

Marcelo Arenas, Pablo Barceló, Leonid Libkin,  
Wim Martens, Andreas Pieris

# Principles of Databases

(Preliminary Version)

July 8, 2022

Santiago Paris  
Bayreuth Edinburgh



---

## Contents

<b>1</b>	<b>Introduction</b>	1
<b>2</b>	<b>Background</b>	5
<hr/>		
<b>Part I The Relational Model: The Classics</b>		
<hr/>		
<b>3</b>	<b>First-Order Logic</b>	23
<b>4</b>	<b>Relational Algebra</b>	31
<b>5</b>	<b>Relational Algebra and SQL</b>	39
<b>6</b>	<b>Equivalence of Logic and Algebra</b>	47
<b>7</b>	<b>First-Order Query Evaluation</b>	51
<b>8</b>	<b>Static Analysis</b>	57
<b>9</b>	<b>Homomorphisms</b>	63
<b>10</b>	<b>Functional Dependencies</b>	67
<b>11</b>	<b>Inclusion Dependencies</b>	75
	<b>Exercises</b>	85
	<b>Bibliographic Comments</b>	87
<hr/>		
<b>Part II Conjunctive Queries</b>		
<hr/>		
<b>12</b>	<b>Syntax and Semantics</b>	91

<b>13 Homomorphisms and Expressiveness</b>	99
<b>14 Query Evaluation</b>	107
<b>15 Containment and Equivalence</b>	115
<b>16 Minimization</b>	121
<b>17 Containment Under Integrity Constraints</b>	129
<b>Exercises</b>	139
<b>Bibliographic Comments</b>	143
<hr/>	
<b>Part III Fast Conjunctive Query Evaluation</b>	
<hr/>	
<b>18 Acyclicity of Conjunctive Queries</b>	149
<b>19 Efficiently Evaluating Boolean ACQs</b>	157
<b>20 Efficiently Evaluating General ACQs</b>	165
<b>21 Treewidth</b>	169
<b>22 Generalized Hypertreewidth</b>	179
<b>23 The Necessity of Bounded Treewidth</b>	187
<b>24 Approximations of Conjunctive Queries</b>	195
<b>25 Bounding the Join Size</b>	203
<b>26 Worst-Case Optimal Join Algorithms</b>	219
<b>27 Leapfrog Triejoin</b>	227
<b>Exercises</b>	241
<b>Bibliographic Comments</b>	245
<hr/>	
<b>Part IV Expressive Languages</b>	
<hr/>	
<b>28 Unions of Conjunctive Queries</b>	249
<b>29 Static Analysis of Unions of Conjunctive Queries</b>	259

<b>30 Unions of Conjunctive Queries with Inequalities</b>	269
<b>31 The Limits of First-Order Queries: Recursion</b>	281
<b>32 The Limits of First-Order Queries: Counting</b>	287
<b>33 Adding Aggregates and Grouping</b>	295
<b>34 Aggregates, Grouping, and Locality</b>	305
<b>35 Adding Recursion: Datalog</b>	317
<b>36 Expressiveness of Datalog Queries</b>	327
<b>37 Datalog Query Evaluation</b>	335
<b>38 Static Analysis of Datalog Queries</b>	343
<b>Exercises</b>	355
<b>Bibliographic Comments</b>	361
<hr/>	
<b>Part V Uncertainty</b>	
<hr/>	
<b>39 Incomplete Databases</b>	365
<b>40 Computing Certain Answers</b>	375
<b>41 Probabilistic Databases</b>	385
<b>42 Inconsistent Databases</b>	395
<b>43 Knowledge-Enriched Databases</b>	405
<b>Exercises</b>	407
<b>Bibliographic Comments</b>	411
<hr/>	
<b>Part VI Query Answering Paradigms</b>	
<hr/>	
<b>44 Bag Semantics</b>	415
<b>45 Incremental Maintenance of Queries</b>	417
<b>46 Provenance Computation</b>	419

<b>47 Top-k Algorithms</b>	421
<b>48 Distributed Evaluation with One Round</b>	423
<b>49 Enumeration and Constant Delay</b>	425
<b>Exercises</b>	427
<b>Bibliographic Comments</b>	429
<hr/>	
<b>Part VII Mappings and Views</b>	
<hr/>	
<b>50 Query Answering using Views</b>	433
<b>51 Determinacy and Rewriting</b>	435
<b>52 Mappings and Data Exchange</b>	437
<b>53 Query Answering for Data Exchange</b>	439
<b>54 Ontology-Based Data Access</b>	441
<b>Exercises</b>	443
<b>Bibliographic Comments</b>	445
<hr/>	
<b>Background for Tree- and Graph-Structured Data</b>	
<hr/>	
<b>55 Background For Tree and Graph Structured Data</b>	449
<hr/>	
<b>Part VIII Tree-Structured Data</b>	
<hr/>	
<b>56 Data Model</b>	463
<b>57 First-Order Logic over Trees</b>	473
<b>58 XPath</b>	479
<b>59 Expressiveness of XPath</b>	485
<b>60 Static Analysis of XPath</b>	491
<b>61 Tree Pattern Queries</b>	495
<b>62 Tree Pattern Query Containment and Equivalence</b>	501

<b>63 Tree Pattern Query Minimization</b>	509
<b>Exercises</b>	515
<b>Bibliographic Comments</b>	517
<hr/>	
<b>Part IX Expressive Languages for Tree-Structured Data</b>	
<b>64 Monadic Second Order Logic</b>	523
<b>65 Tree Automata</b>	527
<b>66 Monadic Datalog</b>	535
<b>67 Schemas for XML</b>	537
<b>68 Static Analysis Under Schema Constraints</b>	547
<b>69 Data Trees</b>	555
<b>70 Static Analysis on Data Trees</b>	557
<b>Exercises</b>	559
<b>Bibliographic Comments</b>	561
<hr/>	
<b>Part X Graph-Structured Data</b>	
<b>71 Data Model and Queries</b>	567
<b>72 Graph Query Evaluation</b>	573
<b>73 Containment</b>	577
<b>74 Querying Property Graphs</b>	583
<b>75 RDF and SPARQL</b>	589
<b>Exercises</b>	591
<b>Bibliographic Comments</b>	593
<b>References</b>	595

---

**Part XI Appendix: Theory of Computation**

---

<b>Big-O Notation</b> .....	599
<b>Turing Machines and Complexity Classes</b> .....	601
<b>Input Encodings</b> .....	611
<b>Tiling Problems</b> .....	615
<b>Formal Languages</b> .....	623



## Introduction

This is a very preliminary partial release of the upcoming book “Principles of Databases”, which will be about the foundational and mathematical principles of databases in their various forms. The early parts focus on an overview of the relational model, and on processing some of the most commonly occurring relational queries. Forthcoming parts will focus on additional aspects of the relational model and will cover tree-structured and graph-structured data as well.

The general philosophy of the book is the following:

- We planned the book such that large parts of it are suitable for teaching. A chapter roughly corresponds to the contents of a single lecture.
- For the ease of teaching and understanding the material, we may sometimes cut corners. If we want to give the reader a relatively quick insight of a particular result, this sometimes means that we present a weaker form of the result than the most general result known in the literature.

We have been teaching from this book ourselves, but the present version will undoubtedly still have errors. If you find any errors in the book, or places that you find particularly unclear, please let us know through the repository: <https://github.com/pdm-book/community>. The new versions of the book, including corrections, will be published in this repository.

## What is Planned

The finished book will consist of the following parts:

- (I) The Relational Model: The Classics
- (II) Conjunctive Queries
- (III) Fast Conjunctive Query Evaluation

- Includes material on acyclic queries, treewidth and hypertreewidth, and worst-case optimal join algorithms.

## (IV) Expressive Languages

- Includes material on adding features found in most commonly used query languages: union, negation, aggregates, and recursion.

## (V) Uncertainty

- Includes material on incomplete information, probabilistic databases, consistent query answering, and query answering in the presence of ontologies.

## (VI) Query Answering Paradigms

- Includes material on bag semantics, incremental maintenance, provenance, top-k queries, distributed evaluation, and constant delay query evaluation.

## (VII) Mappings and Views

- Includes material on determinacy, data exchange, and ontology-based data access.

## (VIII) Tree-Structured Data

- Includes material on first-order logic on trees, XPath, and tree pattern queries, and their evaluation- and static analysis problems.

## (IX) Expressive Languages for Tree-Structured Data

- Includes material on MSO, tree automata, monadic datalog, schema languages, and their static analysis.

## (X) Graph-Structured Data

- Includes material on various types of graph queries, their evaluation and containment, property graphs, RDF, and SPARQL.

We will continue to release parts, not necessarily in the order presented here. Furthermore, the ordering and contents of the chapters is preliminary and may change in future versions.

## What is Still Missing, Even From Released Parts

Let's start by saying what is mainly there: in every chapter that we release, we believe that the technical content is relatively stable. For every part that we have released, two things still need work though:

**Exercises.** We have generated some initial ideas for exercises, but we are aware that the exercises for the currently released parts still need work. In fact, we are open to exercise suggestions.

**Bibliography.** We plan to accompany each part with references and a bibliographic discussion. These are not implemented yet, even for Parts II-IV.

## Proofreaders

The present version has benefited from valuable comments of (in alphabetical order): Johannes Doleschal, Matthias Niewerth, Thomas Schwentick



## Background

In this chapter, we introduce the mathematical concepts and terminology that will be used throughout the book. These include:

- the relational model,
- queries and query languages, and
- computational problems central in the study of principles of databases.

### Basic Notions and Notation

We begin with a brief discussion of the very basic mathematical notions and notation that we are going to use in this book.

#### *Sets*

A *set* contains a finite or infinite number of elements (e.g., numbers, symbols, other sets), without repetition or respect to order. The elements in a set  $S$  are the *members* of  $S$ . We use the symbols  $\in$  and  $\notin$  to denote set membership and nonmembership, respectively. For a finite set  $S$ , we write  $|S|$  for its *cardinality*, that is, the number of elements in it. The set without elements is called the *empty set*, written as  $\emptyset$ .

Given two (finite or infinite) sets  $S$  and  $T$ , we write:

- $S \cup T$  for their *union*  $\{a \mid a \in S \text{ or } a \in T\}$ ,
- $S \cap T$  for their *intersection*  $\{a \mid a \in S \text{ and } a \in T\}$ , and
- $S - T$  for their *difference*  $\{a \mid a \in S \text{ and } a \notin T\}$ .

We further say that

- $S$  is *equal* to  $T$ , written  $S = T$ , when  $x \in S$  if and only if  $x \in T$ ,

- $S$  is a *subset* of  $T$ , written  $S \subseteq T$ , when  $x \in S$  implies  $x \in T$ , and
- $S$  is a *proper* (or *strict*) subset of  $T$ , written  $S \subsetneq T$ , if  $S \subseteq T$  and  $S \neq T$ .

We write  $\mathcal{P}(S)$  for the *powerset* of  $S$ , that is, the set consisting of all the subsets of  $S$ . Analogously, we write  $\mathcal{P}_{\text{fin}}(S)$  for the *finite powerset* of  $S$ , namely the set consisting of all the finite subsets of  $S$ .

We write  $\mathbb{N}$  for the set  $\{0, 1, 2, \dots\}$  of natural numbers. For  $i, j \in \mathbb{N}$ , we denote by  $[i, j]$  the set  $\{k \in \mathbb{N} \mid i \leq k \text{ and } k \leq j\}$ . We simply write  $[i]$  for  $[1, i]$ . We write  $\mathbb{Q}$  for the set of rational numbers and  $\mathbb{Q}_{\geq 0}$  for the set of nonnegative rational numbers.

### Sequences and Tuples

A *sequence* of elements is a list of these elements in some order. We typically identify a sequence by writing the list within parentheses. Recall that in a set the order does not matter, but in a sequence it does. Hence, the sequence  $(1, 2, 3)$  is not the same as  $(3, 2, 1)$ . Similarly, repetition does not matter in a set, but it does matter in a sequence. Thus, the sequence  $(1, 1, 2, 3)$  is different than  $(1, 2, 3)$ , while the set  $\{1, 1, 2, 3\}$  is the same as  $\{1, 2, 3\}$ . Finite sequences are called *tuples*. A sequence with  $k \in \mathbb{N}$  elements is a tuple of *arity*  $k$ , called *k-ary tuple* (or simply *k-tuple*). Note that when  $k = 0$  we get the empty tuple  $()$ . We often abbreviate a  $k$ -ary tuple  $(a_1, \dots, a_k)$  as  $\bar{a}$ . Moreover, for a  $k$ -ary tuple  $\bar{a}$ , we usually assume that its elements are  $(a_1, \dots, a_k)$ . We say that  $\bar{a} = (a_1, \dots, a_k)$  has the *positions*  $1, \dots, k$  and that an element  $b$  *occurs* at position  $i$  if  $b = a_i$ . For example, 1 occurs at positions 1 and 3 in the tuple  $(1, 2, 1, 4)$ . Conversely,  $\bar{a}$  *mentions*  $a$  if  $a \in \{a_1, \dots, a_k\}$ .

For two sets  $S, T$ , we write  $S \times T$  for the set of all pairs  $(a, b)$ , where  $a \in S$  and  $b \in T$ , called the *Cartesian product* or *cross product* of  $S$  and  $T$ . We can also define the Cartesian product of  $k \geq 1$  sets  $S_1, \dots, S_k$ , known as the *k-fold Cartesian product*, which is the set of all tuples  $(a_1, \dots, a_k)$ , where  $a_i \in S_i$  for each  $i \in [k]$ . For the  $k$ -fold Cartesian product of a set  $S$  with itself we write

$$S^k = \underbrace{S \times \dots \times S}_k.$$

A *k-ary mathematical relation* is a set of  $k$ -ary tuples. Let  $T$  be a binary mathematical relation and let  $S = \{a \mid \text{there exists } b \text{ such that } (a, b) \in T \text{ or } (b, a) \in T\}$ . We say that

- $T$  is *transitive* if, for all  $a, b, c \in S$ , we have that  $(a, b) \in T$  and  $(b, c) \in T$  implies that  $(a, c) \in T$ ;
- $T$  is *reflexive* if  $(a, a) \in T$  for every  $a \in S$ ;
- $T$  is *antisymmetric* if, for all  $a, b \in S$ , we have that  $(a, b) \in T$  and  $(b, a) \in T$  implies that  $a = b$ ;
- $T$  is *total*, if, for all  $a, b \in S$ , we have that  $(a, b) \in T$  or  $(b, a) \in T$ ;

- $T$  is a *partial order* if it is transitive, reflexive, and antisymmetric; and
- $T$  is a *total order* if it is a partial order and total.

We also define the following standard operations on  $T$ .

- The *transitive closure* of  $T$ , denoted by  $T^+$ , is defined as  $\{(a, b) \mid (a, b) \in T \text{ or there exist } a_1, \dots, a_n \text{ such that } n \geq 1, (a, a_1) \in T, (a_n, b) \in T \text{ and } (a_i, a_{i+1}) \in T \text{ for every } i \in [n-1]\}$ .
- The *reflexive transitive closure* of  $T$ , denoted by  $T^*$ , is defined as  $T^+ \cup \{(a, a) \mid a \in S\}$ .

Finally, we say that  $T$  is a *successor relation* if

- for every  $a \in S$ , there exists at most one  $b \in S$  such that  $(a, b) \in T$ , and
- $T^*$  is a total order.

### Functions

Consider two (finite or infinite) sets  $S$  and  $T$ . A *function*  $f$  from  $S$  to  $T$ , written  $f : S \rightarrow T$ , is a mapping from (all or some) elements of  $S$  to elements of  $T$ , i.e., for every  $a \in S$ , either  $f(a) \in T$ , in which case we say  $f$  is defined on  $a$ , or  $f(a)$  is undefined, such that the following holds: for every  $a, b \in S$  on which  $f$  is defined,  $a = b$  implies  $f(a) = f(b)$ . We call  $f$  *total* if it is defined on every element of  $S$ ; otherwise, it is called *partial*. By default, we assume functions to be total. When a function  $f$  is partial, we explicitly say this, and write  $\text{Dom}(f)$  for the set of elements from  $S$  on which  $f$  is defined.

We say that a function  $f : S \rightarrow T$  is

- *injective* (or *one-to-one*) if  $a \neq b$  implies  $f(a) \neq f(b)$  for every  $a, b \in S$ ,
- *surjective* (or *onto*) if, for every  $b \in T$ , there is  $a \in S$  such that  $f(a) = b$ ,
- *bijective* (or *one-to-one correspondence*) if it is injective and surjective.

A useful notion is that of composition of functions. Given two functions  $f : S \rightarrow T$  and  $g : T \rightarrow U$ , the *composition* of  $f$  and  $g$ , denoted  $g \circ f$ , is the function from  $S$  to  $U$  defined as follows:  $g \circ f(a) = g(f(a))$  for every  $a \in S$ . Another useful notion is that of union of functions. Given two functions  $f : S \rightarrow T$  and  $g : S' \rightarrow T'$  with  $f(a) = g(a)$  for every  $a \in S \cap S'$ , the *union* of  $f$  and  $g$ , denoted  $f \cup g$ , is the function from  $S \cup S'$  to  $T \cup T'$  defined as follows:  $f \cup g(a) = f(a)$  for every  $a \in S$ , and  $f \cup g(a) = g(a)$  for every  $a \in S'$ .

Given a function  $f : S \rightarrow T$ , for brevity, we will use the same letter  $f$  to denote extensions of  $f$  on more complex objects (such as tuples of elements of  $S$ , sets of elements of  $S$ , etc.). More precisely, if  $\bar{a} = (a_1, \dots, a_k) \in S^k$ , then  $f(\bar{a}) = (f(a_1), \dots, f(a_k))$ . If  $R \subseteq S$ , then  $f(R) = \{f(a) \mid a \in R\}$ . Notice that this convention also extends further, e.g., to sets of sets of tuples.

*Graphs and Trees*

A *graph* is a pair  $(V, E)$ , where  $V$  is a finite set of *nodes*, and  $E$  is a finite set of *edges*. We distinguish between *directed* and *undirected* graphs, which only differ in the way their set  $E$  of edges is defined:

- In *directed graphs* we have that  $E \subseteq V \times V$ .
- In *undirected graphs* we have that  $E \subseteq \{S \subseteq V \mid |S| = 1 \text{ or } |S| = 2\}$ .

In simple words, in directed graphs edges are pairs of nodes, and the order in which the nodes are given indicates the direction of the edge, that is,  $(u, v)$  is an edge from  $u$  to  $v$ , whereas  $(v, u)$  is an edge from  $v$  to  $u$ . In undirected graphs, edges are simply subsets of  $V$  of cardinality 2, and the order in which the nodes are given does not matter, i.e.,  $\{u, v\}$  and  $\{v, u\}$  is the same (undirected) edge between the nodes  $u$  and  $v$ .

A graph  $G_1 = (V_1, E_1)$  is a *subgraph* of a graph  $G_2 = (V_2, E_2)$  if  $V_1 \subseteq V_2$  and  $E_1 \subseteq E_2$ . Moreover, if both  $G_1$  and  $G_2$  are directed graphs, then  $G_1$  is the *subgraph of  $G_2$  induced by  $V_1$*  if  $E_1 = E_2 \cap (V_1 \times V_1)$ , and if both  $G_1$  and  $G_2$  are undirected graphs, then  $G_1$  is the *subgraph of  $G_2$  induced by  $V_1$*  if  $E_1 = E_2 \cap \{S \subseteq V_1 \mid |S| = 1 \text{ or } |S| = 2\}$ .

An *undirected path* in a (directed or undirected) graph  $G = (V, E)$  is a non-empty sequence of nodes  $\pi = u_0 u_1 \cdots u_k$  for  $k \geq 0$ , where

- if  $G$  is a directed graph, then  $(u_{i-1}, u_i) \in E$  or  $(u_i, u_{i-1}) \in E$  for every  $i \in [k]$ , and
- if  $G$  is an undirected graph, then  $\{u_{i-1}, u_i\} \in E$  for every  $i \in [k]$ .

We say that  $\pi$  is *path from  $u_0$  to  $u_k$*  and has *length  $k$* . (The path of length zero is from  $u_0$  to  $u_0$ .) A graph is *connected* if, for every pair of nodes  $u, v \in V$ , there is an undirected path from  $u$  to  $v$ .

A *path* in a directed graph  $G = (V, E)$  is a non-empty sequence of nodes  $\pi = u_0 u_1 \cdots u_k$  for  $k \geq 0$ , where  $(u_{i-1}, u_i) \in E$  for every  $i \in [k]$ . Similar to undirected paths, we call  $\pi$  is a *path from  $u_0$  to  $u_k$*  and its *length* is  $k$ .

A connected directed graph  $T = (V, E)$  is a *tree* if

- for every node  $u$ , there is at most one node  $v$  with  $(v, u) \in E$ , called the *parent* of  $u$ , and
- there is exactly one node  $u$ , called the *root* of  $T$ , without a parent.

The notions of subtree and induced subtree are naturally derived from the notions of subgraph and induced subgraph. Moreover, if  $v$  is the parent of  $u$  in a tree  $T$ , then  $u$  is called a child of  $v$ , and if  $v$  does not have any children, then  $v$  is called a leaf of  $T$ . Finally, if there exists a path from  $v$  to  $u$  in  $T$ , then  $v$  is said to be an *ancestor* of  $u$ , and  $u$  is said to be a *descendant* of  $v$ .



It is common in Computer Science (unlike other fields such as Biology) to graphically depict trees with their root on top and the edges directed downwards. Therefore, even if we deal with directed trees, we may omit the direction of the edges in figures, which is typically specified via an arrow.

## The Relational Model

To define tables in real-life databases, for example, by the `create table` statements of SQL, one needs to specify their names and names of their attributes. Therefore, to model databases, we need two disjoint sets

Rel of *relation names*      and      Att of *attribute names*.

We assume that these sets are countably infinite in order to ensure that we never run out of ways to name new tables and their attributes. In practice, of course, these sets are finite but extremely large: they are strings that can be so large that one never really runs out of names. Theoretically, we model this by assuming that these sets are countably infinite.

In `create table` declarations, one specifies types of attributes as well, for example, integer, Boolean, string. In the study of the theoretical foundations of databases, one typically does not make this distinction, and assumes that all elements populating databases come from another countably infinite set

Const of *values*.

This simplifying assumption does not affect the various results on the complexity of query evaluation, expressiveness of languages, equivalence of queries, and many other subjects studied in this book. At the same time, it brings the setting closer to that of mathematical logic, allowing us to borrow many tools from it. It also allows us to significantly streamline notations.

### *The Named and Unnamed Perspective*

There exist two standard perspectives from which databases can be defined, called the *named* and the *unnamed* perspectives. While the named perspective is closer to how databases appear in database management systems, and therefore more natural when giving examples, the unnamed perspective provides a clean mathematical model that is easier to use for studying the principles of databases. Importantly, the modeling power of those two perspectives is exactly the same, which allows us to go back and forth between the two.

**Named Perspective.** Under the *named perspective*, attribute names are viewed as an explicit part of a database. More precisely, a *database tuple* is a function  $t : U \rightarrow \text{Const}$ , where  $U = \{A_1, \dots, A_k\}$  is a finite subset of Att. The *sort* of  $t$  is  $U$ , and its *arity* is the cardinality  $|U|$  of  $U$ ; we

say that  $t$  is  $k$ -ary if  $|U| = k$ . We usually do not use the function notation for database tuples in the named perspective, and denote them as  $t = (A_1 : a_1, \dots, A_k : a_k)$ , meaning that  $t(A_i) = a_i$  for every  $i \in [k]$ . Notice that, according to this notation,  $(A_1 : a_1, A_2 : a_2)$  and  $(A_2 : a_2, A_1 : a_1)$  represent the same function  $t$ . A *relation instance* in the named perspective is a *finite* set  $S$  of database tuples of the same sort  $U$ , which we also call the sort of the relation instance  $S$  and denote by  $\text{sort}(S)$ . By  $\text{nRI}$  (for *named relational instances*) we denote the set of all such relation instances. A *possibly infinite relation instance* in the named perspective is defined as the notion of relation instance, but without forcing it to be finite. We write  $\text{nRI}^\infty$  for the set of all possibly infinite relation instances in the named perspective.

Database systems usually use a *database schema* that associates attribute names to relation names. This can be formalized as follows.

**Definition 2.1: Named Database Schema**

A *named (database) schema* is a partial function

$$\mathbf{S} : \text{Rel} \rightarrow \mathcal{P}_{\text{fin}}(\text{Att})$$

such that  $\text{Dom}(\mathbf{S})$  is finite. For  $R \in \text{Dom}(\mathbf{S})$ , the *sort of  $R$  under  $\mathbf{S}$*  is the set  $\mathbf{S}(R)$ . The *arity of  $R$  under  $\mathbf{S}$* , denoted  $\text{ars}(\mathbf{S}(R))$ , is  $|\mathbf{S}(R)|$ .

In other words, a named database schema  $\mathbf{S}$  provides a finite set of relation names, together with their (finitely many) attribute names. These attribute names form the sort of the relation names under  $\mathbf{S}$ , and their number specifies the arity of the relation names under  $\mathbf{S}$ . For arities 1, 2, and 3, we speak of unary, binary, and ternary relation names, respectively. We now introduce the notion of database instance of a named schema.

**Definition 2.2: Database Instance (The Named Case)**

A *database instance*  $D$  of a named schema  $\mathbf{S}$  is a function

$$D : \text{Dom}(\mathbf{S}) \rightarrow \text{nRI}$$

such that  $\text{sort}(D(R)) = \mathbf{S}(R)$ , for every  $R \in \text{Dom}(\mathbf{S})$ .

We can also talk about possibly infinite database instances. Formally, a *possibly infinite database instance*  $D$  of a named schema  $\mathbf{S}$  is a function

$$D : \text{Dom}(\mathbf{S}) \rightarrow \text{nRI}^\infty$$

such that  $\text{sort}(D(R)) = \mathbf{S}(R)$ , for every  $R \in \text{Dom}(\mathbf{S})$ . This means that  $D$  is either finite as in Definition 2.2, where each relation name of  $\text{Dom}(\mathbf{S})$  is

mapped to a finite relation instance, or infinite in the sense that at least one relation name of  $\text{Dom}(\mathbf{S})$  is mapped to an infinite relation instance. Infinite database instances are obviously not a real-life concept, and we are not interested in studying them per se. Having said that, they are a very useful mathematical tool as they allow us to prove some results in a more elegant way. In other words, infinite database instances are considered for purely technical reasons, which will be revealed later in the book.

To avoid heavy notation, and because the name  $\mathbf{S}$  of a schema is often not important, we usually provide schema information without explicitly using the symbol  $\mathbf{S}$ . We write  $R[A_1, \dots, A_k]$  instead of  $\mathbf{S}(R) = \{A_1, \dots, A_k\}$  for the schema  $\mathbf{S}$  in question. For example, we write

`City[city_id, name, country]`

to refer to a relation name `City` with attribute names `city_id`, `name`, and `country`. Likewise, we write  $\text{ar}(R)$  instead of  $\text{arg}_{\mathbf{S}}(R)$ . We may even write  $R[k]$  to indicate that the arity of  $R$  under the schema in question is  $k$ .

**Unnamed Perspective.** Under the *unnamed perspective*, a *database tuple* is an element of  $\text{Const}^k$  for some  $k \in \mathbb{N}$ . We denote such tuples using lowercase letters from the beginning of the alphabet, that is, as  $(a_1, \dots, a_k)$ ,  $(b_1, \dots, b_k)$ , etc., or even more succinctly as  $\bar{a}, \bar{b}$ , etc. A *relation instance* in the unnamed perspective is a *finite* set  $S$  of database tuples of the same arity  $k$ . We say that  $k$  is the *arity* of  $S$ , denoted by  $\text{ar}(S)$ . By  $\text{uRI}$  (for *unnamed relation instances*) we denote the set of all such relation instances. A *possibly infinite relation instance* in the unnamed perspective is defined as the notion of relation instance, but without forcing it to be finite. We write  $\text{uRI}^\infty$  for the set of all possibly infinite relation instances in the unnamed perspective. The notion of unnamed database schema follows.

### Definition 2.3: Unnamed Database Schema

An *unnamed (database) schema* is a partial function

$$\mathbf{S} : \text{Rel} \rightarrow \mathbb{N}$$

such that  $\text{Dom}(\mathbf{S})$  is finite. For a relation name  $R \in \text{Dom}(\mathbf{S})$ , the *arity of  $R$  under  $\mathbf{S}$* , denoted  $\text{arg}_{\mathbf{S}}(R)$ , is defined as  $\mathbf{S}(R)$ .

In simple words, an unnamed databases schema  $\mathbf{S}$  provides a finite set of relation names from  $\text{Rel}$ , together with their arity. We proceed to introduce the notion of database instance of an unnamed database schema.

### Definition 2.4: Database Instance (The Unnamed Case)

A *database instance*  $D$  of an unnamed schema  $\mathbf{S}$  is a function

$$D : \text{Dom}(\mathbf{S}) \rightarrow \text{uRI}$$

such that  $\text{ar}(D(R)) = \text{ar}_{\mathbf{S}}(R)$ , for every  $R \in \text{Dom}(\mathbf{S})$ .

Analogously, a *possibly infinite database instance*  $D$  of an unnamed schema  $\mathbf{S}$  is defined as a function of the form

$$D : \text{Dom}(\mathbf{S}) \rightarrow \text{uRI}^\infty$$

such that  $\text{ar}(D(R)) = \text{ar}_{\mathbf{S}}(R)$ , for every  $R \in \text{Dom}(\mathbf{S})$ . Recall that infinite database instances are considered for purely technical reasons. As in the named perspective, in order to avoid heavy notation, we write  $\text{ar}(R)$  instead of  $\text{ar}_{\mathbf{S}}(R)$  for the arity of  $R$  under  $\mathbf{S}$ . We may even write  $R[k]$  to indicate that the arity of  $R$  under the schema in question is  $k$ .

For a (named or unnamed) schema  $\mathbf{S}$ , we write  $\text{Inst}(\mathbf{S})$  for the set of all database instances of  $\mathbf{S}$ . Notice that  $\text{Inst}(\mathbf{S})$  does not contain infinite database instances. We also need the crucial notion of the active domain of a (possibly infinite) database instance, which is, roughly speaking, the set of constants that occur in it. Under the named perspective, we say that a database tuple  $t : U \rightarrow \text{Const}$  *mentions* a constant  $a \in \text{Const}$  if there exists  $A \in U$  such that  $t(A) = a$ . Under the unnamed perspective, a database tuple  $(a_1, \dots, a_k) \in \text{Const}^k$  *mentions*  $a \in \text{Const}$  if there exists  $i \in [k]$  such that  $a_i = a$ . The *active domain* of a (possibly infinite) database instance  $D$  of  $\mathbf{S}$  is defined as the set

$$\{a \in \text{Const} \mid \text{there exists } R \in \text{Dom}(\mathbf{S}) \text{ such that } D(R) \text{ contains a database tuple that mentions } a\}.$$

Henceforth, for brevity, we simply refer to the domain instead of the active domain of  $D$ , and denote it  $\text{Dom}(D)$ . Let us stress that this simplification leads to a clash of terminology. Indeed, given a function  $f : S \rightarrow T$ , the set  $S$  is typically called the domain of  $f$ . Therefore, since  $D$  is formally defined as a function, one may think that the domain of  $D$  is the domain of  $D$  seen as a function, which is a finite set of relation names. We will *never* use the term domain, and the notation  $\text{Dom}(D)$ , to refer to the domain of the function  $D$ .

### *Simplified Terminology and Notation*

We will refer to a (possibly infinite) database instance as a *(possibly infinite) database*, to a relation instance as a *relation*, and to a database tuple as a *tuple*. By abuse of terminology, we will also refer to a mathematical relation as a *relation*, but it will always be clear from the context whether we mean a relation instance or a mathematical relation.

In both the named and the unnamed perspectives, we will write  $R_i^D$  instead of  $D(R_i)$ . When it is clear from the context, we shall omit the superscript  $D$ ,

and simply write  $R_i$  instead of  $R_i^D$ . This means that we will effectively use the same notation for relation names and for relation instances. This is a common practice that is used to simplify notation, and it will never lead to confusion; when the instance is important, we will make it explicit.

Although database schemas are formally defined as partial functions, with their domain being a finite subset of  $\mathbf{Rel}$ , it is often convenient to treat them as sets of relation names. Thus, we will usually treat a schema  $\mathbf{S}$  as the finite set  $\text{Dom}(\mathbf{S})$ . This means that whenever we write  $\mathbf{S} = \{R_1, \dots, R_n\}$ , we actually mean that  $\text{Dom}(\mathbf{S}) = \{R_1, \dots, R_n\}$ . In the unnamed case, we may also write

$$\mathbf{S} = \{R_1[k_1], \dots, R_n[k_n]\}$$

for the fact that  $\text{Dom}(\mathbf{S}) = \{R_1, \dots, R_n\}$  and  $\mathbf{S}(R_i) = k_i$ , for each  $i \in [n]$ . Having this notation for schemas, we can then take, e.g., the union  $\mathbf{S}_1 \cup \mathbf{S}_2$  of two schemas  $\mathbf{S}_1$  and  $\mathbf{S}_2$  (providing that  $\text{Dom}(\mathbf{S}_1)$  and  $\text{Dom}(\mathbf{S}_2)$  are disjoint).

Analogously, databases can be seen as sets, in particular, as sets of facts. In the unnamed perspective, for a  $k$ -ary relation name  $R$ , and a tuple  $\bar{a} \in \text{Const}^k$ , we call  $R(\bar{a})$  a *fact*. Since a fact is always a statement about a single tuple, we simplify the notation  $R((a_1, \dots, a_k))$  to  $R(a_1, \dots, a_k)$ . We will usually treat a (possibly infinite) database  $D$  of an unnamed schema  $\mathbf{S}$  as the set of facts

$$\{R(\bar{a}) \mid R \in \mathbf{S} \text{ and } \bar{a} \in R^D\}.$$

For example, we can write  $D = \{R_1(a, b), R_1(b, c), R_2(a, c, d)\}$  as a shorthand for  $R_1^D = \{(a, b), (b, c)\}$  and  $R_2^D = \{(a, c, d)\}$ . Note that the active domain of  $D$  is precisely the set of constants occurring in  $\{R(\bar{a}) \mid R \in \mathbf{S} \text{ and } \bar{a} \in R^D\}$ .

#### *Named versus Unnamed Perspective*

There is clearly a close connection between the two perspectives, which is not surprising since both are mathematical abstractions of the same concept. A (possibly infinite) database of a named schema can be transformed into a semantically equivalent one of an unnamed schema, and vice versa. By semantically equivalent, we mean databases that are essentially the same modulo representation details. It is instructive to properly formalize this connection, which will be used throughout the book. We do this for databases, but the exact same constructions work also for possibly infinite databases.

**From Named to Unnamed.** Consider a named schema  $\mathbf{S}$ , and assume that there is an ordering  $\prec$  on the set of relation-attribute pairs  $\{(R, A) \mid R \in \text{Dom}(\mathbf{S}) \text{ and } A \in \mathbf{S}(R)\}$ . We define the unnamed schema  $\mathbf{S}' : \mathbf{Rel} \rightarrow \mathbb{N}$  as follows:  $\text{Dom}(\mathbf{S}') = \text{Dom}(\mathbf{S})$ , and  $\mathbf{S}'(R) = \arg(R)$  for every  $R \in \text{Dom}(\mathbf{S})$ . Moreover, for every database  $D$  of  $\mathbf{S}$ , a semantically equivalent database  $D' : \text{Dom}(\mathbf{S}') \rightarrow \mathbf{uRI}$  of  $\mathbf{S}'$  is defined as follows: for every  $R \in \text{Dom}(\mathbf{S}')$ ,

$$D'(R) = \{(a_1, \dots, a_k) \mid (A_1 : a_1, \dots, A_k : a_k) \in D(R) \\ \text{such that } (R, A_1) \prec (R, A_2) \prec \dots \prec (R, A_k)\}.$$

**From Unnamed to Named.** Consider an unnamed database schema  $\mathbf{S}$ . We assume that  $\text{Att}$  contains an attribute name  $\#_i$  for each  $i \geq 1$ . We define the named schema  $\mathbf{S}' : \text{Rel} \rightarrow \mathcal{P}_{\text{fin}}(\text{Att})$  as follows:  $\text{Dom}(\mathbf{S}') = \text{Dom}(\mathbf{S})$ , and  $\mathbf{S}'(R) = \{\#_1, \dots, \#_{\text{ars}(\mathbf{S}(R))}\}$  for every  $R \in \text{Dom}(\mathbf{S})$ . Moreover, for every database  $D$  of  $\mathbf{S}$ , a semantically equivalent database  $D' : \text{Dom}(\mathbf{S}') \rightarrow \text{nRI}$  of  $\mathbf{S}'$  is defined as follows: for every  $R \in \text{Dom}(\mathbf{S}')$ ,

$$D'(R) = \{(\#_1 : a_1, \dots, \#_k : a_k) \mid (a_1, \dots, a_k) \in D(R)\} \quad \boxed{1}$$

Since the above connection between the two perspectives is useful in many places in the book, we assume from now on that, whenever a named database schema is used, the ordering  $<$  on relation-attribute pairs is available.

The unnamed perspective is usually mathematically more elegant, while the named perspective is closer to practice. Therefore, we often define notions in the book using the unnamed perspective, but illustrate them with examples using the named perspective. When we do so, we use the following convention. When we denote a relation name as  $R[A, B, \dots]$  of a named database schema  $\mathbf{S}$  in an example, we assume that the ordering of attributes in  $\mathbf{S}$  is consistent with how we write it in the example, that is,  $(R, A) < (R, B)$ , etc. This allows us to easily switch between the named and unnamed perspective in examples, e.g., by being able to say that the “first” attribute of  $R$  is  $A$ .

## Queries and Query Languages

Queries will appear throughout the book as both *semantic* and *syntactic* objects. As a semantic object, a query  $q$  over a schema  $\mathbf{S}$  is a function that maps databases of  $\mathbf{S}$  to *finite* sets of tuples of the same arity over  $\text{Const}$ .

### Definition 2.5: Queries and Query Languages

Consider a database schema  $\mathbf{S}$ . A *query of arity  $k \geq 0$*  (or simply a  *$k$ -ary query*) over  $\mathbf{S}$  is a function of the form

$$q : \text{Inst}(\mathbf{S}) \rightarrow \mathcal{P}_{\text{fin}}(\text{Const}^k).$$

A *query language* is a set of queries.

An important subject, which will be considered in the book, is to classify query languages according to their expressive power. Two query languages  $\mathcal{L}_1$  and  $\mathcal{L}_2$  are *equally expressive* if  $\mathcal{L}_1 = \mathcal{L}_2$ . Furthermore,  $\mathcal{L}_1$  is *more expressive* than  $\mathcal{L}_2$  if  $\mathcal{L}_2 \subseteq \mathcal{L}_1$ , and  $\mathcal{L}_1$  is *strictly more expressive* than  $\mathcal{L}_2$  if  $\mathcal{L}_2 \subsetneq \mathcal{L}_1$ .

<sup>1</sup> Notice that under the assumption that  $(R, \#_i) < (R, \#_{i+1})$  for every relation name  $R \in \mathbf{S}$  and  $i \in [\mathbf{S}(R) - 1]$ , one can translate a database  $D$  from the unnamed perspective to the named perspective and back, and obtain  $D$  again.

Of course, queries as semantic objects must be given in some syntax. The syntax of queries could be SQL, relational algebra, first-order logic, and Datalog, to name a few. We proceed to explain some of our notational conventions for queries. For the sake of the discussion, we focus on query languages that are based on logic. To this end, we assume a countably infinite set

$\text{Var}$  of *variables*,

disjoint from  $\text{Const}$ ,  $\text{Rel}$ , and  $\text{Att}$ . If  $\varphi$  is a logical formula and  $\bar{x} = (x_1, \dots, x_k) \in \text{Var}^k$  is a tuple of variables, we will denote queries as  $\varphi(\bar{x})$ . We will also use a letter such as  $q$  to refer to the entire query, that is,  $q = \varphi(\bar{x})$ . The purpose of  $\bar{x}$  is to make clear what is the *output* of the query; we will also write  $q(\bar{x})$  to emphasise that  $q$  has the output tuple  $\bar{x}$ . More precisely, we will always define for a database  $D$  and tuple  $\bar{a} = (a_1, \dots, a_k) \in \text{Const}^k$  whether  $D$  *satisfies*  $\varphi$  *using the values*  $\bar{a}$ , denoted by  $D \models \varphi(\bar{a})$ . Then, with the syntactic object  $q = \varphi(\bar{x})$ , we associate a semantic object that produces an *output*, i.e., a set of  $k$ -ary tuples over  $\text{Const}$ , for each database  $D$ , defined as:

$$q(D) = \{\bar{a} \in \text{Const}^k \mid D \models \varphi(\bar{a})\}.$$

This semantic object will always be a query in the sense of Definition [2.5](#). In other words, we will use the letter  $q$  to refer to both

- the syntactic object denoting a query (for example, a logical formula together with an output tuple), and
- the query itself (i.e., the function that maps databases to finite sets of tuples of the same arity over  $\text{Const}$ ).

A query of arity 0 is called *Boolean*. In this case, there are only two possible outputs: either the singleton set  $\{()\}$  containing the empty tuple, or the empty set  $\{\}$ . We interpret  $\{()\}$  as the Boolean value **true**, and  $\{\}$  as **false**. For readability, we write  $q(D) = \mathbf{true}$  in place of  $q(D) = \{()\}$ , and  $q(D) = \mathbf{false}$  in place of  $q(D) = \{\}$ . When denoting Boolean queries, we will often omit the empty tuple  $()$  from the notation, i.e., write  $q = \varphi$  instead of  $q = \varphi()$ .

Next, we introduce relational atoms, which will be useful throughout the book and, in particular, for defining the syntax of query languages that are based on logic. When  $R$  is a  $k$ -ary relation name and  $\bar{u} \in (\text{Const} \cup \text{Var})^k$ ,  $R(\bar{u})$  is a *relational atom*. Observe that the only difference between a fact and a relational atom is that the former mentions only constants, whereas the latter can mention both constants and variables. As for facts, since a relational atom is always a statement about a single tuple, we simply write  $R(u_1, \dots, u_k)$  instead of  $R((u_1, \dots, u_k))$ . Given a set of atoms  $S$ , we write  $\text{Dom}(S)$  for the set of constants and variables in  $S$ . For example,  $\text{Dom}(\{R(a, x, b), R(x, a, y)\}) = \{a, b, x, y\}$ . We also write  $R^S$  for the set of tuples  $\{\bar{u} \mid R(\bar{u}) \in S\}$ .

## Key Problems: Query Evaluation and Query Analysis

Much of what we do in databases boils down to running queries on a database, or statically analyzing queries. The latter is the basis of query optimization: we need to be able to reason about queries, and to be able to replace a query with a better behaved one that has the same output. We proceed to introduce the main algorithmic problems associated with the above tasks. In their most common form, they are parameterized by a query language  $\mathcal{L}$ .

### *Query Evaluation*

We start with the *query evaluation problem*, or simply the *evaluation problem*, that has the following form:

#### **Problem: $\mathcal{L}$ -Evaluation**

**Input:** A query  $q$  from  $\mathcal{L}$ , a database  $D$ , a tuple  $\bar{a}$  over  $\text{Const}$   
**Output:** **true** if  $\bar{a} \in q(D)$ , and **false** otherwise

Note that the evaluation problem is presented as a *decision* problem, that is, a problem whose output is either **true** or **false**. Although in practice the goal is to compute the output of  $q$  on  $D$ , in the study of the principles of databases we are mainly interested in understanding the inherent complexity of a query language. This can be achieved by studying the complexity of the decision version of the evaluation problem, which in turn allows us to employ well established tools from complexity theory such as the standard complexity classes that can be found in Appendix B.

The complexity of the problem as stated above is referred to as *combined complexity* of query evaluation. The term combined reflects the fact that both the query  $q$  and the database  $D$  are part of the input.

Very often we shall deal with a different kind of complexity of query evaluation, where the query  $q$  is *fixed*. This is referred to as *data complexity* since we measure the complexity only in terms of the size of the database  $D$ , which in practice, almost invariably, is much bigger than the size of the query  $q$ . More precisely, when we talk about data complexity, we are actually interested in the complexity of the problem  *$q$ -Evaluation* for some query  $q$ :

#### **Problem: $q$ -Evaluation**

**Input:** A database  $D$ , and a tuple  $\bar{a}$  over  $\text{Const}$   
**Output:** **true** if  $\bar{a} \in q(D)$ , and **false** otherwise

Thus, when we talk about the data complexity of  $\mathcal{L}$ -Evaluation, we actually refer to a family of problems, one for each query  $q$  from  $\mathcal{L}$ . Nonetheless, we



shall apply the standard notions of complexity theory, such as membership in a complexity class, or hardness and completeness for a class, to data complexity. We proceed to precisely explain what we mean by that.

**Definition 2.6: Data Complexity**

Let  $\mathcal{L}$  be a query language, and  $\mathcal{C}$  a complexity class.  $\mathcal{L}$ -Evaluation is

- *in  $\mathcal{C}$  in data complexity* if, for every  $q$  from  $\mathcal{L}$ ,  $q$ -Evaluation is in  $\mathcal{C}$ ,
- *$\mathcal{C}$ -hard in data complexity* if there exists a query  $q$  from  $\mathcal{L}$  such that  $q$ -Evaluation is  $\mathcal{C}$ -hard, and
- *$\mathcal{C}$ -complete in data complexity* if  $\mathcal{L}$ -Evaluation is in  $\mathcal{C}$  in data complexity, and  $\mathcal{C}$ -hard in data complexity.

To reiterate, as we shall use these concepts many times in this book:

**Combined Complexity** of query evaluation refers to the complexity of the  $\mathcal{L}$ -Evaluation problem when all of  $q$ ,  $D$ , and  $\bar{a}$  are inputs, and

**Data Complexity** refers to the complexity of  $\mathcal{L}$ -Evaluation when its input consists only of  $D$  and  $\bar{a}$ , whereas  $q$  is fixed. In other words, it refers to the complexity of the family of problems  $\{q\text{-Evaluation} \mid q \text{ is a query from } \mathcal{L}\}$  in the sense of Definition [2.6](#)

*Query Containment and Equivalence*

The basis of static analysis of queries is the *containment* problem. We say that a query  $q$  is *contained* in a query  $q'$ , written as  $q \subseteq q'$ , if  $q(D) \subseteq q'(D)$  for every database  $D$ ; note that since queries return sets of tuples, the notion of subset is applicable to query outputs. This is the most basic task of reasoning about queries; note that containment is one part of equivalence. Indeed,  $q$  is *equivalent* to  $q'$ , denoted  $q \equiv q'$ , if  $q \subseteq q'$  and  $q' \subseteq q$ . The equivalence problem is the most basic one in query optimization, whose goal is to transform a query  $q$  into an equivalent, and more efficient, query  $q'$ .

In relation to containment and equivalence, we consider the following decision problems, again parameterized by a query language  $\mathcal{L}$ .

**Problem:  $\mathcal{L}$ -Containment**

**Input:** Two queries  $q$  and  $q'$  from  $\mathcal{L}$   
**Output:** true if  $q \subseteq q'$ , and false otherwise

**Problem:  $\mathcal{L}$ -Equivalence****Input:** Two queries  $q$  and  $q'$  from  $\mathcal{L}$ **Output:** `true` if  $q \equiv q'$ , and `false` otherwise

Observe that for the previous problems, the input consists of two queries. Typically, queries are much smaller objects than databases. Therefore, for the containment and equivalence problems, we shall in general tolerate higher complexity than for query evaluation; even intractable complexity will often be reasonable, given the small size of the input.

## Computational Complexity Analysis

We will use two different cost models for analyzing the computational complexity of algorithmic problems such as the ones introduced above.

**Turing Machine models** are typically associated with complexity classes such as PTIME, NP, PSPACE, SPACE( $\log n$ ),  $\Pi_2^P$ , etc.

**Random-Access Machine models** are usually used for analyzing the run-time of efficient algorithms. For instance, if we say that  $n$  numbers can be sorted in time  $O(n \log n)$ , then the intended underlying computational model is a random-access machine model.

Throughout the book, we will assume the Turing Machine model for analyzing the complexity of problems in terms of complexity classes, whereas we will assume random-access models when it comes to proving that algorithms have low complexity. The difference between the two will always be clear in the formal statement, where we will always either

- refer to a concrete complexity class, such as PTIME, NLOGSPACE, or DLOGSPACE<sup>2</sup> in which case we assume Turing Machines, or
- say that the problem “is solvable in” time or space  $O(f(n))$  for some function  $f$ , in which case we assume a random access machine model.

### *Size of the Input*

The complexity of algorithms is always analyzed in terms of the size of their input. To this end, we define the size  $\|o\|$  of syntactic objects  $o$  that we will consider for complexity analysis. In particular:

<sup>2</sup> In this book, we usually denote complexity classes in small caps, see Appendix B. Exceptions to this convention are well-known complexity classes for which fonts are irrelevant, such as  $\Pi_2^P$  and #P.

- $\|\emptyset\| = \|()\| = 1$ .
- $\|u\| = 1$  for each  $u \in \text{Const} \cup \text{Var}$ .
- For a nonempty set  $S = \{e_1, \dots, e_n\}$ , we define  $\|S\| = \sum_{i=1}^n \|e_i\|$ .
- For a tuple  $\bar{u} = (u_1, \dots, u_k)$  or a relational atom  $R(\bar{u}) = R(u_1, \dots, u_k)$  with  $k \geq 1$ , we define  $\|\bar{u}\| = \|R(\bar{u})\| = \sum_{i=1}^k \|u_i\|$ .

Hence, if  $D$  is a nonempty database instance of schema  $\{R_1, \dots, R_n\}$ , then

$$\|D\| = \sum_{i=1}^n (|D(R_i)| \cdot (\text{ar}(R_i))) ,$$

assuming that the arities of  $R_1, \dots, R_n$  are nonzero.

The size  $\|q\|$  of a query  $q$  will depend on the formalisms that we will use throughout the book for representing  $q$ , such as first-order logic or Datalog. This means that the size of a query specified using first-order logic will be defined differently than the size of a query specified in Datalog, because the sizes of the associated first-order formula or Datalog program will be different. Throughout the book, we will therefore define the size of queries at the point where we define the class of queries in question.

Turing Machines and random-access machines perceive their inputs differently. Whereas a random-access machine can store a natural number  $n \in \mathbb{N}$  in a single register, a Turing Machine will store  $n$  as a word of  $O(\log n)$  symbols from its finite alphabet. In Appendix [C](#), we discuss how databases and queries are encoded for Turing Machines. For instance, storing a database  $D$  on a Turing Machine costs space  $O(\|D\| \cdot \log \|D\|)$ . Intuitively, the encoding uses  $O(\|D\|)$  many constants, and we need  $O(\log \|D\|)$  space to encode each such constant using the Turing Machine's finite alphabet.

Since it is well known that the random access model and the Turing Machine model are equally efficient as long as polynomial differences do not matter, we sometimes do a random-access-style analysis for easier presentation, even if the underlying computational model is a Turing Machine.

## Further Background Reading

Should the reader find herself/himself in a situation “that she/he does not have the prerequisites for reading the prerequisites” [\[8\]](#), rather than being discouraged she/he is advised to continue with the main material, as it is still very likely to be understood completely or almost completely. Should the latter happen, the prerequisites can be supplemented by information from many standard sources, some of which are listed below.

The book [\[1\]](#) covers the basics of database theory, while many database systems texts cover design, querying, and building real-life databases, for example, [\[6, 14, 16\]](#). The basic mathematical background needed is covered in a

standard undergraduate “discrete mathematics for computer science” course; moreover, a good source for this material is the book [15]. For additional information about computability theory, we provide a primer in Appendix B. Furthermore, we refer the reader to [9, 11, 17]; standard texts on complexity theory are [2, 13, 20]. For the foundations of finite model theory and descriptive complexity, the reader is referred to [7, 10, 12].

**The Relational Model: The Classics**



## First-Order Logic

Database query languages are either *declarative* or *procedural*. In a declarative language, one provides a specification of what a query result should be, typically by means of logical formulae (sometimes presented in a specialized programming syntax). In the case of relational databases, such languages are usually based on *first-order logic*, which often appears in the literature under the name *relational calculus*. In a procedural language, on the other hand, one specifies how the data is manipulated to produce the desired result. The most commonly used one for relational databases is *relational algebra*. We present these languages next, starting with first-order logic.

### Syntax of First-Order Logic

Recall that a schema  $\mathbf{S}$  can be seen as a finite set of relation names, and each relation name of  $\mathbf{S}$  has an arity under  $\mathbf{S}$ . Recall also that we assume a countably infinite set of values  $\mathbf{Const}$  called constants, and a countably infinite set of variables  $\mathbf{Var}$ . Constants will be typically denoted by  $a, b, c, \dots$ , and variables by  $x, y, z, \dots$  (possibly with subscripts and superscripts). Constants and variables are called *terms*. Formulae of first-order logic are inductively defined using terms, conjunction ( $\wedge$ ), disjunction ( $\vee$ ), negation ( $\neg$ ), existential quantification ( $\exists$ ), and universal quantification ( $\forall$ ).

#### Definition 3.1: Syntax of First-Order Logic

We define *formulae of first-order logic* (FO) over a schema  $\mathbf{S}$  as follows:

- If  $a$  is a constant from  $\mathbf{Const}$ , and  $x, y$  are variables from  $\mathbf{Var}$ , then  $x = a$  and  $x = y$  are atomic formulae.
- If  $u_1, \dots, u_k$  are terms (not necessarily distinct), and  $R$  is a  $k$ -ary relation name from  $\mathbf{S}$ , then  $R(u_1, \dots, u_k)$  is an atomic formula.

- If  $\varphi_1$  and  $\varphi_2$  are formulae, then  $(\varphi_1 \wedge \varphi_2)$ ,  $(\varphi_1 \vee \varphi_2)$ , and  $(\neg\varphi_1)$  are formulae.
- If  $\varphi$  is a formula and  $x \in \text{Var}$ , then  $(\exists x \varphi)$  and  $(\forall x \varphi)$  are formulae.

The *size*  $\|\varphi\|$  of  $\varphi$  is defined to be the total number of occurrences constants, variables, and symbols from  $\{\wedge, \vee, \neg, =, \exists, \forall\}$  occurring in  $\varphi$ . For example, the size of  $(x = a \vee x = b)$  is seven.

Formulae of the form  $x = a$  and  $x = y$  are called *equational atoms*. Furthermore, as already mentioned in Chapter 2, formulae of the form  $R(\bar{u})$  are called *relational atoms*. Note that we allow repetition of variables in relational atoms, for example, we may write  $R(x, x, y)$ . We shall use the standard shorthand  $(\varphi \rightarrow \psi)$  for  $((\neg\varphi) \vee \psi)$  and  $(\varphi \leftrightarrow \psi)$  for  $((\varphi \rightarrow \psi) \wedge (\psi \rightarrow \varphi))$ . To reduce notational clutter, we will often omit the outermost brackets of formulae.

A crucial notion is that of free variables of a formula, which are essentially the variables in a formula that are not quantified. Given an FO formula  $\varphi$ , the set of *free* variables of  $\varphi$ , denoted  $\text{FV}(\varphi)$ , is inductively defined as follows:

- $\text{FV}(x = y) = \{x, y\}$ .
- $\text{FV}(x = a) = \{x\}$ .
- $\text{FV}(R(u_1, \dots, u_k)) = \{u_1, \dots, u_k\} \cap \text{Var}$ .
- $\text{FV}(\varphi_1 \vee \varphi_2) = \text{FV}(\varphi_1 \wedge \varphi_2) = \text{FV}(\varphi_1) \cup \text{FV}(\varphi_2)$ .
- $\text{FV}(\neg\varphi) = \text{FV}(\varphi)$ .
- $\text{FV}(\exists x \varphi) = \text{FV}(\forall x \varphi) = \text{FV}(\varphi) - \{x\}$ .

If  $x \in \text{FV}(\varphi)$ , we call it a *free variable* (of  $\varphi$ ); otherwise,  $x$  is called *bound*. An FO formula  $\varphi$  without free variables is called a *sentence*.

### Example 3.2: First-Order Formulae

Consider the following (named) database schema:

```

Person [ pid, pname, cid ]
Profession [ pid, prname ]
City [ cid, cname, country ]

```

The Person relation stores internal person IDs (**pid**), names (**pname**), and the ID of their city of birth (**cid**). The Profession relation contains the professions of persons by storing their person ID (**pid**) and profession name (**prname**). Finally, City contains a bit of geographic information by storing IDs (**cid**) and names (**cname**) of cities, together with the country they are located in (**country**). In what follows, we give some examples of FO formulae over this schema. Consider first the FO formula:



$$\exists y \exists z \exists u_1 \exists u_2 (\text{Person}(x, y, z) \wedge \text{Profession}(x, u_1) \wedge \text{Profession}(x, u_2) \wedge \neg(u_1 = u_2)). \quad (3.1)$$

This formula has one free variable, that is,  $x$ . Consider now the formula

$$\exists z (\text{Person}(x, y, z) \wedge \forall r \forall s (\neg \text{City}(z, r, s))). \quad (3.2)$$

The free variables of this formula are  $x, y$ . Finally, consider the formula

$$\exists x \exists z (\text{Person}(x, y, z) \wedge (\text{Profession}(x, \text{'author'}) \vee \text{Profession}(x, \text{'actor'}))). \quad (3.3)$$

This formula has one free variable, that is,  $y$ .

## Semantics of First-Order Logic

Given a database  $D$  of a schema  $\mathbf{S}$ , we inductively define the notion of satisfaction of a formula  $\varphi$  over  $\mathbf{S}$  in  $D$  with respect to an *assignment*  $\eta$  for  $\varphi$  over  $D$ . Such an assignment is a function from  $\text{FV}(\varphi)$  to  $\text{Dom}(D) \cup \text{Dom}(\varphi) \subseteq \text{Const}$ , where  $\text{Dom}(\varphi)$  is the set of constants mentioned in  $\varphi$ . For example, for the formula  $R(x, y, a)$ ,  $\eta$  is the function  $\{x, y\} \rightarrow \text{Dom}(D) \cup \{a\}$ . In the following definition (and also later in the book), we write  $\eta[x/u]$ , for a variable  $x$  and term  $u$ , for the assignment that modifies  $\eta$  by setting  $\eta(x) = u$ . Furthermore, to avoid heavy notation, we extend  $\eta$  to be the identity on  $\text{Const}$ .

### Definition 3.3: Semantics of First-Order Logic

Given a database  $D$  of a schema  $\mathbf{S}$ , a formula  $\varphi$  over  $\mathbf{S}$ , and an assignment  $\eta$  for  $\varphi$  over  $D$ , we inductively define when  $\varphi$  is *satisfied* in  $D$  under  $\eta$ , written  $(D, \eta) \models \varphi$ , as follows:

- If  $\varphi$  is  $x = y$ , then  $(D, \eta) \models \varphi$  if  $\eta(x) = \eta(y)$ .
- If  $\varphi$  is  $x = a$ , then  $(D, \eta) \models \varphi$  if  $\eta(x) = a$ .
- If  $\varphi$  is  $R(u_1, \dots, u_k)$ , then  $(D, \eta) \models \varphi$  if  $R(\eta(u_1), \dots, \eta(u_k)) \in D$ .
- If  $\varphi = \varphi_1 \wedge \varphi_2$ , then  $(D, \eta) \models \varphi$  if  $(D, \eta) \models \varphi_1$  and  $(D, \eta) \models \varphi_2$ .
- If  $\varphi = \varphi_1 \vee \varphi_2$ , then  $(D, \eta) \models \varphi$  if  $(D, \eta) \models \varphi_1$  or  $(D, \eta) \models \varphi_2$ .
- If  $\varphi = \neg \psi$ , then  $(D, \eta) \models \varphi$  if  $(D, \eta) \models \psi$  does not hold.
- If  $\varphi = \exists x \psi$ , then  $(D, \eta) \models \varphi$  if  $(D, \eta[x/a]) \models \psi$  for *some* constant  $a \in \text{Dom}(D) \cup \text{Dom}(\varphi)$ .
- If  $\varphi = \forall x \psi$ , then  $(D, \eta) \models \varphi$  if  $(D, \eta[x/a]) \models \psi$  for *each* constant  $a \in \text{Dom}(D) \cup \text{Dom}(\varphi)$ .

An assignment  $\eta$  for a sentence  $\varphi$  has an empty domain (since the domain of  $\eta$  is  $\text{FV}(\varphi)$ ), and thus it is unique. For this unique  $\eta$ , it is either the case that  $(D, \eta) \models \varphi$  or not. If the former is true, then we simply write  $D \models \varphi$  and say that  $D$  *satisfies*  $\varphi$ .

#### Example 3.4: Semantics of First-Order Formulae

We provide an intuitive description of the semantic meaning of the formulae given in Example 3.2:

- Formula (3.1) is satisfied by all  $x$  such that  $x$  is the ID of a person with two different professions.
- Formula (3.2) is satisfied by all  $x, y$  such that  $x$  and  $y$  are the ID and name of persons for which their city of birth is not in the database.
- Formula (3.3) is satisfied by all  $y$  such that  $y$  is the name of a person who is an author or an actor.

It is crucial to say that the semantics of FO are defined in a way that is well-suited for database applications, but slightly departs from the logic literature. In particular, the range of quantifiers is the set of constants  $\text{Dom}(D) \cup \text{Dom}(\varphi)$  (see the last two items of Definition 3.3), whereas in the standard definition is the set of values  $\text{Const}$ . This is why  $\eta$  associates elements of  $\text{Dom}(D) \cup \text{Dom}(\varphi)$  to variables, while in the standard definition one would allow  $\eta$  to associate arbitrary elements of  $\text{Const}$  to variables. The set  $\text{Dom}(D) \cup \text{Dom}(\varphi)$  is called the *active domain of  $D$  and  $\varphi$* . Therefore, Definition 3.3 actually defines the so-called *active domain semantics*, which is standard in the database literature. The importance of the active domain semantics is revealed below where we use FO to define database queries.

## Notational Conventions

We introduce some notational conventions concerning FO formulae that would significantly improve readability:

- Since conjunction is associative, we will omit brackets in long conjunctions and write, for example,  $x_1 \wedge x_2 \wedge x_3 \wedge x_4$  instead of  $((x_1 \wedge x_2) \wedge x_3) \wedge x_4$ . We follow the same convention for disjunction. We also omit brackets within sequences of quantifiers.
- We often write  $\exists \bar{x} \varphi$  for  $\exists x_1 \exists x_2 \dots \exists x_m \varphi$ , where  $\bar{x} = (x_1, \dots, x_m)$ , and likewise for universal quantifiers  $\forall \bar{x}$ .
- We assume that  $\neg$  binds the strongest, followed by  $\wedge$ , then  $\vee$ , and finally quantifiers. For example, by  $\exists x \neg R(x) \wedge S(x)$  we mean the formula  $\exists x ((\neg R(x)) \wedge S(x))$ . We will, however, add brackets to formulae when

we feel that it improves their readability. Notice that this precedence of operators also influences the range of variables; e.g., by  $\forall x R(x) \wedge S(x)$  we mean the formula  $\forall x (R(x) \wedge S(x))$ , as opposed to  $(\forall x R(x)) \wedge S(x)$ .

- Finally, we write  $x \neq y$  instead of  $\neg(x = y)$ , and likewise for  $(x = a)$ .

## Equivalences

In the way FO is defined in Definition 3.1 some constructors are redundant. For instance, we know by De Morgan's laws that  $\neg(\varphi \vee \psi)$  is equivalent to  $\neg\varphi \wedge \neg\psi$ , and  $\neg(\varphi \wedge \psi)$  is equivalent to  $\neg\varphi \vee \neg\psi$ . Furthermore, the formula  $\neg\forall x \varphi$  is equivalent to  $\exists x \neg\varphi$  and  $\neg\exists x \varphi$  is equivalent to  $\forall x \neg\varphi$ . These equivalences mean that the full set of Boolean connectives and quantifiers is not necessary to define all of FO. For example, one can just use  $\vee, \neg$ , and  $\exists$ , or  $\wedge, \neg$ , and  $\exists$ , and this will capture the full expressive power of FO. This is useful for proofs that proceed by induction on the structure of FO formulae.

For some proofs in Part I of the book it will be convenient to assume that constants do not appear in relational atoms. We can always rewrite FO formulae to such a form via equalities, at the expense of a linear blow-up. For instance, we can write  $R(x, a, b)$  as  $\exists x_a \exists x_b R(x, x_a, x_b) \wedge (x_a = a) \wedge (x_b = b)$ .

## First-Order Queries

Recall that a  $k$ -ary query  $q$  produces a finite set of  $k$ -ary tuples  $q(D) \subseteq \text{Const}^k$ , for every database  $D$ . FO formulae can be used to define database queries. In order to do this, we specify together with the formula  $\varphi$  a tuple  $\bar{x}$  of variables that indicates how the output of the query is formed. As a simple example, consider an atomic formula  $\varphi = R(x, y)$  and the tuple  $(x, y)$ . Then the query  $\varphi(x, y)$  would return the entire relation  $R$  from the database. Notice that the query is actually  $R(x, y)(x, y)$ , where the first occurrence of  $(x, y)$  is part of the relational atom  $R(x, y)$ , and the second occurrence specifies how the output of the query is formed. To consider a few other examples, if  $\varphi = R(x, y)$ , then the query  $\varphi(x, x, y)$  returns all tuples  $(a, a, b)$  such that  $(a, b)$  is in the relation  $R$ . Finally, if  $\varphi = R(x, x)$ , then the query  $\varphi(x)$  returns all tuples  $(a)$  such that  $(a, a)$  is in the relation  $R$ . The definition of FO queries follows.

### Definition 3.5: First-Order Queries

A *first-order query* over a schema  $\mathbf{S}$  is an expression of the form  $\varphi(\bar{x})$ , where  $\varphi$  is an FO formula over  $\mathbf{S}$ , and  $\bar{x}$  is a tuple of free variables of  $\varphi$  such that each free variable of  $\varphi$  occurs in  $\bar{x}$  at least once.

For a first-order query  $q = \varphi(\bar{x})$ , we define its *size*  $\|q\|$  as  $\|\varphi\| + \|\bar{x}\|$ .

Let  $\varphi(\bar{x})$  be an FO query over  $\mathbf{S}$ . Given a database  $D$  of  $\mathbf{S}$ , and a tuple  $\bar{a}$  of elements from  $\mathbf{Const}$ , we say that  $D$  *satisfies the query*  $\varphi(\bar{x})$  *using the values*  $\bar{a}$ , denoted by  $D \models \varphi(\bar{a})$ , if there exists an assignment  $\eta$  for  $\varphi$  over  $D$  such that  $\eta(\bar{x}) = \bar{a}$  and  $(D, \eta) \models \varphi$ . Having this notion in place, we can now define what is the output of an FO query on a database.

**Definition 3.6: Evaluation of First-Order Queries**

Given a database  $D$  of a schema  $\mathbf{S}$ , and an FO query  $q = \varphi(x_1, \dots, x_k)$  over  $\mathbf{S}$ , where  $k \geq 0$ , the *output* of  $q$  on  $D$  is defined as the set of tuples

$$q(D) = \{\bar{a} \in \mathbf{Const}^k \mid D \models \varphi(\bar{a})\}.$$

It is clear that  $q(D) \in \mathcal{P}(\mathbf{Const}^k)$ . However, to be able to say that  $q$  defines a  $k$ -ary query over  $\mathbf{S}$  in the sense of Definition 2.5, we need to ensure that  $q(D) \in \mathcal{P}_{\text{fin}}(\mathbf{Const}^k)$ , i.e., the output of  $q$  on  $D$  is finite. This is guaranteed by the following result, which is an immediate consequence of the active domain semantics of FO (see Definition 3.3).

**Proposition 3.7**

Given a database  $D$  of a schema  $\mathbf{S}$ , and an FO query  $q = \varphi(x_1, \dots, x_k)$  over  $\mathbf{S}$ , where  $k \geq 0$ , it holds that

$$q(D) = \{\bar{a} \in (\text{Dom}(D) \cup \text{Dom}(\varphi))^k \mid D \models \varphi(\bar{a})\}.$$

Since, by definition, the set of values  $\text{Dom}(D) \cup \text{Dom}(\varphi)$  is finite, Proposition 3.7 implies that  $q(D) \in \mathcal{P}_{\text{fin}}(\mathbf{Const}^k)$ , and thus,  $q$  defines a  $k$ -ary query over  $\mathbf{S}$  in the sense of Definition 2.5.

Before we proceed further, let us stress that if we adopt the standard semantics of FO from logic textbooks, which uses assignments  $\eta$  that associate arbitrary elements of  $\mathbf{Const}$  to variables, then there is no guarantee that  $q(D)$  is finite. Consider, for example, the query  $q = \varphi(x)$  with  $\varphi = \neg R(x)$ , and the database  $D = \{R(a), P(b)\}$ . Under the standard FO semantics, the output of  $q$  on  $D$  would be the set  $\{c \in \mathbf{Const} - \{a\}\}$ , and thus infinite. On the other hand, under the active domain semantics we have that  $q(D) = \{b\}$ .

**Example 3.8: Evaluation of First Order Queries**

A database  $D$  of the schema in Example 3.2 is depicted in Figure 3.1. We proceed to evaluate the FO queries obtained from the FO formulae given in Example 3.2 on  $D$ :

- Let  $q_1$  be the query  $\varphi_1(x)$ , where  $\varphi_1$  is the formula (3.1). Then

$$q_1(D) = \{('1'), ('3'), ('4')\}.$$

Person			Profession	
pid	pname	cid	pid	prname
1	Aretha	MPH	1	singer
2	Billie	BLT	1	songwriter
3	Bob	DLT	1	actor
4	Freddie	ST	2	singer
			3	singer
			3	songwriter
			3	author
			4	singer
			4	songwriter

City		
cid	cname	country
MPH	Memphis	United States
DLT	Duluth	United States
ST	Stone Town	Tanzania

Fig. 3.1: A database of the schema in Example 3.2

- Let  $q_2$  be the query  $\varphi_2(x, y)$ , where  $\varphi_2$  is the formula (3.2). Then

$$q_2(D) = \{('2', 'Billie')\}.$$

- Let  $q_3$  be the query  $\varphi_3(y)$ , where  $\varphi_3$  is the formula (3.3). Then

$$q_3(D) = \{('Aretha'), ('Bob')\}.$$

## Boolean First-Order Queries

FO sentences, that is, FO formulae without free variables, are used to define Boolean queries, i.e., queries that return **true** or **false**, and hence the name *Boolean FO queries*. By definition, the output of a query  $q$  on a database  $D$  corresponds to a set of tuples, and thus, Boolean FO queries will be no exception to this. We consider such queries to be of the form  $q = \varphi()$ , where  $\varphi$  is an FO sentence, and  $()$  denotes the empty tuple. There are two cases:

- either  $q(D)$  consists of the empty tuple, that is,  $q(D) = \{()\}$ , which happens precisely when  $D \models \varphi$ , or
- $q(D)$  is the empty set, which happens precisely when  $D \models \neg\varphi$ .

By convention, we write  $q(D) = \mathbf{true}$  if  $D \models \varphi$ , and  $q(D) = \mathbf{false}$  otherwise.



## Relational Algebra

Queries expressed in FO are *declarative* and tell us *what* the output of a query should be. In this chapter, we introduce *relational algebra*, abbreviated RA, which contrasts itself with FO because it is *procedural*, i.e., it specifies *how* the output of queries can be obtained via a sequence of operations on the data. Relational algebra is of significant practical importance in databases, since database systems typically use relational-algebra-like representations of queries to do query optimization, that is, to discover methods in which a given query can be evaluated efficiently.

We present relational algebra in its most elementary form, in both the unnamed and the named perspective. The following table gives a quick overview of the operators in the unnamed and named relational algebra.

<i>Operator Name</i>	<i>(Unnamed) RA Symbol</i>	<i>Named RA Symbol</i>
selection	$\sigma_\theta$	$\sigma_\theta$
projection	$\pi_\alpha$	$\pi_\alpha$
Cartesian product	$\times$	
rename		$\rho$
union	$\cup$	$\cup$
difference	$-$	$-$
join	$\bowtie_\theta$	$\bowtie$

We explain these operators and their semantics next, in the definitions of the *unnamed* and *named RA*. Since we will usually be working with the unnamed perspective in this book, we will often abbreviate “unnamed RA” as “RA”.

### Syntax of the Unnamed Relational Algebra

Under the unnamed perspective, RA consists of five primitive operations: selection, projection, Cartesian product, union, and difference. Before giving the

formal definitions of those operations, we first introduce conditions over sets of integers, which are needed for defining the selection operation. A *condition*  $\theta$  over  $\{1, \dots, k\}$ , for some  $k \geq 0$ , is a Boolean combination of statements of the form  $i \doteq j$ ,  $i \doteq a$ ,  $i \not\doteq j$ , and  $i \not\doteq a$ , where for  $a \in \mathbf{Const}$  and  $i, j \in [k]$ . Intuitively, a condition  $i \doteq j$  is used to indicate that in a tuple the values of the  $i$ -th attribute and the  $j$ -th attribute must be the same, while  $i \not\doteq j$  is used to indicate that these values must be different. Moreover, a condition  $i \doteq a$  is used to indicate that in a tuple the value of the  $i$ -th attribute must be the constant  $a$ , while  $i \not\doteq a$  is used to indicate that this value must be different than  $a$ . Let us clarify that we use the symbols  $\doteq$  and  $\not\doteq$ , instead of  $=$  and  $\neq$ , to avoid writing statements such as “ $1 = 2$ ”, which are likely to confuse the reader. Notice that by using De Morgan’s laws to propagate negation, we can define conditions as *positive* Boolean combinations of statements  $i \doteq j$  and  $i \not\doteq j$ , i.e., Boolean combinations using only conjunction  $\wedge$  and disjunction  $\vee$ . For example,  $\neg((1 \doteq 2) \vee (2 \not\doteq 3))$  is equivalent to  $(1 \not\doteq 2) \wedge (2 \doteq 3)$ .

#### Definition 4.1: Syntax of Unnamed Relational Algebra

We inductively define *RA expressions* over a schema  $\mathbf{S}$ , and their associated arities, as follows:

**Base Expressions.** If  $R$  is a  $k$ -ary relation name from  $\mathbf{S}$ , then  $R$  is an atomic RA expression over  $\mathbf{S}$  of arity  $k$ . If  $a \in \mathbf{Const}$ , then  $\{a\}$  is an RA expression over  $\mathbf{S}$  of arity 1.

**Selection.** If  $e$  is an RA expression over  $\mathbf{S}$  of arity  $k \geq 0$  and  $\theta$  is a condition over  $[k]$ , then  $\sigma_\theta(e)$  is an RA expression over  $\mathbf{S}$  of arity  $k$ .

**Projection.** If  $e$  is an RA expression over  $\mathbf{S}$  of arity  $k \geq 0$  and  $\alpha = (i_1, \dots, i_m)$ , for  $m \geq 0$ , is a list of numbers from  $[k]$ , then  $\pi_\alpha(e)$  is an RA expression over  $\mathbf{S}$  of arity  $m$ .

**Cartesian Product.** If  $e_1, e_2$  are RA expressions over  $\mathbf{S}$  of arity  $k \geq 0$  and  $m \geq 0$ , respectively, then their Cartesian product  $(e_1 \times e_2)$  is an RA expression over  $\mathbf{S}$  of arity  $k + m$ .

**Union.** If  $e_1, e_2$  are RA expressions over  $\mathbf{S}$  of the same arity  $k \geq 0$ , then their union  $(e_1 \cup e_2)$  is an RA expression over  $\mathbf{S}$  of arity  $k$ .

**Difference.** If  $e_1, e_2$  are RA expressions over  $\mathbf{S}$  of the same arity  $k \geq 0$ , then their difference  $(e_1 - e_2)$  is an RA expression over  $\mathbf{S}$  of arity  $k$ .

The *size*  $\|e\|$  of an RA expression  $e$  is the total number of occurrences of relation names, constants, natural numbers, and symbols from  $\{\sigma, \pi, \times, \cup, -, \wedge, \vee, \neg, \doteq, \not\doteq\}$  in  $e$ . For instance, the size of  $\sigma_{1 \doteq 2}(\pi_{(1,2)}(R))$  is eight.

Notice that in the definition of the projection operation, we allow  $m$  to be 0, in which case the list of integers  $\alpha = (i_1, \dots, i_m)$  is the empty list  $()$ . This is useful for expressing Boolean queries.



## Semantics of Unnamed Relational Algebra

We proceed to define the semantics of RA expressions. We first need to define the operation of projection over tuples. For a tuple  $\bar{a} = (a_1, \dots, a_k) \in \text{Const}^k$ , and a list  $\alpha = (i_1, \dots, i_m)$  of numbers from  $[k]$ , the projection  $\pi_\alpha(\bar{a})$  is defined as the tuple  $(a_{i_1}, a_{i_2}, \dots, a_{i_m})$ .<sup>1</sup> Here are some simple examples:

$$\pi_{(1,3)}(a, b, c, d) = (a, c) \quad \pi_{(1,3,3)}(a, b, c, d) = (a, c, c) \quad \pi_{()}(a, b, c, d) = ()$$

We also need the notion of satisfaction of conditions over tuples. We inductively define when a tuple  $\bar{a}$  *satisfies the condition*  $\theta$ , denoted  $\bar{a} \models \theta$ :

$$\begin{aligned} \bar{a} \models i \doteq j & \text{ if } a_i = a_j & \bar{a} \models i \doteq a & \text{ if } a_i = a \\ \bar{a} \models i \not\doteq j & \text{ if } a_i \neq a_j & \bar{a} \models i \not\doteq a & \text{ if } a_i \neq a \\ \bar{a} \models \theta \wedge \theta' & \text{ if } \bar{a} \models \theta \text{ and } \bar{a} \models \theta' & \bar{a} \models \theta \vee \theta' & \text{ if } \bar{a} \models \theta \text{ or } \bar{a} \models \theta' \\ \bar{a} \models \neg\theta & \text{ if } \bar{a} \models \theta \text{ does not hold} \end{aligned}$$

We are now ready to define the semantics of RA expressions.

### Definition 4.2: Semantics of Unnamed RA Expressions

Let  $D$  be a database of a schema  $\mathbf{S}$ , and  $e$  an RA expression over  $\mathbf{S}$ . We inductively define the *output*  $e(D)$  of  $e$  on  $D$  as follows:

- If  $e = R$ , where  $R$  is a relation name from  $\mathbf{S}$ , then  $e(D) = R^D$ .
- If  $e = \{a\}$ , for  $a \in \text{Const}$ , then  $e(D) = \{a\}$ .
- If  $e = \sigma_\theta(e_1)$ , where  $e_1$  is an RA expression of arity  $k \geq 0$  and  $\theta$  is a condition over  $[k]$ , then  $e(D) = \{\bar{a} \mid \bar{a} \in e_1(D) \text{ and } \bar{a} \models \theta\}$ .
- If  $e = \pi_\alpha(e_1)$ , where  $e_1$  is an RA expression of arity  $k \geq 0$  and  $\alpha = (i_1, \dots, i_m)$ , for  $m \geq 0$ , is a list of numbers from  $[k]$ , then  $e(D)$  is the  $m$ -ary relation  $\{\pi_\alpha(\bar{a}) \mid \bar{a} \in e_1(D)\}$ .
- If  $e = (e_1 \times e_2)$ , where  $e_1$  and  $e_2$  are RA expressions of arity  $k \geq 0$  and  $\ell \geq 0$ , respectively, then  $e(D) = e_1(D) \times e_2(D)$ .
- If  $e = (e_1 \cup e_2)$ , where  $e_1$  and  $e_2$  are RA expressions of the same arity  $k \geq 0$ , then  $e(D) = e_1(D) \cup e_2(D)$ .
- If  $e = (e_1 - e_2)$ , where  $e_1$  and  $e_2$  are RA expressions of the same arity  $k \geq 0$ , then  $e(D) = e_1(D) - e_2(D)$ .

We sometimes use derived operations, one of them of special importance:

<sup>1</sup> The projection  $\pi_\alpha(\bar{u})$ , where  $\bar{u}$  is tuple from  $(\text{Const} \cup \text{Var})^k$ , is defined in the same way. For example,  $\pi_{(1,3)}(a, x, y, d) = (a, y)$  and  $\pi_{(1,3,3)}(a, x, y, d) = (a, y, y)$ . We are going to apply the projection operator over tuples of constants and variables in subsequent chapters such as Chapters [10](#) and [11](#).

**Join.** Given a  $k$ -ary RA expression  $e_1$ , an  $m$ -ary RA expression  $e_2$ , and a condition  $\theta$  over  $\{1, \dots, k+m\}$ , the  $\theta$ -join of  $e_1$  and  $e_2$  is denoted  $e_1 \bowtie_{\theta} e_2$ . Its output on a database  $D$  is defined as

$$(e_1 \bowtie_{\theta} e_2)(D) = \sigma_{\theta}(e_1(D) \times e_2(D)).$$

We note that RA expressions readily define queries on databases. Indeed, if  $e$  is a RA expression, then the *output* of  $e$  on a database  $D$  is  $e(D)$ . In the remainder of the book, we will therefore sometimes also refer to  $e$  as a *query*.

#### Example 4.3: Unnamed RA Queries

Consider again the (named) database schema:

```
Person [ pid, pname, cid ]
Profession [ pid, pname ]
City [ cid, cname, country ]
```

The RA expression

$$\pi_{(1)}(\sigma_{5 \neq 7}((\text{Person} \bowtie_{1 \doteq 4} \text{Profession}) \bowtie_{1 \doteq 6} \text{Profession}))$$

returns the IDs of persons with at least two professions. The expression

$$\pi_{(1,2)}(\text{Person}) - \pi_{(1,2)}(\text{Person} \bowtie_{3 \doteq 4} \text{City})$$

returns the ID and name of persons whose city of birth does not appear in the database. Finally, the expression

$$\pi_{(2)}(\sigma_{(5 \doteq \text{'author'}) \vee (5 \doteq \text{'actor'})}(\text{Person} \bowtie_{1 \doteq 4} \text{Profession}))$$

returns the names of persons that are author or actors.

## Syntax of the Named Relational Algebra

Under the named perspective, the presentation changes a bit. Before giving the formal definition, let us first note that the notion of condition, needed for defining the selection operation, is now over a set of attributes, and not a set of integers as in the case of unnamed RA. More precisely, a *condition*  $\theta$  over a set of attributes  $U \subseteq \text{Att}$  is a Boolean combination of statements of the form  $A \doteq B$ ,  $A \doteq a$ ,  $A \neq B$ , and  $A \neq a$ , where  $a \in \text{Const}$  and  $A, B \in U$ .

#### Definition 4.4: Syntax of Named Relational Algebra

We inductively define *named RA expressions* over a schema  $\mathbf{S}$ , and their associated sorts, as follows:

**Base Expressions.** If  $R \in \mathbf{S}$ , then  $R$  is an atomic named RA expression over  $\mathbf{S}$  of sort  $\mathbf{S}(R)$ . If  $a \in \mathbf{Const}$  and  $A \in \mathbf{Att}$ , then  $\{(A: a)\}$  is a named RA expression of sort  $\{A\}$ .

**Selection.** If  $e$  is a named RA expression of sort  $U$  and  $\theta$  is a *condition over  $U$* , then  $\sigma_\theta(e)$  is a named RA expression of sort  $U$ .

**Projection.** If  $e$  is a named RA expression of sort  $U$  and  $\alpha \subseteq U$ , then  $\pi_\alpha(e)$  is a named RA expression of sort  $\alpha$ .

**Join.** If  $e_1, e_2$  are named RA expressions of sort  $U_1$  and  $U_2$ , respectively, then their join  $(e_1 \bowtie e_2)$  is a named RA expression of sort  $U_1 \cup U_2$ .

**Rename.** If  $e$  is a named RA expression of sort  $U$ , then  $\rho_{A \rightarrow B}(e)$ , where  $A \in U$  and  $B \in \mathbf{Att} - U$ , is a named RA expression of sort  $(U - \{A\}) \cup \{B\}$ .

**Union.** If  $e_1, e_2$  are named RA expressions of the same sort  $U$ , then their union  $(e_1 \cup e_2)$  is a named RA expression of sort  $U$ .

**Difference.** If  $e_1, e_2$  are named RA expressions of the same sort  $U$ , then their difference  $(e_1 - e_2)$  is a named RA expression of sort  $U$ .

The *size*  $\|e\|$  of a named RA expression  $e$  is the total number of occurrences of relation names, constants, attribute names, and symbols from  $\{\sigma, \pi, \bowtie, \rho, \cup, -, \wedge, \vee, \neg, \dot{=}, \neq\}$  in  $e$ . For instance, the size of  $\sigma_{A \dot{=} B}(\pi_{(A,B)}(R))$  is eight.

Notice in the definition of the projection operation the contrast with the unnamed perspective, where  $\alpha$  is a list of numbers with repetitions.

## Semantics of the Named Relational Algebra

The semantics of named RA expressions is defined similarly to the unnamed case, with the main difference that  $e(D)$  is now a named relation instance. Therefore, we only discuss rename and join, and leave the others as exercises.

**Rename.** If  $e = \rho_{A \rightarrow B}(e_1)$ , where  $e_1$  is a named RA expression of sort  $U$ ,  $A \in U$ , and  $B \in \mathbf{Att} - U$ , then  $e(D)$  is the relation

$$\{t \mid t(B) = t_1(A) \text{ and } t(C) = t_1(C) \text{ for } t_1 \in e_1(D) \text{ and } C \in U - \{A\}\}.$$

Note that renaming does not change the data at all, it only changes names of attributes. Nonetheless, this operation is necessary under the named perspective. For instance, consider two relations,  $R$  and  $S$ , the former with a single attribute  $A$  and the latter with a single attribute  $B$ . Suppose we want to find their union in relational algebra. The problem is that the union is only defined if the sorts of  $R$  and  $S$  are the same, which is not the case. To take their union, we can therefore rename the attribute of  $S$  to be  $A$ , and complete the task by writing the expression  $(R \cup \rho_{B \rightarrow A}(S))$ .

**Join.** The other new primitive operator in the named perspective is *join* (also known in the literature as *natural join*). It is simply a join of two relations on the condition that their common attributes are the same. Formally, if  $e = e_1 \bowtie e_2$ , where  $e_1$  and  $e_2$  are named RA expressions of sorts  $U_1$  and  $U_2$ , then  $e(D)$  is the set of tuples  $t$  such that

$$t(A) = \begin{cases} t_1(A) & \text{if } A \in U_1, \\ t_2(A) & \text{if } A \in U_2 - U_1, \end{cases}$$

where  $t_1 \in e_1(D)$ ,  $t_2 \in e_2(D)$ , and  $t_1(A) = t_2(A)$  for all  $A \in U_1 \cap U_2$ . To give an example, consider the relations  $R[A, B]$  and  $S[B, C]$ . Their join  $R \bowtie S$  has attributes  $A, B, C$ , and consists of triples  $(a, b, c)$  such that  $R(a, b)$  and  $S(b, c)$  are both facts in the database. Notice that, if  $R$  and  $S$  have no common attributes, their join is their Cartesian product. For this reason, we do not have the operator  $\times$  in the named RA.

Similarly to the unnamed perspective, we can interpret named RA expressions  $e$  as queries over databases  $D$ . However, since queries return tuples in  $\text{Const}^k$  according to Definition 2.5, and since  $e(D)$  is a named relation instance, we still need to explain how we go from  $e(D)$  to a finite set of tuples over  $\text{Const}$ . To this end, we will assume that the order  $\leq$  that we introduced in Chapter 2 for translating between databases from the named to the unnamed perspective, is also an order on  $\text{Att}$ , i.e., we assume that it is an order on the set  $\text{Att} \cup (\text{Rel} \times \text{Att})$ .<sup>2</sup> We can now associate to  $e$  a query  $q_e$  by defining that, on database  $D$ , the *output of  $q_e$  on  $D$*  is the set

$$q_e(D) = \{(a_1, \dots, a_k) \mid (A_1: a_1, \dots, A_k: a_k) \in e(D) \text{ such that } A_1 \leq A_2 \leq \dots \leq A_k\}.$$

In the remainder of the book, we will usually not formally distinguish between the RA expression  $e$  and the query  $q_e$ . In particular, if we talk about the *query*  $e$ , then we mean the query  $q_e$  that we just defined.

#### Example 4.5: Named RA Queries

We provide named RA versions for the expressions given in Example 4.3. The expression

$$\pi_{\{\text{pid}\}}(\sigma_{\text{prname} \neq \text{prname2}}((\text{Person} \bowtie \text{Profession}) \bowtie \rho_{\text{prname} \rightarrow \text{prname2}}(\text{Profession})))$$

returns the IDs of persons with at least two professions. The expression

<sup>2</sup> We can assume that  $A \leq (R, B)$  for all  $A, B \in \text{Att}$  and  $R \in \text{Rel}$ , although this is inconsequential.

$$\pi_{\{\text{pid}, \text{pname}\}}(\text{Person}) - \pi_{\{\text{pid}, \text{pname}\}}(\text{Person} \bowtie \text{City})$$

returns the ID and name of persons whose city of birth does not appear in the database. Finally, the expression

$$\pi_{\{\text{pname}\}}(\sigma_{(\text{pname} \doteq \text{'author'}) \vee (\text{pname} \doteq \text{'actor'})}(\text{Person} \bowtie \text{Profession}))$$

returns the names of persons that are authors or actors.

## Expressiveness of Named and Unnamed RA

We often use named RA in examples since it is closer to how we think about real-life databases. On the other hand, many results are easier to state and prove in unnamed RA. This comes at no cost since, as we discuss below, every named RA query can be expressed in unnamed RA, and vice versa.

Let  $f$  be the function that converts a database  $D$  from the named to the unnamed perspective, as presented in Chapter 2. Recall that this converts each tuple  $t = (A_1: a_1, \dots, A_k: a_k)$  in  $D(R)$ , for a relation name  $R$  of sort  $\{A_1, \dots, A_k\}$  (with  $(R, A_1) \leq \dots \leq (R, A_k)$ ), into a tuple  $t' = (a_1, \dots, a_k)$  in  $f(D)(R)$ . Let  $q_n$  be a named RA query,  $q_u$  an unnamed RA query, and  $\mathbf{S}$  a named database schema. We say that  $q_n$  is *equivalent to  $q_u$  under  $\mathbf{S}$*  if, for every database  $D$  of  $\mathbf{S}$ , we have that  $q_n(D) = q_u(f(D))$ .

Note that two queries can be equivalent under one schema but not equivalent under another one. This is unavoidable since the order inside an unnamed tuple depends on the names of the attributes (the order is defined by  $\leq$ ). For instance, if  $R$  is a binary relation name, then  $\pi_{(1)}(R)$  is equivalent to  $\pi_{\{B\}}(R)$  over  $\mathbf{S}_1 = \{R(B, C)\}$  but not over  $\mathbf{S}_2 = \{R(A, B)\}$ .

The following theorem establishes that each named RA query can be translated into an equivalent unnamed RA query. We leave the statement of the reverse direction and its proof as an exercise.

### Theorem 4.6

Consider a named database schema  $\mathbf{S}$ , and a named RA query  $q_n$ . There exists an unnamed RA query  $q_u$  that is equivalent to  $q_n$  under  $\mathbf{S}$ .

*Proof.* We prove this by induction on the structure of  $q_n$ . Assume that  $q_n$  has sort  $\{A_1, \dots, A_k\}$  with  $A_1 \leq \dots \leq A_k$ , where  $\leq$  is the ordering we used in the definition of named RA queries. We proceed to explain how to obtain an unnamed RA query  $q_u$  that is equivalent to  $q_n$ , which means that the  $i$ -th attribute in the output of  $q_u$  corresponds to the  $A_i$ -attribute in the output of  $q_n$ . In the remainder of the proof, whenever we write a set of attributes as a set  $\{A_1, \dots, A_k\}$ , we assume that  $A_1 \leq \dots \leq A_k$ . The base cases are:

- If  $q_n = R$ , for a relation name  $R \in \mathbf{S}$  of sort  $\{A_1, \dots, A_k\}$ , then  $q_u = R$ .
- If  $q_n = \{(A : a)\}$ , then  $q_u = \{a\}$ .

For the inductive step, assume that  $q'_n$  and  $q''_n$  are named RA expressions of sort  $U' = \{A'_1, \dots, A'_k\}$  and  $U'' = \{A''_1, \dots, A''_\ell\}$ , respectively, and assume that they are equivalent to the unnamed RA expressions  $q'_u$  and  $q''_u$ , respectively.

- Let  $q_n = \sigma_\theta(q'_n)$ . Then  $q_u = \sigma_{\theta'}(q'_u)$ , where  $\theta'$  is the condition that is obtained from  $\theta$  by replacing each occurrence of attribute  $A'_i$  with  $i$ , for every  $i \in [k]$ . For example, if  $\theta$  is the condition  $(A'_1 \doteq A'_3) \wedge (A'_2 \neq b)$ , then  $\theta'$  is the condition  $(1 \doteq 3) \wedge (2 \neq b)$ .
- Let  $q_n = \pi_\alpha(q'_n)$  and  $\alpha \subseteq U'$ . Then  $q_u = \pi_{\alpha'}(q'_u)$ , where  $\alpha'$  is the list of all  $i \in [k]$  with  $A'_i \in \alpha$ .
- Let  $q_n = (q'_n \bowtie q''_n)$ . Then  $q_u = \pi_\alpha(q'_u \bowtie_\theta q''_u)$ , where  $\theta$  is the conjunction of all conditions  $i = j$  such that  $A'_i = A''_j$ , for  $i \in [k]$  and  $j \in [\ell]$ . To define  $\alpha$ , let  $\{A_1, \dots, A_m\} = U' \cup U''$  and let  $g : [m] \rightarrow [k + \ell]$  be such that

$$g(i) = \begin{cases} j & \text{if } A_i = A'_j, \\ k + j & \text{if } A_i = A''_j \text{ and } A''_j \in U'' - U'. \end{cases}$$

We now define  $\alpha = (g(1), \dots, g(m))$ . Therefore,  $\theta$  allows us to mimic the natural join on  $q'_n$  and  $q''_n$ , while  $\pi_\alpha$  is used for getting rid of redundant attributes and putting the attributes in an ordering that conforms to  $\prec$ .

- Let  $q_n = \rho_{A \rightarrow B}(q'_n)$ , where  $A = A'_i$  for some  $i \in [k]$ . Let  $j = |\{i \mid A'_i \prec B\}|$ . Then  $q_u = \pi_\alpha(q'_u)$ , where  $\alpha$  is obtained from  $(1, \dots, k)$  by deleting  $i$  and reinserting it right after  $j$  if  $j > 0$ , and at the beginning of the list if  $j = 0$ .
- Finally, if  $q_n = q'_n \cup q''_n$ , then  $q_u = q'_u \cup q''_u$ , where  $q'_u$  and  $q''_u$  are the unnamed RA expressions that are obtained by the induction hypothesis for  $q'_n$  and  $q''_n$ , respectively. The case when  $q_n = q'_n - q''_n$  is analogous.  $\square$

## Relational Algebra and SQL

---

In this chapter, we shed light on the relationship between relational algebra and SQL, the dominant query language in the relational database world. It is a complex language (the full description takes many hundreds of pages), and thus here we focus our attention on its core fragment.

### A Core of SQL

We assume that the reader by virtue of being interested in the principles of databases has some basic familiarity with relational databases and thus, by necessity, with SQL. Of course SQL is a language with a multitude of features, but its very core captures relational algebra queries. For now, we concentrate on that core part of the language and demonstrate its correspondence with relational algebra. The set of queries we consider are of the form

$$Q_1, Q_2 := \begin{array}{l} \text{SELECT } [\text{DISTINCT}] \text{ } \langle \text{select\_list} \rangle \\ \text{FROM } \langle \text{from\_list} \rangle \\ \text{WHERE } \langle \text{condition} \rangle \\ | \quad Q_1 \text{ UNION } Q_2 \quad | \quad Q_1 \text{ EXCEPT } Q_2 \end{array}$$

The `from_list` provides the list of relation names and subqueries used in the query, and also their *aliases*. For example, we can put `R AS R1` on the list, in which case `R1` is used as a new name for `R`. This could be used to shorten the name, e.g.,

`RelationWithAVeryLongName AS ShortName`

or to use the same relation more than once, in which case different aliases are needed. We can also put `Q AS Name` on the list, in which case the subquery `Q` is evaluated and the result of it, which is a relation, is given name `Name`. We shall do both in the examples very soon.

The `select_list` contains attributes of relation names mentioned in `FROM` or constants. For example, if we had `R AS R1` in `FROM` and `R` has an attribute `A`, we can have a reference to `R1.A` in that list. Likewise, if `from_list` contained `Q AS Name` and the output of `Q` contained attribute `B`, we can refer to `Name.B`. Constants can be output as well, e.g., for example, `5 AS C` will output the constant 5 as value of attribute `C`.

The keyword `DISTINCT` is to instruct the query to perform duplicate elimination. In general, SQL tables and query results are allowed to contain duplicates. For example, in a database containing two facts,  $R(a, b)$  and  $R(a, c)$ , projecting on the first column would result in *two* copies of  $a$ . We shall discuss duplicates in Chapter 44. In this Chapter, we will always assume that SQL queries only return sets, and omit `DISTINCT` from queries used in examples.

As *conditions* in this basic fragment we shall consider:

- equalities between attributes, e.g., `R.A = S.B`,
- equalities between attributes and constants, e.g., `Person.name = 'John'`,
- complex conditions built from these basic ones by using `AND`, `OR`, and `NOT`.

### Example 5.1: SQL Queries

Consider the FO query  $\varphi_1(x)$ , where  $\varphi_1$  is the FO formula (3.1). This can be written as the SQL query

```
SELECT P.pid
FROM Person AS P, Profession AS Pr1, Profession AS Pr2
WHERE P.pid = Pr1.pid
      AND P.pid = Pr1.pid
      AND NOT (Pr1.pname = Pr2.pname)
```

The formula  $\varphi_1$  mentions the relation name `Person` once, and the relation name `Profession` twice, and so does the above SQL query in the `FROM` clause (assigning different names to different occurrences, to avoid ambiguity). The first two conditions in the `WHERE` clause capture the use of the same variable  $x$  in three atomic subformulae of  $\varphi_1$ , whereas the last condition corresponds to the subformula  $\neg(u_1 = u_2)$ .

This query could alternatively be written as

```
SELECT T.id
FROM (SELECT P.pid AS id, Pr1.pname AS prof1,
           Pr2.pname AS prof2
      FROM Person AS P, Profession AS Pr1, Profession AS Pr2
      WHERE P.pid = Pr1.pid AND P.pid = Pr1.pid) AS T
WHERE NOT (T.prof1 = T.prof2)
```

that has a subquery in `FROM`.



Consider now the query  $\varphi_2(x, y)$ , where  $\varphi_2$  is the FO formula (3.2), which asks for IDs and names of people whose cities of birth were not recorded in the *City* relation. This can be expressed as the SQL query:

```
SELECT Person.pid, Person.pname
FROM Person
EXCEPT
SELECT Person.pid, Person.pname
FROM Person, City
WHERE Person.cid = City.cid
```

The first subquery asks for all people, the second subquery for those that have a city of birth recorded, and **EXCEPT** is their difference. This query returns people as pairs, consisting of their ID and their name.

## Core SQL to Relational Algebra

This section gives an intuition as to what happens when an SQL query is executed on a DBMS. A declarative query is translated into a procedural query to be executed. The real translation of SQL into RA is significantly more complex and, of course, captures many more features of SQL (and thus, the algebra implemented in DBMSs goes beyond the algebra we consider here). Nonetheless, the translation we outline presents the key ideas of the real-life translation.

Assume that we start with the query

```
SELECT  $\alpha_1$  AS  $B_1$ , ...,  $\alpha_n$  AS  $B_n$ 
FROM  $Q_1$  AS  $S_1$ , ...,  $Q_m$  AS  $S_m$ 
WHERE condition
```

where all relation names and subqueries in **FROM** have been renamed so they are different, and each  $\alpha_i$  is of the form  $S_j.A_p$ , that is, one of the attributes of the relation names in the **FROM** clause. Let  $\rho_i$  be the sequence of renaming operators that rename each attribute **A** of  $R_i$  to  $S_i.A$ . Let  $\rho_{\text{out}}$  be the sequence of renaming operators that forms the output, i.e., it renames each  $\alpha_i$  as  $B_i$ . Then, the translated query in relational algebra follows:

$$\rho_{\text{out}} \left( \pi_{\{\alpha_1, \dots, \alpha_n\}} \left( \sigma_{\text{condition}} \left( \rho_1(Q'_1) \bowtie \dots \bowtie \rho_m(Q'_m) \right) \right) \right),$$

where each  $Q'_i$  is the translation of  $Q_i$  into relational algebra. When  $Q_i$  is a relation  $R$  in the database, then  $Q'_i = R$ .

Essentially the **FROM** defines the join, **WHERE** provides the condition for selection, and **SELECT** is the final projection (hence, some clash of the naming conventions in SQL and RA).

The translation is then supplemented by translating **UNION** to RA's union  $\cup$  and **EXCEPT** to RA's difference  $-$ .

## Relational Algebra to Core SQL

We now show that (named) relational algebra queries can always be written as Core SQL queries. Let  $e$  be a named RA expression. We inductively translate  $e$  into an equivalent SQL query  $Q_e$  as follows.

**Base Expressions.** If  $e = R$ , and  $R$  has attributes  $A_1, \dots, A_n$ , then  $Q_e$  is

```
SELECT R.A1 AS A1, ... R.An AS An FROM R AS R
```

Of course in real SQL one can omit **AS R** and also use shorthand **\*** for listing all attributes, but here we keep the complete notation for the inductive construction.

If  $e = \{(A : a)\}$ , then  $Q_e$  is simply

```
SELECT a AS A
```

For the induction, we will assume that we can write all queries  $Q_e$  with a **SELECT** statement of the form

```
SELECT X1 AS A1, ..., Xn AS An
```

where

- each  $X_i$  is either a constant  $a$  or of the form  $Q.A$  or  $R.A$  for a query or a relation in **FROM**, and
- all attribute names  $A_i$  for  $i \in [n]$  are different

and a **FROM** statement that, if present, is of the form

```
FROM Q1 AS T1, ..., Qm AS Tm
```

where each  $Q_i$  is either a relation name, and all the names  $T_j$  for  $j \in [m]$  are different, and also different from names of relations in the database.

**Selection and Projection.** Assume that  $e$  is translated into

```
SELECT X1 AS A1, ..., Xn AS An
FROM Q1 AS T1, ..., Qm AS Tm
WHERE condition
```

- Then,  $\sigma_\theta(e)$  is translated into

```
SELECT X1 AS A1, ..., Xn AS An
FROM Q1 AS T1, ..., Qm AS Tm
WHERE condition AND Cθ
```

where  $C_\theta$  expresses the condition  $\theta$  in SQL syntax, but uses the names  $X_i$  instead of  $A_i$  due to the ordering in which SQL applies aliases. For instance, if  $\theta$  is  $(A_1 \doteq A_2) \wedge \neg(A_3 \doteq 1)$  then  $C_\theta$  is  $(X_1 = X_2)$  AND NOT  $(X_3 = 1)$ .

To illustrate this with an example, consider the relation  $R[A, B]$  and the query  $Q_e$  given as

```
SELECT R.A AS A, R.B AS C, 5 AS D
FROM R AS R
WHERE R.A = R.B
```

and assume that  $\theta$  is  $(A \doteq 3) \vee \neg(C \doteq D)$ . Then  $\sigma_\theta(e)$  is translated into

```
SELECT R.A AS A, R.B AS C, 5 AS D
FROM R AS R
WHERE R.A = R.B AND ((R.A = 3) OR NOT (R.B = 5))
```

- Furthermore,  $\pi_\alpha(e)$  is translated into

```
SELECT Xi1 AS Ai1, ..., Xik AS Aik
FROM Q1 AS T1, ..., Qm AS Tm
WHERE condition
```

where  $A_{i_1}, \dots, A_{i_k}$  are the elements from the set  $\alpha$ .

**Rename.** Assume now that  $e$  is translated into

```
SELECT ..., Qi.Aj AS A, ...
FROM Q1 AS T1, ..., Qm AS Tm
WHERE condition
```

Then,  $\rho_{A \rightarrow B}(e)$  is translated into

```
SELECT ..., Qi.Aj AS B, ...
FROM Q1 AS T1, ..., Qm AS Tm
WHERE condition
```

**Join** Assume now that we have expressions  $e_1$  and  $e_2$  that are translated into queries  $Q_{e_1}$  and  $Q_{e_2}$ , respectively. Then their natural join  $e_1 \bowtie e_2$  is translated into

```
SELECT T1.* AS  $\bar{A}$ , {T2.C | C ∈ sort(Qe2) − sort(Qe1)} AS  $\bar{B}$ 
FROM Qe1 AS T1, Qe2 AS T2
WHERE T1.X1 = T2.X1 AND ... AND T1.Xk = T2.Xk
```

where  $X_1, \dots, X_k$  are the attributes in  $\text{sort}(Q_{e_1}) \cap \text{sort}(Q_{e_2})$  and  $T_1.*$  abbreviates the list of all attributes of  $Q_{e_1}$ , with each of them assigned a name from the tuple  $\bar{A}$ . That is, all attributes of  $Q_{e_1}$  are in the output, together with attributes of  $Q_{e_2}$  that do not occur in  $Q_{e_1}$ , appropriately renamed, but only if the common attributes of  $Q_{e_1}$  and  $Q_{e_2}$  are equal.

**Difference** If  $e = e_1 - e_2$ , then  $Q_e$  is

$$(Q_{e_1}) \text{ EXCEPT } (Q_{e_2})$$

**Union** Finally, if  $e = e_1 \cup e_2$ , then  $Q_e$  is

$$(Q_{e_1}) \text{ UNION } (Q_{e_2})$$

This completes the translation from (named) RA to Core SQL.

## Other SQL Features Captured by RA

A very important feature of SQL is using *subqueries*. We have seen them in **FROM**, but they can be used in conditions in **WHERE** as well. They are very convenient for a declarative presentation of queries (although from the point of view of expressiveness of the language, they can be omitted). Consider, for example, the query that computes the difference of two relations  $R$  and  $S$  with one attribute  $A$ . We could use **EXCEPT**, but using subqueries we can also write

```
SELECT R.A
FROM R
WHERE R.A NOT IN (SELECT S.A FROM S)
```

saying that we need to return elements of  $R$  that are not present in  $S$ , or

```
SELECT R.A
FROM R
WHERE NOT EXISTS (SELECT S.A FROM S WHERE S.A = R.A)
```

which asks for elements  $a$  of  $R$  such that there is no  $b$  in  $S$  satisfying  $a = b$ . Both queries express the difference.

### Example 5.2: Subqueries in SQL

Consider the query  $\varphi(x, y)$ , where  $\varphi$  is the FO formula (3.2), which asks for IDs and names of people whose cities of birth were not recorded in the **City** relation. This can also be written as the SQL query:

```
SELECT P.pid, P.pname
```

```
FROM Person AS P
WHERE P.cid NOT IN (SELECT City.cid FROM City)
```

The above two forms of subqueries, using **IN** and **EXISTS**, potentially with negation **NOT**, correspond to adding the following two types of selection conditions to RA, which, nevertheless, do not increase the expressiveness of RA; see Exercise [1.5](#).

- $\bar{a} \in e$ , where  $\bar{a}$  is a tuple of terms and  $e$  is an expression, checking whether  $\bar{a}$  belongs to the result of the evaluation of  $e$ , and
- $\text{empty}(e)$ , checking if the result of the evaluation of  $e$  is empty.

Such an addition does not increase expressiveness (Exercise [1.6](#)) but makes writing queries easier.

## Other SQL Features Not Captured by RA

**Bag Semantics.** As mentioned already, SQL's data model is based on bags, i.e., the same tuple may occur multiple times in a database or output of a query. Here we tacitly assumed that all relations are sets and each **SELECT** is followed by **DISTINCT** to ensure that duplicates are eliminated. To see how RA operations change in the presence of duplicates, see Chapter [44](#).

**Grouping and Aggregation.** An extremely common feature of SQL queries is the use of aggregation and grouping. Aggregation allows numerical functions to be applied to entire columns, for example, to find the total salary of all employees in a company. Grouping allows such columns to be split according to a value of some attribute; an example of this is a query that returns the total salary of each department in a company. These features will be discussed in more detail in Chapter [33](#).

**Nulls.** SQL databases permit missing values in tuples. To handle this, they allow a special element **null** to be placed as a value. The handling of nulls is very different though from the handling of values from **Const**, and even the notion of query output changes in this case. These issues are discussed in detail in Chapters [39](#) and [40](#).

**Types.** In SQL databases, attributes must be typed, i.e., all values in a column must have the same type. There are standard types such as numbers (integers, floats), strings of various length, fixed or varying, date, time, and many others. With the exception of the consideration of arithmetic operations (Chapter [33](#)), this is a subject that we do study in this book.



## Equivalence of Logic and Algebra

In this chapter, we prove that the declarative query language based on FO, and the procedural query language RA have the same expressive power, which is a fundamental result of relational database theory. Recall that we focus on the unnamed version of RA for reasons that we explained earlier.

### Theorem 6.1

The languages of RA queries and of FO queries are equally expressive.

The proof of Theorem 6.1 boils down to showing that, for a schema  $\mathbf{S}$ , the following statements hold:

- (a) For every RA expression  $e$  over  $\mathbf{S}$ , there exists an FO query  $q_e$  such that  $q_e(D) = e(D)$ , for every database  $D$  of  $\mathbf{S}$ .
- (b) For every FO query  $q$  over  $\mathbf{S}$ , there exists an RA expression  $e_q$  such that  $e_q(D) = q(D)$ , for every database  $D$  of  $\mathbf{S}$ .

In the proof of the above, we need a mechanism that allows us to substitute variables in formulae. For an FO formula  $\varphi$  and variables  $\{x_1, \dots, x_n\}$ , we denote by  $\varphi[x_1/y_1, \dots, x_n/y_n]$  the formula obtained from  $\varphi$  by simultaneously replacing each  $x_i$  with  $y_i$ . We also use the notation  $\exists\{x_1, \dots, x_n\}\varphi$  for a set of variables  $\{x_1, \dots, x_n\}$  as an abbreviation for  $\exists x_1 \dots \exists x_n \varphi$ . Notice that the ordering of quantification is irrelevant for the semantics of this formula.

### From RA to FO

We first show (a) by induction on the structure of  $e$ . The base cases are:

- If  $e = R$  for  $R \in \text{Dom}(\mathbf{S})$ , then the FO query is  $\varphi_e(x_1, \dots, x_{\text{ar}(R)})$ , where

$$\varphi_e = R(x_1, \dots, x_{\text{ar}(R)})$$

with all the variables  $x_1, \dots, x_{\text{ar}(R)}$  being different.

- If  $e$  is  $\{a\}$  with  $a \in \text{Const}$ , then the FO query is  $\varphi_e(x)$ , where

$$\varphi_e = (x = a).$$

We now proceed with the induction step. Assume that  $e$  and  $e'$  are RA expressions over  $\mathbf{S}$  for which we have equivalent FO queries  $\varphi_e(x_1, \dots, x_k)$  and  $\varphi_{e'}(y_1, \dots, y_\ell)$ , respectively. By renaming variables, we can assume, without loss of generality, that  $\{x_1, \dots, x_k\}$  and  $\{y_1, \dots, y_\ell\}$  are disjoint.

- Let  $\theta$  be a condition over  $\{1, \dots, k\}$ . Taking  $\bar{x} = (x_1, \dots, x_k)$ , we inductively define the formula  $\theta[\bar{x}]$  as follows:
  - if  $\theta$  is  $i \doteq j$ ,  $i \doteq a$ ,  $i \not\doteq j$ , or  $i \not\doteq a$ , then  $\theta[\bar{x}]$  is  $x_i = x_j$ ,  $x_i = a$ ,  $x_i \neq x_j$ , or  $x_i \neq a$ , respectively,
  - if  $\theta = \theta_1 \wedge \theta_2$ , then  $\theta[\bar{x}] = \theta_1[\bar{x}] \wedge \theta_2[\bar{x}]$ ,
  - if  $\theta = \theta_1 \vee \theta_2$ , then  $\theta[\bar{x}] = \theta_1[\bar{x}] \vee \theta_2[\bar{x}]$ , and
  - if  $\theta = \neg\theta_1$ , then  $\theta[\bar{x}] = \neg\theta_1[\bar{x}]$ .

Then, the FO query equivalent to  $\sigma_\theta(e)$  is  $\varphi_{\sigma_\theta(e)}(\bar{x}) = \varphi_e(\bar{x}) \wedge \theta[\bar{x}]$ .

- Let  $\alpha = (i_1, \dots, i_p)$  be a list of numbers from  $\{1, \dots, k\}$ . The FO query equivalent to  $\pi_\alpha(e)$  is  $\varphi_{\pi_\alpha(e)}(x_{i_1}, \dots, x_{i_p})$ , where  $\varphi_{\pi_\alpha(e)}$  is the formula

$$\exists(\{x_1, \dots, x_n\} - \{x_{i_1}, \dots, x_{i_p}\}) \varphi_e.$$

Notice that, if  $\alpha$  has repetitions, then  $(x_{i_1}, \dots, x_{i_p})$  has repeated variables. For example, if  $e = R$ , where  $R$  is binary, and  $\alpha = (1, 1)$ , then the FO query is  $\varphi_e(x_1, x_1)$  with  $\varphi_e = \exists x_2 R(x_1, x_2)$ .

- The FO query equivalent to  $e \times e'$  is  $\varphi_{e \times e'}(x_1, \dots, x_k, y_1, \dots, y_\ell)$ , where  $\varphi_{e \times e'}$  is the formula

$$\varphi_e \wedge \varphi_{e'}.$$

- Let  $e \cup e'$  be an RA expression, which is only well-defined if  $k = \ell$ . The equivalent FO query is  $\varphi_{e \cup e'}(x_1, \dots, x_k)$ , where  $\varphi_{e \cup e'}$  is

$$\varphi_e \vee (\varphi_{e'}[y_1/x_1, \dots, y_k/x_k]).$$

- Let  $e - e'$  be an RA expression, which is only well-defined if  $k = \ell$ . The equivalent FO query is  $\varphi_{e - e'}(x_1, \dots, x_k)$ , where  $\varphi_{e - e'}$  is

$$\varphi_e \wedge \neg(\varphi_{e'}[y_1/x_1, \dots, y_k/x_k]).$$

We leave the verification of the construction, that is, the inductive proof of the equivalence of  $e$  and  $\varphi_e(\bar{x})$ , to the reader. This concludes part (a).



## From FO to RA

For proving (b), we assume that relational atoms do not mention constants, which we observed in Chapter 3 is always possible. We also consider a slight generalization of FO queries that will simplify the induction:  $\varphi(x_1, \dots, x_n)$  is an FO query even if the free variables of  $\varphi$  are a subset of  $\{x_1, \dots, x_n\}$ . The semantics of such a query  $\varphi(x_1, \dots, x_n)$  is the usual semantics of the FO query  $\varphi'(x_1, \dots, x_n)$ , where  $\varphi'$  is the formula  $\varphi \wedge (x_1 = x_1) \wedge \dots \wedge (x_n = x_n)$ .

Let  $q$  be an FO query of the form  $\varphi(x_1, \dots, x_n)$ . We can assume, without loss of generality, that  $\varphi$  is in *prenex normal form*, that is, of the form

$$Q_k \cdots Q_1 \varphi_{\text{qf}} ,$$

where

- each  $Q_j$  is of the form  $\exists y_j$  or  $\neg \exists y_j$ ,
- $\varphi_{\text{qf}}$  is quantifier-free and has (free) variables  $y_1, \dots, y_m$ ,
- $\{x_1, \dots, x_n\} = \{y_{k+1}, \dots, y_m\}$ , and
- $\varphi_{\text{qf}}$  only uses the Boolean operators  $\vee$  and  $\neg$ .

Let  $\text{Dom}(\varphi) = \{a_1, \dots, a_\ell\}$ . First, we build an RA expression  $\text{Adom}$  for the active domain, that is,

$$\text{Adom} = \bigcup_{i=1}^{\ell} \{a_i\} \cup \bigcup_{R[n] \in \mathbf{S}} (\pi_1(R) \cup \dots \cup \pi_n(R)) .$$

In the following, we denote by  $\text{Adom}^i$ , for  $i \in \mathbb{N}$ , the  $i$ -fold Cartesian product

$$\underbrace{\text{Adom} \times \dots \times \text{Adom}}_i .$$

We construct for each subformula  $\psi$  of  $\varphi$  an RA query  $e_\psi$ . The induction hypothesis consists of two parts.

- (1) For each subformula  $\psi$  of  $\varphi_{\text{qf}}$ , the expression  $e_\psi$  has arity  $m$  and is equivalent to the FO query  $\psi(y_1, \dots, y_m)$ .
- (2) For all the other subformulae  $\psi$  of  $\varphi$ , it holds that  $\psi = Q_j \cdots Q_1 \varphi_{\text{qf}}$ , for  $j \in [k]$ ,  $\text{FV}(\psi) = \{y_{j+1}, \dots, y_m\}$ , and the expression  $e_\psi$ , which has arity  $m - j$ , is equivalent to the FO query  $Q_j \cdots Q_1 \varphi_{\text{qf}}(y_{j+1}, \dots, y_m)$ .

The inductive construction defines the expression

$$e_\psi = \begin{cases} \pi_{1,\dots,m}(\sigma_{i_1=m+1,\dots,i_j=m+j}(\text{Adom}^m \times R)) & \text{if } \psi \text{ is } R(y_{i_1}, \dots, y_{i_j}) \\ \sigma_{i=j}(\text{Adom}^m) & \text{if } \psi \text{ is } y_i = y_j \\ \sigma_{i=a}(\text{Adom}^m) & \text{if } \psi \text{ is } y_i = a \\ e_{\psi_1} \cup e_{\psi_2} & \text{if } \psi \text{ is } (\psi_1 \vee \psi_2) \\ \text{Adom}^m - e_{\psi'} & \text{if } \psi \text{ is } \neg\psi', \text{ and} \\ & \psi \text{ is a subformula of } \varphi_{\text{qf}} \\ \pi_{2,\dots,m-j+1}(e_{\psi'}) & \text{if } \psi \text{ is } \exists y_j \psi' \text{ and } Q_j = \exists y_j \\ \text{Adom}^{m-j} - \pi_{2,\dots,m-j+1}(e_{\psi'}) & \text{if } \psi \text{ is } \neg\exists y_j \psi' \end{cases}$$

We leave the proof that the inductive construction gives an expression that is equivalent to  $\varphi(y_{k+1}, \dots, y_m)$  to the reader. To obtain an expression equivalent to  $\varphi(x_1, \dots, x_n)$ , observe that  $x_i \in \{y_{k+1}, \dots, y_m\}$  for every  $i \in [n]$ . Therefore, there exists a function  $f : [n] \rightarrow [m - k]$  such that  $x_i = y_{f(i)}$  for every  $i \in [n]$ . This means that the expression  $\pi_{(f(1), \dots, f(n))} e_\varphi$  is equivalent to  $\varphi(x_1, \dots, x_n)$ .

## First-Order Query Evaluation

In this chapter, we study the complexity of evaluating first-order queries, that is, **FO-Evaluation**. Recall that this is the problem of checking whether  $\bar{a} \in q(D)$  for an FO query  $q$ , a database  $D$ , and a tuple  $\bar{a}$  over **Const**.

### Combined Complexity

We first concentrate on the combined complexity of the problem, that is, when the input consists of the query  $q$ , the database  $D$ , and the tuple  $\bar{a}$ .

#### Theorem 7.1

FO-Evaluation is PSPACE-complete.

*Proof.* For simplicity, we consider the case where the query does not use constant values and leave the extension to arbitrary FO queries as an exercise. We start with the upper bound. Consider an FO query  $q = \varphi(\bar{x})$  without constants, a database  $D$ , and a tuple  $\bar{a}$ . First, we observe that, if  $\bar{a}$  contains values  $a'_1, \dots, a'_\ell$  that are not in  $\text{Dom}(D)$ , we can replace them with arbitrary elements  $a''_1, \dots, a''_\ell$  in  $\text{Const} - \text{Dom}(D)$  that can be represented using  $O(\|D\|)$  many bits. We can also assume that the tuples  $\bar{x} = (x_1, \dots, x_n)$  and  $\bar{a} = (a_1, \dots, a_m)$  are *compatible*, that is, they have the same length (i.e.,  $n = m$ ), and  $x_i = x_j$  implies  $a_i = a_j$  for every  $i, j \in [n]$ . Indeed, if  $\bar{x}$  and  $\bar{a}$  are not compatible, which can be easily checked using logarithmic space, then  $\bar{a} \notin q(D)$  holds trivially. We can also assume that  $\varphi$  uses only  $\neg$ ,  $\vee$ , and  $\exists$  (see Exercise 1.1).

By Definition 3.6,  $\bar{a} \in q(D)$  if and only if  $(D, \eta) \models \varphi$  with  $\eta$  being the assignment for  $\varphi$  over  $D$  such that  $\eta(\bar{x}) = \bar{a}$ . Therefore, to establish that **FO-Evaluation** is in PSPACE, it suffices to show that the problem of checking whether  $(D, \eta) \models \varphi$  is in PSPACE. This is done by exploiting the recursive procedure **EVALUATION**, depicted in Algorithm 1. Notice that the algorithm

performs simple Boolean tests for determining its output values, like testing if  $R(\eta(\bar{x}))$  is an element of  $D$  in line [1](#) or whether  $\eta(x_i) = \eta(x_j)$  in line [2](#). It is not difficult to verify that  $(D, \eta) \models \varphi$  if and only if  $\text{EVALUATION}(\varphi, D, \eta) = \text{true}$ . It remains to argue that  $\text{EVALUATION}(\varphi, D, \eta)$  uses polynomial space.

---

**Algorithm 1**  $\text{EVALUATION}(\varphi, D, \eta)$ 


---

**Input:** An FO formula  $\varphi$ , a database  $D$ , and an assignment  $\eta$  for  $\varphi$  over  $D$ .

**Output:** **true** if  $(D, \eta) \models \varphi$ , and **false** otherwise.

---

```

1: if  $\varphi$  is of the form  $R(\bar{x})$  then return  $R(\eta(\bar{x})) \in D$ 
2: else if  $\varphi$  is of the form  $(x_i = x_j)$  then return  $\eta(x_i) = \eta(x_j)$ 
3: else if  $\varphi$  is of the form  $(x_i = a)$  then return  $\eta(x_i) = a$ 
4: else if  $\varphi$  is of the form  $\neg\varphi'$  then return  $\neg\text{EVALUATION}(\varphi', D, \eta)$ 
5: else if  $\varphi$  is of the form  $\varphi' \vee \varphi''$  then
6:   return  $\text{EVALUATION}(\varphi', D, \eta) \vee \text{EVALUATION}(\varphi'', D, \eta)$ 
7: else if  $\varphi$  is of the form  $\exists x \varphi'$  then
8:   return  $\bigvee_{a \in \text{Dom}(D)} \text{EVALUATION}(\varphi', D, \eta[x/a])$ 
9:    $\triangleright \eta[x/a]$  extends  $\eta$  by setting  $\eta(x) = a$ .
```

---

**Lemma 7.2.**  $\text{EVALUATION}(\varphi, D, \eta)$  runs in  $\text{SPACE}(O(\|\varphi\|^2 \cdot \log \|D\|))$ .

*Proof.* Observe that the total space used by  $\text{EVALUATION}(\varphi, D, \eta)$  is its recursion depth times the space needed by each recursive call. It is clear that the recursion depth is  $O(\|\varphi\|)$ . We proceed to argue, by induction on the structure of  $\varphi$ , that each recursive call uses  $O(\|\varphi\| \cdot \log \|D\|)$  space on a Turing Machine, which in turn implies that the total space used by  $\text{EVALUATION}(\varphi, D, \eta)$  is  $O(\|\varphi\|^2 \cdot \log \|D\|)$ .

- Assume first that  $\varphi = R(\bar{x})$ . In this case, the algorithm checks whether  $R(\eta(\bar{x})) \in D$ . The space needed to store  $\eta(\bar{x})$  on the work tape (adopting the encoding discussed in Appendix [C](#)) is  $O(\|\varphi\| \cdot \log \|D\|)$ . Furthermore, as shown in Appendix [C](#) (see Lemma [C.1](#)), for a tuple  $\bar{t}$  over  $\text{Dom}(D)$ , we can check whether  $R(\bar{t}) \in D$  using  $O(\text{ar}(R) \cdot \log \|D\|)$  space if  $\text{ar}(R) > 0$ , and  $O(\log \|D\|)$  space if  $\text{ar}(R) = 0$ . Therefore, in the worst-case where  $\text{ar}(R) > 0$ , we can check whether  $R(\eta(\bar{x})) \in D$  using space

$$O(\|\varphi\| \cdot \log \|D\|) + O(\text{ar}(R) \cdot \log \|D\|).$$

Since  $\text{ar}(R) \leq \|\varphi\|$ , the total space used is  $O(\|\varphi\| \cdot \log \|D\|)$ .

- When  $\varphi = (x_i = x_j)$ , the algorithm checks whether  $\eta(x_i) = \eta(x_j)$ , which can be done using  $O(\|\varphi\| \cdot \log \|D\|)$  space by simply storing the tuples  $\eta(x_i)$  and  $\eta(x_j)$  (adopting the encoding from Appendix [C](#)) on the work tape, and then check that they are equal. The case  $\varphi = (x_i = a)$  is analogous.

- When  $\varphi = \neg\varphi'$ , the algorithm computes the value  $\neg\text{EVALUATION}(\varphi', D, \eta)$ , which, by induction hypothesis, can be done using  $O(\|\varphi\| \cdot \log \|D\|)$  space.
- When  $\varphi = \varphi' \vee \varphi''$ , the algorithm computes  $\text{EVALUATION}(\varphi', D, \eta) \vee \text{EVALUATION}(\varphi'', D, \eta)$ , which, by induction hypothesis, can be done using  $O(\|\varphi\| \cdot \log \|D\|)$  space.
- Finally, assume that  $\varphi = \exists x \varphi'$ . In this case, the algorithm computes  $\bigvee_{a \in \text{Dom}(D)} \text{EVALUATION}(\varphi', D, \eta[x/a])$ . This is done by iterating over the constants of  $\text{Dom}(D)$  in the order provided by the encoding of  $D$  (see Appendix [C](#)), and reusing the space used by the previous iteration. Thus, it suffices to argue that computing the value  $\text{EVALUATION}(\varphi', D, \eta[x/a])$ , for some value  $a \in \text{Dom}(D)$ , can be done using  $O(\|\varphi\| \cdot \log \|D\|)$  space. The latter clearly holds by induction hypothesis, and the claim follows.  $\square$

For the lower bound, we provide a reduction from QSAT, which we know is PSPACE-complete (see Appendix [B](#)). Consider an input to QSAT given by

$$\psi = \exists \bar{x}_1 \forall \bar{x}_2 \exists \bar{x}_3 \dots Q_n \bar{x}_n \psi' \langle \bar{x}_1, \dots, \bar{x}_n \rangle,$$

where  $Q_n = \forall$  if  $n$  is even, and  $Q_n = \exists$  if  $n$  is odd. We assume that  $\psi'$  is in negation normal form, which means that negation is only applied to variables, since QSAT remains PSPACE-hard. We construct the database

$$D = \{\text{Zero}(0), \text{One}(1)\}$$

and the Boolean FO query

$$q_\psi = \exists \bar{x}_1 \forall \bar{x}_2 \exists \bar{x}_3 \dots Q_n \bar{x}_n \psi'',$$

where  $\psi''$  is obtained from  $\psi'$  by replacing each occurrence of the literal  $x$  by  $\text{One}(x)$ , and each occurrence of the literal  $\neg x$  by  $\neg \text{One}(x)$ . The only reason why we add  $\text{Zero}(0)$  to the database is to ensure that 0 is in the active domain. For example, if  $\psi'(x_1, x_2, x_3) = (x_1 \wedge x_2) \vee (\neg x_1 \wedge x_3)$ , then  $\psi'' = (\text{One}(x_1) \wedge \text{One}(x_2)) \vee (\neg \text{One}(x_1) \wedge \text{One}(x_3))$ . It is not hard to verify that  $\psi$  is satisfiable if and only if  $D \models q_\psi$  (we leave the proof as an exercise).  $\square$

Note that  $q(D)$ , for an FO query  $q = \varphi(\bar{x})$  and a database  $D$ , can also be computed in polynomial space as follows: iterate over all tuples  $\bar{a}$  over  $\text{Dom}(D)$  that are compatible with  $\bar{x}$ , and output  $\bar{a}$  if and only if  $\text{EVALUATION}(\varphi, D, \eta) = \text{true}$  with  $\eta$  being the assignment for  $\varphi$  over  $D$  such that  $\eta(\bar{x}) = \bar{a}$ . It is easy to show that this procedure runs in polynomial space. This, of course, relies on the fact that the running space of a Turing Machine with output is defined without considering the output tape; see Appendix [B](#) for details.

## Data Complexity

How can it be that databases are so successful in practice, even though Theorem [7.1](#) proves that the most essential database problem is PSPACE-complete,

a complexity class that we consider to be intractable? If we take a closer look at the lower bound proof of Theorem 7.1, we see that the entire difficulty of the problem is encoded in the query. In fact, the database  $D = \{\text{Zero}(0), \text{One}(1)\}$  consists of only two atoms, whereas the query  $q$  can be arbitrarily large. This is in contrast to what we typically experience in practice, where databases are orders of magnitude larger than queries, which means that databases and queries contribute in different ways to the complexity of evaluation. This brings us to the data complexity of FO query evaluation.

As discussed in Chapter 2, when we study the data complexity of query evaluation, we essentially consider the query to be fixed, and only the database and the candidate output are considered as input. Formally, we are interested in the complexity of the problem  $q$ -Evaluation for an FO query  $q$ , which takes as input a database  $D$  and a tuple  $\bar{a}$  over  $\text{Dom}(D)$ , and asks whether  $\bar{a} \in q(D)$ . Recall that, by convention, we say that FO-Evaluation is in a complexity class  $\mathcal{C}$  in data complexity if  $q$ -Evaluation is in  $\mathcal{C}$  for every FO query  $q$ .

### Theorem 7.3

FO-Evaluation is in DLOGSPACE in data complexity.

*Proof.* Fix an FO query  $q = \varphi(\bar{x})$ . Our goal is to show that  $q$ -Evaluation is in DLOGSPACE. As for Theorem 7.1, we prove the result for the case where  $\bar{a}$  is over  $\text{Dom}(D)$ , and leave the extension to tuples over **Const** as an exercise.

Consider a database  $D$ , and a tuple  $\bar{a}$  over  $\text{Dom}(D)$ . Observe that the input word encoding  $D$  on a Turing Machine has length  $O(\|D\| \log \|D\|)$ . We therefore need to prove that  $q$ -Evaluation can be solved in  $\text{SPACE}(O(\log(\|D\| \log \|D\|))) = \text{SPACE}(O(\log \|D\|))$ . As explained in the proof of Theorem 7.1, we can assume that the relational atoms in  $\varphi$  do not contain constants, the tuples  $\bar{x} = (x_1, \dots, x_n)$  and  $\bar{a} = (a_1, \dots, a_m)$  are compatible, and that  $\varphi$  uses only  $\neg$ ,  $\vee$ , and  $\exists$ . To prove our claim it suffices to show that checking whether  $(D, \eta) \models \varphi$  with  $\eta$  being the assignment for  $\varphi$  over  $D$  such that  $\eta(\bar{x}) = \bar{a}$  is in DLOGSPACE. This is done by exploiting the procedure  $\text{EVALUATION}_\varphi$ , which takes as input  $D$  and  $\eta$ , and it simply calls the procedure  $\text{EVALUATION}$  given in Algorithm 1 with input  $\varphi$ ,  $D$  and  $\eta$ . More precisely,  $\text{EVALUATION}_\varphi(D, \eta)$  does the following: if  $\text{EVALUATION}(\varphi, D, \eta) = \text{true}$ , then return **true**; otherwise, return **false**. From the correctness of  $\text{EVALUATION}$ , it is straightforward to see that  $(D, \eta) \models \varphi$  if and only if  $\text{EVALUATION}_\varphi(D, \eta) = \text{true}$ . Moreover, from the complexity analysis of  $\text{EVALUATION}$  performed in the proof of Theorem 7.1 and the fact that  $\varphi$  is fixed, we conclude that  $\text{EVALUATION}_\varphi(D, \eta)$  runs in space  $O(\log \|D\|)$ , and the claim follows.  $\square$

Theorem 7.3 essentially tells us that fixing the query indeed has a big impact to the complexity of evaluation, which goes from PSPACE to DLOGSPACE. Actually, FO-Evaluation is in  $\text{AC}_0$  in data complexity, a class that is properly contained in DLOGSPACE. The class  $\text{AC}_0$  consists of those languages that are

accepted by polynomial-size circuits of constant depth and unbounded fan-in (the number of inputs to their gates). This is the reason why **FO-Evaluation** is often regarded as an “embarrassingly parallel” task.





## Static Analysis

We now study central static analysis tasks for FO queries. We focus on satisfiability, containment, and equivalence, which are key ingredients for query optimization. As we shall see, these problems are undecidable for FO queries. This in turn implies that, given an FO query, computing an optimal equivalent FO query is, in general, algorithmically impossible.

### Satisfiability

A query  $q$  is *satisfiable* if there is a database  $D$  such that  $q(D)$  is non-empty. It is clear that a query that is not satisfiable it is not a useful query since its output on a database is always empty. In relation to satisfiability, we consider the following problem, parameterized by a query language  $\mathcal{L}$ .

**Problem:  $\mathcal{L}$ -Satisfiability**

**Input:** A query  $q$  from  $\mathcal{L}$

**Output:** **true** if there is a database  $D$  such that  $q(D) \neq \emptyset$ , and **false** otherwise

Notice that satisfiability is, in a sense, the most elementary static analysis question one can ask about a query: “does there exist a database *at all* for which the query returns an answer?” Indeed, if there does not, then optimizing the query is extremely simple: one can just always return the empty set of answers, independently of the input database.

We are asking the satisfiability question focussing on finite databases. In the case of possibly infinite databases, we know from a classical result in logic that goes back in the 1930s, known as Church’s Theorem (sometimes called Church-Turing Theorem), that checking for satisfiability is undecidable. The problem remains undecidable even for finite databases, a result proved by

Trakhtenbrot in the 1950s, i.e., several years after Church's Theorem. In what follows we present Trakhtenbrot's Theorem.

**Theorem 8.1: Trakhtenbrot's Theorem**

FO-Satisfiability is undecidable.

*Proof.* The proof is by reduction from the halting problem for Turing Machines; details on Turing Machines can be found in Appendix B. It is well-known that the problem of deciding whether a (deterministic) Turing Machine  $M = (Q, \Sigma, \delta, s)$  halts on the empty word is undecidable. Our goal is to construct a Boolean FO query  $q_M$  such that the following are equivalent:

1.  $M$  halts on the empty word.
2. There exists a database  $D$  such that  $q_M(D) = \text{true}$ .

The Boolean FO query  $q_M$  will be over the schema

$$\{\prec[2], \text{First}[1], \text{Succ}[2]\} \cup \{\text{Symbol}_a[2] \mid a \in \Sigma\} \cup \{\text{Head}[2], \text{State}[2]\}.$$

The intuitive meaning of the above relation names is the following:

- $\prec(\cdot, \cdot)$  encodes a strict linear order over the underlying domain, which will be used to simulate the time steps of the computation of  $M$  on the empty word, and the tape cells of  $M$ .
- $\text{First}(\cdot)$  contains the first element from the linear order  $\prec$ .
- $\text{Succ}(\cdot, \cdot)$  encodes the successor relation over the linear order  $\prec$ .
- $\text{Symbol}_a(t, c)$ : at time instant  $t$ , the tape cell  $c$  contains the symbol  $a$ .
- $\text{Head}(t, c)$ : at time instant  $t$ , the head points at cell  $c$ .
- $\text{State}(t, p)$ : at time instant  $t$ , the machine  $M$  is in state  $p$ .

Having the above schema in place, we can now proceed with the definition of the Boolean FO query  $q_M$ , which is of the form

$$\varphi_{\prec} \wedge \varphi_{\text{first}} \wedge \varphi_{\text{succ}} \wedge \varphi_{\text{comp}},$$

where  $\varphi_{\prec}$ ,  $\varphi_{\text{first}}$  and  $\varphi_{\text{succ}}$  are FO sentences that are responsible for defining the relations  $\prec$ ,  $\text{First}$  and  $\text{Succ}$ , respectively, while  $\varphi_{\text{comp}}$  is an FO sentence responsible for mimicking the computation of  $M$  on the empty word. The definitions of the above FO sentences follow. For the sake of readability, we write  $x \prec y$  instead of the formal  $\prec(x, y)$ .

*The Sentence  $\varphi_{\prec}$*

This sentence simply expresses that the binary relation  $\prec$  over the underlying domain is total, irreflexive, and transitive:

$$\begin{aligned} \forall x \forall y \left( \neg(x = y) \rightarrow (x \prec y \vee y \prec x) \right) \wedge \\ \forall x \neg(x \prec x) \wedge \\ \forall x \forall y \forall z \left( (x \prec y \wedge y \prec z) \rightarrow x \prec z \right). \end{aligned}$$

Note that irreflexivity and transitivity together imply that the relation  $\prec$  is also asymmetric, i.e.,  $\forall x \forall y \neg(x \prec y \wedge y \prec x)$ .

*The Sentence  $\varphi_{first}$*

This sentence expresses that  $\text{First}(\cdot)$  contains the smallest element over  $\prec$ :

$$\forall x \left( \text{First}(x) \leftrightarrow \forall y (x = y \vee x \prec y) \right).$$

*The Sentence  $\varphi_{succ}$*

It simply defines the successor relation over  $\prec$  as expected:

$$\forall x \forall y \left( \text{Succ}(x, y) \leftrightarrow (x \prec y \wedge \neg \exists z (x \prec z \wedge z \prec y)) \right).$$

*The Sentence  $\varphi_{comp}$*

Assume that the set of states of  $M$  is  $Q = \{p_1, \dots, p_k\}$ , where  $p_1 = s$  is the start state,  $p_2 = \text{"yes"}$  is the accepting state, and  $p_3 = \text{"no"}$  is the rejecting state. The key idea is to associate to each state of  $M$  a distinct element of the underlying domain, which in turn will allow us to refer to the states of  $M$ . Thus,  $\varphi_{comp}$  is defined as the following FO sentence; for a subformula  $\psi$  of  $\varphi_{comp}$ , we write  $\psi\langle\bar{x}\rangle$  to indicate that  $\text{FV}(\psi)$  consists of the variables in  $\bar{x}$ :

$$\begin{aligned} \exists x_1 \cdots \exists x_k \left( \bigwedge_{i,j \in [k]: i < j} \neg(x_i = x_j) \wedge \varphi_{\text{start}}\langle x_1 \rangle \wedge \varphi_{\text{consistent}}\langle x_1, \dots, x_k \rangle \wedge \right. \\ \left. \varphi_{\delta}\langle x_1, \dots, x_k \rangle \wedge \varphi_{\text{halt}}\langle x_2, x_3 \rangle \right), \end{aligned}$$

where

- $\varphi_{\text{start}}$  defines the start configuration  $sc(\varepsilon)$ ,
- $\varphi_{\text{consistent}}$  performs several consistency checks to ensure that the computation of  $M$  on the empty word is faithfully described,
- $\varphi_{\delta}$  encodes the transition function of  $M$ , and
- $\varphi_{\text{halt}}$  checks whether  $M$  halts.

The definitions of the subformulae of  $\varphi_{\text{comp}}$  follow.

**The Formula  $\varphi_{\text{start}}$ .** It is defined as the conjunction of the following FO formulae, expressing that the first tape cell contains the left marker

$$\forall x (\text{First}(x) \rightarrow \text{Symbol}_{\triangleright}(x, x)),$$

the rest of tape cells contain the blank symbol

$$\forall x \forall y ((\text{First}(x) \wedge \neg \text{First}(y)) \rightarrow \text{Symbol}_{\sqcup}(x, y)),$$

the head points to the first cell

$$\forall x (\text{First}(x) \rightarrow \text{Head}(x, x)),$$

and the machine  $M$  is in state  $s$

$$\forall x (\text{First}(x) \rightarrow \text{State}(x, x_1)).$$

Note that we refer to the start state  $s = p_1$  via the variable  $x_1$ .

**The Formula  $\varphi_{\text{consistent}}$ .** It is defined as the conjunction of the following FO formulae, expressing that, at any time instant  $x$ ,  $M$  is in exactly one state

$$\forall x \left( \left( \bigvee_{i=1}^k \text{State}(x, x_i) \right) \wedge \bigwedge_{i,j \in [k] : i < j} \neg (\text{State}(x, x_i) \wedge \text{State}(x, x_j)) \right),$$

each tape cell  $y$  contains exactly one symbol

$$\forall x \forall y \left( \left( \bigvee_{a \in \Sigma} \text{Symbol}_a(x, y) \right) \wedge \bigwedge_{a,b \in \Sigma : a \neq b} \neg (\text{Symbol}_a(x, y) \wedge \text{Symbol}_b(x, y)) \right),$$

and the head points at exactly one cell

$$\forall x \left( \exists y \text{Head}(x, y) \wedge \forall y \forall z \left( (\text{Head}(x, y) \wedge \text{Head}(x, z)) \rightarrow y = z \right) \right).$$

**The Formula  $\varphi_{\delta}$ .** It is defined as the conjunction of the following FO formulae: for each pair  $(p_i, a) \in (Q - \{\text{"yes"}, \text{"no"}\}) \times \Sigma$  with  $\delta(p_i, a) = (p_j, b, \text{dir})$ ,

$$\begin{aligned} & \forall x \forall y \left( (\text{State}(x, x_i) \wedge \text{Head}(x, y) \wedge \text{Symbol}_a(x, y) \wedge \exists t (x \prec t)) \rightarrow \right. \\ & \left. \exists z \exists w \left( \text{Succ}(x, z) \wedge \text{Move}(y, w) \wedge \text{Head}(z, w) \wedge \text{Symbol}_b(z, w) \wedge \text{State}(z, x_j) \wedge \right. \right. \\ & \quad \left. \left. \forall u \left( \neg(y = u) \rightarrow \bigwedge_{c \in \Sigma} (\text{Symbol}_c(x, u) \rightarrow \text{Symbol}_c(z, u)) \right) \right) \right), \end{aligned}$$

where

$$\text{Move}(y, w) = \begin{cases} \text{Succ}(y, w) & \text{if } \text{dir} = \rightarrow, \\ \text{Succ}(w, y) & \text{if } \text{dir} = \leftarrow, \text{ and} \\ y = w & \text{if } \text{dir} = -. \end{cases}$$

**The Formula  $\varphi_{\text{halt}}$ .** Finally, this formula checks whether  $M$  has reached an accepting or a rejecting configuration

$$\exists x (\text{State}(x, x_2) \vee \text{State}(x, x_3)).$$

Recall that, by assumption,  $p_2 = \text{“yes”}$  and  $p_3 = \text{“no”}$ . Thus, the states “yes” and “no” can be accessed via the variables  $x_2$  and  $x_3$ , respectively.

This completes the construction of the Boolean FO query  $q_{\text{comp}}$ , and thus of  $q_M$ . It is not hard to verify that  $M$  halts on the empty word if and only if there exists a database  $D$  such that  $q(D) = \text{true}$ , and the claim follows.  $\square$

The proof of Theorem 8.1 relies on the *finiteness* of databases; it does not work for possibly infinite databases. Assuming that the Turing Machine  $M$  does not halt on the empty word, we can construct an infinite database  $D$  such that  $q_M(D) = \text{true}$  (we leave this as an exercise).<sup>1</sup> As mentioned earlier, Church’s Theorem shows the undecidability of the satisfiability problem for FO queries over possibly infinite databases (see also Exercise 1.11).

We have seen in Chapter 6 that FO and RA have the same expressive power (Theorem 6.1). This fact and Theorem 8.1 immediately imply the following.

### Corollary 8.2

RA-Satisfiability is undecidable.

## Containment and Equivalence

We now focus on the problems of containment and equivalence for FO queries: given two FO queries  $q$  and  $q'$ , is it the case that  $q \subseteq q'$  and  $q \equiv q'$ , respectively. By exploiting Theorem 8.1 it is easy to show the following.

### Theorem 8.3

FO-Containment and FO-Equivalence are undecidable.

<sup>1</sup> The output of an FO query on an infinite database  $D$  is defined in the same way as for databases (see Definition 3.6).

*Proof.* The proof is by an easy reduction from FO-Satisfiability. Consider an FO query  $q$ . From the proof of Theorem 8.1, we know that FO-Satisfiability is undecidable even for Boolean FO queries. Consider the Boolean FO query

$$q' = \exists x (R(x) \wedge \neg R(x)),$$

which is trivially unsatisfiable. It is easy to verify that  $q$  is unsatisfiable if and only if  $q \equiv q'$  (or even  $q \subseteq q'$ ), and the claim follows.  $\square$

The following is an easy consequence of the fact that FO and RA have the same expressive power, and Theorem 8.3

**Corollary 8.4**

RA-Containment and RA-Equivalence are undecidable.

## Homomorphisms

Homomorphisms are a fundamental tool that plays a very prominent role in various aspects of relational databases. We introduce them here, because we will use them in Chapter 10 to reason about functional dependencies. In this chapter, we define homomorphisms and provide some simple examples.

### Definition of Homomorphism

Homomorphisms are structure-preserving functions between two objects of the same type. In our setting, the objects that we are interested in are (possibly infinite) databases and queries. To talk about them as one we define homomorphisms among (possibly infinite) sets of relational atoms. Recall that relational atoms are of the form  $R(\bar{u})$ , where  $\bar{u}$  is a tuple that can mix variables and constants, e.g.,  $R(a, x, 2, b)$ . Recall also that we write  $\text{Dom}(S)$  for the set of constants and variables occurring in a set of relational atoms  $S$ ; for example,  $\text{Dom}(\{R(a, x, b), R(x, a, y)\}) = \{a, b, x, y\}$ .

The way that the notion of homomorphism is defined between sets of atoms is slightly different from the standard notion of mathematical homomorphism, namely constant values of  $\text{Const}$  should be mapped to themselves. The reason for this is that, in general, a value  $a \in \text{Const}$  represents an object different from the one represented by  $b \in \text{Const}$  with  $a \neq b$ , and homomorphisms, as structure preserving functions, should preserve this information as well.

#### Definition 9.1: Homomorphism

Let  $S, S'$  be sets of relational atoms over the same schema. A *homomorphism from  $S$  to  $S'$*  is a function  $h : \text{Dom}(S) \rightarrow \text{Dom}(S')$  such that:

1.  $h(a) = a$  for every  $a \in \text{Dom}(S) \cap \text{Const}$ , and
2. if  $R(\bar{u})$  is an atom in  $S$ , then  $R(h(\bar{u}))$  is an atom in  $S'$ .

If  $h(\bar{u}) = \bar{v}$ , where  $\bar{u}, \bar{v}$  are tuples of the same length over  $\text{Dom}(S)$  and  $\text{Dom}(S')$ , respectively, then  $h$  is a *homomorphism from  $(S, \bar{u})$  to  $(S', \bar{v})$* . We write  $S \rightarrow S'$  if there exists a homomorphism from  $S$  to  $S'$ , and  $(S, \bar{u}) \rightarrow (S', \bar{v})$  if there exists a homomorphism from  $(S, \bar{u})$  to  $(S', \bar{v})$ .

### Example 9.2: Homomorphism

Assume that  $S$  and  $S'$  are sets of relational atoms over the schema  $\{R[2]\}$ . In this way, we can view both  $S$  and  $S'$  as a directed graph: the set of nodes is the set of constants and variables occurring in the relational atoms, and  $R(u, v)$  means that there exists an edge from  $u$  to  $v$ . Unless stated otherwise, the elements in  $S$  and  $S'$  are variables.

**A homomorphism always exists.** Let  $S' = \{R(z, z)\}$ . The function  $h : \text{Dom}(S) \rightarrow \text{Dom}(S')$  such that  $h(x) = z$ , for each  $x \in \text{Dom}(S)$ , is a homomorphism from  $S$  to  $S'$  since  $R(h(x), h(y)) = R(z, z)$  is an atom of  $S'$ , for every  $x, y \in \text{Dom}(S)$ .

**A homomorphism does not exist.** Let  $S = \{R(a, x)\}$  and  $S' = \{R(z, z)\}$ , where  $a \in \text{Const}$ . In contrast to the previous example, there is no homomorphism  $h$  from  $S$  to  $S'$  since, by definition,  $h(a)$  must be equal to  $a$ , while  $a \notin \text{Dom}(S')$ .

**A homomorphism is easy to find.** Let now  $S' = \{R(x, y), R(y, x)\}$ . Assume that a homomorphism  $h$  from  $S$  to  $S'$  exists. As usual,  $h^{-1}$  stands for the inverse, i.e.,  $h^{-1}(x) = \{z \in \text{Dom}(S) \mid h(z) = x\}$ , and likewise for  $h^{-1}(y)$ . The sets  $h^{-1}(x)$  and  $h^{-1}(y)$  are disjoint since  $x \neq y$ . If we have an edge  $(z, w)$  in  $S$ , we know that the variables  $z$  and  $w$  cannot belong to the same set  $h^{-1}(x)$  or  $h^{-1}(y)$ ; otherwise, either  $R(x, x)$  or  $R(y, y)$  would be an atom in  $S$  by the definition of the homomorphism. This means that  $S$ , viewed as a directed graph, is *bipartite*: its nodes are partitioned into two sets such that edges can only connect vertices in different sets. In other words, the nodes of the directed graph given by  $S$  can be colored with two colors  $x$  and  $y$ . Thus, in this case, checking for the existence of a homomorphism witnessing  $S \rightarrow S'$  is the same as checking for the existence of a 2-coloring of  $S$ , which can be done in polynomial time (by using, for example, a coloring version of depth-first search).

**A homomorphism is hard to find.** We now add  $z$  to  $\text{Dom}(S')$ , and let  $S' = \{R(x, y), R(y, x), R(x, z), R(z, x), R(y, z), R(z, y)\}$ . Then, as before, if  $h : \text{Dom}(S) \rightarrow \text{Dom}(S')$  is a homomorphism from  $S$  to  $S'$ , and  $R(z, w)$  is an edge in  $S$ , then  $h(z) \neq h(w)$ . In other words, the nodes of the directed graph given by  $S$  can be colored with three colors  $x, y$  and  $z$ . Therefore, in this case, checking for the existence



of a homomorphism witnessing  $S \rightarrow S'$  is the same as checking for the existence of a 3-coloring of  $S$ , which is an NP-complete problem.

## Grounding Sets of Atoms

In several chapters, it will be convenient to have a mechanism viewing sets of atoms as databases. This is done by converting a set of atoms  $S$  into a possibly infinite database by replacing the variables occurring in  $S$  by new constants not already in  $S$ .<sup>1</sup> This process is called *grounding*, and can be easily defined via homomorphisms.

### Definition 9.3: Grounding

Let  $S$  be a set of relational atoms over a schema  $\mathbf{S}$ . A possibly infinite database  $D$  of  $\mathbf{S}$  is called a *grounding* of  $S$  if there exists a homomorphism from  $S$  to  $D$  that is a bijection.

Note that, in general, there is no unique grounding for a set of atoms. Consider, for example, the set of atoms

$$S = \{R(x, a, y), P(y, b, x, z)\},$$

where  $a, b$  are constants and  $x, y, z$  are variables. The databases

$$D_1 = \{R(c_1, a, d_1), R(d_1, b, c_1, e_1)\} \text{ and } D_2 = \{R(c_2, a, d_2), R(d_2, b, c_2, e_2)\}$$

with  $c_1 \neq c_2$ ,  $d_1 \neq d_2$ , and  $e_1 \neq e_2$ , are both groundings of  $S$ . On the other hand,  $D_1$  and  $D_2$  are isomorphic databases, that is, they are the same up to renaming of constants. This simple observation can be generalized to any set of atoms. In particular, for a set of atoms  $S$ , it is straightforward to show that, for every two groundings  $D_1$  and  $D_2$  of  $S$ , there is a bijection  $\rho : \text{Const} \rightarrow \text{Const}$  such that  $\rho(D_1) = D_2$ . Therefore, we can refer to:

- the grounding of  $S$ , denoted  $S^\downarrow$ , and
- the unique bijective homomorphism  $G_S$  from  $S$  to  $S^\downarrow$ .

We conclude the chapter with a note on the difference between  $\text{Dom}(S)$  and  $\text{Dom}(S^\downarrow)$  to avoid confusion later in the book. If  $S$  is a set of atoms, then  $\text{Dom}(S) \subseteq \text{Const} \cup \text{Var}$ , that is, it may contain *both* constants and variables. On the other hand, by definition,  $\text{Dom}(S^\downarrow)$  contains only constants. Similarly,  $R^S$  is a set of tuples that may mention constants and variables, while  $R^{S^\downarrow}$  is a set of tuples that mention only constants.

<sup>1</sup> Converting a database into a set of atoms by replacing constants with variables is needed less often; this is discussed in Chapter 13.



## Functional Dependencies

In a relational database system, it is possible to specify semantic properties that should be satisfied by all databases of a certain schema, such as “*every person should have at most one social security number*”. Such properties are crucial in the development of transparent and usable database schemas for complex applications, as well as for optimizing the evaluation of queries. However, the relational model as presented in Chapter 2 is not powerful enough to express such semantic properties. This can be achieved by incorporating *integrity constraints*, also known as *dependencies*.

One of the most important classes of dependencies supported by relational systems is the class of *functional dependencies*, which can express that the values of some attributes of a tuple uniquely (or functionally) determine the values of other attributes of that tuple.

### Example 10.1: Functional Dependencies

Consider the (named) database schema

`Person [ pid, name, cid ]`

We can express that the id of a person uniquely determines that person via the functional dependency

$$\text{Person} : \{1\} \rightarrow \{1, 2, 3\},$$

which states that whenever two tuples of the relation Person agree on the first attribute, the id, they should also agree on all the other attributes.

Note that the form of dependency used in Example 10.1, where a set of attributes determines the *entire tuple*, is of particular interest and is called a *key dependency*. We may also say that the id attribute is a *key* of Person.

## Syntax and Semantics

We start with the syntax of functional dependencies.

### Definition 10.2: Syntax of Functional Dependencies

A *functional dependency* (FD)  $\sigma$  over a schema  $\mathbf{S}$  is an expression

$$R : U \rightarrow V$$

where  $R \in \mathbf{S}$  and  $U, V \subseteq \{1, \dots, \text{ar}(R)\}$ . If  $V = \{1, \dots, \text{ar}(R)\}$ , then  $\sigma$  is called a *key dependency*, and we simply write  $\text{key}(R) = U$ .

Intuitively, an FD  $R : U \rightarrow V$  expresses that the values of the attributes  $U$  of  $R$  functionally determine the values of the attributes  $V$  of  $R$ , while a key dependency  $\text{key}(R) = U$  states that the values of the attributes  $U$  of  $R$  functionally determine the values of *all* the attributes of  $R$ . We proceed to formally define the semantics of FDs. Note that in the following definition, by abuse of notation, we write  $U$  and  $V$  in the projection expressions  $\pi_U(\cdot)$  and  $\pi_V(\cdot)$  for the lists consisting of the elements of  $U$  and  $V$  in ascending order.

### Definition 10.3: Semantics of FDs

A database  $D$  of a schema  $\mathbf{S}$  *satisfies* an FD  $\sigma$  of the form  $R : U \rightarrow V$  over  $\mathbf{S}$ , denoted  $D \models \sigma$ , if for each pair of tuples  $\bar{a}, \bar{b} \in R^D$ ,

$$\pi_U(\bar{a}) = \pi_U(\bar{b}) \text{ implies } \pi_V(\bar{a}) = \pi_V(\bar{b}).$$

$D$  *satisfies* a set  $\Sigma$  of FDs, written  $D \models \Sigma$ , if  $D \models \sigma$  for each  $\sigma \in \Sigma$ .

Note that the notion of satisfaction for FDs can be easily transferred to finite sets of atoms by exploiting the notion of grounding of sets of atoms. In particular, a finite set of atoms  $S$  satisfies an FD  $\sigma$ , denoted  $S \models \sigma$ , if  $S^\downarrow \models \sigma$ , while  $S$  satisfies a set  $\Sigma$  of FDs, written  $S \models \Sigma$ , if  $S^\downarrow \models \Sigma$ .

## Satisfaction of Functional Dependencies

A central task is checking whether a database  $D$  satisfies a set  $\Sigma$  of FDs.

### Problem: FD-Satisfaction

**Input:** A database  $D$  of a schema  $\mathbf{S}$ , and a set  $\Sigma$  of FDs over  $\mathbf{S}$

**Output:** true if  $D \models \Sigma$ , and false otherwise

It is not difficult to show the following result:

**Theorem 10.4**

FD-Satisfaction is in PTIME.

*Proof.* Consider a database  $D$  of a schema  $\mathbf{S}$ , and a set  $\Sigma$  of FDs over  $\mathbf{S}$ . Let  $\sigma$  be an FD from  $\Sigma$  of the form  $R : U \rightarrow V$ . To check whether  $D \models \sigma$  we need to check that, for every  $\bar{a}, \bar{b} \in R^D$ ,  $\pi_U(\bar{a}) = \pi_U(\bar{b})$  implies  $\pi_V(\bar{a}) = \pi_V(\bar{b})$ . It is easy to verify that this can be done in time  $O(\|D\|^2)$ . Therefore, we can check whether  $D \models \Sigma$  in time  $O(\|\Sigma\| \cdot \|D\|^2)$ , and the claim follows.  $\square$

**The Chase for Functional Dependencies**

Another crucial task in connection with dependencies is that of (logical) implication, which allows us to discover new dependencies from existing ones. A natural problem that arises in this context is, given a set of dependencies  $\Sigma$  and a dependency  $\sigma$ , to determine whether  $\Sigma$  implies  $\sigma$ . This means checking if, for every database  $D$  such that  $D \models \Sigma$ , it holds that  $D \models \sigma$ . Before formalizing and studying this problem, we first introduce a fundamental algorithmic tool for reasoning about dependencies known as *the chase*. Actually, the chase should be understood as a family of algorithms since, depending on the class of dependencies in question, we may get a different variant. However, all the chase variants have the same objective, that is, given a finite set of relational atoms  $S$ , and a set  $\Sigma$  of dependencies, to transform  $S$  as dictated by  $\Sigma$  into a set of relational atoms that satisfies  $\Sigma$ .

Consider a finite set  $S$  of relational atoms over a schema  $\mathbf{S}$ , and an FD  $\sigma = R : U \rightarrow V$  over  $\mathbf{S}$ . We say that  $\sigma$  is *applicable to  $S$  with  $(\bar{u}, \bar{v})$* , where  $\bar{u}, \bar{v} \in R^S$ <sup>1</sup> if  $\pi_U(\bar{u}) = \pi_U(\bar{v})$  and  $\pi_V(\bar{u}) \neq \pi_V(\bar{v})$ . Let  $\pi_V(\bar{u}) = (u_1, \dots, u_k)$  and  $\pi_V(\bar{v}) = (v_1, \dots, v_k)$ . For technical convenience, we assume that there is a strict total order  $<$  on the elements of the set  $\mathbf{Const} \cup \mathbf{Var}$  such that  $a < x$ , for each  $a \in \mathbf{Const}$  and  $x \in \mathbf{Var}$ , i.e., constants are smaller than variables according to  $<$ . Let  $h_{\bar{u}, \bar{v}} : \text{Dom}(S) \rightarrow \text{Dom}(S)$  be a function such that

$$h_{\bar{u}, \bar{v}}(w) = \begin{cases} u_i & \text{if } w = v_i \text{ and } u_i < v_i, \text{ for some } i \in [k], \\ v_i & \text{if } w = u_i \text{ and } v_i < u_i, \text{ for some } i \in [k], \\ w & \text{otherwise.} \end{cases}$$

The *result of applying  $\sigma$  to  $S$  with  $(\bar{u}, \bar{v})$*  is defined as

$$S' = \begin{cases} \perp & \text{if there is an } i \in [k] \text{ with } u_i \neq v_i \text{ and } u_i, v_i \in \mathbf{Const}, \\ h_{\bar{u}, \bar{v}}(S) & \text{otherwise.} \end{cases}$$

<sup>1</sup> Recall that tuples in  $R^S$  can contain both constants and variables.

Intuitively, the application of  $\sigma$  to  $S$  with  $(\bar{u}, \bar{v})$  fails, indicated by  $\perp$ , whenever we have two distinct constants from **Const** that are supposed to be equal to satisfy  $\sigma$ . In case of non-failure,  $S'$  is obtained from  $S$  by simply replacing  $u_i$  and  $v_i$  by the smallest of the two, for every  $i \in [k]$ . Recall that, by our assumption on  $<$ , if one of  $u_i, v_i$  is a variable and the other one is a constant, then the variable is always replaced by the constant. The application of  $\sigma$  to  $S$  with  $(\bar{u}, \bar{v})$ , which results to  $S'$ , is denoted by  $S \xrightarrow{\sigma, (\bar{u}, \bar{v})} S'$ .

We are now ready to introduce the notion of chase sequence of a finite set  $S$  of relational atoms under a set  $\Sigma$  of FDs, which formalizes the objective of transforming  $S$  as dictated by  $\Sigma$  into a set of atoms that satisfies  $\Sigma$ .

**Definition 10.5: The Chase for FDs**

Consider a finite set  $S$  of relational atoms over a schema **S**, and a set  $\Sigma$  of FDs over **S**.

- A *finite chase sequence* of  $S$  under  $\Sigma$  is a finite sequence  $s = S_0, \dots, S_n$  of sets of relational atoms, where  $S_0 = S$ , and
  - for each  $i \in [0, n-1]$ , there is an FD  $\sigma = R : U \rightarrow V$  in  $\Sigma$  and atoms  $R(\bar{u}), R(\bar{v}) \in S_i$  such that  $S_i \xrightarrow{\sigma, (\bar{u}, \bar{v})} S_{i+1}$ , and
  - either  $S_n = \perp$ , in which case we say that  $s$  is *failing*, or, for every FD  $\sigma = R : U \rightarrow V$  in  $\Sigma$  and atoms  $R(\bar{u}), R(\bar{v}) \in S_n$ ,  $\sigma$  is not applicable to  $S_n$  with  $(\bar{u}, \bar{v})$ , in which case  $s$  is called *successful*.
- An *infinite chase sequence* of  $S$  under  $\Sigma$  is an infinite sequence  $S_0, S_1, \dots$  of sets of relational atoms, where  $S_0 = S$ , and for each  $i \geq 0$ , there is an FD  $\sigma = R : U \rightarrow V$  in  $\Sigma$  and atoms  $R(\bar{u}), R(\bar{v}) \in S_i$  such that  $S_i \xrightarrow{\sigma, (\bar{u}, \bar{v})} S_{i+1}$ .

We proceed to present some fundamental properties of the chase for FDs.<sup>[2]</sup> In what follows, let  $S$  be a finite set of relational atoms, and  $\Sigma$  a finite set of FDs, both over the same schema **S**. It is not hard to see that there are no infinite chase sequences under FDs.<sup>[3]</sup> This is a consequence of the fact that each non-failing chase application does not introduce new terms but only equalizes them. Therefore, in the worst-case, the chase either will fail, or will produce after finitely many steps a set of relational atoms with only one term, which trivially satisfies every functional dependency.

**Lemma 10.6.** *There is no infinite chase sequence of  $S$  under  $\Sigma$ .*

<sup>2</sup> Formal proofs are omitted since in Chapter 52 we are going to present the chase for a more general class of dependencies than FDs, known as equality-generating dependencies, and provide proofs there for all the desired properties.

<sup>3</sup> As we discuss in Chapter 11, this is not the case for other types of dependencies, in particular, inclusion dependencies.

Although there could be several finite chase sequences of  $S$  under  $\Sigma$ , depending on the application order of the FDs in  $\Sigma$ , we can show that all those sequences either fail or end in exactly the same set of relational atoms.

**Lemma 10.7.** *Let  $S_0, \dots, S_n$  and  $S'_0, \dots, S'_m$  be two finite chase sequences of  $S$  under  $\Sigma$ . Then it holds that  $S_n = S'_m$ .*

The above lemma allows us to refer to *the* result of the chase of  $S$  under  $\Sigma$ , denoted by  $\text{Chase}(S, \Sigma)$ , which is defined as  $S_n$  for some (any) finite chase sequence  $S_0, \dots, S_n$  of  $S$  under  $\Sigma$ . Notice that we do not need to define the result of infinite chase sequences under FDs since, by Lemma 10.6, they do not exist. Hence,  $\text{Chase}(S, \Sigma)$  is either the symbol  $\perp$ , or a finite set of relational atoms. It is not difficult to verify that in the latter case,  $\text{Chase}(S, \Sigma)$  satisfies  $\Sigma$ . Actually, this follows from the definition of successful chase sequences.

**Lemma 10.8.** *If  $\text{Chase}(S, \Sigma) \neq \perp$ , then  $\text{Chase}(S, \Sigma) \models \Sigma$ .*

A central notion is that of chase homomorphism, which essentially computes the result of a successful finite chase sequence of  $S$  under  $\Sigma$ . Consider such a chase sequence  $s = S_0, S_1, \dots, S_n$  of  $S$  under  $\Sigma$  such that

$$S_0 \xrightarrow{\sigma_0, (\bar{u}_0, \bar{v}_0)} S_1 \xrightarrow{\sigma_1, (\bar{u}_1, \bar{v}_1)} S_2 \cdots S_{n-1} \xrightarrow{\sigma_{n-1}, (\bar{u}_{n-1}, \bar{v}_{n-1})} S_n.$$

Recall that  $S_i = h_{\bar{u}_{i-1}, \bar{v}_{i-1}}(S_{i-1})$ , for each  $i \in [n]$ . The *chase homomorphism* of  $s$ , denoted  $h_s$ , is defined as the composition of functions

$$h_{\bar{u}_{n-1}, \bar{v}_{n-1}} \circ h_{\bar{u}_{n-2}, \bar{v}_{n-2}} \circ \cdots \circ h_{\bar{u}_0, \bar{v}_0}.$$

It is clear that  $h_s(S_0) = h_s(S) = S_n$ . Since, by Lemma 10.7, different finite chase sequences have the same result, we get the following.

**Lemma 10.9.** *Let  $s$  and  $s'$  be successful finite chase sequences of  $S$  under  $\Sigma$ . It holds that  $h_s(S) = h_{s'}(S)$ .*

Therefore, assuming that  $\text{Chase}(S, \Sigma) \neq \perp$ , we can refer to *the* chase homomorphism of  $S$  under  $\Sigma$ , denoted  $h_{S, \Sigma}$ . It should be clear that  $\text{Chase}(S, \Sigma) \neq \perp$  implies  $h_{S, \Sigma}(S) = \text{Chase}(S, \Sigma)$ .

By Lemma 10.6,  $\text{Chase}(S, \Sigma)$  can be computed after finitely many steps. Furthermore, assuming that  $\text{Chase}(S, \Sigma) \neq \perp$ , also the chase homomorphism  $h_{S, \Sigma}$  can be computed after finitely many steps. In fact, as the next lemma states, this is even possible after polynomially many steps.

**Lemma 10.10.**  *$\text{Chase}(S, \Sigma)$  can be computed in polynomial time. Furthermore, if  $\text{Chase}(S, \Sigma) \neq \perp$ , then  $h_{S, \Sigma}$  can be computed in polynomial time.*

The last main property of the chase states that, if  $\text{Chase}(S, \Sigma) \neq \perp$ , then it acts as a representative of all the sets of atoms  $S'$  that satisfy  $\Sigma$  and  $S \rightarrow S'$ , that is, there exists a homomorphism from  $S$  to  $S'$ .

**Lemma 10.11.** *Let  $S'$  be a set of atoms over  $\mathbf{S}$  such that  $(S, \bar{u}) \rightarrow (S', \bar{v})$  and  $S' \models \Sigma$ . If  $\text{Chase}(S, \Sigma) \neq \perp$ , then  $(\text{Chase}(S, \Sigma), h_{S, \Sigma}(\bar{u})) \rightarrow (S', \bar{v})$ .*

Note that the definition of the chase for FDs, as well as its main properties, would be technically simpler if we focus on sets of constant-free atoms since in this case there are no failing chase sequences. As we shall see, this suffices for studying the implication problem for FDs. Nevertheless, we consider sets of atoms with constants since the chase is also used in Chapter 17 for studying a different problem for which the proper treatment of constants is vital.

## Implication of Functional Dependencies

We now proceed to study the implication problem for FDs, which we define next. Given a set  $\Sigma$  of FDs over a schema  $\mathbf{S}$  and a single FD  $\sigma$  over  $\mathbf{S}$ , we say that  $\Sigma$  *implies*  $\sigma$ , denoted  $\Sigma \models \sigma$ , if, for every database  $D$  of  $\mathbf{S}$ , we have that  $D \models \Sigma$  implies  $D \models \sigma$ . The main problem of concern is the following:

### Problem: FD-Implication

**Input:** A set  $\Sigma$  of FDs over a schema  $\mathbf{S}$ , and an FD  $\sigma$  over  $\mathbf{S}$   
**Output:** true if  $\Sigma \models \sigma$ , and false otherwise

We proceed to show the following result:

### Theorem 10.12

FD-Implication is in PTIME.

To show Theorem 10.12, we first show how implication of FDs can be characterized via the chase for FDs. This is done by showing that checking whether a set of FDs  $\Sigma$  implies an FD  $\sigma$  boils down to checking whether the result of the chase of the prototypical set of relational atoms  $S_\sigma$  that violates  $\sigma$  is a set of atoms that satisfies  $\sigma$ . Given an FD  $\sigma$  of the form  $R : U \rightarrow V$ , the set  $S_\sigma$  is defined as  $\{R(x_1, \dots, x_{\text{ar}(R)}), R(y_1, \dots, y_{\text{ar}(R)})\}$ , where

- $x_1, \dots, x_{\text{ar}(R)}, y_1, \dots, y_{\text{ar}(R)}$  are variables,
- for each  $i, j \in \{1, \dots, \text{ar}(R)\}$  with  $i \neq j$ , it holds that  $x_i \neq x_j$  and  $y_i \neq y_j$ , and
- for each  $i \in \{1, \dots, \text{ar}(R)\}$ , we have  $x_i = y_i$  if and only if  $i \in U$ .

We can now show the following useful characterization:



**Proposition 10.13**

Consider a set  $\Sigma$  of FDs over a schema  $\mathbf{S}$ , and an FD  $\sigma$  over  $\mathbf{S}$ . Then:

$$\Sigma \models \sigma \quad \text{if and only if} \quad \text{Chase}(S_\sigma, \Sigma) \models \sigma.$$

*Proof.* ( $\Rightarrow$ ) By hypothesis, for every finite set of relational atoms  $S$ , it holds that  $S \models \Sigma$  implies  $S \models \sigma$ . Observe that  $\text{Chase}(S_\sigma, \Sigma) \neq \perp$  since  $S_\sigma$  contains only variables. Therefore, by Lemma 10.8, we have that  $\text{Chase}(S_\sigma, \Sigma) \models \Sigma$ . Since, by Lemma 10.6,  $\text{Chase}(S_\sigma, \Sigma)$  is finite, we get that  $\text{Chase}(S_\sigma, \Sigma) \models \sigma$ .

( $\Leftarrow$ ) Consider now a database  $D$  of  $\mathbf{S}$  such that  $D \models \Sigma$ , and with  $\sigma$  being of the form  $R : \{i_1, \dots, i_k\} \rightarrow \{j_1, \dots, j_\ell\}$ , assume that there are tuples  $(a_1, \dots, a_{\text{ar}(R)}), (b_1, \dots, b_{\text{ar}(R)}) \in R^D$  such that  $(a_{i_1}, \dots, a_{i_k}) = (b_{i_1}, \dots, b_{i_k})$ . Recall also that  $S_\sigma$  is of the form  $\{R(x_1, \dots, x_{\text{ar}(R)}), R(y_1, \dots, y_{\text{ar}(R)})\}$ . Let  $\bar{z} = (x_{j_1}, \dots, x_{j_\ell}, y_{j_1}, \dots, y_{j_\ell})$  and  $\bar{c} = (a_{j_1}, \dots, a_{j_\ell}, b_{j_1}, \dots, b_{j_\ell})$ . It is clear that  $(S_\sigma, \bar{z}) \rightarrow (D, \bar{c})$ . Since  $D \models \Sigma$  and  $\text{Chase}(S_\sigma, \Sigma) \neq \perp$ , by Lemma 10.11

$$(\text{Chase}(S_\sigma, \Sigma), h_{S_\sigma, \Sigma}(\bar{z})) \rightarrow (D, \bar{c}).$$

Since, by hypothesis,  $\text{Chase}(S_\sigma, \Sigma) \models \sigma$ , we can conclude that

$$(h_{S_\sigma, \Sigma}(x_{j_1}), \dots, h_{S_\sigma, \Sigma}(x_{j_\ell})) = (h_{S_\sigma, \Sigma}(y_{j_1}), \dots, h_{S_\sigma, \Sigma}(y_{j_\ell})),$$

which in turn implies that

$$(a_{j_1}, \dots, a_{j_\ell}) = (b_{j_1}, \dots, b_{j_\ell}).$$

Therefore,  $D \models \sigma$ , and the claim follows.  $\square$

By Proposition 10.13, we get a simple procedure for checking whether a set  $\Sigma$  of FDs implies an FD  $\sigma$  that runs in polynomial time:

if  $\text{Chase}(S_\sigma, \Sigma) \models \sigma$ , then return **true**; otherwise, return **false**.

We know that the set of atoms  $\text{Chase}(S_\sigma, \Sigma)$  can be constructed in polynomial time (Lemma 10.10), and we also know that  $\text{Chase}(S_\sigma, \Sigma) \models \sigma$  can be checked in polynomial time (Theorem 10.4), and Theorem 10.12 follows.



## Inclusion Dependencies

In this chapter, we concentrate on another central class of constraints supported by relational database systems, called *inclusion dependencies* (also known as *referential constraints*). With this type of constraints we can express relationships among attributes of different relations, which is not possible using functional dependencies that can talk only about one relation.

### Example 11.1: Inclusion Dependencies

Having the (named) database schema

```
Person [ pid, pname, cid ]  
Profession [ pid, pname ]
```

we would like to express that the values occurring in the first attribute of Profession are person ids. This can be done via the inclusion dependency

$$\text{Profession}[1] \subseteq \text{Person}[1].$$

This dependency simply states that the set of values occurring in the first attribute of the relation Profession should be a subset of the set of values appearing in the first attribute of the relation Person.

## Syntax and Semantics

We start with the syntax of inclusion dependencies.

### Definition 11.2: Syntax of Inclusion Dependencies

An *inclusion dependency* (IND)  $\sigma$  over a schema  $\mathbf{S}$  is an expression

$$R[i_1, \dots, i_k] \subseteq P[j_1, \dots, j_k]$$

where  $k \geq 1$ ,  $R, P$  belong to  $\mathbf{S}$ , and  $(i_1, \dots, i_k)$  and  $(j_1, \dots, j_k)$  are lists of distinct integers from  $\{1, \dots, \text{ar}(R)\}$  and  $\{1, \dots, \text{ar}(P)\}$ , respectively.

Intuitively, an IND  $R[i_1, \dots, i_k] \subseteq P[j_1, \dots, j_k]$  states that if  $R(\bar{a})$  belongs to a database  $D$ , then in the same database an atom  $P(\bar{b})$  should exist such that the  $i_\ell$ -th element of  $\bar{a}$  coincides with the  $j_\ell$ -th element of  $\bar{b}$ , for  $\ell \in [k]$ . The formal definition of the semantic meaning of INDs follows.

### Definition 11.3: Semantics of INDs

A database  $D$  of a schema  $\mathbf{S}$  *satisfies* an IND  $\sigma$  of the form  $R[i_1, \dots, i_k] \subseteq P[j_1, \dots, j_k]$  over  $\mathbf{S}$ , denoted  $D \models \sigma$ , if for every tuple  $\bar{a} \in R^D$ , there exists a tuple  $\bar{b} \in P^D$  such that

$$\pi_{(i_1, \dots, i_k)}(\bar{a}) = \pi_{(j_1, \dots, j_k)}(\bar{b}).$$

$D$  *satisfies* a set  $\Sigma$  of INDs, denoted  $D \models \Sigma$ , if  $D \models \sigma$  for each  $\sigma \in \Sigma$ .

## Satisfaction of Inclusion Dependencies

A central task is checking whether a database  $D$  satisfies a set  $\Sigma$  of INDs.

### Problem: IND-Satisfaction

**Input:** A database  $D$  over a schema  $\mathbf{S}$ , and a set  $\Sigma$  of INDs over  $\mathbf{S}$

**Output:** **true** if  $D \models \Sigma$ , and **false** otherwise

It is not difficult to show the following result:

### Theorem 11.4

IND-Satisfaction is in PTIME.

*Proof.* Consider a database  $D$  of a schema  $\mathbf{S}$ , and a set  $\Sigma$  of INDs over  $\mathbf{S}$ . Let  $\sigma$  be an IND from  $\Sigma$  of the form  $R[i_1, \dots, i_k] \subseteq P[j_1, \dots, j_k]$ . To check whether  $D \models \sigma$  we need to check that, for every tuple  $(a_1, \dots, a_{\text{ar}(R)}) \in R^D$ , there exists a tuple  $(b_1, \dots, b_{\text{ar}(P)}) \in P^D$  such that  $(a_{i_1}, \dots, a_{i_k}) = (b_{j_1}, \dots, b_{j_k})$ . It is not difficult to verify that this can be done in time  $O(\|D\|^2)$ . Therefore, we can check whether  $D \models \Sigma$  in time  $O(\|\Sigma\| \cdot \|D\|^2)$ , and the claim follows.  $\square$

## The Chase for Inclusion Dependencies

As for FDs, the other crucial task of interest in connection with INDs is (logical) implication. Unsurprisingly, the main tool for studying the implication problem for INDs is the chase for INDs, which we introduce next.

Consider a finite set  $S$  of atoms over  $\mathbf{S}$ , and an IND  $\sigma = R[i_1, \dots, i_m] \subseteq P[j_1, \dots, j_m]$  over  $\mathbf{S}$ . We say that  $\sigma$  is *applicable to  $S$  with  $\bar{u} = (u_1, \dots, u_{\text{ar}(R)})$*  if  $\bar{u} \in R^S$ . Let  $\text{new}(\sigma, \bar{u}) = P(v_1, \dots, v_{\text{ar}(P)})$ , where, for each  $\ell \in [\text{ar}(P)]$ ,

$$v_\ell = \begin{cases} u_{i_k} & \text{if } \ell = j_k, \text{ for } k \in [m], \\ x_\ell^{\sigma, \pi(i_1, \dots, i_m)(\bar{u})} & \text{otherwise,} \end{cases}$$

with  $x_\ell^{\sigma, \pi(i_1, \dots, i_m)(\bar{u})} \in \text{Var} - \text{Dom}(S)$ .<sup>1</sup> The *result of applying  $\sigma$  to  $S$  with  $\bar{u}$*  is the set of atoms  $S' = S \cup \{\text{new}(\sigma, \bar{u})\}$ . In simple words,  $S'$  is obtained from  $S$  by adding the new atom  $\text{new}(\sigma, \bar{u})$ , which is uniquely determined by  $\sigma$  and  $\bar{u}$ . The application of  $\sigma$  to  $S$  with  $\bar{u}$ , which results in  $S'$ , is denoted  $S \xrightarrow{\sigma, \bar{u}} S'$ .

We are now ready to introduce the notion of chase sequence of a finite set  $S$  of relational atoms under a set  $\Sigma$  of INDs, which formalizes the objective of transforming  $S$  as dictated by  $\Sigma$  into a set of atoms that satisfies  $\Sigma$ .

### Definition 11.5: The Chase for INDs

Consider a finite set  $S$  of relational atoms over a schema  $\mathbf{S}$ , and a set  $\Sigma$  of INDs over  $\mathbf{S}$ .

- A *finite chase sequence* of  $S$  under  $\Sigma$  is a finite sequence  $s = S_0, \dots, S_n$  of sets of relational atoms, where  $S_0 = S$ , and
  1. for each  $i \in [0, n-1]$ , there is  $\sigma = R[\alpha] \subseteq P[\beta]$  in  $\Sigma$  and  $\bar{u} \in R^{S_i}$  such that  $\text{new}(\sigma, \bar{u}) \notin S_i$  and  $S_i \xrightarrow{\sigma, \bar{u}} S_{i+1}$ , and
  2. for each IND  $\sigma = R[\alpha] \subseteq P[\beta]$  in  $\Sigma$  and  $\bar{u} \in R^{S_n}$ ,  $\text{new}(\sigma, \bar{u}) \in S_n$ .

The *result* of  $s$  is defined as the set of atoms  $S_n$ .

- An *infinite chase sequence* of  $S$  under  $\Sigma$  is an infinite sequence  $s = S_0, S_1, \dots$  of sets of atoms, where  $S_0 = S$ , and
  1. for each  $i \geq 0$ , there is  $\sigma = R[\alpha] \subseteq P[\beta]$  in  $\Sigma$  and  $\bar{u} \in R^{S_i}$  such that  $\text{new}(\sigma, \bar{u}) \notin S_i$  and  $S_i \xrightarrow{\sigma, \bar{u}} S_{i+1}$ , and
  2. for each  $i \geq 0$ , and for each  $\sigma = R[\alpha] \subseteq P[\beta]$  in  $\Sigma$  and  $\bar{u} \in R^{S_i}$  such that  $\sigma$  is applicable to  $S_i$  with  $\bar{u}$ , there exists  $j > i$  such that  $\text{new}(\sigma, \bar{u}) \in S_j$ .

<sup>1</sup> One could adopt a simpler naming scheme for these newly introduced variables. For example, for each  $\ell \in [\text{ar}(P)] - \{j_1, \dots, j_m\}$ , we could simply name the new variable  $x_\ell^{\sigma, \bar{u}}$ . For further details on this matter see the comments for Part II.

The *result* of  $s$  is defined as the set infinite set of atoms  $\bigcup_{i \geq 0} S_i$ .

In the case of finite chase sequences, the first condition in Definition 11.5 simply says that  $S_{i+1}$  is obtained from  $S_i$  by applying  $\sigma$  to  $S_i$  with  $\bar{u}$ , while  $\sigma$  has not been already applied to some  $S_j$ , for  $j < i$ , with  $\bar{u}$ . The second condition states that no new atom, which is not already in  $S_n$ , can be derived by applying an IND of  $\Sigma$  to  $S_n$ . Now, in the case of infinite chase sequences, the first condition in Definition 11.5, as in the finite case, says that  $S_{i+1}$  is obtained from  $S_i$  by applying  $\sigma$  to  $S_i$  with  $\bar{u}$ , while  $\sigma$  has not been already applied before. The second condition is known as the *fairness condition*, and it ensures that all the INDs that are applicable eventually will be applied.

We proceed to show some fundamental properties of the chase for INDs.<sup>2</sup> In what follows, let  $S$  be a finite set of relational atoms, and  $\Sigma$  a finite set of INDs, both over the same schema  $\mathbf{S}$ . Recall that in the case of FDs we know that there are no infinite chase sequences since a chase application does not introduce new terms, but only equalizes terms. However, in the case of INDs, a chase step may introduce new variables not occurring in the given set of atoms, which may lead to infinite chase sequences. Indeed, this can happen even for simple sets of atoms and INDs. For example, it is not hard to verify that the single chase sequence of  $\{R(a, b)\}$  under  $\{R[2] \subseteq R[1]\}$  is infinite.

Although we may have infinite chase sequences, we can still establish some favourable properties. It is clear that there are several chase sequences of  $S$  under  $\Sigma$  depending on the order that we apply the INDs of  $\Sigma$ . However, the adopted naming scheme of new variables ensures that, no matter when we apply an IND  $\sigma$  with a tuple  $\bar{u}$ , the newly generated atom  $\text{new}(\sigma, \bar{u})$  is always the same, which in turn allows us to show that all those chase sequences have the same result. At this point, let us stress that the result of an infinite chase sequence  $s = S_0, S_1, \dots$  of  $S$  under  $\Sigma$  always exists.<sup>3</sup> This can be shown by exploiting classical results of fixpoint theory. By using Kleene's Theorem, we can show that  $\bigcup_{i \geq 0} S_i$  coincides with the least fixpoint of a continuous operator (which corresponds to a single chase step) on the complete lattice  $(\text{Inst}(\mathbf{S}), \subseteq)$ , which we know that always exists by Knaster-Tarski's Theorem (we leave the proof as an exercise). We can now state the announced result.

**Lemma 11.6.** *The following hold:*

1. *There exists a finite chase sequence of  $S$  under  $\Sigma$  if and only if there is no infinite chase sequence of  $S$  under  $\Sigma$ .*
2. *Let  $S_0, \dots, S_n$  and  $S'_0, \dots, S'_m$  be two finite chase sequences of  $S$  under  $\Sigma$ . Then, it holds that  $S_n = S'_m$ .*

<sup>2</sup> Formal proofs are omitted since in Chapter 43 we are going to present the chase for a more general class of dependencies than INDs, known as tuple-generating dependencies, and provide proofs there for all the desired properties.

<sup>3</sup> This statement trivially holds for finite chase sequences.

3. Let  $S_0, S_1, \dots$  and  $S'_0, S'_1, \dots$  be two infinite chase sequences of  $S$  under  $\Sigma$ . Then, it holds that  $\bigcup_{i \geq 0} S_i = \bigcup_{i \geq 0} S'_i$ .

Lemma 11.6 allows us to refer to *the* unique result of the chase of  $S$  under  $\Sigma$ , denoted  $\text{Chase}(S, \Sigma)$ , which is defined as the result of some (any) finite or infinite chase sequence of  $S$  under  $\Sigma$ . At this point, the reader may expect that the next key property is that  $\text{Chase}(S, \Sigma)$  satisfies  $\Sigma$ . However, it should not be overlooked that  $\text{Chase}(S, \Sigma)$  is a possibly infinite set of atoms, and thus, we cannot directly apply the notion of satisfaction from Definition 11.3. Nevertheless, Definition 11.3 can be readily applied to possibly infinite databases, which in turn allows us to transfer the notion of satisfaction for INDs to sets of atoms via the notion of grounding. In particular, a set of atoms  $S$  satisfies an IND  $\sigma$ , denoted  $S \models \sigma$ , if  $S^\downarrow \models \sigma$ , while  $S$  satisfies a set  $\Sigma$  of INDs, written  $S \models \Sigma$ , if  $S^\downarrow \models \Sigma$ . We can now formally state that  $\text{Chase}(S, \Sigma)$  satisfies  $\Sigma$ . Let us clarify, though, that in the case where only infinite chase sequences exist, this result heavily relies on the fairness condition.

**Lemma 11.7.** *It holds that  $\text{Chase}(S, \Sigma) \models \Sigma$ .*

The last crucial property states that  $\text{Chase}(S, \Sigma)$  acts as a representative of all the finite or infinite sets of atoms  $S'$  that satisfy  $\Sigma$ , and such that there exists a homomorphism from  $S$  to  $S'$ , that is,  $S \rightarrow S'$ .

**Lemma 11.8.** *Let  $S'$  be a set of atoms over  $\mathbf{S}$  such that  $(S, \bar{u}) \rightarrow (S', \bar{v})$  and  $S' \models \Sigma$ . It holds that  $(\text{Chase}(S, \Sigma), \bar{u}) \rightarrow (S', \bar{v})$ .*

## Implication of Inclusion Dependencies

We now proceed to study the implication problem for INDs. The notion of implication for INDs is defined in the same way as for functional dependencies. More precisely, given a set  $\Sigma$  of INDs over a schema  $\mathbf{S}$  and a single IND  $\sigma$  over  $\mathbf{S}$ , we say that  $\Sigma$  *implies*  $\sigma$ , denoted  $\Sigma \models \sigma$ , if, for every database  $D$  of  $\mathbf{S}$ , we have that  $D \models \Sigma$  implies  $D \models \sigma$ . This leads to the following problem:

### Problem: IND-Implication

**Input:** A set  $\Sigma$  of INDs over a schema  $\mathbf{S}$ , and an IND  $\sigma$  over  $\mathbf{S}$

**Output:** true if  $\Sigma \models \sigma$ , and false otherwise

Although for FDs the implication problem is tractable (Theorem 10.12), for INDs it turns out to be an intractable problem:

**Theorem 11.9**

IND-Implication is PSPACE-complete.

We first concentrate on the upper bound. We are going to establish a result, analogous to Proposition 10.13 for FDs, that characterizes implication of INDs via the chase. However, since the chase for INDs may build an infinite set of atoms, we can only characterize implication under possibly infinite databases. Given a set  $\Sigma$  of INDs over a schema  $\mathbf{S}$  and a single IND  $\sigma$  over  $\mathbf{S}$ , we say that  $\Sigma$  *implies without restriction*  $\sigma$ , denoted  $\Sigma \models^\infty \sigma$ , if, for every possibly infinite database  $D$  of  $\mathbf{S}$ , we have that  $D \models \Sigma$  implies  $D \models \sigma$ .

Given an IND  $\sigma$  of the form  $R[i_1, \dots, i_k] \subseteq P[j_1, \dots, j_k]$ , the set  $S_\sigma$  is defined as the singleton  $\{R(x_1, \dots, x_{\text{ar}(R)})\}$ , where  $x_1, \dots, x_{\text{ar}(R)}$  are distinct variables. We can now show the following auxiliary lemma.

**Lemma 11.10.** *Consider a set  $\Sigma$  of INDs over schema  $\mathbf{S}$ , and an IND  $\sigma$  over  $\mathbf{S}$ . It holds that  $\Sigma \models^\infty \sigma$  if and only if  $\text{Chase}(S_\sigma, \Sigma) \models \sigma$ .*

*Proof.* ( $\Rightarrow$ ) By hypothesis, for every possibly infinite set of relational atoms  $S$ , it holds that  $S \models \Sigma$  implies  $S \models \sigma$ . By Lemma 11.7,  $\text{Chase}(S_\sigma, \Sigma) \models \Sigma$ , and therefore,  $\text{Chase}(S_\sigma, \Sigma) \models \sigma$ .

( $\Leftarrow$ ) Consider now a possibly infinite database  $D$  of  $\mathbf{S}$  such that  $D \models \Sigma$ , and with  $\sigma$  being of the form  $R[i_1, \dots, i_k] \subseteq P[j_1, \dots, j_k]$ , assume that there exists a tuple  $(a_1, \dots, a_{\text{ar}(R)}) \in R^D$ . Recall also that  $S_\sigma$  is of the form  $\{R(x_1, \dots, x_{\text{ar}(R)})\}$ . Let  $\bar{y} = (x_{i_1}, \dots, x_{i_k})$  and  $\bar{b} = (a_{i_1}, \dots, a_{i_k})$ . It is clear that  $(S_\sigma, \bar{y}) \rightarrow (D, \bar{b})$ . Since  $D \models \Sigma$ , by Lemma 11.8

$$(\text{Chase}(S_\sigma, \Sigma), \bar{y}) \rightarrow (D, \bar{b}).$$

Since, by hypothesis,  $\text{Chase}(S_\sigma, \Sigma) \models \sigma$ , we can conclude that there exists a tuple  $(z_1, \dots, z_{\text{ar}(P)}) \in P^{\text{Chase}(S_\sigma, \Sigma)}$  such that

$$(x_{i_1}, \dots, x_{i_k}) = (z_{j_1}, \dots, z_{j_k}),$$

which in turn implies that there exists  $(c_1, \dots, c_{\text{ar}(P)}) \in P^D$  such that

$$(a_{i_1}, \dots, a_{i_k}) = (c_{j_1}, \dots, c_{j_k}).$$

Therefore,  $D \models \sigma$ , and the claim follows.  $\square$

Lemma 11.10 alone is of little use since it characterizes implication of INDs under possibly infinite databases, whereas we are interested only in (finite) databases. However, we can show that implication of INDs is *finitely controllable*, which means that implication under finite databases ( $\models$ ) and implication under possibly infinite databases ( $\models^\infty$ ) coincide.



**Theorem 11.11: Finite Controllability of Implication**

Consider a set  $\Sigma$  of INDs over a schema  $\mathbf{S}$ , and an IND  $\sigma$  over  $\mathbf{S}$ . Then:

$$\Sigma \models \sigma \quad \text{if and only if} \quad \Sigma \models^\infty \sigma.$$

Although the above theorem is crucial for our analysis, we do not discuss its proof here (see Exercise 1.15). An immediate consequence of Lemma 11.10 and Theorem 11.11 is the following:

**Corollary 11.12**

Consider a set  $\Sigma$  of INDs over a schema  $\mathbf{S}$ , and an IND  $\sigma$  over  $\mathbf{S}$ . Then:

$$\Sigma \models \sigma \quad \text{if and only if} \quad \text{Chase}(S_\sigma, \Sigma) \models \sigma.$$

Due to Corollary 11.12 the reader may think that the procedure for checking whether  $\Sigma \models \sigma$ , which will lead to the PSPACE upper bound claimed in Theorem 11.9 is simply to construct the set of atoms  $\text{Chase}(S_\sigma, \Sigma)$ , and then check whether it satisfies  $\sigma$ , which can be achieved due to Theorem 11.4. However, it should not be forgotten that  $\text{Chase}(S_\sigma, \Sigma)$  may be infinite. Therefore, we need to rely on a finer procedure that avoids the explicit construction of  $\text{Chase}(S_\sigma, \Sigma)$ . We proceed to present a technical lemma that is the building block of this refined procedure, but first we need some terminology.

Given an IND  $\sigma = R[i_1, \dots, i_m] \subseteq P[j_1, \dots, j_m]$  and a tuple of variables  $\bar{x} = (x_1, \dots, x_{\text{ar}(R)})$ , we define the atom  $\text{new}^*(\sigma, \bar{x})$  as the atom obtained from  $\text{new}(\sigma, \bar{x})$  after replacing the newly introduced variables with the special variable  $\star \notin \{x_1, \dots, x_{\text{ar}(R)}\}$ , which should be understood as a placeholder for new variables. Formally,  $\text{new}^*(\sigma, \bar{x}) = P(y_1, \dots, y_{\text{ar}(P)})$ , where, for  $\ell \in [\text{ar}(P)]$ ,

$$y_\ell = \begin{cases} x_{i_k} & \text{if } \ell = j_k, \text{ for } k \in [m], \\ \star & \text{otherwise.} \end{cases}$$

Given a set  $\Sigma$  of INDs, a *witness of  $\sigma$  relative to  $\Sigma$*  is a sequence of atoms  $R_1(\bar{x}_1), \dots, R_n(\bar{x}_n)$ , for  $n \geq 1$ , such that:

- $S_\sigma = \{R_1(\bar{x}_1)\}$ ,
- for each  $i \in [2, n]$ , there is  $\sigma_i = R_{i-1}[\alpha_{i-1}] \subseteq R_i[\alpha_i]$  in  $\Sigma$  that is applicable to  $\{R_{i-1}(\bar{x}_{i-1})\}$  with  $\bar{x}_{i-1}$  such that  $R_i(\bar{x}_i) = \text{new}^*(\sigma_i, \bar{x}_{i-1})$ ,
- $R_n = P$ , and
- $\pi_{(i_1, \dots, i_m)}(\bar{x}_1) = \pi_{(j_1, \dots, j_m)}(\bar{x}_n)$ .

A witness of  $\sigma$  relative to  $\Sigma$  is essentially a compact representation, which uses only  $\text{ar}(R) + 1$  variables, of a sequence of atoms of  $\text{Chase}(S_\sigma, \Sigma)$  that

**Algorithm 2** IMPLICATIONWITNESS( $\Sigma, \sigma$ )**Input:** A set  $\Sigma$  of INDs over  $\mathbf{S}$  and  $\sigma = R[i_1, \dots, i_k] \subseteq P[j_1, \dots, j_k]$  over  $\mathbf{S}$ .**Output:** **true** if there is a witness of  $\sigma$  relative  $\Sigma$ , and **false** otherwise.

---

```

1: if  $R = P$  and  $(i_1, \dots, i_k) = (j_1, \dots, j_k)$  then
2:   return true
3:  $S_\nabla := \{R(\bar{x})\}$ , where  $\bar{x} = (x_1, \dots, x_{\text{ar}(R)})$  consists of distinct variables
4:  $S_\triangleright := \emptyset$ 
5: repeat
6:   if  $\sigma' = T[\alpha] \subseteq T'[\beta] \in \Sigma$  is applicable to  $S_\nabla$  with  $\bar{y} \in \text{Dom}(S_\nabla)^{\text{ar}(T')}$  then
7:      $S_\triangleright := \{\text{new}^*(\sigma', \bar{y})\}$ 
8:   if  $S_\triangleright = \emptyset$  then
9:     return false
10:   $S_\nabla := S_\triangleright$ 
11:   $S_\triangleright := \emptyset$ 
12:   $Check := b$ , where  $b \in \{0, 1\}$ 
13: until  $Check = 1$ 
14: return  $(T' = P \wedge \pi_{(i_1, \dots, i_k)}(\bar{x}) = \pi_{(j_1, \dots, j_k)}(\bar{z}))$ 

```

---

witnesses the following: starting from  $S_\sigma = \{R(x_1, \dots, x_{\text{ar}(R)})\}$ , an atom  $P(y_1, \dots, y_{\text{ar}(P)})$  with  $\pi_{(i_1, \dots, i_m)}(x_1, \dots, x_{\text{ar}(R)}) = \pi_{(j_1, \dots, j_m)}(y_1, \dots, y_{\text{ar}(P)})$  can be obtained via chase applications, which means that  $\text{Chase}(S_\sigma, \Sigma) \models \sigma$ . It is also not difficult to see that if  $\text{Chase}(S_\sigma, \Sigma) \models \sigma$ , then a witness of  $\sigma$  relative to  $\Sigma$  can be extracted from  $\text{Chase}(S_\sigma, \Sigma)$ . This discussion is summarized in the following technical lemma, whose proof is left as an exercise.

**Lemma 11.13.** *Consider a set  $\Sigma$  of INDs over a schema  $\mathbf{S}$ , and an IND  $\sigma$  over  $\mathbf{S}$ . Then,  $\text{Chase}(S_\sigma, \Sigma) \models \sigma$  iff there is a witness of  $\sigma$  relative to  $\Sigma$ .*

By Corollary 11.12 and Lemma 11.13, we have that the problem of checking whether a set  $\Sigma$  of INDs over a schema  $\mathbf{S}$  implies a single IND  $\sigma$  over  $\mathbf{S}$ , boils down to checking whether a witness of  $\sigma$  relative to  $\Sigma$  exists. This is done via the nondeterministic procedure depicted in Algorithm 2. Assume that  $\sigma$  is of the form  $R[i_1, \dots, i_k] \subseteq P[j_1, \dots, j_k]$ . The algorithm first checks whether  $R[i_1, \dots, i_k] = P[j_1, \dots, j_k]$ , in which case a witness of  $\sigma$  relative to  $\Sigma$  trivially exists, and returns **true**. Otherwise, it proceeds to nondeterministically construct a witness of  $\sigma$  relative to  $\Sigma$  (if one exists). This is done by constructing one atom after the other via chase steps, without having to store more than two consecutive atoms. The algorithm starts from  $S_\nabla = \{R(x_1, \dots, x_{\text{ar}(R)})\}$ ;  $S_\nabla$  should be understood as the “current atom”, which at the beginning is  $S_\sigma$ , from which we construct the “next atom”  $S_\triangleright$  in the sequence. The repeat-until loop is responsible for constructing  $S_\triangleright$  from  $S_\nabla$ . This is done by guessing an IND  $\sigma' \in \Sigma$ , and adding to  $S_\triangleright$  the atom  $\text{new}^*(\sigma', \bar{y})$  if  $\sigma'$  is applicable to  $S_\nabla$  with  $\bar{y}$ ; note that  $\bar{y}$  is the single tuple occurring in  $S_\nabla$ . This is repeated until the algorithm chooses to exit the loop by setting  $Check$  to 1, and check whether

$S_{\triangleright}$  consists of an atom  $T'(\bar{z})$  with  $T' = P$  and  $\pi_{(i_1, \dots, i_k)}(\bar{x}) = \pi_{(j_1, \dots, j_k)}(\bar{z})$ , in which case it returns **true**; otherwise, it returns **false**.

It is easy to verify that Algorithm 2 uses polynomial space. This heavily relies on the fact that the atoms generated during its computation contain only variables from  $\{x_1, \dots, x_{\text{ar}(R)}\}$  and the special variable  $\star$ , which in turn implies that  $S_{\nabla}$  and  $S_{\triangleright}$  can be represented using polynomial space. It also takes polynomial space to check if  $R[i_1, \dots, i_k] = P[j_1, \dots, j_k]$  (see line 1), to check if an IND is applicable to  $S_{\nabla}$  with  $\bar{y}$  (see line 6), and to check if  $T' = P$  and  $\pi_{(i_1, \dots, i_k)}(\bar{x}) = \pi_{(j_1, \dots, j_k)}(\bar{z})$  (see line 14). Therefore, **IND-Implication** is in NPSpace, and thus in PSPACE since  $\text{NPSpace} = \text{PSPACE}$ .

A PSPACE lower bound for **IND-Implication** can be shown via a reduction from the following PSPACE-hard problem: given 2-TM  $M$  that runs in linear space, and a word  $w$  over the alphabet of  $M$ , decide whether  $M$  accepts input  $w$ . The formal proof is left as Exercise 1.17.



---

## Exercises

**Exercise 1.1.** Let  $q$  be an FO query. Prove that one can compute in polynomial time an FO query  $q'$  that uses only  $\neg$ ,  $\vee$ , and  $\exists$  such that  $q \equiv q'$ .

**Exercise 1.2.** We say that a query  $q$  from a database schema  $\mathbf{S}$  to a relation schema  $\mathbf{S}'$  is *C-generic*, for some  $C \subseteq \text{Const}$ , if for every database  $D$  of  $\mathbf{S}$ , and for every bijection  $\rho : \text{Const} \rightarrow \text{Const}$  that is the identity on  $C$ ,  $q(\rho(D)) = \rho(q(D))$ . Show that an FO query  $\varphi(\bar{x})$  over a schema  $\mathbf{S}$  is  $\text{Dom}(\varphi)$ -generic.

**Exercise 1.3.** The semantics of the rename and join operations in the named RA has been defined in Chapter 4. Provide formal definitions for the semantics of the other operations, i.e., selection, projection, union, and difference.

**Exercise 1.4.** State and prove the converse of Theorem 4.6

**Exercise 1.5.** Prove that allowing conditions of the form  $\bar{a} \in e$  and  $\text{empty}(e)$  in selection conditions of RA does not increase its expressiveness. In particular, show that selections with these new conditions can be expressed using standard operations of RA.

**Exercise 1.6.** Prove that adding nested subqueries in the FROM clause does not increase expressiveness. In particular, extend the translation from basic SQL to RA that handles nested subqueries in FROM.

**Exercise 1.7.** The proofs of Theorems 7.1 and 7.3 only consider the special case of FO-Evaluation where the query does not use constants, i.e., elements from  $\text{Const}$ . How can the proof be extended to FO-Evaluation in general, i.e., allowing for general FO queries?

**Exercise 1.8.** For showing that FO-Evaluation is PSPACE-hard, we provided a reduction from QSAT. In particular, for an input to QSAT given by  $\psi$ , we constructed a database  $D$  and an FO query  $q_\psi$  (see the proof of Theorem 7.1). Show that  $\psi$  is satisfiable if and only if  $D \models q_\psi$ .

**Exercise 1.9.** For an integer  $k > 0$ , we write  $\text{FO}_k$  for the class of FO queries that can mention at most  $k$  variables. The evaluation problem for the class of  $\text{FO}_k$  queries, for some fixed  $k > 0$ , is defined as expected: given an  $\text{FO}_k$  query  $q$ , a database  $D$ , and a tuple  $\bar{a}$ , decide whether  $\bar{a} \in q(D)$ . Show that the evaluation problem for  $\text{FO}_k$  queries, for a fixed  $k > 0$ , is in PTIME.

**Exercise 1.10.** Let  $q_M$  be the Boolean FO query constructed in the proof of Theorem 8.1. Prove that if the Turing machine  $M$  on the empty word does not halt, then there exists an infinite database  $D$  such that  $q(D) = \text{true}$ .

**Exercise 1.11.** Let FO-Unrestricted-Satisfiability be the unrestricted version of FO-Satisfiability where we consider possibly infinite databases. In other words, FO-Unrestricted-Satisfiability is defined as follows: given an FO query  $q$ , is there a possibly infinite database  $D$  such that  $q(D) \neq \emptyset$ ? Show that FO-Unrestricted-Satisfiability is undecidable by adapting the proof of Theorem 8.1.

**Exercise 1.12.** Prove that FO-Containment remains undecidable even if the left hand-side query is a Boolean query  $q = \exists \bar{x} \varphi$ , where  $\varphi$  is a conjunction of relational atoms or the negation of relational atoms.

**Exercise 1.13.** The algorithms underlying Theorems 10.4 and 11.4 for checking whether a database satisfies a set of FDs and INDs, respectively, were designed with simplicity instead of efficiency in mind. Provide more efficient algorithms for the problems FD-Satisfaction and IND-Satisfaction.

**Exercise 1.14.** Prove that the result of an infinite chase sequence of a finite set of relational atoms under a set of INDs always exists.

**Exercise 1.15.** Prove Theorem 11.11. The non-trivial task is to show that if  $\Sigma \models^\infty \sigma$  does not hold, then also  $\Sigma \models \sigma$  does not hold. One can exploit Lemma 11.10, which states that if  $\Sigma \models^\infty \sigma$  does not hold, then  $\text{Chase}(S_\sigma, \Sigma)$  does not satisfy  $\sigma$ . If  $\text{Chase}(S_\sigma, \Sigma)$  is finite, then we have that  $\Sigma \models \sigma$  does not hold. The main task is, when  $\text{Chase}(S_\sigma, \Sigma)$  is infinite, to convert  $\text{Chase}(S_\sigma, \Sigma)$  into a finite set  $S$  such that  $S \models \Sigma$ , but  $S$  does not satisfy  $\sigma$ .

**Exercise 1.16.** Prove Lemma 11.13.

**Exercise 1.17.** Prove that IND-Implication is PSPACE-hard. To this end, provide a reduction from the following PSPACE-hard problem: given 2-TM  $M$  that runs in linear space, and a word  $w$  over the alphabet of  $M$ , decide whether  $M$  accepts input  $w$ .

---

## Bibliographic Comments

To be done.

## Part II

---

### Conjunctive Queries





## Syntax and Semantics

Conjunctive queries are of special importance to databases. They express relational joins, which correspond to the operation that is most commonly performed by relational database engines. This is because data is typically spread over multiple relations, and thus, to answer queries, one needs to join such relations. Actually, conjunctive queries have the power of select-project-join RA queries, which means that they correspond to a very common type of queries written in Core SQL. The goal of this chapter is to introduce the syntax and semantics of conjunctive queries.

### Syntax of Conjunctive Queries

We start with the syntax of conjunctive queries.

#### Definition 12.1: Syntax of Conjunctive Queries

A *conjunctive query* (CQ) over a schema  $\mathbf{S}$  is an FO query  $\varphi(\bar{x})$  over  $\mathbf{S}$  with  $\varphi$  being a formula of the form

$$\exists \bar{y} (R_1(\bar{u}_1) \wedge \cdots \wedge R_n(\bar{u}_n))$$

for  $n \geq 1$ , where  $R_i(\bar{u}_i)$  is a relational atom, and  $\bar{u}_i$  a tuple of constants and variables mentioned in  $\bar{x}$  and  $\bar{y}$ , for every  $i \in [n]$ .

It is very common to represent CQs via a rule-like syntax, which is reminiscent of the syntax of logic programming rules. In particular, the CQ  $\varphi(\bar{x})$  given in Definition 12.1 can be written as the *rule*

$$\text{Answer}(\bar{x}) \text{ :- } R_1(\bar{u}_1), \dots, R_n(\bar{u}_n),$$

where  $\text{Answer}$  is a relation name not in  $\mathbf{S}$ , and its arity (under the singleton schema  $\{\text{Answer}\}$ ) is equal to the arity of  $q$ . The relational atom  $\text{Answer}(\bar{x})$

that appears on the left of the  $\text{:-}$  symbol is called the *head* of the rule, while the expression  $R_1(\bar{u}_1), \dots, R_n(\bar{u}_n)$  that appears on the right of the  $\text{:-}$  symbol is called the *body* of the rule. In general, throughout the book, we use the rule-like syntax for CQs. Nevertheless, for convenience, we will freely interpret a CQ as a first-order query or as a rule.

### Example 12.2: Conjunctive Queries

Consider again the relational schema from Example 3.2:

```

Person [ pid, pname, cid ]
Profession [ pid, pname ]
City [ cid, cname, country ]

```

The following CQ can be used to retrieve the list of names of computer scientists that were born in the city of Athens in Greece:

$$\exists x \exists z \left( \text{Person}(x, y, z) \wedge \text{Profession}(x, \text{'computer scientist'}) \wedge \text{City}(z, \text{'Athens'}, \text{'Greece'}) \right).$$

In rule-like representation, this query is expressed as follows:

```

Answer(y) :- Person(x, y, z), Profession(x, 'computer scientist'),
              City(z, 'Athens', 'Greece').

```

A CQ  $q$  is *Boolean* if it has no output variables, i.e.,  $\bar{x}$  is the empty tuple. When we write a Boolean CQ as a rule, we simply write Answer as the head, instead of Answer(). For example, the following Boolean CQ checks whether there exists a computer scientist that was born in the city of Athens in Greece:

```

Answer :- Person(x, y, z), Profession(x, 'computer scientist'),
           City(z, 'Athens', 'Greece').

```

## Semantics of Conjunctive Queries

Since CQs are FO queries, the definition of their output on a database can be inherited from Definition 3.6. More precisely, given a database  $D$  of a schema  $\mathbf{S}$ , and a  $k$ -ary CQ  $q = \varphi(\bar{x})$  over  $\mathbf{S}$ , where  $k \geq 0$ , the *output* of  $q$  on  $D$  is

$$q(D) = \{ \bar{a} \in \text{Dom}(D)^k \mid D \models \varphi(\bar{a}) \}.$$

Notice that  $q(D)$  consists of tuples over  $\text{Dom}(D)$ , not over  $\text{Dom}(D) \cup \text{Dom}(\varphi)$ . This is because CQs do not allow for equational atoms, and thus, there is no way for a constant of  $\text{Dom}(\varphi) - \text{Dom}(D)$  to appear in the output.

Interestingly, there is a more intuitive (and equivalent) way of defining the semantics of CQs when they are viewed as rules. The body of a CQ  $q$  of the form  $\text{Answer}(\bar{x}) :- \text{body}$  can be seen as a pattern that must be matched with the database  $D$  via an assignment  $\eta$  that maps the variables in  $q$  to  $\text{Dom}(D)$ . For each such assignment  $\eta$ , if  $\eta$  applied to this pattern produces only facts of  $D$ , it means that the pattern matches with  $D$  via  $\eta$ , and the tuple  $\eta(\bar{x})$  is an output of  $q$  on  $D$ . We proceed to formalize this informal description.

Consider a database  $D$  and a CQ  $q$  of the form

$$\text{Answer}(\bar{x}) :- R_1(\bar{u}_1), \dots, R_n(\bar{u}_n).$$

An *assignment* for  $q$  over  $D$  is a function  $\eta$  from the set of variables in  $q$  to  $\text{Dom}(D)$ . We say that  $\eta$  is *consistent* with  $D$  if

$$\{R_1(\eta(\bar{u}_1)), \dots, R_n(\eta(\bar{u}_n))\} \subseteq D,$$

where, for  $i \in [n]$ , the fact  $R_i(\eta(\bar{u}_i))$  is obtained by replacing each variable  $x$  in  $\bar{u}_i$  with  $\eta(x)$ , and leaving the constants in  $\bar{u}_i$  untouched. The consistency of  $\eta$  with  $D$  essentially means that the body of  $q$  matches with  $D$  via  $\eta$ . Having this notion in place, we can define what is the output of a CQ on a database.

#### Definition 12.3: Evaluation of CQs

Given a database  $D$  of a schema  $\mathbf{S}$ , and a CQ  $q(\bar{x})$  over  $\mathbf{S}$ , the *output* of  $q$  on  $D$  is defined as the set of tuples

$$q(D) = \{\eta(\bar{x}) \mid \eta \text{ is an assignment for } q \text{ over } D \text{ consistent with } D\}.$$

It is an easy exercise to show that the semantics of CQs inherited from the semantics of FO queries in Definition 3.6 and the semantics of CQs given in Definition 12.3, are equivalent, i.e., for a CQ  $q = \varphi(\bar{x})$  and a database  $D$ ,

$$\{\bar{a} \in \text{Dom}(D)^k \mid D \models \varphi(\bar{a})\} = \{\eta(\bar{x}) \mid \eta \text{ is an assignment for } q \text{ over } D \text{ consistent with } D\}.$$

#### Example 12.4: Evaluation of CQs

Let  $\mathbf{S}$  be the schema from Example 3.2, which has been also used in Example 12.2. Let  $D$  be the database of  $\mathbf{S}$  shown in Figure 3.1; we recall the relations *Person* and *Profession* in Figure 12.1. The following CQ  $q$  can be used to retrieve the ids and names of actors:

$$\text{Answer}(x, y) :- \text{Person}(x, y, z), \text{Profession}(x, \text{'actor'}).$$

Observe that the assignment  $\eta$  for  $q$  over  $D$  such that

$$\eta(x) = \text{'1'} \quad \eta(y) = \text{'Aretha'} \quad \eta(z) = \text{'MPH'}$$

Person			Profession	
pid	pname	cid	pid	pname
1	Aretha	MPH	1	singer
2	Billie	BLT	1	songwriter
3	Bob	DLT	1	actor
4	Freddie	ST	2	singer
			3	singer
			3	songwriter
			3	author
			4	singer
			4	songwriter

Fig. 12.1: The relations Person and Profession for Example 12.4

is consistent with  $D$ . Indeed, when applied to the body of  $q$  it produces the facts  $\text{Person}('1', 'Aretha', 'MPH')$  and  $\text{Profession}('1', 'actor')$ , both of which are facts of  $D$ . On the other hand, the assignment  $\eta'$  such that

$$\eta'(x) = '2' \quad \eta'(y) = 'Billie' \quad \eta'(z) = 'BLT'$$

is not consistent with  $D$ . When applied to the body of  $q$ , it generates the fact  $\text{Profession}('2', 'actor')$  that is not in  $D$ . It is straightforward to verify that  $\eta$  is the only assignment for  $q$  over  $D$  that is consistent with  $D$ , which in turn implies that the output of  $q$  on  $D$  is

$$q(D) = \{('1', 'Aretha')\}.$$

If  $q$  is a Boolean CQ, then  $q(D) = \mathbf{true}$  if and only if there is an assignment for  $q$  over  $D$  that is consistent with  $D$ . In other words,  $q(D) = \mathbf{true}$  if and only if the body of the CQ matches with  $D$  via at least one assignment for  $q$  over  $D$ . For instance, if in Example 12.4 we consider also the Boolean CQ  $q'$

$$\text{Answer} \text{ :- } \text{Person}(x, y, z), \text{Profession}(x, 'actor'),$$

which is the Boolean version of  $q$  in Example 12.4, then  $q'(D) = \mathbf{true}$  since the assignment  $\eta$  is consistent with  $D$ . On the other hand, given the CQ  $q''$

$$\text{Answer} \text{ :- } \text{Person}(x, y, z), \text{Profession}(x, 'nurse'),$$

$q''(D) = \mathbf{false}$  since there is no assignment  $\eta$  such that  $\text{Person}(\eta(x), \eta(y), \eta(z))$  and  $\text{Profession}(\eta(x), 'nurse')$  are both facts of  $D$ .

## Conjunctive Queries as a Fragment of FO

When CQs are seen as FO queries they use only relational atoms, conjunction ( $\wedge$ ), and existential quantification ( $\exists$ ). Thus, every CQ can be expressed using

formulae from the fragment of FO that corresponds to the closure of relational atoms under  $\exists$  and  $\wedge$ ; we refer to this fragment of FO as  $\text{FO}^{\text{rel}}[\wedge, \exists]$ . Actually, the converse is also true. Consider a query  $\varphi(\bar{x})$  with  $\varphi$  being an  $\text{FO}^{\text{rel}}[\wedge, \exists]$  formula. It is easy to show that  $\varphi(\bar{x})$  is equivalent to a CQ. We first rename variables in order to ensure that bound variables do not repeat (which leads to an equivalent query), and then push the existential quantifiers outside. This conversion can be easily illustrated via a simple example.

**Example 12.5: From  $\text{FO}^{\text{rel}}[\wedge, \exists]$  Queries to CQs**

Consider the  $\text{FO}^{\text{rel}}[\wedge, \exists]$  query  $\varphi(x)$  with

$$\varphi = (\exists y R(x, a, y)) \wedge (\exists y S(y, x, b)).$$

We first rename the second occurrence of  $y$ , and get the query  $\varphi'(x)$  with

$$\varphi' = (\exists y R(x, a, y)) \wedge (\exists z S(z, x, b)).$$

We then push all the quantifiers outside, and get the CQ  $\varphi''(x)$  with

$$\varphi'' = \exists y \exists z (R(x, a, y) \wedge S(z, x, b)).$$

From the above discussion, we immediately get that:

**Theorem 12.6**

The languages of CQs and of  $\text{FO}^{\text{rel}}[\wedge, \exists]$  queries are equally expressive.

Notice that  $\text{FO}^{\text{rel}}[\wedge, \exists]$  is *not* the same as  $\text{FO}[\wedge, \exists]$ , that is, the fragment of FO that allows only for conjunction ( $\wedge$ ) and existential quantification ( $\exists$ ). Fragments defined by listing a set of features of FO are assumed to be the closure of *all* atomic formulae (including equational atoms) under those features. Therefore, the fragment  $\text{FO}[\wedge, \exists]$  allows also for equational atoms, which means that the query  $\varphi(x, y)$  with  $\varphi = (x = y)$  is an  $\text{FO}[\wedge, \exists]$  query. As we shall see in the next chapter, though,  $\varphi(x, y)$  is not equivalent to a CQ.

## Conjunctive Queries as a Fragment of RA

The class of CQs has the same expressive power as the fragment of RA that is built from base expressions  $R \in \text{Rel}$  and allows for selection, projection, and Cartesian product. Furthermore, conditions in selections are conjunctions of equalities. Note that base expressions of the form  $\{a\}$  with  $a \in \text{Const}$  are not included. This fragment of RA is called the *select-project-join* (SPJ) fragment. Henceforth, we simply refer to the associated queries as SPJ queries. Recall

that the join operation is actually a selection from the Cartesian product on a condition that is a conjunction of equalities. We proceed to show that:

**Theorem 12.7**

The languages of CQs and of SPJ queries are equally expressive.

*Proof.* We first show how to translate a CQ  $q$  of the form

$$\text{Answer}(\bar{x}) \text{ :- } R_1(\bar{u}_1), \dots, R_n(\bar{u}_n)$$

into an SPJ query. In fact,  $q$  can be expressed as the query

$$\pi_\alpha(\sigma_{\theta(q)}(\sigma_{\theta(\bar{u}_1)}(R_1) \times \sigma_{\theta(\bar{u}_2)}(R_2) \times \dots \times \sigma_{\theta(\bar{u}_n)}(R_n))),$$

where conditions in selections, as well as the list of positions in the projections are defined as follows:

- For each  $i \in [n]$ ,  $\theta(\bar{u}_i)$  is a conjunction of statements  $j \doteq a$  and  $j \doteq k$ , where  $a \in \text{Const}$  and  $j, k \in [\text{ar}(R_i)]$ , such that  $j \doteq a$  is a conjunct of  $\theta(\bar{u}_i)$  if and only if the  $j$ -th component of  $\bar{u}_i$  is the constant  $a$ , and  $j \doteq k$  is a conjunct of  $\theta(\bar{u}_i)$  if and only if the  $j$ -th and the  $k$ -th components of  $\bar{u}_i$  are the same variable. If no constant occurs in  $\bar{u}_i$ , and  $\bar{u}_i$  consists of distinct variables, then the selection is omitted; we have  $R_i$  instead of  $\sigma_{\theta(\bar{u}_i)}(R_i)$ .
- The condition  $\theta(q)$  is a conjunction of statements of the form  $j \doteq k$ , where  $j, k \in [\text{ar}(R_1) + \dots + \text{ar}(R_n)]$ , such that  $j \doteq k$  is a conjunct of  $\theta(q)$  if and only if the following hold:
  - (i) if  $j = \text{ar}(R_1) + \dots + \text{ar}(R_\ell) + \ell'$ , for some  $\ell \in [0, n-1]$  and  $\ell' \in [\text{ar}(R_{\ell+1})]$ , then  $k > \text{ar}(R_1) + \dots + \text{ar}(R_{\ell+1})$ , and
  - (ii) the  $j$ -th and the  $k$ -th components of  $\bar{u}_1 \bar{u}_2 \dots \bar{u}_n$  are the same variable.
 Item (i) states that  $j$  and  $k$  should be positions from different  $\bar{u}_i$  tuples.
- Finally,  $\alpha$  is a list of positions among  $\bar{u}_1 \bar{u}_2 \dots \bar{u}_n$  that form the output tuple of variables  $\bar{x}$ .

The correctness of the above translation is left as an exercise. Note that instead of using the condition  $\theta(q)$ , one can replace the Cartesian products by  $\theta$ -joins (recall that the  $\theta$ -join of relations  $R$  and  $S$  is defined as  $R \bowtie_\theta S = \sigma_\theta(R \times S)$ ). Here is a simple example that illustrates the above translation.

**Example 12.8: From CQs to SPJ Queries**

Consider the CQ  $q$  defined as

$$\text{Answer}(x, x, y) \text{ :- } R_1(\underbrace{x, z, z, a, x}_{\bar{u}_1}), R_2(\underbrace{a, y, z, a, b}_{\bar{u}_2}), R_3(\underbrace{x, y, z}_{\bar{u}_3}).$$

It is easy to verify that

$$\begin{aligned}\theta(\bar{u}_1) &= (4 \doteq a) \wedge (1 \doteq 5) \wedge (2 \doteq 3) \\ \theta(\bar{u}_2) &= (1 \doteq a) \wedge (4 \doteq a) \wedge (5 \doteq b),\end{aligned}$$

while the selection operation  $\sigma_{\theta(\bar{u}_3)}$  is omitted since neither a constant nor a repetition of variables occurs in  $\bar{u}_3$ .

The condition  $\theta(q)$  essentially has to specify that in

$$\bar{u}_1 \bar{u}_2 \bar{u}_3 = (x, z, z, a, x, a, y, z, a, b, x, y, z)$$

the variable  $x$  in  $\bar{u}_1$  and the variable  $x$  in  $\bar{u}_3$  are the same, the variable  $z$  in  $\bar{u}_1$  and the variable  $z$  in both  $\bar{u}_2$  and  $\bar{u}_3$  are the same, and that the variable  $y$  in  $\bar{u}_2$  and the variable  $y$  in  $\bar{u}_3$  are the same. This results in

$$\begin{aligned}\theta(q) &= (1 \doteq 11) \wedge (5 \doteq 11) \wedge (2 \doteq 8) \wedge (2 \doteq 13) \wedge \\ &\quad (3 \doteq 8) \wedge (3 \doteq 13) \wedge (8 \doteq 13) \wedge (7 \doteq 12).\end{aligned}$$

Finally,  $\alpha$  corresponds to variable  $x$  repeated twice and variable  $y$ , i.e.,  $\alpha = (1, 1, 7)$ . Summing up, the CQ  $q$  is expressed as

$$\begin{aligned}\pi_{(1,1,7)} \Big( &\sigma_{(1 \doteq 11) \wedge (2 \doteq 8) \wedge (2 \doteq 13) \wedge (7 \doteq 12)} \Big( \sigma_{(4 \doteq a) \wedge (1 \doteq 5) \wedge (2 \doteq 3)}(R_1) \times \\ &\sigma_{(1 \doteq a) \wedge (4 \doteq a) \wedge (5 \doteq b)}(R_2) \times R_3 \Big) \Big).\end{aligned}$$

For the sake of readability, we have eliminated  $(5 \doteq 11)$  from  $\theta(q)$  since it can be derived from  $(1 \doteq 11)$  in  $\theta(q)$  and  $(1 \doteq 5)$  in  $\theta(\bar{u}_1)$ , and likewise for conditions  $(3 \doteq 8)$ ,  $(3 \doteq 13)$  and  $(8 \doteq 13)$  in  $\theta(q)$ .

We now proceed with the other direction, and show that every SPJ query  $e$  can be expressed as a CQ  $q_e$ . The proof is by induction on the structure of  $e$ . We can assume that in  $e$  all selections are either of the form  $\sigma_{i \doteq a}$  or  $\sigma_{i \doteq j}$  (because more complex selections can be obtained by applying a sequence of simple selections). We also assume that all projections are of the form  $\pi_{\bar{i}}$  that exclude the  $i$ -th component; for instance,  $\pi_{\bar{2}}(R)$  applied to a ternary relation  $R$  will transform each tuple  $(a, b, c)$  into  $(a, c)$  by excluding the second component (again, more complex projections are simply sequences of these simple ones).

- If  $e = R$ , where  $R$  is a  $k$ -ary relation, then  $q_e$  is the CQ  $\varphi(\bar{x}) = R(\bar{x})$ , where  $\bar{x}$  is a  $k$ -ary tuple of pairwise distinct and fresh variables.
- If  $e$  is of arity  $k$  with  $q_e = \varphi(x_1, \dots, x_k)$ , where the  $x_i$ 's are not necessarily distinct, then
  - $q_{\sigma_{i \doteq a}(e)}$  is the CQ obtained from  $q_e$  by replacing each occurrence of the variable  $x_i$  by the constant  $a$ ,
  - $q_{\sigma_{i \doteq j}(e)}$  is the CQ obtained from  $q_e$  by replacing each occurrence of the variable  $x_j$  with the variable  $x_i$ , and



- $q_{\pi_i(e)}$  is the CQ  $\varphi(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_k)$  if  $x_i$  occurs among the  $x_j$ 's with  $j \neq i$ , and  $\exists x_i \varphi(x_1, \dots, x_k)$  otherwise.
- If  $e_1$  is  $k$ -ary with  $q_{e_1} = \varphi_1(x_1, \dots, x_k)$  and  $\varphi_1 = \exists \bar{z} \psi_1$ , and  $e_2$  is  $m$ -ary with  $q_{e_2} = \varphi_2(y_1, \dots, y_m)$  and  $\varphi_2 = \exists \bar{w} \psi_2$ , then  $q_{(e_1 \times e_2)}$  is the CQ  $\varphi(x_1, \dots, x_k, y_1, \dots, y_m)$  with  $\varphi = \exists \bar{z} \exists \bar{w} \psi_1 \wedge \psi_2$ ; we assume that  $\psi_1$  and  $\psi_2$  do not share variables.

This completes the construction of the CQ  $q_e$ . The correctness of the above translation is left as an exercise to the reader.

We conclude by explaining further the difference between the two cases of handling projection. Consider the unary relations  $U, V$  and an RA expression  $e = \pi_1(\sigma_{1 \dot{=} 2}(U \times V))$ . First, notice that  $U \times V$  is translated as  $\varphi(x, y) = U(x) \wedge V(y)$ , since the expression  $U$  has to be translated as a relational atom of the form  $U(z)$  where the variable  $z$  is fresh, and likewise for the expression  $V$ ; thus, the occurrences of  $U$  and  $V$  in  $e$  have to be translated considering distinct variables, in this case  $x$  and  $y$ . Then  $\sigma_{1 \dot{=} 2}(U \times V)$  is translated as  $\varphi(x, x) = U(x) \wedge V(x)$ , since  $y$  is replaced with  $x$ . Finally,  $\pi_1(\sigma_{1 \dot{=} 2}(U \times V))$  is obtained by eliminating the first occurrence of  $x$  as an output variable: the CQ defining  $e$  is  $\psi(x) = U(x) \wedge V(x)$ . On the other hand, the correct way to define  $e' = \pi_2(U \times V)$  as a CQ is to existentially quantify over  $y$  in  $\varphi(x, y)$  that defines  $U \times V$ , that is, the CQ  $\psi'(x)$  with  $\psi = \exists y (U(x) \wedge V(y))$ .  $\square$

The following is an immediate corollary of Theorems 12.6 and 12.7 that relates the languages of  $\text{FO}^{\text{rel}}[\wedge, \exists]$  and SPJ queries.

**Corollary 12.9**

The language of  $\text{FO}^{\text{rel}}[\wedge, \exists]$  queries and the language of SPJ queries are equally expressive.

## Homomorphisms and Expressiveness

As already discussed in Chapter 9, homomorphisms are a fundamental tool that plays a key role in various aspects of relational databases. In this chapter, we discuss how homomorphisms emerge in the context of CQs. In particular, we show that they provide an alternative way to describe the evaluation of CQs, and also use them as a tool to understand the expressiveness of CQs.

### CQ Evaluation and Homomorphisms

We can recast the semantics of CQs using the notion of homomorphism. The key observation is that the body of a CQ, written as a rule, can be viewed as a set of atoms. More precisely, given a CQ  $q$  of the form

$$\text{Answer}(\bar{x}) :- R_1(\bar{u}_1), \dots, R_n(\bar{u}_n)$$

we define the set of relational atoms

$$A_q = \{R_1(\bar{u}_1), \dots, R_n(\bar{u}_n)\}.$$

Thus, we can naturally talk about homomorphisms from CQs to databases.

#### Definition 13.1: Homomorphisms from CQs to Databases

Consider a CQ  $q(\bar{x})$  over a schema  $\mathbf{S}$ , and a database  $D$  of  $\mathbf{S}$ . We say that there is a *homomorphism from  $q$  to  $D$* , written as  $q \rightarrow D$ , if  $A_q \rightarrow D$ . We also say that there is a *homomorphism from  $(q, \bar{x})$  to  $(D, \bar{a})$* , written as  $(q, \bar{x}) \rightarrow (D, \bar{a})$ , if  $(A_q, \bar{x}) \rightarrow (D, \bar{a})$ .

To define the output of a CQ  $q(\bar{x})$  on a database  $D$  (see Definition 12.3), we used the notion of assignment for  $q$  over  $D$ , which is a function from the set of variables in  $q$  to  $\text{Dom}(D)$ . The output of  $q$  on  $D$  consists of all the tuples  $\eta(\bar{x})$ , where  $\eta$  is an assignment for  $q$  over  $D$  that is consistent with  $D$ , i.e.,

$$\{R_1(\eta(\bar{u}_1)), \dots, R_n(\eta(\bar{u}_n))\} \subseteq D.$$

Since, for  $i \in [n]$ ,  $R_i(\eta(\bar{u}_i))$  is the fact obtained after replacing each variable  $x$  in  $\bar{u}_i$  with  $\eta(x)$ , and leaving the constants in  $\bar{u}_i$  untouched, such an assignment  $\eta$  corresponds to a function  $h : \text{Dom}(A_q) \rightarrow \text{Dom}(D)$ , which is the identity on the constants occurring in  $q$ , such that  $R(h(\bar{u}_i)) = R(\eta(\bar{u}_i))$ . But, of course, this is the same as saying that  $h$  is a homomorphism from  $q$  to  $D$ . Therefore,  $q(D)$  is the set of all tuples  $h(\bar{x})$ , where  $h$  is a homomorphism from  $q$  to  $D$ , i.e., the set of all tuples  $\bar{a}$  over  $\text{Dom}(D)$  with  $(q, \bar{x}) \rightarrow (D, \bar{a})$ . This leads to an alternative characterization of CQ evaluation in terms of homomorphisms.

**Theorem 13.2**

Given a database  $D$  of a schema  $\mathbf{S}$ , and a CQ  $q(\bar{x})$  of arity  $k \geq 0$  over  $\mathbf{S}$ ,

$$q(D) = \{\bar{a} \in \text{Dom}(D)^k \mid (q, \bar{x}) \rightarrow (D, \bar{a})\}.$$

Here is a simple example that illustrates the above characterization.

**Example 13.3: CQ Evaluation via Homomorphisms**

Let  $D$  and  $q$  be the database and the CQ, respectively, that have been considered in Example 12.4. We know that  $q(D) = \{('1', 'Aretha')\}$ . By the characterization given in Theorem 13.2 we conclude that  $(q, (x, y)) \rightarrow (D, ('1', 'Aretha'))$ . To verify that this is the case, recall that we need to check whether  $(A_q, (x, y)) \rightarrow (D, ('1', 'Aretha'))$ , where

$$A_q = \{\text{Person}(x, y, z), \text{Profession}(x, \text{'actor'})\}.$$

Consider the function  $h : \text{Dom}(A_q) \rightarrow \text{Dom}(D)$  such that

$$h(x) = '1' \quad h(y) = 'Aretha' \quad h(z) = 'MPH' \quad h(\text{'actor'}) = \text{'actor'}.$$

It is clear that the following facts belong to  $D$ :

$$\begin{aligned} \text{Person}(h(x), h(y), h(z)) &= \text{Person}('1', 'Aretha', 'MPH') \\ \text{Profession}(h(x), h(\text{'actor'})) &= \text{Profession}('1', \text{'actor'}) \end{aligned}$$

Moreover,  $h((x, y)) = ('1', 'Aretha')$ . Thus,  $h$  is a homomorphism from  $(A_q, (x, y))$  to  $(D, ('1', 'Aretha'))$ , witnessing that

$$(A_q, (x, y)) \rightarrow (D, ('1', 'Aretha')).$$

## Preservation Results for CQs

Some particularly useful properties of CQs are their preservation under various operations, such as application of homomorphisms, or taking direct products. These properties will provide a precise explanation of the expressiveness of CQs as a subclass of FO queries.

### *Preservation Under Homomorphisms*

By saying that a query  $q$  is preserved under homomorphisms, we essentially mean the following: if a tuple  $\bar{a}$  belongs to the output of  $q$  on a database  $D$ , and  $(D, \bar{a}) \rightarrow (D', \bar{b})$ , then  $\bar{b}$  should belong to the output of  $q$  on  $D'$ . Although we can naturally talk about homomorphisms among databases (since databases are sets of relational atoms), there is a caveat that is related to the fact that homomorphisms are the identity on constant values. Since  $\text{Dom}(D) \subseteq \text{Const}$  for every database  $D$ , it follows that  $D \rightarrow D'$  if and only if  $D \subseteq D'$ . Thus, the notion of homomorphism among databases is actually subset inclusion. However, the intention underlying the notion of homomorphism is to preserve the structure, possibly by leaving some constants unchanged.

To overcome this mismatch, we need a mechanism that allows us to convert a database into a set of relational atoms by replacing constant values with variables.<sup>1</sup> To this end, for a finite set of constants  $C \subseteq \text{Const}$ , we define an injective function  $V_C : \text{Const} \rightarrow \text{Const} \cup \text{Var}$  such that

- $V_C$  is the identity on  $C$  and
- $V_C(a) \in \text{Var}$  for every  $a \notin C$ .

We then write  $(D, \bar{a}) \rightarrow_C (D', \bar{b})$  if  $(V_C(D), V_C(\bar{a})) \rightarrow (D', \bar{b})$ . Note that in  $V_C(D)$  and  $V_C(\bar{a})$  all constants, except for those in  $C$ , have been replaced by variables, so the definition of homomorphism no longer trivializes to being a subset.

---

<sup>1</sup> This is essentially the opposite of grounding a set of atoms discussed in Chapter 9.

**Example 13.4: Homomorphisms Among Databases**

Consider the databases

$$D_1 = \{R(a, b), R(b, a)\} \quad D_2 = \{R(c, c)\}.$$

If  $C_1 = \emptyset$ , then we have that  $V_{C_1}(a)$  and  $V_{C_1}(b)$  are distinct elements of  $\text{Var}$ , let say  $V_{C_1}(a) = x$  and  $V_{C_1}(b) = y$ . Hence,

$$V_{C_1}(D_1) = \{R(x, y), R(y, x)\} \quad V_{C_1}((a, b)) = (x, y),$$

from which we conclude that  $(D_1, (a, b)) \rightarrow_{C_1} (D_2, (c, c))$  since

$$(V_{C_1}(D_1), V_{C_1}((a, b))) \rightarrow (D_2, (c, c)).$$

On the other hand, if  $C_2 = \{a, b\}$ , then

$$V_{C_2}(D_1) = \{R(a, b), R(b, a)\} \quad V_{C_2}((a, b)) = (a, b).$$

Therefore, it does not hold that  $(D_1, (a, b)) \rightarrow_{C_2} (D_2, (c, c))$ , since it does not hold that  $(V_{C_2}(D_1), V_{C_2}((a, b))) \rightarrow (D_2, (c, c))$ .

We can now define the notion of preservation under homomorphisms.

**Definition 13.5: Preservation Under Homomorphisms**

Consider a  $k$ -ary FO query  $q = \varphi(\bar{x})$  over a schema  $\mathbf{S}$ . We say that  $q$  is *preserved under homomorphisms* if, for every two databases  $D$  and  $D'$  of  $\mathbf{S}$ , and tuples  $\bar{a} \in \text{Dom}(D)^k$  and  $\bar{b} \in \text{Dom}(D')^k$ , it holds that

$$(D, \bar{a}) \rightarrow_{\text{Dom}(\varphi)} (D', \bar{b}) \text{ and } \bar{a} \in q(D) \text{ implies } \bar{b} \in q(D').$$

We then show the following for CQs.

**Proposition 13.6**

Every CQ is preserved under homomorphisms.

*Proof.* Consider a  $k$ -ary CQ  $q(\bar{x})$  over a schema  $\mathbf{S}$ , and let  $C$  be the set of constants in  $q$ . Assume that  $(D, \bar{a}) \rightarrow_C (D', \bar{b})$  for some databases  $D, D'$  of  $\mathbf{S}$ , and tuples  $\bar{a} \in \text{Dom}(D)^k$  and  $\bar{b} \in \text{Dom}(D')^k$ . Assume also that  $\bar{a} \in q(D)$ . Let  $h$  be a homomorphism witnessing  $(V_C(D), V_C(\bar{a})) \rightarrow (D', \bar{b})$ . By Theorem 13.2,  $(q, \bar{x}) \rightarrow (D, \bar{a})$  via some  $h'$ . It holds that  $h_q = V_C \circ h'$  is a homomorphism

witnessing  $(q, \bar{x}) \rightarrow (\mathbf{V}_C(D), \mathbf{V}_C(\bar{a}))$  since  $h_q$  is the identity on  $C$ ; indeed, for  $a \in C$ ,  $\mathbf{V}_C(h'(a)) = a$  by definition. Observe that  $h \circ h_q$  is a homomorphism from  $(q, \bar{x})$  to  $(D', \bar{b})$ , and thus, by Theorem 13.2,  $\bar{b} \in q(D')$ , as needed.  $\square$

Another key property is that of monotonicity. A query  $q$  over a schema  $\mathbf{S}$  is *monotone* if, for every two databases  $D$  and  $D'$  of  $\mathbf{S}$ , we have that

$$D \subseteq D' \text{ implies } q(D) \subseteq q(D').$$

We show that homomorphism preservation implies monotonicity of CQs.

**Corollary 13.7**

Every CQ is monotone.

*Proof.* Let  $q$  be a CQ over  $\mathbf{S}$ , and  $C$  be the set of constants occurring in  $q$ . Consider the databases  $D, D'$  of  $\mathbf{S}$  such that  $D \subseteq D'$ , and assume that  $\bar{a} \in q(D)$ . It is clear that  $\mathbf{V}_C^{-1}$  is a homomorphism from  $(\mathbf{V}_C(D), \mathbf{V}_C(\bar{a}))$  to  $(D', \bar{a})$  and thus,  $(D, \bar{a}) \rightarrow_C (D', \bar{a})$ . By Proposition 13.6, we get that  $\bar{a} \in q(D')$ .  $\square$

*Preservation under Direct Products*

The second preservation result stated here concerns *direct products*. Given two directed graphs  $G_1 = (V_1, E_1)$  and  $G_2 = (V_2, E_2)$ , their direct product  $G_1 \otimes G_2$  has  $V_1 \times V_2$  as the set of vertices, i.e., each vertex is a pair  $(v_1, v_2)$  with  $v_1 \in V_1$  and  $v_2 \in V_2$ . In  $G_1 \otimes G_2$  there is an edge between  $(v_1, v_2)$  and  $(v'_1, v'_2)$  if there is an edge from  $v_1$  to  $v'_1$  in  $E_1$  and from  $v_2$  to  $v'_2$  in  $E_2$ . Note that the notion of direct product is different from that of Cartesian product. Indeed, the Cartesian product of two binary relations is a 4-ary relation, while their direct product is still binary.

The definition of direct products for databases is essentially the same, modulo one small technical detail. Elements of databases come from  $\mathbf{Const}$ . For two constants  $a_1$  and  $a_2$ , the pair  $(a_1, a_2)$  is not an element of  $\mathbf{Const}$ , but we can think of it as such. Indeed, since  $\mathbf{Const}$  is countably infinite, there is a *pairing* function, i.e., a bijection  $\tau : \mathbf{Const} \times \mathbf{Const} \rightarrow \mathbf{Const}$ . A typical example, assuming that  $\mathbf{Const}$  is enumerated as  $c_0, c_1, c_2, \dots$ , is to define  $\tau(c_n, c_m) = c_k$  for  $k = (n + m)(n + m + 1)/2 + m$ . Given a pairing function, we can think of  $(a_1, a_2)$  as being in  $\mathbf{Const}$ , represented by  $\tau(a_1, a_2)$ , and then simply extend the previous definition to arbitrary databases as follows. Given two databases  $D$  and  $D'$  of a schema  $\mathbf{S}$ , their direct product  $D \otimes D'$  is a database of  $\mathbf{S}$  that, for each  $n$ -ary relation name  $R$  in  $\mathbf{S}$ , contains the following facts:

$$R(\tau(a_1, a'_1), \dots, \tau(a_n, a'_n)) \text{ where } R(a_1, \dots, a_n) \in D \text{ and } R(a'_1, \dots, a'_n) \in D'.$$

Technically speaking, this definition depends on the choice of a pairing function, but this choice is irrelevant for FO queries (see Exercise 2.4).

We proceed to define the notion of preservation under direct products. We do this for Boolean queries without constants, as this suffices to understand the limitations of CQs. Exercises 2.6 and 2.7 explain how these results can be extended to queries with constants and free variables, respectively.

**Definition 13.8: Preservation under Direct Products**

A Boolean FO query  $q$  over a schema  $\mathbf{S}$  is *preserved under direct products* if, for every two databases  $D$  and  $D'$  of  $\mathbf{S}$ , it holds that

$$D \models q \text{ and } D' \models q \text{ implies } D \otimes D' \models q.$$

We then show the following for CQs.

**Proposition 13.9**

Every Boolean CQ is preserved under direct products.

*Proof.* As stated earlier, for technical clarity, we only consider CQs that do not mention constants, but the result holds even for CQs with constants (see Exercise 2.6). Let  $q$  be a Boolean CQ without constants over a schema  $\mathbf{S}$ , and let  $D, D'$  be databases of  $\mathbf{S}$  such that  $D \models q$  and  $D' \models q$ . By Theorem 13.2 there are homomorphisms  $h, g$  witnessing  $q \rightarrow D$  and  $q \rightarrow D'$ , respectively. Define now  $f(x) = \tau(h(x), g(x))$ . Assume that  $R(u_1, \dots, u_n)$  is an atom in  $q$ . Then  $R(h(u_1), \dots, h(u_n)) \in D$  and  $R(g(u_1), \dots, g(u_n)) \in D'$ . Hence,

$$R(f(u_1), \dots, f(u_n)) = R(\tau(h(u_1), g(u_1)), \dots, \tau(h(u_n), g(u_n)))$$

belongs to  $D \otimes D'$ , proving that  $f$  is a homomorphism from  $q$  to  $D \otimes D'$ . Thus, by Theorem 13.2  $D \otimes D' \models q$ , as needed.  $\square$

## Expressiveness of CQs

The above preservation results allow us to delineate the expressiveness boundaries of CQs. By Theorem 12.7, CQs and SPJ queries, that is, RA queries that do *not* have inequality in selections, union (and disjunction in selection conditions), and difference, are equally expressive. We prove that none of these is expressible as a CQ. Also notice that in the definition of CQs we disallow explicit equality: CQs correspond to  $\text{FO}^{\text{rel}}[\wedge, \exists]$  queries, i.e., FO queries based on the fragment of FO that is the closure of relational atoms under  $\exists$  and  $\wedge$ . Implicit equality is, of course, allowed by reusing variables. We show that by adding explicit equality one obtains queries that cannot be expressed as CQs.

**CQs cannot express inequality.** This is because CQs with inequality are not preserved under homomorphisms.<sup>2</sup> Consider, e.g., the FO query

$$q_1 = \exists x \exists y (R(x, y) \wedge x \neq y).$$

For  $D = \{R(a, b)\}$  and  $D' = \{R(c, c)\}$ , we have that  $D \rightarrow_\emptyset D'$ . However,  $D \models q_1$  while  $D' \not\models q_1$ . As a second example, consider the FO query

$$q_2 = \exists x (S(x) \wedge x \neq a),$$

where  $a$  is a constant. Given  $D = \{S(b)\}$  and  $D' = \{S(a)\}$ , we have that  $D \rightarrow_{\{a\}} D'$ . However,  $D \models q_2$  while  $D' \not\models q_2$ .

**CQs cannot express negative relational atoms.** The reason is because such queries are not monotone. Consider, for example, the FO query

$$q = \neg P(a),$$

where  $a$  is a constant. If we take  $D = \emptyset$  and  $D' = \{P(a)\}$ , then  $D \subseteq D'$  but  $D \models q$  while  $D' \not\models q$ .

**CQs cannot express difference.** This is because difference is not monotone. Consider, for example, the FO query

$$q = \exists x (P(x) \wedge \neg Q(x)).$$

For  $D = \{P(a)\} \subseteq D' = \{P(a), Q(a)\}$ , we have that  $D \models q$  while  $D' \not\models q$ .

**CQs cannot express union.** This is because such queries are not preserved under direct products. Consider, for example, the FO query

$$q = \exists x (R(x) \vee S(x)).$$

Let  $D = \{R(a)\}$  and  $D' = \{S(a)\}$ . Then,  $D \models q$  and  $D' \models q$ , but  $D \otimes D'$  is empty, and thus,  $D \otimes D' \not\models q$ .

**CQs cannot express explicit equality.** This is because such queries are not preserved under direct products. Consider, for example, the FO query

$$q = \exists x \exists y (x = y).$$

Let  $D = \{R(a)\}$  and  $D' = \{S(a)\}$ . Observe that  $D \models q$  and  $D' \models q$ , but  $D \otimes D' \not\models q$  since  $D \otimes D'$  is empty.

---

<sup>2</sup> Conjunctive queries with inequality are studied in-depth in Chapter 30.





## Query Evaluation

In this chapter, we study the complexity of evaluating conjunctive queries, that is, **CQ-Evaluation**. Recall that this is the problem of checking whether  $\bar{a} \in q(D)$  for a CQ query  $q$ , a database  $D$ , and a tuple  $\bar{a}$  over  $\text{Dom}(D)$ . Recall that for FO queries the same problem is PSPACE-complete (Theorem 7.1). As we show next, the complexity for CQs lies in NP.

### Theorem 14.1

CQ-Evaluation is NP-complete.

*Proof.* We start with the upper bound. Consider a CQ  $q(\bar{x})$ , a database  $D$ , and a tuple  $\bar{a} \in \text{Dom}(D)$ . By Theorem 13.2,  $\bar{a} \in q(D)$  if and only if  $(q, \bar{x}) \rightarrow (D, \bar{a})$ . Therefore, we need to show that checking whether there exists a homomorphism from  $(q, \bar{x})$  to  $(D, \bar{a})$  is in NP. This is done by guessing a function  $h : \text{Dom}(A_q) \rightarrow \text{Dom}(D)$ , and then verifying that  $h$  is a homomorphism from  $(A_q, \bar{x})$  to  $(D, \bar{a})$ , i.e.,  $h$  is the identity on  $\text{Dom}(A_q) \cap \text{Const}$ , and  $R(\bar{u}) \in A_q$  implies  $R(h(\bar{u})) \in D$ . Since both steps are feasible in polynomial time, we conclude that checking whether  $(q, \bar{x}) \rightarrow (D, \bar{a})$  is in NP, as needed.

For the lower bound, we provide a reduction from a graph-theoretic problem, called **Clique**, which is NP-complete. Recall that a *clique* in an undirected graph  $G = (V, E)$  is a complete subgraph  $G' = (V', E')$  of  $G$ , i.e., every two distinct nodes of  $V'$  are connected via an edge of  $E'$ . We say that such a clique is of size  $k \geq 1$  if  $V'$  consists of  $k$  nodes. The problem **Clique** follows:

### Problem: Clique

**Input:** An undirected graph  $G$ , and an integer  $k \geq 1$   
**Output:** **true** if  $G$  has a clique of size  $k$ , and **false** otherwise

Consider an input to **Clique** given by  $G = (V, E)$  and  $k \geq 1$ . The goal is to construct in polynomial time a database  $D$  and a Boolean CQ  $q$  such that  $G$

has a clique of size  $k$  if and only if  $D \models q$ . We construct the database

$$D = \{\text{Node}(v) \mid v \in V\} \cup \{\text{Edge}(v, u) \mid \{v, u\} \in E \text{ and } v \neq u\},$$

which essentially stores the undirected graph  $G$ , but without loops of the form  $\{v\}$  that may occur in  $E$  (that is,  $\{v\} \in E$ ). We can eliminate loops, which is crucial for the correctness of the CQ that we construct next, since they do not affect the existence of a clique of size  $k$  in  $G$ , i.e.,  $G$  has a clique of size  $k$  if and only if  $G'$  obtained from  $G$  after eliminating the loops has a clique of size  $k$ . We also construct

$$q = \exists x_1 \cdots \exists x_k \left( \bigwedge_{i=1}^k \text{Node}(x_i) \wedge \bigwedge_{i,j \in [k] : i \neq j} \text{Edge}(x_i, x_j) \right),$$

which asks whether  $G$  has a clique of size  $k$ . It is clear that  $D$  and  $q$  can be constructed in polynomial time from  $G$  and  $k$ . Moreover, it is easy to see that  $G$  has a clique of size  $k$  if and only if  $D \models q$ , and the claim follows.  $\square$

The data complexity of CQ-Evaluation is immediately inherited from FO-Evaluation (see Theorem 7.3) since CQs are FO queries. Recall that, by convention, CQ-Evaluation is in a complexity class  $\mathcal{C}$  in data complexity if, for every CQ query  $q$ , the problem  $q$ -Evaluation, which takes as input a database  $D$  and a tuple  $\bar{a}$  over  $\text{Dom}(D)$ , and asks whether  $\bar{a} \in q(D)$ , is in  $\mathcal{C}$ .

#### Corollary 14.2

CQ-Evaluation is in DLOGSPACE in data complexity.

Actually, as discussed in Chapter 7, FO-Evaluation, and thus CQ-Evaluation, is in  $\text{AC}_0$  in data complexity, a class that is properly contained in DLOGSPACE. Recall that  $\text{AC}_0$  consists of those languages that are accepted by polynomial-size circuits of constant depth and unbounded fan-in.

## Parameterized Complexity

As discussed in Chapter 2, queries are typically much smaller than databases in practice. This motivated the notion of data complexity, where the cost of evaluation is measured only in terms of the size of the database, while the query is considered to be fixed. However, an algorithm that runs, for example, in time  $O(\|D\|^{\|q\|})$ , although is tractable in terms of data complexity since  $\|q\|$  is a constant, it cannot be considered to be really practical when the database  $D$  is very large, even if the query  $q$  is small. This suggests that we need to rely on a finer notion of complexity than data complexity for classifying query evaluation algorithms as practical or impractical.

This finer notion of complexity is *parameterized complexity*, which is relevant whenever we need to classify the complexity of a problem depending on some central parameters. In the context of query evaluation, it is sensible to consider the size of the database and the size of the query as separate parameters when designing evaluation algorithms, and target algorithms that take less time on the former parameter. For example, a query evaluation algorithm that runs in time  $O(\|D\| \cdot \|q\|^2)$  is expected to perform better in practice than an algorithm that runs in time  $O(\|D\|^2 \cdot \|q\|)$ . Moreover, if the difference between  $\|D\|$  and  $\|q\|$  is significant, as it usually happens in real-life, then even an algorithm that runs in time  $O(\|D\| \cdot 2^{\|q\|})$  could perform better in practice than an algorithm that runs in time  $O(\|D\|^2 \cdot \|q\|)$ .

#### Background on Parameterized Complexity

Before studying the parameterized complexity of CQ-Evaluation when considering the size of the database and the size of the query as separate parameters, we first need to introduce some fundamental notions of parameterized complexity. We start with the notion of parameterized problem (or language).

##### Definition 14.3: Parameterized Problem

Consider a finite alphabet  $\Sigma$ . A *parameterization* of  $\Sigma^*$  is a polynomial time computable function  $\kappa : \Sigma^* \rightarrow \mathbb{N}$ . A *parameterized problem (over  $\Sigma$ )* is a pair  $(L, \kappa)$ , where  $L \subseteq \Sigma^*$ , and  $\kappa$  is a parameterization of  $\Sigma^*$ .

A typical example of such a problem is the parameterized version of Clique.

##### Example 14.4: Parameterized Clique

Recall that **Clique** is the set of pairs  $(G, k)$ , where  $G$  is an undirected graph that contains a clique of size  $k \geq 1$ . Assume that graph-integer pairs are encoded as words over some finite alphabet  $\Sigma$ . Let  $\kappa : \Sigma^* \rightarrow \mathbb{N}$  be the parameterization of  $\Sigma^*$  defined by

$$\kappa(w) = \begin{cases} k & \text{if } w \text{ is the encoding of a graph-integer pair } (G, k), \\ 1 & \text{otherwise,} \end{cases}$$

for  $w \in \Sigma^*$ . We denote the parameterized problem  $(\text{Clique}, \kappa)$  as **p-Clique**.

The input to a parameterized problem  $(L, \kappa)$  over the alphabet  $\Sigma$  is a word  $w \in \Sigma^*$ , and the numbers  $\kappa(w)$  are the corresponding *parameters*. Similarly to (non-parameterized) problems that are represented in the form input-output, we will represent parameterized problems in the form input-parameter-output. For example, **p-Clique** is represented as follows:

**Problem: p-Clique**

**Input:** An undirected graph  $G$ , and an integer  $k \geq 1$   
**Parameter:**  $k$   
**Output:** **true** if  $G$  has a clique of size  $k$ , and **false** otherwise

Analogously, we can talk about the parameterized version of CQ-Evaluation, where the parameter is the size of the query:

**Problem: p-CQ-Evaluation**

**Input:** A CQ  $q(\bar{x})$ , a database  $D$ , and a tuple  $\bar{a}$  over  $\text{Dom}(D)$   
**Parameter:**  $\|q\|$   
**Output:** **true** if  $\bar{a} \in q(D)$ , and **false** otherwise

Recall that the motivation underlying parameterized complexity is to have a finer notion of complexity that allows us to classify algorithms as practical or impractical. But when an algorithm in the realm of parameterized complexity is considered to be practical? This brings us to *fixed-parameter tractability*.

**Definition 14.5: Fixed-Parameter Tractability**

Consider a finite alphabet  $\Sigma$ , and a parameterization  $\kappa : \Sigma^* \rightarrow \mathbb{N}$  of  $\Sigma^*$ . An algorithm  $A$  with input alphabet  $\Sigma$  is an *fpt-algorithm with respect to  $\kappa$*  if there exists a computable function  $f : \mathbb{N} \rightarrow \mathbb{R}_0^+$ , and a polynomial  $p(\cdot)$  such that, for every  $w \in \Sigma^*$ ,  $A$  on input  $w$  runs in time

$$O(p(|w|) \cdot f(\kappa(w))).$$

A parameterized problem  $(L, \kappa)$  is *fixed-parameter tractable* if there is an fpt-algorithm with respect to  $\kappa$  that decides  $L$ . We write FPT for the class of all fixed-parameter tractable problems.

In simple words,  $(L, \kappa)$  is fixed-parameter tractable if there is an algorithm that decides whether  $w \in L$  in time arbitrarily large in the parameter  $\kappa(w)$ , but polynomial in the size of the input  $w$ . This reflects the assumption that  $\kappa(w)$  is much smaller than  $|w|$ , and thus, an algorithm that runs, e.g., in time  $O(|w| \cdot 2^{\kappa(w)})$  is preferable than one that runs in time  $O(|w|^{\kappa(w)})$ .

Whenever we deal with an intractable problem, e.g., the problem of concern of this chapter, i.e., CQ-Evaluation, it would be ideal to be able to show that its parameterized version is in FPT. The reader may be tempted to think that p-CQ-Evaluation is in FPT, and that this can be easily shown by exploiting the algorithm for proving that CQ-Evaluation is in NP. It turns out that this is not true. Consider a CQ  $q(\bar{x})$ , a database  $D$ , and a tuple  $\bar{a}$  over  $\text{Dom}(D)$ .

To check if  $\bar{a} \in q(D)$ , we can iterate over all functions  $h : \text{Dom}(A_q) \rightarrow \text{Dom}(D)$  until we find one that is a homomorphism from  $(A_q, \bar{x})$  to  $(D, \bar{a})$ , in which case we return **true**; otherwise, we return **false**. Since there are  $|\text{Dom}(D)|^{|\text{Dom}(A_q)|}$  such functions, we conclude that this algorithm runs in time

$$O(\|D\|^{\|q\|} \cdot r(\|D\| + \|q\|))$$

for some polynomial  $r(\cdot)$ ; note that the size of  $\bar{a}$  is not included in the bound since it is polynomially bounded by  $\|D\|$  and  $\|q\|$ . Therefore, we cannot conclude that **p-CQ-Evaluation** is in FPT since the expression that describes the running time of the above algorithm is not of the form  $O(p(\|D\|) \cdot f(\|q\|))$ , for some polynomial  $p(\cdot)$  and computable function  $f : \mathbb{N} \rightarrow \mathbb{R}_0^+$ , as required by fixed-parameter tractability in Definition 14.5.

It is widely believed that there is no fpt-algorithm that decides the parameterized version of **CQ-Evaluation**. But then the natural question that comes up is the following: how can we prove that a parameterized problem is not in FPT? Several complexity classes have been defined in the context of parameterized complexity in order to prove that a parameterized problem is not in FPT. Such classes are widely believed to properly contain FPT. This means that if a parameterized problem is complete for one of those classes, then this is a strong indication that the problem in question is not in FPT. Notice here the analogy with classes such as NP and PSPACE: it is not known whether these classes properly contain PTIME, but if a problem is complete for any of them, then this is considered as a strong evidence that the problem is not tractable. We proceed to define one of such classes, namely W[1], which will allow us to pinpoint the exact complexity of **p-CQ-Evaluation**.

To define the class W[1], we need to introduce some auxiliary terminology. Consider a schema **S**. Let  $X$  be a relation name of arity  $m \geq 0$  that does not belong to **S**, and  $\varphi$  an FO sentence over  $\mathbf{S} \cup \{X\}$ . For a database  $D$  of **S**, and a relation  $S \subseteq \text{Dom}(D)^m$ , we write  $D \models \varphi(S)$  to indicate that  $D' \models \varphi$ , where  $D' = D \cup \{X(\bar{a}) \mid \bar{a} \in S\}$ . We further define the problem **p-WD $_{\varphi}$**  as follows:

**Problem: p-WD $_{\varphi}$**

**Input:** A database  $D$  of the schema **S**, and  $k \in \mathbb{N}$   
**Parameter:**  $k$   
**Output:** **true** if there exists  $S \subseteq \text{Dom}(D)^m$  such that  $|S| = k$  and  $D \models \varphi(S)$ , and **false** otherwise

Notice that the sentence  $\varphi$  is fixed in the definition of **p-WD $_{\varphi}$** . Therefore, a different FO sentence  $\psi$  of the form described above gives rise to a different parameterized problem, dubbed **p-WD $_{\psi}$** . The last notion that we need before introducing the class W[1] is that of FPT-reduction.

An *FPT-reduction* from a parameterized problem  $(L_1, \kappa_1)$  over  $\Sigma_1$  to a parameterized problem  $(L_2, \kappa_2)$  over  $\Sigma_2$  is a function  $\Phi : \Sigma_1^* \rightarrow \Sigma_2^*$  such that

the following holds: there are computable functions  $f, g : \mathbb{N} \rightarrow \mathbb{R}_0^+$ , and a polynomial  $p(\cdot)$ , such that, for every word  $w \in \Sigma_1^*$ :

1.  $w \in L_1$  if and only if  $\Phi(w) \in L_2$ ,
2.  $\Phi(w)$  can be computed in time  $p(|w|) \cdot f(\kappa_1(w))$ , and
3.  $\kappa_2(\Phi(w)) \leq g(\kappa_1(w))$ .

The first and the second conditions are natural. The third condition is needed to ensure the crucial property that FPT is closed under FPT-reductions: if there exists an FPT-reduction from  $(L_1, \kappa_1)$  to  $(L_2, \kappa_2)$ , and  $(L_2, \kappa_2) \in \text{FPT}$ , then  $(L_1, \kappa_1) \in \text{FPT}$ ; the proof is left as an exercise.

We now have all the ingredients needed for introducing the class  $\text{W}[1]$ . Recall that *universal* FO sentences are FO sentences of the form  $\forall x_1 \cdots \forall x_n \psi$ , where  $\psi$  is quantifier free and  $\text{FV}(\psi) = \{x_1, \dots, x_n\}$ .

**Definition 14.6: The Class  $\text{W}[1]$**

A parameterized problem  $(L, \kappa)$  is in  $\text{W}[1]$  if there exists a schema  $\mathbf{S}$ , a relation name  $X$  not in  $\mathbf{S}$ , and a universal FO sentence  $\varphi$  over  $\mathbf{S} \cup \{X\}$ , such that there exists an FPT-reduction from  $(L, \kappa)$  to  $\text{p-WD}_\varphi$ .

To give some intuition about the definition of  $\text{W}[1]$ , we show that  $\text{p-Clique}$  is in  $\text{W}[1]$ . We first define a universal FO sentence  $\varphi$ , and then show that there exists an FPT-reduction from  $\text{p-Clique}$  to  $\text{p-WD}_\varphi$ . Assume that  $\mathbf{S}$  consists of the relation names  $\text{Node}[1]$  and  $\text{Edge}[2]$ . Let also  $\text{Elem}[1]$  be a relation name not in  $\mathbf{S}$ . We define the universal FO sentence  $\varphi$  over  $\mathbf{S} \cup \{\text{Elem}\}$

$$\forall x \forall y ((\text{Elem}(x) \wedge \text{Elem}(y) \wedge x \neq y) \rightarrow \text{Edge}(x, y)).$$

We proceed to show that there is an FPT-reduction from  $\text{p-Clique}$  to  $\text{p-WD}_\varphi$ . Consider an input to  $\text{p-Clique}$  given by  $G = (V, E)$  and  $k \geq 1$ . Let

$$D = \{\text{Node}(v) \mid v \in V\} \cup \{\text{Edge}(v, u) \mid \{v, u\} \in E \text{ and } v \neq u\}.$$

The sentence  $\varphi$  checks whether the nodes in the relation  $\text{Elem}$  form a clique. Thus,  $G$  has a clique of size  $k$  if and only if there exists  $S \subseteq \text{Dom}(D)$  such that  $|S| = k$  and  $D \models \varphi(S)$ . It is also clear that  $(D, k)$  can be computed in polynomial time. Therefore, the above reduction from  $\text{p-Clique}$  to  $\text{p-WD}_\varphi$  is an FPT-reduction, which in turn implies that  $\text{p-Clique} \in \text{W}[1]$ .

Before we proceed with the parameterized complexity of  $\text{CQ-Evaluation}$ , let us comment on the nomenclature of  $\text{W}[1]$ . The class  $\text{W}[1]$  is the first level of a hierarchy of complexity classes  $\text{W}[t]$ , for each  $t \geq 1$ ; hence the number 1. More specifically, the class  $\text{W}[t]$  is defined in the same way as the class  $\text{W}[1]$ , but allowing the FO sentence  $\varphi$  in  $\text{p-WD}_\varphi$  to be of the form  $\forall \bar{x}_1 \exists \bar{x}_2 \cdots Q \bar{x}_t \psi$ , where  $\psi$  is quantifier free,  $Q = \exists$  if  $t$  is even, and  $Q = \forall$  if  $t$  is odd. The  $\text{W}$ -hierarchy is defined as the union of all the classes  $\text{W}[t]$ , that is,  $\bigcup_{t \geq 1} \text{W}[t]$ .

### Parameterized Complexity of CQ-Evaluation

We know that **p-Clique** is  $W[1]$ -complete, which means that **p-Clique**  $\in W[1]$  (this has been shown above), and every parameterized problem in  $W[1]$  can be reduced via an FPT-reduction to **p-Clique**. We also know that  $FPT \subseteq W[1]$ , and it is widely believed that this inclusion is strict (the status of the question whether  $FPT \neq W[1]$  is comparable to that of  $PTime \neq NP$ ). Thus, it is unlikely that **p-Clique**  $\in FPT$  (as  $FPT$  is closed under FPT-reductions). We use this result to prove that the same holds for **p-CQ-Evaluation**, thus providing strong evidence that this problem is not fixed-parameter tractable.

#### Theorem 14.7

**p-CQ-Evaluation** is  $W[1]$ -complete.

*Proof.* For the lower bound, we show that there exists an FPT-reduction from **p-Clique** to **p-CQ-Evaluation**. We use the same reduction as for the lower bound in Theorem 14.1 which we recall here for the sake of readability. Consider an input to **p-Clique** given by  $G = (V, E)$  and  $k \geq 1$ . The database is

$$D = \{\text{Node}(v) \mid v \in V\} \cup \{\text{Edge}(v, u) \mid \{v, u\} \in E \text{ and } v \neq u\},$$

and the Boolean CQ is

$$q = \exists x_1 \cdots \exists x_k \left( \bigwedge_{i=1}^k \text{Node}(x_i) \wedge \bigwedge_{i,j \in [k]: i \neq j} \text{Edge}(x_i, x_j) \right).$$

As discussed in the proof of Theorem 14.1,  $G$  has a clique of size  $k$  if and only if  $D \models q$ , and  $D$  and  $q$  can be constructed in polynomial time from  $G$  and  $k$ . To conclude that this is an FPT-reduction, it remains to show that the third condition in the definition of FPT-reductions holds, i.e.,  $\|q\| \leq g(k)$  for some computable function  $g : \mathbb{N} \rightarrow \mathbb{R}_0^+$ . It is easy to verify that  $\|q\| \leq c \cdot \log k \cdot k^2$  for some constant  $c \in \mathbb{R}^+$ , and thus, **p-CQ-Evaluation** is  $W[1]$ -hard.

We now focus on the upper bound. For technical clarity, we consider only constant-free Boolean CQs over a schema consisting of a single binary relation name **Edge**. We leave the prove for the general case, where no restrictions are imposed to the query and its schema, as an exercise.

We first define a universal FO sentence  $\varphi$ , and then show that there exists an FPT-reduction from **p-CQ-Evaluation** to **p-WD $_{\varphi}$** . Consider the schema

$$\mathbf{S} = \{\text{Const}[1], \text{Var}[1], \text{Edge}_1[2], \text{Edge}_2[2]\}.$$

Consider also the relation name **Hom**[2] that does not belong to  $\mathbf{S}$ . We define the universal FO sentence  $\varphi$  over  $\mathbf{S} \cup \{\text{Hom}\}$  as follows:



$$\begin{aligned}
& \forall x \forall y \forall z ((\text{Hom}(x, y) \wedge \text{Hom}(x, z)) \rightarrow y = z) \wedge \\
& \forall x \forall y (\text{Hom}(x, y) \rightarrow (\text{Var}(x) \wedge \text{Const}(y))) \wedge \\
& \forall x_1 \forall y_1 \forall x_2 \forall y_2 ((\text{Edge}_1(x_1, y_1) \wedge \text{Hom}(x_1, x_2) \wedge \text{Hom}(y_1, y_2)) \rightarrow \text{Edge}_2(x_2, y_2)).
\end{aligned}$$

We show that there is an FPT-reduction from **p-CQ-Evaluation** to **p-WD<sub>φ</sub>**. Consider an input to **p-CQ-Evaluation** given by a constant-free Boolean CQ  $q$  over the schema  $\{\text{Edge}[2]\}$ , and a database  $D$  of  $\{\text{Edge}[2]\}$ . Assuming that  $\{x_1, \dots, x_n\}$  are the variables occurring in  $q$ , we define the database  $D'$  as

$$\begin{aligned}
D' \cup & \{\text{Const}(a) \mid a \in \text{Dom}(D)\} \cup \{\text{Var}(a_{x_1}), \dots, \text{Var}(a_{x_n})\} \\
& \cup \{\text{Edge}_1(a_{x_i}, a_{x_j}) \mid \text{Edge}(x_i, x_j) \text{ is an atom occurring in } q\} \\
& \cup \{\text{Edge}_2(a, b) \mid \text{Edge}(a, b) \in D\}.
\end{aligned}$$

Roughly, the relation  $\text{Const}$  stores the constants occurring in  $D$ , the relation  $\text{Var}$  stores the variables occurring in  $q$ , the relation  $\text{Edge}_1$  stores the atoms of  $q$ , and the relation  $\text{Edge}_2$  stores the facts of  $D$ . We further define  $n = k$ , that is,  $k$  is the number of variables occurring in  $q$ .

With the definitions of  $D'$  and  $k$  in place, we can now explain the meaning of the FO sentence  $\varphi$ . The first conjunct  $\forall x \forall y \forall z ((\text{Hom}(x, y) \wedge \text{Hom}(x, z)) \rightarrow y = z)$  states that  $\text{Hom}$  represents a function, as only one value can be associated to  $x$ . The second conjunct states that  $\text{Hom}$  maps variables of  $q$  to constants of  $D$ . Finally, the third conjunct states that  $\text{Hom}$  represents a homomorphism from  $q$  to  $D$ . Notice, however, that  $\varphi$  does not impose the restriction that every variable occurring in  $q$  has to be mapped to a constant of  $D$ , as this requires a non-universal FO sentence of the form

$$\forall x (\text{Var}(x) \rightarrow \exists y (\text{Const}(y) \wedge \text{Hom}(x, y))).$$

Instead, the parameter  $k = n$  is used to force  $\text{Hom}$  to map every variable in  $q$  to a constant of  $D$ , as  $n$  is the number of variables occurring in  $q$ .

Summing up,  $q(D) = \text{true}$  if and only if there is  $S \subseteq \text{Dom}(D')^2$  with  $|S| = k$  and  $D' \models \varphi(S)$ . It is also clear that  $D'$  and  $k$  can be constructed from  $D$  and  $q$  in polynomial time, and  $k \leq \|q\|$ . Thus, we have provided an FPT-reduction from **p-CQ-Evaluation** to **p-WD<sub>φ</sub>**, which shows that **p-CQ-Evaluation**  $\in \text{W}[1]$  (for constant-free Boolean CQs over a single binary relation).  $\square$

## Containment and Equivalence

We have seen in Chapter 8 that the satisfiability problem for FO and RA is undecidable. In terms of query optimization, satisfiability is arguably the most elementary task one can think of, since it simply asks whether a query has a non-empty output on at least one database. Indeed, if a query is not satisfiable, then we do not even need to access the database in order to compute its output, which is trivially empty. Furthermore, for FO and RA, undecidability of other static analysis tasks such as containment and equivalence immediately follow from the undecidability of satisfiability.

On the other hand, the satisfiability problem for CQs is trivial. Indeed, given a CQ  $q$ , there is always a database on which  $q$  has a non-empty output, that is, the grounding  $A_q^\downarrow$  of  $A_q$  (see Definition 9.3). This means that static analysis for CQs is drastically different than for FO and RA, which in turn indicates that we need to revisit the problems of containment and equivalence in the case of CQs. This is the goal of this chapter.

### Optimizing A Simple Query

We start by first illustrating the role of containment and equivalence for CQs in query optimization by means of a simple example.

#### Example 15.1: A CQ with Redundancy

Consider again the relational schema

```
Person [ pid, pname, cid ]  
Profession [ pid, prname ]  
City [ cid, cname, country ]
```

from Chapter 3, and the CQ

$$q = \text{Answer}(y) :- \text{Person}(x, y, z), \text{Profession}(x, \text{'actor'}), \text{Profession}(x, w)$$

over this schema. The query  $q$  asks for names of persons who are actors and who have some profession. It is clear that  $q$  contains some redundancy since, if a person is an actor, then this person also has a profession (namely, being an actor). In fact, the CQ

$$q' = \text{Answer}(y) :- \text{Person}(x, y, z), \text{Profession}(x, \text{'actor'})$$

asks the same query, but in smarter way in the sense that it mentions fewer relational atoms in its body. We make two observations:

- (a) The query  $q'$  is a part of  $q$ , that is, all atoms in the body of  $q'$  belong also to the body of  $q$ .
- (b) In order to test if  $q$  and  $q'$  are equivalent, we only need to test if  $q' \subseteq q$ . The other inclusion immediately follows from (a).

The above example suggests that the following simple strategy may be useful for optimizing a CQ  $q$ . We write  $(q - R(\bar{u}))$  for the CQ obtained by deleting from the body of  $q$  the relational atom  $R(\bar{u})$ .

---

**Algorithm 3** OPTIMIZE-BY-CONTAINMENT( $q$ )

---

**Input:** A CQ  $q(\bar{x})$

**Output:** A CQ  $q^*(\bar{x})$  that is equivalent to  $q(\bar{x})$ , and may mention fewer atoms

- 1: **while** there exists an atom  $R(\bar{u})$  in the body of  $q$  such that  $(q - R(\bar{u})) \subseteq q$  **do**
  - 2:      $q := (q - R(\bar{u}))$
  - 3: **return**  $q(\bar{x})$
- 

The approach in Algorithm 3 captures a very natural idea for optimizing CQs: *keep removing atoms from the body of the CQ as long as the resulting CQ is equivalent to the original one*. In order to carry out this strategy (and numerous other, more intricate, optimization strategies), it is crucial that we are able to effectively test containment, and thus equivalence, between CQs. We therefore study in this chapter the closely related problems CQ-Containment and CQ-Equivalence. We will retake Algorithm 3 in Chapter 16.

## Containment

We first concentrate on CQ-Containment. We start by illustrating the notion of containment for CQs via a simple example.

**Example 15.2: CQ Containment**

Consider the CQ

$$q_1 = \text{Answer}(y_1) :- \text{Person}(x_1, y_1, z_1), \text{Profession}(x_1, \text{'actor'}), \\ \text{City}(z_1, \text{'Los Angeles'}, \text{'United States'})$$

asking for names of actors who live in Los Angeles, and the CQ

$$q_2 = \text{Answer}(y_2) :- \text{Person}(x_2, y_2, z_2), \text{Profession}(x_2, w_2)$$

asking for persons who have a profession. It is easy to verify that  $q_1 \subseteq q_2$  since  $q_1$  imposes the extra conditions that the returned persons are actors who live in Los Angeles.

We proceed to show that checking for containment in the case of CQs is decidable, but an intractable problem.

**Theorem 15.3**

CQ-Containment is NP-complete.

The proof of Theorem 15.3 relies on a useful characterization of containment of CQs in terms of homomorphisms, which we present below. Given two CQs  $q(\bar{x})$  and  $q'(\bar{x}')$ , we write  $(q', \bar{x}') \rightarrow (q, \bar{x})$  for the fact that there exists a homomorphism from  $(A_{q'}, \bar{x}')$  to  $(A_q, \bar{x})$ ; we also write  $q \rightarrow q'$  to indicate that  $A_q \rightarrow A_{q'}$ . Recall that  $A_q$  and  $A_{q'}$  are the sets of atoms occurring in the body of  $q$  and  $q'$ , respectively, when seen as rules.

We also remind the reader that for a set of atoms  $S$ , we write  $S^\downarrow$  for the grounding of  $S$ , which allows us to view  $S$  as a database. Such a grounding is given by the bijective homomorphism  $G_S$  from  $S$  to  $S^\downarrow$  that replaces variables in  $S$  by new constants; in particular,  $G_S(S) = S^\downarrow$ .

**Theorem 15.4: Homomorphism Theorem**

Let  $q(\bar{x})$  and  $q'(\bar{x}')$  be CQs. Then:

$$q \subseteq q' \quad \text{if and only if} \quad (q', \bar{x}') \rightarrow (q, \bar{x}).$$

*Proof.* ( $\Rightarrow$ ) Assume that  $q \subseteq q'$ . Since  $G_{A_q}$  is a homomorphism, Theorem 13.2 implies  $G_{A_q}(\bar{x}) \in q(G_{A_q}(A_q))$ . Since  $q \subseteq q'$ , we have  $G_{A_q}(\bar{x}) \in q'(G_{A_q}(A_q))$ . Applying Theorem 13.2 again, we conclude that there exists a homomorphism  $h$  from  $(A_{q'}, \bar{x}')$  to  $(G_{A_q}(A_q), G_{A_q}(\bar{x}))$ . Since  $G_{A_q}$  is bijective,  $G_{A_q}^{-1} \circ h$  is a homomorphism from  $(A_{q'}, \bar{x}')$  to  $(A_q, \bar{x})$ , as needed.

( $\Leftarrow$ ) Conversely, assume that  $(q', \bar{x}') \rightarrow (q, \bar{x})$ , and let  $h$  be a homomorphism from  $(A_{q'}, \bar{x}')$  to  $(A_q, \bar{x})$ . Given a database  $D$ , assume that  $\bar{a} \in q(D)$ . By Theorem 13.2 there exists a homomorphism  $g$  from  $(A_q, \bar{x})$  to  $(D, \bar{a})$ . Since homomorphisms compose,  $g \circ h$  is a homomorphism from  $(A_{q'}, \bar{x}')$  to  $(D, \bar{a})$  and, thus,  $\bar{a} \in q'(D)$  by Theorem 13.2. Therefore, we have that  $q(D) \subseteq q'(D)$ , from which we conclude that  $q \subseteq q'$ .  $\square$

The next example shows the usefulness of the Homomorphism Theorem.

#### Example 15.5: Homomorphism Theorem

Consider again the CQs  $q_1$  and  $q_2$  from Example 15.2 and recall that  $q_1 \subseteq q_2$ . This is confirmed by the Homomorphism Theorem since

$$(q_2, y_2) \rightarrow (q_1, y_1).$$

This is the case since the function  $h : \text{Dom}(A_{q_2}) \rightarrow \text{Dom}(A_{q_1})$  defined as

$$h(x_2) = x_1 \quad h(y_2) = y_1 \quad h(z_2) = z_1 \quad h(w_2) = \text{'actor'}$$

is a homomorphism from  $(A_{q_2}, y_2)$  to  $(A_{q_1}, y_1)$ .

An easy consequence of the Homomorphism Theorem is that the problem CQ-Containment can be reduced to CQ-Evaluation.

#### Corollary 15.6

Let  $q(\bar{x})$  and  $q'(\bar{x}')$  be CQs. Then:

$$q \subseteq q' \text{ if and only if } \mathbf{G}_{A_q}(\bar{x}) \in q'(\mathbf{G}_{A_q}(A_q)).$$

*Proof.* By Theorem 15.4 we conclude that

$$q \subseteq q' \text{ if and only if } (A_{q'}, \bar{x}') \rightarrow (A_q, \bar{x}).$$

We can also show that

$$(A_{q'}, \bar{x}') \rightarrow (A_q, \bar{x}) \text{ if and only if } (A_{q'}, \bar{x}') \rightarrow (\mathbf{G}_{A_q}(A_q), \mathbf{G}_{A_q}(\bar{x})).$$

Indeed, if  $(A_{q'}, \bar{x}') \rightarrow (A_q, \bar{x})$  is witnessed via  $h$ , then we have that  $\mathbf{G}_{A_q} \circ h$  is a homomorphism from  $(A_{q'}, \bar{x}')$  to  $(\mathbf{G}_{A_q}(A_q), \mathbf{G}_{A_q}(\bar{x}))$ . Conversely, assuming that  $(A_{q'}, \bar{x}') \rightarrow (\mathbf{G}_{A_q}(A_q), \mathbf{G}_{A_q}(\bar{x}))$  is witnessed via  $g$ ,  $\mathbf{G}_{A_q}^{-1} \circ g$  is a homomorphism from  $(A_{q'}, \bar{x}')$  to  $(A_q, \bar{x})$ . By Theorem 13.2 we get that

$$(A_{q'}, \bar{x}') \rightarrow (\mathbf{G}_{A_q}(A_q), \mathbf{G}_{A_q}(\bar{x})) \text{ if and only if } \mathbf{G}_{A_q}(\bar{x}) \in q'(\mathbf{G}_{A_q}(A_q)).$$

Consequently, we get that  $q \subseteq q'$  if and only if  $\mathbf{G}_{A_q}(\bar{x}) \in q'(\mathbf{G}_{A_q}(A_q))$ .  $\square$

By exploiting the Homomorphism Theorem, we can further show that the problem **CQ-Evaluation** can be reduced to **CQ-Containment**, i.e., the opposite of what Corollary 15.6 shows. In the proof of Corollary 15.6, we essentially convert the CQ  $q$  into a database via the bijective homomorphism  $G_{A_q}$ . Now we are going to do the opposite, i.e., convert a database into a CQ. As discussed in Chapter 13, we can convert a database  $D$  into a set of relational atoms via the injective function  $V_C : \text{Const} \rightarrow \text{Const} \cup \text{Var}$ , where  $C$  is a finite set of constants. Recall that  $V_C(D)$  is the set of relational atoms obtained from  $D$  by replacing constants, except for those in  $C$ , with variables. The following corollary, which establishes that **CQ-Evaluation** can be reduced to **CQ-Containment**, is stated for Boolean CQs, as this suffices for the purpose of pinpointing the complexity of **CQ-Containment**, but it can be easily generalized to arbitrary CQs.

#### Corollary 15.7

Let  $q$  be a Boolean CQ,  $D$  a database, and  $q_D$  the Boolean CQ such that  $A_{q_D} = V_C(D)$ , where  $C = \text{Dom}(A_q) \cap \text{Const}$ . Then:

$$D \models q \text{ if and only if } q_D \subseteq q.$$

*Proof.* By Theorem 13.2, we conclude that

$$D \models q \text{ if and only if } q \rightarrow D.$$

It is easy to show that

$$q \rightarrow D \text{ if and only if } q \rightarrow q_D.$$

Indeed, if  $q \rightarrow D$  is witnessed via  $h$ , then we get that  $V_C \circ h$  is a homomorphism from  $q$  to  $q_D$ . Conversely, assuming that  $q \rightarrow q_D$  is witnessed via  $g$ ,  $V_C^{-1} \circ g$  is a homomorphism from  $q$  to  $D$ . By Theorem 15.4, we conclude that

$$q \rightarrow q_D \text{ if and only if } q_D \subseteq q.$$

From the above equivalences, we get that  $D \models q$  if and only if  $q_D \subseteq q$ .  $\square$

By Theorem 14.1, **CQ-Evaluation** is in NP, and thus, Corollary 15.6 implies that also **CQ-Containment** is in NP. Moreover, since **CQ-Evaluation** is NP-hard even for Boolean CQs (this is because the CQ that the reduction from **Clique** to **CQ-Evaluation** builds in the proof of Theorem 14.1 is Boolean), Corollary 15.7 implies that **CQ-Containment** is NP-hard. Therefore, **CQ-Containment** is NP-complete, which establishes Theorem 15.3.

## Equivalence

We now focus on the equivalence problem: given two CQs  $q, q'$ , check whether  $q \equiv q'$ , i.e., whether  $q(D) = q'(D)$  for every database  $D$ . We show that:

**Theorem 15.8**

CQ-Equivalence is NP-complete.

*Proof.* Concerning the upper bound, it suffices to observe that

$$q \equiv q' \text{ if and only if } q \subseteq q' \text{ and } q' \subseteq q,$$

which implies that CQ-Equivalence is in NP since, by Theorem 15.3 the problem of deciding whether  $q \subseteq q'$  and  $q' \subseteq q$  is in NP.

Concerning the lower bound, we provide a reduction from CQ-Containment. In fact, CQ-Containment is NP-hard even if we consider Boolean CQs (this is a consequence of the proof of Theorem 15.3). Consider two Boolean CQs

$$q = \text{Answer} :- R_1(\bar{u}_1), \dots, R_n(\bar{u}_n) \quad q' = \text{Answer} :- R'_1(\bar{u}'_1), \dots, R'_m(\bar{u}'_m),$$

We assume that  $q, q'$  do not share variables since we can always rename variables without affecting the semantics of a query. Let  $q_\cap$  be the Boolean CQ

$$\text{Answer} :- R_1(\bar{u}_1), \dots, R_n(\bar{u}_n), R'_1(\bar{u}'_1), \dots, R'_m(\bar{u}'_m),$$

which essentially computes the intersection of  $q$  and  $q'$ . In other words, for every database  $D$ ,  $q(D) \cap q'(D) = q_\cap(D)$ . It is straightforward to see that

$$q \subseteq q' \text{ if and only if } q \equiv q_\cap,$$

which in turn implies that CQ-Equivalence is NP-hard, as needed.  $\square$

## Minimization

Query optimization is the task of transforming a query into an equivalent one that is easier to evaluate. Since joins are expensive operations, we typically consider an equivalent version of a CQ  $q$  with fewer atoms in its body, and thus, with fewer joins to perform. Ideally, we would like to compute a CQ  $q'$  that is equivalent to  $q$ , and is also minimal, i.e., it has the minimum number of atoms. This brings us to the notion of minimization of CQs.

### Definition 16.1: Minimization of CQs

Consider a CQ  $q$  over a schema  $\mathbf{S}$ . A CQ  $q'$  over  $\mathbf{S}$  is a *minimization* of  $q$  if the following hold:

1.  $q \equiv q'$ , and
2. for every CQ  $q''$  over  $\mathbf{S}$ ,  $q' \equiv q''$  implies  $|A_{q'}| \leq |A_{q''}|$ .

In other words,  $q'$  is a minimization of  $q$  if it is equivalent to  $q$  and has the smallest number of atoms among all the CQs that are equivalent to  $q$ . It is straightforward to see that every CQ  $q$  over a schema  $\mathbf{S}$  has a minimization, which is actually a query from the finite set (up to variable renaming)

$$M_q = \{q' \mid q' \text{ is a CQ over } \mathbf{S} \text{ and } |A_{q'}| \leq |A_q|\}$$

that collects all the CQs over  $\mathbf{S}$  (up to variable renaming) with at most  $|A_q|$  atoms. Hence, to compute a minimization of  $q$ , we could, e.g., iterate over all CQs of  $M_q$  in increasing order with respect to the number of body atoms, until we find one that is equivalent to  $q$ . But now the following questions arise:

1. Is there a smarter procedure for computing a minimization of  $q$  instead of naively iterating over the exponentially many CQs of  $M_q$ ? In particular, does the strategy of removing atoms from  $q$  as long as the resulting query is equivalent to  $q$  (see Algorithm 3) lead to a minimization of  $q$ ?



2. Which minimization of  $q$  should be computed? Is there one that stands out as the best?

The above questions have neat answers, which we discuss in detail in the rest of the chapter. In a nutshell, one can indeed find minimizations of a CQ  $q$  by removing atoms from its body. Moreover, although  $q$  may have several minimizations, they are all the same (up to variable renaming). This implies that no matter in which order we remove atoms from the body of  $q$ , we will always compute the same minimization of  $q$  (up to variable renaming).

#### Minimization via Atom Removals

Consider a CQ  $q$  of the form  $\text{Answer}(\bar{x}) :- R_1(\bar{u}_1), \dots, R_n(\bar{u}_n)$ . The CQ  $q'$  obtained from  $q$  by removing the atom  $R_i(\bar{u}_i)$ , for some  $i \in [n]$ , is

$$\text{Answer}(\bar{x}') :- R_1(\bar{u}_1), \dots, R_{i-1}(\bar{u}_{i-1}), R_{i+1}(\bar{u}_{i+1}), \dots, R_n(\bar{u}_n),$$

where  $\bar{x}'$  is obtained from  $\bar{x}$  by removing every variable that is only mentioned in the atom  $R_i(\bar{u}_i)$ . For example, if we remove the atom  $R(x)$  from the CQ  $\text{Answer}(x, y) :- R(x), S(y)$ , then we obtain the CQ  $\text{Answer}(y) :- S(y)$  as the variable  $x$  is only mentioned in  $R(x)$ . On the other hand, if we remove the atom  $R(x)$  from the CQ  $\text{Answer}(x, y) :- R(x), T(x, y)$ , then we obtain the CQ  $\text{Answer}(x, y) :- T(x, y)$  since  $x$  occurs also in  $T(x, y)$ .

The building block of minimization via atom removals is as follows: given a CQ  $q(\bar{x})$ , construct a CQ  $q'(\bar{x})$  by removing an atom  $R(\bar{u})$  from the body of  $q$  such that  $(q, \bar{x}) \rightarrow (q', \bar{x})$ . Notice that the output tuple  $\bar{x}$  remains the same, which means that the atom  $R(\bar{u})$  either it does not contain a variable of  $\bar{x}$ , or it contains only variables of  $\bar{x}$  that occur also in atoms of  $A_q - \{R(\bar{u})\}$ . In this way, we actually construct a CQ that is equivalent to  $q$ . Indeed, since  $(q, \bar{x}) \rightarrow (q', \bar{x})$ , we get that  $q' \subseteq q$  (by Theorem 15.4). Moreover,  $(q', \bar{x}) \rightarrow (q, \bar{x})$  holds trivially due to the identity homomorphism from  $A_{q'}$  to  $A_q$ , and thus,  $q \subseteq q'$  (again by Theorem 15.4). We then iteratively remove atoms as above until we reach a CQ  $q''(\bar{x})$  that is minimal, i.e., any CQ  $q'''(\bar{x})$  that can be obtained by removing an atom from the body of  $q''$  is such that  $(q'', \bar{x}) \rightarrow (q''', \bar{x})$  does not hold. The CQ  $q''$  is typically called a *core* of  $q$ . The formal definition follows.

#### Definition 16.2: Core of a CQ

Consider a CQ  $q(\bar{x})$ . A CQ  $q'(\bar{x})$  is a *core* of  $q$  if the following hold:

1.  $A_{q'} \subseteq A_q$ ,
2.  $(q, \bar{x}) \rightarrow (q', \bar{x})$ , and
3. for every CQ  $q''(\bar{x})$  with  $A_{q''} \subsetneq A_{q'}$ ,  $(q', \bar{x}) \rightarrow (q'', \bar{x})$  does not hold.

The first condition in Definition 16.2 expresses that either  $q = q'$ , or  $q'$  is obtained by removing atoms from  $q$  but without altering the output tuple  $\bar{x}$ ,

the second condition ensures that  $q \equiv q'$ , and the third condition states that  $q'$  is minimal. Here is an example that illustrates the notion of core of a CQ.

**Example 16.3: Core of a CQ**

Consider the Boolean CQ  $q_1$  defined as

$$\text{Answer} \text{ :- } R(x, y), R(x, z).$$

The function  $h$  defined as  $h(x) = x$ ,  $h(y) = y$  and  $h(z) = y$  is a homomorphism from  $\{R(x, y), R(x, z)\}$  to  $\{R(x, y)\}$ . Therefore,  $q_1 \rightarrow q'_1$ , where  $q'_1$  is the Boolean CQ defined as

$$\text{Answer} \text{ :- } R(x, y).$$

Since, by definition, a CQ must have at least one atom in its body, we conclude that  $q'_1$  is a core of  $q_1$ . Observe that the Boolean CQ  $q''_1$

$$\text{Answer} \text{ :- } R(x, z)$$

is also a core of  $q_1$  due to the homomorphism  $h'$  defined as  $h(x) = x$ ,  $h(y) = z$  and  $h(z) = z$ . Therefore, a CQ may have several cores that are syntactically different, depending on the order that atoms are removed.

Consider now the Boolean CQ  $q_2$  defined as

$$\text{Answer} \text{ :- } R(x, y), R(y, z).$$

Observe that there is neither a homomorphism from  $\{R(x, y), R(y, z)\}$  to  $\{R(x, y)\}$ , nor a homomorphism from  $\{R(x, y), R(y, z)\}$  to  $\{R(y, z)\}$ . This means that there is no way to remove an atom from  $q_2$  and get an equivalent CQ. Therefore, we conclude that  $q_2$  is its own core.

Finally, consider the CQ  $q_3$  defined as

$$\text{Answer}(x, y, z) \text{ :- } R(x, y), R(x, z),$$

which is actually  $q_1$  with all the variables in the output tuple. By removing the atom  $R(x, z)$  from  $q_3$ , we obtain the CQ  $q'_3$

$$\text{Answer}(x, y) \text{ :- } R(x, y).$$

In this case, there is no homomorphism from  $(A_{q_3}, (x, y, z))$  to  $(A_{q'_3}, (x, y))$  since there is no way to map the ternary tuple  $(x, y, z)$  to the binary tuple  $(x, y)$ . Hence,  $q'_3$  is not equivalent to  $q_3$ . The case where we remove the atom  $R(x, y)$  from  $q_3$  is analogous. Therefore,  $q_3$  is its own core.

We proceed to show that the notion of core captures our original intention, that is, the construction of a minimization of a CQ.

**Proposition 16.4**

Every CQ  $q$  has at least one core, and every core of  $q$  is a minimization of  $q$ .

*Proof.* We first show that a CQ  $q(\bar{x})$  has a core. If  $q$  is a core of itself, then the claim follows. Assume now that this is not the case. This means that condition (3) in the definition of core (Definition 16.2) is violated, which in turn implies that there is a CQ  $q'(\bar{x})$  with  $A_{q'} \subsetneq A_q$  such that  $(q, \bar{x}) \rightarrow (q', \bar{x})$ . If  $q'$  is a core of itself, then it is clear that  $q'$  is a core of  $q$ . Otherwise, we iteratively apply the above argument until we reach a core of  $q$ .

We now proceed to show that a core of  $q(\bar{x})$  is a minimization of it. We first show a useful technical lemma:

**Lemma 16.5.** *Consider a CQ  $q_1(\bar{y}_1)$ , and assume that there is a CQ  $q_2(\bar{y}_2)$  such that  $q_1 \equiv q_2$  and  $|A_{q_2}| < |A_{q_1}|$ . Then, there is a CQ  $q_3(\bar{y}_1)$  such that*

$$(q_1, \bar{y}_1) \rightarrow (q_3, \bar{y}_1) \quad \text{and} \quad A_{q_3} \subsetneq A_{q_1}.$$

*Proof.* By Theorem 15.4, we conclude that

$$(q_1, \bar{y}_1) \rightarrow (q_2, \bar{y}_2) \quad \text{and} \quad (q_2, \bar{y}_2) \rightarrow (q_1, \bar{y}_1).$$

Assume that these statements are witnessed via the homomorphisms  $h_1$  and  $h_2$ , respectively. Let  $q_3(\bar{y}_3)$  be the CQ such that

$$A_{q_3} = h_2(A_{q_2}) \quad \text{and} \quad \bar{y}_3 = h_2(\bar{y}_2).$$

It is clear that  $\bar{y}_3 = \bar{y}_1$  and  $A_{q_3} \subseteq A_{q_1}$ . Furthermore, since  $|A_{q_3}| \leq |A_{q_2}|$  and  $|A_{q_2}| < |A_{q_1}|$ , we conclude that  $|A_{q_3}| < |A_{q_1}|$ , and thus,  $A_{q_3} \subsetneq A_{q_1}$ . It remains to show that  $(q_1, \bar{y}_1) \rightarrow (q_3, \bar{y}_1)$ . Since homomorphisms compose, the latter is witnessed via the homomorphism  $h_2 \circ h_1$ .  $\square$

Consider now a CQ  $q'(\bar{x})$  that is a core of  $q(\bar{x})$ . Towards a contradiction, assume that  $q'$  is not a minimization of  $q$ . This implies that there exists a CQ  $q''$  such that  $q' \equiv q''$  and  $|A_{q''}| < |A_{q'}|$ . By Lemma 16.5, we conclude that there exists a CQ  $q'''(\bar{x})$  such that  $(q', \bar{x}) \rightarrow (q''', \bar{x})$  and  $A_{q'''} \subsetneq A_{q'}$ . This contradicts our hypothesis that  $q'$  is a core of  $q$ , and the claim follows.  $\square$

By Proposition 16.4, to compute a minimization of a CQ  $q$ , we simply need to compute a core of it. This can be done via the simple iterative procedure COMPUTECORE, given in Algorithm 4. Notice that this algorithm is a more detailed reformulation of Algorithm 3. It is straightforward to show that, for a CQ  $q$ , COMPUTECORE( $q$ ) terminates after finitely many steps. It is also not difficult to show that the procedure COMPUTECORE is correct.

**Lemma 16.6.** *Given a CQ  $q$ , COMPUTECORE( $q$ ) is a core of  $q$ .*

**Algorithm 4** COMPUTECORE( $q$ )**Input:** A CQ  $q(\bar{x})$ **Output:** A CQ  $q^*(\bar{x})$  that is a core of  $q(\bar{x})$ 


---

```

1:  $S := A_q$ 
2: while there exists  $R(\bar{u}) \in S$  such that each variable in  $\bar{x}$ 
3:   occurs in  $\text{Dom}(S - \{R(\bar{u})\})$  and  $(S, \bar{x}) \rightarrow (S - \{R(\bar{u})\}, \bar{x})$  do
4:    $S := S - \{R(\bar{u})\}$ 
5: return  $q^*(\bar{x}) := R_1(\bar{u}_1), \dots, R_n(\bar{u}_n)$ , where  $S = \{R_1(\bar{u}_1), \dots, R_n(\bar{u}_n)\}$ 

```

---

*Proof.* At each iteration of the while-loop, the CQ  $q'(\bar{x})$  with  $A_{q'} = S$  (which is indeed a CQ since, by construction, every variable in  $\bar{x}$  occurs in  $A_{q'}$ ) is such that  $A_{q'} \subseteq A_q$  and  $(q, \bar{x}) \rightarrow (q', \bar{x})$ . Therefore, the CQ  $q^*(\bar{x})$  returned by the algorithm is such that  $A_{q^*} \subseteq A_q$  and  $(q, \bar{x}) \rightarrow (q^*, \bar{x})$ . Furthermore, by construction, for every CQ  $q''(\bar{x})$  with  $A_{q''} \subsetneq A_{q^*}$ ,  $(q^*, \bar{x}) \rightarrow (q'', \bar{x})$  does not hold. Therefore,  $q^*$  satisfies all the three conditions given in the definition of core (Definition 16.2), and thus, it is a core of  $q$ , as needed.  $\square$

Note that COMPUTECORE is a nondeterministic algorithm. Observe that there may be several atoms  $R(\bar{y}) \in S$  satisfying the condition of the while loop (in particular, the condition  $(S, \bar{x}) \rightarrow (S - \{R(\bar{y})\}, \bar{x})$ ), but we do not specify how such an atom is selected. In fact, the atom  $R(\bar{y})$  of  $S$  that is eventually removed from  $S$  at step 4 is chosen nondeterministically. Therefore, the final result computed by the algorithm depends on how the atoms to be removed from  $S$  are chosen, and thus, different executions of COMPUTECORE( $q$ ) may compute cores of  $q$  that are syntactically different. This fact should not be surprising as it has been already illustrated in Example 16.3 (see the queries  $q'_1$  and  $q''_1$  that are cores of  $q_1$ ). This leads to the second main question raised above: is there a core of  $q$  that stands out as the best?

*Uniqueness of Minimizations*

It turns out that such a concept as the best core does not exist since a CQ has a *unique core* (up to variable renaming). This is a consequence of the fact that every CQ has a *unique minimization* (up to variable renaming). We proceed to show the latter statement.

We say that two CQs  $q(\bar{x}), q'(\bar{x}')$  are *isomorphic* if one can be turned into the other via renaming of variables, i.e., if there is a bijection  $\rho : \text{Dom}(A_q) \rightarrow \text{Dom}(A_{q'})$  that is a homomorphism from  $(A_q, \bar{x})$  to  $(A_{q'}, \bar{x}')$ , and its inverse  $\rho^{-1}$  is a homomorphism from  $(A_{q'}, \bar{x}')$  to  $(A_q, \bar{x})$ . (Recall from Chapter 9 that homomorphisms between sets of atoms are always the identity on constants.)

**Proposition 16.7**

Consider a CQ  $q(\bar{x})$ , and let  $q'(\bar{x}')$  and  $q''(\bar{x}'')$  be minimizations of  $q$ . Then  $q'$  and  $q''$  are isomorphic.

*Proof.* We need to show that there is a bijection  $\rho : \text{Dom}(A_{q'}) \rightarrow \text{Dom}(A_{q''})$  that is a homomorphism from  $(A_{q'}, \bar{x}')$  to  $(A_{q''}, \bar{x}'')$ , and its inverse  $\rho^{-1}$  is a homomorphism from  $(A_{q''}, \bar{x}'')$  to  $(A_{q'}, \bar{x}')$ . Since both  $q'$  and  $q''$  are minimizations of  $q$ , we get that  $q \equiv q'$  and  $q \equiv q''$ , and thus,  $q' \equiv q''$ . By Theorem 15.4

$$(q', \bar{x}') \rightarrow (q'', \bar{x}'') \quad \text{and} \quad (q'', \bar{x}'') \rightarrow (q', \bar{x}').$$

Assume that these statements are witnessed via the homomorphisms  $h$  and  $g$ , respectively. We proceed to show a useful statement concerning  $h$  and  $g$ :

**Lemma 16.8.** *The functions  $h$  and  $g$  are bijections.*

*Proof.* We concentrate on  $h$ , and show that it is both surjective and injective; the proof for  $g$  is analogous. We give a proof by contradiction:

- Assume first that  $h$  is not surjective. This implies that there is a variable  $z \in \text{Dom}(A_{q''})$  such that there is no variable  $y \in \text{Dom}(A_{q'})$  with  $h(y) = z$ . Let  $R(\bar{u}) \in A_{q''}$  be an atom that mentions  $z$ . We have that  $R(\bar{u}) \notin h(A_{q'})$ . We define  $q'''(\bar{x}''')$  as the CQ with  $A_{q'''} = h(A_{q'})$ . It is clear that  $(q', \bar{x}') \rightarrow (q''', \bar{x}''')$  via  $h$ , and  $(q''', \bar{x}''') \rightarrow (q'', \bar{x}'')$  via  $g$ . Therefore, by Theorem 15.4,  $q' \equiv q'''$ . Since  $q' \equiv q''$ , we conclude that  $q'' \equiv q'''$ . Observe also that  $A_{q'''} \subsetneq A_{q''}$ , which implies that  $|A_{q'''}| < |A_{q''}|$ . But this contradicts the fact that  $q''$  is a minimization of  $q$ , and thus,  $h$  is surjective.
- Assume now that  $h$  is not injective. This implies that there are two distinct variables  $y, z \in \text{Dom}(A_{q'})$  such that  $h(y) = h(z)$ . Hence,  $g(h(y)) = g(h(z))$ , which implies that  $g \circ h$  is a homomorphism from  $(q', \bar{x}') \rightarrow (q'', \bar{x}'')$  that is not surjective. Therefore, there exists a variable  $u \in \text{Dom}(A_{q''})$  such that there is no variable  $v \in \text{Dom}(A_{q'})$  with  $g(h(v)) = u$ . Let  $R(\bar{u}) \in A_{q''}$  be an atom that mentions  $u$ . We have that  $R(\bar{u}) \notin g(h(A_{q'}))$ . We define  $q'''(\bar{x}''')$  as the CQ with  $A_{q'''} = g(h(A_{q'}))$ . It is clear that  $(q', \bar{x}') \rightarrow (q''', \bar{x}''')$  via  $g \circ h$ . Observe also that  $A_{q'''} \subsetneq A_{q''}$ . Hence,  $(q''', \bar{x}''') \rightarrow (q'', \bar{x}'')$  via the identity homomorphism, which means that  $q' \equiv q'''$  due to Theorem 15.4 and  $|A_{q'''}| < |A_{q''}|$ . But this contradicts the fact that  $q'$  is a minimization of  $q$ , which in turn implies that  $h$  is injective.

Since  $h$  is both surjective and injective, the claim follows.  $\square$

We are now ready to define the bijection  $\rho : \text{Dom}(A_{q'}) \rightarrow \text{Dom}(A_{q''})$ . Let  $f = g \circ h$ . It is clear that  $f$  is a homomorphism from  $(A_{q'}, \bar{x}')$  to  $(A_{q'}, \bar{x}')$ . Since, by Lemma 16.8, both  $h$  and  $g$  are bijections, we can further conclude that  $f$  is a bijection. This implies that there exists  $k \geq 0$  such that the function

$$f^k = \underbrace{f \circ \dots \circ f}_k$$

is the identity homomorphism from  $(A_{q'}, \bar{x}')$  to  $(A_{q'}, \bar{x}')$ . Let  $\rho = h \circ f^{k-1}$ . Since both  $h$  and  $f^{k-1}$  are bijections, we get that also  $\rho$  is a bijection. It is also clear that  $\rho$  is a homomorphism from  $(A_{q'}, \bar{x}')$  to  $(A_{q''}, \bar{x}'')$ . Notice also that  $g \circ \rho = f^k$  is the identity, which means that  $g$  is the inverse of  $\rho$ . Thus, the inverse of  $\rho$  is a homomorphism from  $(A_{q''}, \bar{x}'')$  to  $(A_{q'}, \bar{x}')$ . Therefore,  $\rho$  witnesses the fact that  $q'$  and  $q''$  are isomorphic, and the claim follows.  $\square$

From Proposition 16.4, which tells us that a core of a CQ  $q$  is a minimization of  $q$ , and Proposition 16.7, we immediately get the following corollary:

**Corollary 16.9**

Consider a CQ  $q$ , and let  $q'$  and  $q''$  be cores of  $q$ . It holds that  $q'$  and  $q''$  are isomorphic.

Recall that different executions of the nondeterministic procedure COMPUTECORE on some input CQ  $q$ , may compute cores of  $q$  that are syntactically different. However, Corollary 16.9 tells us that those cores differ only on the names of their variables. In other words, cores of  $q$  computed by different executions of COMPUTECORE( $q$ ) are actually the same up to variable renaming.



## Containment Under Integrity Constraints

As discussed in Chapters 10 and 11, relational systems support the specification of semantic properties that should be satisfied by all databases of a certain schema. This is achieved via integrity constraints, also called dependencies. The question that arises is how static analysis, and in particular the notion of containment of CQs, studied in Chapter 16, is affected in the presence of constraints. In this chapter, we study this question concentrating on functional dependencies (FDs) and inclusion dependencies (INDs).

### Functional Dependencies

We start with FDs, and illustrate via an example how containment of CQs is affected if we focus on databases that satisfy a given set of FDs.

#### Example 17.1: Containment of CQs Under FDs

Consider the CQs  $q_1$  and  $q_2$  defined as

$$\begin{aligned}\text{Answer}(x_1, y_1) &:- R(x_1, y_1), R(y_1, z_1), R(x_1, z_1) \\ \text{Answer}(x_2, y_2) &:- R(x_2, y_2), R(y_2, y_2),\end{aligned}$$

respectively. It is easy to verify that  $(q_2, (x_2, y_2)) \rightarrow (q_1, (x_1, y_1))$  does not hold, and thus, we have that  $q_1 \not\subseteq q_2$  by the Homomorphism Theorem. For example, if we consider the database

$$D = \{R(1, 2), R(2, 3), R(1, 3)\},$$

then  $q_1(D) = \{(1, 2)\}$  and  $q_2(D) = \emptyset$ , so that  $q_1(D) \not\subseteq q_2(D)$ . Suppose now that  $q_1, q_2$  will be evaluated only over databases that satisfy the FD

$$\sigma = R : \{1\} \rightarrow \{2\}.$$



In particular,  $q_1$  and  $q_2$  will not be evaluated over the database  $D$  since it does not satisfy  $\sigma$ . We can show that, for every database  $D'$ ,

$$D' \models \sigma \text{ implies } q_1(D') \subseteq q_2(D').$$

To see this, consider an arbitrary database  $D'$  that satisfies  $\sigma$ , and assume that  $(a, b) \in q_1(D')$ . By Theorem 13.2, we have that

$$(q_1, (x_1, y_1)) \rightarrow (D', (a, b))$$

via a homomorphism  $h_1$ . Since  $D' \models \sigma$  and

$$\{R(h_1(x_1), h_1(y_1)), R(h_1(x_1), h_1(z_1))\} \subseteq D',$$

it holds that  $h_1(y_1) = h_1(z_1)$ . Since  $R(h_1(y_1), h_1(z_1)) \in D'$ , we get that

$$(q_2, (x_2, y_2)) \rightarrow (D', (a, b))$$

via  $h_2$  such that  $h_2(x_2) = h_1(x_1)$  and  $h_2(y_2) = h_1(y_1) = h_1(z_1)$ .

Our goal is to revisit the problem of containment for CQs in the presence of FDs. More precisely, given two CQs  $q$  and  $q'$ , and a set  $\Sigma$  of FDs, we say that  $q$  is *contained in  $q'$  under  $\Sigma$* , denoted by  $q \subseteq_{\Sigma} q'$ , if for every database  $D$  that satisfies  $\Sigma$ , it holds that  $q(D) \subseteq q'(D)$ . The problem of interest follows:

**Problem: CQ-Containment-FD**

**Input:** Two CQs  $q$  and  $q'$ , and a set  $\Sigma$  of FDs

**Output:** **true** if  $q \subseteq_{\Sigma} q'$ , and **false** otherwise

We proceed to show the following result:

**Theorem 17.2**

CQ-Containment-FD is NP-complete.

It is clear that the NP-hardness is inherited from CQ containment without constraints (see Theorem 15.3). Recall that, by the Homomorphism Theorem, checking whether a CQ  $q(\bar{x})$  is contained in a CQ  $q'(\bar{x}')$  in the absence of constraints boils down to checking whether  $(q', \bar{x}') \rightarrow (q, \bar{x})$ . Even though this is not enough in the presence of FDs, we can adopt a similar approach providing that we first transform, by identifying terms as dictated by the FDs, the set of atoms  $A_q$  in  $q$  into a new set of atoms  $S$  that satisfies the FDs, and the tuple of variables  $\bar{x}$  into a new tuple  $\bar{u}$ , which may contain also constants, and then check whether  $(A_{q'}, \bar{x}') \rightarrow (S, \bar{u})$ . This simple idea has

been already illustrated by Example 17.1. Unsurprisingly, the transformation of  $A_q$  and  $\bar{x}$  into  $S$  and  $\bar{u}$ , respectively, can be done by exploiting the chase for FDs, which has been introduced in Chapter 10. For brevity, we simply write  $\text{Chase}(q, \Sigma)$  instead of  $\text{Chase}(A_q, \Sigma)$ , and  $h_{q, \Sigma}$  instead of  $h_{A_q, \Sigma}$ . We now show the following result by providing a proof similar to that of the Homomorphism Theorem:

**Theorem 17.3**

Let  $q(\bar{x})$  and  $q'(\bar{x}')$  be CQs over a schema  $\mathbf{S}$ , and  $\Sigma$  a set of FDs over  $\mathbf{S}$ . The following are equivalent:

1.  $q \subseteq_{\Sigma} q'$ .
2.  $\text{Chase}(q, \Sigma) \neq \perp$  implies  $(A_{q'}, \bar{x}') \rightarrow (\text{Chase}(q, \Sigma), h_{q, \Sigma}(\bar{x}))$ .

*Proof.* For brevity, let  $S = \text{Chase}(q, \Sigma)$  and  $\bar{u} = h_{q, \Sigma}(\bar{x})$ .

We first show that (1) implies (2). By hypothesis,  $q \subseteq_{\Sigma} q'$ . It is clear that, if  $S \neq \perp$ , then  $\mathbf{G}_S(\bar{u}) \in q(\mathbf{G}_S(S))$ . Since, by Lemma 10.8,  $S \models \Sigma$ , which means that  $\mathbf{G}_S(S) \models \Sigma$ , we have that  $\mathbf{G}_S(\bar{u}) \in q'(\mathbf{G}_S(S))$ . By Theorem 13.2, there exists a homomorphism  $h$  from  $(A_{q'}, \bar{x}')$  to  $(\mathbf{G}_S(S), \mathbf{G}_S(\bar{u}))$ . Clearly,  $\mathbf{G}_S^{-1} \circ h$  is a homomorphism from  $(A_{q'}, \bar{x}')$  to  $(S, \bar{u})$ , as needed.

For showing that (2) implies (1) we proceed by case analysis:

- Assume first that  $S = \perp$ . This implies that, for every database  $D$  of  $\mathbf{S}$  such that  $D \models \Sigma$ , there is no homomorphism from  $q$  to  $D$ ; otherwise, there is a successful finite chase sequence of  $q$  under  $\Sigma$ , which contradicts the fact that  $S = \perp$ . Therefore, for every database  $D$  of  $\mathbf{S}$  such that  $D \models \Sigma$ ,  $q(D) = \emptyset$ , which in turn implies that  $q \subseteq_{\Sigma} q'$ .
- Assume now that  $S \neq \perp$ . By hypothesis, we get that  $(A_{q'}, \bar{x}') \rightarrow (S, \bar{u})$  via a homomorphism  $h$ . Let  $D$  be an arbitrary database of  $\mathbf{S}$  such that  $D \models \Sigma$ , and assume that  $\bar{a} \in q(D)$ . By Theorem 13.2,  $(q, \bar{x}) \rightarrow (D, \bar{a})$ . Since  $D \models \Sigma$ , Lemma 10.11 implies that  $(S, \bar{u}) \rightarrow (D, \bar{a})$  via a homomorphism  $g$ . Since homomorphisms compose,  $g \circ h$  is a homomorphism from  $(q', \bar{x}')$  to  $(D, \bar{a})$ . By Theorem 13.2,  $\bar{a} \in q'(D)$ , which implies that  $q \subseteq_{\Sigma} q'$ .

Since in both cases we get that  $q \subseteq_{\Sigma} q'$ , the claim follows.  $\square$

The following is an easy consequence of Theorem 17.3 and Theorem 13.2

**Corollary 17.4**

Let  $q(\bar{x})$  and  $q'(\bar{x}')$  be CQs over a schema  $\mathbf{S}$ , and  $\Sigma$  a set of FDs over  $\mathbf{S}$ . With  $S = \text{Chase}(q, \Sigma)$ , the following are equivalent:

1.  $q \subseteq_{\Sigma} q'$ .
2.  $S \neq \perp$  implies  $G_S(h_{q,\Sigma}(\bar{x})) \in q'(G_S(S))$ .

By Lemma 10.10,  $\text{Chase}(q, \Sigma)$  can be computed in polynomial time. Moreover, if  $\text{Chase}(q, \Sigma) \neq \perp$ , then the chase homomorphism  $h_{q,\Sigma}$  can be also computed in polynomial time. Since CQ-Evaluation is in NP (see Theorem 14.1), we conclude that CQ-Containment-FD is also in NP, and Theorem 17.2 follows.

## Inclusion Dependencies

We now focus on INDs. We first illustrate via an example how containment of CQs is affected if we focus on databases that satisfy a set of INDs.

### Example 17.5: Containment of CQs Under INDs

Consider the CQs  $q_1$  and  $q_2$  defined as

$$\begin{aligned} \text{Answer}(x_1, y_1) &:- R(x_1, y_1), R(y_1, z_1), P(z_1, y_1) \\ \text{Answer}(x_2, y_2) &:- R(x_2, y_2), R(y_2, z_2), S(x_2, y_2, z_2), \end{aligned}$$

respectively. It is clear that  $(q_2, (x_2, y_2)) \rightarrow (q_1, (x_1, y_1))$  does not hold, and thus, we have that  $q_1 \not\subseteq q_2$  by the Homomorphism Theorem. Suppose now that  $q_1$  and  $q_2$  will be evaluated only over databases that satisfy

$$\sigma_1 = R[1, 2] \subseteq S[1, 2] \quad \text{and} \quad \sigma_2 = S[2, 3] \subseteq R[1, 2].$$

We can show that, for every database  $D$ ,

$$D \models \{\sigma_1, \sigma_2\} \text{ implies } q_1(D) \subseteq q_2(D).$$

Consider an arbitrary database  $D$  that satisfies  $\{\sigma_1, \sigma_2\}$ , and assume that  $(a, b) \in q_1(D)$ , or, equivalently,  $(q_1, (x_1, y_1)) \rightarrow (D, (a, b))$  via a homomorphism  $h_1$ . This implies that  $R(h_1(x_1), h_1(y_1)) \in D$ . Since  $D \models \sigma_1$ , we get that  $D$  contains an atom of the form  $S(h_1(x_1), h(y_1), c)$ . But since  $D \models \sigma_2$ , we also get that  $D$  contains the atom  $R(h_1(y_1), c)$ . Hence,

$$\{R(h_1(x_1), h_1(y_1)), R(h_1(y_1), c), S(h_1(x_1), h_1(y_1), c)\} \subseteq D.$$

This implies that  $(q'_1, (x_1, y_1)) \rightarrow (D, (a, b))$ , where  $q'_1$  is obtained from  $q_1$  by adding certain atoms according to  $\sigma_1$  and  $\sigma_2$ , i.e.,  $q'_1$  is defined as

$$\text{Answer}(x_1, y_1) :- R(x_1, y_1), R(y_1, z_1), P(z_1, y_1), S(x_1, y_1, w_1), R(y_1, w_1),$$

where  $w_1$  is a new variable not in  $q_1$ . Now observe that  $(q_2, (x_2, y_2)) \rightarrow (q'_1, (x_1, y_1))$ , which implies that  $(q_2, (x_2, y_2)) \rightarrow (D, (a, b))$ . By the Homomorphism Theorem,  $(a, b) \in q_2(D)$ , and thus,  $q_1(D) \subseteq q_2(D)$ .

Our goal is to revisit the problem of CQ containment in the presence of INDs. Given two CQs  $q$  and  $q'$ , and a set  $\Sigma$  of INDs,  $q$  is contained in  $q'$  under  $\Sigma$ , denoted  $q \subseteq_{\Sigma} q'$ , if for every database  $D$  that satisfies  $\Sigma$ ,  $q(D) \subseteq q'(D)$ . The problem of interest is defined as expected:

**Problem: CQ-Containment-IND**

**Input:** Two CQs  $q$  and  $q'$ , and a set  $\Sigma$  of INDs

**Output:** **true** if  $q \subseteq_{\Sigma} q'$ , and **false** otherwise

Although the complexity of CQ containment in the presence of FDs remains NP-complete (Theorem 17.2), this is not true for INDs:

**Theorem 17.6**

CQ-Containment-IND is PSPACE-complete.

We first focus on the upper bound. Recall again that, by the Homomorphism Theorem, checking whether a CQ  $q(\bar{x})$  is contained in a CQ  $q'(\bar{x}')$  in the absence of constraints boils down to checking whether  $(q', \bar{x}') \rightarrow (q, \bar{x})$ . Although this is not enough in the presence of INDs, we can adopt a similar approach providing that we first transform, by adding atoms as dictated by the INDs, the set of atoms  $A_q$  occurring in  $q$  into a new set of atoms  $S$  that satisfies the INDs, and then check whether  $(A_{q'}, \bar{x}') \rightarrow (S, \bar{x})$ . This simple idea has been already illustrated by Example 17.5. As expected, the transformation of  $A_q$  into  $S$  can be achieved by exploiting the chase for INDs, which has been already introduced in Chapter 11.

We are going to establish a statement analogous to Theorem 17.3. However, since the chase for INDs may build an infinite set of atoms, we can only characterize CQ containment under possibly infinite databases. Notice that here we refer to the output of a CQ over a possibly infinite database. Although this is defined in the same way as for databases (Definition 12.3), we proceed to give the formal definition for the sake of completeness.

Consider a possibly infinite database  $D$  and a CQ  $q$  of the form

$$\text{Answer}(\bar{x}) \text{ :- } R_1(\bar{u}_1), \dots, R_n(\bar{u}_n).$$

An *assignment* for  $q$  over  $D$  is a function  $\eta$  from the set of variables in  $q$  to  $\text{Dom}(D)$ . We say that  $\eta$  is *consistent* with  $D$  if

$$\{R_1(\eta(\bar{u}_1)), \dots, R_n(\eta(\bar{u}_n))\} \subseteq D,$$

where, for  $i \in [n]$ ,  $R_i(\eta(\bar{u}_i))$  is the fact obtained after replacing each variable  $x$  in  $\bar{u}_i$  with  $\eta(x)$ , and leave the constants in  $\bar{u}_i$  untouched. Having this notion, we can define what is the output of a CQ on a possibly infinite database.

**Definition 17.7: Evaluation on Possibly Infinite Databases**

Given a possibly infinite database  $D$  of a schema  $\mathbf{S}$ , and a CQ  $q(\bar{x})$  over  $\mathbf{S}$ , the *output* of  $q$  on  $D$  is defined as the set of tuples

$$q(D) = \{\eta(\bar{x}) \mid \eta \text{ is an assignment for } q \text{ over } D \text{ consistent with } D\}.$$

We can naturally talk about homomorphisms from CQs to possibly infinite databases. Actually, Definition 13.1 merely extends to possibly infinite databases, which allows us to state a result analogous to Theorem 13.2:

**Theorem 17.8**

Given a possibly infinite database  $D$  of a schema  $\mathbf{S}$ , and a CQ  $q(\bar{x})$  of arity  $k \geq 0$  over  $\mathbf{S}$ , it holds that

$$q(D) = \{\bar{a} \in \text{Dom}(D)^k \mid (q, \bar{x}) \rightarrow (D, \bar{a})\}.$$

Consider two CQs  $q$  and  $q'$ , and a set  $\Sigma$  of INDs. We say that  $q$  is *contained without restriction in  $q'$  under  $\Sigma$* , denoted  $q \subseteq_{\Sigma}^{\infty} q'$ , if for every possibly infinite database  $D$  that satisfies  $\Sigma$ ,  $q(D) \subseteq q'(D)$ . For brevity, we write  $\text{Chase}(q, \Sigma)$  instead of  $\text{Chase}(A_q, \Sigma)$ . The next result is shown as Theorem 17.3:

**Theorem 17.9**

Let  $q(\bar{x}), q'(\bar{x}')$  be CQs over schema  $\mathbf{S}$ , and  $\Sigma$  a set of INDs over  $\mathbf{S}$ . Then:

$$q \subseteq_{\Sigma}^{\infty} q' \quad \text{if and only if} \quad (A_{q'}, \bar{x}') \rightarrow (\text{Chase}(q, \Sigma), \bar{x}).$$

The above statement alone is of little use since we are interested in finite databases. However, combined with the following result, known as the *finite controllability* of CQ containment under INDs, we get the desired characterization of CQ containment under finite databases via the chase.

**Theorem 17.10: Finite Controllability of Containment**

Let  $q$  and  $q'$  be CQs over a schema  $\mathbf{S}$ , and  $\Sigma$  a set of INDs over  $\mathbf{S}$ . Then:

$$q \subseteq_{\Sigma} q' \quad \text{if and only if} \quad q \subseteq_{\Sigma}^{\infty} q'.$$

The above theorem is a deep result that is extremely useful for our analysis, but whose proof is out of the scope of this book. An easy consequence of Theorems 17.9 and 17.10 combined with Theorem 17.8 is the following:

**Corollary 17.11**

Let  $q(\bar{x})$  and  $q'(\bar{x}')$  be CQs over a schema  $\mathbf{S}$ , and  $\Sigma$  a set of INDs over  $\mathbf{S}$ . With  $S = \text{Chase}(q, \Sigma)$ , the following holds:

$$q \subseteq_{\Sigma} q' \text{ if and only if } G_S(\bar{x}) \in q'(G_S(S)).$$

Due to Corollary 17.11, the reader may be tempted to think that the procedure for checking whether  $q \subseteq_{\Sigma} q'$ , which in turn will lead to the PSPACE upper bound claimed in Theorem 17.6, is to check whether  $G_S(\bar{x})$  belongs to the evaluation of  $q'$  over  $S^{\downarrow}$ , where  $S = \text{Chase}(q, \Sigma)$ . However, it should not be forgotten that  $\text{Chase}(q, \Sigma)$  may be infinite. Hence, we need a finer procedure that avoids the explicit construction of  $\text{Chase}(q, \Sigma)$ . We present a lemma that is the building block of this procedure, but first we need some terminology.

For an IND  $\sigma = R[i_1, \dots, i_m] \subseteq P[j_1, \dots, j_m]$ , a tuple  $\bar{u} = (u_1, \dots, u_{\text{ar}(R)})$ , and a set of variables  $V$ ,  $\text{new}^V(\sigma, \bar{u})$  is the atom obtained from  $\text{new}(\sigma, \bar{u})$  after replacing each newly introduced variable with a distinct variable from  $V$ . Formally,  $\text{new}^V(\sigma, \bar{u}) = P(v_1, \dots, v_{\text{ar}(P)})$ , where, for each  $\ell \in [\text{ar}(P)]$ ,

$$v_{\ell} = \begin{cases} u_{i_k} & \text{if } \ell = j_k, \text{ for } k \in [m], \\ x \in V & \text{otherwise,} \end{cases}$$

such that, for each  $i, j \in [\text{ar}(P)] - \{j_1, \dots, j_m\}$ ,  $i \neq j$  implies  $v_i \neq v_j$ .<sup>1</sup> Given two CQs  $q(\bar{x}), q'(\bar{x}')$  over a schema  $\mathbf{S}$ , and a set  $\Sigma$  of INDs over  $\mathbf{S}$ , a *witness of  $q'$  relative to  $q$  and  $\Sigma$*  is a triple  $(\mathcal{V}, \mathcal{S}, Q)$ , where  $\mathcal{V}$  is a sequence of (not necessarily disjoint) sets of variables  $V_1, \dots, V_n$ , for  $n \geq 0$ ,  $\mathcal{S}$  is a sequence of disjoint sets of relational atoms  $S_0, \dots, S_n$ , and  $Q \subseteq \bigcup_{i \in [0, n]} S_i$ , such that:

- $|\bigcup_{i \in [n]} V_i| \leq 3 \cdot |A_{q'}| \cdot \max_{R \in \mathbf{S}} \{\text{ar}(R)\}$ ,
- for each  $i \in [n]$ ,  $V_i \cap (\text{Dom}(S_{i-1}) \cup \text{Dom}(S)) = \emptyset$ ,
- for each  $i \in [0, n]$ ,  $|S_i| \leq |A_{q'}|$ ,
- $S_0 \subseteq A_q$ ,
- for each  $i \in [n]$  and  $P(\bar{v}) \in S_i$ , there exists  $\sigma = R[\alpha] \subseteq P[\beta]$  in  $\Sigma$  that is applicable on  $S_{i-1}$  with some  $\bar{u} \in R^{S_{i-1}}$  such that  $P(\bar{v}) = \text{new}^{V_i}(\sigma, \bar{u})$ ,
- for each  $i \in [n]$  and  $x \in \text{Dom}(S_i) - \text{Dom}(S_{i-1})$ , there is only one occurrence of  $x$  in  $S_i$ , i.e., it is mentioned only once by exactly one atom of  $S_i$ ,
- $|Q| \leq |A_{q'}|$ , and
- $G_Q(\bar{x}) \in q'(G_Q(Q))$ .

<sup>1</sup> We assume some fixed mechanism that chooses the variable  $v_{\ell}$  from the set  $V$  whenever  $\ell \in [\text{ar}(P)] - \{j_1, \dots, j_m\}$ .

Let  $S = \text{Chase}(q, \Sigma)$ . Notice that  $G_S(\bar{x}) \in q'(G_S(S))$  holds due to the existence of a set  $A \subseteq \text{Chase}(q, \Sigma)$  such that  $(A_{q'}, \bar{x}') \rightarrow (A, \bar{x})$ . It is also not difficult to see that the construction of  $A$  can be witnessed via a sequence  $A_0, A_1, \dots, A_n$  of disjoint subsets of  $\text{Chase}(q, \Sigma)$ , where each such set consists of at most  $|A_{q'}|$  atoms,  $A_0 \subseteq A_{q'}$ ,  $A_n = A$ , and for each  $i \in [n]$ , the atoms of  $A_i$  are obtained from the atoms of  $A_{i-1}$  via chase applications using INDs of  $\Sigma$ . A witness of  $q'$  relative to  $q$  and  $\Sigma$  should be understood as a compact representation, which uses only polynomially many variables, of such a sequence  $A_0, A_1, \dots, A_n$  of disjoint subsets of  $\text{Chase}(q, \Sigma)$ . Therefore, the existence of a witness of  $q'$  relative to  $q$  essentially implies that  $G_S(\bar{x}) \in q'(G_S(S))$ . Furthermore, if  $G_S(\bar{x}) \in q'(G_S(S))$ , then a witness of  $q'$  relative to  $q$  and  $\Sigma$  can be extracted from  $\text{Chase}(q, \Sigma)$ . The above informal discussion is summarized in the following technical lemma, whose proof is left as an exercise.

---

**Algorithm 5** CONTAINMENTWITNESS( $q, q', \Sigma$ )

---

**Input:** Two CQs  $q(\bar{x})$  and  $q'(\bar{x}')$  over  $\mathbf{S}$ , and a set  $\Sigma$  of INDs over  $\mathbf{S}$ .

**Output:** **true** if there is a witness for  $q'$  relative to  $q$  and  $\Sigma$ , and **false** otherwise.

```

1:  $A_\nabla := A$ , where  $A \subseteq A_{q'}$  with  $|A| \leq |A_{q'}|$ 
2:  $A_\triangleright := \emptyset$ 
3:  $Q := A$ , where  $A \subseteq A_\nabla$ 
4:  $V := \{y_1, \dots, y_m\} \subset \text{Var} - \text{Dom}(A_q)$  for some  $m \in [3 \cdot |A_{q'}| \cdot \max_{R \in \Sigma} \{\text{ar}(R)\}]$ 
5: repeat
6:   repeat
7:     if  $\sigma = R[\alpha] \subseteq P[\beta] \in \Sigma$  is applicable on  $A_\nabla$  with  $\bar{u} \in \text{Dom}(A_\nabla)^{\text{ar}(R)}$ 
8:       then
9:          $N := \text{new}^V(\sigma, \bar{u})$ 
10:         $V := V - \text{Dom}(\{N\})$ 
11:         $A_\triangleright := A_\triangleright \cup \{N\}$ 
12:        if  $|A_\triangleright| < |A_{q'}|$  then
13:           $\text{Next} := b$ , where  $b \in \{0, 1\}$ 
14:        else
15:           $\text{Next} := 1$ 
16:    until  $\text{Next} = 1$ 
17:    if  $A_\triangleright = \emptyset$  then
18:      return false
19:     $V := V \cup ((\text{Dom}(A_\nabla) \cap \text{Var}) - (\text{Dom}(A_\triangleright) \cup \text{Dom}(Q)))$ 
20:     $A_\nabla := A_\triangleright$ 
21:     $A_\triangleright := \emptyset$ 
22:     $Q := Q \cup A$ , where  $A \subseteq A_\nabla$ 
23:    if  $|Q| < |A_{q'}|$  then
24:       $\text{Evaluate} := b$ , where  $b \in \{0, 1\}$ 
25:    else
26:       $\text{Evaluate} := 1$ 
27:  until  $\text{Evaluate} = 1$ 
28: return  $G_Q(\bar{x}) \in q'(G_Q(Q))$ 

```

---

**Lemma 17.12.** *Let  $q(\bar{x})$  and  $q'(\bar{x}')$  be CQs over a schema  $\mathcal{S}$ , and  $\Sigma$  a set of INDs over  $\mathcal{S}$ . With  $S = \text{Chase}(q, \Sigma)$ , it holds that  $G_S(\bar{x}) \in q'(G_S(S))$  if and only if there exists a witness of  $q'$  relative to  $q$  and  $\Sigma$ .*

By Corollary 17.11 and Lemma 17.12, we conclude that the problem of checking whether a CQ  $q(\bar{x})$  is contained in a CQ  $q'(\bar{x}')$  under a set  $\Sigma$  of INDs boils down to checking whether a witness of  $q'$  relative to  $q$  and  $\Sigma$  exists. This is done via the nondeterministic procedure shown in Algorithm 5. It essentially constructs the sequence of sets of variables  $V_1, \dots, V_n$ , and the sequence of sets of atoms  $S_0, \dots, S_n$ , required by a witness for  $q'$  relative to  $q$  and  $\Sigma$ , one after the other (if they exist), without storing more than two consecutive sets of a sequence during its computation. It also constructs on the fly the set of atoms  $Q$ . This is done by storing some of the atoms of a set  $S_i$  (possibly none) into  $Q$  before discarding it. Finally, the algorithm checks whether  $G_Q(\bar{x}) \in q'(G_Q(Q))$ , in which case it returns **true**; otherwise, it returns **false**. We proceed to give a bit more detailed description of Algorithm 5.

**Initialization.** The algorithm starts by guessing a subset of  $A_q$  with at most  $|A_{q'}|$  atoms, which is stored in  $A_{\nabla}$  (see line 1);  $A_{\nabla}$  should be seen as the “current set” from which we construct the “next set”  $A_{\triangleright}$  in the sequence. It also guesses a subset of  $A_{\nabla}$  that is stored in  $Q$  (see line 3); this step is part of the “on the fly” construction of the set  $Q$ . It also collects  $3 \cdot |A_{q'}| \cdot \max_{R \in \mathcal{S}} \{\text{ar}(R)\}$  variables not occurring in  $A_q$  in the set  $V$  (see line 4).

**Inner repeat-until loop.** The inner repeat-until loop (see lines 6 - 15) is responsible for constructing the set  $A_{\triangleright}$  from  $A_{\nabla}$ . This is done by guessing an IND  $\sigma \in \Sigma$  and a tuple  $\bar{u}$  over  $\text{Dom}(A_{\nabla})$ , and adding to  $A_{\triangleright}$  the atom  $\text{new}^V(\sigma, \bar{u})$  if  $\sigma$  is applicable on the current set  $A_{\nabla}$  with  $\bar{u}$ . It also removes from  $V$  the variables that has been used in  $\text{new}^V(\sigma, \bar{u})$  since they should not be reused in any other atom of  $A_{\triangleright}$  that will be generated by a subsequent iteration. This is repeated until  $A_{\triangleright}$  contains exactly  $|A_{q'}|$  atoms, which means that its construction has been completed, or the algorithm nondeterministically chooses that its construction has been completed, even if it contains less than  $|A_{q'}|$  atoms, by setting *Next* to 1. Once  $A_{\triangleright}$  is in place, the algorithm updates  $V$  by adding to it the variables that occur in the current set  $A_{\nabla}$ , but have not been propagated to  $A_{\triangleright}$  and do not occur in  $Q$  (see line 18). This essentially gives rise to the next set of variables in the sequence of sets of variable under construction. Then  $A_{\nabla}$  is not needed further, and we can reuse the space that it occupies. The set  $A_{\triangleright}$  becomes the current set  $A_{\nabla}$  (see line 19), while  $A_{\triangleright}$  becomes empty (see line 20). Then the algorithm guesses a subset of  $A_{\nabla}$  that is stored in  $Q$  (see line 21); this step is part of the “on the fly” construction of  $Q$ .

**Outer repeat-until loop.** The above is repeated until  $Q$  contains more than  $|A_{q'}|$  atoms (in the worst-case,  $2 \cdot |A_{q'}|$  atoms), which means that its construction has been completed, or the algorithm nondeterministically chooses that its construction has been completed, even if it contains less



than  $|A_{q'}|$  atoms, by setting *Evaluate* to 1. The algorithm returns **true** if  $G_{S'}(\bar{x}) \in q'(G_{S'}(S'))$ ; otherwise, it returns **false**.

It is not difficult to verify that Algorithm 5 uses polynomial space, which is actually the space needed to represent the sets  $A_{\nabla}$ ,  $A_{\triangleright}$ ,  $Q$  and  $V$ , as well as the space needed to check whether an IND is applicable on  $A_{\nabla}$  with some tuple  $\bar{u} \in \text{Dom}(A_{\nabla})^{\text{ar}(R)}$  (see line 7), and the space needed to check whether  $G_Q(\bar{x}) \in q'(G_Q(Q))$  (see line 27). This shows that **CQ-Containment-IND** is in NPSpace, and thus in PSpace since  $\text{NPSpace} = \text{PSpace}$ .

The PSpace-hardness of **CQ-Containment-IND** is shown via a reduction from **IND-Implication**, which is PSpace-hard (see Theorem 11.9). Recall that the **IND-Implication** problem takes as input a set  $\Sigma$  of INDs over a schema  $\mathbf{S}$ , and an IND  $\sigma$  over  $\mathbf{S}$ , and asks whether  $\Sigma \models \sigma$ , i.e., whether for every database over  $\mathbf{S}$ ,  $D \models \Sigma$  implies  $D \models \sigma$ . We are going to construct two CQs  $q$  and  $q'$  such that  $\Sigma \models \sigma$  if and only if  $q \subseteq_{\Sigma} q'$ .

Assume that  $\sigma = R[i_1, \dots, i_k] \subseteq P[j_1, \dots, j_k]$ . The CQ  $q$  is defined as

$$\text{Answer}(x_{i_1}, \dots, x_{i_k}) \quad :- \quad R(x_1, \dots, x_{\text{ar}(R)}),$$

while the CQ  $q'$  is defined as

$$\text{Answer}(x_{i_1}, \dots, x_{i_k}) \quad :- \quad R(x_1, \dots, x_{\text{ar}(R)}), P(x_{f(1)}, \dots, x_{f(\text{ar}(R))}),$$

where, for each  $m \in [\text{ar}(P)]$ ,

$$f(m) = \begin{cases} i_{\ell} & \text{if } m = j_{\ell}, \text{ where } \ell \in [k], \\ \text{ar}(R) + m & \text{otherwise.} \end{cases}$$

The function  $f$  ensures that the variable at position  $j_{\ell}$  in the  $P$ -atom of  $q'$  is  $x_{i_{\ell}}$ , i.e., the same as the one at position  $i_{\ell}$  in the  $R$ -atom of  $q'$ , while all the variables in the  $P$ -atom occurring at a position not in  $\{j_1, \dots, j_k\}$  are new variables occurring only once in the  $P$ -atom, and not occurring in the  $R$ -atom. It is an easy exercise to show that indeed  $\Sigma \models \sigma$  if and only if  $q \subseteq_{\Sigma} q'$ .

---

## Exercises

**Exercise 2.1.** Let  $q$  be the CQ given in Example 12.8. Express  $q$  as an RA query using  $\theta$ -joins instead of Cartesian product.

**Exercise 2.2.** Prove the correctness of the translation of a CQ into an SPJ query, and the translation of an SPJ query into a CQ, given in the proof of Theorem 12.7 which establishes that the languages of CQs and of SPJ queries are equally expressive.

**Exercise 2.3.** For a CQ  $q$ , let  $e_q$  be the equivalent SPJ query obtained by applying the translation in the proof of Theorem 12.7. What is the size of  $e_q$  with respect to the size of  $q$ ? Conversely, assuming that  $q_e$  is the CQ obtained after translating an SPJ query  $e$  into a CQ according to the translation in the proof of Theorem 12.7, what is the size of  $q_e$  with respect to the size of  $e$ ?

**Exercise 2.4.** Prove that the choice of a pairing function in the definition of direct product does not matter. More precisely, let  $\otimes_\tau$  be the direct product defined using a pairing function  $\tau$ . Then, for every Boolean FO query  $q$ , every two databases  $D$  and  $D'$ , and every two pairing functions  $\tau$  and  $\tau'$ , show that  $D \otimes_\tau D' \models q$  if and only if  $D \otimes_{\tau'} D' \models q$ .

**Exercise 2.5.** Let  $q$  be a Boolean FO query without constants over a schema  $\mathbf{S}$ . Prove that the following are equivalent:

1. There exists a CQ  $q'$  over  $\mathbf{S}$  such that  $q \equiv q'$ , i.e.,  $q(D) = q'(D)$  for every database  $D$  of  $\mathbf{S}$ .
2.  $q$  is preserved under homomorphisms and direct products.

**Exercise 2.6.** The goal of this exercise is to extend the notion of preservation under direct products to queries with constants. To this end, we first refine the definition of a pairing function. Let  $C \subseteq \text{Const}$  be a finite set of constants, and  $\tau_C$  a pairing function such that  $\tau_C(a, a) = a$  for each  $a \in C$ . First, prove that such a pairing function exists. Then, prove that for any two databases  $D$

and  $D'$  of the same schema  $\mathbf{S}$ , and for a Boolean CQ  $q$  over  $\mathbf{S}$  that mentions only constants from  $C$ , if  $D \models q$  and  $D' \models q$ , then  $D \otimes D' \models q$ , where the definition of  $\otimes$  uses the pairing function  $\tau_C$ .

**Exercise 2.7.** The goal is to extend further the notion of preservation under direct products to queries with constants that are not Boolean. For a finite set of constants  $C \subseteq \text{Const}$ , let  $\tau_C$  be a pairing function defined as in Exercise 2.6. Then, given two tuples  $\bar{a} = (a_1, \dots, a_n)$  and  $\bar{b} = (b_1, \dots, b_n)$ , define the  $n$ -ary tuple  $\bar{a} \otimes \bar{b}$  as  $(\tau_C(a_1, b_1), \dots, \tau_C(a_n, b_n))$ . Consider now an  $n$ -ary CQ  $q(\bar{x})$  that mentions only constants from  $C$ . Show that if  $\bar{a} \in q(D)$  and  $\bar{b} \in q(D')$ , then  $\bar{a} \otimes \bar{b} \in q(D \otimes D')$ , where  $\otimes$  is defined with the pairing function  $\tau_C$ .

**Exercise 2.8.** Use Exercise 2.6 to prove that the Boolean query  $q = \exists x (x = a)$ , where  $a$  is a constant, cannot be expressed as a CQ.

**Exercise 2.9.** Use Exercise 2.7 to prove that the query  $q = \varphi(x, y)$ , where  $\varphi$  is the equational atom  $(x = y)$ , cannot be expressed as a CQ.

**Exercise 2.10.** Consider a parameterized problem  $(L_1, \kappa_1)$  over  $\Sigma_1$ , and a parameterized problem  $(L_2, \kappa_2)$  over  $\Sigma_2$ . Show that if there is an FPT-reduction from  $(L_1, \kappa_1)$  to  $(L_2, \kappa_2)$ , and  $(L_2, \kappa_2) \in \text{FPT}$ , then  $(L_1, \kappa_1) \in \text{FPT}$ .

**Exercise 2.11.** Recall that in the proof of the fact that **p-CQ-Evaluation** is in  $\text{W}[1]$  (see Theorem 14.7), for technical clarity, we consider only constant-free Boolean CQs over a schema consisting of a single binary relation name. Prove that **p-CQ-Evaluation** is in  $\text{W}[1]$  even for arbitrary CQs.

**Exercise 2.12.** Show Corollary 15.7 for arbitrary (non-Boolean) CQs.

**Exercise 2.13.** Show that the binary relation  $\equiv$  over CQs is an equivalence relation, i.e., is reflexive, symmetric, and transitive. Show also that the binary relation  $\subseteq$  over CQs is reflexive and transitive, but not necessarily symmetric.

**Exercise 2.14.** Answer the following questions about CQs and their cores.

- (i) Consider the Boolean CQ  $q_1$  over the schema  $\{E[2]\}$  defined as

$$\text{Answer} \text{ :- } E(x_1, y_1), E(y_1, z_1), E(z_1, w_1), E(w_1, x_1), E(x_2, y_2), E(y_2, x_2).$$

Assume that  $E$  is used to represent the edge relation of a directed graph  $G$ . What does  $q_1$  check for  $G$ ? Compute the core of  $q_1$ .

- (ii) Consider the Boolean CQ  $q_2$  over the schema  $\{R[1], S[1]\}$  defined as

$$\text{Answer} \text{ :- } R(x), S(x), R(y), S(y).$$

Compute the core of  $q_2$ .

(iii) Consider the CQ  $q_3$  over the schema  $\{R[1], S[1]\}$  defined as

$$\text{Answer}(x, y) \text{ :- } R(x), S(x), R(y), S(y).$$

Prove that  $q_3$  is a core of itself.

**Exercise 2.15.** Let  $q(\bar{x})$  be a CQ, and  $q'(\bar{x})$  a core of  $q(\bar{x})$ . Prove that there is a homomorphism from  $(q, \bar{x})$  to  $(q', \bar{x})$  that is the identity on  $\text{Dom}(S_{q'})$ .

**Exercise 2.16.** Recall that COMPUTECORE (see Algorithm 4) is nondeterministic. Devise a deterministic algorithm that computes the core of a CQ, and show that it runs in exponential time in the size of the input query.

**Exercise 2.17.** (a) Let CQ-Minimization be the problem where, given a Boolean CQ  $q$  and integer  $k \in \mathbb{N}$ , the question is if there exists a CQ  $q'$  such that  $q' \equiv q$  and  $|q'| \leq k$ . Prove that CQ-Minimization is NP-complete.  
 (b) Let CQ-Minimality be the problem where, given a Boolean CQ  $q$ , the question is to answer **true** if  $q$  is minimal and **false** otherwise. Prove that CQ-Minimality is coNP-complete.

**Exercise 2.18.** Let  $D$  be a database, and  $T = \{\bar{a}_1, \dots, \bar{a}_n\}$  a set of  $m$ -ary tuples over  $\text{Dom}(D)$ , for  $m > 0$ . Show that there exists a CQ  $q(\bar{x})$  such that  $q(D) = T$  if and only if the following hold:

1.  $\prod_{i \in [n]} \bar{a}_i$  appears in  $\prod_{i \in [n]} D$ , and
2. there is no tuple  $\bar{b} \in \text{Dom}(D)^m - T$  such that  $\prod_{i \in [n]} (D, \bar{a}_i) \rightarrow (D, \bar{b})$ .

**Exercise 2.19.** The purpose of this exercise is to understand what happens if we allow equalities of the form  $x = y$  or  $x = a$  in CQs. We define a *conjunctive query with equalities* (CQ<sup>=</sup>) similarly to a CQ, but we additionally allow equational atoms. Such queries can therefore be written as rules

$$\text{Answer}(\bar{x}) \text{ :- } R_1(\bar{u}_1), \dots, R_n(\bar{u}_n), y_1 = v_1, \dots, y_k = v_k,$$

where  $\{v_1, \dots, v_k\} \subseteq \text{Var} \cup \text{Const}$ .

1. Why are the queries  $\text{Answer}(x) \text{ :- } x = y$  and  $\text{Answer}(x) \text{ :- } x = a$  not expressible as CQs?
2. Prove that (i) Theorem 14.1 (ii) Theorem 15.3, and (iii) Theorem 15.8 also hold for CQ<sup>=</sup>.
3. Prove that queries in CQ<sup>=</sup> can be minimized with a variant of Algorithm 4.

**Exercise 2.20.** Prove that the following problem is CONEXPTIME-complete: given a database  $D$ , and a set  $T = \{\bar{a}_1, \dots, \bar{a}_n\}$  of  $m$ -ary tuples over  $\text{Dom}(D)$ , for  $m > 0$ , check whether there exists a CQ  $q$  such that  $q(D) = T$ .

**Exercise 2.21.** Prove that FO-Containment remains undecidable even if one of the two input queries is a CQ.

**Exercise 2.22.** Prove Lemma 17.12

**Exercise 2.23.** Prove that the reduction at the end of Chapter 17 from IND-Implication to CQ-Containment-IND, which establishes that the latter is PSPACE-hard, is correct.

---

## Bibliographic Comments

To be done.

## **Fast Conjunctive Query Evaluation**





---

## Motivation

Here we are interested in understanding when CQ evaluation can be solved efficiently in combined complexity. In Theorem [14.1](#), we have shown that CQ evaluation is NP-complete by reducing from an NP-complete problem over graphs. It is known, on the other hand, that several NP-complete problems over graphs become tractable if they are restricted to be *nearly acyclic*. As we show in this part of the book, similar ideas can be applied to prove that CQ evaluation is tractable when CQs are nearly acyclic. This is highly relevant from a practical point of view, as such CQs appear often in real-world applications.



## Acyclicity of Conjunctive Queries

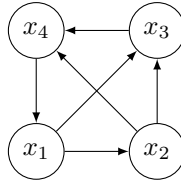
We start by studying the notion of *acyclicity* for CQs, which has received considerable attention in the database literature since the early 1980s. In this chapter, we define acyclicity and present an algorithm to recognize it. In the next chapter, we will present two algorithms that show that acyclic CQs can be evaluated efficiently.

### Conjunctive Queries and Hypergraphs

We have seen in Theorem 14.1 that the evaluation problem for CQs is NP-complete. So, the reader may wonder why databases are so successful in practice, even though the most fundamental database problem on the most common class of queries is NP-complete. The crux is that the *syntactic shape* of a CQ plays a key role on how complex is its evaluation. For instance, assume that we have a database  $D$  of the schema  $\{E[2]\}$ , i.e.,  $D$  can be understood as a directed graph with  $E$  being the edge relation. Evaluating the CQ

$$\text{Answer} := E(x_1, x_2), E(x_2, x_3), E(x_3, x_4), E(x_4, x_1), E(x_1, x_3), E(x_2, x_4)$$

can be seen as matching a variant of the 4-clique, namely a graph of the form



in  $D$ . In fact, this correspondence between evaluation of CQs and graph matching is precisely what we used in Theorem 14.1 to reduce the Clique problem into the CQ-Evaluation problem.

However, CQs in practice are usually not shaped as cliques. Instead, *tree-shaped* CQs are much more common. Since it is well-known that finding cliques in graphs is computationally difficult, whereas finding tree-like structures in graphs is much easier, it makes sense to study the evaluation problem of CQs for which the associated graph is acyclic.

To make this precise, however, we need to consider a generalization of undirected graphs that can deal with relations of arity three or more. Such graphs are called *hypergraphs*.

#### Definition 18.1: Hypergraph

A *hypergraph* is a pair  $H = (V, E)$ , where

- $V$  is a finite set of *nodes* and
- $E$  is a set of subsets of  $V$ , called *hyperedges*.

### Acyclicity of Hypergraphs

Defining the notion of acyclicity for hypergraphs is not as simple as it is for graphs. In fact, several natural, non-equivalent notions of acyclicity for hypergraphs exist. We work here with one such a notion, often referred to as  $\alpha$ -*acyclicity*, which has received considerable attention in database theory.

We will call a hypergraph  $H$  *acyclic* if it admits a *join tree*, that is, if its hyperedges can be arranged in the form of a tree (recall the definition of tree from Chapter 2), while preserving the connectivity of elements that occur in different hyperedges.

#### Definition 18.2: Join Tree and Acyclic Hypergraph

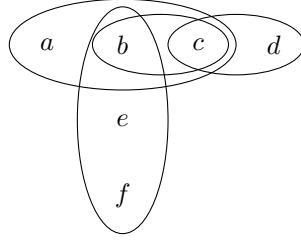
Given a hypergraph  $H = (V, E)$ , a tree  $T$  is a *join tree* of  $H$  if

- the nodes of  $T$  are precisely the hyperedges in  $E$  and,
- for each node  $v \in V$ , the set of nodes of  $T$  in which  $v$  is an element forms a connected subtree of  $T$ .

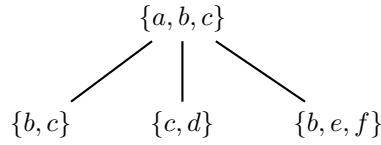
Moreover,  $H$  is *acyclic* if  $H$  admits a join tree.

#### Example 18.3: Acyclic and Non-Acyclic Hypergraphs

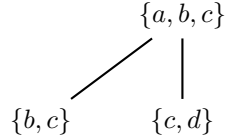
Consider the following hypergraph  $H_1 = (V_1, E_1)$ :



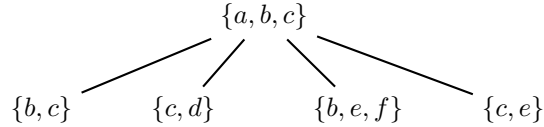
Thus, we have that  $V_1 = \{a, b, c, d, e, f\}$  and  $E_1 = \{\{a, b, c\}, \{b, c\}, \{c, d\}, \{b, e, f\}\}$ . It holds that  $H_1$  is an acyclic hypergraph, as the following tree  $T_1$  is a join tree for  $H_1$ :



Recall that by convention, the root of  $T_1$  is depicted on top and its edges are directed downwards. We have that  $T_1$  is a join tree of  $H_1$  as the nodes of  $T_1$  are precisely the hyperedges in  $E_1$ , and for each  $v \in V_1$ , the set of nodes of  $T_1$  in which  $v$  occurs defines a connected subtree of  $T_1$ . As an example of this latter condition, if we consider  $v = c$ , then we obtained the following subtree of  $T_1$  that is connected:



On the other hand, consider a hypergraph  $H_2$  that extends  $H_1$  with the hyperedge  $\{c, e\}$ . We have that  $H_2$  is not acyclic, as it is not possible to construct a join tree for it. For instance, consider the extension  $T_2$  of  $T_1$  that is obtained by adding a node  $\{c, e\}$  as a children of  $\{a, b, c\}$ .



Then we have that  $T_2$  is not a join tree for  $H_2$  as the set of nodes of  $T_2$  in which  $e$  occurs do not define a connected subtree of  $T_2$ :

$$\{b, e, f\} \quad \{c, e\}$$

It is not hard to see that for undirected graphs, the notion of  $\alpha$ -acyclicity coincides with the usual notion of acyclicity that stems from graph theory (i.e., tree-shaped or forest-shaped graphs).

## Acyclicity of Conjunctive Queries

The notion of acyclic hypergraph is the key concept in the definition of acyclic CQs. Each CQ  $q$  is naturally associated with a hypergraph  $H_q$  that represents the structure of joins among its variables. In particular, if  $q$  is of the form

$$\text{Answer}(\bar{x}) \text{ :- } R_1(\bar{u}_1), \dots, R_n(\bar{u}_n),$$

where  $\bar{u}_i$  is a tuple of constants and variables for every  $i \in [n]$ , then  $H_q = (V, E)$  is a hypergraph such that

- its set  $V$  of vertices contains all variables mentioned in  $q$  and
- the hyperedges in  $E$  are precisely the sets of variables appearing in the atoms of  $q$ , i.e.,  $E = \{X_i \mid i \in [n]\}$ , where  $X_i$  is the set of variables occurring in  $\bar{u}_i$ .

### Definition 18.4: Acyclicity of CQs

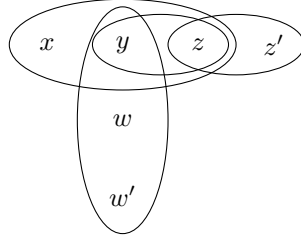
An *acyclic conjunctive query* (ACQ) is a conjunctive query  $q$  such that its associated hypergraph  $H_q$  is acyclic.

### Example 18.5: Acyclic and Non-Acyclic CQs

Consider the following CQ  $q_1$ :

$$\text{Answer}(x, y) \text{ :- } R(x, y, z), T(y, z), S(y, w, w'), T(z, z').$$

Then we have that  $H_{q_1}$  is the following hypergraph:



We know from Example 18.3 that  $H_{q_1}$  is an acyclic hypergraph, as  $H_{q_1}$  can be obtained from the hypergraph  $H_1$  in Example 18.3 by renaming the nodes of  $H_1$ . Therefore, we have that  $q_1$  is an acyclic CQ. In the same way, we obtain that the following CQ  $q_2$ :

$$\text{Answer}(x, y) :- R'(x, y, z, a), T'(y, z, a), S(y, w, w'), T'(z, z', b)$$

is acyclic as  $H_{q_2} = H_{q_1}$ . In particular, notice that constants  $a, b$  in  $q_2$  do not play any role in the construction of  $H_{q_2}$ . On the other hand, the following CQ  $q_3$ :

$$\text{Answer}(x, y) :- R(x, y, z), T(y, z), S(y, w, w'), T(z, z'), T(z, w)$$

is not acyclic as the hypergraph  $H_{q_3}$  is not acyclic. Notice that this latter fact is also obtained from Example 18.3, as  $H_{q_3}$  can be obtained from the hypergraph  $H_2$  in Example 18.3 by renaming the nodes of  $H_2$ . Finally, consider the following CQ  $q_4$ :

$$\text{Answer}(x, y) :- R(x, y, z), R(y, y, z), T(y, z), S(y, w, w'), T(z, z').$$

Then we have that  $q_3$  is an acyclic CQ since  $H_{q_3} = H_{q_1}$ . Notice that this latter condition holds as the set of variable occurring in  $R(y, y, z)$  is  $\{y, z\}$ , which is the same as the set of variables occurring in  $T(y, z)$ .

## Acyclicity Recognition

In the following chapter, we will show that ACQs can be evaluated efficiently. But before doing so, it is important to explain why acyclicity itself can be efficiently recognized. This follows from the existence of an equivalent definition of acyclicity in terms of an iterative process described in the following proposition.

**Proposition 18.6: GYO Algorithm**

A hypergraph  $H = (V, E)$  is acyclic if and only if all of its vertices can be deleted by repeatedly applying the following two operations (in no particular order):

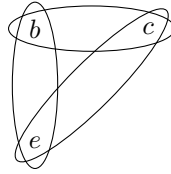
1. Delete a vertex that appears in at most one hyperedge.
2. Delete a hyperedge that is contained in another hyperedge.

The proof of Proposition 18.6 is left as an exercise for the reader (see Exercise 3.1). The characterization given in this proposition leads directly to a polynomial-time algorithm for checking acyclicity of hypergraphs, and thus of CQs: Given a CQ  $q$ , we apply operations (1) and (2) in the statement of Proposition 18.6 over  $H_q$  until a fixpoint is reached. The query  $q$  is acyclic if and only if we are left with no vertices. Interestingly, a simple extension of this algorithm also constructs in polynomial time a join tree of  $H_q$  when  $q$  is acyclic (see Exercise 3.2).

**Example 18.7: Application of GYO Algorithm**

Consider the hypergraph  $H_1$  in Example 18.3. As expected, by using the previous algorithm we obtain that  $H_1$  is acyclic. In fact, all vertices of  $H_1$  are deleted by applying the following sequence of operations: delete vertex  $d$  (that appears only in hyperedge  $\{c, d\}$ ), delete vertices  $e$  and  $f$ , delete hyperedges  $\{b\}$  and  $\{c\}$  (that are contained in hyperedge  $\{b, c\}$ ), delete hyperedge  $\{b, c\}$  (that is contained in hyperedge  $\{a, b, c\}$ ), and delete vertices  $a$ ,  $b$  and  $c$ .

On the other hand, and also as expected, by applying the previous algorithm on hypergraph  $H_2$  from Example 18.3, we obtain that  $H_2$  is not acyclic. In fact, no matter what order is used when applying the two operations of the algorithm, we reach the following fixed point:



Notice that no operation can be applied to reduce this hypergraph, which is intuitively correct as this hypergraph represents the canonical example of an undirected graph that is not acyclic.



It is important to mention that there are more sophisticated algorithms that check whether a CQ  $q$  is acyclic and construct a join tree of  $H_q$  if the latter is the case, in time  $O(\|H_q\|)$ , that is, linear time. We summarize this result in the following proposition.

**Proposition 18.8: CQ Acyclicity Checking**

There exists a linear-time algorithm that, given a CQ  $q$ , checks whether  $q$  is acyclic, and if this is the case constructs a join tree of  $H_q$ .

We will exploit this proposition later in the presentation of efficient evaluation algorithms for acyclic conjunctive queries.



## Efficiently Evaluating Boolean ACQs

We will present two algorithms that show that acyclic CQs can be evaluated efficiently. The first one, known as Yannakakis's algorithm, makes use of the decomposition of an acyclic CQ as a join tree, which was defined in the previous chapter, while the second one is based on a simple *consistency* criterion. Yannakakis's algorithm achieves a relatively efficient running time of  $O(\|D\| \cdot \log \|D\| \cdot \|q\|)$ , where  $D$  is the database and  $q$  is the query. On the other hand, the algorithm based on the consistency criterion has the advantage that it does not require the CQ itself to be acyclic, only its *core* (as defined in Chapter 16).

### Semijoins and Acyclic CQs

The evaluation of acyclic CQs is tightly related to a particular relational algebra operation, known as *semijoin*, which we describe next. In the named relational algebra, one would define the semijoin as follows:

**Semijoin.** If  $e_1, e_2$  are named RA expressions of sort  $U_1$  and  $U_2$ , respectively, then their *semijoin*  $(e_1 \bowtie e_2)$  is a named RA expression of sort  $U_1$ . It is defined as  $(e_1 \bowtie e_2) := \pi_{U_1}(e_1 \Join e_2)$ .

In this chapter, we use the semijoin operator on outputs of conjunctive queries. Let  $D$  be a database. Given CQs  $q(\bar{x})$  and  $q'(\bar{x}')$  and tuples  $\bar{a} \in q(D)$  and  $\bar{b} \in q'(D)$ , we call  $\bar{a}$  and  $\bar{b}$  *consistent* if they have the same value on each position that contains a common variable of  $\bar{x}$  and  $\bar{x}'$ . We then define the semijoin of  $q(D)$  and  $q'(D)$ , denoted by  $q(D) \bowtie q'(D)$ , as the set of tuples  $\bar{a} \in q(D)$  that are consistent with some tuple  $\bar{b} \in q'(D)$ .

**Example 19.1: Semijoin of Conjunctive Queries**

Assume that  $D = \{R(a, b, c), R(d, d, d), S(c, b, e), S(d, e, e)\}$ , and that  $q$  and  $q'$  are the following CQs:

$$\begin{aligned} q &= \text{Answer}(x, y, z) \text{ :- } R(x, y, z) \\ q' &= \text{Answer}(z, y, w) \text{ :- } S(z, y, w). \end{aligned}$$

Then we have that  $q(D) = \{(a, b, c), (d, d, d)\}$ ,  $q'(D) = \{(c, b, e), (d, e, e)\}$ , and  $q(D) \bowtie q'(D) = \{(a, b, c)\}$ . In particular, we have that  $(a, b, c)$  is in  $q(D) \bowtie q'(D)$  since  $(a, b, c)$  belongs to  $q(D)$  and  $(a, b, c)$  is consistent with the tuple  $(c, b, e) \in q'(D)$ , as these two tuples have the same value  $b$  in the position that corresponds to the variable  $y$  shared by  $(x, y, z)$  and  $(z, y, w)$ , and have the same value  $c$  in the position that corresponds to the variable  $z$  shared by  $(x, y, z)$  and  $(z, y, w)$ . Moreover, we have that  $(d, d, d)$  is not in  $q(D) \bowtie q'(D)$  as this tuple is not consistent with any tuple in  $q'(D)$ . Finally, notice that  $q'(D) \bowtie q(D) = \{(c, b, e)\}$ , which shows that, as opposed to the case of the join operator,  $\bowtie$  is not commutative.

We now explain the relationship between the CQs in ACQ and the semijoin operator. Let  $q = \text{Answer} \text{ :- } R_1(\bar{u}_1), \dots, R_n(\bar{u}_n)$  be a Boolean ACQ, and consider an arbitrary join tree  $T$  of  $H_q$ . Recall that the set of nodes of  $T$  is  $\{X_i \mid i \in [n]\}$ , where each  $X_i$  is the set of variables occurring in  $\bar{u}_i$ . For every node  $s$  of  $T$ , we define the following CQs for some  $i \in [n]$ , assuming that  $s = X_i$  and that  $\bar{y}_i$  is a tuple of pairwise distinct variables consisting exactly of the variables in  $X_i$ :

- A CQ  $q_s = \text{Answer}(\bar{y}_i) \text{ :- } R_{j_1}(\bar{u}_{j_1}), \dots, R_{j_p}(\bar{u}_{j_p})$ , where  $\{R_{j_1}(\bar{u}_{j_1}), \dots, R_{j_p}(\bar{u}_{j_p})\}$  is the set of atoms of  $q$  such that  $X_{j_\ell} = X_i$  for each  $\ell \in [p]$ .
- A CQ  $Q_s(\bar{y}_i)$  whose set atoms is the union of those that appear in CQs  $q_{s'}$ , where  $s'$  is a descendant of  $s$  in  $T$  (including  $s$  itself).

Notice that  $Q_s \subseteq q_s$  for each node  $s$  of  $T$ . Moreover, if  $s$  is a (non-leaf) node of  $T$  with children  $s_1, \dots, s_p$ , then the set of atoms of  $Q_s$  is the union of the atoms of  $Q_{s_1}, \dots, Q_{s_p}$ , and the atoms of  $q_s$ .

In what follows, we present a fundamental connection between the evaluation of acyclic CQs and the semijoin operator. To understand this connection, assume that  $q = \text{Answer} \text{ :- } R_1(\bar{u}_1), \dots, R_n(\bar{u}_n)$  is a Boolean ACQ, and suppose that  $T$  is a join tree of  $H_q$  with root  $r = X_\ell$ , for some  $\ell \in [n]$ . Then we have that  $Q_r$  is a CQ of the form  $Q_r = \text{Answer}(\bar{y}_\ell) \text{ :- } R_1(\bar{u}_1), \dots, R_n(\bar{u}_n)$ , which means that  $Q_r$  has the same body as  $q$ . Thus, for every database  $D$ , it holds that

$$q(D) = \text{true} \text{ if and only if } Q_r(D) \neq \emptyset$$

and, therefore, an efficient algorithm for the evaluation of  $Q_r$  can also be used to evaluate  $q$ . We show in the next section that the following proposition gives us such an algorithm. The proposition tells us that  $Q_r$  can be inductively evaluated by computing semijoins while traversing  $T$  in a bottom-up manner, provided that the CQ  $q_s$  has been previously evaluated for every node  $s$  of  $T$ .

### Proposition 19.2

Let  $q$  be a Boolean ACQ,  $T$  a join tree of  $H_q$  and  $D$  a database. Then for every node  $s$  of  $T$ ,

- if  $s$  is a leaf of  $T$ , then  $Q_s(D) = q_s(D)$  and
- otherwise, if the children of  $s$  in  $T$  are  $s_1, \dots, s_p$ , then

$$Q_s(D) = \bigcap_{i=1}^p (q_s(D) \bowtie Q_{s_i}(D)) .$$

*Proof.* Assume that  $q = \text{Answer} \text{ :- } R_1(\bar{u}_1), \dots, R_n(\bar{u}_n)$  is the Boolean ACQ. Therefore, the set of nodes of  $T$  is  $\{X_i \mid i \in [n]\}$ , where each  $X_i$  is the set of variables occurring in  $\bar{u}_i$ . If  $s$  is a leaf of  $T$ , then  $q_s$  and  $Q_s$  are the same CQ and, thus,  $Q_s(D) = q_s(D)$ .

Let us assume then that  $s$  is a non-leaf node of  $T$  with children  $s_1, \dots, s_p$ . Moreover, assume  $s = X_\ell$ ,  $s_1 = X_{k_1}$ ,  $\dots$ ,  $s_p = X_{k_p}$ , where  $\ell, k_1, \dots, k_p$  are pairwise distinct numbers in the set  $[n]$ . Then we have that  $\bar{y}_\ell, \bar{y}_{k_1}, \dots, \bar{y}_{k_p}$  are the tuples of free variables of CQs  $Q_s, Q_{s_1}, \dots, Q_{s_p}$ , respectively. Let us consider first an arbitrary tuple in  $Q_s(D)$ . By definition, such a tuple is of the form  $h(\bar{y}_\ell)$  for some homomorphism  $h$  from  $Q_s$  to  $D$ . It is not hard to see that  $h(\bar{y}_\ell) \in q_s(D) \bowtie Q_{s_i}(D)$  for every  $i \in [p]$ . Indeed,  $h(\bar{y}_\ell) \in q_s(D)$  since  $Q_s \subseteq q_s$ , and  $h(\bar{y}_{k_i}) \in Q_{s_i}(D)$  since the atoms of  $Q_{s_i}$  are contained in those of  $Q_s$ . Moreover,  $h(\bar{y}_\ell)$  and  $h(\bar{y}_{k_i})$  are consistent by definition. We conclude that  $h(\bar{y}_\ell) \in \bigcap_{i=1}^p (q_s(D) \bowtie Q_{s_i}(D))$ .

Let us consider now an arbitrary tuple in  $\bigcap_{i=1}^p (q_s(D) \bowtie Q_{s_i}(D))$ . By definition, such a tuple is of the form  $h(\bar{y}_\ell)$  for some homomorphism  $h$  from  $q_s$  to  $D$ . Moreover, for each  $i \in [p]$ , there is a homomorphism  $h_i$  from  $Q_{s_i}$  to  $D$  such that  $h(\bar{y}_\ell)$  and  $h_i(\bar{y}_{k_i})$  are consistent; i.e., they have the same values on positions where  $\bar{y}_\ell$  and  $\bar{y}_{k_i}$  have common variables. We claim that  $h' = h \cup h_1 \cup \dots \cup h_p$  is a well-defined homomorphism from  $Q_s$  to  $D$ . Since  $h'(\bar{y}_\ell) = h(\bar{y}_\ell)$ , this shows that  $h(\bar{y}_\ell) \in Q_s(D)$  as desired.

We first show that  $h'$  is well defined. Take an arbitrary variable  $y$  in  $Q_s$ . If  $y$  occurs only in  $q_s$  but not in any of the CQs  $Q_{s_i}$  (for  $i \in [p]$ ), or if  $y$  occurs only in one of the CQs  $Q_{s_i}$  (for  $i \in [p]$ ) but not in  $q_s$ , then clearly  $h'(y)$  is well defined. There are two other possibilities:  $y$  occurs in  $q_s$  and in  $Q_{s_i}$ , for some  $i \in [p]$ , or  $y$  occurs in  $Q_{s_i}$  and  $Q_{s_j}$ , for some  $i, j \in [p]$  with  $i \neq j$ . We only consider the latter case since the former can be handled analogously. By

definition of join trees, the nodes in  $T$  that contain  $y$  are connected, which means that  $y \in s \cap s_i \cap s_j$ . Therefore, we conclude that  $h_i(y) = h_j(y) = h(y)$  since  $h(\bar{y}_\ell)$  is consistent with both  $h_i(\bar{y}_{k_i})$  and  $h_j(\bar{y}_{k_j})$ .

We now prove that  $h'$  is a homomorphism from  $Q_s$  to  $D$ . Take an arbitrary atom  $R(\bar{z})$  in  $Q_s$ . Then,  $R(h'(\bar{z})) = R(h(\bar{z}))$  or  $R(h'(\bar{z})) = R(h_i(\bar{z}))$  for some  $i \in [p]$ . Thus,  $R(h'(\bar{z})) \in D$  because  $h$  and  $h_i$  are homomorphisms. This concludes the proof of the proposition.  $\square$

## Yannakakis's Algorithm

Yannakakis's algorithm uses the conditions in Proposition 19.2 to evaluate a Boolean ACQ, as shown in Algorithm 6. The correctness of the algorithm follows from Proposition 19.2—which justifies the correctness of the inductive computation carried out in the while loop—and the fact that the atoms of  $Q_r$  are precisely those of  $q$ , from which we conclude that  $Q_r(D) \neq \emptyset$  if and only if  $q(D) = \text{true}$ .

---

### Algorithm 6 YANNAKAKIS( $q, D$ )

---

**Input:** A Boolean ACQ  $q$  and a database  $D$

**Output:**  $q(D)$

```

1:  $T :=$  a join tree of  $H_q$ 
2:  $N :=$  the set of nodes of  $T$ 
3:  $r :=$  the root of  $T$ 
4: while  $N \neq \emptyset$  do
5:   Choose  $s \in N$  such that no child of  $s$  is in  $N$ 
6:   Compute  $q_s(D)$ 
7:   if  $s$  is a leaf of  $T$  then
8:      $Q_s(D) := q_s(D)$ 
9:   else
10:    Let  $s_1, \dots, s_p$  be the children of  $s$  in  $T$ 
11:     $Q_s(D) := \bigcap_{i=1}^p (q_s(D) \times Q_{s_i}(D))$ 
12:     $N := N - \{s\}$ 
13: if  $Q_r(D) \neq \emptyset$  then
14:   return true
15: else
16:   return false

```

---

We now analyze the complexity of the algorithm. We first notice that Proposition 18.8 tell us that a join tree  $T$  of  $H_q$  can be computed in time  $O(\|q\|)$ . We show next that the remainder of the algorithm can be implemented in time  $O(\|D\| \cdot \log \|D\| \cdot \|q\|)$ . To see why this is the case, we need the following observation (see Exercise 3.3): the time needed to compute  $q(D) \times q'(D)$ , given  $q(D)$  and  $q'(D)$ , is  $O(N \log N)$  with  $N = \|q(D)\| + \|q'(D)\|$ . In

particular, then, each  $q_s(D)$ , for a node  $s$  in  $T$ , can be computed in time  $O(\|D\| \cdot \log \|D\| \cdot \|q_s\|)$ . Therefore, the collection of all queries  $q_s(D)$ , for  $s$  a node in  $T$ , can be computed in time  $O(\|D\| \cdot \log \|D\| \cdot \|q\|)$ .

Now, if  $s$  is a node of  $T$  with children  $s_1, \dots, s_p$ , we can compute  $Q_s(D) = \bigcap_{1 \leq i \leq p} q_{s_i}(D) \times Q_{s_i}(D)$  in time  $O(\|D\| \cdot \log \|D\| \cdot p)$ . This follows from the fact that  $\|q_s(D)\| \leq \|D\|$  and  $\|Q_{s_i}(D)\| \leq \|q_{s_i}(D)\| \leq \|D\|$ , for each  $i \in [p]$ . Therefore, we can inductively compute the collection of all queries  $Q_s$ , for  $s$  a node in  $T$ , in time  $O(\|D\| \cdot \log \|D\| \cdot \|T\|) = O(\|D\| \cdot \log \|D\| \cdot \|q\|)$ .

In summary, we obtain the following result:

**Theorem 19.3**

ACQ-Evaluation can be solved in time  $O(\|D\| \cdot \log \|D\| \cdot \|q\|)$ .

*Proof.* We already proved in the preceding analysis that the theorem holds for Boolean ACQs. Assume now that we are given a non-Boolean ACQ  $q = \text{Answer}(\bar{x}) :- R_1(\bar{u}_1), \dots, R_n(\bar{u}_n)$ , a database  $D$ , and a tuple  $\bar{a}$  over  $\text{Const}$  of the same arity than  $\bar{x}$ . We want to check whether  $\bar{a} \in q(D)$ . We start by turning  $q(\bar{x})$  into a Boolean ACQ by simultaneously replacing in  $q(\bar{x})$  each free variable  $x_i$  in  $\bar{x} = (x_1, \dots, x_k)$  by its corresponding value  $a_i$  in  $\bar{a} = (a_1, \dots, a_k)$ . We denote this Boolean CQ as  $q_{\bar{a}}$ . Clearly,  $q_{\bar{a}}$  is acyclic and, in addition,  $\bar{a} \in q(D)$  if and only if  $q_{\bar{a}}(D) = \text{true}$ .  $\square$

## The Consistency Algorithm

While Yannakakis's algorithm uses a join tree of an acyclic CQ  $q$  in order to evaluate  $q$  over a database  $D$  in time  $O(\|D\| \cdot \log \|D\| \cdot \|q\|)$ , if we only aim for tractability then there is no need for such a join tree to be computed. In fact, below we present an algorithm that evaluates  $q$  on  $D$  in polynomial time, only by holding the *promise* that  $q$  is acyclic (i.e., that a join tree of  $q$  exists). The design of such an algorithm is based on a simple consistency criterion, established in the following proposition, which characterizes when  $q(D) = \text{true}$  for a Boolean ACQ  $q$  and a database  $D$ .

**Proposition 19.4: Consistency Property**

Let  $q = \text{Answer} :- R_1(\bar{u}_1), \dots, R_n(\bar{u}_n)$  be a Boolean ACQ and  $D$  a database, and  $q_i = \text{Answer}(\bar{x}_i) :- R_i(\bar{u}_i)$  be a CQ such that  $\bar{x}_i$  is the tuple obtained from  $\bar{u}_i$  by removing constants, for each  $i \in [n]$ . Then the following are equivalent:

1.  $q(D) = \text{true}$ .
2. There are nonempty sets  $S_1 \subseteq q_1(D), \dots, S_n \subseteq q_n(D)$  such that

$$S_i = S_i \times S_j \text{ for all } i, j \in [n].$$

That is, each tuple in  $S_i$  is consistent with some tuple in  $S_j$  for all  $i, j \in [n]$ .

*Proof.* Assume first that  $q(D) = \mathbf{true}$ , i.e., there is a homomorphism  $h$  from  $q$  to  $D$ . In this case, we can choose  $S_i$  to be  $\{h(\bar{x}_i)\}$ , for each  $i \in [n]$ .

For the other direction, assume that nonempty sets  $S_1, \dots, S_n$  as described in Item 2 exist. Let  $T$  be an arbitrary join tree of  $H_q$ . One can then prove by induction the following for each node  $s$  of  $T$  (recall the notation in Algorithm 6):

If  $s$  is the set  $X_i$  of variables occurring in  $\bar{x}_i$ , for  $i \in [n]$ , then  $S_i \subseteq Q_s(D)$  (see Exercise 3.5).

In particular, if the root  $r$  of  $T$  is the set  $X_j$  of variables occurring in  $\bar{x}_j$ , for  $j \in [n]$ , then  $S_j \subseteq Q_r(D)$ . Therefore, since  $S_j$  is nonempty, we conclude that  $Q_r(D)$  is also nonempty. This implies that there is at least one homomorphism from  $Q_r$  to  $D$ . But the atoms of  $Q_r$  and  $q$  are the same by definition, and thus  $q(D) = \mathbf{true}$ .  $\square$

We are ready to present the consistency algorithm, which can be understood as a greatest fixed-point computation that checks for the existence of nonempty sets  $S_1, \dots, S_n$  as described in Item 2 of Proposition 19.4.

---

**Algorithm 7** CONSISTENCY( $q, D$ )

---

**Input:** A Boolean ACQ  $q = \text{Answer} :- R_1(\bar{u}_1), \dots, R_n(\bar{u}_n)$  and a database  $D$

**Output:**  $q(D)$

```

1:  $S_i := q_i(D)$ , for each  $i \in [n]$ 
2: while  $S_i \neq S_i \times S_j$  for some  $i, j \in [n]$  do
3:    $S_i := S_i \times S_j$ 
4: if  $S_i \neq \emptyset$  for every  $i \in [n]$  then
5:   return true
6: else
7:   return false

```

---

The algorithm initializes  $S_i$  to be  $q_i(D)$ , for each  $i \in [n]$ . It then iteratively deletes every tuple in  $S_i$  that is not consistent with a tuple in  $S_j$ , for some  $j \in [n]$ . If some  $S_i$  becomes empty during this procedure, the algorithm declares  $q(D) = \mathbf{false}$ . Otherwise,  $q(D) = \mathbf{true}$ . The algorithm runs in polynomial time, but not in time  $O(\|D\| \cdot \log \|D\| \cdot \|q\|)$  as Yannakakis's algorithm. Next, we establish that it is correct.



**Proposition 19.5**

Given a Boolean ACQ  $q$  and a database  $D$ , we have that  $q(D) = \text{true}$  if and only if  $\text{CONSISTENCY}(q, D) = \text{true}$ .

We leave the proof of Proposition 19.5 as an exercise for the reader (see Exercise 3.6).

**Acyclicity of the Core**

It is known that there are Boolean CQs that are not acyclic, yet their core is acyclic. (The reader is asked to prove this fact in Exercise 3.7). Interestingly, the consistency algorithm continues being correct for the evaluation problem of such a class of CQs.

**Proposition 19.6**

Let  $q$  be a Boolean CQ whose core is acyclic and  $D$  a database. Then we have that  $q(D) = \text{true}$  if and only if  $\text{CONSISTENCY}(q, D) = \text{true}$ .

*Proof.* Let  $q$  be a CQ whose core  $q'$  is acyclic. It is then the case that:

$$\begin{array}{lll} q(D) = \text{true} & \text{if and only if} & q'(D) = \text{true} \\ & \text{if and only if} & \text{CONSISTENCY}(q', D) = \text{true} \\ & \text{if and only if} & \text{CONSISTENCY}(q, D) = \text{true}. \end{array}$$

The first equivalence holds since  $q \equiv q'$ , the second one since  $q'$  is acyclic (based on Proposition 19.5), and the last one given the fact that if there exists a homomorphism from  $q$  to  $q'$  and  $\text{CONSISTENCY}(q', D) = \text{true}$ , then  $\text{CONSISTENCY}(q, D) = \text{true}$  (see Exercise 3.9).  $\square$

As a corollary, we obtain the following, for the class **ACoreCQ** of CQs that have an acyclic core.

**Theorem 19.7**

**ACoreCQ-Evaluation** is in PTIME.



## Generalized Hypertreewidth

A significant number of real-world CQs are not acyclic, but are in some sense “nearly-acyclic”. Bounded generalized hypertreewidth provides a natural formalization of the notion of near-acyclicity that retains several of the good properties of acyclicity. In particular, CQs of bounded generalized hypertreewidth can be evaluated in polynomial time. Importantly, most of the CQs found in real-world situations are of small generalized hypertreewidth, thus establishing the practical relevance of the concept.

### The notion of generalized hypertreewidth

The definition of generalized hypertreewidth is based on the important notions of *tree decompositions* and *generalized hypertree decompositions*. Let  $H = (V, E)$  be a hypergraph. A *tree decomposition* of  $H = (V, E)$  is a pair  $(T, \chi)$ , formed by a tree  $T$  and a mapping  $\chi$  that assigns a subset of the nodes in  $V$  to each node  $s \in T$ , for which the following statements hold:

1. For each edge  $e \in E$ , there is a node  $s \in T$  such that  $e \subseteq \chi(s)$ .
2. For each node  $v \in V$ , the set of nodes  $s \in T$  for which  $v$  occurs in  $\chi(s)$  is connected.

#### Definition 22.1: Generalized Hypertree Decomposition

A *generalized hypertree decomposition* of  $H = (V, E)$  is a triple  $(T, \chi, \lambda)$  such that:

1.  $(T, \chi)$  is a tree decomposition of  $H$ .
2.  $\lambda$  is a mapping that assigns a subset of the hyperedges in  $E$  to each node  $s \in T$ .

3. For each node  $s \in T$ , it is the case that  $\chi(s) \subseteq \bigcup_{e \in \lambda(s)} e$ .

In other words, a generalized hypertree decomposition  $(T, \chi, \lambda)$  of  $H$  extends the tree decomposition  $(T, \chi)$  by *covering* each set  $\chi(s)$  of nodes, for  $s \in T$ , with a set  $\lambda(s)$  of hyperedges from  $H$ .

The *width* of a node  $s$  in the generalized hypertree decomposition  $(T, \chi, \lambda)$  is the number of atoms in  $\lambda(s)$ . The width of  $(T, \chi, \lambda)$  is the maximal width of the nodes of  $T$ . The *generalized hypertreewidth* of a hypergraph is the minimum width of its generalized hypertree decompositions.

The notion of generalized hypertreewidth of a CQ is defined as follows:

**Definition 22.2: Generalized hypertreewidth**

The generalized hypertreewidth of a CQ  $q$  corresponds to the generalized hypertreewidth of its associated hypergraph  $H_q$ .

For each fixed  $k \geq 1$ , we write  $\text{GHW}(k)$  for the set of CQs whose generalized hypertreewidth is at most  $k$ . That is, CQs in  $\text{GHW}(k)$  are those which admit generalized hypertree decompositions of the form  $(T, \chi, \lambda)$  in which each set of variables  $\chi(s)$ , for  $s \in T$ , is covered by a set  $\lambda(s)$  of at most  $k$  edges in  $H_q$ .

It is easy to prove that  $\text{GHW}(k) \subsetneq \text{GHW}(k+1)$ , for each  $k \geq 1$ . As an example, let us recall the CQ

$$q(x, y) :- R(x, y, z), S(y, z), S(y, w, w'), T(z, z'), T(z, w),$$

which was introduced in Chapter 18. The CQ  $q$  is in  $\text{GHW}(2)$ . This is witnessed by the generalized hypertree decomposition  $(T, \chi, \lambda)$  such that  $T$  consists of two nodes  $r$  and  $t$  for which  $\chi(r) = \{x, y, z, w, w'\}$ ,  $\lambda(r) = \{\{x, y, z\}, \{y, w, w'\}\}$ ,  $\chi(t) = \{z, z'\}$ , and  $\lambda(t) = \{\{z, z'\}\}$ . On the other hand,  $q(x, y)$  is not in  $\text{GHW}(1)$ . This is because CQs in  $\text{GHW}(1)$  are acyclic (see Proposition 22.3 below) and we know from Example 18.5 that  $q$  is not acyclic.

By slightly abusing notation, in the rest of the section we assume for simplicity that if  $q$  is a CQ and  $(T, \lambda, \chi)$  is a generalized hypertree decomposition of  $H_q$ , then  $\chi(s)$  corresponds to a set of atoms from  $q$ , for each node  $s \in T$ .

## Bounded generalized hypertreewidth and acyclicity

The notion of bounded generalized hypertreewidth properly subsumes acyclicity. More in particular, acyclic CQs coincide with the CQs of generalized hypertreewidth one:

**Proposition 22.3**

GHW(1) is the class of acyclic CQs.

*Proof.* It is easy to see that each acyclic CQ is in GHW(1). In fact, let  $q := R_1(\bar{u}_1), \dots, R_m(\bar{u}_m)$  be an acyclic CQ and  $T$  an arbitrary join tree of  $q$ . Then  $T$  can be turned into a generalized hypertree decomposition of  $q$  of width one as follows. For each node  $s \in T$  which is associated with the set  $X_i$  of variables that are mentioned in  $\bar{u}_i$ , for  $1 \leq i \leq m$ , we define  $\chi(s) = X_i$  and  $\lambda(s) = \{R_i(\bar{u}_i)\}$ . On the other hand, assume that  $q$  is in GHW(1) and  $(T, \chi, \lambda)$  is a generalized hypertree decomposition of  $q$  of width one. From  $(T, \chi, \lambda)$  one can construct a join tree of  $q$  as follows. Each node  $s \in T$  is associated with the set of variables mentioned in the single atom in  $\lambda(s)$ . If two such nodes end up being associated with the same set of variables, we simply delete one of them (provided that it is not the root). Clearly, then, each node  $s \in T$  is associated with the set  $X_i$  of variables mentioned in some tuple  $\bar{u}_i$ , for  $1 \leq i \leq m$ , and no two distinct nodes are associated with the same  $X_i$ . Still, there might be several  $1 \leq i \leq m$  such that  $X_i$  is associated with no node in  $T$ . However, by definition of generalized hypertree decomposition, for each such an  $1 \leq i \leq m$  there must be a node  $s \in T$  such that  $X_i \subseteq \chi(s) \subseteq X_j$ , assuming that  $s$  is associated with the set  $X_j$  of variables mentioned in  $\bar{u}_j$ , for  $1 \leq j \leq m$ . One can then create a new children  $s'$  of  $s$  which is associated with  $X_i$ . This construction yields a join tree of  $q$ .  $\square$

**Tractable evaluation based on consistency**

It is shown next that CQs of bounded generalized hypertreewidth can be evaluated in polynomial time by extending the consistency criterion for acyclic CQs developed in Proposition 19.4. For reasons explained in Chapter 18, we concentrate on Boolean CQs.

Let  $q := R_1(\bar{u}_1), \dots, R_m(\bar{u}_m)$  be a Boolean CQ. For each  $S \subseteq \{1, \dots, m\}$ , let us define  $q_S(\bar{x}_S)$  to be the CQ whose set of atoms is  $\{R_i(\bar{u}_i) \mid i \in S\}$  and  $\bar{x}_S$  is a tuple that consists precisely of the variables mentioned in such atoms. Our consistency criterion establishes the following:

**Proposition 22.4**

Fix  $k \geq 1$ . Let  $q := R_1(\bar{u}_1), \dots, R_m(\bar{u}_m)$  be a Boolean CQ in GHW( $k$ ) and  $D$  a database. Then the following are equivalent:

1.  $q(D) = \text{true}$ .

2. For each  $S \subseteq \{1, \dots, m\}$  with at most  $k$  elements, there is a nonempty set  $\text{Cons}(q_S(D)) \subseteq q_S(D)$  such that, for each  $S' \subseteq \{1, \dots, m\}$  with at most  $k$  atoms it is the case that

$$\text{Cons}(q_S(D)) = \text{Cons}(q_S(D)) \times \text{Cons}(q_{S'}(D)).$$

That is, each tuple in  $\text{Cons}(q_S(D))$  is consistent with some tuple in  $\text{Cons}(q_{S'}(D))$ , for every  $S, S' \subseteq \{1, \dots, m\}$  with at most  $k$  atoms.

*Proof.* Suppose first that  $q(D) = \text{true}$ ; i.e., there is a homomorphism  $h$  from  $q$  to  $D$ . It is not hard to see then that one can choose  $\text{Cons}(q_S(D))$  to be  $\{h(\bar{x}_S)\}$ , for each  $S \subseteq \{1, \dots, m\}$  with at most  $k$  elements. Suppose, on the other hand, that for each  $S \subseteq \{1, \dots, m\}$  with at most  $k$  elements a nonempty set  $\text{Cons}(q_S(D))$  as described in (2) exists. Let  $(T, \chi, \lambda)$  be an arbitrary rooted and directed generalized hypertree decomposition of  $q$  of width  $k$ . This implies that  $\text{Cons}(q_{\lambda(s)}(D))$  is well-defined for each node  $s \in T$ , as  $|\lambda(s)| \leq k$  by definition. One can then prove by induction the following for each node  $s \in T$  (see Exercise 3.14): Let  $X_s$  be the set of variables mentioned in sets of the form  $\chi(s')$ , for  $s'$  a descendant of  $s$  in  $T$  (including  $s$  itself) and  $S(X_s)$  the set of all atoms in  $q$  that only mention variables in  $X_s$ . Then  $\text{Cons}(q_{\lambda(s)}(D)) \subseteq q'_{S(X_s)}(D)$ , where  $q'_{S(X_s)}$  is the CQ that has the same set of atoms than  $q_{S(X_s)}$  but its tuple of free variables is  $\bar{x}_{\lambda(s)}$ . In particular, for the root  $r$  of  $T$  it is the case that  $\text{Cons}(q_{\lambda(r)}(D)) \subseteq q'_{S(X_r)}(D)$ . Therefore, since  $\text{Cons}(q_{\lambda(r)}(D))$  is nonempty one can conclude that  $q'_{S(X_r)}(D)$  is also nonempty. This implies that there is at least one homomorphism from  $q'_{S(X_r)}$  to  $D$ . But the atoms of  $q'_{S(X_r)}$  and  $q$  are the same by definition, and thus  $q(D) = \text{true}$ .  $\square$

As for acyclic CQs, this result allows us to construct a simple greatest-fixed point procedure that checks for the existence of sets  $\text{Cons}(q_S(D))$  as described in item (2) of Proposition 22.4. The algorithm, called  $k$ -CONSISTENCY, is presented next:

---

**Algorithm 9**  $k$ -CONSISTENCY( $q, D$ )

---

**Input:** A Boolean CQ  $q := R_1(\bar{u}_1), \dots, R_m(\bar{u}_m)$  in  $\text{GHW}(k)$  and a database  $D$ .

**Output:** If  $q(D) = \text{true}$ , then  $\text{Cons}(q_S(D)) \neq \emptyset$  for each  $S \subseteq \{1, \dots, m\}$  with at most  $k$  elements. Otherwise, **fail**.

- 1:  $\text{Cons}(q_S(D)) := q_S(D)$ , for each  $S \subseteq \{1, \dots, m\}$  with at most  $k$  elements
  - 2: **while**  $\text{Cons}(q_S(D)) \neq \text{Cons}(q_S(D)) \times \text{Cons}(q_{S'}(D))$  for  $S, S' \subseteq \{1, \dots, m\}$  with at most  $k$  elements **do**
  - 3:      $\text{Cons}(q_S(D)) := \text{Cons}(q_S(D)) \times \text{Cons}(q_{S'}(D))$
  - 4: **if**  $\text{Cons}(q_S(D)) = \emptyset$  for some  $S \subseteq \{1, \dots, m\}$  with at most  $k$  elements **then**
  - 5:     **fail**
-

It is not hard to see that  $k$ -CONSISTENCY runs in polynomial time, for each fixed  $k \geq 1$ . In fact, a naïve analysis shows that  $k$ -CONSISTENCY( $q, D$ ) can be implemented in time  $O(\|D\|^{2k} \cdot \|q\|^{2k})$  based on the following observations:

1. For each  $S \subseteq \{1, \dots, m\}$  with at most  $k$  elements, the value of  $q_S(D)$  can be computed in time  $O(\|D\|^k)$ . Therefore, the initialization step of the algorithm in which  $\text{Cons}(q_S(D))$  is set to be  $q_S(D)$ , for each  $S \subseteq \{1, \dots, m\}$  with at most  $k$  elements, runs in time  $O(\|D\|^k \cdot \|q\|^k)$ .
2. Each further step of the algorithm in which  $\text{Cons}(q_S(D))$  is set to be  $\text{Cons}(q_S(D)) \times \text{Cons}(q_{S'}(D))$ , for each  $S, S' \subseteq \{1, \dots, m\}$  with at most  $k$  elements, takes time  $O(\|D\|^k \cdot \|q\|^k)$ . Since each such a step deletes at least one tuple from some  $\text{Cons}(q_S(D))$ , for an  $S \subseteq \{1, \dots, m\}$  with at most  $k$  elements, the maximum number of steps performed by the algorithm is bounded by  $O(\|D\|^k \cdot \|q\|^k)$ .

The soundness and completeness of  $k$ -CONSISTENCY is established next:

**Proposition 22.5**

Let  $q$  be a Boolean CQ in  $\text{GHW}(k)$  and  $D$  a database. Then:

$$q(D) = \text{true} \iff k\text{-CONSISTENCY}(q, D) \neq \text{fail}.$$

The proof of Proposition 22.5 is left as an exercise for the reader (see Exercise 3.15). As a corollary, one obtains the following:

**Theorem 22.6**

Fix  $k \geq 1$ .  $\text{GHW}(k)$ -Evaluation can be solved in polynomial time  $O(\|D\|^{2k} \cdot \|q\|^{2k})$ .

## Bounded generalized hypertreewidth of the core

Recall that the CONSISTENCY procedure for evaluating acyclic CQs, presented in Chapter 18, continues being sound for the class of CQs whose core is acyclic. The  $k$ -CONSISTENCY algorithm presented above preserves this good behavior, this time with respect to the class of CQs whose core is in  $\text{GHW}(k)$ :

**Proposition 22.7**

Fix  $k \geq 1$ . Let  $q$  be a Boolean CQ whose core is in  $\text{GHW}(k)$  and  $D$  a database. Then:

$$q(D) = \text{true} \iff k\text{-CONSISTENCY}(q, D) \neq \text{fail}.$$

The proof of this result is similar to the proof of Proposition 19.6. As a corollary, one then obtains the following:

**Theorem 22.8**

Fix  $k \geq 1$ . The evaluation problem for the class of CQs whose core is in  $\text{GHW}(k)$  can be solved in polynomial time.

### Computing generalized hypertree decompositions for faster evaluation

CQs in  $\text{GHW}(k)$  can be evaluated in polynomial time, for each fixed  $k \geq 1$ , via the  $k$ -CONSISTENCY procedure. Such a procedure assumes that the input CQ  $q$  is in  $\text{GHW}(k)$ , i.e., that a generalized hypertree decomposition of  $q$  of width  $k$  exists, but no such a decomposition is required to be computed. On the other hand, as we will see later, having access to a generalized hypertree decomposition of width  $k$  with good properties helps improving the cost of evaluation for CQs in  $\text{GHW}(k)$ . This is in line with the case of acyclic CQs, for which we know that computing a join tree allows us to perform evaluation in linear time using Yannakakis's algorithm.

Having access to a generalized hypertree decomposition is not a problem for the case  $k = 1$ ; in fact,  $\text{GHW}(1)$  corresponds to the class of acyclic CQs, and from an acyclic CQ we can always compute a join tree (or, equivalently, a generalized hypertree decomposition), in linear time. Unfortunately, for  $k > 1$  this good property no longer holds as the following result shows:

**Theorem 22.9**

Fix  $k > 1$ . Assuming  $\text{PTIME} \neq \text{NP}$ , there is no polynomial time algorithm that given an input CQ  $q$  computes a generalized hypertree decomposition of width  $k$  whenever  $q \in \text{GHW}(k)$ .

*Proof.* The proof relies on the following difficult result:

**Proposition 22.10**

Fix  $k > 1$ . The problem of checking if a given CQ is in  $\text{GHW}(k)$  is NP-complete.

In fact, assume for the sake of contradiction that for some  $k > 1$  there is a polynomial time algorithm  $\mathcal{A}$  that given an input CQ  $q$  computes a generalized



hypertree decomposition of width  $k$  whenever  $q \in \text{GHW}(k)$ . We show then that there is a polynomial time algorithm  $\mathcal{A}'$  that checks whether a given CQ is in  $\text{GHW}(k)$ , thus contradicting Proposition 22.10. Take an arbitrary CQ  $q$  and run  $\mathcal{A}$  on  $q$ . The algorithm  $\mathcal{A}'$  accepts iff  $\mathcal{A}$  outputs a generalized hypertree decomposition of  $q$  of width  $k$  (the latter can be checked in polynomial time). It is easy to see, then, that  $\mathcal{A}'$  accepts iff  $q$  is in  $\text{GHW}(k)$ .  $\square$

Does this result completely rule out the possibility of using generalized hypertree decompositions for query evaluation? Not necessarily, for the following reasons. It can be proved that if  $q$  is in  $\text{GHW}(k)$ , then there is a generalized hypertree decomposition of  $q$  of width  $k$  with at most  $n$  nodes, where  $n$  is the number of variables in  $q$ . Therefore, in order to check if  $q \in \text{GHW}(k)$ , and, if so, compute a generalized hypertree decomposition of  $q$  of width  $k$  with at most  $n$  nodes, we can do the following: Iterate over all tuples of the form  $(T, \chi, \lambda)$ , where  $T$  is a tree with at most  $n$  nodes,  $\chi$  is a mapping that assigns a subset of the variables in  $q$  to each node  $s \in T$ , and  $\lambda$  is a mapping that assigns a subset of at most  $k$  atoms from  $q$  to each node  $s \in T$ . Then check if any of them is a generalized hypertree decomposition of  $q$ . This takes time  $2^{\|q\|^c}$ , for some integer  $c \geq 1$ . While this algorithm exhibits an exponential behavior, it is not completely impractical as the problem corresponds to a static analysis task for which the input, the CQ  $q$ , is often small.

We can then establish the following:

**Theorem 22.11**

Fix  $k \geq 1$ . Then  $\text{GHW}(k)$ -Evaluation can be solved in time  $2^{\|q\|^c} + O(\|D\|^k \cdot \|q\|)$ , for some integer  $c \geq 1$ .

*Proof.* First compute in time  $2^{\|q\|^c}$ , for some integer  $c \geq 1$ , a generalized hypertree decomposition of  $q$  of width  $k$  with at most  $n$  nodes, where  $n$  is the number of variables in  $q$ . Recall that for each node  $s \in T$  the CQ  $q_{\lambda(s)}$  is defined as follows: the atoms of  $q_{\lambda(s)}$  are precisely those in  $\lambda(s)$  and the free variables of  $q_{\lambda(s)}$  are all the variables that are mentioned in such atoms. Compute then the value of  $q_{\lambda(s)}(D)$ , for each  $s \in T$ . Each  $q_{\lambda(s)}(D)$  can be computed in time  $O(\|D\|^k)$ , and thus computing them all takes time  $O(\|D\|^k \cdot \|T\|)$ , which is  $O(\|D\|^k \cdot \|q\|)$ . By mimicking Yannakakis's algorithm, we inductively compute the values  $Q_{\lambda(s)}(D)$ 's, for  $s$  a node in  $T$ , which are defined as follows:

- If  $s$  is a leaf of  $T$ , then  $Q_{\lambda(s)}(D) = q_{\lambda(s)}(D)$ .
- If  $s$  has children  $s_1, \dots, s_p$ , then  $Q_{\lambda(s)}(D) = \bigcap_{1 \leq i \leq p} q_{\lambda(s)}(D) \bowtie Q_{\lambda(s_i)}(D)$ .

This takes time  $O(\|D\|^k \cdot \|q\|)$  since  $\|q_{\lambda(s)}\|$  is  $O(\|D\|^k)$ , for each node  $s \in T$ , and  $T$  has at most  $n$  nodes. It can be proved then (see Exercise 1) that for each node  $s \in T$  we have that:

$$Q_{\lambda(s)}(D) = q'_{S(X_s)}(D),$$

where  $q'_{S(X_s)}$  is as defined in the proof of Proposition 22.4. In particular, the atoms of  $q'_{S(X_s)}$  are precisely the atoms of  $q$  that only mention variables that appear in the subtree of  $T$  rooted in  $s$ , and the free variables of  $q'_{S(X_s)}$  are those that are mentioned in  $\lambda(s)$ . The algorithm then accepts if  $Q_{\lambda(r)}(D) \neq \emptyset$ , where  $r$  is the root of  $T$ . In fact, from the previous observation it is the case that  $Q_{\lambda(r)}(D) = q'_{S(X_r)}(D)$ . Moreover, clearly the atoms of  $q'_{S(X_r)}$  and  $q$  are the same. It follows then that  $Q_{\lambda(r)}(D) \neq \emptyset$  if and only if there is a homomorphism from  $q$  to  $D$ , i.e.,  $q(D) = \mathbf{true}$ .  $\square$

The bound for evaluation of CQs in  $\text{GHW}(k)$  obtained in Theorem 22.11 is better, in many practical situations, than the one offered by the consistency algorithm. In fact, from Theorem 22.11 we obtain that the problem can be solved in time  $2^{\|q\|^c} + O(\|D\|^k \cdot \|q\|)$ . This is better than the  $O(\|D\|^{2k} \cdot \|q\|^{2k})$  obtained by applying the  $k$ -CONSISTENCY algorithm whenever  $2^{\|q\|^c}$  is  $O(\|D\|^{2k})$ . But this is not uncommon, as often the database  $D$  is very large and the CQ  $q$  is orders of magnitude smaller.

## The Necessity of Bounded Treewidth

If a class  $\mathcal{C}$  of CQs has bounded generalized hypertreewidth, then the evaluation problem for  $\mathcal{C}$  can be solved in polynomial time. Moreover, this positive behavior continues to hold even if  $\mathcal{C}$  itself does not have bounded generalized hypertreewidth, but the class  $\mathcal{C}_{\text{core}}$  of all cores of CQs in  $\mathcal{C}$  does. More formally, it follows from Proposition 22.7 that if  $\mathcal{C}_{\text{core}}$  satisfies that there is a  $k \geq 1$  such that every CQ  $q \in \mathcal{C}_{\text{core}}$  is in  $\text{GHW}(k)$ , then the evaluation problem for  $\mathcal{C}$  can be solved in polynomial time.

A crucial question at this stage is whether this notion exhausts the space of tractability for CQ evaluation. That is, whether for every class  $\mathcal{C}$  of CQs the following are equivalent:

1. Evaluation for  $\mathcal{C}$  can be solved in polynomial time.
2.  $\mathcal{C}_{\text{core}}$  has bounded generalized hypertreewidth.

Perhaps not surprisingly, it can be shown that this is not the case in general; in fact, there are more general notions of bounded CQ-width, e.g., bounded *fractional hypertreewidth*, that lead to tractability of CQ evaluation and properly extend the notion of bounded generalized hypertreewidth.

On the other hand, it is shown in this section that there is one important scenario in which conditions 1 and 2 expressed above are equivalent (at least under widely-held complexity theoretical assumptions); namely, when the arity of the underlying schemas of the CQs in  $\mathcal{C}$  is fixed in advance. This includes the important case in which all CQs in  $\mathcal{C}$  come from the same schema. In other words, notions such as bounded fractional hypertreewidth, that ensure tractability of CQ evaluation, properly extend bounded generalized hypertreewidth only when schemas of unbounded arity are allowed.

The equivalence of conditions 1 and 2 over schemas of fixed arity is obtained by proving that if  $\mathcal{C}$  is a class of CQs over bounded arity schemas such that  $\mathcal{C}_{\text{core}}$  is not of bounded generalized hypertreewidth, then the evaluation problem for  $\mathcal{C}$  is W[1]-complete. Thus, under the standard assumption that W[1]-complete problems are not tractable, one can conclude the evaluation

problem for  $\mathcal{C}$  cannot be solved in polynomial time. But not only that, under the assumption that  $\text{FPT} \neq \text{W}[1]$  one also obtains a stronger result: if  $\mathcal{C}$  is a class of CQs over fixed arity schemas, then the evaluation problem for  $\mathcal{C}$  is tractable if and only if the evaluation problem for  $\mathcal{C}$  is fixed-parameter tractable if and only if  $\mathcal{C}_{\text{core}}$  has bounded generalized hypertreewidth. In other words, at least in this restricted scenario the notion of fixed parameter tractability does not add to the usual notion of tractability.

## CQs of bounded treewidth

Over schemas of fixed arity, the notion of bounded generalized hypertreewidth collapses to the notion of bounded *treewidth*. The concept of treewidth has received vast attention in the context of graph theory, where it is used as a measure for the degree of cyclicity of a graph. In particular, the class of trees corresponds precisely to the class of undirected graphs of treewidth one, while the class of all  $k$ -cliques, for  $k \geq 1$ , has unbounded treewidth (in particular, the  $k$ -clique has treewidth  $k - 1$  for each  $k > 1$ ). The notion of treewidth of a CQ corresponds to a natural generalization of the notion of treewidth of a graph to the case when schemas are arbitrary. The main result this section will be presented in terms of bounded treewidth, instead of bounded generalized hypertreewidth, as the former notion more naturally fits in the scenario in which schemas are of fixed arity.

Recall that a tree decomposition of a hypergraph  $H = (V, E)$  is a pair  $(T, \chi)$ , formed by a tree  $T$  and a mapping  $\chi$  that assigns a subset of the nodes in  $V$  to each node  $s \in T$ , for which the following statements hold:

1. For each hyperedge  $e \in E$ , there is a node  $s \in T$  such that  $e \subseteq \chi(s)$ .
2. For each node  $v \in V$ , the set of nodes  $s \in T$  for which  $v$  occurs in  $\chi(s)$  is connected.

The *width* of  $(T, \chi)$  is  $(\max_{s \in T} |\chi(s)|) - 1$ . The treewidth of  $H$  is the minimum width of its tree decompositions. The treewidth of a CQ  $q$  is the treewidth of its associated hypergraph  $H_q$ . We denote by  $\text{TW}(k)$  the set of all CQs of treewidth at most  $k$ , for  $k \geq 1$ .

The treewidth of a CQ can equivalently be defined as the treewidth of the graph that represents which variables occur together in some atom of  $q$ . Formally, let  $q$  be a CQ and  $G_q$  the *Gaifman graph* of  $q$ ; this is the simple graph whose nodes are the variables of  $q$  and there is an edge between distinct variables  $x$  and  $y$  if and only if there is some atom in  $q$  that mentions both  $x$  and  $y$ . As stated next,  $q$  and  $G_q$  have the same treewidth (the proof of this result is left as an exercise):

**Lemma 23.1.** *Let  $q$  be a CQ. The treewidth of  $q$  is precisely the treewidth of its Gaifman graph  $G_q$ .*

It is easy to see that the notion of bounded generalized hypertreewidth properly extends the notion of bounded treewidth. As an example, consider the class formed by all boolean CQs  $q_n$ , for  $n \geq 3$ , defined as follows:

$$q_n \leftarrow \bigwedge_{1 \leq i, j \leq n} E(x_i, x_j), T_n(x_1, \dots, x_n),$$

where  $\bigwedge_{1 \leq i, j \leq n} E(x_i, x_j)$  is a shortening for the fact that all atoms of the form  $E(x_i, x_j)$ , for  $1 \leq i, j \leq n$ , are in  $q_n$ . Notice that  $\{q_n \mid n \geq 3\} \subseteq \text{GHW}(1)$  as every CQ  $q_n$  in  $\mathcal{C}$  is acyclic: this is witnessed by the generalized hypertree decomposition that contains a single node labeled with all variables in  $\{x_1, \dots, x_n\}$  that is covered by the single atom  $T_n(x_1, \dots, x_n)$ . On the other hand, it is the case that  $q_{n+1} \notin \text{TW}(n)$ , for each  $n > 1$ . This follows directly from the fact that the treewidth of the  $(n+1)$ -clique is exactly  $n$ .

Notice, on the other hand, that there is no bound on the arity of the schemas over which the CQs in  $\{q_n \mid n \geq 3\}$  are defined (as  $q_n$  contains the  $n$ -ary atom  $T_n(x_1, \dots, x_n)$ ). As mentioned above, this is necessary since over fixed arity schemas the notions of bounded treewidth and bounded generalized hypertreewidth coincide. This is formally stated below:

**Lemma 23.2.** *Let  $q$  be a CQ defined over a schema all of whose relation symbols have arity at most  $c$ , for  $c \geq 1$ . Then for every  $k \geq 1$ :*

- $q \in \text{GHW}(k)$  implies  $q \in \text{TW}(ck - 1)$ .
- $q \in \text{TW}(k)$  implies  $q \in \text{GHW}(k + 1)$ .

*In particular, if  $\mathcal{C}$  is a class of CQs defined over fixed arity schemas, then  $\mathcal{C}$  is of bounded treewidth iff  $\mathcal{C}$  is of bounded generalized hypertreewidth.*

The proof of this result is left as an easy exercise for the reader.

## The main result

The main result of this section, which is the characterization of the tractable classes of CQs over fixed arity schemas in terms of the notion of bounded treewidth, is stated next. Only the version for boolean CQs is presented, but there is a simple extension of the result that characterizes tractability for arbitrary CQs (see Exercise 5):

### Theorem 23.3

Assume that  $\text{FPT} \neq \text{W}[1]$ . Let  $\mathcal{C}$  be a recursively enumerable class of boolean CQs such that there is a bound on the arity of the relation symbols that are mentioned by CQs in  $\mathcal{C}$ . Then the following are equivalent:

1.  $\mathcal{C}$ -Evaluation can be solved in polynomial time.
2.  $\mathcal{C}$ -Evaluation is FPT.

### 3. $\mathcal{C}_{\text{core}}$ has bounded treewidth.

The restriction on  $\mathcal{C}$  to be recursively enumerable can be removed if one assumes a stronger complexity theoretical assumption; namely, that the *non-uniform* versions of FPT and W[1] are also different.

## Overall idea behind the proof of Theorem 23.3

The implication from (1) to (2) is straightforward. The implication from (3) to (1) follows from Proposition 22.7 and Lemma 23.2. The implication from (2) to (3) is proved next via the contrapositive. To do this, it is shown that if  $\mathcal{C}_{\text{core}}$  does not have bounded treewidth, then the evaluation problem for  $\mathcal{C}_{\text{core}}$  is W[1]-hard under fpt-reductions and, therefore, not in FPT based on the assumption that  $\text{FPT} \neq \text{W}[1]$ .

Let  $\mathcal{C}$  be a recursively enumerable class of CQs over fixed arity schemas such that  $\mathcal{C}_{\text{core}}$  is not of bounded treewidth. The proof constructs an fpt-reduction from the parameterized problem  $p\text{-CLIQUE}$  to  $p\text{-}\mathcal{C}\text{-Evaluation}$ . Recall that  $p\text{-CLIQUE}$  is the problem of given a simple graph  $G = (V, E)$  and an integer  $k \geq 1$ , decide if  $G$  has a  $k$ -clique using  $k$  as a parameter. As mentioned in Chapter 14, the problem  $p\text{-CLIQUE}$  is W[1]-complete under fpt-reductions.

The construction makes use of a deep result in graph theory presented below. The  $(n \times m)$ -grid is the undirected graph whose vertices are the pairs  $(i, j)$ , for  $1 \leq i \leq n$  and  $1 \leq j \leq m$ , and whose set of edges is:

$$\begin{aligned} \{ \{(i, j), (i + 1, j)\} \mid 1 \leq i < n, 1 \leq j \leq m \} \cup \\ \{ \{(i, j), (i, j + 1)\} \mid 1 \leq i \leq n, 1 \leq j < m \} \}. \end{aligned}$$

It can be proved that the  $(n \times n)$ -grid has treewidth  $n$ , for each  $n > 1$ . A *minor* of a simple graph  $G$  is a graph  $H$  can be obtained from a subgraph  $G'$  of  $G$  by contracting edges. Then:

### Theorem 23.4: Excluded Grid Theorem

There is a function  $t : \mathbb{N} \rightarrow \mathbb{N}$ , such that for every  $k \geq 1$  and simple graph  $G$  of treewidth at least  $t(k)$  it is the case that  $G$  contains a  $(k \times K)$ -grid as a minor, for  $K = \binom{k}{2}$ .

Let  $G = (V, E)$  and  $H = (V', E')$  be simple graphs. A *minor map* from  $H$  to  $G$  is a mapping  $\mu : V' \rightarrow 2^V$  that satisfies the following:

- (M1) If  $n$  is a node of  $H$ , then  $\mu(n)$  is a connected and nonempty subset of the nodes of  $G$ .
- (M2) The sets of the form  $\mu(n)$ , for  $n$  a node of  $H$ , are pairwise disjoint.

- (M3) For every edge  $(n_1, n_2) \in E'$ , there exist nodes  $n'_1 \in \mu(n_1)$  and  $n'_2 \in \mu(n_2)$  such that  $(n'_1, n'_2) \in E$ .

A minor map is said to be onto if, in addition, the sets of the form  $\mu(n)$ , for  $n$  a node of  $H$ , define a partition of  $V$ .

The following is left as an easy exercise for the reader (see Exercise 4):

**Lemma 23.5.**  *$H$  is a minor of  $G$  if and only if there is a minor map  $\mu$  from  $H$  to  $G$ . If, in addition,  $G$  is connected, then there exists an onto minor map  $\mu$  from  $H$  to  $G$ .*

The reduction from  $p$ -CLIQUE to  $p$ -C-Evaluation is explained next. Consider an input for  $p$ -CLIQUE given by the pair  $(G, k)$ , where  $G$  is a simple graph and  $k \geq 1$  is an integer. Since  $\mathcal{C}_{\text{core}}$  is not of bounded treewidth, there is some CQ  $q \in \mathcal{C}$  whose core  $q'$  has treewidth at least  $t(k)$ . This means that there is at least one connected component  $q^*$  of  $q'$  whose treewidth is at least  $t(k)$ . By Lemma 23.2 then, the Gaifman graph  $G_{q^*}$  of  $q^*$  has treewidth at least  $t(k)$ , and from the Excluded Grid Theorem it is the case that the  $(k \times K)$ -grid is a minor of  $G_{q^*}$ , where  $K = \binom{k}{2}$ . Notice that  $q^*$  is also a core. Moreover, since the  $(k \times K)$ -grid is a minor of  $G_{q^*}$  and  $G_{q^*}$  is connected, Lemma 23.5 implies that there exists an onto minor map  $\mu$  from the  $(k \times K)$ -grid to  $G_{q^*}$ .

Based on  $G$ ,  $q^*$ , and  $\mu$ , the proof constructs a database  $D = D(G, q^*, \mu)$  over the schema of  $q^*$  such that:

$$q^*(D) = \text{true} \iff G \text{ contains a } k\text{-clique.}$$

Let us assume that  $q' \setminus q^*$  is the CQ that is obtained from  $q'$  by removing all atoms in the connected component  $q^*$ . Let us denote by  $D'$  the disjoint union of  $D$  and (the canonical database of)  $q' \setminus q^*$ . Notice that because  $q'$  is a core and  $q^*$  is connected, it must be the case that if  $h$  is a homomorphism from  $q'$  to  $D'$  then  $h$  maps  $q^*$  to  $D$ . Therefore:

$$q'(D') = \text{true} \iff q^*(D) = \text{true} \iff G \text{ contains a } k\text{-clique.}$$

Since  $q'$  is the core of  $q$ , one concludes that:

$$q(D') = \text{true} \iff G \text{ contains a } k\text{-clique.}$$

The construction of  $D = D(G, q^*, \mu)$  and the proof that  $q^*(D) = \text{true}$  if and only if  $G$  contains a  $k$ -clique are presented later. To conclude that the reduction is indeed an fpt-reduction, it is necessary to further establish the following facts:

1. There is a computable function  $g : \mathbb{N} \rightarrow \mathbb{N}$  such that  $\|q\| \leq g(k)$ .
2. There is a computable function  $f : \mathbb{N} \rightarrow \mathbb{N}$  such that the pair  $(D', q)$  can be constructed in time  $f(k) \cdot \|G\|^{O(1)}$  from  $(G, k)$ .

Condition 1 holds since  $q$  depends exclusively on  $k$  and  $\mathcal{C}$  is recursively enumerable. To show that condition 2 also holds, it is sufficient to show the following:

3. There is a computable function  $f : \mathbb{N} \rightarrow \mathbb{N}$  such that  $D = D(G, q^*, \mu)$  can be constructed in time  $f(k) \cdot \|G\|^{O(1)}$  from  $(G, k)$ .

This is because  $D'$  is the disjoint union of  $D$  and  $q' \setminus q^*$  and the latter is computable from  $q$ . Condition 3 is established during the construction of  $D$ .

### The construction of $D(G, q^*, \mu)$

Recall that the set  $\{1, \dots, n\}$  is denoted  $[n]$ , for each  $n \geq 1$ . Let  $(G, k)$  be an input to  $p$ -CLIQUE, and assume that  $G = (V, E)$  and  $K = \binom{k}{2}$ . It is convenient to interchangeably interpret the columns of the  $(k \times K)$ -grid as elements of  $[K]$  and unordered pairs of elements over  $[k]$ . For that, let us define an arbitrary bijection  $\varphi$  from the set  $[K]$  to the set of all unordered pairs of elements over  $[k]$ . The notation  $i \in \varphi(p)$ , for  $i \in [k]$  and  $p \in [K]$ , is just a shortening for the fact that the integer  $i$  is contained in the pair  $\varphi(p)$ .

As explained before,  $q^*$  is a CQ in  $\mathcal{C}$  that satisfies the following: It is a core, it is connected, and the  $(k \times K)$ -grid is a minor of  $G_{q^*}$ . In addition,  $\mu : [k] \times [K] \rightarrow 2^V$  is a minor map satisfying the aforementioned conditions M1, M2, and M3. The database  $D = D(G, q^*, \mu)$  is then defined over the alphabet of  $q^*$  as follows:

**The domain** The domain of  $D$  consists of all tuples:

$$(v, e, i, p, x) \in (V \times E \times [k] \times [K] \times \text{Dom}(q^*)),$$

such that the following holds:  $v \in e \iff i \in \varphi(p)$  and, in addition,  $x \in \mu(i, p)$ .

**The facts** Let us define the projection  $\Pi : D \rightarrow \text{Dom}(q^*)$  such that  $\Pi(v, e, i, p, x) = x$ . One can assume that  $\Pi$  extends to tuples by defining it component-wise. Then for every fact of the form  $R(\bar{x}) \in D_{q^*}$ , it is the case that  $D$  contains all facts of the form  $R(\bar{b})$  such that  $\Pi(\bar{b}) = \bar{x}$  and the following conditions hold for any two elements  $b = (v, e, i, p, x)$  and  $b' = (v', e', i', p', x')$  in the domain of  $D$  that are mentioned in  $\bar{b}$ :

- (C1) If  $i = i'$  then  $v = v'$ .
- (C2) If  $p = p'$  then  $e = e'$ .

Before establishing the correctness of the construction, let us analyze the cost of computing  $D = D(G, q^*, \mu)$ . First of all,  $q^*$  and  $\mu$  can be computed from  $q$ , which in turn depends only on  $k$ . Once they are computed, it is possible to construct  $D$  in time:



$$O(|V| \cdot |E| \cdot k \cdot K \cdot \|q^*\|)^r,$$

where  $r$  is the maximum arity of a relation symbol mentioned in  $q^*$ . But such a maximum arity is fixed by assumption, implying that there is a computable function  $f : \mathbb{N} \rightarrow \mathbb{N}$  such that  $D = D(G, q^*, \mu)$  can be constructed in time  $f(k) \cdot \|G\|^{O(1)}$  from  $(G, k)$ .

## Correctness of the construction

It is finally shown that  $q^*(D) = \mathbf{true}$  if and only if  $G$  contains a  $k$ -clique. This is proved in the two lemmas that follow:

**Lemma 23.6.** *If  $G$  contains a  $k$ -clique, then  $q^*(D) = \mathbf{true}$ .*

*Proof.* Let  $\{v_1, \dots, v_k\}$  be a set of vertices that defines a  $k$ -clique in  $G$ . For  $p \in [K]$  with  $\varphi(p) = \{i, j\}$ , let us define  $e_p$  to be the edge  $\{v_i, v_j\}$ . Therefore, it is possible to define a mapping  $h : q^* \rightarrow \text{Dom}(D)$  in such a way that:

$$h(x) = (v_i, e_p, i, p, x),$$

where  $i \in [k]$  and  $p \in [K]$  are the elements that satisfy that  $x \in \mu(i, p)$ . In fact,  $h(x)$  belongs to  $D$  as by definition it is the case that  $v_i \in e_p \iff i \in \varphi(p)$ .

Next it is shown that  $h : q^* \rightarrow \text{Dom}(D)$  is a homomorphism, thus implying  $q^*(D) = \mathbf{true}$ . Consider a fact of the form  $R(\bar{x})$  in  $q^*$ , where  $\bar{x} = (x_1, \dots, x_r)$ . Let  $i_1, \dots, i_r$  and  $p_1, \dots, p_r$  be such that  $x_j \in \mu(i_j, p_j)$ , for each  $1 \leq j \leq r$ . Then:

$$h(\bar{x}) = ((v_{i_1}, e_{p_1}, i_1, p_1, x_1), \dots, (v_{i_r}, e_{p_r}, i_r, p_r, x_r)).$$

Conditions (C1) and (C2) described above are trivially satisfied, and thus  $R(h(\bar{x})) \in D$ .  $\square$

**Lemma 23.7.** *If  $q^*(D) = \mathbf{true}$ , then  $G$  contains a  $k$ -clique.*

*Proof.* Since  $q^*(D) = \mathbf{true}$ , there is a homomorphism  $h : q^* \rightarrow D$ . Notice that, by definition,  $\Pi$  is a homomorphism from  $D$  to  $q^*$ . Then  $f = \Pi \circ h$  is a homomorphism from  $q^*$  to  $q^*$ , and thus it is also an isomorphism since  $q^*$  is a core. One can assume, without loss of generality, that  $f$  is the identity. If not, simply consider  $h \circ f^{-1}$  as the homomorphism instead of  $h$ .

As  $f = \Pi \circ h$  is the identity, for each element  $x \in \text{Dom}(q^*)$  such that  $x \in \mu(i, p)$ , for  $i \in [k]$  and  $p \in [K]$ , it must be the case that

$$h(x) = (v_x, e_x, i, p, x),$$

for some  $v_x \in V$  and  $e_x \in E$  such that  $v_x \in e_x \iff i \in \varphi(p)$ . To prove the existence of a  $k$ -clique in  $G$  it is necessary to establish several properties of the  $h(x)$ 's, for  $x \in \text{Dom}(q^*)$ , as stated in the following claims:

**Claim 23.8.** *For each  $i \in [k]$ ,  $p \in [K]$ , and  $x, x' \in \mu(i, p)$ , it is the case that  $v_x = v_{x'}$  and  $e_x = e_{x'}$ .*

*Proof.* Since  $\mu(i, p)$  is connected in  $q^*$ , it suffices to prove the claim for  $x, x'$  such that there is an edge between  $x$  and  $x'$  in  $G_{q^*}$ . Let  $R(\bar{x})$  be a fact in  $q^*$  such that both  $x$  and  $x'$  are mentioned in  $\bar{x}$ . Then  $R(h(\bar{x})) \in D$ , and thus  $v_x = v_{x'}$  and  $e_x = e_{x'}$  since conditions (C1) and (C2) hold.  $\square$

**Claim 23.9.** *For each  $i, i' \in [k]$ ,  $p \in [K]$ ,  $x \in \mu(i, p)$ , and  $x' \in \mu(i', p)$ , it is the case that  $e_x = e_{x'}$ .*

*Proof.* Let us assume without loss of generality that  $i \leq i'$ . If  $i = i'$  then the result holds from Claim 23.8. Let us suppose then that  $i < i'$ . It is sufficient to establish the result for the case  $i' = i + 1$ , as all other cases follow easily by induction.

Since there is an edge between  $\{i, p\}$  and  $\{i+1, p\} = \{i', p\}$  in the  $(k \times K)$ -grid and  $\mu$  is a minor map, there is an edge between an element  $y \in \mu(i, p)$  and an element  $y' \in \mu(i', p)$  in  $G_{q^*}$ . Thus, there is a fact  $R(\bar{y})$  in  $q^*$  such that both  $y$  and  $y'$  are mentioned in  $\bar{y}$ . Let us assume without loss of generality that  $\bar{y} = (y, y', \dots)$ . Since  $R(h(\bar{y})) \in D$ , it must be the case that  $e_y = e_{y'}$  by condition (C2). But by Claim 23.8 it is the case that  $e_y = e_x$  and  $e_{y'} = e_{x'}$ . This finishes the proof of the claim.  $\square$

**Claim 23.10.** *For each  $i \in [k]$ ,  $p, p' \in [K]$ ,  $x \in \mu(i, p)$ , and  $x' \in \mu(i, p')$ , it is the case that  $v_x = v_{x'}$ .*

*Proof.* Analogous to the proof of Claim 23.9 and left as an exercise for the reader.  $\square$

Summing up, the previous claims imply that there are vertices  $v_1, \dots, v_k$  and edges  $e_1, \dots, e_K$  in  $G$  satisfying the following:

- For each  $i \in [k]$  and element  $x \in \text{Dom}(q^*)$  such that  $x \in \mu(i, p)$ , for some  $p \in [K]$ , it is the case that  $h(x)$  is of the form  $(v_i, e, i, p, x)$ .
- For each  $p \in [K]$  and element  $x \in \text{Dom}(q^*)$  such that  $x \in \mu(i, p)$ , for some  $i \in [k]$ , it is the case that  $h(x)$  is of the form  $(v, e_p, i, p, x)$ .

Fix an arbitrary pair  $\{i, j\}$  with  $1 \leq i < j \leq k$  and consider the pair  $p \in [K]$  such that  $\varphi(p) = \{i, j\}$ . It is possible to prove that  $e_p = \{v_i, v_j\}$ , and, therefore, that  $\{v_1, \dots, v_k\}$  defines a  $k$ -clique in  $G$ . In fact, since  $\mu$  is a minor map, there are elements  $x \in \mu(i, p)$  and  $x' \in \mu(j, p)$ . From the previous remarks,  $h(x) = (v_i, e_p, i, p, x)$  and  $h(x') = (v_j, e_p, j, p, x)$ . Now, since  $v_i \in e_p \iff i \in \varphi(p)$  and  $v_j \in e_p \iff j \in \varphi(p)$  by definition, we conclude that  $e_p = \{v_i, v_j\}$ . This finishes the proof of Lemma 23.7.  $\square$

## Bounding the Join Size

So far, the search for efficient CQ evaluation methods has focused on the so-called *structural approach*. The underlying idea is to exploit structural properties of the input CQs, such as bounded generalized hypertreewidth, to develop tractable evaluation algorithms. This approach, however, disregards quantitative aspects such as the cardinalities of different relations in the database, which play a fundamental role in query evaluation.

The goal of this chapter is to present a result that brings together both structural and quantitative aspects. It does so by providing a tight bound on the size of the answer  $q(D)$  of a CQ  $q$  over a database  $D$  in terms of the cardinalities of the relations in  $D$  and a sophisticated notion of width for  $q$  based on *fractional covers*. As we will see in the next chapter, this result also provides the tools for developing worst-case optimal evaluation algorithms for CQs.

### Join Queries

We consider in this chapter the named perspective for relational algebra, as defined in Chapter 4. Recall that, under this perspective, the join  $R_1 \bowtie R_2$  of relations  $R_1[U_1]$  and  $R_2[U_2]$  is the set of tuples  $t$  of sort  $U_1 \cup U_2$  such that there exist  $t_1 \in R_1$  and  $t_2 \in R_2$  with  $t(A) = t_1(A)$  for every  $A \in U_1$  and  $t(A) = t_2(A)$  for every  $A \in U_2$ . That is, if  $A$  is a common attribute name in  $R_1$  and  $R_2$ , then  $t_1$  and  $t_2$  have the same value for their  $A$ -attribute.

For the sake of presentation, this chapter only deals with CQs that represent *join queries*, which we define next.

#### Definition 25.1: Join Query

A *join query* over a named schema  $\mathbf{S}$  is a named RA query over  $\mathbf{S}$  of the form

$$R_1 \bowtie \cdots \bowtie R_n .$$

We will assume throughout the chapter that  $R_1, \dots, R_n$  are pairwise different. Notice that this is possible because the join operator  $\bowtie$  is commutative, associative, and idempotent (that is,  $(R \bowtie R)(D) = R(D)$  for every database  $D$ ). Join queries therefore correspond to CQs of the form  $\text{Answer}(\bar{x}) :- R_1(\bar{y}_1), \dots, R_n(\bar{y}_n)$  where

- (1) each relation name  $R_i$  with  $i \in [n]$  occurs exactly once;
- (2) for every  $i \in [n]$ , no variables are repeated in the tuple  $\bar{y}_i$ ; and
- (3) every variable that appears in some tuple  $\bar{y}_i$  for  $i \in [n]$  also appears in  $\bar{x}$ .

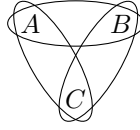
Condition (3) states that join queries do not have bound variables, i.e., they are projection free. For several of the results presented next this restriction is not essential, and in fact such results continue to hold for arbitrary CQs at the cost of more complicated proofs.

## A Gentle Introduction

As a gentle introduction, consider the join query

$$q_\Delta = R[A, B] \bowtie S[B, C] \bowtie T[C, A],$$

and the hypergraph



that visualizes the structure of  $q_\Delta$ . For simplicity, let us assume that we evaluate  $q_\Delta$  on a database  $D$  where  $R$ ,  $S$ , and  $T$  have the same number  $n$  of tuples. How many tuples can there be in  $q_\Delta(D)$ ?

Trivially,  $|q_\Delta(D)|$  is at most  $n^3$ , because  $|q_\Delta(D)|$  is at most as large as the size  $|R \times S \times T|$  of the Cartesian product of  $R$ ,  $S$ , and  $T$ . But we can also see that  $|q_\Delta(D)|$  is at most  $n^2 = |R \times S|$ . This is because we can obtain  $q_\Delta(D)$  as follows. We first compute  $R \times S$  and select those tuples in  $R \times S$  that agree on the  $B$ -attribute. Finally, we observe that joining the result of the previous step with  $T$  can only remove further tuples: it selects those tuples that agree with some tuple in  $T$  on their  $A$ - and  $C$ -attribute.

This idea can be generalized when we consider *edge covers* of the hypergraph of a join query  $q$  of the form  $R_1 \bowtie \dots \bowtie R_n$ . Let  $\{A_1, \dots, A_m\}$  be the set of attributes used by  $q$ , i.e.,  $\cup_{i \in [n]} \mathbf{S}(R_i)$ . We define the hypergraph  $H_q = (V_q, E_q)$  of  $q$  to consist of the set of nodes  $V_q = \{A_1, \dots, A_m\}$  and hyperedges  $E_q = \{\mathbf{S}(R_i) \mid i \in [n]\}$ .

**Definition 25.2: Edge Cover**

An *edge cover* of a hypergraph  $H = (V, E)$  is a subset  $E' \subseteq E$  of its edges such that, for each node  $v \in V$ , there is an edge  $e \in E'$  with  $v \in e$ .

The intuition of edge covers is that the set of selected edges in  $E'$  “cover” each node of  $H$ . For example, the hyperedges  $\{A, B\}$  and  $\{B, C\}$  (corresponding to  $R[A, B]$  and  $S[B, C]$  in  $q_\Delta$ ) are an edge cover of  $H_{q_\Delta}$ . The principle that we used to give the  $n^2$  upper bound for  $|q_\Delta(D)|$  can be generalized as follows.

**Theorem 25.3**

Let  $q = R_1 \bowtie \cdots \bowtie R_n$  be a join query and let  $E$  be an edge cover of its hypergraph  $H_q$ . Then, for every database  $D$ , we have that

$$|q(D)| \leq \prod_{S(R_i) \in E} |R_i^D| \quad \in O(|D|^{|E|}).$$

We will prove a generalization of this result as Theorem 25.8.

Whereas Theorem 25.3 gives an upper bound for  $|q(D)|$ , what can we say about lower bounds? Consider again the query  $q_\Delta$ . We can see that, for every  $k \in \mathbb{N}$ , the database  $D_k$ , in which  $R[A, B] = \{(A : i, B : 1) \mid i \in [k]\}$ ,  $S[B, C] = \{(B : 1, C : 1)\}$ , and  $T[C, A] = \{(C : 1, A : i) \mid i \in [k]\}$  returns  $k$  answers to  $q$ . Indeed, we have that  $q(D_k) = \{(A : i, B : 1, C : 1) \mid i \in [k]\}$ .

Interestingly, also this idea can be generalized using a standard graph-theoretical notion, namely *independent sets*.

**Definition 25.4: Independent Set**

An *independent set* of a hypergraph  $H = (V, E)$  is a subset  $V' \subseteq V$  of its nodes such that, for each edge  $e \in E$ , it holds that  $|e \cap V'| \leq 1$ .

The intuition behind independent sets is that all distinct node pairs in  $V'$  are “independent”, i.e., not connected by an edge. For example, the node  $A$  is an independent set of  $H_{q_\Delta}$  which, incidentally, we used to construct the databases  $D_k$ . Similar to before, we can also generalize this idea.

**Proposition 25.5**

Let  $q = R_1 \bowtie \cdots \bowtie R_n$  be a join query and let  $V$  be an independent set of its hypergraph  $H_q$ . Then, for every  $k \in \mathbb{N}$ , there is a database  $D$  with

$$k^{|V|} \leq |q(D)|$$

and every relation in  $D$  has at most  $k$  tuples.

Again, we will prove a generalization of this result as Proposition 25.20.

## The AGM Bound

We now embark on generalizing the results from the gentle introduction and start with focusing on upper bounds. We consider the following generalization of the edge covers from Definition 25.2

### Definition 25.6: Fractional Edge Cover

A *fractional edge cover* of a hypergraph  $H = (V, E)$  is a function

$$f : E \rightarrow \mathbb{Q}_{\geq 0}$$

such that, for each node  $v \in V$ , it holds that  $\sum_{v \in e} f(e) \geq 1$ . The *weight* of  $f$  is the value  $\sum_{e \in E} f(e)$ .

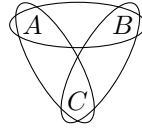
It is easy to see that each edge cover is indeed also a fractional edge cover. Let us illustrate the difference between edge covers and their fractional variant on an example.

### Example 25.7

Consider again the join query

$$q_{\Delta} = R[A, B] \bowtie S[B, C] \bowtie T[C, A],$$

and its hypergraph  $H_{q_{\Delta}}$



By slight abuse of notation, let us denote the hyperedge  $\{A, B\}$  as  $R$ , the hyperedge  $\{B, C\}$  as  $S$ , and  $\{C, A\}$  as  $T$ . We can then write the conditions for a fractional edge cover  $f$  of  $H_{q_{\Delta}}$  as the following system of equations.

$$\begin{aligned} f(R), f(S), f(T) &\geq 0, \\ f(R) + f(T) &\geq 1, \quad f(R) + f(S) \geq 1, \quad f(S) + f(T) \geq 1. \end{aligned}$$

The equation  $f(R) + f(T) \geq 1$ , for example, is obtained by considering attribute  $A$  and the fact that  $R$  and  $T$  are the edges in  $H_{q_\Delta}$  that contain  $A$ . One solution to this system of equations, i.e., a fractional edge cover, is  $f(R) = f(S) = f(T) = 1/2$ . The weight of this fractional edge cover is  $3/2$ .

Recall that in a database  $D$ , we denote by  $R^D$  the relation instance associated to relation name  $R$ , that is,  $R^D$  is the set of tuples  $\bar{a}$  such that  $R(\bar{a}) \in D$ . With this notation, we can state the AGM bound, which received its name from the last names of Albert Atserias, Martin Grohe, and Dániel Marx.

**Theorem 25.8: AGM Bound**

Consider a join query  $q = R_1 \bowtie \cdots \bowtie R_n$  over schema  $\mathbf{S}$  and a fractional edge cover  $f$  of  $q$ . Then, for every database  $D$ , we have that

$$|q(D)| \leq \prod_{i=1}^n |R_i^D|^{f(\mathbf{S}(R_i))}.$$

Before proving the theorem, it is worth illustrating the application of the AGM bound on a specific query.

**Example 25.9**

Applying the AGM bound to the fractional cover of Example 25.7 we obtain that, over every database  $D$ ,

$$|q(D)| \leq \sqrt{|R^D| \cdot |S^D| \cdot |T^D|}.$$

We could also consider another fractional cover, namely  $f(R) = f(S) = 1$  and  $f(T) = 0$ . Applying the AGM bound to this solution implies that

$$|q(D)| \leq |R^D| \cdot |S^D|.$$

Which bound is better depends on the underlying database  $D$ .

- Consider first the case when  $R^D = S^D = T^D$ , i.e., when the interpretations of all three relations over  $D$  coincide. Moreover, assume that  $R^D$  does not contain any tuple of the form  $(a, a)$ . Then  $D$  can be naturally seen as a graph without loops and with  $M = |R^D|$  edges, and, thus,  $|q(D)|$  is equal to three times the number of directed triangles in such a graph. (Notice that each directed triangle is counted three times as a tuple  $(A: a, B: b, C: c)$  is considered to be different from a tuple  $(A: b, B: c, C: a)$ ). In this case, the first bound is bet-

ter, as it implies that  $|q(D)| \leq M^{3/2}$ , while the second one implies that  $|q(D)| \leq M^2$ .

The fact that the maximum number of directed triangles in a graph with  $M$  edges is bounded by

$$\frac{M^{3/2}}{3}$$

is non-trivial, which illustrates the power of Theorem 25.8 as a tool for obtaining meaningful bounds on the size of the evaluation of a join query.

- Consider now the case when  $|R^D| = |S^D| = 1$  and  $|T^D| = M$ . Then the second bound is tighter, as it implies that  $|q(D)| = 1$  while the first one implies that  $|q(D)| \leq \sqrt{M}$ .

## Proof of the AGM Bound

We will now prove the AGM bound. The proof makes use of a key *query decomposition lemma*, which is in turn proved by applying the following version of Hölder's inequality, which we state without proof.

### Theorem 25.10: Hölder's Inequality

Let  $p, r$  be positive integers,  $y_1, \dots, y_r$  be non-negative real numbers such that  $y_1 + \dots + y_r \geq 1$ , and  $a_{i,j}$  be a non-negative real number for every  $i \in [p]$  and  $j \in [r]$ . Then it holds that

$$\sum_{i=1}^p \prod_{j=1}^r a_{i,j}^{y_j} \leq \prod_{j=1}^r \left( \sum_{i=1}^p a_{i,j} \right)^{y_j}.$$

From now on, fix a join query  $q = R_1 \bowtie \dots \bowtie R_n$  over schema  $\mathbf{S}$ , and assume that  $X = \{A_1, \dots, A_m\}$  is the set of attributes occurring in  $q$ . Given  $Y \subseteq X$ , define  $\mathcal{R}(Y)$  as  $\bigcup_{A_j \in Y} \mathcal{R}(A_j)$ . That is,  $\mathcal{R}(Y)$  is the set of those  $R_i$ 's that mention some attribute in  $Y$ , for  $i \in [n]$ . Moreover, assuming that  $\mathcal{R}(Y) = \{R_{i_1}, \dots, R_{i_\ell}\}$ , where  $1 \leq i_1 < \dots < i_\ell \leq n$ , define query

$$q_Y = \pi_{Y \cap \mathbf{S}(R_{i_1})}(R_{i_1}) \bowtie \dots \bowtie \pi_{Y \cap \mathbf{S}(R_{i_\ell})}(R_{i_\ell}). \quad (25.1)$$

Recall that  $\mathbf{S}(R)$  is the set of attributes associated to relation name  $R$  under schema  $\mathbf{S}$ . Thus, query  $q_Y$  defines the join of the projection over  $Y$  of those  $R_i$ 's that mention at least some attribute in  $Y$ .<sup>1</sup>

<sup>1</sup> Recall that, in order for a projection  $\pi_\alpha(e)$  to be well defined, it is required that  $\alpha$  is a subset of the sort  $U$  of  $e$ . This is why we intersect with  $\mathbf{S}(R_{i_j})$  in every projection.



Finally, to state the query decomposition lemma, recall the definition of the semijoin operator  $\bowtie$  from Chapter 18. In particular, recall that two tuples  $\bar{a}$  and  $\bar{b}$  are said to be consistent if they have the same value in each shared attribute. Thus, the term  $R \bowtie \{\bar{a}\}$  is used in the lemma to denote the set of tuples in  $R$  that are consistent with  $\bar{a}$ .

**Lemma 25.11 (Query Decomposition Lemma).** *Assume that  $f$  is a fractional edge cover for the hypergraph  $H_q$  of  $q$  and consider an arbitrary partition  $\{Y, Z\}$  of  $X$ . Then for every database  $D$  it holds that*

$$\sum_{\bar{a} \in q_Y(D)} \prod_{R_i \in \mathcal{R}(Z)} |R_i^D \bowtie \{\bar{a}\}|^{f(\mathbf{S}(R_i))} \leq \prod_{i=1}^n |R_i^D|^{f(\mathbf{S}(R_i))}.$$

Here, we assume that all relations  $R_i^D$  are non-empty.

*Proof.* Let  $A$  be an arbitrary attribute in  $Y$ . Define  $Y' = Y - \{A\}$  and  $Z' = Z \cup \{A\}$ . Next we show that

$$\begin{aligned} \sum_{\bar{a} \in q_Y(D)} \prod_{R_i \in \mathcal{R}(Z)} |R_i^D \bowtie \{\bar{a}\}|^{f(\mathbf{S}(R_i))} \\ \leq \sum_{\bar{a}' \in q_{Y'}(D)} \prod_{R_i \in \mathcal{R}(Z')} |R_i^D \bowtie \{\bar{a}'\}|^{f(\mathbf{S}(R_i))}. \end{aligned} \quad (25.2)$$

The lemma is then obtained by repeatedly applying (25.2) until  $Y'$  is empty, in which case the right-hand side of the inequality is precisely  $\prod_{i=1}^n |R_i^D|^{f(\mathbf{S}(R_i))}$ .

Each tuple  $\bar{a} \in q_Y(D)$  can be decomposed as a pair  $\bar{a}', v$  such that  $\bar{a}' \in q_{Y'}(D)$  and  $v$  is the value of  $\bar{a}$  for attribute  $A$ . In what follows, we use  $(\bar{a}', v)$  as an alternative notation for tuple  $\bar{a}$ . Then the left-hand side of equation (25.2) can be expressed as

$$\sum_{\bar{a}' \in q_{Y'}(D)} \sum_{v : (\bar{a}', v) \in q_Y(D)} \prod_{R_i \in \mathcal{R}(Z)} |R_i^D \bowtie \{(\bar{a}', v)\}|^{f(\mathbf{S}(R_i))},$$

which in turn can be rewritten as

$$\begin{aligned} \sum_{\bar{a}' \in q_{Y'}(D)} \sum_{v : (\bar{a}', v) \in q_Y(D)} \left( \left( \prod_{R_i \in \mathcal{R}(Z)} |R_i^D \bowtie \{(\bar{a}', v)\}|^{f(\mathbf{S}(R_i))} \right) \right. \\ \left. \left( \prod_{R_i \in \mathcal{R}(Z') - \mathcal{R}(Z)} 1^{f(\mathbf{S}(R_i))} \right) \right). \end{aligned} \quad (25.3)$$

Notice that equation (25.3) is equivalent to

$$\sum_{\bar{a}' \in q_{Y'}(D)} \sum_{v : (\bar{a}', v) \in q_Y(D)} \prod_{R_i \in \mathcal{R}(Z')} |R_i^D \bowtie \{(\bar{a}', v)\}|^{f(\mathbf{S}(R_i))}, \quad (25.4)$$

since for those  $R_i$  in  $\mathcal{R}(Z') - \mathcal{R}(Z)$ , it holds that the set of attributes of  $R_i$  is contained in  $Y$  and, thus,  $|R_i^D \bowtie \{(\bar{a}', v)\}| = 1$ . Moreover, given that  $R_i^D \bowtie \{(\bar{a}', v)\} = R_i^D \bowtie \{\bar{a}'\}$  for each  $R_i \in \mathcal{R}(Z') - \mathcal{R}(A)$ , equation (25.4) can be expressed as

$$\sum_{\bar{a}' \in q_{Y'}(D)} \prod_{R_i \in \mathcal{R}(Z') - \mathcal{R}(A)} |R_i^D \bowtie \{\bar{a}'\}|^{f(\mathbf{S}(R_i))} \sum_{v : (\bar{a}', v) \in q_Y(D)} \prod_{R_i \in \mathcal{R}(A)} |R_i^D \bowtie \{(\bar{a}', v)\}|^{f(\mathbf{S}(R_i))}. \quad (25.5)$$

This is the moment in which Hölder's inequality stated in Claim 25.10 is applied to obtain that the value expressed in equation (25.5) is bounded by

$$\sum_{\bar{a}' \in q_{Y'}(D)} \prod_{R_i \in \mathcal{R}(Z') - \mathcal{R}(A)} |R_i^D \bowtie \{\bar{a}'\}|^{f(\mathbf{S}(R_i))} \prod_{R_i \in \mathcal{R}(A)} \left( \sum_{v : (\bar{a}', v) \in q_Y(D)} |R_i^D \bowtie \{(\bar{a}', v)\}| \right)^{f(\mathbf{S}(R_i))}. \quad (25.6)$$

Observe that it is possible to apply Claim 25.10 because: (i)  $f$  is a fractional cover of  $q$  and, thus,  $\sum_{R_i \in \mathcal{R}(A)} f(\mathbf{S}(R_i)) \geq 1$ ; and (ii)  $|R_i^D \bowtie \{(\bar{a}', v)\}| \geq 0$  for every  $R_i \in \mathcal{R}(A)$  and  $v$  such that  $(\bar{a}', v) \in q_Y(D)$ .

To conclude, the value expressed in Equation (25.6) is bounded by

$$\sum_{\bar{a}' \in q_{Y'}(D)} \prod_{R_i \in \mathcal{R}(Z') - \mathcal{R}(A)} |R_i^D \bowtie \{\bar{a}'\}|^{f(\mathbf{S}(R_i))} \prod_{R_i \in \mathcal{R}(A)} |R_i^D \bowtie \{\bar{a}'\}|^{f(\mathbf{S}(R_i))}. \quad (25.7)$$

This is because for every  $R_i \in \mathcal{R}(A)$ : (i)  $R_i^D \bowtie \{(\bar{a}', v)\} \subseteq R_i^D \bowtie \{\bar{a}'\}$ ; and (ii) if  $(\bar{a}', v_1)$  and  $(\bar{a}', v_2)$  are in  $q_Y(D)$ , with  $v_1 \neq v_2$ , then  $R_i^D \bowtie \{(\bar{a}', v_1)\}$  and  $R_i^D \bowtie \{(\bar{a}', v_2)\}$  are disjoint given that  $A \in \mathbf{S}(R_i)$ . Finally, equation (25.7) can be simplified as

$$\sum_{\bar{a}' \in q_{Y'}(D)} \prod_{R_i \in \mathcal{R}(Z')} |R_i^D \bowtie \{\bar{a}'\}|^{f(\mathbf{S}(R_i))},$$

which finishes the proof of the lemma.  $\square$

We now move to the proof of Theorem 25.8, which is done by induction on the size of the set  $X = \{A_1, \dots, A_m\}$  of attributes occurring in  $q$ . Notice that this theorem trivially holds if  $|R_i^D| = 0$  for some  $i \in [n]$ . Thus, we assume in the following proof that  $|R_i^D| \geq 1$  for every  $i \in [n]$ .

- **Base case.** We have that  $|X| = 1$ , and hence each relation  $R_i$  is unary for each  $i \in [n]$ . Then we have that

$$\begin{aligned}
|q(D)| &\leq \min_{\ell \in [n]} |R_\ell^D| \\
&\leq \left( \min_{\ell \in [n]} |R_\ell^D| \right)^{\sum_{i=1}^n f(\mathbf{S}(R_i))} \\
&= \prod_{i=1}^n \left( \min_{\ell \in [n]} |R_\ell^D| \right)^{f(\mathbf{S}(R_i))} \\
&\leq \prod_{i=1}^n |R_i^D|^{f(\mathbf{S}(R_i))},
\end{aligned}$$

where the second inequality holds given that  $\min_{\ell \in [n]} |R_\ell^D| \geq 1$  and that  $\sum_{i=1}^n f(\mathbf{S}(R_i)) \geq 1$ .

- **Inductive case.** We have that  $|X| = m$  with  $m > 1$ . Let  $Y = \{A_1, \dots, A_{m-1}\}$ . Let  $\mathbf{S}'$  be the schema such that  $\mathbf{S}'(R) = \{A_m\}$  for every  $R \in \mathcal{R}(A_m)$ . Let  $q'$  be the natural join query over  $\mathbf{S}'$ , that is,

$$q' = \bowtie_{R \in \mathbf{S}'} R.$$

Notice that  $q'$  is in fact the query  $q_{\{A_m\}} = q_{X-Y}$  as in (25.1), which computes the intersection over all the unary relations  $R$  over schema  $\mathbf{S}'$ . For every tuple  $\bar{a} \in q_Y(D)$ , let  $D_{\bar{a}}$  be the database over  $\mathbf{S}'$  such that  $R^{D_{\bar{a}}} = \pi_{\{A_m\}}(R^D \times \{\bar{a}\})$  for each  $R \in \mathcal{R}(A_m)$ . Then we have that

$$q(D) = \bigcup_{\bar{a} \in q_Y(D)} (q'(D_{\bar{a}}) \times \{\bar{a}\}) \quad (25.8)$$

and, therefore,

$$|q(D)| \leq \sum_{\bar{a} \in q_Y(D)} |q'(D_{\bar{a}})|.$$

Since  $f$  satisfies  $\sum_{A_m \in \mathbf{S}(R_i)} f(\mathbf{S}(R_i)) \geq 1$ , we obtain by induction hypothesis for every  $\bar{a} \in q_Y(D)$  that

$$|q'(D_{\bar{a}})| \leq \prod_{R_i \in \mathcal{R}(A_m)} |R_i^{D_{\bar{a}}}|^{f(\mathbf{S}(R_i))} = \prod_{R_i \in \mathcal{R}(A_m)} |R_i^D \times \{\bar{a}\}|^{f(\mathbf{S}(R_i))}.$$

It is possible then to conclude that

$$\begin{aligned}
|q(D)| &\leq \sum_{\bar{a} \in q_Y(D)} |q'(D_{\bar{a}})| \leq \\
&\sum_{\bar{a} \in q_Y(D)} \prod_{R_i \in \mathcal{R}(A_m)} |R_i^D \times \{\bar{a}\}|^{f(\mathbf{S}(R_i))} \leq \prod_{i=1}^n |R_i^D|^{f(\mathbf{S}(R_i))},
\end{aligned}$$

where the last inequality holds by Lemma 25.11.

This finishes the proof of Theorem 25.8.

## Fractional Independent Sets

In the remainder of this chapter, we will prove that the AGM bound is tight. To this end, we will consider a fractional version of *independent sets*.

### Definition 25.12: Fractional Independent Set

A *fractional independent set* of a hypergraph  $H = (V, E)$  is a function

$$g : V \rightarrow \mathbb{Q}_{\geq 0}$$

such that, for each edge  $e \in E$ , it holds that  $\sum_{v \in e} g(v) \leq 1$ . The *weight* of  $g$  is the sum  $\sum_{v \in V} g(v)$ .

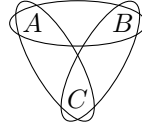
Fractional independent sets assign a nonnegative weight to each node in the hypergraph, such that the sum of the weights of the nodes in a single hyperedge should not exceed 1. This generalizes the condition of the “classical” independent set problem that states that at most one node per edge is allowed in the independent set. We illustrate the notion on an example.

### Example 25.13

Consider again the join query

$$q_{\Delta} = R[A, B] \bowtie S[B, C] \bowtie T[C, A] ,$$

and its hypergraph  $H_{q_{\Delta}}$



We can write the conditions for a fractional independent set  $g$  of  $H_{q_{\Delta}}$  as the following system of equations.

$$\begin{aligned} g(A), g(B), g(C) &\geq 0, \\ g(A) + g(B) &\leq 1, \quad g(B) + g(C) \leq 1, \quad g(C) + g(A) \leq 1. \end{aligned}$$

The equation  $g(A) + g(B) \leq 1$  is obtained by considering the hyperedge  $\{A, B\}$ . One solution to this system of equations, i.e., a fractional independent set, is  $g(A) = g(B) = g(C) = 1/2$ . The weight of this fractional independent set is  $3/2$ .

## Connection to Linear Programming

Fractional edge covers and fractional independent sets are closely related to each other, which will become clear when we explore their connection to linear programming. To this end, we first need to introduce some background on linear programs.

### Definition 25.14: Linear Program

Let  $a_1, \dots, a_n$  be real numbers and  $x_1, \dots, x_n$  be variables. A *linear function* (over  $x_1, \dots, x_n$ ) is a function  $f$  of the form

$$f(x_1, \dots, x_n) = a_1x_1 + a_2x_2 + \dots + a_nx_n .$$

If  $f$  is a linear function and  $b$  a real number, then the equations

$$f(x_1, \dots, x_n) \leq b \quad \text{and} \quad f(x_1, \dots, x_n) \geq b$$

are *linear inequalities*. A *linear program* (over  $x_1, \dots, x_n$ ) consists of a linear function  $f$ , a set of linear inequalities, and a *task*, which is to *minimize* or to *maximize*.

Let us illustrate how we will denote linear programs by means of a few examples.

### Example 25.15

When we denote  $f(R)$ ,  $f(S)$ , and  $f(T)$  as variables  $x_R$ ,  $x_S$ , and  $x_T$ , respectively, then we can associate a linear program to the system of equations in Example 25.7, namely

$$\begin{aligned} &\text{minimize} && x_R + x_S + x_T \\ &\text{subject to} && x_R + x_T \geq 1 \\ & && x_R + x_S \geq 1 \\ & && x_S + x_T \geq 1 \\ &\text{and} && x_R \geq 0 \quad x_S \geq 0 \quad x_T \geq 0 . \end{aligned}$$

Here, the linear function is  $x_R + x_S + x_T$ , the task is to minimize, and we have a set of six linear equations.

Likewise, when we denote  $g(A)$ ,  $g(B)$ , and  $g(C)$  as variables  $y_A$ ,  $y_B$ , and  $y_C$ , respectively, then we can associate a linear program to the system of equations in Example 25.13, namely

$$\begin{array}{ll}
\text{maximize} & y_A + y_B + y_C \\
\text{subject to} & y_A + y_B \leq 1 \\
& y_B + y_C \leq 1 \\
& y_C + y_A \leq 1 \\
\text{and} & y_A \geq 0 \quad y_B \geq 0 \quad y_C \geq 0 .
\end{array}$$

Here, the task is to maximize.

### Definition 25.16: Feasible Solutions and Optimal Values

Let  $L$  be a linear program with linear inequalities

$$\begin{array}{l}
f_1(x_1, \dots, x_n) \theta_1 b_1 \\
\vdots \\
f_k(x_1, \dots, x_n) \theta_k b_k ,
\end{array}$$

where  $\theta_i \in \{\leq, \geq\}$  for every  $i \in [k]$ . A *feasible solution* to  $L$  is an assignment  $\rho : \{x_1, \dots, x_n\} \rightarrow \mathbb{R}$  such that  $f_i(\rho(x_1), \dots, \rho(x_n)) \theta_i b_i$  holds for every  $i \in [k]$ .

Let  $L$  additionally have the linear function  $f$ . If the task of  $L$  is to minimize (resp., to maximize), then its *optimal value*, if it exists, is the smallest (resp., largest) value of  $f(\rho(x_1), \dots, \rho(x_n))$  such that  $\rho$  is a feasible solution to  $L$ . In this case, we call  $\rho$  an *optimal solution* to  $L$ .

### Example 25.17

Both  $\rho_1$  with  $\rho_1(x_R) = \rho_1(x_S) = \rho_1(x_T) = 1/2$  and  $\rho_2$  with  $\rho_2(x_R) = \rho_2(x_S) = 1$  and  $\rho_2(x_T) = 0$  are feasible solutions to the first linear program in Example 25.15. One can show that the optimal value is  $3/2$  and therefore  $\rho_1$  is an optimal solution.

## Linear Programming Duality

We will use the Strong Duality Theorem of linear programming, which we state next. To this end, the *dual* of a linear program of the form

$$\begin{array}{ll}
\text{minimize} & \sum_{i=1}^n c_i x_i \\
\text{subject to} & \sum_{i=1}^n a_{i,j} x_i \geq b_j \quad \text{for each } j \in [m], \\
& \text{and} \quad x_i \geq 0 \quad \text{for each } i \in [n]
\end{array}$$

over the variables  $x_1, \dots, x_n$  is the linear program

$$\begin{array}{ll}
\text{maximize} & \sum_{j=1}^m b_j y_j \\
\text{subject to} & \sum_{j=1}^m a_{i,j} y_j \leq c_i \quad \text{for each } i \in [n], \\
& \text{and} \quad y_j \geq 0 \quad \text{for each } j \in [m]
\end{array}$$

over the variables  $y_1, \dots, y_m$ .

**Theorem 25.18: Strong Duality Theorem**

Let  $\rho^*$  be the optimal value of a linear program  $P$  and let  $\tau^*$  be the optimal value of its dual program  $P_D$ . Then  $\rho^* = \tau^*$ .

We are now ready to show the key observation that ties the previous notions in this chapter together with linear programs. The observation is that, if we consider a join query  $q$ , then the linear program that maximizes the weight of the fractional independent set of  $H_q$  is the dual of the linear program that minimizes the weight of the fractional edge cover of  $H_q$ . For instance, in Example 25.15 the second linear program is the dual of the first one. We make this more precise next.

Let  $q = R_1 \bowtie \dots \bowtie R_n$  be a join query over some schema  $\mathbf{S}$  and assume that  $\{A_1, \dots, A_m\}$  is the set of all attributes occurring in  $q$ . Recall that we use  $\mathcal{R}(A_j)$  to denote the set of relation names  $R_i$  in  $\{R_1, \dots, R_n\}$  that use attribute  $A_j$  under schema  $\mathbf{S}$ , i.e., such that  $A_j \in \mathbf{S}(R_i)$ .

Let us denote the linear program

$$\begin{array}{ll}
\text{minimize} & \sum_{i=1}^n x_i \\
\text{subject to} & \sum_{R_i \in \mathcal{R}(A_j)} x_i \geq 1 \quad \text{for each } j \in [m], \\
& \text{and} \quad x_i \geq 0 \quad \text{for each } i \in [n],
\end{array}$$

that minimizes the weight of the fractional edge cover of  $H_q$ , as  $\text{EC}(H_q)$ . Likewise, let us denote the linear program

$$\begin{array}{ll}
\text{maximize} & \sum_{j=1}^m y_j \\
\text{subject to} & \sum_{A_j \in \mathbf{S}(R_i)} y_j \leq 1 \quad \text{for each } i \in [n] \\
\text{and} & y_j \geq 0 \quad \text{for each } j \in [m] .
\end{array}$$

that maximizes the weight of the fractional independent set of  $H_q$ , as  $\text{IS}(H_q)$ . Notice that  $\text{IS}(H_q)$  is indeed the dual program of  $\text{EC}(H_q)$ .

Let us call a fractional edge cover  $f$  of  $H_q$  *minimal* if its weight equals the optimal value of  $\text{EC}(H_q)$ ; and a fractional independent set  $g$  *maximal* if its weight equals the optimal value of  $\text{IS}(H_q)$ . Then Theorem 25.18 implies the following Corollary.

**Corollary 25.19**

Let  $q$  be a join query and  $H_q$  be its hypergraph. Let  $f$  be a minimal fractional edge cover and  $g$  be a maximal fractional independent set of  $H_q$ . Then the weight of  $f$  equals the weight of  $g$ .

## Optimality of the AGM Bound

We show next that the bound stated in Theorem 25.8 is strict.

**Proposition 25.20**

Let  $q = R_1 \bowtie \cdots \bowtie R_n$  be a join query over schema  $\mathbf{S}$  and let  $f$  be a minimal edge cover of the hypergraph  $H_q$  of  $q$ . Then there exist arbitrarily large databases  $D$  such that  $|R_i^D| \geq 1$  for each  $i \in [n]$  and

$$|q(D)| = \prod_{i=1}^n |R_i^D|^{f(\mathbf{S}(R_i))} .$$

*Proof.* Assume that  $\{A_1, \dots, A_m\}$  is the set of attributes occurring in  $q$ . Let  $f$  be a minimal fractional edge cover and  $g$  be a maximal fractional independent set of  $H_q$ . Let  $f^*$  be the weight of  $f$  and  $g^*$  be the weight of  $g$ .

We now define a database  $D_k$  for each  $k \in \mathbb{N}$ . For each  $i \in [n]$ , let  $\sigma_i$  be a function such that  $\mathbf{S}(R_i) = \{A_{\sigma_i(1)}, \dots, A_{\sigma_i(r)}\}$ , where  $r = \text{arg}(\mathbf{S}(R_i))$ . We define the database  $D_k$  as

$$\begin{aligned}
R_i^{D_k} = \{ & (A_{\sigma_i(1)} : j_1, \dots, A_{\sigma_i(r)} : j_r) \mid \\
& j_\ell \in \{1, \dots, \lfloor k^{g(A_{\sigma_i(\ell)})} \rfloor\} \text{ for every } \ell \in [r] \} .
\end{aligned}$$



It is straightforward to show that

$$q(D_k) = \{1, \dots, \lfloor k^{g(A_1)} \rfloor\} \times \dots \times \{1, \dots, \lfloor k^{g(A_m)} \rfloor\}.$$

Since  $g(A_j) \in \mathbb{Q}$  for every  $j \in [m]$ , there are arbitrarily large  $k$  such that every value  $k^{g(A_j)}$  is an integer. In this case, all relations  $R_i^{D_k}$  have size  $k$  and we have that  $|q(D_k)| = k^{\sum_{j=1}^m g(A_j)}$ . We then obtain that

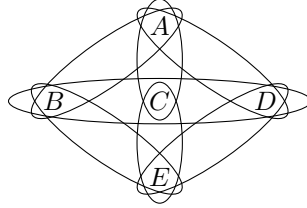
$$\prod_{i=1}^n |R_i^{D_k}|^{f(\mathbf{S}(R_i))} = k^{\sum_{i=1}^n f(\mathbf{S}(R_i))} = k^{f^*} = k^{g^*}.$$

Here, the first equality holds because  $|R_i^{D_k}| = k$ . The second equality is by definition of  $f^*$ , and the third is because  $f^* = g^*$  (Corollary 25.19). This concludes the proof.  $\square$

We illustrate the construction of Proposition 25.20 with an example.

#### Example 25.21

Consider a join query with the following hypergraph  $H$ :



This hypergraph has a fractional edge cover  $f$  that assigns weight  $\frac{1}{3}$  to every edge, amounting to a total weight of  $\frac{7}{3}$ . It has a fractional independent set  $g$  that assigns

- weight  $\frac{1}{3}$  to  $B$ ,  $C$ , and  $D$ ; and
- weight  $\frac{2}{3}$  to  $A$  and  $E$ ,

also amounting to a total weight of  $\frac{7}{3}$ . Therefore,  $f$  is minimal and  $g$  is maximal. We now use the construction in the proof of Proposition 25.20 to construct a database from  $g$ . Let us choose a value of  $k$  such that  $k^{g(v)}$  is an integer for every node  $v$  of  $H$ . This is the case for  $k = 8$ , since  $\sqrt[3]{8} = 2$  and  $8^{2/3} = 4$ . The construction in Proposition 25.20 now defines the database  $D_k$  as

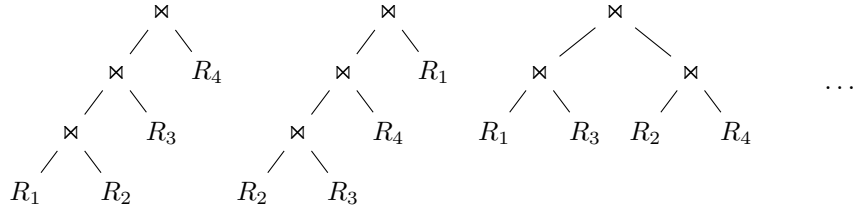
$A$	$B$			
$E$	$C$			
$D$		$B$	$C$	$D$
1	1	1	1	1
1	2	1	1	2
2	1	1	2	1
2	2	1	2	2
3	1	2	1	1
3	2	2	1	2
4	1	2	2	1
4	2	2	2	2

where every binary relation, i.e., over attributes  $\{A, B\}$ ,  $\{A, C\}$ ,  $\{A, D\}$ ,  $\{E, B\}$ ,  $\{E, C\}$ , and  $\{E, D\}$  is like the relation on the left and the ternary relation like the relation on the right.

## Worst-Case Optimal Join Algorithms

### Join Plans

In a database system, join queries are typically evaluated using *join plans*, which aim at finding a clever ordering in which to perform the joins. Assume, for example, that we have a join query  $q = R_1 \bowtie R_2 \bowtie R_3 \bowtie R_4$ . Due to the commutativity and associativity of the join operator, we can for example consider the trees



Each of these trees represent a join plan: if we have a database  $D$  with relations  $R_1, \dots, R_4$  and evaluate the joins of  $q$  in a bottom-up fashion over these trees, we end up with  $q(D)$  at the root of each tree.

Since joins are very expensive operators in database systems, it makes sense to think about which join plan leads to the fastest computation of the result. Here, it is important to know that, in general, joining big relations typically takes more time than joining small relations. Query optimizers therefore try to perform joins that produce small results early, since it is easier and faster to continue the computation with a small intermediate result than with a large one. For instance, if  $|R_1 \bowtie R_2(D)|$  is much smaller than  $|R_2 \bowtie R_3(D)|$ , we may prefer the leftmost join plan to the middle one. Since the join plans we considered here always take the join of two relations, they are called *pairwise join plans*.

### Are Pairwise Join Plans Sub-Optimal?

Perhaps surprisingly, there exist cases in which every possible pairwise join plan produces after the first step an intermediate result that contains more tuples than the entire output. Consider the join query

$$q(A, B, C) :- R(A, B), S(B, C), T(C, A)$$

and, for  $n \in \mathbb{N}$ , the database  $D_n$  in which  $R[A, B] = S[B, C] = T[C, A]$  is the relation

1	1
1	2
$\vdots$	$\vdots$
1	$n/2$
2	1
$\vdots$	$\vdots$
$n/2$	1

which consists of  $n - 1$  tuples.

Any join plan that starts with joining two relations will produce an intermediate result that is quadratic in  $n$ , because, if we consider  $q_{RS} = R \bowtie S$ ,  $q_{TR} = T \bowtie R$ , and  $q_{ST} = S \bowtie T$ , then we have

$$|q_{RS}(D_n)| = |q_{TR}(D_n)| = |q_{ST}(D_n)| = \left(\frac{n}{2}\right)^2 = \frac{n^2}{4}.$$

But, according to the AGM bound, we have that  $|q(D)| \leq n^{3/2}$ .

This example shows that no query evaluation algorithm that uses join plans in the way that we explained in the beginning of this chapter has a runtime proportional to the number of tuples in the output. Furthermore, it raises the question whether such an algorithm even exists. We will show in this chapter that there exists an algorithm that can evaluate join queries  $q$  on databases  $D$  in time  $\tilde{O}(nm|q(D)|)$ , where  $n$  and  $m$  are the number of atoms and attributes in  $q$ , respectively. Here, the  $\tilde{O}$ -notation hides logarithmic factors, see Appendix [A](#).

### Worst-Case Optimal Join By Example

Consider again the join query  $q(A, B, C) :- R(A, B), S(B, C), T(C, A)$  and the database  $D$  in Figure [26.1](#). Here, we represent atoms  $R(a, b)$ ,  $S(a, b)$ , and  $T(a, b)$  with green, blue, and orange arrows from  $a$  to  $b$ , respectively.

The algorithm that we present in this section initializes three sets  $L_1$ ,  $L_2$ , and  $L_3$  as empty and computes  $q(D)$  in the following steps, each of which focuses on an attribute  $A$ ,  $B$ , or  $C$ :

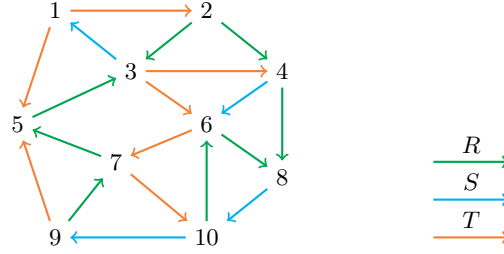


Fig. 26.1: Visualization of a database for illustrating worst-case optimal join

- (1) Compute  $L_1 := \pi_A(R \bowtie T)$ .
- (2) For each  $a \in L_1$ ,
  - compute the values  $b$  in  $\pi_B(R \bowtie S)$  such that  $(a, b) \in R$  and
  - add the pairs  $(a, b)$  to  $L_2$ .
- (3) For each  $(a, b) \in L_2$ ,
  - compute the values  $c$  in  $\pi_C(S \bowtie T)$  such that  $(b, c) \in S$  and  $(c, a) \in T$
  - add the triples  $(a, b, c)$  to  $L_3$ .
- (4) Return  $L_3$ .

On the database in Figure 26.1, the algorithm therefore computes

- (1)  $L_1 = \{5, 2, 6, 7, 4, 10\}$ ;
- (2)  $L_2 = \{(5, 3), (2, 3), (2, 4), (6, 8), (4, 8)\}$ ;
- (3)  $L_3 = \{(5, 3, 1), (2, 3, 1)\}$ .

## A Worst-Case Optimal Join Algorithm

We now describe a join algorithm that is worst-case optimal. As in Chapter 25, fix a join query  $q = R_1 \bowtie \cdots \bowtie R_n$  over schema  $\mathbf{S}$ , and assume that  $X = \{A_1, \dots, A_m\}$  is the set of attributes occurring in  $q$ . Given  $Y \subseteq X$ , recall that  $\mathcal{R}(Y) = \bigcup_{A_j \in Y} \mathcal{R}(A_j)$  is the set of those  $R_i$ 's that mention some attribute in  $Y$ , for  $i \in [n]$ . Moreover, assuming that  $\mathcal{R}(Y) = \{R_{i_1}, \dots, R_{i_\ell}\}$ , where  $1 \leq i_1 < \cdots < i_\ell \leq n$ , recall that query  $q_Y$  is

$$\pi_{Y \cap \mathbf{S}(R_{i_1})}(R_{i_1}) \bowtie \cdots \bowtie \pi_{Y \cap \mathbf{S}(R_{i_\ell})}(R_{i_\ell}) .$$

Algorithm 10, called AE-Join (short for “Attribute Elimination Join”) describes a worst-case optimal join algorithm by repeated elimination of attributes. The algorithm considers the attributes in  $\{A_1, \dots, A_m\}$  one by one

**Algorithm 10** AEJOIN( $q, D$ )**Input:** Join query  $q$  using attributes  $A_1, \dots, A_m$  and database  $D$ **Output:**  $q(D)$ 

```

1:  $L_0 := \{\emptyset\}$   $\triangleright L_0$  is the set with the empty tuple
2: for  $i = 1, \dots, m$  do
3:    $L_i := \emptyset$ 
4:   for each tuple  $\bar{a} \in L_{i-1}$  do
5:      $V := \bigcap_{R \in \mathcal{R}(A_i)} \pi_{\{A_i\}}(R^D \times \{\bar{a}\})$ 
6:      $L_i := L_i \cup (\{\bar{a}\} \times V)$ 
7: return  $L_m$ 

```

and completely deals with it before going to the next. As such, it iteratively computes sets  $L_0, \dots, L_m$ , where each  $L_i$  can be seen as the result of the query, *considering only the attributes*  $A_1, \dots, A_i$ . As we show next,  $L_m$  is the result of  $q$  on  $D$ .

**Proposition 26.1**

Given a join query  $q$  and a database  $D$ , we have that AEJOIN( $q, D$ ) returns  $q(D)$ .

*Proof.* We show by induction on  $i \in [m]$  that  $L_i = q_Z(D)$ , with  $Z = \{A_1, \dots, A_i\}$ . The result then follows since  $q_Z = q$  when  $Z = \{A_1, \dots, A_m\}$ . For the base case  $i = 1$ , this holds trivially, as we have that

$$L_1 = \bigcap_{R \in \mathcal{R}(A_1)} \pi_{\{A_1\}} R^D = q_{\{A_1\}}(D)$$

by definition. Consider now the inductive case  $i$ , for  $i > 1$ , and define  $Y = \{A_1, \dots, A_{i-1}\}$ . Notice that, by definition of  $q_{\{A_i\}}$ , it holds that  $V = q_{\{A_i\}}(R^{D_{\bar{a}}})$ , where  $D_{\bar{a}}$  is the database over  $\mathbf{S}'$  such that  $R^{D_{\bar{a}}} = \pi_{\{A_i\}}(R^D \times \{\bar{a}\})$  for each  $R \in \mathcal{R}(A_i)$ . Hence,

$$L_i = \bigcup_{\bar{a} \in q_Y(D)} (q_{\{A_i\}}(R^{D_{\bar{a}}}) \times \{\bar{a}\}).$$

But, as explained in Equation (25.8) and the analysis that precedes it, the right-hand side of this expression coincides with  $q_Z(D)$ .  $\square$

We invite the reader to compare how Algorithm 10 implements the procedure on the triangle query that we described before. Furthermore, notice that the exact execution of algorithm is highly dependent on the ordering  $A_1, \dots, A_m$  of the attributes in the input. Indeed, the for-loop considers increasing values of  $i$  from 1 to  $m$ , which corresponds to first considering  $A_1$ , then  $A_2$  etc. However, we will show that the algorithm is worst-case optimal independent of the ordering  $A_1, \dots, A_m$  of the attributes.

## Complexity Analysis

We will prove that Algorithm [10](#) runs in time

$$\tilde{O}\left(n \cdot m \cdot \prod_{j=1}^n |R_j^D|^{x_j}\right),$$

where  $(x_1, \dots, x_n)$  is a fractional edge cover of  $q$ .

In order to reach this bound, we will rely on a couple of assumptions. To facilitate the analysis, we assume that all relations  $R^D$  mentioned in  $q$  are non-empty. Notice that this assumption is without loss of generality, since it can be checked in constant time before evaluating the query. Second, we assume that there exists an ordering on the values in the database, since we will work with sorted lists.

Furthermore, we rely on the following complexity assumption on how the relations in  $D$  can be accessed:

**Claim 26.2.** *We can construct a data structure such that, given a relation  $R$  containing  $n$  tuples, an attribute name  $A_i$  of  $R$ , and a tuple  $\bar{a} = (A_1 : a_1, \dots, A_{i-1} : a_{i-1})$ , we can get access to an increasingly sorted list  $\lambda$  of values  $b$  such that  $(\bar{a}, b) \in \pi_{\{A_1, \dots, A_i\}} R$  in time  $O(i \log n)$ .*

Notice that not all attributes  $A_1, \dots, A_i$  have to be attributes of  $R$ . Finally, we will rely on the following claim about sorted lists:

**Claim 26.3.** *Given  $n$  lists  $\lambda_1, \dots, \lambda_n$  of increasingly sorted values, we can compute the set  $V = \{j \mid j \text{ occurs in each list } \lambda_1, \dots, \lambda_n\}$  in time  $\tilde{O}(n \cdot \min_{i \in [n]} |\lambda_i|)$ .*

We will prove Claims [26.2](#) and [26.3](#) in Chapter [27](#).

Using these claims, we show that the algorithm can be implemented such that for every  $i \in [m]$ ,

- (1) at the end of the  $i$ th for-loop in line [2](#), we have computed  $L_i$  in time

$$\tilde{O}\left(n \cdot i \cdot \prod_{j=1}^n |R_j^D|^{x_j}\right)$$

since the start of the algorithm and

- (2) given a tuple  $\bar{a}$ , the set  $V$  can be computed in time

$$\tilde{O}\left(n \cdot \prod_{R_j \in \mathcal{R}(A_i)} |R_j^D \ltimes \bar{a}|^{x_j}\right).$$

Notice that we measure the time in (1) and (2) differently. In (1) we count the number of steps since the start of the algorithm and in (2) we consider

the tuple  $\bar{a}$  as input, i.e., we only measure from the beginning of an iteration of the inner for-loop. Our desired bound then follows from (1), taking  $i = m$ .

We proceed by induction on  $i$ . To this end, assume that  $i = 1$ . In this case, notice that  $L_1 = V = \cap_{R \in \mathcal{R}(A_1)} \pi_{\{A_1\}} R^D$ . According to Claim 26.2 for each  $R \in \mathcal{R}(A_1)$ , we can get access to a sorted list  $\lambda_R$  containing the values  $\pi_{A_1} R$ . Using Claim 26.3, we can compute the set  $V$ , which are the values that occur in each of these lists, in time

$$\tilde{O}\left(n \cdot \min_{R \in \mathcal{R}(A_1)} |R^D|\right).$$

Furthermore, we have that

$$\begin{aligned} n \cdot \min_{R \in \mathcal{R}(A_1)} |R^D| &\leq n \cdot \prod_{R_j \in \mathcal{R}(A_1)} \left( \min_{R \in \mathcal{R}(A_1)} |R^D| \right)^{x_j} \\ &\leq n \cdot \prod_{R_j \in \mathcal{R}(A_1)} |R_j^D|^{x_j} \\ &\leq n \cdot \prod_{j=1}^n |R_j^D|^{x_j}, \end{aligned}$$

where the first inequality holds because  $(x_1, \dots, x_n)$  is a fractional edge cover and the last line holds because we assumed (without loss of generality) that all relations  $R_j^D$  are non-empty.

For general  $i$ , we first note that, given a tuple  $\bar{a}$ , we can compute the set  $V$ , analogously to the  $i = 1$  case, in time

$$\tilde{O}\left(n \cdot \min_{R \in \mathcal{R}(A_i)} |R^D \bowtie \bar{a}|\right)$$

using Claims 26.2 and 26.3. Furthermore,

$$n \cdot \min_{R \in \mathcal{R}(A_i)} |R^D \bowtie \bar{a}| \leq n \cdot \prod_{R_j \in \mathcal{R}(A_i)} |R_j^D \bowtie \bar{a}|^{x_j},$$

again analogously to the  $i = 1$  case.

We will now prove that, at the end of an iteration of the for-loop in line 2, we can compute  $L_i$  in

$$\tilde{O}\left(n \cdot i \cdot \prod_{j=1}^n |R_j^D|^{x_j}\right)$$

steps since the beginning of the entire algorithm. By induction, we can finish computing  $L_{i-1}$  after

$$\tilde{O}\left(n \cdot (i-1) \cdot \prod_{j=1}^n |R_j^D|^{x_j}\right)$$



steps, which we then use to compute  $L_i$ . So we need to show that we can compute  $L_i$  using an additional

$$\tilde{O}\left(n \cdot \prod_{j=1}^n |R_j^D|^{x_j}\right)$$

steps, making use of  $L_{i-1}$ . Indeed, given  $L_{i-1}$ , we can compute  $L_i$  in

$$\tilde{O}\left(\sum_{\bar{a} \in L_{i-1}} \left(n \cdot \prod_{R_j \in \mathcal{R}(A_i)} |R_j^D \ltimes \bar{a}|^{x_j}\right)\right)$$

steps and, taking  $Y = \{A_1, \dots, A_{i-1}\}$  we have that

$$\begin{aligned} \sum_{\bar{a} \in L_{i-1}} \left(n \cdot \prod_{R_j \in \mathcal{R}(A_i)} |R_j^D \ltimes \bar{a}|^{x_j}\right) &= n \cdot \sum_{\bar{a} \in L_{i-1}} \left(\prod_{R_j \in \mathcal{R}(A_i)} |R_j^D \ltimes \bar{a}|^{x_j}\right) \\ &= n \cdot \sum_{\bar{a} \in q_Y(D)} \left(\prod_{R_j \in \mathcal{R}(A_i)} |R_j^D \ltimes \bar{a}|^{x_j}\right) \\ &\leq n \cdot \sum_{\bar{a} \in q_Y(D)} \left(\prod_{R_j \in \mathcal{R}(\{A_i, \dots, A_m\})} |R_j^D \ltimes \bar{a}|^{x_j}\right) \\ &\leq n \cdot \prod_{j=1}^n |R_j^D|^{x_j} \end{aligned}$$

where the last inequality is by Lemma [25.11](#), the Query Decomposition Lemma.



---

## Bibliographic Comments

To be done.

## Part IV

---

### Expressive Languages



## Unions of Conjunctive Queries

The first, and simplest, addition to conjunctive queries is *union*, which leads to the language of union of conjunctive queries.

### Definition 28.1: Union of Conjunctive Queries

A *union of conjunctive queries* (UCQ) over a schema  $\mathbf{S}$  is an FO query  $\varphi(\bar{x})$  over  $\mathbf{S}$  where  $\varphi$  is a formula of the form

$$\varphi_1 \vee \cdots \vee \varphi_n$$

for  $n \geq 1$ , with  $\text{FV}(\varphi) = \text{FV}(\varphi_i)$  and  $\varphi_i(\bar{x})$  being a CQ, for every  $i \in [n]$ .

For notational convenience, we denote a UCQ  $q = \varphi(\bar{x})$  with  $\varphi = \varphi_1 \vee \cdots \vee \varphi_n$  as  $q_1 \cup \cdots \cup q_n$ , where  $q_i$  is the CQ  $\varphi_i(\bar{x})$ , for each  $i \in [n]$ . It is easy to verify that, given a database  $D$  of a schema  $\mathbf{S}$ , and a UCQ  $q = q_1 \cup \cdots \cup q_n$  over  $\mathbf{S}$ , it holds that  $q(D) = q_1(D) \cup \cdots \cup q_n(D)$ .

### Example 28.2: Union of Conjunctive Queries

Consider the relational schema from Example 3.2:

```

Person [ pid, pname, cid ]
Profession [ pid, pname ]
City [ cid, cname, country ]

```

The UCQ  $\varphi(y)$ , where  $\varphi = \varphi_1 \vee \varphi_2$  with

$$\varphi_1 = \exists x \exists z (\text{Person}(x, y, z) \wedge \text{Profession}(x, \text{'computer scientist'}) \wedge \text{City}(z, \text{'Athens'}, \text{'Greece'})).$$

and

$$\varphi_2 = \exists x \exists z (\text{Person}(x, y, z) \wedge \text{Profession}(x, \text{'computer scientist'}) \wedge \text{City}(z, \text{'Putú'}, \text{'Chile'})).$$

can be used to retrieve the list of names of computer scientists that were born in the city of Athens in Greece, or in the city of Putú in Chile.

## Union of Conjunctive Queries as a Fragment of FO

By definition, UCQ s use only relational atoms, conjunction ( $\wedge$ ), disjunction ( $\vee$ ), and existential quantification ( $\exists$ ). Therefore, every UCQ can be expressed using formulae from the fragment of FO that corresponds to the closure of relational atoms under  $\wedge$ ,  $\vee$  and  $\exists$ ; we refer to this fragment of FO as  $\text{FO}^{\text{rel}}[\wedge, \vee, \exists]$ . Interestingly, we can show that the converse is also true, which leads to the following expressive power result:

### Theorem 28.3

The language of UCQ s and the language of  $\text{FO}^{\text{rel}}[\wedge, \vee, \exists]$  queries are equally expressive.

*Proof.* As discussed above, by definition, a UCQ is trivially an  $\text{FO}^{\text{rel}}[\wedge, \vee, \exists]$  query. The interesting task is to show that an  $\text{FO}^{\text{rel}}[\wedge, \vee, \exists]$  query  $\varphi(\bar{x})$  can be equivalently expressed as a UCQ. This is done in three main steps:

- First, we propagate disjunction by using the following simple rules:

$$\chi \wedge (\xi \vee \psi) \rightsquigarrow (\chi \wedge \xi) \vee (\chi \wedge \psi) \quad \text{and} \quad \exists x (\chi \vee \psi) \rightsquigarrow \exists x \chi \vee \exists x \psi,$$

where  $\chi, \xi$  and  $\psi$  are FO formulae. By applying the above rules, we can convert the formula  $\varphi$  into an equivalent formula of the form

$$\varphi_1 \vee \cdots \vee \varphi_n$$

for  $n \geq 1$ , where  $\varphi_i$  is a formula from  $\text{FO}^{\text{rel}}[\wedge, \exists]$ , for each  $i \in [n]$ .

- We then convert  $\varphi_i$ , for each  $i \in [n]$ , into a formula of the form  $\exists \bar{x}_i \varphi'_i$ , where  $\varphi'_i$  is a quantifier-free conjunction of relational atoms. This is done in the same way as the transformation of an  $\text{FO}^{\text{rel}}[\wedge, \exists]$  query into a CQ (see Example 12.5): we first rename variables in order to ensure that bound variables do not repeat, and then push the existential quantifiers outside.
- After applying the above steps, we end up with a formula  $\psi$  of the form

$$\psi_1 \vee \cdots \vee \psi_n$$

where  $\psi_i = \exists \bar{x}_i \varphi'_i$  and  $\varphi'_i$  is a quantifier-free conjunction of relational atoms, for each  $i \in [n]$ . Observe also that during the above two steps we

have not altered the set of free variables of  $\varphi$ , i.e.,  $\text{FV}(\varphi) = \text{FV}(\psi)$ . Hence,  $\psi(\bar{x})$  is a syntactically valid FO query that is equivalent to  $\varphi(\bar{x})$ . However, it should not be overlooked that  $\psi(\bar{x})$  is not yet a UCQ since there is no guarantee that  $\text{FV}(\psi) = \text{FV}(\psi_i)$ , for each  $i \in [n]$ . In this final step, we explain how  $\psi(\bar{x})$  can be converted into an equivalent UCQ.

Assume that we have access to a unary relation **dom** that stores all the values in the given database. In other words, we assume that every database  $D$  comes with a unary relation **dom** such that, for every  $a \in \text{Dom}(D)$ ,  $\text{dom}(a) \in D$ . In this case, it is easy to see that  $\psi'(\bar{x})$  with

$$\psi' = \bigvee_{i=1}^n \left( \psi_i \wedge \bigwedge_{y \in \text{FV}(\psi) - \text{FV}(\psi_i)} \text{dom}(y) \right)$$

is a syntactically valid UCQ that is equivalent to  $\psi(\bar{x})$ . Indeed,  $\text{FV}(\psi) = \text{FV}(\psi')$ , each disjunct  $\psi'_i$  of  $\psi$  is such that  $\text{FV}(\psi') = \text{FV}(\psi'_i)$ , and  $\psi'_i(\bar{x})$  is a CQ. Moreover,  $\psi(\bar{x})(D) = \psi'(\bar{x})(D)$  for every database  $D$  equipped with the unary relation **dom**. It remains to show that the relational atoms of the form **dom**( $\cdot$ ) in  $\psi'$  can be eliminated with the help of disjunction.

Consider a formula of the form  $\exists \bar{y} \chi \wedge \text{dom}(u)$ , where  $\chi$  is an arbitrary formula over a schema **S**. It is easy to verify that it is equivalent to

$$\bigvee_{R \in \mathbf{S}} \bigvee_{i \in \{1, \dots, \text{ar}(R)\}} \exists \bar{y} \exists z_1 \dots \exists z_{\text{ar}(R)-1} \chi \wedge R(z_1, \dots, z_{i-1}, u, z_i, \dots, z_{\text{ar}(R)-1})$$

where  $z_1, \dots, z_{i-1}, z_{i+1}, \dots, z_{\text{ar}(R)-1}$  are new variables not occurring in  $\chi$ . Indeed, if  $a$  is a value that occurs in the input database, then  $a$  must occur in a tuple of some relation  $R$  at some position  $i \in \{1, \dots, \text{ar}(R)\}$ . Using this transformation, we can eventually eliminate all the relational atoms of the form **dom**( $\cdot$ ) in  $\psi'$ , and obtain a formula  $\psi'' = \psi'_1 \vee \dots \vee \psi'_m$  such that  $\psi'(\bar{x})$  and  $\psi''(\bar{x})$  are equivalent queries, and  $\psi''_i(\bar{x})$  is a syntactically valid CQ for each  $i \in [m]$ , which in turn implies that  $\psi''(\bar{x})$  is a UCQ.  $\square$

Let  $\text{FO}[\wedge, \vee, \exists]$  be the fragment of FO that corresponds to the closure of relational atoms *and* equational atoms under  $\wedge$ ,  $\vee$  and  $\exists$ . This is known as the *existential positive* fragment of FO, and is typically denoted as  $\exists\text{FO}^+$ ; hence, from now on, by  $\exists\text{FO}^+$  we actually mean  $\text{FO}[\wedge, \vee, \exists]$ . In other words,  $\exists\text{FO}^+$  is the fragment of FO obtained by explicitly adding equality to  $\text{FO}^{\text{rel}}[\wedge, \vee, \exists]$ . It is easy to show that the language of  $\exists\text{FO}^+$  queries is strictly more expressive than the language of UCQs, and thus, by Theorem 28.3, also the language of  $\text{FO}^{\text{rel}}[\wedge, \vee, \exists]$  queries. Consider, for example, the  $\exists\text{FO}^+$  query  $q = \varphi(x)$  with  $\varphi = (x = a)$ , where  $a \in \text{Const}$ . Clearly,  $q(D) = \{(a)\}$  for every database  $D$ , even if  $a \notin \text{Dom}(D)$ . However, given a database  $D'$  such that  $a \notin \text{Dom}(D')$ , for every unary UCQ  $q'$ , it holds that  $\{(a)\} \not\subseteq q'(D')$  since the output of a CQ, and thus of a UCQ, on  $D'$  consists of tuples of constants from  $\text{Dom}(D')$ .



Observe that the  $\exists\text{FO}^+$  query  $q$  above uses equality among a variable and a constant. It turns out that this is crucial for showing that  $\exists\text{FO}^+$  queries form a strictly more expressive language than UCQ s. Interestingly, the language of queries based on  $\text{FO}^{\text{rel}, \text{var}=[\wedge, \vee, \exists]}$ , that is, the fragment of FO that corresponds to the closure of relational atoms and equational atoms of the form  $x = y$ , where both  $x$  and  $y$  are variables, under  $\wedge$ ,  $\vee$  and  $\exists$ , has the same expressive power as the language of UCQ s.

**Theorem 28.4**

The language of UCQ s and the language of  $\text{FO}^{\text{rel}, \text{var}=[\wedge, \vee, \exists]}$  queries are equally expressive.

*Proof.* By definition, a UCQ is trivially an  $\text{FO}^{\text{rel}, \text{var}=[\wedge, \vee, \exists]}$  query. It remains to show that an  $\text{FO}^{\text{rel}, \text{var}=[\wedge, \vee, \exists]}$  query  $\varphi(\bar{x})$  can be equivalently expressed as a UCQ. This is done by first observing that equational atoms that mention only variables can be eliminated by using atoms of the form  $\text{dom}(\cdot)$ , where as before  $\text{dom}$  is a unary relation that stores all the values in the given database. Indeed, for each equational atom  $x = y$  in  $\varphi$ , we replace  $y$  by  $x$  everywhere in  $\varphi$  and  $\bar{x}$ , and add the atom  $\text{dom}(x)$  to the conjunction. To see why the latter is needed, consider, for example, the query  $\psi(x, y)$  with  $\psi = (x = y)$ . We cannot just throw away the equational atom; instead, this query is equivalent to  $\psi'(x, x)$  with  $\psi' = \text{dom}(x)$ . Hence, after the above transformation, we obtain an  $\text{FO}^{\text{rel}}[\wedge, \vee, \exists]$  query  $\varphi'(\bar{x}')$  that is equivalent to  $\varphi(\bar{x})$  over databases equipped with the unary relation  $\text{dom}$ . Since, as discussed in the proof of Theorem 28.3,  $\text{dom}(\cdot)$  atoms can be eliminated with the help of disjunction, we can convert  $\varphi'(\bar{x}')$  into an  $\text{FO}^{\text{rel}}[\wedge, \vee, \exists]$  query  $\varphi''(\bar{x}')$  that is equivalent to  $\varphi(\bar{x})$  over *all* databases. Finally, by Theorem 28.3, we know that there exists a UCQ that is equivalent to  $\varphi''(\bar{x}')$ , and thus to  $\varphi(\bar{x})$ , and the claim follows.  $\square$

We have seen that UCQ s are not powerful enough for expressing every  $\exists\text{FO}^+$  query. We have also seen that the key reason for this is the fact that UCQ s, although can express equality among variables, cannot express equality among variables and constants. The question that comes up is whether the addition of equality among variables and constants to UCQ s leads to a language that can express every  $\exists\text{FO}^+$  query. A UCQ *with variable-constant equality*  $\varphi(\bar{x})$  is defined as a UCQ with the only difference that a disjunct of  $\varphi$  can be a conjunction of relational atoms and equational atoms of the form  $(x = a)$ , where  $x$  is a variable and  $a$  is a constant. By using the same ideas as in the proofs of Theorems 28.3 and 28.4, it is easy to show the following:

**Theorem 28.5**

The language of UCQ s with variable-constant equality and the language of  $\exists\text{FO}^+$  queries are equally expressive.

It is important to stress that the transformations described in the proofs of the above expressive power results can be costly. Already, the transformation of an  $\text{FO}^{\text{rel}}[\wedge, \vee, \exists]$  query into a UCQ may lead to an exponentially sized query. Consider, e.g., an  $\text{FO}^{\text{rel}}[\wedge, \vee, \exists]$  query  $\varphi(\bar{x})$ , where  $\varphi$  is of the form

$$(\varphi_1 \vee \varphi'_1) \wedge \cdots \wedge (\varphi_n \vee \varphi'_n)$$

and  $\varphi_i(\bar{x})$ ,  $\varphi'_i(\bar{x})$  are CQs, for every  $i \in [n]$ . Representing  $\varphi(\bar{x})$  as a UCQ requires transforming an FO formula in conjunctive normal form into an FO formula in disjunctive normal form, resulting in a UCQ consisting of  $2^n$  CQs. Consequently, even though UCQ s and  $\text{FO}^{\text{rel}}[\wedge, \vee, \exists]$  (or  $\text{FO}^{\text{rel, var}}[\wedge, \vee, \exists]$ ) queries have the same expressive power, some problems related to them will have different complexity (whenever the size of the query matters). The same holds for UCQ s with variable-constant equality and  $\exists\text{FO}^+$  queries.

## Union of Conjunctive Queries as a Fragment of RA

We know, by Theorem 12.7, that the language of CQs has the same expressive power as the language of SPJ queries. Recall that SPJ is the fragment of RA that is built from base expressions  $R \in \text{Rel}$  (crucially, base expressions of the form  $\{a\}$  with  $a \in \text{Const}$  are not included), and allows for selection, projection, and Cartesian product. Furthermore, conditions in selections are conjunctions of equalities. It should not come as a surprise the fact that by adding union to SPJ we get a fragment of RA, called *select-project-join-union* (SPJU), that has the same expressive power as UCQ s.

### Theorem 28.6

The language of UCQ s and the language of SPJU queries are equally expressive.

*Proof.* The fact that every UCQ can be expressed as an SPJU query immediately follows from Theorem 12.7, which shows that every CQ can be expressed as an SPJ query. Indeed, a UCQ  $q_1 \cup \cdots \cup q_n$  is equivalent to the SPJU query  $e_1 \cup \cdots \cup e_n$ , where  $e_i$  is an SPJ query that is equivalent to  $q_i$ , for each  $i \in [n]$ .

Consider now an SPJU  $k$ -ary query  $e$ . We proceed to show that  $e$  can be expressed as a UCQ. This is done in three main steps:

- First, we propagate union through other operations to become the outermost operation by applying the following simple rules:

$$\begin{aligned} \sigma_\theta(e_1 \cup e_2) &\rightsquigarrow \sigma_\theta(e_1) \cup \sigma_\theta(e_2) \\ \pi_\alpha(e_1 \cup e_2) &\rightsquigarrow \pi_\alpha(e_1) \cup \pi_\alpha(e_2) \\ e_1 \times (e_2 \cup e_3) &\rightsquigarrow (e_1 \times e_2) \cup (e_1 \times e_3). \end{aligned}$$

By applying the above rules, we get an SPJU query

$$e' = e_1 \cup \dots \cup e_n$$

where  $e_i$  is an SPJ query, for each  $i \in [n]$ .

- Let  $\varphi_i(x_i^1, \dots, x_i^k)$  be the CQ that is equivalent to the SPJ query  $e_i$ , for each  $i \in [n]$ ; such a CQ always exists due to Theorem 12.7. Let

$$\varphi = \varphi_1 \vee \dots \vee \varphi_n.$$

- At this point, one may think that the above step leads to the desired UCQ that is equivalent to the SPJU query  $e$ . However, there is no guarantee that, for each  $i, j \in [n]$  with  $i \neq j$ , it holds that  $\text{FV}(\varphi_i) = \text{FV}(\varphi_j)$ . We proceed to convert  $\varphi$  into an FO formula  $\psi$  such that  $\psi(z_1, \dots, z_k)$ , where  $z_1, \dots, z_k$  are distinct variables not occurring in  $\varphi$ , is an  $\text{FO}^{\text{rel}, \text{var}=[\wedge, \vee, \exists]}$  query that is equivalent to the query  $e'$ , and thus, to the query  $e$ . This suffices to show our claim since, by Theorem 28.4, we get that  $\psi(z_1, \dots, z_k)$  can be equivalently expressed as a UCQ.

Consider an arbitrary disjunct  $\varphi_i$  of  $\varphi$ . Let  $P_{\varphi_i} = \{P_1, \dots, P_\ell\}$ , where  $\ell \leq k$  is the number of distinct variables occurring in  $(x_i^1, \dots, x_i^k)$ , be the partition of the set of integers  $[k]$  such that, for every  $j, j' \in [k]$ ,  $j, j'$  belong to the same set of  $P_{\varphi_i}$  if and only if  $x_i^j = x_i^{j'}$ . For example, with  $k = 5$  and  $(x_i^1, \dots, x_i^5) = (x, y, x, z, y)$ ,  $P_{\varphi_i} = \{\{1, 3\}, \{2, 5\}, \{4\}\}$ . Let  $\psi_i$  be the formula obtained from  $\varphi_i$  as follows: for every set  $\{j_1, \dots, j_m\} \in P_{\varphi_i}$ , replace in  $\varphi_i$  the variable  $x_i^{j_1}$  (note that  $x_i^{j_1} = x_i^{j_2} = \dots = x_i^{j_m}$ ) with the variable  $z_{j_1}$ , and add as a conjunct the conjunction of equational atoms

$$\bigwedge_{j \in \{j_2, \dots, j_m\}} (z_{j_1} = z_j).$$

It is easy to verify that  $\psi(z_1, \dots, z_k)$  with

$$\psi = \psi_1 \vee \dots \vee \psi_n$$

is an  $\text{FO}^{\text{rel}, \text{var}=[\wedge, \vee, \exists]}$  query that is equivalent to  $e'$ , as needed.  $\square$

The following is an immediate corollary of Theorems 28.3, 28.4 and 28.6.

#### Corollary 28.7

The language of  $\text{FO}^{\text{rel}}[\wedge, \vee, \exists]$  (or even  $\text{FO}^{\text{rel}, \text{var}=[\wedge, \vee, \exists]}$ ) queries and the language of SPJU queries are equally expressive.

It should be clear that the inability of SPJ queries to state base expressions of the form  $\{a\}$  with  $a \in \text{Const}$  it is crucial for the validity of Theorem 28.6. Indeed, the addition of such base expressions to the SPJU fragment of RA leads

to the strictly more expressive language of *positive relational algebra* ( $\text{RA}^+$ ) queries. Interestingly, by providing a proof similar to that of Theorem 28.6, we can show that adding base expressions of the form  $\{a\}$  with  $a \in \text{Const}$  to SPJU corresponds to the addition of variable-constant equality to UCQ s.

**Theorem 28.8**

The language of UCQ s with variable-constant equality and the language of  $\text{RA}^+$  queries are equally expressive.

The following is an immediate corollary of Theorems 28.5 and 28.8 that relates the languages of  $\exists\text{FO}^+$  queries and  $\text{RA}^+$  queries.

**Corollary 28.9**

The language of  $\exists\text{FO}^+$  queries and the language of  $\text{RA}^+$  queries are equally expressive.

Note that the transformations described in the proofs of the above expressive power results (in particular, in the proof of Theorem 28.6) can be costly. For example, given an SPJU query of the form

$$(e_1 \cup e'_1) \times \cdots \times (e_n \cup e'_n)$$

where  $e_i, e'_i$  are SPJ queries, for each  $i \in [n]$ , after propagating the union during the first step of the transformation in the proof of Theorem 28.6, we get an SPJU query that is the union of  $2^n$  SPJ queries.

## Preservation Under Homomorphisms

We have already seen that CQs are preserved under homomorphisms (Proposition 13.6). In other words, given a  $k$ -ary CQ  $q = \varphi(\bar{x})$  over a schema  $\mathbf{S}$ , for every two databases  $D$  and  $D'$  of  $\mathbf{S}$ , and tuples  $\bar{a} \in \text{Dom}(D)^k$  and  $\bar{b} \in \text{Dom}(D')^k$ ,

$$(D, \bar{a}) \rightarrow_{\text{Dom}(\varphi)} (D', \bar{b}) \text{ and } \bar{a} \in q(D) \text{ implies } \bar{b} \in q(D').$$

It is not difficult to show that preservation under homomorphisms extends to UCQ s with variable-constant equality (and thus, to  $\exists\text{FO}^+$  queries).

**Proposition 28.10**

Every UCQ with variable-constant equality is preserved under homomorphisms.

*Proof.* Let  $q = \varphi(\bar{x})$  be a  $k$ -ary UCQ with variable-constant equality over a schema  $\mathbf{S}$ . Assume that  $\varphi = \varphi_1 \vee \cdots \vee \varphi_n$ , and let  $q_i$  be the query  $\varphi_i(\bar{x})$ , for each  $i \in [n]$ ; note that these queries are not CQs since they use equational atoms. Consider two databases  $D$  and  $D'$  of  $\mathbf{S}$ , and tuples  $\bar{a} \in \text{Dom}(D)^k$  and  $\bar{b} \in \text{Dom}(D')^k$  such that  $(D, \bar{a}) \rightarrow_{\text{Dom}(\varphi)} (D', \bar{b})$  and  $\bar{a} \in q(D)$ . Since  $q(D) = \bigcup_{i=1}^n q_i(D)$ , it is clear that  $\bar{a} \in q_i(D)$  for some  $i \in [n]$ . By providing a proof similar to that of Proposition 13.6, which shows that CQs are preserved under homomorphisms, we can show that  $q_i$  is preserved under homomorphisms. Therefore,  $\bar{b} \in q_i(D')$ , which in turn implies that  $\bar{b} \in q(D')$ , as needed.  $\square$

It is far more remarkable, though, that the converse is true, that is, every FO query that is preserved under homomorphisms can be expressed as a UCQ with variable-constant equality.

### Theorem 28.11

Consider an FO query  $q$  that is preserved under homomorphisms. There exists a UCQ with variable-constant equality that is equivalent to  $q$ .

This is a deep result whose proof is beyond the scope of this book, but it is very important and found many applications in the foundations of databases. An immediate corollary of Proposition 28.10 and Theorem 28.11 is that:

### Corollary 28.12

The language of UCQs with variable-constant equality and the language of FO queries preserved under homomorphisms are equally expressive.

It is important to stress that many preservation results known in logic are true on arbitrary (finite or infinite) structures, but fail on finite structures. Preservation results of this kind can be transferred to the database setting only for possibly infinite databases, but not for (finite) databases, which is of course what is of interest to us. Remarkably, Theorem 28.11 is a rare exception that holds in the case of (finite) databases.

## Query Evaluation

A general rule of thumb is that whatever is true about the evaluation of CQs, is true about the evaluation of their unions. We illustrate this by analyzing the combined complexity of evaluating UCQs and acyclic UCQs; concerning the data complexity, both problems are in DLOGSPACE due to Theorem 7.3.

### Evaluation of UCQs

We first concentrate on UCQ-Evaluation. Recall that this is the problem of checking whether  $\bar{a} \in q(D)$  for a UCQ  $q$ , a database  $D$ , and a tuple  $\bar{a}$  over  $\text{Dom}(D)$ . We show that indeed UCQ-Evaluation has the same (combined) complexity as CQ-Evaluation.

**Theorem 28.13**

UCQ-Evaluation is NP-complete.

*Proof.* It is clear that the NP-hardness is inherited from CQ-Evaluation, which we know is NP-hard (Theorem 14.1). We proceed to show the upper bound. Consider a UCQ  $q(\bar{x})$  of the form  $q_1 \cup \dots \cup q_n$ , a database  $D$ , and a tuple  $\bar{a} \in \text{Dom}(D)$ . By Theorem 13.2 for  $i \in [n]$ ,  $\bar{a} \in q_i(D)$  if and only if  $(q_i, \bar{x}) \rightarrow (D, \bar{a})$ . Thus, we need to show that checking whether there exists an integer  $i \in [n]$  and a homomorphism from  $(q_i, \bar{x})$  to  $(D, \bar{a})$  is in NP. This is done by guessing an integer  $i \in [n]$  and a function  $h : \text{Dom}(A_{q_i}) \rightarrow \text{Dom}(D)$ , and then verifying that  $h$  is a homomorphism from  $(A_{q_i}, \bar{x})$  to  $(D, \bar{a})$ , i.e.,  $h$  is the identity on  $\text{Dom}(A_{q_i}) \cap \text{Const}$ ,  $R(\bar{u}) \in A_{q_i}$  implies  $R(h(\bar{u})) \in D$ , and  $h(\bar{x}) = \bar{a}$ . Since all the above steps are feasible in polynomial time, the claim follows.  $\square$

It is not difficult to show, by providing a proof similar to that of Theorem 28.13, that evaluating UCQs with variable-constant equality remains in NP. What is more interesting is the fact that evaluating  $\exists\text{FO}^+$  queries, as well as  $\text{RA}^+$  queries, is also in NP. We have seen that every  $\exists\text{FO}^+$  can be converted into an equivalent UCQ with variable-constant equality (Theorem 28.5); the same holds for  $\text{RA}^+$  queries (Theorem 28.8). However, the conversion can be very costly; it may take, in general, exponential time. Therefore, knowing that evaluating UCQs with variable-constant equality is in NP does not immediately imply that evaluating  $\exists\text{FO}^+$  and  $\text{RA}^+$  queries is also in NP. Hence, one has to adopt a more refined procedure than simply converting the given  $\exists\text{FO}^+$  or  $\text{RA}^+$  query into a UCQ with variable-constant equality (see Exercise 4.1).

### Evaluation of Acyclic UCQs

We now consider the problem of evaluating *acyclic* UCQs, that is, UCQs of the form  $q_1 \cup \dots \cup q_n$ , where, for each  $i \in [n]$ , the CQ  $q_i$  is acyclic. Recall that a CQ is acyclic if its associated hypergraph is acyclic (Definition 18.4). We further know that acyclic CQs can be efficiently evaluated. More precisely, by Theorem 19.3, checking whether  $\bar{a} \in q(D)$  for an acyclic CQ query  $q$ , a database  $D$ , and a tuple  $\bar{a}$  over  $\text{Dom}(D)$  is feasible in time  $O(\|D\| \cdot \log \|D\| \cdot \|q\|)$ . It is easy to show that the same holds for acyclic UCQs.

**Theorem 28.14**

Consider an acyclic UCQ  $q$ , a database  $D$ , and a tuple  $\bar{a}$  over  $\text{Dom}(D)$ . Checking whether  $\bar{a} \in q(D)$  is feasible in time  $O(\|D\| \cdot \log \|D\| \cdot \|q\|)$ .

*Proof.* We need to check whether  $\bar{a} \in \bigcup_{i=1}^n q_i(D)$ . By Theorem 19.3, for every  $i \in [n]$ , checking whether  $\bar{a} \in q_i(D)$  is feasible in time  $O(\|D\| \cdot \log \|D\| \cdot \|q_i\|)$ . This implies that checking whether  $\bar{a} \in \bigcup_{i=1}^n q_i(D)$  can be done in time

$$O\left(\|D\| \cdot \log \|D\| \cdot \sum_{i \in [n]} \|q_i\|\right) = O(\|D\| \cdot \log \|D\| \cdot \|q\|)$$

and the claim follows.  $\square$

## Static Analysis of Unions of Conjunctive Queries

In this chapter, we study the containment and equivalence problems for unions of conjunctive queries, as well as the task of minimizing such queries.

### Containment and Equivalence

We first focus on UCQ-Containment, the problem of deciding whether a UCQ  $q$  is contained in a UCQ  $q'$ , that is, whether  $q(D) \subseteq q'(D)$  for every database  $D$ . We show that it has the same complexity as CQ-Containment. But first we present a useful result that characterizes when  $q$  is contained in  $q'$  in terms of containment of the individual CQs occurring in  $q$  and  $q'$ .

#### Proposition 29.1

Consider two UCQs  $q = q_1 \cup \dots \cup q_n$  and  $q' = q'_1 \cup \dots \cup q'_m$ . The following are equivalent:

1.  $q \subseteq q'$ .
2. For every  $i \in [n]$ , there exists  $j \in [m]$  such that  $q_i \subseteq q'_j$ .

*Proof.* We first show that (1) implies (2). We assume that the output tuple of  $q$  and  $q'$  is  $\bar{x}$  and  $\bar{x}'$ , respectively. This means that, for each  $i \in [n]$  and  $j \in [m]$ , the output tuple of  $q_i$  and  $q'_j$  is  $\bar{x}$  and  $\bar{x}'$ , respectively. Consider an arbitrary integer  $i \in [n]$ . Recall that  $G_{A_{q_i}}$  is a homomorphism from  $A_{q_i}$  to the grounding of  $A_{q_i}$ . Thus, by Theorem 13.2, we get that  $G_{A_{q_i}}(\bar{x}) \in q_i(G_{A_{q_i}}(A_{q_i}))$ . Since  $q \subseteq q'$ , we have that  $G_{A_{q_i}}(\bar{x}) \in q'(G_{A_{q_i}}(A_{q_i}))$ , and hence,  $G_{A_{q_i}}(\bar{x}) \in q'_j(G_{A_{q_i}}(A_{q_i}))$  for some  $j \in [m]$ . As in the proof of the Homomorphism Theorem (Theorem 15.4), we can show that there exists a homomorphism  $h$  from  $(q'_j, \bar{x}')$  to  $(q_i, \bar{x})$ . From the Homomorphism Theorem itself, we conclude then that  $q_i \subseteq q'_j$ .

For showing that (2) implies (1), assume that  $\bar{a} \in q(D)$  for some database  $D$  and tuple  $\bar{a}$  over  $\text{Dom}(D)$ . Clearly, there exists  $i \in [n]$  such that  $\bar{a} \in q_i(D)$ .



By hypothesis, there exists  $j \in [m]$  such that  $q_i \subseteq q'_j$ , and thus,  $\bar{a} \in q'_j(D)$ . The latter implies that  $\bar{a} \in q'(D)$ , which in turn shows that  $q \subseteq q'$ , as needed.  $\square$

We are now ready to pinpoint the complexity of UCQ-Containment.

**Theorem 29.2**

UCQ-Containment is NP-complete.

*Proof.* It is clear that the NP-hardness is inherited from CQ-Containment, which we know is NP-hard (Theorem 15.3). Consider now two UCQs  $q = q_1 \cup \dots \cup q_n$  and  $q' = q'_1 \cup \dots \cup q'_m$ . We proceed to show that checking whether  $q \subseteq q'$  is in NP. We assume that, for  $i \in [n]$  and  $j \in [m]$ , the output tuple of  $q_i$  and  $q'_j$  is  $\bar{x}$  and  $\bar{x}'$ , respectively. By the Homomorphism Theorem (Theorem 15.4) and Proposition 29.1, to check whether  $q \subseteq q'$  it suffices to do the following:

- for each  $i \in [n]$ , guess an integer  $j_i \in [m]$ , and a function  $h : \text{Dom}(A_{q'_{j_i}}) \rightarrow A_{q_i}$ , and
- for each  $i \in [n]$ , verify that  $h_i$  is a homomorphism from  $(q'_{j_i}, \bar{x}')$  to  $(q_i, \bar{x})$ .

Since both steps are feasible in polynomial time, we conclude that deciding whether  $q \subseteq q'$  is in NP, and the claim follows.  $\square$

An immediate corollary of Theorem 29.2 is that the equivalence problem for UCQs, that is, given two UCQs  $q, q'$ , check whether  $q \equiv q'$ , is in NP since it boils down to two containment checks:  $q \subseteq q'$  and  $q' \subseteq q$ . The NP-hardness is inherited from CQ-Equivalence (Theorem 15.8).

**Corollary 29.3**

UCQ-Equivalence is NP-complete.

Recall that query evaluation remains NP-complete even if we consider  $\text{FO}^{\text{rel}}[\wedge, \vee, \exists]$  and SPJU queries (or even  $\exists\text{FO}^+$  and  $\text{RA}^+$  queries), despite the fact that these languages allow us to express UCQs in a more succinct way. However, this is not true in the case of containment, where we can show that the complexity increases. We illustrate this for SPJU-Containment, that is, the problem of deciding whether an SPJU query is contained in another SPJU query. The treatment for  $\text{FO}^{\text{rel}}[\wedge, \vee, \exists]$ , and the more expressive languages  $\exists\text{FO}^+$  and  $\text{RA}^+$ , is similar and is left as an exercise (see Exercise 4.3).

We proceed to show that SPJU-Containment is  $\Pi_2^P$ -complete. This essentially tells us that, given two SPJU queries  $e$  and  $e'$ , the problem of deciding whether  $e \not\subseteq e'$  is in  $\Sigma_2^P = \text{NP}^{\text{NP}}$ , i.e., it can be solved via a nondeterministic algorithm that runs in polynomial time assuming that it has access to an

oracle that can solve any problem in NP. In other words, the complement of SPJU-Containment is  $\Sigma_2^P$ -complete.<sup>1</sup> To show this we need some preparation.

We associate to an SPJU query  $e$  a set of SPJ queries, denoted  $\text{SPJ}(e)$ . This is done by induction on the structure of  $e$ ; essentially, for every union  $e_1 \cup e_2$  that occurs in  $e$ , we look at possible ways of resolving this union, i.e., choosing  $e_1$  or  $e_2$ . The set  $\text{SPJ}(e)$  is formally defined as follows:

$$\text{SPJ}(e) = \begin{cases} \{R\} & \text{if } e = R \\ \{\sigma_\theta(e'') \mid e'' \in \text{SPJ}(e')\} & \text{if } e = \sigma_\theta(e') \\ \{\pi_\alpha(e'') \mid e'' \in \text{SPJ}(e')\} & \text{if } e = \pi_\alpha(e') \\ \{e'_1 \times e'_2 \mid e'_1 \in \text{SPJ}(e_1) \text{ and } e'_2 \in \text{SPJ}(e_2)\} & \text{if } e = e_1 \times e_2 \\ \text{SPJ}(e_1) \cup \text{SPJ}(e_2) & \text{if } e = e_1 \cup e_2. \end{cases}$$

The following is easily shown by structural induction.

#### Proposition 29.4

Consider an SPJU query  $e$ , and assume that  $\text{SPJ}(e) = \{e_1, \dots, e_n\}$ . For every  $i \in [n]$ ,  $e_i$  is an SPJ query, and  $e \equiv e_1 \cup \dots \cup e_n$ .

It is clear that Proposition 29.4 provides an algorithm for solving the problem SPJU-Containment: given two SPJU queries  $e_1$  and  $e_2$ , compute the sets  $\text{SPJ}(e_1)$  and  $\text{SPJ}(e_2)$ , and then check whether, for every  $e'_1 \in \text{SPJ}(e_1)$ , there exists  $e'_2 \in \text{SPJ}(e_2)$  such that  $e'_1 \subseteq e'_2$ ; the latter is essentially a containment check among two CQs, since an SPJ query can be easily converted into a CQ, which can be performed using the algorithm underlying Theorem 15.3. However, this is a very naive algorithm, which only shows that SPJU-Containment is in EXPTIME. Indeed, it explicitly constructs the sets  $\text{SPJ}(e_1)$  and  $\text{SPJ}(e_2)$ , which are in general of exponential size, and then performs exponentially many containment checks among SPJ queries. To establish the desired  $\Pi_2^P$  upper bound, we need to rely on a refined version of the above algorithm.

The key observation towards such a refined procedure is that finding a query  $e' \in \text{SPJ}(e)$ , for an SPJU query  $e$ , amounts to “resolving” each union in  $e$ . In other words, having the parse tree  $T_e$  of the expression  $e$ , for every union node in  $T_e$ , with two subtrees under it, we keep only one of those subtrees, while the other one is replaced by  $\perp$ . The obtained tree  $T'_e$  is essentially the parse tree of an SPJ query from  $\text{SPJ}(e)$ . Consider, for example, the query

$$e = (e_1 \cup e_2) \times (e_3 \cup (e_4 \cup e_5)),$$

<sup>1</sup> This is actually the first time in the book that we encounter the complexity classes  $\Sigma_2^P$  and  $\Pi_2^P$ , which contain NP. For further details see Appendix B.

where  $e_1, \dots, e_5$  are  $\cup$ -free. One way to resolve the union nodes in  $T_e$  is

$$(e_1 \cup \perp) \times (e_3 \cup \perp),$$

which leads to  $e_1 \times e_3 \in \text{SPJ}(e)$ . Another way is

$$(\perp \cup e_2) \times (\perp \cup (\perp \cup e_5)),$$

which leads to  $e_2 \times e_5 \in \text{SPJ}(e)$ . If union occurs  $k$  times in  $e$ , this gives us  $2^k$  expressions that can result from resolving those union nodes, as each one gives us two choices. However, each way to resolve the union nodes in  $T_e$  can be carried out in polynomial time, or, in other words, we can guess any SPJ query from  $\text{SPJ}(e)$  in polynomial time. This fact allows us to devise the refined procedure for **SPJU-Containment**.

Before doing this, we need to establish an intermediate complexity result, which will be crucial in the complexity analysis of this refined procedure. Furthermore, it illustrates the fact that restricting the language of the left-hand side query in the containment check has an impact on the complexity.

### Proposition 29.5

Consider an SPJ query  $e$ , and an SPJU query  $e'$ . The problem of deciding whether  $e \subseteq e'$  is in NP.

*Proof.* By Proposition 29.4 it suffices to show that checking whether there exists an SPJ query  $e'' \in \text{SPJ}(e')$  such that  $e \subseteq e''$  is in NP. This is done by guessing an SPJ query  $e''$  from  $\text{SPJ}(e')$ , and a mapping  $h : \text{Dom}(A_{q_{e''}}) \rightarrow \text{Dom}(A_{q_e})$ , where  $q_e(\bar{x})$  and  $q_{e''}(\bar{y})$  are the CQs obtained after converting  $e$  and  $e''$  into CQs by applying the translation given in Theorem 12.7. It is easy to see that such translation can be carried out in polynomial time. We finally verify that  $h$  is a homomorphism from  $(q_{e''}, \bar{y})$  to  $(q_e, \bar{x})$ . The correctness of this procedure is guaranteed by the Homomorphism Theorem (Theorem 15.4). Since both steps are feasible in polynomial time (recall that an SPJ query from  $\text{SPJ}(e)$  can be guessed in polynomial time), we conclude that checking whether  $e \subseteq e'$  is in NP, and the claim follows.  $\square$

Proposition 29.5 essentially tells us that **SPJU-Containment** remains NP-complete whenever the left-hand side query in the containment check does not use union. However, as already mentioned, the complexity increases when considering the problem in its general form without any assumptions.

### Theorem 29.6

SPJU-Containment is  $\Pi_2^P$ -complete.

*Proof.* Consider two SPJU queries  $e$  and  $e'$ . We proceed to show that checking whether  $e \not\subseteq e'$  is in  $\Sigma_2^P = \text{NP}^{\text{NP}}$ , which in turn implies that **SPJU-Containment** is in  $\Pi_2^P$ . By Proposition 29.4, it suffices to show that the problem of checking whether there exists an SPJ query  $\hat{e} \in \text{SPJ}(e)$  such that  $\hat{e} \not\subseteq e'$  is in  $\text{NP}^{\text{NP}}$ . This is done by simply guessing an SPJ query  $\hat{e}$  from  $\text{SPJ}(e)$ , and verifying that  $\hat{e} \not\subseteq e'$ . It is clear that the “guess” step can be performed in polynomial time. Concerning the “verify” step, by Proposition 29.5, it can be performed in constant time assuming that we have access to an oracle that can solve any problem in NP. In particular, the oracle takes as input the queries  $\hat{e}$  and  $e'$ , and does the following: if  $\hat{e} \subseteq e'$ , then return **false**; otherwise, return **true**. Therefore, checking whether  $e \not\subseteq e'$  is in  $\text{NP}^{\text{NP}}$ , as needed.

To prove the  $\Pi_2^P$ -hardness one can provide a polynomial-time reduction from the problem  $\forall\exists\text{QSAT}$  (Exercise 4.2).  $\square$

An immediate corollary of Theorem 29.6 is that the equivalence problem for SPJU queries, that is, given two SPJU queries  $e, e'$ , check whether  $e \equiv e'$ , is in  $\Pi_2^P$  since it boils down to two containment checks:  $e \subseteq e'$  and  $e' \subseteq e$ . The  $\Pi_2^P$ -hardness is shown via an easy reduction from **SPJU-Containment**. Given two SPJU queries  $e, e'$  of arity  $k$ , it holds that  $e \subseteq e'$  iff  $\pi_{(1,\dots,k)}(e \bowtie_{\theta} e') \equiv e$ , where  $\theta = (1 \dot{=} k + 1) \wedge (2 \dot{=} k + 2) \wedge \dots \wedge (k \dot{=} 2k)$ . Therefore:

#### Corollary 29.7

SPJU-Equivalence is  $\Pi_2^P$ -complete.

## Minimization

In Chapter 16, we studied the notion of minimization of CQs, which aims to provide equivalent CQs that are also minimal. More precisely, given a CQ  $q$  over a schema  $\mathbf{S}$ , a CQ  $q'$  over  $\mathbf{S}$  is a minimization of  $q$  if  $q \equiv q'$ , and for every CQ  $q''$  over  $\mathbf{S}$  with  $q' \equiv q''$  it holds that  $|A_{q'}| \leq |A_{q''}|$ . We have also seen how the minimization of a CQ, which is unique (up to variable renaming), can be computed by simply removing atoms from its body. We proceed to discuss how the ideas developed around CQ minimization can be extended to UCQs. We start by defining the notion of minimization for UCQs.

#### Definition 29.8: Minimization of UCQs

Consider a UCQ  $q = q_1 \cup \dots \cup q_n$  over a schema  $\mathbf{S}$ . A UCQ  $q' = q'_1 \cup \dots \cup q'_m$ , for  $m \leq n$ , over  $\mathbf{S}$  is a *minimization* of  $q$  if the following hold:

1.  $q \equiv q'$ ,
2. for every  $i \in [m]$ , and CQ  $p$  over  $\mathbf{S}$ , if  $q'_i \equiv p$  then  $|A_{q'_i}| \leq |A_p|$ , and

3. for every UCQ  $q'' = q_1'' \cup \dots \cup q_\ell''$ , if  $q' \equiv q''$  then  $m \leq \ell$ .

In simple words,  $q'$  is a minimization of  $q$  if it is equivalent to  $q$ , each CQ of  $q'$  has the smallest number of atoms among all the CQs that are equivalent to it, and  $q'$  has the smallest number of CQs among all the UCQs that are equivalent to it. We proceed to show that minimizations of a UCQ  $q$  can be found by simply removing atoms from the body of its CQs (i.e., by computing a core of its CQs), as well as removing CQs from it. Moreover, although  $q$  may have several minimizations, they are all the same (up to variable renaming).

#### *Minimization via Atom and CQ Removals*

We start by defining the notion of core for UCQs, which is a generalization of the notion of core for CQs given in Definition 16.2.

#### **Definition 29.9: Core of a UCQ**

Consider a UCQ  $q(\bar{x}) = q_1 \cup \dots \cup q_n$ . A UCQ  $q'(\bar{x}) = q'_1 \cup \dots \cup q'_m$  is a *core* of  $q$  if the following hold:

1.  $m \leq n$ , and there is a set of integers  $\{i_1, \dots, i_m\} \subseteq [n]$  such that, for every  $j \in [m]$ ,  $q'_j$  is a core of  $q_{i_j}$ ,
2. for every  $i \in [n]$ , there exists  $j \in [m]$  such that  $q_i \subseteq q'_j$ , and
3. for every  $i \in [m]$ , there is no  $j \in [m] - \{i\}$  such that  $q'_i \subseteq q'_j$ .

The first condition in Definition 29.9 states that  $q'$  can be obtained from  $q$  by eliminating some of its CQs, and then computing a core of the remaining CQs, the second condition ensures that  $q \equiv q'$ , and the third condition states that  $q'$  is minimal with respect to the number of CQs occurring in it. We show that the notion of core captures our intention of constructing a minimization.

#### **Proposition 29.10**

Every UCQ  $q$  has at least one core, and every core of  $q$  is a minimization of  $q$ .

*Proof.* We first show that we can always construct a core of  $q$ . Assuming that  $q = q_1 \cup \dots \cup q_n$ , we first replace  $q_i$  with a core of it, for each  $i \in [n]$ , which we know always exists by Proposition 16.4, and get a UCQ  $q' = q'_1 \cup \dots \cup q'_n$ . If  $q'$  is a core of itself, then the claim follows. Assume now that this is not the case. This means that condition (3) in the definition of core (Definition 29.9) is violated, which in turn implies that there exists  $i \in [n]$  such that  $q'_i \subseteq q'_j$  for some  $j \in [n] - \{i\}$ . If  $q''$  obtained from  $q'$  by removing the CQ  $q'_i$  is a core

**Algorithm 16** COMPUTECOREUCQ( $q$ )**Input:** A UCQ  $q(\bar{x}) = q_1 \cup \dots \cup q_n$ **Output:** A UCQ  $q^*(\bar{x})$  that is a core of  $q(\bar{x})$ 


---

```

1:  $Q := \emptyset$ 
2: for  $i = 1$  to  $n$  do
3:    $Q := Q \cup \{\text{COMPUTECORE}(q_i)\}$ 
4: while there are distinct CQs  $q', q'' \in Q$  such that  $q' \subseteq q''$  do
5:    $Q := Q - \{q'\}$ 
6: return  $q^* = \bigcup_{q' \in Q} q'$ 

```

---

of itself, then it is clear that  $q''$  is a core of  $q$ . Otherwise, we iteratively apply the above argument until we reach a core of  $q$ .

We now proceed to show that a core of  $q(\bar{x}) = q_1 \cup \dots \cup q_n$  is a minimization of it. Towards a contradiction, assume that  $q'(\bar{x}) = q'_1 \cup \dots \cup q'_m$  is a core of  $q$  but not a minimization of  $q$ . This implies that one of the following holds:

1. there exists  $i \in [n]$  and a CQ  $p$  such that  $q'_i \equiv p$  and  $|A_p| < |A_{q'_i}|$ , or
2. there exists a UCQ  $q'' = q''_1 \cup \dots \cup q''_\ell$  such that  $q' \equiv q''$  and  $\ell < m$ .

Assuming (1) holds, we have that  $q'_i$  is not a core of itself, and hence  $q'_i$  cannot be a core of a CQ in  $q$ . This contradicts the fact that  $q'$  is a core of  $q$ . Assume now (2). Since  $q' \subseteq q''$ , by Proposition 29.1 we get that, for every  $i \in [m]$ , there exists  $j \in [\ell]$  such that  $q'_i \subseteq q''_j$ . Then, from the fact that  $\ell < m$ , we obtain that there exist  $i, i' \in [m]$ , with  $i \neq i'$ , and  $j^* \in [\ell]$  such that  $q'_i \subseteq q''_{j^*}$  and  $q'_{i'} \subseteq q''_{j^*}$ . Now, since  $q'' \subseteq q'$ , by Proposition 29.1 we get that, for each  $j \in [\ell]$ , there exists  $i \in [m]$  such that  $q''_j \subseteq q'_i$ . Therefore, there exists  $i^* \in [m]$  such that  $q''_{j^*} \subseteq q'_{i^*}$ . This implies that  $q'_i \subseteq q'_{i^*}$  and  $q'_{i'} \subseteq q'_{i^*}$ . Consequently, there exist  $i \in [m]$  and  $j \in [m] - \{i\}$  such that  $q'_i \subseteq q'_j$ , which again contradicts our hypothesis that  $q'$  is a core of  $q$ , and the claim follows.  $\square$

By Proposition 29.10, computing a minimization of a UCQ boils down to computing a core of it. This can be done by applying the iterative procedure COMPUTECOREUCQ, given in Algorithm 16, which in turn relies on COMPUTECORE, given in Algorithm 4, that computes the core of a CQ. It is clear that, for a UCQ  $q$ , COMPUTECOREUCQ( $q$ ) terminates after finitely many steps. It is also easy to verify that COMPUTECOREUCQ is correct.

**Lemma 29.11.** *Given a UCQ  $q$ , COMPUTECOREUCQ( $q$ ) is a core of  $q$ .*

Let us clarify that COMPUTECOREUCQ is a nondeterministic algorithm since the procedure COMPUTECORE is nondeterministic. Moreover, there may be several CQs satisfying the condition of the while loop (in particular, there may be several CQs that must be removed from the set  $Q$ ), but we do not

specify how such a CQ is selected. Actually, the CQ  $q'$  of  $Q$  that is eventually removed from  $Q$  at step 5 is chosen nondeterministically. Therefore, the final result computed by the algorithm depends on the computation of COMPUTECORE at step 3, as well as how the CQs to be removed from  $Q$  are chosen at step 5. Consequently, different executions of COMPUTECORE( $q$ ) may compute cores of  $q$  that are syntactically different. However, as we discuss next, different minimizations of a UCQ  $q$  are isomorphic, which in turn implies, due to Proposition 29.10, that different cores of  $q$  are isomorphic.

#### Uniqueness of Minimization

We say that two UCQs  $q(\bar{x}) = q_1 \cup \dots \cup q_n$  and  $q'(\bar{x}') = q'_1 \cup \dots \cup q'_m$  are *isomorphic* if one can be turned into the other via renaming of variables, i.e., there is a bijection  $\delta : \{q_1, \dots, q_n\} \rightarrow \{q'_1, \dots, q'_m\}$ , which means that  $n = m$ , such that for every  $i \in [n]$  the CQs  $q_i$  and  $\delta(q_i)$  are isomorphic.

#### Proposition 29.12

Consider a UCQ  $q(\bar{x})$ , and let  $q'(\bar{x}')$  and  $q''(\bar{x}'')$  be minimizations of  $q$ . Then  $q'$  and  $q''$  are isomorphic.

*Proof.* Assume that  $q' = q'_1 \cup \dots \cup q'_n$  and  $q'' = q''_1 \cup \dots \cup q''_n$ . We need to show that there is a bijection  $\delta : \{q'_1, \dots, q'_n\} \rightarrow \{q''_1, \dots, q''_n\}$  such that for every  $i \in [n]$  the CQs  $q'_i$  and  $\delta(q'_i)$  are isomorphic. We first show an auxiliary lemma:

**Lemma 29.13.** *There is a bijection  $\tau : [n] \rightarrow [n]$  such that  $q'_i \equiv q''_{\tau(i)}$ , for each  $i \in [n]$ .*

*Proof.* To prove the claim, it suffices to show that, for each  $i \in [n]$ :

1. there exists  $j \in [n]$  such that  $q'_i \equiv q''_j$ , and
2. for each  $j, k \in [n]$ , if  $q'_i \equiv q''_j$  and  $q'_i \equiv q''_k$  then  $j = k$ .

We first show claim (1). Since  $q$  and  $q'$  are minimizations of  $q$ , we conclude that  $q \equiv q'$  and  $q \equiv q''$ , and thus,  $q' \equiv q''$ . By Proposition 29.1, we get that, for each  $i \in [n]$ , there exists  $j \in [n]$  such that  $q'_i \subseteq q''_j$ . But, again by Proposition 29.1, there exists  $k \in [n]$  such that  $q''_j \subseteq q'_k$ . Necessarily,  $q'_i = q'_k$ ; otherwise,  $q'$  is equivalent to the UCQ obtained from  $q'$  after eliminating  $q'_i$  (since  $q'_i \subseteq q'_k$ ), which contradicts the fact that  $q'$  is a minimization of  $q$ . We conclude then that  $q'_i \equiv q''_j$ .

We now prove claim (2). Towards a contradiction, assume that there exists  $i \in [n]$ , and distinct integers  $j, k \in [n]$  such that  $q'_i \equiv q''_j$  and  $q'_i \equiv q''_k$ . This implies that the UCQ obtained from  $q''$  after eliminating one of the CQs  $q''_j$  or  $q''_k$  (since  $q''_j \equiv q''_k$ ) is equivalent to  $q''$ , which contradicts the fact that  $q''$  is a minimization of  $q$ . This completes the proof of Lemma 29.13.  $\square$

Having the bijection  $\tau : [n] \rightarrow [n]$  provided by Lemma 29.13, we define the bijection  $\delta : \{q'_1, \dots, q'_n\} \rightarrow \{q''_1, \dots, q''_n\}$  as  $\delta(q'_i) = q''_{\tau(i)}$ , for each  $i \in [n]$ . It remains to argue that, for every  $i \in [n]$ , the CQs  $q'_i$  and  $\delta(q'_i)$  are isomorphic. Since  $q'$  and  $q''$  are minimizations of  $q$ , both  $q'_i$  and  $\delta(q'_i)$  must be cores of themselves. Hence,  $\delta(q'_i)$  is a minimization of  $q'_i$ , and thus, by Proposition 16.7, we get that  $q'_i$  and  $\delta(q'_i)$  are isomorphic.  $\square$

From Proposition 29.10, which tells us that a core of a UCQ  $q$  is a minimization of  $q$ , and Proposition 29.12, we get the following corollary:

**Corollary 29.14**

Consider a UCQ  $q$ , and let  $q'$  and  $q''$  be cores of  $q$ . It holds that  $q'$  and  $q''$  are isomorphic.

Recall that different executions of the nondeterministic procedure COMPUTECOREUCQ on some UCQ  $q$ , may compute cores of  $q$  that are syntactically different. However, Corollary 29.14 tells us that those cores differ only on the names of their variables. In other words, cores of  $q$  computed by different executions of COMPUTECOREUCQ( $q$ ) are the same up to variable renaming.





## Unions of Conjunctive Queries with Inequalities

It is not difficult to show that UCQ s, or even UCQ s with variable-constant equality that are equally expressive to  $\exists\text{FO}^+$  queries, are not powerful enough for expressing simple queries that involve negation such as the query

$$q = \exists x \exists y (\text{Edge}(x, y) \wedge \neg(x = y)),$$

which essentially asks whether a graph has an edge  $(v, u)$  that is not a loop, i.e.,  $v$  and  $u$  are different nodes. Observe that  $q$  is not preserved under homomorphisms: for  $D = \{R(a, b)\}$  and  $D' = \{R(c, c)\}$ , we have that  $D \rightarrow_\emptyset D'$ , but  $D \models q$  while  $D' \not\models q$ . On the other hand, we know by Proposition 28.10 that UCQ s (even with variable-constant equality) are preserved under homomorphisms, which immediately implies that  $q$  cannot be expressed as a UCQ (even with variable-constant equality). This raises the question whether there are languages obtained by adding a tamed negation to UCQ s without increasing the complexity of query evaluation, and, more importantly, without losing the decidability of containment and equivalence. The best known such language is that of UCQ s *with inequality*, that is, UCQ s that can also use expressions of the form  $\neg(v = u)$  as the query  $q$  given above. In this chapter, we introduce the language of UCQ s with inequality, and study the main computational problems: evaluation and containment.

### Conjunctive Queries with Inequality

Before we introduce and study UCQ s with inequality, it is important to study and understand first CQs with inequalities. Note that in the rest of the chapter we write  $v \neq u$  as an abbreviation for  $\neg(v = u)$ .

#### *Syntax of Conjunctive Queries with Inequality*

We start with the syntax of conjunctive queries with inequalities.

**Definition 30.1: Syntax of CQs with Inequality**

A *conjunctive query with inequality* ( $\text{CQ}^\neq$ ) over a schema  $\mathbf{S}$  is an FO query  $\varphi(\bar{x})$  over  $\mathbf{S}$  where  $\varphi$  is a formula of the form

$$\exists \bar{y} (R_1(\bar{u}_1) \wedge \cdots \wedge R_n(\bar{u}_n) \wedge v_1 \neq v'_1 \wedge \cdots \wedge v_m \neq v'_m)$$

for  $n \geq 1$  and  $m \geq 0$ , where

- for each  $i \in [n]$ ,  $R_i(\bar{u}_i)$  is a relational atom, and  $\bar{u}_i$  a tuple of constants and variables mentioned in  $\bar{x}$  and  $\bar{y}$ , and
- for each  $i \in [m]$ ,  $v_i$  is a variable mentioned in  $\bar{u}_k$  for some  $k \in [n]$ , and  $v'_i$  is a variable mentioned in  $\bar{u}_k$  for some  $k \in [n]$  or a constant.

Note that the second item of the definition requires that each variable that participates in at least one inequality appears also in at least one relational atom. This is a common assumption in CQs with inequality that, as we discuss below, allows us to show that homomorphisms provide an alternative way to describe the evaluation of CQs with inequality in the same way as for CQs.

It is common to represent CQs with inequality via a rule-like syntax. In particular, the  $\text{CQ}^\neq$   $\varphi(\bar{x})$  given in Definition 12.1 can be written as the *rule*

$$\text{Answer}(\bar{x}) \text{ :- } R_1(\bar{u}_1), \dots, R_n(\bar{u}_n), v_1 \neq v'_1, \dots, v_m \neq v'_m,$$

where  $\text{Answer}$  is a relation name not in  $\mathbf{S}$ , and its arity (under the singleton schema  $\{\text{Answer}\}$ ) is equal to the arity of  $q$ . The relational atom  $\text{Answer}(\bar{x})$  that appears on the left of the  $\text{:-}$  symbol is the *head* of the rule, while the expression that appears on the right of the  $\text{:-}$  symbol is the *body* of the rule. In general, we use the rule-like syntax. Nevertheless, for convenience, we will freely interpret a CQ with inequality as an FO query or as a rule.

*Semantics of Conjunctive Queries with Inequality*

Since a  $\text{CQ}^\neq$  is an FO query, the definition of its output on a database can be inherited from Definition 3.6. More precisely, given a database  $D$  of a schema  $\mathbf{S}$ , and a  $k$ -ary  $\text{CQ}^\neq$   $q = \varphi(\bar{x})$  over  $\mathbf{S}$ , where  $k \geq 0$ , the *output* of  $q$  on  $D$  is

$$q(D) = \{\bar{a} \in \text{Dom}(D)^k \mid D \models \varphi(\bar{a})\}.$$

Recall that every variable that occurs in an inequality of  $q$  it also occurs in a relational atom of  $q$ . For this reason, the output of  $q$  only consists of tuples of constants from  $\text{Dom}(D)$ . It is easy to verify that if we drop this condition, then the output may mention constants that occur in the query but not in the database. Consider, for example, the FO query  $q = \varphi(y)$  with

$$\varphi = \exists x (R(x) \wedge x \neq y \wedge a \neq b),$$

where  $x, y$  are variables and  $a, b$  are constants, which is not a CQ with inequality since  $y$  does not occur in a relational atom. Clearly, for  $D = \{R(a)\}$ ,  $q(D) = \{(b)\}$ , while the constant  $b$  does not belong to  $\text{Dom}(D)$ .

As for plain CQs, there is a more intuitive (and equivalent) way of defining the semantics of CQs with inequality when they are viewed as rules. The body of a  $\text{CQ}^\neq q$  of the form  $\text{Answer}(\bar{x}) :- \text{body}$  can be seen as a pattern that must be matched with the database  $D$  via an assignment  $\eta$  that maps the variables in  $q$  to  $\text{Dom}(D)$ . For each such assignment  $\eta$ , if  $\eta$  applied to this pattern produces only facts of  $D$ , and at the same time respects all the inequalities, it means that the pattern matches with  $D$  via  $\eta$ , and the tuple  $\eta(\bar{x})$  is an output of  $q$  on  $D$ . We proceed to formalize this informal description.

Consider a database  $D$  and a  $\text{CQ}^\neq q$  of the form

$$\text{Answer}(\bar{x}) :- R_1(\bar{u}_1), \dots, R_n(\bar{u}_n), v_1 \neq v'_1, \dots, v_m \neq v'_m.$$

An *assignment* for  $q$  over  $D$  is a function  $\eta$  from the set of variables in  $q$  to  $\text{Dom}(D)$ . We say that  $\eta$  is *consistent* with  $D$  if

$$R_i(\eta(\bar{u}_i)) \in D \quad \text{and} \quad \eta(v_j) \neq \eta(v'_j)$$

for each  $i \in [n]$  and  $j \in [m]$ , where the fact  $R_i(\eta(\bar{u}_i))$  is obtained by replacing each variable  $x$  in  $\bar{u}_i$  with  $\eta(x)$ , and leaving the constants in  $\bar{u}_i$  untouched. The consistency of  $\eta$  with  $D$  essentially means that the body of  $q$  matches with  $D$  via  $\eta$ . We can now define what is the output of a  $\text{CQ}^\neq$  on a database.

**Definition 30.2: Evaluation of CQs with Inequality**

Given a database  $D$  of a schema  $\mathbf{S}$ , and a  $\text{CQ}^\neq q(\bar{x})$  over  $\mathbf{S}$ , the *output* of  $q$  on  $D$  is defined as the set of tuples

$$q(D) = \{\eta(\bar{x}) \mid \eta \text{ is an assignment for } q \text{ over } D \text{ consistent with } D\}.$$

It is an easy exercise to show that the semantics of  $\text{CQ}^\neq$  inherited from the semantics of FO queries in Definition 3.6 and the semantics of CQs given in Definition 30.2 are equivalent, i.e., for a CQ  $q = \varphi(\bar{x})$  and a database  $D$ ,

$$\{\bar{a} \in \text{Dom}(D)^k \mid D \models \varphi(\bar{a})\} = \{\eta(\bar{x}) \mid \eta \text{ is an assignment for } q \text{ over } D \text{ consistent with } D\}.$$

*Evaluation and Homomorphisms*

We proceed to discuss how homomorphisms emerge in the context of CQs with inequality. In particular, we show that they provide an alternative way to describe the evaluation of CQs with inequality. Given a  $\text{CQ}^\neq q$  of the form

$$\text{Answer}(\bar{x}) :- R_1(\bar{u}_1), \dots, R_n(\bar{u}_n), v_1 \neq v'_1, \dots, v_m \neq v'_m,$$

we define the sets of relational atoms and inequalities

$$A_q^+ = \{R_1(\bar{u}_1), \dots, R_n(\bar{u}_n)\} \quad \text{and} \quad A_q^- = \{v_1 \neq v'_1, \dots, v_m \neq v'_m\},$$

respectively. As usual,  $\text{Dom}(A_q^+)$  collects all the variables and constants occurring in the relational atoms of  $A_q^+$ . We further write  $\text{Dom}(A_q^-)$  for the set of constants occurring in  $A_q^-$ . Notice that there may be constants in  $\text{Dom}(A_q^-)$  that do not occur in  $\text{Dom}(A_q^+)$  since, according to Definition 30.1, only the variables (not the constants) that occur in an inequality must also occur in a relational atom. Having the above sets in place, we can naturally talk about homomorphisms from CQs with inequality to databases.

**Definition 30.3: Homomorphisms from  $\text{CQ}^\neq$  to Databases**

Consider a  $\text{CQ}^\neq q(\bar{x})$  over a schema  $\mathbf{S}$ , and a database  $D$  of  $\mathbf{S}$ .

- A *homomorphism from  $q$  to  $D$*  is a function  $h : \text{Dom}(A_q^+) \cup \text{Dom}(A_q^-) \rightarrow \text{Dom}(D) \cup \text{Dom}(A_q^-)$  that is a homomorphism from  $A_q^+$  to  $D$ , is the identity on  $\text{Dom}(A_q^-)$ , and  $h(v) \neq h(u)$  for every  $v \neq u \in A_q^-$ . We write  $q \rightarrow D$  if such a homomorphism exists.
- A *homomorphism from  $(q, \bar{x})$  to  $(D, \bar{a})$*  is a homomorphism  $h$  from  $q$  to  $D$  such that  $h(\bar{x}) = \bar{a}$ . We write  $(q, \bar{x}) \rightarrow (D, \bar{a})$  if such a homomorphism exists.

To define the output of a  $\text{CQ}^\neq q(\bar{x})$  on a database  $D$  (see Definition 30.2), we used the notion of assignment for  $q$  over  $D$ , which is a function from the set of variables in  $q$  to  $\text{Dom}(D)$ . The output of  $q$  on  $D$  consists of all the tuples  $\eta(\bar{x})$ , where  $\eta$  is an assignment for  $q$  over  $D$  that is consistent with  $D$ , i.e.,  $R_i(\eta(\bar{u}_i)) \in D$  and  $\eta(v_j) \neq \eta(v'_j)$  for each  $i \in [n]$  and  $j \in [m]$ . Since  $R_i(\eta(\bar{u}_i))$  is the fact obtained after replacing each variable  $x$  in  $\bar{u}_i$  with  $\eta(x)$  and leaving the constants in  $\bar{u}_i$  untouched, for  $i \in [n]$ , such an assignment  $\eta$  corresponds to a function  $h : \text{Dom}(A_q^+) \cup \text{Dom}(A_q^-) \rightarrow \text{Dom}(D) \cup \text{Dom}(A_q^-)$ , which is the identity on the set of constants occurring in  $q$ , such that  $R(h(\bar{u}_i)) = R(\eta(\bar{u}_i))$ . But, of course, this is the same as saying that  $h$  is a homomorphism from  $q$  to  $D$ . Therefore,  $q(D)$  is the set of all tuples  $h(\bar{x})$ , where  $h$  is a homomorphism from  $q$  to  $D$ , i.e., the set of all tuples  $\bar{a}$  over  $\text{Dom}(D)$  with  $(q, \bar{x}) \rightarrow (D, \bar{a})$ . This leads to an alternative characterization of  $\text{CQ}^\neq$  evaluation.

**Theorem 30.4**

For a database  $D$  of a schema  $\mathbf{S}$ , and a  $\text{CQ}^\neq q(\bar{x})$  of arity  $k \geq 0$  over  $\mathbf{S}$ ,

$$q(D) = \{\bar{a} \in \text{Dom}(D)^k \mid (q, \bar{x}) \rightarrow (D, \bar{a})\}.$$

Unlike plain CQs, CQs with inequality are not preserved under homomorphisms. This has been already illustrated at the beginning of the chapter via

the  $\text{CQ}^\neq$  Answer  $\text{:- Edge}(x, y), x \neq y$ , which asks for the existence of an edge in a graph that is not a loop. On the other hand, we can easily show that CQs with inequality remain monotone, i.e., given a  $\text{CQ}^\neq$   $q$  over a schema  $\mathbf{S}$ , and two databases  $D$  and  $D'$  of  $\mathbf{S}$ , if  $D \subseteq D'$  then  $q(D) \subseteq q(D')$ .

**Proposition 30.5**

Every  $\text{CQ}^\neq$  is monotone.

*Proof.* Let  $q(\bar{x})$  be a  $\text{CQ}^\neq$  over a schema  $\mathbf{S}$ . Consider the databases  $D, D'$  of  $\mathbf{S}$  such that  $D \subseteq D'$ , and assume that  $\bar{a} \in q(D)$ . By Theorem 30.4, we get that  $(q, \bar{x}) \rightarrow (D, \bar{a})$ . Since  $D \subseteq D'$ , we immediately get that  $(q, \bar{x}) \rightarrow (D', \bar{a})$ , and thus, again by Theorem 30.4,  $\bar{a} \in q(D')$ , as needed.  $\square$

The fact that CQs with inequality are monotone allows us to clarify the expressiveness boundaries of  $\text{CQ}^\neq$ . In particular, we can show that the negation allowed in  $\text{CQ}^\neq$  is indeed quite restricted. Note that the following arguments have been already used in Chapter 13 to show that CQs cannot express negative relational atoms and difference.

**$\text{CQ}^\neq$  cannot express negative relational atoms.** The reason is because such queries are not monotone. Consider, for example, the FO query

$$q = \neg P(a),$$

where  $a$  is a constant. If we take  $D = \emptyset$  and  $D' = \{P(a)\}$ , then  $D \subseteq D'$  but  $D \models q$  while  $D' \not\models q$ .

**$\text{CQ}^\neq$  cannot express difference.** This is because difference is not monotone. Consider, for example, the FO query

$$q = \exists x(P(x) \wedge \neg Q(x)).$$

For  $D = \{P(a)\} \subseteq D' = \{P(a), Q(a)\}$ , we have that  $D \models q$  while  $D' \not\models q$ .

*Query Evaluation*

We now proceed to study the complexity of  $\text{CQ}^\neq$ -Evaluation, that is, the problem of checking whether  $\bar{a} \in q(D)$  for a  $\text{CQ}^\neq$  query  $q$ , a database  $D$ , and a tuple  $\bar{a}$  over  $\text{Dom}(D)$ . We actually show that it has the same (combined) complexity as CQ-Equivalence; concerning the data complexity, it is clear that  $\text{CQ}^\neq$ -Evaluation is in DLOGSPACE due to Theorem 7.3.

**Theorem 30.6**

$\text{CQ}^\neq$ -Evaluation is NP-complete.

*Proof.* It is clear that the NP-hardness is inherited from **CQ-Evaluation**, which we know is NP-hard (Theorem 14.1). We proceed to show the upper bound. Consider a  $\text{CQ}^\neq q(\bar{x})$ , a database  $D$ , and a tuple  $\bar{a} \in \text{Dom}(D)$ . By Theorem 30.4  $\bar{a} \in q(D)$  if and only if  $(q, \bar{x}) \rightarrow (D, \bar{a})$ . Therefore, we need to show that checking whether there exists a homomorphism from  $(q, \bar{x})$  to  $(D, \bar{a})$  is in NP. This is done by guessing a function  $h : \text{Dom}(A_q^+) \cup \text{Dom}(A_q^-) \rightarrow \text{Dom}(D) \cup \text{Dom}(A_q^-)$  that is the identity on constants, and then verifying that (i)  $h$  is a homomorphism from  $(A_q^+, \bar{x})$  to  $(D, \bar{a})$ , and (ii) for every  $v \neq u \in A_q^-$ ,  $h(v) \neq h(u)$ . Since both steps are feasible in polynomial time, we conclude that checking whether  $(q, \bar{x}) \rightarrow (D, \bar{a})$  is in NP, and the claim follows.  $\square$

### Containment

We now focus on **CQ $^\neq$ -Containment**, the problem of deciding whether a  $\text{CQ}^\neq$  is contained in another  $\text{CQ}^\neq$ . We show that this problem is decidable, but its complexity is higher than **CQ-Containment**, that is,  $\Pi_2^P$ -complete.

Recall that for CQs (without inequality) the building block underlying the procedure for checking containment is the Homomorphism Theorem (Theorem 15.4) that provides a useful characterization of containment in terms of homomorphisms: given two CQs  $q(\bar{x})$  and  $q'(\bar{x})$ ,  $q \subseteq q'$  iff  $(q', \bar{x}') \rightarrow (q, \bar{x})$ . It is not difficult to see, however, that this is no longer true once inequalities are allowed. This can be easily shown via a simple example.

#### Example 30.7: The Homomorphism Theorem Fails for $\text{CQ}^\neq$

Consider the (Boolean) queries

$$\begin{aligned} q = \text{Answer} & :- \text{Edge}(x, y) \\ q' = \text{Answer} & :- \text{Edge}(x', y'), x' \neq y'. \end{aligned}$$

It is clear that  $q' \rightarrow q$ , that is, there is a homomorphism  $h$  from  $A_{q'}^+$  to  $A_q^+$  with  $h(x') \neq h(y')$ , but  $q \not\subseteq q'$ : for  $D = \{R(a, a)\}$  we have that  $D \models q$  but  $D \not\models q'$ . This should not be surprising since the existence of an edge in a graph does immediately imply that there is a non-loop edge.

As the above example illustrates, the key reason why a result similar to the Homomorphism Theorem fails for CQs with inequality is because homomorphisms, as defined for these queries, do not compose. In other words, the fact that  $q' \rightarrow q$  and  $q \rightarrow D$ , for  $q, q'$  CQs with inequalities, does not allow us to conclude that  $q' \rightarrow D$  by simply composing homomorphisms. Observe that after composing  $h'$ , where  $h'(x') = x$  and  $h'(y') = y$ , which witnesses the fact that  $q' \rightarrow q$ , with  $h$ , where  $h(x) = h(y) = a$ , which witnesses that  $q \rightarrow D$ , we get the function  $g$  with  $g(x') = g(y') = a$  that is not a homomorphism from  $q'$  to  $D$  since the inequality is not preserved. The above discussion indicates that we need a version of the Homomorphism Theorem that somehow ensures

the following: no matter how the homomorphism  $h$ , which witnesses the fact that  $q \rightarrow D$ , looks like,  $h \circ h'$  is a homomorphism from  $q'$  to  $D$ .

For a CQ (without inequality)  $q(\bar{x})$ , and a CQ $^\neq$   $q'(\bar{x}')$ , we write  $(q', \bar{x}') \rightarrow (q, \bar{x})$  for the fact that there exists a homomorphism from  $(A_{q'}, \bar{x}')$  to  $(A_q, \bar{x})$ , which in turn is defined in exactly the same way as the notion of homomorphism from a CQ $^\neq$  to a database (see Definition 30.3).

Given a CQ $^\neq$   $q(\bar{x})$  of the form

$$\text{Answer}(\bar{x}) \text{ :- } R_1(\bar{u}_1), \dots, R_n(\bar{u}_n), v_1 \neq v'_1, \dots, v_m \neq v'_m,$$

let  $H_q$  be the set of all functions  $h : \text{Dom}(A_q^+) \cup \text{Dom}(A_q^-) \rightarrow \text{Dom}(A_q^+) \cup \text{Dom}(A_q^-)$  such that (i)  $h$  is the identity on constants, and (ii) for every  $v \neq u \in A_q^-$ ,  $h(v) \neq h(u)$ . Intuitively,  $H_q$  collects all the possible ways that  $q$  can be mapped into a homomorphic image of itself. For a function  $h \in H_q$ , we write  $h(q^+)$  for the CQ (without inequality)

$$\text{Answer}(h(\bar{x})) \text{ :- } R_1(h(\bar{u}_1)), \dots, R_n(h(\bar{u}_n)),$$

i.e., the CQ obtained from  $q$  by eliminating the inequalities and then replacing each term  $u$  with the term  $h(u)$ . We are now ready to establish the version of the Homomorphism Theorem for CQs with inequality.

### Theorem 30.8

Let  $q(\bar{x})$  and  $q'(\bar{x}')$  be CQ $^\neq$ s. The following are equivalent:

1.  $q \subseteq q'$ .
2.  $(q', \bar{x}') \rightarrow (h(q^+), h(\bar{x}))$ , for each  $h \in H_q$ .

*Proof.* We first establish that (1) implies (2). Consider an arbitrary function  $h \in H_q$ ; for brevity, let  $p = h(q^+)$ . Note that  $p$  is a CQ without inequalities. By definition of  $H_q$ , we have that  $h$  is a homomorphism from  $(q, \bar{x})$  to  $(p, h(\bar{x}))$ . Since  $\mathbb{G}_{A_p}$  is a bijective homomorphism from  $(p, h(\bar{x}))$  to  $(\mathbb{G}_{A_p}(A_p), \mathbb{G}_{A_p}(h(\bar{x})))$ , we conclude that  $(\mathbb{G}_{A_p} \cup \mu) \circ h$ , where  $\mu$  is the identity on the set of constants  $\text{Dom}(A_q^-) - \text{Dom}(A_p)$ , is a homomorphism from  $(q, \bar{x})$  to  $(\mathbb{G}_{A_p}(A_p), \mathbb{G}_{A_p}(h(\bar{x})))$ . By Theorem 30.4,  $\mathbb{G}_{A_p}(h(\bar{x})) \in q(\mathbb{G}_{A_p}(A_p))$ . Since, by hypothesis,  $q \subseteq q'$ , we conclude that  $\mathbb{G}_{A_p}(h(\bar{x})) \in q'(\mathbb{G}_{A_p}(A_p))$ . By applying again Theorem 30.4, we get that there exists a homomorphism  $g$  from  $(q', \bar{x}')$  to  $(\mathbb{G}_{A_p}(A_p), \mathbb{G}_{A_p}(h(\bar{x})))$ . Since  $\mathbb{G}_{A_p}$  is a bijection,  $(\mathbb{G}_{A_p}^{-1} \cup \mu') \circ g$ , where  $\mu'$  is the identity on the set of constants  $\text{Dom}(A_q^-) - \text{Dom}(A_p)$ , is a homomorphism from  $(q', \bar{x}')$  to  $(p, h(\bar{x}))$ .

We now proceed to show that (2) implies (1). Given a database  $D$ , assume that  $\bar{a} \in q(D)$ . By Theorem 30.4, there exists a homomorphism  $g$  from  $(q, \bar{x})$  to  $(D, \bar{a})$ . Let  $h$  be a function from  $\text{Dom}(A_q^+) \cup \text{Dom}(A_q^-)$  to itself such that

- $h$  is the identity on constants,



- for every two variables  $x, y \in \text{Dom}(A_q^+)$ ,  $h(x) = h(y)$  iff  $g(x) = g(y)$ , and
- $h(x)$  is a variable, for every variable  $x \in \text{Dom}(A_q^+)$ .

In simple words,  $h$  unifies the variables in  $q^+$  that are mapped by  $g$  to the same constant of  $\text{Dom}(D)$ . It is easy to verify that such a function always exists and belongs to  $H_q$ . It is also clear that  $(h(q^+), h(\bar{x})) \rightarrow (D, \bar{a})$  is witnessed via a bijective homomorphism  $g'$ . Since  $h \in H_q$ , by hypothesis, there exists a homomorphism  $h'$  from  $(q', \bar{x}')$  to  $(h(q^+), h(\bar{x}))$ . Since  $g'$  is a bijection, we get that  $(g' \cup \mu) \circ h'$ , where  $\mu$  is the identity on the set of constants  $\text{Dom}(A_q^-) - \text{Dom}(A_{h(q^+)})$ , is a homomorphism from  $(q', \bar{x}')$  to  $(D, \bar{a})$ . By Theorem 30.4, we get that  $\bar{a} \in q'(D)$ , and the claim follows.  $\square$

The next example illustrates how Theorem 30.8 is used in order to confirm what has been discussed in Example 30.7.

**Example 30.9: Application of Theorem 30.8**

Consider again the CQs  $q$  and  $q'$  given in Example 30.7 and recall that  $q \not\subseteq q'$ . This is confirmed by Theorem 30.8 since, for the function  $h \in H_q$  with  $h(x) = h(y) = x$ , we can conclude that there is no homomorphism from  $q'$  to  $h(q)$ . Indeed, the only way to map  $q'$  to  $h(q)$  is via the function  $g$  with  $g(x') = g(y') = x$ , which is not a homomorphism from  $q'$  to  $h(q)$ .

From Theorem 30.8, we immediately get a procedure for  $\text{CQ}^\neq$ -Containment: given two CQs with inequality  $q(\bar{x})$  and  $q'(\bar{x}')$ , return **true** if  $(q', \bar{x}) \rightarrow (h(q^+), h(\bar{x}))$ , for every function  $h \in H_q$ ; otherwise, return **false**. However, this naive approach only shows that  $\text{CQ}^\neq$ -Containment is in EXPTIME since  $H_q$  consists, in general, of exponentially many functions, i.e., we need to perform exponentially many homomorphism checks, while each one takes exponential time. By providing a more clever procedure, we can establish the following.

**Theorem 30.10**

$\text{CQ}^\neq$ -Containment is  $\Pi_2^P$ -complete.

*Proof.* Consider two  $\text{CQ}^\neq$   $q(\bar{x})$  and  $q'(\bar{x}')$ . We proceed to show that checking whether  $q \not\subseteq q'$  is in  $\Sigma_2^P = \text{NP}^{\text{NP}}$ , which in turn implies that  $\text{CQ}^\neq$ -Containment is in  $\Pi_2^P$ . By Theorem 30.8, it suffices to show that the problem of checking whether there exists a function  $h \in H_q$  such that there is no homomorphism from  $(q', \bar{x}')$  to  $(h(q^+), h(\bar{x}))$  is in  $\text{NP}^{\text{NP}}$ . This can be done by simply guessing a function  $h$  from  $H_q$ , and the verifying that indeed there is no homomorphism from  $(q', \bar{x}')$  to  $(h(q^+), h(\bar{x}))$ . It is clear that the “guess” step is feasible in polynomial time. Concerning the “verify” step, we first observe the following, which can be easily shown via a simple guess-and-check algorithm:

**Lemma 30.11.** *Given a  $CQ^\neq$   $q_1(\bar{x}_1)$  and a  $CQ$  (without inequality)  $q_2(\bar{x}_2)$ , the problem of deciding whether  $(q_1, \bar{x}_1) \rightarrow (q_2, \bar{x}_2)$  is in NP.*

By Lemma 30.11, we conclude that the “verify” step can be performed in constant time assuming we have access to an oracle that can solve any problem in NP. In particular, the oracle takes as input the queries  $q'$  and  $h(q^+)$ , and does the following: if  $(q', \bar{x}') \rightarrow (h(q^+), h(\bar{x}))$ , then return **false**; otherwise, return **true**. Therefore, checking whether  $q \not\subseteq q'$  is in  $NP^{NP}$ , as needed.

For the  $\Pi_2^P$ -hardness of  $CQ^\neq$ -Containment see Exercise 4.4.  $\square$

With Theorem 30.10 in place, we can easily pinpoint the complexity of the equivalence problem for  $CQ^\neq$ : given two  $CQ^\neq$   $q$  and  $q'$ , check whether  $q \equiv q'$ , i.e., whether  $q(D) = q'(D)$  for every database  $D$ . We show that:

**Theorem 30.12**

$CQ^\neq$ -Equivalence is  $\Pi_2^P$ -complete.

*Proof.* For the upper bound, it suffices to observe that  $q \equiv q'$  iff  $q \subseteq q'$  and  $q' \subseteq q$ . This implies that  $CQ^\neq$ -Equivalence is in  $\Pi_2^P$  since, by Theorem 30.10, the problem of deciding whether  $q \subseteq q'$  and  $q' \subseteq q$  is in  $\Pi_2^P$ .

For the lower bound, we provide an easy reduction from  $CQ^\neq$ -Containment that is  $\Pi_2^P$ -hard (Theorem 30.10); in fact, this holds even for Boolean queries. Given two Boolean  $CQ^\neq$   $q, q'$ , we can easily construct a  $CQ^\neq$   $q_\cap$  that computes the intersection of  $q$  and  $q'$ , i.e., for every database  $D$ ,  $q(D) \cap q'(D) = q_\cap(D)$ , by merging the bodies of  $q$  and  $q'$ ; a similar construction has been already used in the proof of Theorem 15.8 that deals with  $CQ$ -Equivalence. Since  $q \subseteq q'$  iff  $q \equiv q_\cap$ , we get that  $CQ^\neq$ -Equivalence is  $\Pi_2^P$ -hard, and the claim follows.  $\square$

## Adding Union to $CQ^\neq$

We now proceed to add the union operator to CQs with inequality.

**Definition 30.13: Union of Conjunctive Queries with Inequality**

A *union of conjunctive queries with inequality* ( $UCQ^\neq$ ) over a schema  $\mathbf{S}$  is an FO query  $\varphi(\bar{x})$  over  $\mathbf{S}$ , where  $\varphi$  is of the form

$$\varphi_1 \vee \cdots \vee \varphi_n$$

for  $n \geq 1$ ,  $FV(\varphi) = FV(\varphi_i)$ , and  $\varphi_i(\bar{x})$  is a  $CQ^\neq$  for every  $i \in [n]$ .

As for  $UCQ$  s without inequality, for notational convenience, we denote a  $UCQ^\neq$   $q = \varphi(\bar{x})$  with  $\varphi = \varphi_1 \vee \cdots \vee \varphi_n$  as  $q_1 \cup \cdots \cup q_n$ , where  $q_i$  is the  $CQ^\neq$   $\varphi_i(\bar{x})$ , for each  $i \in [n]$ . It is easy to verify that, given a database  $D$  of a schema  $\mathbf{S}$ , and a  $UCQ^\neq$   $q = q_1 \cup \cdots \cup q_n$  over  $\mathbf{S}$ , we have that  $q(D) = q_1(D) \cup \cdots \cup q_n(D)$ .

**Example 30.14: Union of Conjunctive Queries with Inequality**

Consider the relational schema from Example 3.2:

Person [ pid, pname, cid ]  
 Profession [ pid, pname ]  
 City [ cid, cname, country ]

The UCQ  $\varphi(y)$ , where  $\varphi = \varphi_1 \vee \varphi_2$  with

$$\varphi_1 = \exists x \exists z (\text{Person}(x, y, z) \wedge \text{Profession}(x, \text{'computer scientist'}) \wedge \text{City}(z, w, \text{'Greece'}) \wedge w \neq \text{'Athens'}).$$

and

$$\varphi_2 = \exists x \exists z (\text{Person}(x, y, z) \wedge \text{Profession}(x, \text{'computer scientist'}) \wedge \text{City}(z, w, \text{'Chile'}) \wedge w \neq \text{'Santiago'}).$$

can be used to retrieve the list of names of computer scientists that were born in Greece or Chile, but not in the capital city of the country.

*Query Evaluation*

We proceed with **UCQ<sup>≠</sup>-Evaluation**, the problem of checking whether  $\bar{a} \in q(D)$  for a UCQ<sup>≠</sup>  $q$ , a database  $D$ , and a tuple  $\bar{a}$  over  $\text{Dom}(D)$ . We show that it has the same (combined) complexity as **CQ<sup>≠</sup>-Evaluation**; concerning the data complexity, both problems are in DLOGSPACE due to Theorem 7.3

**Theorem 30.15**

UCQ<sup>≠</sup>-Evaluation is NP-complete.

*Proof.* It is clear that the NP-hardness is inherited from **CQ-Evaluation**, which we know is NP-hard (Theorem 14.1). We proceed to show the upper bound. Consider a UCQ<sup>≠</sup>  $q(\bar{x})$  of the form  $q_1 \cup \dots \cup q_n$ , a database  $D$ , and a tuple  $\bar{a} \in \text{Dom}(D)$ . By Theorem 30.4, for  $i \in [n]$ ,  $\bar{a} \in q_i(D)$  if and only if  $(q_i, \bar{x}) \rightarrow (D, \bar{a})$ . Thus, we need to show that checking whether there exists an integer  $i \in [n]$  and a homomorphism from  $(q_i, \bar{x})$  to  $(D, \bar{a})$  is in NP. This is done by guessing an integer  $i \in [n]$  and a function  $h : \text{Dom}(A_{q_i}^+) \cup \text{Dom}(A_{q_i}^-) \rightarrow \text{Dom}(D) \cup \text{Dom}(A_{q_i}^-)$ , and then verifying that  $h$  is indeed a homomorphism from  $(A_{q_i}, \bar{x})$  to  $(D, \bar{a})$ , i.e., is a homomorphism from  $A_{q_i}^+$  to  $D$ , is the identity on  $\text{Dom}(A_{q_i}^-)$ , for every  $v \neq u \in A_{q_i}^-$ ,  $h(v) \neq h(u)$ , and  $h(\bar{x}) = \bar{a}$ . Since all the above steps are feasible in polynomial time, the claim follows.  $\square$

### Containment

We finally concentrate on  $\text{UCQ}^\neq\text{-Containment}$ , the problem of deciding whether a  $\text{UCQ}^\neq q$  is contained in a  $\text{UCQ}^\neq q'$ , that is, whether  $q(D) \subseteq q'(D)$  for every database  $D$ . We show that it has the same complexity as  $\text{CQ}^\neq\text{-Containment}$ . To this end, we first observe that Proposition 29.1 for UCQ s (without inequality) holds also for UCQ s with inequality. More precisely, the proof of Proposition 29.1 relies on the fact that, for a database  $D$ , and a UCQ  $q = q_1 \cup \dots \cup q_n$ ,  $q(D) = \bigcup_{i \in [n]} q_i(D)$ . Since the latter property holds even for UCQ s with inequality, the same proof shows the following:

#### Proposition 30.16

Consider two  $\text{UCQ}^\neq q = q_1 \cup \dots \cup q_n$  and  $q' = q'_1 \cup \dots \cup q'_m$ . The following are equivalent:

1.  $q \subseteq q'$ .
2. For every  $i \in [n]$ , there exists  $j \in [m]$  such that  $q_i \subseteq q'_j$ .

We are now ready to pinpoint the complexity of  $\text{UCQ}^\neq\text{-Containment}$ .

#### Theorem 30.17

$\text{UCQ}^\neq\text{-Containment}$  is  $\Pi_2^P$ -complete.

*Proof.* It is clear that the  $\Pi_2^P$ -hardness is inherited from  $\text{CQ}^\neq\text{-Containment}$ , which we know is  $\Pi_2^P$ -hard (Theorem 30.10). Consider now two  $\text{UCQ}^\neq q = q_1 \cup \dots \cup q_n$  and  $q' = q'_1 \cup \dots \cup q'_m$ . We proceed to show that checking whether  $q \not\subseteq q'$  is in  $\Sigma_2^P = \text{NP}^{\text{NP}}$ , which in turn implies that  $\text{UCQ}^\neq\text{-Containment}$  is in  $\Pi_2^P$ . We assume that, for  $i \in [n]$  and  $j \in [m]$ , the output tuple of  $q_i$  and  $q'_j$  is  $\bar{x}$  and  $\bar{x}'$ , respectively. By Theorem 30.8 and Proposition 30.16 to check whether  $q \not\subseteq q'$  it suffices to do the following:

- guess an  $i \in [n]$ , and, for every  $j \in [m]$ , guess a function  $h_j \in H_{q_i}$ , and
- for each  $j \in [m]$ , verify that there is no homomorphism from  $(q'_j, \bar{x}')$  to  $(h_j(q_i^+), h_j(\bar{x}))$ .

It is clear that the “guess” step is feasible in polynomial time. Moreover, by Lemma 30.11, the “verify” step can be performed in polynomial time assuming access to an oracle that can solve any problem in NP. In particular, for each  $j \in [m]$ , the oracle is called with input the queries  $q'_j$  and  $h_j(q_i^+)$ , and does the following: if  $(q'_j, \bar{x}') \rightarrow (h_j(q_i^+), h_j(\bar{x}))$ , then return **false**; otherwise, return **true**. Therefore, deciding whether  $q \not\subseteq q'$  is in  $\text{NP}^{\text{NP}}$ , as needed.  $\square$

An immediate corollary of Theorem 30.17 is that the equivalence problem for UCQs with inequality, that is, given two UCQ $^\neq$   $q, q'$ , check whether  $q \equiv q'$ , is in  $\Pi_2^p$  since it boils down to two containment checks:  $q \subseteq q'$  and  $q' \subseteq q$ . The  $\Pi_2^p$ -hardness is inherited from CQ $^\neq$ -Equivalence (Theorem 30.12).

**Corollary 30.18**

UCQ $^\neq$ -Equivalence is  $\Pi_2^p$ -complete.

## The Limits of First-Order Queries: Recursion

We have seen that the language of UCQs has the same expressive power as the language of  $\text{FO}^{\text{rel}}[\wedge, \vee, \exists]$  queries (Theorem 28.3), and that adding equational atoms  $x = a$  between variables and constants leads to the strictly more expressive language of  $\exists\text{FO}^+$  queries (Theorem 28.5). If we further add negation to  $\exists\text{FO}^+$  queries, we then get the full power of FO queries. Universal quantification can be expressed by means of negation and existential quantification:  $\forall x \varphi$  is equivalent to  $\neg \exists x \neg \varphi$ . In fact, this is how SQL typically expresses universal quantification, as it lacks explicit statements for it.

We have already learned some interesting facts about the language of FO queries. Theorem 6.1 states that it has the same expressive power as the language of RA queries. This tells us that adding negation to  $\exists\text{FO}^+$  queries is the same as extending  $\text{RA}^+$  queries with the difference operation, and allowing inequalities in conditions. The problem of evaluating FO queries is PSPACE-complete in combined complexity (Theorem 7.1), and in DLOGSPACE in data complexity (Theorem 7.3). On the other hand, static analysis for FO queries, unlike  $\exists\text{FO}^+$  queries or UCQs, is undecidable (Theorems 8.1 and 8.3).

In the next two chapters, we address a different type of questions about the language of FO queries: what are the limitations of its expressive power? In other words, what kinds of queries FO cannot express. Knowing this justifies what practical languages need to add on top of the language of FO queries.

In this chapter, we study the inexpressibility in FO, and thus in RA, of queries requiring recursive computation. Note that recursion is the subject of Chapters 35 – 38. A canonical query of this type is reachability in directed graphs. Given a directed graph  $G$ , this query computes all the pairs of nodes  $(u, v)$  in  $G$  such that  $v$  is reachable from  $u$ . It can be computed by the following simple algorithm: a node  $v$  is reachable from a node  $u$  if

1. there is an edge from  $u$  to  $v$ , or
2. there is an edge from  $u$  to some node  $w$  such that  $v$  is reachable from  $w$ .

This algorithm is indeed recursive since the second item defines reachability in terms of itself. Such a description allows us to extract arbitrarily long paths from the input graph. For example, in the graph with nodes  $\{0, \dots, n\}$  and edges  $(i, i + 1)$  for all  $i \in [n - 1]$ , a node  $j$  is reachable from  $i$  iff  $i < j$ . Thus, the reachability query will extract paths of any length from 1 to  $n$ .

To show that such queries cannot be defined in FO, we present a fundamental property of FO queries, that is, *locality*. Intuitively, locality means that FO queries can only talk about objects that they see in a fixed neighbourhood of their free variables (this will be made precise below). This will let us prove that we cannot express queries such as reachability in directed graphs using FO queries. In particular, since seeing if there exists a path from one node to another may involve paths of arbitrary length, such a query is not local.

## Notion of Locality and Its Use

To build towards the notion of locality, we need to define some auxiliary terminology. Given a database  $D$ , its *Gaifman graph*, denoted  $G_D$ , is an undirected graph whose nodes are the elements of  $\text{Dom}(D)$ , and whose edges are the sets  $\{a, b\}$  such that  $a$  and  $b$  appear together in some fact of  $D$ , i.e., there is a fact of the form  $R(\dots, a, \dots, b, \dots)$  in  $D$ . Given  $a, b \in \text{Dom}(D)$ , the *distance* of  $a$  and  $b$  in  $D$ , denoted by  $d_D(a, b)$ , is the length of the shortest path in  $G_D$  from  $a$  to  $b$ ; by convention,  $d_D(a, a) = 0$ , and  $d_D(a, b) = \infty$  if there is no path in  $G_D$  from  $a$  to  $b$ . For a tuple  $\bar{a} = (a_1, \dots, a_n)$  over  $\text{Dom}(D)$ , the distance of  $\bar{a}$  and  $b$  in  $D$ , denoted  $d_D(\bar{a}, b)$ , is  $\min_{i \in [n]} \{d_D(a_i, b)\}$ . The *radius- $r$  ball* of  $\bar{a}$  in  $D$ , for  $r \geq 0$ , denoted  $B_r^D(\bar{a})$ , is the set  $\{b \in \text{Dom}(D) \mid d_D(a_i, b) \leq r \text{ for } i \in [n]\}$ . The *radius- $r$  neighborhood* (or simply  *$r$ -neighborhood*) of  $\bar{a}$  in  $D$ , denoted  $N_r^D(\bar{a})$ , is the set of facts  $\{R(\bar{b}) \in D \mid \bar{b} \text{ is over } B_r^D(\bar{a})\}$ , i.e., the set of all facts of  $D$  that contain only elements of  $B_r^D(\bar{a})$ . The tuple  $\bar{a}$  should be understood as a tuple of distinguished elements in  $N_r^D(\bar{a})$ . Finally, we say that two  $r$ -neighborhoods  $N_r^D(\bar{a})$  and  $N_r^D(\bar{b})$  are *isomorphic* if there exists a bijection  $h : \text{Dom}(N_r^D(\bar{a})) \rightarrow \text{Dom}(N_r^D(\bar{b}))$  such that  $h(\bar{a}) = \bar{b}$ , and  $R(\bar{c}) \in N_r^D(\bar{a})$  if and only if  $R(h(\bar{c})) \in N_r^D(\bar{b})$ . We can now define locality of queries.

### Definition 31.1: Locality of Queries

A  $k$ -ary query  $q$  over schema  $\mathbf{S}$  is  *$r$ -local*, for  $r \geq 0$ , if, for every database  $D$  of  $\mathbf{S}$ , and every two tuples  $\bar{a}, \bar{b} \in \text{Dom}(D)^k$  such that  $N_r^D(\bar{a})$  and  $N_r^D(\bar{b})$  are isomorphic,  $\bar{a} \in q(D)$  iff  $\bar{b} \in q(D)$ . A query is called *local* if it is  $r$ -local for some  $r \geq 0$ .

Interestingly, we can show that FO queries are local. In fact, we can show more than that; we can further connect the number  $r$  witnessing the locality of an FO query to the depth of quantifier nesting in formulae, defined below.

**Definition 31.2: Quantifier Rank of FO Formulae**

The *quantifier rank* of an FO formula is inductively defined as follows:

- $\text{rank}(\varphi) = 0$  if  $\varphi$  is atomic.
- $\text{rank}(\varphi \vee \psi) = \text{rank}(\varphi \wedge \psi) = \max\{\text{rank}(\varphi), \text{rank}(\psi)\}$ .
- $\text{rank}(\neg\varphi) = \text{rank}(\varphi)$ .
- $\text{rank}(\exists x \varphi) = \text{rank}(\forall x \varphi) = \text{rank}(\varphi) + 1$ .

Note that  $\text{rank}(\varphi)$  is not the total number of quantifiers, but rather the depth of their nesting, that is, the largest number of quantifiers one encounters in a branch of a parse tree of  $\varphi$ . For example,  $\text{rank}(\exists x S(x) \vee \forall y \neg R(y, y))$  is 1 and not 2. The formal statement about the locality of FO queries follows.

**Theorem 31.3**

Every FO query  $q = \varphi(\bar{x})$  is  $r$ -local for  $r = (3^{\text{rank}(\varphi)} - 1)/2$ .

Before we give the proof of Theorem 31.3, let us explain how it is used in order to obtain inexpressibility results for FO queries. Specifically, we look at the already seen *reachability* query, defined formally below.

**Definition 31.4: The Reachability Query**

The *reachability query*, denoted  $q_{\text{reach}}$ , over a schema  $\mathbf{S} = \{E[2]\}$  is such that, for every database  $D$  of  $\mathbf{S}$  and  $a, b \in \text{Dom}(D)$ ,  $(a, b) \in q_{\text{reach}}(D)$  if

- $E(a, b) \in D$ , or
- there are constants  $c_1, \dots, c_n \in \text{Dom}(D)$ , for  $n > 0$ , such that  $\{E(a, c_1), E(c_1, c_2), \dots, E(c_{n-1}, c_n), E(c_n, b)\} \subseteq D$ .

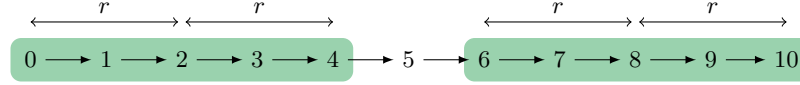
With Theorem 31.3 in place, to show that the reachability query cannot be expressed as an FO query, that is, there is no FO query  $q$  over  $\{E[2]\}$  such that  $q_{\text{reach}} \equiv q$ , it suffices to show that  $q_{\text{reach}}$  is not local.

**Proposition 31.5**

The reachability query  $q_{\text{reach}}$  over  $\mathbf{S} = \{E[2]\}$  is not local.

*Proof.* By contradiction, we assume that  $q_{\text{reach}}$  is  $r$ -local, for some  $r > 0$ . Consider a database  $D$  of  $\mathbf{S}$  consisting of the facts  $E(i, i + 1)$  for each  $i \in [0, 5r - 1]$ . In the output of  $q_{\text{reach}}$  on  $D$  we have all the pairs  $(i, j)$  with  $i < j$ . Note that the  $r$ -neighborhoods  $N_r^D(r, 4r)$  and  $N_r^D(4r, r)$  are isomorphic.



Fig. 31.1: Illustration of locality (with  $r = 2$ )

Indeed, each of these  $r$ -neighborhoods is a disjoint union of two chains of length  $2r$ , with the distinguished elements in the middle of those chains. This is illustrated in Figure 31.1, where  $E(i, j)$  is visualized using an arrow from  $i$  to  $j$ . Therefore, by locality,

$$(r, 4r) \in q_{\text{reach}}(D) \iff (4r, r) \in q_{\text{reach}}(D).$$

However,  $(r, 4r) \in q_{\text{reach}}(D)$  but  $(4r, r) \notin q_{\text{reach}}(D)$ , since  $r < 4r$ , which leads to a contradiction. This implies that  $q_{\text{reach}}$  is not local, as needed.  $\square$

## Proving Locality of FO

We proceed with the proof of Theorem 31.3. For clarity of the presentation, we assume that the FO query  $q = \varphi(\bar{x})$  is over a schema  $\mathbf{S}$  that consists of a binary relation  $R$ , but of course the proof generalizes to arbitrary schemas. We proceed by induction on the quantifier rank of  $\varphi$ .

Assume first that  $\text{rank}(\varphi) = 0$ . This implies that  $\varphi$  is an atomic formula  $x = y$  or  $R(x, y)$ . For a database  $D$  of  $\mathbf{S}$ , and tuples  $\bar{a}, \bar{b}$  over  $\text{Const}$ , the fact that  $N_0^D(\bar{a})$  and  $N_0^D(\bar{b})$  are isomorphic means that the same atomic formulae are true about  $\bar{a}$  and  $\bar{b}$  since the neighborhood does not contain any other elements, and thus,  $\bar{a} \in q(D)$  iff  $\bar{b} \in q(D)$ . Therefore,  $q$  is 0-local, as needed.

Assume now that  $\text{rank}(\varphi) = k$  for  $k > 0$ . We first observe that conjunction and disjunction, as well as negation, do not alter locality. More precisely, if both  $\varphi$  and  $\psi$  are  $r$ -local, then so are  $\varphi \vee \psi$  and  $\varphi \wedge \psi$ , as well as  $\neg\varphi$ . Notice also that conjunction, disjunction and negation do not alter the quantifier rank. Thus, the induction step should treat the case where  $\varphi(\bar{x}) = \exists y \psi(\bar{x}, y)$ ; we do not need to explicitly treat the case  $\varphi(\bar{x}) = \forall y \psi(\bar{x}, y)$  since  $\forall y \psi(\bar{x}, y)$  is equivalent to  $\neg \exists y \neg \psi(\bar{x}, y)$ , which has the same quantifier rank. Since  $\text{rank}(\varphi) = k$ , we get that  $\text{rank}(\psi) = k - 1$ . By induction hypothesis, the query  $\psi(\bar{x}, y)$  is  $r$ -local, where  $r = (3^{k-1} - 1)/2$ . We proceed to show that  $\varphi(\bar{x})$  is  $(3^k - 1)/2$ -local, or, equivalently,  $3r + 1$ -local.

**Lemma 31.6.** *It holds that  $\varphi(\bar{x})$  is  $(3r + 1)$ -local.*

*Proof.* Consider a database  $D$  of  $\mathbf{S}$ , and two tuples  $\bar{a}, \bar{b}$  over  $\text{Const}$  such that  $N_{3r+1}^D(\bar{a})$  and  $N_{3r+1}^D(\bar{b})$  are isomorphic. We need to establish that  $\bar{a} \in \varphi(\bar{x})(D)$  iff  $\bar{b} \in \varphi(\bar{x})(D)$ . To this end, it suffices to show that there exists a bijection  $f$  from  $\text{Dom}(D)$  to  $\text{Dom}(D)$  such that  $N_r^D(\bar{a}, c)$  and  $N_r^D(\bar{b}, f(c))$  are isomorphic,

for every  $c \in \text{Dom}(D)$ . Indeed, if such a bijection exists, the claim follows. Assuming that  $\bar{a} \in \varphi(\bar{x})(D)$ , we can find a witness  $c \in \text{Dom}(D)$  of the existential quantifier  $\exists y$  such that  $(\bar{a}, c) \in \psi(\bar{x}, y)(D)$ . But since  $N_r^D(\bar{a}, c)$  and  $N_r^D(\bar{b}, f(c))$  are isomorphic,  $f(c)$  is a witness that  $\bar{b} \in \varphi(\bar{x})(D)$ , since  $(\bar{b}, f(c))$  is in  $\psi(\bar{x}, y)(D)$ . For the converse, we use the bijection  $f^{-1}$  instead of  $f$ . The rest of the proof is devoted to showing the existence of the bijection  $f$ .

Assume that  $h$  is an isomorphism between  $N_{3r+1}^D(\bar{a})$  and  $N_{3r+1}^D(\bar{b})$ . Consider an arbitrary  $c \in B_{2r+1}^D(\bar{a})$ . Since  $B_r^D(c) \subseteq B_{3r+1}^D(\bar{a})$  and  $B_r^D(h(c)) \subseteq B_{3r+1}^D(\bar{b})$ , we thus obtain, from the fact that  $h$  is an isomorphism, that  $N_r^D(\bar{a}, c)$  and  $N_r^D(\bar{b}, h(c))$  are isomorphic; in fact the same isomorphism  $h$  witnesses this. Towards proving the existence of  $f$ , we define a *pointed* database as a pair  $(D', d)$  for  $d \in \text{Dom}(D')$ . An *isomorphism type* of pointed databases is the set of pointed databases that are all isomorphic to each other, and no pointed database from outside the set is isomorphic to them. In other words, it is an equivalence class with respect to the equivalence relation of being isomorphic. Note that each neighborhood  $N_r^D(c)$  is a pointed database. Let  $T_1, \dots, T_k$  be all the isomorphism types of  $r$ -neighborhoods present in  $D$ , and let  $n_i$  be the number of elements  $c \in \text{Dom}(D)$  such that  $N_r^D(c)$  belongs to  $T_i$ , for  $i \in [k]$ . Let  $m_i(\bar{a})$  be the number of elements  $c \in B_{2r+1}^D(\bar{a})$  such that  $N_r^D(c) \in T_i$ , and let  $m_i(\bar{b})$  be similarly defined for  $B_{2r+1}^D(\bar{b})$ . Since  $N_r^D(c)$  and  $N_r^D(h(c))$  are isomorphic for  $c \in B_{2r+1}^D(\bar{a})$ , we obtain that  $m_i(\bar{a}) = m_i(\bar{b})$ , for  $i \in [k]$ . Now, consider  $n_i - m_i(\bar{a})$ ; it is the number of elements  $c \in \text{Dom}(D) - B_{2r+1}^D(\bar{a})$  such that  $N_r^D(c) \in T_i$ , and let  $n_i - m_i(\bar{b})$  similarly count the number of elements  $c \in \text{Dom}(D) - B_{2r+1}^D(\bar{b})$  such that  $N_r^D(c) \in T_i$ . Since  $m_i(\bar{a}) = m_i(\bar{b})$ , we see that  $n_i - m_i(\bar{a}) = n_i - m_i(\bar{b})$ , for  $i \in [k]$ . Hence,

$$|\{c \notin B_{2r+1}^D(\bar{a}) \mid N_r^D(c) \in T_i\}| = |\{c \notin B_{2r+1}^D(\bar{b}) \mid N_r^D(c) \in T_i\}|$$

and thus, there is a bijection  $g : \text{Dom}(D) - B_{2r+1}^D(\bar{a}) \rightarrow \text{Dom}(D) - B_{2r+1}^D(\bar{b})$  so that  $N_r(c)$  and  $N_r(g(c))$  belong to the same set  $T_i$ , i.e., they are isomorphic.

We now define  $f : \text{Dom}(D) \rightarrow \text{Dom}(D)$ : on  $B_{2r+1}^D(\bar{a})$  we let  $f(c) = h(c)$ , and on  $\text{Dom}(D) - B_{2r+1}^D(\bar{a})$  we let  $f(c) = g(c)$ . For  $c \in B_{2r+1}^D(\bar{a})$ , we have already seen that  $N_r^D(\bar{a}, c)$  and  $N_r^D(\bar{b}, f(c))$  are isomorphic. If  $c \notin B_{2r+1}^D(\bar{a})$ , then  $N_r^D(\bar{a}, c)$  is the disjoint union of  $N_r^D(\bar{a})$  and  $N_r^D(c)$ , and  $N_r^D(\bar{b}, f(c))$  is the disjoint union of  $N_r^D(\bar{b})$  and  $N_r^D(f(c))$ . Thus,  $N_r^D(\bar{a}, c)$  and  $N_r^D(\bar{b}, f(c))$  are isomorphic since they are the disjoint union of isomorphic sets.  $\square$

This completes the proof of Theorem [31.3](#).

## Other Forms of Locality

Locality as seen in Definition [31.1](#) is not the only manifestation of this property for FO. A different one, called *Hanf-locality*, is defined in Exercise [4.18](#). Essentially, it says that if we have two databases,  $D$  and  $D'$ , and a bijection  $f$  between their domains such that  $N_r^D(a)$  and  $N_r^{D'}(f(a))$  are isomorphic,

for an appropriately chosen  $r$  (say,  $3^k$ ), then no Boolean FO query given by a formula  $\varphi$  of quantifier rank  $k$  distinguishes  $D$  from  $D'$ : it is either true in both, or false in both.

A different type of property makes it especially easy to show that some queries are inexpressible in FO. We present it again for graphs, that is, for schemas consisting of a single binary relation name, but it can be stated for arbitrary schemas too. Given a graph, i.e., a database  $D$  of the schema  $\{E[2]\}$ , *degrees* of  $D$  are numbers of the form  $|\{b \mid E(a, b) \in D\}|$  and  $|\{b \mid E(b, a) \in D\}|$ . These are essentially out-degrees and in-degrees of graph nodes. We write  $\deg(D)$  for the set of all degrees that occur in  $D$ . We can then show the following property for binary FO queries over  $\{E[2]\}$ , i.e., FO queries that take a graph as their input and return a graph as well.

**Proposition 31.7**

Let  $q$  be a binary FO query over  $\{E[2]\}$ . There is a function  $f : \mathbb{N} \rightarrow \mathbb{N}$  such that, for every database  $D$  of  $\{E[2]\}$  and  $k > 0$ , if all elements of  $\deg(D)$  are bounded by  $k$ , then  $|\deg(q(D))| \leq f(k)$ .

To use this to show that  $q_{\text{reach}}$  cannot be defined as an FO query, we refer again to the database  $D$  in the proof of Proposition 31.5 consisting of the facts  $E(i, i+1)$  for  $i \in [0, n-1]$ . In this case,  $\deg(D) = \{0, 1\}$ . Assume  $q_{\text{reach}}$  can be expressed as an FO query. By Proposition 31.7, there exists a function  $f : \mathbb{N} \rightarrow \mathbb{N}$  such that  $|\deg(q_{\text{reach}}(D))| \leq f(1)$ . However, if  $n > f(1)$ , this does not hold: the output of  $q_{\text{reach}}(D)$  has all edges  $(i, j)$  for  $0 \leq i < j \leq n$ , and thus,  $n+1$  different degrees, from 0 to  $n$ , are present in it.

## The Limits of First-Order Queries: Counting

In this chapter, we continue the study of the limitations of FO queries and what these limitations tell us in terms of extending the capabilities of practical query languages, and present two fundamental inexpressibility results concerning constant-free FO queries:

- They cannot express nontrivial statements about cardinalities of sets (for example, is the cardinality of a set even?).
- They cannot compare cardinalities of relations.

Let us stress that the results presented in this chapter do not hold for FO queries with constants. This is discussed further in the comments for Part [IV](#).

### An Easy Expressiveness Bound

We start by providing a rather preliminary result on the expressive power of constant-free FO queries (Theorem [32.3](#)), which in turn allows us to conclude that such queries cannot express nontrivial statements about the cardinalities of sets. Note that an FO sentence  $\varphi$  is called *constant-free* if it does not mention any constants, that is, the set  $\text{Dom}(\varphi)$  is empty. We further call an FO query  $\varphi(\bar{x})$  constant-free if  $\varphi$  is constant-free.

Let  $\mathbf{S} = \{R[1]\}$ , that is,  $\mathbf{S}$  is a schema consisting of a single unary relation name  $R$ . Databases of  $\mathbf{S}$  are essentially sets, i.e., they store the elements of a set  $R$ . It is easy to see that two databases  $D$  and  $D'$  of  $\mathbf{S}$  with  $|D| = |D'|$  satisfy exactly the same constant-free FO sentences over  $\mathbf{S}$ , i.e., for every such FO sentence  $\varphi$  over  $\mathbf{S}$ ,  $D \models \varphi$  iff  $D' \models \varphi$ . This is because  $D$  and  $D'$  are the same up to renaming of constants. We can thus define

$$\mu_n(\varphi) = \begin{cases} 1 & \text{if } D \models \varphi, \text{ for every } D \in \text{Inst}(\mathbf{S}) \text{ with } |D| = n \\ 0 & \text{if } D \not\models \varphi, \text{ for every } D \in \text{Inst}(\mathbf{S}) \text{ with } |D| = n. \end{cases}$$

We can then show the following useful technical result:

**Proposition 32.1**

Consider a constant-free FO sentence  $\varphi$  over  $\mathbf{S} = \{R[1]\}$ . There is  $k \in \mathbb{N}$  such that either  $\mu_n(\varphi) = 1$  for all  $n \geq k$ , or  $\mu_n(\varphi) = 0$  for all  $n \geq k$ .

To prove the above result we need a few basic notions and facts about first-order logic. A (possibly infinite) set  $T$  of first-order sentences over a schema  $\mathbf{S}$  is called a *first-order theory* over  $\mathbf{S}$ , or simply a *theory* over  $\mathbf{S}$ . The notion of satisfaction of an FO sentence by a database (see Definition 3.3) can be naturally extended to possibly infinite databases. We call a possibly infinite database of a schema  $\mathbf{S}$  a *model* of a theory  $T$  over  $\mathbf{S}$  if  $D \models \varphi$  for every  $\varphi \in T$ . We further say that  $T$  is *consistent* if it has at least one model. We know that a consistent theory  $T$  over  $\mathbf{S}$  has always a countably infinite model since  $\mathbf{S}$  is finite and Const is countably infinite; the latter is a consequence of a basic result in logic known as the Löwenheim-Skolem Theorem. A sentence  $\varphi$  is a consequence of a theory  $T$ , written  $T \models \varphi$ , if every model of  $T$  satisfies  $\varphi$ . We further know that if  $T \models \varphi$ , then there exists a finite subset  $T_0$  of  $T$  such that  $T_0 \models \varphi$ ; this is known as the Compactness Theorem of first-order logic. We are now ready to give the proof of Proposition 32.1.

*Proof (of Proposition 32.1).* Consider the theory  $T = \{\psi_n \mid n \in \mathbb{N}\}$ , where

$$\psi_n = \exists x_1 \cdots \exists x_n \bigwedge_{i,j \in [n]: i < j} \neg(x_i = x_j),$$

i.e., it states that there are  $n$  distinct elements. Clearly,  $T$  is consistent since any possibly infinite database of  $\mathbf{S}$  is a model of  $T$ . It is easy to show that:

**Lemma 32.2.** *Either  $T \models \varphi$  or  $T \models \neg\varphi$ .*

*Proof.* Notice that  $T \models \varphi$  and  $T \models \neg\varphi$  cannot be both true since in this case there exists a possibly infinite database  $D$  of  $\mathbf{S}$  such that  $D \models \varphi$  and  $D \models \neg\varphi$ , which cannot be the case. Assume now that  $T \not\models \varphi$  and  $T \not\models \neg\varphi$ . This implies that both theories  $T \cup \{\varphi\}$  and  $T \cup \{\neg\varphi\}$  are consistent, and thus, they have countably infinite models. Since there is only one countably infinite model, up to isomorphism, we get a contradiction as it cannot satisfy both  $\varphi, \neg\varphi$ .  $\square$

We now proceed to show the claim by considering the two cases provided by Lemma 32.2: either  $T \models \varphi$  or  $T \models \neg\varphi$ . Assume first that  $T \models \varphi$ . By the Compactness Theorem, there exists a finite subset  $T_0$  of  $T$  such that  $T_0 \models \varphi$ . Let  $\psi_k$  be the sentence with the largest index  $k$  among the sentences of  $T_0$ . It is clear that  $\psi_k \models \varphi$  since  $\psi_k \models \psi_m$  whenever  $m \leq k$ . Therefore, for every database  $D$  of  $\mathbf{S}$  with  $|D| \geq k$  it is the case that  $D \models \varphi$ . This implies that  $\mu_n(\varphi) = 1$  for every  $n \geq k$ . Analogously, if  $T \models \neg\varphi$ , then we can show that there is  $k \in \mathbb{N}$  such that  $\mu_n(\varphi) = 0$  for every  $n \geq k$ , and the claim follows.  $\square$

We can use Proposition 32.1 to show that only simple properties of cardinalities of sets can be expressed using constant-free FO queries. We proceed to make this more precise. Given a set of integers  $C \subseteq \mathbb{N}$ , let  $q_C$  be a Boolean query over the schema  $\mathbf{S} = \{R[1]\}$  that asks whether the cardinality of the input database is equal to an integer of  $C$ . In other words, for every database  $D$  of  $\mathbf{S}$ ,  $D \models q_C$  iff  $|D| \in C$ . Interestingly, we can precisely characterize when  $q_C$  is expressible as a constant-free FO query.

**Theorem 32.3**

Let  $C \subseteq \mathbb{N}$ , and  $\mathbf{S}$  be the schema  $\{R[1]\}$ . The following are equivalent:

1. There is a constant-free FO query  $q$  over  $\mathbf{S}$  such that  $q_C \equiv q$ .
2. Either  $C$  is a finite set, or  $\mathbb{N} - C$  is a finite set.

*Proof.* We first prove that (1) implies (2). Since, by hypothesis,  $q_C$  can be expressed as a constant-free FO query, Proposition 32.1 implies that there is an integer  $k \in \mathbb{N}$  such that one of the following statements hold:

- (i) For every  $D \in \text{Inst}(\mathbf{S})$  with  $|D| \geq k$  it is the case that  $D \models q_C$ .
- (ii) For every  $D \in \text{Inst}(\mathbf{S})$  with  $|D| \geq k$  it is the case that  $D \not\models q_C$ .

Assuming that (i) holds, there are finitely many integers, let us say  $i_1, \dots, i_m$ , for  $m \geq 0$ , such that, given a database  $D'$  of  $\mathbf{S}$  with  $|D'| \in \{i_1, \dots, i_m\}$ ,  $D' \not\models q_C$ . This in turn implies that  $\mathbb{N} - C$  is finite. Analogously, when (ii) holds, we can show that  $C$  is finite, and statement (2) follows.

We now show that (2) implies (1). This is shown by constructing a Boolean constant-free FO query  $q$  over  $\mathbf{S}$  such that  $q_C \equiv q$ . We first observe that, given an integer  $k \in \mathbb{N}$ , it is easy to construct an FO sentence  $\varphi_k$  over  $\mathbf{S}$  that is satisfied only by databases  $D$  over  $\mathbf{S}$  with  $|D| = k$ ; in particular,  $\varphi_k$  is

$$\begin{aligned} \exists x_1 \cdots \exists x_k \left( \bigwedge_{i=1}^k R(x_i) \wedge \bigwedge_{i,j \in [k] : i < j} \neg(x_i = x_j) \right) \wedge \\ \forall x_1 \cdots \forall x_{k+1} \left( \bigwedge_{i=1}^{k+1} R(x_i) \rightarrow \bigvee_{i,j \in [k+1] : i < j} x_i = x_j \right), \end{aligned}$$

where the first conjunct states that  $|D| \geq k$ , while the second conjunct states that  $|D| \leq k$ . By exploiting the fact that either  $C$  is finite, or  $\mathbb{N} - C$  is finite, the desired Boolean constant-free FO query  $q$  is defined as  $\varphi()$ , where

$$\varphi = \begin{cases} \bigvee_{i \in C} \varphi_i & \text{if } C \text{ is finite} \\ \neg \left( \bigvee_{i \in \mathbb{N} - C} \varphi_i \right) & \text{if } \mathbb{N} - C \text{ is finite.} \end{cases}$$

It is easy to verify that  $q_C \equiv q$ , and the claim follows.  $\square$

According to Theorem 32.3 it is impossible to check using a constant-free FO query whether the cardinality of a set is even, or, more generally, whether is divisible by some number  $n \in \mathbb{N}$ . Such a query is of the form  $q_C$  where  $C$  is an infinite set, and thus, not expressible as a constant-free FO query.

## Zero-One Law

Although Theorem 32.3 allows us to conclude that constant-free FO queries cannot express nontrivial statements about the cardinalities of sets, it is not powerful enough to tell us something about comparisons of cardinalities of relations. We proceed to establish a stronger property of FO sentences than the one established by Proposition 32.1 above, known as *zero-one law*, which will allow us to derive inexpressibility results concerning cardinality comparisons.

We start by reformulating the definition of the function  $\mu_n$  used above. We assume that the values occurring in a database are integers in order to be able to enumerate them. Then, for an FO sentence  $\varphi$  over a schema  $\mathbf{S}$ , let

$$\mu_n(\varphi) = \frac{|\{D \in \text{Inst}(\mathbf{S}) \mid \text{Dom}(D) = [n] \text{ and } D \models \varphi\}|}{|\{D \in \text{Inst}(\mathbf{S}) \mid \text{Dom}(D) = [n]\}|}.$$

In simple words,  $\mu_n(\varphi)$  is the proportion of databases of  $\mathbf{S}$  with  $\text{Dom}(D) = [n]$  that satisfy  $\varphi$ . Notice that this new definition of the function  $\mu_n$  applies to arbitrary schemas, not only to those with a single unary relation name. The intuition behind the quantity  $\mu_n(\varphi)$  can be described in purely probabilistic terms. Consider the finite set of databases  $D$  of  $\mathbf{S}$  with  $\text{Dom}(D) = [n]$ . Then  $\mu_n(\varphi)$  is the probability that a database one picks uniformly at random from this set satisfies  $\varphi$ . Now, by taking the limit  $\lim_{n \rightarrow \infty} \mu_n(\varphi)$ , we essentially describe the asymptotic behavior of the sequence  $(\mu_i(\varphi))_{i>0}$ ; intuitively, it defines the probability that a randomly picked database satisfies  $\varphi$ .

### Definition 32.4: 0–1 Law

We say that an FO sentence  $\varphi$  over a schema  $\mathbf{S}$  enjoys the 0–1 law if

$$\lim_{n \rightarrow \infty} \mu_n(\varphi) \in \{0, 1\}.$$

Intuitively, if an FO sentence  $\varphi$  over a schema  $\mathbf{S}$  enjoys the 0–1 law, then it is either satisfied by almost all the databases of  $\mathbf{S}$ , or violated by almost all the databases of  $\mathbf{S}$ . This is the case for constant-free FO sentences over a schema with a single unary relation name. Observe that there exists only one database  $D$  of  $\mathbf{S} = \{R[1]\}$  with  $\text{Dom}(D) = [n]$ . Therefore, for a constant-free FO sentence  $\varphi$  over  $\mathbf{S}$ , Proposition 32.1 says that the sequence  $(\mu_i(\varphi))_{i>0}$

eventually stabilizes, namely  $\lim_{n \rightarrow \infty} \mu_n(\varphi) \in \{0, 1\}$ . Interestingly, this can be shown for every constant-free FO sentence over an arbitrary schema.

Before showing this, let us stress that not all logical sentences enjoy the 0–1 law. There are, for example, FO sentences that mention constants that do not enjoy the 0–1 law; this is discussed further in the comments for Part [IV](#). Another example is the query  $\varphi_{\text{EVEN}}$ , expressed in a logical formalism that goes beyond first-order logic, that checks whether the cardinality of a database of the schema  $\mathbf{S} = \{R[1]\}$  is even, i.e., for every database  $D$  of  $\mathbf{S}$ ,  $D \models \varphi_{\text{EVEN}}$  iff  $|D|$  is even. Thus,  $\mu_n(\varphi_{\text{EVEN}})$  is 1 when  $n$  is even, and 0 when  $n$  is odd, which means that the limit  $\lim_{n \rightarrow \infty} \mu_n(\varphi)$  does not even exist.

**Theorem 32.5: 0–1 Law**

Every constant-free FO sentence enjoys the 0–1 law.

*Proof.* We shall not prove the result in its full generality, but instead consider the special case of constant-free FO sentences over a schema  $\mathbf{S}$  with two unary relation names, i.e.,  $\mathbf{S} = \{R[1], S[1]\}$ . Some ideas on how to extend the proof to graphs, i.e., to schemas with a single binary relation name, are explained in Exercise [4.13](#). This special case we study here builds on the proof of Proposition [32.1](#) illustrates key elements in the proof of the 0–1 law, and allows us to derive corollaries about cardinality comparisons.

The key elements are the same as in the general proof of the 0–1 law, and are summarized in the next technical lemma:

**Lemma 32.6.** *There exists a first-order theory  $T$  over  $\mathbf{S}$  such that:*

1.  $\lim_{n \rightarrow \infty} \mu_n(\psi) = 1$  for each  $\psi \in T$ , and
2.  $T$  has a unique, up to isomorphism, countably infinite model.

Before we give the proof of the lemma, let us explain how it can be used to show that every constant-free FO sentence enjoys the 0–1 law. Let  $T$  be the theory provided by Lemma [32.6](#). Condition (2) of the lemma, and the same argument as in the proof of Lemma [32.2](#) show that for every constant-free FO sentence  $\varphi$ , either  $T \models \varphi$  or  $T \models \neg\varphi$ . Consider now a constant-free FO sentence  $\varphi$  over  $\mathbf{S}$ . We proceed by case analysis:

- Assume that  $T \models \varphi$ . By the Compactness Theorem,  $\varphi$  is a consequence of a finite subset  $\{\psi_1, \dots, \psi_m\}$  of sentences of  $T$ . Since  $\lim_{n \rightarrow \infty} \mu_n(\psi_i) = 1$  for each  $i \in [m]$ , we get that  $\lim_{n \rightarrow \infty} \mu_n(\bigwedge_{i=1}^m \psi_i) = 1$ . Therefore, we have that  $\lim_{n \rightarrow \infty} \mu_n(\varphi) = 1$  since  $\mu_n(\varphi) \geq \mu_n(\bigwedge_{i=1}^m \psi_i)$ .
- Assume now that  $T \models \neg\varphi$ . We apply the same argument to  $\neg\varphi$ , and conclude that  $\lim_{n \rightarrow \infty} \mu_n(\neg\varphi) = 1$ , which implies that  $\lim_{n \rightarrow \infty} \mu_n(\varphi) = 0$ .

Hence,  $\lim_{n \rightarrow \infty} \mu_n(\varphi) \in \{0, 1\}$ . We proceed with the proof of Lemma [32.6](#).



*Proof (of Lemma 32.6).* We construct a theory  $T$  over  $\mathbf{S}$ , and then show that it satisfies conditions (1) and (2). For each  $k \in \mathbb{N}$ , the theory  $T$  contains

$$\psi_k(RS) = \exists x_1 \cdots \exists x_k \left( \bigwedge_{i,j \in [k]: i < j} \neg(x_i = x_j) \wedge \bigwedge_{i \in [k]} (R(x_i) \wedge S(x_i)) \right),$$

which states that, for a database  $D$ , the set  $R^D \cap S^D$  has at least  $k$  elements,

$$\psi_k(R\bar{S}) = \exists x_1 \cdots \exists x_k \left( \bigwedge_{i,j \in [k]: i < j} \neg(x_i = x_j) \wedge \bigwedge_{i \in [k]} (R(x_i) \wedge \neg S(x_i)) \right),$$

which states that  $R^D - S^D$  has at least  $k$  elements, and

$$\psi_k(\bar{R}S) = \exists x_1 \cdots \exists x_k \left( \bigwedge_{i,j \in [k]: i < j} \neg(x_i = x_j) \wedge \bigwedge_{i \in [k]} (\neg R(x_i) \wedge S(x_i)) \right),$$

which states that  $S^D - R^D$  has at least  $k$  elements.

We first observe that the theory  $T$  is consistent. Indeed, it has a countably infinite model  $D$  such that  $R^D \cup S^D = \mathbb{N}$  (recall the assumption that the values occurring in a database are integers), and  $R^D \cap S^D$ ,  $R^D - S^D$ , and  $S^D - R^D$  are countably infinite; for example, we can take  $R^D = \{3n, 3n+1 \mid n \in \mathbb{N}\}$  and  $S^D = \{3n, 3n+2 \mid n \in \mathbb{N}\}$ . It is easy to see that any two such infinite databases are isomorphic, and thus, up to isomorphism,  $T$  has only one countably infinite model, satisfying condition (2) of Lemma 32.6.

We now prove that  $T$  satisfies the first condition. We show that, for each integer  $k \in \mathbb{N}$ , the sentences  $\psi_k(RS)$ ,  $\psi_k(R\bar{S})$ , and  $\psi_k(\bar{R}S)$  are true in almost all databases of  $\mathbf{S}$ . We do this as an illustration for the sentence  $\psi_k(RS)$ . To this end, we consider its negation stating that  $|R^D \cap S^D| < k$ , for a database  $D$ . We first provide an upper bound for  $\mu_n(\neg\psi_k(RS))$ :

- The numerator of  $\mu_n(\neg\psi_k(RS))$  coincides with the number of different ways that we can choose two sets  $R$  and  $S$  from  $[n]$  such that  $R \cup S = [n]$  and  $|R \cap S| < k$ . For each  $j < k$ , we have  $\binom{n}{j}$  ways to choose an intersection  $R \cap S$  of cardinality less than  $k$ . Then we need to choose the elements of  $R - S$  from the remaining  $n - j$  elements, and there are  $2^{n-j}$  ways of doing so. The remaining elements must belong to  $S - R$  since  $R \cup S = [n]$ . Thus, there are  $\binom{n}{j} \cdot 2^{n-j}$  ways to choose two sets from  $[n]$  whose intersection has exactly  $j$  elements, and whose union contains all the elements of  $[n]$ . This means that there are at most

$$\sum_{j \in [0, k-1]} \binom{n}{j} \cdot 2^{n-j} \leq n^k \cdot 2^n$$

ways of choosing two sets whose intersection has cardinality less than  $k$ .

- The denominator of  $\mu_n(\neg\psi_k(RS))$  coincides with the number of ways to choose two sets  $R$  and  $S$  from  $[n]$  such that  $R \cup S = [n]$ . There are  $3^n$  ways of doing so since, for each element of  $[n]$ , there are 3 possibilities: it can either belong to  $R \cap S$ , or  $R - S$ , or  $S - R$ .

From the above analysis, we conclude that

$$\mu_n(\neg\psi_k(RS)) \leq n^k \left(\frac{2}{3}\right)^n$$

which means  $\lim_{n \rightarrow \infty} \mu_n(\neg\psi_k(RS)) = 0$ , and therefore

$$\lim_{n \rightarrow \infty} \mu_n(\psi_k(RS)) = 1.$$

The proof for  $\psi_k(R\bar{S})$  and  $\psi_k(\bar{S}R)$  are very similar. This shows that  $T$  satisfies both conditions (1) and (2), and concludes the proof of Lemma 32.6.  $\square$

This completes the proof of Theorem 32.5.  $\square$

Theorem 32.5 allows us to establish inexpressibility results concerning cardinality comparisons. For  $\diamond \in \{<, \leq, =\}$ , let  $q_\diamond$  be a Boolean query over the schema  $\mathbf{S} = \{R[1], S[1]\}$  that compares the cardinalities of the relations  $R^D$  and  $S^D$  for a database  $D$  according to the comparison  $\diamond$ . In other words, for every database  $D$  of  $\mathbf{S}$ , the query  $q_\diamond$  is true in  $D$  iff the comparison  $|R^D| \diamond |S^D|$  is true. We now show via a 0–1 law argument that none of the comparisons  $|R| = |S|$ , or  $|R| < |S|$ , or  $|R| \leq |S|$ , is expressible as a constant-free FO query.

**Theorem 32.7**

Let  $\mathbf{S} = \{R[1], S[1]\}$ . For every  $\diamond \in \{<, \leq, =\}$ , there is no constant-free FO query  $q$  over  $\mathbf{S}$  such that  $q_\diamond \equiv q$ .

*Proof.* We prove the result for the case of  $q_<$  via a 0–1 law argument; a very similar argument works for  $q_\leq$ . The case of  $q_=>$  can be also shown via a 0–1 law argument, and is left as an exercise (see Exercise 4.10).

Let  $F_n^=$  be the number of all databases  $D$  of  $\mathbf{S}$  such that  $\text{Dom}(D) = [n]$  and  $|R^D| = |S^D|$ ; we define  $F_n^<$  and  $F_n^>$  likewise. From the proof of Theorem 32.5 we know that  $F_n^= + F_n^< + F_n^> = 3^n$ . Moreover, by symmetry,  $F_n^< = F_n^>$ , and thus,  $F_n^= + 2F_n^< = 3^n$ . We first estimate the value  $F_n^=$ . To have a database  $D$  in which  $|R^D| = |S^D|$ , for every  $k \leq \lfloor n/2 \rfloor$ , we can pick  $k$  elements to belong to  $R^D - S^D$ , from the remaining  $n - k$  elements we can pick  $k$  elements to belong to  $S^D - R^D$ , and the remaining ones belong to  $R^D \cap S^D$ . Hence

$$F_n^= = \sum_{k \leq \lfloor n/2 \rfloor} \binom{n}{k} \binom{n-k}{k}$$

and one can show (Exercise 4.12) that

$$\lim_{n \rightarrow \infty} \frac{F_n^=}{3^n} = 0.$$

Assume now that there is a constant-free FO query  $q = \varphi()$  such  $q_< \equiv q$ , i.e.,  $q$  expresses the condition  $|R^D| < |S^D|$ . Then

$$\lim_{n \rightarrow \infty} \mu_n(\varphi) = \lim_{n \rightarrow \infty} \frac{F_n^<}{3^n} = \lim_{n \rightarrow \infty} \frac{3^n - F_n^=}{2 \cdot 3^n} = \frac{1}{2} - \lim_{n \rightarrow \infty} \frac{F_n^=}{2 \cdot 3^n} = \frac{1}{2}$$

which contradicts the 0–1 law, and the claim follows. □

## Adding Aggregates and Grouping

When the core of SQL was presented in Chapter 5, a commonly used feature of it was omitted, namely *aggregation*. It is typically used together with *grouping* to apply numerical functions to entire columns of a relation. Recall that one of the relation names in the schema that we usually use in examples is

City [ cid, cname, country ]

If we want to know how many cities each country has, we can write in SQL

```
SELECT country, COUNT(cid) AS city_count
FROM City
GROUP BY country
```

For each value of the country attribute, it groups together all the tuples having this value, and then counts the number of occurrences of cid in such a group and outputs it as the value of the new attribute city\_count. Here, COUNT is an *aggregate function*: it applies to a collection, and produces a single numerical value. The standard aggregates of SQL, in addition to COUNT, are SUM and AVG that compute the sum and the average of a collection of *numbers*, as well as MIN and MAX that compute the minimum and the maximum.

Queries with aggregates are extremely common and useful in practice; for example, “find the average grade for each class” or “find the total cost of products sold to each country”. The addition of aggregate functions, however, takes us out of the realm of FO and RA queries. Indeed, by Theorem 32.7, we know that cardinality comparisons are not expressible via FO queries. However, they are easily expressible with the help of aggregate functions:

```
SELECT DISTINCT 1 FROM R, S
WHERE (SELECT COUNT(*) FROM R) > (SELECT COUNT(*) FROM S)
```

outputs 1 if  $|R| > |S|$ , and nothing otherwise.

In this chapter, we introduce a query language with aggregates and grouping based on RA. Of course one can also define a logical language with aggregates, but this is cumbersome for the reasons that are explained below. Note that in Chapter 34, we analyze queries that cannot be expressed even if we have the powerful features of aggregates and grouping. Let us now discuss the technical issues that arise due to aggregates:

**Numerical Attributes.** So far we assumed that database elements come from a countably infinite set  $\text{Const}$  of values. With the addition of aggregates, however, we can no longer make this assumption, as we need to distinguish attributes that are *numerical*. For example, we can only apply **AVG** over numbers. Thus, in the description of relational schemas, it is no longer sufficient to simply state what the arity of each relation name is. In addition, we need to provide information about attributes that are numerical. The standard approach for solving this technical issue is to consider *two-sorted schemas*: there will be columns populated by the usual values from  $\text{Const}$ , and columns populated by values from a numerical domain  $\text{Num}$ ; e.g., the natural numbers  $\mathbb{N}$ , or the integers  $\mathbb{Z}$ , or the rationals  $\mathbb{Q}$ .

**Infinite Numerical Domains.** The second issue manifests itself when we deal with logical languages for aggregates. In Chapter 3, we defined the satisfaction of a logical formula over the active domain of the database and the formula (recall that Definition 3.3 defines the active domain semantics of first-order logic). However, aggregates can produce new numerical values that do not occur in the active domain. Therefore, the satisfaction of formulae must be defined with respect to the entire infinite numerical domain  $\text{Num}$ . This leads to the situation where a logical formula  $\varphi$  may be satisfied by an infinite number of assignments of values to its free variables, and thus, the expression  $\varphi(\bar{x})$ , for some tuple  $\bar{x}$  over the free variables of  $\varphi$ , may not define a query since its output on a database may be infinite. Although we can define a logical language with aggregates that does not exhibit this problem, its syntax is cumbersome. Therefore, we present the language with aggregates and grouping at the level of relational algebra, where the above problem does not arise. We present an extension of first-order logic with aggregates in Chapter 34 that uses a relatively simple syntax, and show that RA with aggregates and grouping translates to it. However, this logical language will not serve as the basis for defining a query language with aggregates, but rather as a technical tool for analyzing the expressive power of aggregates.

## Two-Sorted Schemas, Databases and Queries

We first revisit the notions of database schema, database instance, and query in order to take into account the fact that relations can now have both ordinary and numerical values. These are actually straightforward adaptations of

the definitions given in Chapter 2, but, for the sake of completeness and readability, we proceed to properly introduce those notions. As usual, for technical clarity, we adopt the unnamed perspective.

Each relation name  $R$  in a schema should come not simply with its arity  $k$ , but rather with a tuple  $\tau$  of arity  $k$  over  $\{\mathbf{o}, \mathbf{n}\}$ , where  $\mathbf{o}$  indicates a column of ordinary type taking values from  $\mathbf{Const}$ , and  $\mathbf{n}$  indicates a column of numerical type taking values from a set of numerical values  $\mathbf{Num}$ . The formal definition of two-sorted database schemas follows.

**Definition 33.1: Two-Sorted Database Schema**

A *two-sorted (database) schema* is a partial function

$$\mathbf{S} : \text{Rel} \rightarrow \{\mathbf{o}, \mathbf{n}\}^k,$$

for  $k \in \mathbb{N}$ , such that  $\text{Dom}(\mathbf{S})$  is finite. For a relation name  $R \in \text{Dom}(\mathbf{S})$ , the *arity of  $R$  under  $\mathbf{S}$* , denoted  $\text{ar}_{\mathbf{S}}(R)$ , is defined as  $k$ .

In order to avoid heavy notation, we write  $\text{ar}(R)$  instead of  $\text{ar}_{\mathbf{S}}(R)$  for the arity of  $R$  under  $\mathbf{S}$ . We may even write  $R : \tau$  to indicate that  $\mathbf{S}(R) = \tau$ . As for plain schemas, a two-sorted schema can be naturally seen as a finite set of relation names. We may also write  $\mathbf{S} = \{R_1 : \tau_1, \dots, R_n : \tau_n\}$  for the fact that  $\text{Dom}(\mathbf{S}) = \{R_1, \dots, R_n\}$  and  $\mathbf{S}(R_i) = \tau_i$  for each  $i \in [n]$ .

The elements of database tuples are coming from the set of values  $\mathbf{Const}$ , and the set of numerical values  $\mathbf{Num}$ , namely a *two-sorted database tuple* is an element of  $(\mathbf{Const} \cup \mathbf{Num})^k$  for some  $k \in \mathbb{N}$ . A *two-sorted relation instance* is a finite set  $S$  of two-sorted database tuples of the same arity  $k$ . We say that  $k$  is the *arity of  $S$* , denoted by  $\text{ar}(S)$ . By **tsRI** (for *two-sorted relation instances*) we denote the set of all such relation instances. The formal definition of a database instance of a two-sorted schema follows.

**Definition 33.2: Two-Sorted Database Instance**

A *database instance*  $D$  of a two-sorted schema  $\mathbf{S}$  is a function

$$D : \text{Dom}(\mathbf{S}) \rightarrow \text{tsRI}$$

such that, for every  $R \in \text{Dom}(\mathbf{S})$ , the following hold:

- $\text{ar}(D(R)) = \text{ar}_{\mathbf{S}}(R)$ , and
- with  $\mathbf{S}(R) = (\tau_1, \dots, \tau_k)$  and  $D(R) = (a_1, \dots, a_k)$ , we have that  $a_i \in \mathbf{Const}$  if  $\tau_i = \mathbf{o}$ , and  $a_i \in \mathbf{Num}$  if  $\tau_i = \mathbf{n}$ , for every  $i \in [k]$ .

We will refer to a database instance of a two sorted schema as a *two-sorted database*, or simply as a database whenever is clear that the underlying schema is two-sorted. The active domain (or simply domain) of a two-sorted database

$D$  is defined in the same way as the active domain of a plain database given in Chapter 2 and is denoted  $\text{Dom}(D)$ ; we will *never* use the term domain, and the notation  $\text{Dom}(D)$ , to refer to the domain of the function  $D$ , i.e.,  $\text{Dom}(\mathbf{S})$ . A two-sorted database can be naturally seen as a finite set of facts.

We can now naturally define two-sorted queries as functions that map two-sorted databases to *finite* sets of tuples of the same type over  $\text{Const} \cup \text{Num}$ . For a two-sorted schema  $\mathbf{S}$ , we write  $\text{Inst}(\mathbf{S})$  for the set of all databases of  $\mathbf{S}$ .

### Definition 33.3: Two-Sorted Queries

Consider a two-sorted database schema  $\mathbf{S}$ . A *query of type*  $(\tau_1, \dots, \tau_k) \in \{\mathbf{o}, \mathbf{n}\}^k$ , for  $k \geq 0$ , over  $\mathbf{S}$  is a function of the form

$$q : \text{Inst}(\mathbf{S}) \rightarrow \mathcal{P}_{\text{fin}}((\text{Const} \cup \text{Num})^k)$$

such that, for every  $D \in \text{Inst}(\mathbf{S})$  with  $q(D) = (a_1, \dots, a_k)$ , it is the case that  $a_i \in \text{Const}$  if  $\tau_i = \mathbf{o}$ , and  $a_i \in \text{Num}$  if  $\tau_i = \mathbf{n}$ , for every  $i \in [k]$ .

Regarding the domain  $\text{Num}$ , we assume that it comes equipped with:

- *(Numerical) predicates*: a predicate  $P$  or arity  $k > 0$  over  $\text{Num}$  is a subset of  $\text{Num}^k$ . For brevity, we say that  $P(a_1, \dots, a_k)$  holds if  $(a_1, \dots, a_k) \in P$ . Examples of binary predicates are  $=$  and  $<$ .
- *(Numerical) functions*: a function of arity  $k > 0$  over  $\text{Num}$  is a function of the form  $f : \text{Num}^k \rightarrow \text{Num}$ . Examples of binary functions are  $+$  and  $\times$ .
- *Aggregate functions or aggregates*: an aggregate  $\mathcal{F}$  over  $\text{Num}$  is a function that maps bags (or multisets) of elements of  $\text{Num}$  to  $\text{Num}$ . In a bag, unlike a set, an element can appear multiple times; e.g.,  $\{1, 1, 2, 4\}$  is a bag that has two occurrences of 1, and one occurrence of 2 and 4. We shall use the brackets  $\{\}$  to distinguish bags from sets.

Let us stress that aggregates must be applied to bags rather than sets since values in databases can repeat. Here is a simple example that illustrates this.

### Example 33.4: Aggregates over Bags

Consider the two-sorted database

$$D = \{R(a, 1), R(b, 1), R(c, 2), R(d, 4)\}.$$

It is clear that the second column of  $R^D$  is the bag

$$\{1, 1, 2, 4\}.$$

Assume now that we are interested in computing the average of the numerical values occurring in the second column of  $R^D$ . Applying the aver-

age aggregate to  $\{1, 1, 2, 4\}$  will correctly produce the value 2. However, if we apply it to the set  $\{1, 2, 4\}$ , we get an incorrect value  $2.333 \dots$ .

## Syntax of RA with Aggregates and Grouping

Assuming a set  $\Omega$  of predicates, functions, and aggregates over the numerical domain  $\text{Num}$ , we are going to define *relational algebra with aggregates and grouping*, denoted  $\text{RA}_{\text{Aggr}}(\Omega)$ . Its expressions, unlike expressions of RA, will be *typed*. The type of an expression is again a tuple over  $\{\mathbf{o}, \mathbf{n}\}$  indicating which attributes of the output are ordinary and which are numerical. In addition to the usual operations, we add the operations of selection based on numerical predicates, applying functions, and applying aggregates.

Before giving the formal definition of  $\text{RA}_{\text{Aggr}}(\Omega)$ , we first need to introduce some auxiliary notions. For a tuple  $\tau = (\tau_1, \dots, \tau_k) \in \{\mathbf{o}, \mathbf{n}\}^k$ , where  $k \in \mathbb{N}$ , we inductively define  $(\Omega, \tau)$ -terms, and their associated *types*, as follows:

- Every  $a \in \text{Const}$  is an  $(\Omega, \tau)$ -term of type  $\mathbf{o}$ .
- Every integer  $i \in [k]$  is an  $(\Omega, \tau)$ -term of type  $\tau_i$ .
- If  $f$  is an  $m$ -ary numerical function from  $\Omega$ , and  $t_1, \dots, t_m$  are  $(\Omega, \tau)$ -terms of type  $\mathbf{n}$ , then  $f(t_1, \dots, t_m)$  is an  $(\Omega, \tau)$ -term of type  $\mathbf{n}$ .

We write  $\tau(t)$  for the type of an  $(\Omega, \tau)$ -term  $t$ , and  $\text{VAR}(t)$  for the set of all integers that appear in  $t$ . An  $(\Omega, \tau)$ -condition  $\theta$  is a Boolean combination of statements of the form  $i \doteq j$  and  $i \neq j$ , for  $i, j \in [k]$ , and  $P(i_1, \dots, i_m)$ , where  $i_1, \dots, i_m \in [k]$  and  $P$  is an  $m$ -ary predicate from  $\Omega$ .

### Definition 33.5: Syntax of RA with Aggregates and Grouping

Consider a set  $\Omega$  of predicates, functions, and aggregates over  $\text{Num}$ . We inductively define  $\text{RA}_{\text{Aggr}}(\Omega)$  *expressions* over a two-sorted schema  $\mathbf{S}$ , and their associated *types*, as follows:

**Base Expression.** If  $R : \tau$  belongs to  $\mathbf{S}$ , then  $R$  is an  $\text{RA}_{\text{Aggr}}(\Omega)$  expression over  $\mathbf{S}$  of type  $\tau$ .

**Selection.** If  $e$  is an  $\text{RA}_{\text{Aggr}}(\Omega)$  expression over  $\mathbf{S}$  of type  $\tau$ , and  $\theta$  is an  $(\Omega, \tau)$ -condition, then  $\sigma_\theta(e)$  is an  $\text{RA}_{\text{Aggr}}(\Omega)$  expression over  $\mathbf{S}$  of type  $\tau$ .

**Projection.** If  $e$  is an  $\text{RA}_{\text{Aggr}}(\Omega)$  expression of type  $\tau = (\tau_1, \dots, \tau_k)$ , for  $k \geq 0$ , and  $\alpha = (t_1, \dots, t_m)$ , for  $m \geq 0$ , is a list of  $(\Omega, \tau)$ -terms, then  $\pi_\alpha(e)$  is an  $\text{RA}_{\text{Aggr}}(\Omega)$  expression over  $\mathbf{S}$  of type  $(\tau(t_1), \dots, \tau(t_m))$ .

**Cartesian Product.** If  $e_1, e_2$  are  $\text{RA}_{\text{Aggr}}(\Omega)$  expressions over  $\mathbf{S}$  of type  $(\tau_1, \dots, \tau_k)$  and  $(\tau'_1, \dots, \tau'_m)$ , for  $k, m \geq 0$ , respectively, then  $(e_1 \times e_2)$  is an  $\text{RA}_{\text{Aggr}}(\Omega)$  expression over  $\mathbf{S}$  of type  $(\tau_1, \dots, \tau_k, \tau'_1, \dots, \tau'_m)$ .



**Aggregation and Grouping.** If  $e$  is an  $\text{RA}_{\text{Aggr}}(\Omega)$  expression over  $\mathbf{S}$  of type  $\tau = (\tau_1, \dots, \tau_k)$ , for  $k \geq 0$ ,  $\alpha = (i_1, \dots, i_m)$ , for  $m \geq 0$ , is a list of integers from  $[k]$ , and  $t_1, \dots, t_\ell$ , for  $\ell \geq 0$ , are  $(\Omega, \tau)$ -terms of type  $\mathbf{n}$  such that  $\text{VAR}(t_i) \cap \{i_1, \dots, i_m\} = \emptyset$ , for each  $i \in [\ell]$ , then

$$\text{Aggr}_\alpha[\mathcal{F}_1(t_1), \dots, \mathcal{F}_\ell(t_\ell)](e),$$

where  $\mathcal{F}_1, \dots, \mathcal{F}_\ell$  are aggregates from  $\Omega$ , is an  $\text{RA}_{\text{Aggr}}(\Omega)$  expression over  $\mathbf{S}$  of type  $(\tau_{i_1}, \dots, \tau_{i_m}, \mathbf{n}, \dots, \mathbf{n})$  with  $\mathbf{n}$  repeated  $\ell$  times.

**Union.** If  $e_1, e_2$  are  $\text{RA}_{\text{Aggr}}(\Omega)$  expressions over  $\mathbf{S}$  of type  $\tau$ , then  $(e_1 \cup e_2)$  is an  $\text{RA}_{\text{Aggr}}(\Omega)$  expression over  $\mathbf{S}$  of type  $\tau$ .

**Difference.** If  $e_1, e_2$  are  $\text{RA}_{\text{Aggr}}(\Omega)$  expressions over  $\mathbf{S}$  of type  $\tau$ , then  $(e_1 - e_2)$  is an  $\text{RA}_{\text{Aggr}}(\Omega)$  expression over  $\mathbf{S}$  of type  $\tau$ .

Concerning the base expression in Definition 33.5, observe the difference with the definition of RA (Definition 4.1). In particular, in ordinary RA we have base expressions  $\{a\}$  of arity 1, where  $a \in \text{Const}$ . Thus, one would expect in Definition 33.5 base expressions of the form  $\{a\}$  of type  $\mathbf{o}$ . However, such base expressions would be redundant since every element of  $\text{Const}$  is a term, and thus, can be expressed via expressions of the form  $\pi_{(a)}(R)$  of type  $\mathbf{o}$  (based on the semantics of  $\text{RA}_{\text{Aggr}}(\Omega)$  expressions as defined next).

## Semantics of RA with Aggregates and Grouping

The semantics of the standard relational algebra operations is the same as it was presented in Chapter 4. For numerical selection conditions,  $P(i_1, \dots, i_m)$  is true in a tuple  $(a_1, \dots, a_k)$  iff  $i_1, \dots, i_m$  correspond to columns of type  $\mathbf{n}$  and  $P(a_{i_1}, \dots, a_{i_m})$  holds. For example,  $<(1, 3)$  is true in a tuple  $(a_1, a_2, a_3)$  iff the first and the third components are numerical and  $a_1 < a_3$ .

It remains to explain the new generalized projection, and aggregation with grouping. We first provide informal explanations by means of SQL examples.

### Example 33.6: Generalized Projection

Projection allows us compute functions on attributes and output them. For example, for  $R[A, B, C]$  with  $(R, A) \triangleleft (R, B) \triangleleft (R, C)$ , the SQL query

```
SELECT R.A+R.C, R.B*R.C
FROM R
```

is translated into the following  $\text{RA}_{\text{Aggr}}(\Omega)$  expression of type  $(\mathbf{n}, \mathbf{n})$ :

$$\pi_{(\text{add}(1,3), \text{mult}(2,3))}(R)$$

assuming that  $\Omega$  contains the binary functions `add` and `mult` such that

$$\text{add}(x, y) = x + y \quad \text{and} \quad \text{mult}(x, y) = x \cdot y.$$

Grouping and aggregation can also be explained intuitively by using SQL queries. In particular,  $\text{Aggr}_{(i_1, \dots, i_m)}[\mathcal{F}_1(t_1), \dots, \mathcal{F}_\ell(t_\ell)](R)$  corresponds to

```
SELECT   $i_1, \dots, i_m, \mathcal{F}_1(t_1), \dots, \mathcal{F}_\ell(t_\ell)$ 
FROM    R
GROUP BY  $i_1, \dots, i_m$ 
```

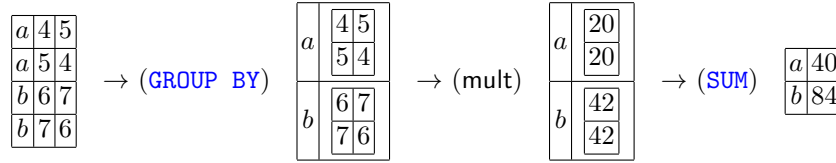
assuming that the attributes of  $R$  are  $1, \dots, k$ .

### Example 33.7: Aggregation and Grouping

For a ternary relation  $R$ , the evaluation of

$$\text{Aggr}_{(1)}[\text{SUM}(\text{mult}(2, 3))](R)$$

is illustrated below:



The **GROUP BY** clause does the grouping relative to the first attribute, keeping duplicates if necessary. Then, the term values are computed, again preserving duplicates; for example, for both tuples (4, 5) and (5, 4), the result of the multiplication is 20, and thus two copies are kept. Finally, the aggregate sums up the values of those terms.

We now define the semantics of generalized projection, and aggregation with grouping. We first need some auxiliary notions. Consider an  $(\Omega, \tau)$ -term  $t$ , where  $\tau = (\tau_1, \dots, \tau_k) \in \{\mathbf{o}, \mathbf{n}\}^k$ . A tuple  $\bar{a} = (a_1, \dots, a_n) \in (\text{Const} \cup \text{Num})^n$ , for  $n \geq k$ , is *compatible* with  $t$  if  $\tau_i = \mathbf{o}$  implies  $a_i \in \text{Const}$ , and  $\tau_i = \mathbf{n}$  implies  $a_i \in \text{Num}$ , for every  $i \in \text{VAR}(t)$ . We define the *evaluation* of  $t$  over a tuple  $\bar{a} = (a_1, \dots, a_n)$  that is compatible with it, denoted  $\text{eval}(t|\bar{a})$ , as follows:

- if  $t = a$  with  $a \in \text{Const}$ , then  $\text{eval}(t|\bar{a}) = a$ ,
- if  $t = i$  with  $i \in [k]$ , then  $\text{eval}(t|\bar{a}) = a_i$ , and
- if  $t = f(t_1, \dots, t_m)$ , then  $\text{eval}(t|\bar{a}) = f(\text{eval}(t_1|\bar{a}), \dots, \text{eval}(t_m|\bar{a}))$ .

Furthermore, given a two-sorted relation  $R$  of arity  $n \in \mathbb{N}$  consisting of tuples that are compatible with  $t$ , a list of integers  $\alpha = (i_1, \dots, i_m)$ , for  $m \geq 0$ , from  $[n]$ , and a tuple  $\bar{a} \in \pi_\alpha(R)$ , we define the bag

$$B(\bar{a}, R, t) = \{ \text{eval}(t|\bar{c}) \mid \bar{c} \in R \text{ and } \bar{a} = \pi_\alpha(\bar{c}) \}.$$

For instance, going back to Example 33.7 for  $t = \text{mult}(2, 3)$ ,

$$B((a), R, t) = \{ \text{eval}(t|(a, 4, 5)), \text{eval}(t|(a, 5, 4)) \} = \{ 20, 20 \}.$$

### Definition 33.8: Semantics of Projection and Aggregation

Consider a set  $\Omega$  of predicates, functions, and aggregates over  $\text{Num}$ . Let  $D$  be a database of a two-sorted schema  $\mathbf{S}$ , and  $e$  an  $\text{RA}_{\text{Aggr}}(\Omega)$  over  $\mathbf{S}$ . We define the *output*  $e(D)$  of  $e$  on  $D$  as follows:

- If  $e = \pi_{(t_1, \dots, t_m)}(e_1)$ , where  $e_1$  is an  $\text{RA}_{\text{Aggr}}(\Omega)$  expression, then

$$e(D) = \{ (\text{eval}(t_1|\bar{a}), \dots, \text{eval}(t_m|\bar{a})) \mid \bar{a} \in e_1(D) \}.$$

- If  $e = \text{Aggr}_{(i_1, \dots, i_m)}[\mathcal{F}_1(t_1), \dots, \mathcal{F}_\ell(t_\ell)](e_1)$ , where  $e_1$  is an  $\text{RA}_{\text{Aggr}}(\Omega)$  expression, then

$$\{ (\bar{a}, \mathcal{F}_1(B(\bar{a}, R, t_1)), \dots, \mathcal{F}_\ell(B(\bar{a}, R, t_\ell))) \mid \bar{a} \in \pi_{(i_1, \dots, i_m)}(e_1(D)) \}.$$

It is clear that  $\text{RA}_{\text{Aggr}}(\Omega)$  expressions readily define queries over two-sorted schemas. Indeed, if  $e$  is an  $\text{RA}_{\text{Aggr}}(\Omega)$  expression, then the *output* of  $e$  on a two-sorted database  $D$  is  $e(D)$ . We thus may refer to  $e$  as a *query*.

Here is another example, slightly more involved than the ones given above, of expressing an SQL query as an RA query with aggregates and grouping.

### Example 33.9: RA with Aggregates and Grouping

For  $R[A, B, C]$  with  $(R, A) < (R, B) < (R, C)$ , the SQL query

```
SELECT R.A, SUM(R.B), AVG(R.B*R.C)
FROM R
GROUP BY R.A
HAVING SUM(R.B) > AVG(R.C)
```

is translated into the  $\text{RA}_{\text{Aggr}}(\Omega)$  query

$$\pi_{(1,2,4)} \left( \sigma_{2>3} \left( \text{Aggr}_{(1)} [\text{SUM}(2), \text{AVG}(3), \text{AVG}(\text{mult}(2, 3))] (R) \right) \right)$$

assuming  $\Omega$  contains the predicate  $<$ , the numerical function **mult**, and the aggregate functions **SUM** and **AVG** with the obvious meaning. The aggregate expression groups by the first attribute, and computes the sum of the second, and the averages of the third and the product of the second and the third attributes, which then become the second, third,

and fourth attributes. The selection  $\sigma_{2>3}$  enforces the condition in the **HAVING** clause, and the projection outputs the attributes listed in **SELECT**.

## Complexity of RA with Aggregates and Grouping

We proceed to study the complexity of evaluating  $\text{RA}_{\text{Aggr}}(\Omega)$  queries for some set  $\Omega$  of predicates, functions, and aggregates over **Num**. Note that the query evaluation problem for  $\text{RA}_{\text{Aggr}}(\Omega)$  is defined in a slightly different way than the query evaluation problem for the query languages that we have seen so far. In particular, the input database is two-sorted, while the candidate tuple mentions both constants from **Const** and values from **Num**.

### Problem: $\text{RA}_{\text{Aggr}}(\Omega)$ -Evaluation

**Input:** A query  $e$  from  $\text{RA}_{\text{Aggr}}(\Omega)$ , a two-sorted database  $D$ , a tuple  $\bar{a}$  over  $\text{Const} \cup \text{Num}$

**Output:** **true** if  $\bar{a} \in e(D)$ , and **false** otherwise

We can also talk about the data complexity of  $\text{RA}_{\text{Aggr}}(\Omega)$ -Evaluation. As discussed in Chapter 2, when we study the data complexity of query evaluation, we essentially consider the query to be fixed, and only the database and the candidate output are part of the input. Formally, we are interested in the complexity of  $e$ -Evaluation, for an  $\text{RA}_{\text{Aggr}}(\Omega)$  query  $e$ , which takes as input a two-sorted database  $D$  and a tuple  $\bar{a}$  over  $\text{Const} \cup \text{Num}$ , and asks whether  $\bar{a} \in e(D)$ . As usual,  $\text{RA}_{\text{Aggr}}(\Omega)$ -Evaluation is in  $\mathcal{C}$  in data complexity, for some complexity class  $\mathcal{C}$ , if  $e$ -Evaluation is in  $\mathcal{C}$  for every  $\text{RA}_{\text{Aggr}}(\Omega)$  query  $e$ .

Given an  $\text{RA}_{\text{Aggr}}(\Omega)$  query  $e$ , to check whether  $\bar{a} \in e(D)$ , we actually need to compute  $e(D)$ . Indeed, the only way to check if a numerical value equals the output of an aggregate function is to compute the entire bag of values to which the aggregate is applied. It is clear that the complexity of computing  $e(D)$ , and therefore, the complexity of  $\text{RA}_{\text{Aggr}}(\Omega)$ -Evaluation, heavily relies on how complex is to compute predicates, functions, and aggregates from  $\Omega$ . We concentrate on predicates, functions, and aggregates that are easily computable, and show that, although evaluating  $\text{RA}_{\text{Aggr}}(\Omega)$  queries is in general intractable, it becomes tractable when we focus on data complexity.

We say that a  $k$ -ary predicate  $P$  is computable in polynomial time if, for a tuple  $\bar{a} \in \text{Num}^k$ , we can check in polynomial time whether  $P(\bar{a})$  holds. We also say that that a  $k$ -ary function  $f$  is computable in polynomial time if  $f(\bar{a})$  can be computed in polynomial time. Analogously, we can talk about aggregates that are computable in polynomial time. By a simple inspection of each operation of  $\text{RA}_{\text{Aggr}}(\Omega)$ , it is not difficult to show the following result:

**Theorem 33.10**

Consider a set  $\Omega$  of predicates, functions, and aggregates over **Num** that are computable in polynomial time. Then  $\text{RA}_{\text{Aggr}}(\Omega)$ -Evaluation is in EXPTIME, and in PTIME in data complexity.

## Aggregates, Grouping, and Locality

Aggregation and grouping are powerful features that allow us to express interesting counting properties. There are some common queries nonetheless that cannot be expressed even if aggregation and grouping are available. These are the same recursive queries, such as reachability, that were already shown to be inexpressible in FO in Chapter 31. In fact, the tool we shall use to show this is also the same, namely locality. It turns out that adding any amount of aggregate and counting power does not break the locality of FO, i.e., queries can still only see within a fixed distance of their free variables. In particular, in this chapter we are going to show the following inexpressibility result:

### Theorem 34.1

Consider a numerical domain  $\text{Num}$ , and an arbitrary set  $\Omega$  of predicates, functions, and aggregates over  $\text{Num}$ . There is no  $\text{RA}_{\text{Aggr}}(\Omega)$  query  $e$  over the schema  $\mathbf{S} = \{E : (\mathbf{o}, \mathbf{o})\}$  that expresses the reachability query  $q_{\text{reach}}$ .

This will motivate extending query languages with explicit recursive computation, as will be done in Chapters 35 – 38.

### Locality of Two-Sorted Queries

We need to lift the main definitions related to locality, seen in Chapter 31, to queries over databases that include values of the numerical type. This is very routine: the definitions of the Gaifman graph, radius- $r$  balls and neighborhoods, and their isomorphism, do not change at all.

### Definition 34.2: Locality of Queries

A query  $q$  of type  $\tau \in \{\mathbf{o}, \mathbf{n}\}^k$  over a two-sorted schema  $\mathbf{S}$  is  $r$ -local, for  $r \geq 0$ , if, for every database  $D$  of  $\mathbf{S}$ , and every two tuples  $\bar{a}, \bar{b} \in$

$(\text{Const} \cup \text{Num})^k$  such that  $N_r^D(\bar{a})$  and  $N_r^D(\bar{b})$  are isomorphic,  $\bar{a} \in q(D)$  iff  $\bar{b} \in q(D)$ . A query is called *local*, if it is  $r$ -local for some  $r \geq 0$ .

We proceed to show a locality result focussing on schemas that consist of relation names of ordinary type. For brevity, we say that a schema  $\mathbf{S}$  is of ordinary type if, for each  $R : \tau \in \mathbf{S}$ ,  $\tau \in \{\mathbf{o}\}^k$  for some  $k \geq 0$ .

### Theorem 34.3

Consider a numerical domain  $\text{Num}$ , and a set  $\Omega$  of predicates, functions, and aggregates over  $\text{Num}$ . Every  $\text{RA}_{\text{Aggr}}(\Omega)$  query of type  $\{\mathbf{o}\}^k$ , for  $k \geq 0$ , over a schema of ordinary type is local.

What if instead we consider arbitrary two-sorted schemas? In this case there is no known result analogous to Theorem 34.3. In fact, proving an analog of Theorem 34.1 for relations of type  $(\mathbf{n}, \mathbf{n})$  would resolve deep open problems in complexity theory (see Exercise 4.19 for further explanations).

The rather long proof of Theorem 34.3 consists of three main steps:

**FO with Aggregates.** We first present an extension of FO with aggregates, and show that RA with aggregates and grouping translates to it. Note that this logical language is a useful technical tool for showing Theorem 34.3 but it cannot serve as the basis for defining a query language with aggregates; further details are given below.

**Counting Logic.** We then express an RA query with aggregates and grouping, via its translation into the extension of FO with aggregates, in an infinitary counting logic that is easier to analyze mathematically.

**Locality of Counting Logic.** We finally prove the locality of this counting logic. In fact, we shall see that the proof given in Chapter 31 applies with a few minor modifications.

For the sake of clarity, the second and third steps are focussing on queries and logical formulae that do not mention ordinary constants from  $\text{Const}$ . Extending the proof to handle constants is the subject of Exercise 4.15.

## First-Order Logic with Aggregates

Considering the set  $\Omega$  of predicates, functions, and aggregates over the numerical domain  $\text{Num}$ , we are going to define *first-order logic with aggregates*, denoted  $\text{FO}_{\text{Aggr}}(\Omega)$ . Similarly to  $\text{RA}_{\text{Aggr}}(\Omega)$ , the logical language  $\text{FO}_{\text{Aggr}}(\Omega)$  must be *typed*, in particular, each variable should come with its type. We therefore assume that the set  $\text{Var}$  is partitioned into two infinite sets  $\text{Var}_{\mathbf{o}}$  and  $\text{Var}_{\mathbf{n}}$  of variables of type  $\mathbf{o}$  and  $\mathbf{n}$ , respectively. We use  $x, y, z, \dots$  for variables

from  $\text{Var}_o$ , which will be ranging over the set of constants  $\text{Const}$ , and  $i, j, \dots$  for variables from  $\text{Var}_n$ , which will be ranging over the numerical domain  $\text{Num}$ .

#### Definition 34.4: Syntax of FO with Aggregates

Consider a two-sorted schema  $\mathbf{S}$ . We define  $\Omega$ -terms (relative to  $\mathbf{S}$ ) with their associated *types*, and  $\text{FO}_{\text{Aggr}}(\Omega)$  formulae over  $\mathbf{S}$ , by mutual induction as follows:

##### $\Omega$ -terms

- Each constant of  $\text{Const}$  and variable of  $\text{Var}_o$  is an  $\Omega$ -term of type  $o$ .
- Each value  $a \in \text{Num}$  is an  $\Omega$ -term of type  $n$ , whose set of free variables  $\text{FV}(a)$  is empty.
- Each variable  $i \in \text{Var}_n$  is an  $\Omega$ -term of type  $n$ , whose set of free variables  $\text{FV}(i)$  is  $\{i\}$ .
- If  $f$  is an  $m$ -ary numerical function from  $\Omega$ , and  $t_1, \dots, t_m$  are  $\Omega$ -terms of type  $n$ , then  $f(t_1, \dots, t_m)$  is an  $\Omega$ -term of type  $n$ , whose set of free variables  $\text{FV}(f(t_1, \dots, t_m))$  is  $\text{FV}(t_1) \cup \dots \cup \text{FV}(t_m)$ .
- If  $\varphi$  is an  $\text{FO}_{\text{Aggr}}(\Omega)$  formula over  $\mathbf{S}$ ,  $t$  is an  $\Omega$ -term of type  $n$ , and  $\mathcal{F}$  is an aggregate of  $\Omega$ , then

$$\text{Aggr}_{\mathcal{F}}(\bar{u})(\varphi, t)$$

where  $\bar{u} = (u_1, \dots, u_k)$  is a tuple of variables over  $\text{Var}$  such that  $\{u_1, \dots, u_k\} \subseteq \text{FV}(t)$  and  $\text{FV}(t) \subseteq \text{FV}(\varphi)$ , is an  $\Omega$ -term of type  $n$ , whose set of free variables  $\text{FV}(\text{Aggr}_{\mathcal{F}}(\bar{u})(\varphi, t))$  is  $\text{FV}(\varphi) - \{u_1, \dots, u_k\}$ .

##### $\text{FO}_{\text{Aggr}}(\Omega)$ Formulae

- If  $a \in \text{Const}$ , and  $x \in \text{Var}_o$ , then  $x = a$  is an atomic formula, whose set of variables  $\text{FV}(x = a)$  is  $\{x\}$ .
- If  $x, y \in \text{Var}_o$ , then  $x = y$  is an atomic formula, whose set of free variables  $\text{FV}(x = y)$  is  $\{x, y\}$ .
- If  $t$  is an  $\Omega$ -term of type  $n$ , and  $i \in \text{Var}_n$ , then  $i = t$  is an atomic formula, whose set of free variables  $\text{FV}(i = t)$  is  $\{i\} \cup \text{FV}(t)$ .
- If  $u_1, \dots, u_k$  are  $\Omega$ -terms (not necessarily distinct) from  $\text{Const} \cup \text{Num} \cup \text{Var}$ , where  $u_i$  is of type  $\tau_i$  for each  $i \in [k]$ , and  $R : (\tau_1, \dots, \tau_k)$  belongs to  $\mathbf{S}$ , then  $R(u_1, \dots, u_k)$  is an atomic formula, whose set of free variables  $\text{FV}(R(u_1, \dots, u_k))$  is  $\{u_1, \dots, u_k\} \cap \text{Var}$ .
- If  $u_1, \dots, u_k$  are  $\Omega$ -terms (not necessarily distinct) from  $\text{Num} \cup \text{Var}_n$ , and  $P$  is a  $k$ -ary numerical predicate from  $\Omega$ , then  $P(u_1, \dots, u_k)$  is



an atomic formula, whose set of free variables  $FV(P(u_1, \dots, u_k))$  is  $\{u_1, \dots, u_k\} \cap \text{Var}_n$ .

- If  $\varphi_1$  and  $\varphi_2$  are formulae, then  $(\varphi_1 \wedge \varphi_2)$  and  $(\varphi_1 \vee \varphi_2)$  are formulae, whose set of free variables  $FV(\varphi_1 \wedge \varphi_2) = FV(\varphi_1 \vee \varphi_2)$  is  $FV(\varphi_1) \cup FV(\varphi_2)$ .
- If  $\varphi$  is a formula, then  $(\neg\varphi)$  is a formula, whose set of free variables  $FV(\neg\varphi)$  is  $FV(\varphi)$ .
- If  $\varphi$  is a formula and  $u \in \text{Var}$ , then  $(\exists u \varphi)$  and  $(\forall u \varphi)$  are formulae, whose set of free variables  $FV(\exists u \varphi) = FV(\forall u \varphi)$  is  $FV(\varphi) - \{u\}$ .

We will omit the outermost brackets of  $\text{FO}_{\text{Aggr}}(\Omega)$  formulae. To define the semantics of  $\text{FO}_{\text{Aggr}}(\Omega)$ , we need the notion of assignment for  $\Omega$ -terms and formulae. Given an  $\Omega$ -term  $t$  of type  $n$ , an *assignment*  $\eta$  for  $t$  is a function from  $FV(t)$  to  $\text{Const} \cup \text{Num}$  such that  $\eta(u) \in \text{Const}$  if  $u \in \text{Var}_o$ , and  $\eta(u) \in \text{Num}$  if  $u \in \text{Var}_n$ . Similarly, given an  $\text{FO}_{\text{Aggr}}(\Omega)$  formula  $\varphi$ , and a database  $D$ , an *assignment*  $\eta$  for  $\varphi$  over  $D$  is a function from  $FV(\varphi)$  to  $\text{Dom}(D) \cup \text{Dom}(\varphi) \cup \text{Num}$ , where  $\text{Dom}(\varphi)$  is the set of constants and numerical values occurring in  $\varphi$ , such that  $\eta(u) \in \text{Const}$  if  $u \in \text{Var}_o$ , and  $\eta(u) \in \text{Num}$  if  $u \in \text{Var}_n$ . We write  $\eta[u/a]$ , for  $u \in \text{Var}$  and  $a \in \text{Const} \cup \text{Num}$ , for the assignment that modifies  $\eta$  by setting  $\eta(u) = a$ . To avoid heavy notation, we extend an assignment to be the identity on  $\text{Const} \cup \text{Num}$ . The semantics of  $\text{FO}_{\text{Aggr}}(\Omega)$  follows.

**Definition 34.5: Semantics of FO with Aggregates**

Consider a two-sorted schema  $\mathbf{S}$ . We define the value of  $\Omega$ -terms (relative to  $\mathbf{S}$ ) of type  $\mathbf{n}$ , and the satisfaction of  $\text{FO}_{\text{Aggr}}(\Omega)$  formulae over  $\mathbf{S}$ , by mutual induction as follows.

**Value of  $\Omega$ -terms**

Let  $t$  be an  $\Omega$ -term of type  $\mathbf{n}$  and  $\eta$  an assignment for  $t$ . The *value* of  $t$  in a database  $D$  of  $\mathbf{S}$  under  $\eta$ , denoted  $t^{D,\eta}$ , is defined as follows:

- If  $t = a$  with  $a \in \text{Num}$ , then  $t^{D,\eta} = a$ .
- If  $t = \iota$  with  $\iota \in \text{Var}_{\mathbf{n}}$ , then  $t^{D,\eta} = \eta(\iota)$ .
- If  $t = f(t_1, \dots, t_m)$ , then  $t^{D,\eta} = f(t_1^{D,\eta}, \dots, t_m^{D,\eta})$ .
- If  $t = \text{Aggr}_{\mathcal{F}}(\bar{u})(\varphi, t_0)$ , with  $H$  being the set of all assignments  $\eta'$  for  $\varphi$  over  $D$  that agree with  $\eta$  on  $\text{FV}(\varphi)$  and satisfy that  $(D, \eta') \models \varphi$ , then  $t^{D,\eta} = \mathcal{F}(\{\bar{u}\})$  if  $H$  is infinite, otherwise  $t^{D,\eta} = \mathcal{F}(\{\bar{u}^{D,\eta'} \mid \eta' \in H\})$ .

**Satisfaction of  $\text{FO}_{\text{Aggr}}(\Omega)$  Formulae**

Let  $\varphi$  be a formula over  $\mathbf{S}$ , and  $\eta$  an assignment for  $\varphi$ . We define when  $\varphi$  is *satisfied* in a database  $D$  of  $\mathbf{S}$  under  $\eta$ , written  $(D, \eta) \models \varphi$ , as follows (omitting Boolean connectives, which are defined in the standard way):

- If  $\varphi$  is  $x = a$ , then  $(D, \eta) \models \varphi$  if  $\eta(x) = a$ .
- If  $\varphi$  is  $x = y$ , then  $(D, \eta) \models \varphi$  if  $\eta(x) = \eta(y)$ .
- If  $\varphi$  is  $\iota = t$ , then  $(D, \eta) \models \varphi$  if  $\eta(\iota) = t^{D,\eta}$ .
- If  $\varphi$  is  $R(u_1, \dots, u_k)$ , for  $R$  a  $k$ -ary relation symbol from  $\mathbf{S}$ , then  $(D, \eta) \models \varphi$  if  $R(\eta(u_1), \dots, \eta(u_k)) \in D$ .
- If  $\varphi$  is  $P(u_1, \dots, u_k)$ , for  $P$  a  $k$ -ary numerical predicate from  $\Omega$ , then  $(D, \eta) \models \varphi$  if  $(\eta(u_1), \dots, \eta(u_k)) \in P$ .
- If  $\varphi = \exists u \psi$ , then  $(D, \eta) \models \varphi$  if  $(D, \eta[u/a]) \models \psi$  for *some*  $a \in (\text{Dom}(D) \cup \text{Dom}(\varphi)) \cap \text{Const}$  if  $u \in \text{Var}_{\mathbf{o}}$ , and  $a \in \text{Num}$  if  $u \in \text{Var}_{\mathbf{n}}$ .
- If  $\varphi = \forall x \psi$ , then  $(D, \eta) \models \varphi$  if  $(D, \eta[x/a]) \models \psi$  for *each*  $a \in \text{Dom}(D) \cup (\text{Dom}(\varphi) \cap \text{Const})$  if  $u \in \text{Var}_{\mathbf{o}}$ , and  $a \in \text{Num}$  if  $u \in \text{Var}_{\mathbf{n}}$ .

We can now observe a crucial difference between first-order logic with aggregates, and first-order logic as defined in Chapter 3. Given an  $\text{FO}_{\text{Aggr}}(\Omega)$  formula  $\varphi$ , we may have infinitely many assignments  $\eta$  for  $\varphi$  over a database  $D$  such that  $(D, \eta) \models \varphi$ . Consider, for example, the formula  $\varphi = R(x) \wedge (j = \iota + 1)$ ,

and the database  $D = \{R(a)\}$  with  $a \in \text{Const}$ . Assuming that  $\text{Num}$  is the set of integers, any of the infinitely many assignments  $\eta : \{x, i, j\} \rightarrow \{a\} \cup \text{Num}$  for which  $\eta(x) = a$  and  $\eta(j) = \eta(i) + 1$  satisfies  $(D, \eta) \models \varphi$ . This implies that first-order logic with aggregates cannot be used to define a query language as we did with ordinary FO in Chapter 3. Nevertheless, given an expression  $\varphi(\bar{u})$ , where  $\varphi$  is an  $\text{FO}_{\text{Aggr}}(\Omega)$  formula, and  $\bar{u}$  is a tuple over  $\text{FV}(\varphi)$  such that each free variable of  $\varphi$  occurs in  $\bar{u}$  at least once, we can define the *output* of  $\varphi(\bar{u})$  on a database  $D$ , denoted  $\varphi(\bar{u})(D)$ , in the obvious way, but we cannot call  $\varphi(\bar{u})$  a query since  $\varphi(\bar{u})(D)$  may be infinite. One can still define a logic with aggregates that can in turn be used to define a query language, but the syntax is much more cumbersome.

We now show that every RA query with aggregates and grouping can be expressed via FO with aggregates.

#### Proposition 34.6

Consider an  $\text{RA}_{\text{Aggr}}(\Omega)$  query  $e$  over a two-sorted schema  $\mathbf{S}$ . There exists an  $\text{FO}_{\text{Aggr}}(\Omega)$  formula  $\varphi_e$  over  $\mathbf{S}$ , and a tuple  $\bar{u}_e$  over  $\text{FV}(\varphi_e)$  that mentions all the variables of  $\text{FV}(\varphi_e)$ , such that  $e(D) = \varphi_e(\bar{u}_e)(D)$ , for every database  $D$  of  $\mathbf{S}$ .

*Proof.* The proof is by induction on the structure of  $e$ . Most of the cases are treated in the same way as in the proof of Theorem 6.1, in particular, the proof that every RA query can be equivalently written as an FO query. We proceed to discuss the two new cases, namely generalized projection and grouping with aggregation. For an  $(\Omega, \tau)$ -term  $t$ , where  $\tau = (\tau_1, \dots, \tau_k) \in \{\mathbf{o}, \mathbf{n}\}^k$ , and a tuple of variables  $\bar{u} = (u_1, \dots, u_k)$  with  $u_i$  being of type  $\tau_i$ , for each  $i \in [k]$ , we write  $t(\bar{u})$  for the  $\Omega$ -term obtained from  $t$  by replacing each  $i \in \text{VAR}(t)$  with  $u_i$ . We are now ready to proceed with the translations:

- Assume first that  $e$  is  $\pi_{(t_1, \dots, t_m)}(e')$ . By induction hypothesis, there exists an  $\text{FO}_{\text{Aggr}}(\Omega)$  formula  $\varphi_{e'}$ , and a tuple  $\bar{u}_{e'} = (u_1, \dots, u_k)$  over  $\text{FV}(\varphi_{e'})$ , such that  $\varphi_{e'}(\bar{u}_{e'})$  expresses  $e'$ . We define the  $\text{FO}_{\text{Aggr}}(\Omega)$  formula

$$\varphi_e = \exists u_1 \cdots \exists u_k \left( \varphi_{e'} \wedge \bigwedge_{i=1}^m v_i = t_i(\bar{u}_{e'}) \right),$$

and the tuple  $\bar{u}_e = (v_1, \dots, v_m)$  over  $\text{FV}(\varphi_e)$ .

- Assume now that  $e$  is  $\text{Aggr}_{(i_1, \dots, i_m)}[\mathcal{F}_1(t_1), \dots, \mathcal{F}_\ell(t_\ell)](e')$ ; for the sake of clarity, we assume that  $(i_1, \dots, i_m) = (1, \dots, m)$ , but the same construction can be applied to any list of integers. By induction hypothesis, there exists an  $\text{FO}_{\text{Aggr}}(\Omega)$  formula  $\varphi_{e'}$ , and a tuple  $\bar{u}_{e'} = (u_1, \dots, u_k)$  over  $\text{FV}(\varphi_{e'})$ , such that  $\varphi_{e'}(\bar{u}_{e'})$  expresses  $e'$ . We define the  $\text{FO}_{\text{Aggr}}(\Omega)$  formula

$$\varphi_e = \exists w_{m+1} \cdots \exists w_k \left( \psi_{e'} \wedge \bigwedge_{i=1}^{\ell} v_i = \text{Aggr}_{\mathcal{F}_i}(w_{m+1}, \dots, w_k) (\psi_{e'}, t_i(u_1, \dots, u_m, w_{m+1}, \dots, w_k)) \right),$$

where  $\psi_{e'}$  is obtained from  $\varphi_{e'}$  by replacing  $u_i$  with  $w_i$  for each  $i \in [m+1, k]$ , and the tuple  $\bar{u}_e = (u_1, \dots, u_m, v_1, \dots, v_\ell)$  over  $\text{FV}(\varphi_e)$ .

It is easy to verify the correctness of the above translations.  $\square$

## An Infinitary Counting Logic

We now proceed with the second step of the proof of Theorem 34.1, where the goal is to express an  $\text{RA}_{\text{Aggr}}(\Omega)$  query, by exploiting Proposition 34.6, into a convenient infinitary counting logic. We proceed to introduce the counting logic  $\mathbf{L}_C$ , and its sublogic  $\mathbb{L}_C$ , which we prove to be equivalent to  $\mathbf{L}_C$ . Recall that, for the sake of clarity, in this step we focus on queries and formulae that do not mention constants from  $\text{Const}$ . In what follows, every value of  $\text{Num}$ , and every variable of  $\text{Var}$ , is a term of the respective type, and of rank 0.

### Definition 34.7: An Infinitary Counting Logic

Consider a two-sorted schema  $\mathbf{S}$ . We define *formulae of  $\mathbf{L}_C$  over  $\mathbf{S}$* , and their associated *rank*, by induction as follows:

- If  $x, y \in \text{Var}_o$ , then  $x = y$  is an atomic formula with  $\text{rank}(x = y) = 0$ .
- If  $a \in \text{Num}$ , and  $\iota \in \text{Var}_n$ , then  $\iota = a$  is an atomic formula with  $\text{rank}(\iota = a) = 0$ .
- If  $\iota, j \in \text{Var}_n$ , then  $\iota = j$  is an atomic formula with  $\text{rank}(\iota = j) = 0$ .
- If  $u_1, \dots, u_k$  are terms (not necessarily distinct), where  $u_i$  is of type  $\tau_i$  for each  $i \in [k]$ , and  $R : (\tau_1, \dots, \tau_k)$  belongs to  $\mathbf{S}$ , then  $R(u_1, \dots, u_k)$  is an atomic formula with  $\text{rank}(R(u_1, \dots, u_k)) = 0$ .
- If  $\Phi$  is a (possibly infinite) set of formulae, and  $k = \sup_{\varphi \in \Phi} \text{rank}(\varphi)$  is finite, then  $\psi = \left( \bigwedge_{\varphi \in \Phi} \varphi \right)$  and  $\psi' = \left( \bigvee_{\varphi \in \Phi} \varphi \right)$  are formulae with  $\text{rank}(\psi) = \text{rank}(\psi') = k$ .
- If  $\varphi$  is a formula, then  $(\neg \varphi)$  is a formula with  $\text{rank}(\neg \varphi) = \text{rank}(\varphi)$ .
- If  $\varphi$  is a formula,  $\bar{u} = (u_1, \dots, u_k)$  is a tuple over  $\text{Var}$ , and  $n \in \mathbb{N}$ , then  $\psi = (\exists^{\geq n} \bar{u} \varphi)$  is a formula with  $\text{rank}(\psi) = \text{rank}(\varphi) + k$ .

We further define *formulae of  $\mathbb{L}_C$  over  $\mathbf{S}$* , and their associated *rank*, in the same way as  $\mathbf{L}_C$  formulae, with the difference that only quantification of the form  $\exists^{\geq n} x$ , where  $x$  is a single variable type  $o$ , is allowed:

- If  $\varphi$  is a formula of  $\mathbb{L}_{\mathbf{C}}$ ,  $x \in \text{Var}_{\mathbf{o}}$ , and  $n \in \mathbb{N}$ , then  $\psi = (\exists^{\geq n} x \varphi)$  is a formula of  $\mathbb{L}_{\mathbf{C}}$  with  $\text{rank}(\psi) = \text{rank}(\varphi) + 1$ .

Let us stress that  $\mathbf{L}_{\mathbf{C}}$ , and thus  $\mathbb{L}_{\mathbf{C}}$ , consists of formulae of finite rank since in the definition of the infinitary conjunctions and disjunctions we explicitly restrict the rank to be finite; otherwise, we may get formulae of infinite rank, e.g., if  $\Phi = \{\varphi_1, \varphi_2, \dots\}$  with  $\text{rank}(\varphi_i) = i$  for each  $i > 0$ , then  $\text{rank}\left(\bigvee_{\varphi \in \Phi} \varphi\right)$  is infinite. The set of free variables of an  $\mathbf{L}_{\mathbf{C}}$  formula  $\varphi$ , denoted  $\text{FV}(\varphi)$ , as well as the semantics of  $\mathbf{L}_{\mathbf{C}}$ , are defined in the expected way. Let us only discuss the details in the case of a formula of the form  $\psi = \exists^{\geq n} \bar{u} \varphi$ , which essentially states that there exist at least  $n$  witnesses for  $\bar{u}$ . For an assignment  $\eta$  for  $\psi$  over a database  $D$ , i.e., a function that maps  $\text{FV}(\psi)$  to  $\text{Dom}(D) \cup \text{Num}$ ,  $\psi$  is satisfied in  $D$  under  $\eta$ , written  $(D, \eta) \models \psi$ , if there are at least  $n$  assignments  $\eta'$  for  $\varphi$  over  $D$  that agree with  $\eta$  on  $\text{FV}(\psi)$  such that  $(D, \eta') \models \varphi$ . We can use the shorthand  $\exists^{\geq n} \bar{u} \varphi$  to say that there are *exactly*  $n$  such assignments for  $\varphi$  over  $D$ , that is,  $\exists^{\geq n} \bar{u} \varphi \wedge \neg \exists^{\geq n+1} \bar{u} \varphi$ , which does not alter the rank. Given an expression  $\varphi(\bar{u})$ , where  $\varphi$  is an  $\mathbf{L}_{\mathbf{C}}$  formula, and  $\bar{u}$  is a tuple over  $\text{FV}(\varphi)$  such that each free variable of  $\varphi$  occurs in  $\bar{u}$  at least once, we can define the output of  $\varphi(\bar{u})$  on a database  $D$ , denoted  $\varphi(\bar{u})(D)$ , in the obvious way.

#### Proposition 34.8

Consider an  $\text{FO}_{\text{Aggr}}(\Omega)$  formula  $\varphi$  over a two-sorted schema  $\mathbf{S}$ , and a tuple  $\bar{u}$  over  $\text{FV}(\varphi)$  that mentions all the variables of  $\text{FV}(\varphi)$ . There exists an  $\mathbb{L}_{\mathbf{C}}$  formula  $\psi$  over  $\mathbf{S}$ , with  $\text{FV}(\varphi) = \text{FV}(\psi)$ , such that  $\varphi(\bar{u})(D) = \psi(\bar{u})(D)$ , for every database  $D$  of  $\mathbf{S}$ .

*Proof.* We first show the statement for the counting logic  $\mathbf{L}_{\mathbf{C}}$ .

**Lemma 34.9.** *There is an  $\mathbf{L}_{\mathbf{C}}$  formula  $\varphi^{\diamond}$  over  $\mathbf{S}$ , with  $\text{FV}(\varphi) = \text{FV}(\varphi^{\diamond})$ , such that  $\varphi(\bar{u})(D) = \varphi^{\diamond}(\bar{u})(D)$ , for every database  $D$  of  $\mathbf{S}$ .*

*Proof.* We translate every  $\Omega$ -term  $t$  of type  $\mathbf{n}$  occurring in  $\varphi$  into an  $\mathbf{L}_{\mathbf{C}}$  formula  $\alpha_t^i$ , where  $i$  is a distinguished free variable of  $\alpha_t^i$ . Intuitively,  $\alpha_t^i$  states that  $i$  is the value of  $t$ , that is,  $(D, \eta) \models \alpha_t^i$  iff  $t^{D, \eta} = \eta(i)$ . We further translate the formula  $\varphi$  into an  $\mathbf{L}_{\mathbf{C}}$  formula  $\varphi^{\diamond}$  with  $\text{FV}(\varphi) = \text{FV}(\varphi^{\diamond})$ .

To ensure that  $\varphi^{\diamond}$  is indeed an  $\mathbf{L}_{\mathbf{C}}$  formula, we need to show that it has finite rank. To this end, we first need to transfer the notion of rank to  $\Omega$ -terms and  $\text{FO}_{\text{Aggr}}(\Omega)$  formulae by mutual induction. Let  $t$  be an  $\Omega$ -term:

- If  $t \in \text{Num} \cup \text{Var}$ , then  $\text{rank}(t) = 0$ .
- If  $t = f(t_1, \dots, t_m)$ , then  $\text{rank}(t) = \max_{i \in [m]} \{\text{rank}(t_i)\}$ .
- If  $t = \text{Aggr}_{\mathcal{F}}(\bar{v})(\varphi', t')$ , then  $\text{rank}(t) = \max\{\text{rank}(\varphi'), \text{rank}(t')\} + k$ , where  $k$  is the arity of the tuple  $\bar{v}$ .

Consider now an  $\text{FO}_{\text{Aggr}}(\Omega)$  formula  $\varphi'$ :

- If  $\varphi'$  is  $x = y$ , then  $\text{rank}(\varphi') = 0$ .
- If  $\varphi'$  is  $\iota = t$ , then  $\text{rank}(\varphi') = \text{rank}(t)$ .
- If  $\varphi'$  is  $R(v_1, \dots, v_k)$ , then  $\text{rank}(\varphi') = \max_{i \in [k]} \{\text{rank}(v_i)\}$ .
- If  $\varphi'$  is  $\varphi_1 \wedge \varphi_2$  or  $\varphi_1 \vee \varphi_2$ , then  $\text{rank}(\varphi') = \max\{\text{rank}(\varphi_1), \text{rank}(\varphi_2)\}$ .
- If  $\varphi'$  is  $\neg \varphi_1$ , then  $\text{rank}(\varphi') = \text{rank}(\varphi_1)$ .
- If  $\varphi'$  is  $\exists v \varphi_1$  or  $\forall v \varphi_1$ , then  $\text{rank}(\varphi') = \text{rank}(\varphi_1) + 1$ .

We are now ready to provide the translation of an  $\Omega$ -term  $t$  of type  $n$  occurring in  $\varphi$ , and of  $\varphi$  itself, into  $\mathbf{L}_{\mathbf{C}}$  by mutual induction. In what follows, given an  $\mathbf{L}_{\mathbf{C}}$  formula  $\chi$ , a variable  $\iota \in \text{FV}(\chi)$ , and a value  $a \in \text{Num}$ , we write  $\chi[\iota/a]$  for the formula obtained from  $\chi$  after replacing  $\iota$  with  $a$ .

*Translation of an  $\Omega$ -term  $t$  occurring in  $\varphi$  into an  $\mathbf{L}_{\mathbf{C}}$  formula  $\alpha_t^i$*

- If  $t = a$  with  $a \in \text{Num}$ , then  $\alpha_t^i = (\iota = a)$  of rank 0.
- If  $t = \iota$  with  $\iota \in \text{Var}_n$ , then  $\alpha_t^i = (\iota = \iota)$  of rank 0.
- If  $t = f(t_1, \dots, t_m)$ , then

$$\alpha_t^i = \bigvee_{\substack{(a_1, \dots, a_m, a_{m+1}) \in \text{Num}^{m+1} : \\ f(a_1, \dots, a_m) = a_{m+1}}} \left( \bigwedge_{i \in [m]} \alpha_{t_i}^{j_i}[\iota_i/a_i] \rightarrow \iota = a_{m+1} \right)$$

of rank  $\max_{i \in [m]} \{\text{rank}(\alpha_{t_i}^{j_i})\} \leq \text{rank}(t)$ , where  $\iota$  is a new numerical variable not occurring in  $\alpha_{t_i}^{j_i}$ , for each  $i \in [m]$ .

- If  $t = \text{Aggr}_{\mathcal{F}}(\bar{v})(\varphi', t')$ , with  $\mathcal{B}$  being the set of all bags (finite or infinite) over  $\text{Num}$ , then

$$\alpha_t^i = \bigvee_{B \in \mathcal{B}} (\chi_B \wedge \zeta_B \wedge \iota = \mathcal{F}(B)),$$

where  $\iota$  is a new numerical variable not occurring in  $\chi_B$  and  $\zeta_B$ , and  $\chi_B$  and  $\zeta_B$  are defined as follows. For  $B \in \mathcal{B}$ , we write  $\text{supp}(B)$  for its *support*, i.e., the set of elements that appear in it, and  $\sharp(a, B)$  for the number of occurrences of  $a$  in  $B$ . We then define

$$\chi_B = \bigwedge_{a \in \text{supp}(B)} \exists^{\sharp(a, B)} \bar{v} \left( (\varphi')^\diamond \wedge \alpha_{t'}^j[\iota/a] \right)$$

stating that the values of  $t'$ , as  $\bar{v}$  ranges over tuples satisfying  $\varphi'$ , have exactly the same multiplicities as in  $B$ , and

$$\zeta_B = \forall \bar{v} \left( (\varphi')^\diamond \rightarrow \bigvee_{a \in \text{supp}(B)} \alpha_{t'}^j[j/a] \right)$$

stating that only elements of  $B$  are values of  $t'$  as  $\bar{v}$  ranges over tuples satisfying  $\varphi'$ . It is easy to verify that  $\alpha_t^j$  is of rank  $\max\{\text{rank}((\varphi')^\diamond), \text{rank}(\alpha_{t'}^j)\} + k \leq \max\{\text{rank}(\varphi'), \text{rank}(t')\} + k = \text{rank}(t)$ , where  $k$  is the arity of the tuple  $\bar{v}$ . Moreover, it should be clear that  $\alpha_t^j$  essentially states that, for some bag  $B \in \mathcal{B}$ , the values of  $t'$  form exactly  $B$ , and the value of  $t$  is the value of the aggregate  $\mathcal{F}$  on  $B$ .

*Translation of the formula  $\varphi$  into an  $\mathbf{L}_C$  formula  $\varphi^\diamond$*

- If  $\varphi$  is the atomic formula  $x = y$  or  $R(\bar{v})$ , then  $\varphi^\diamond$  is precisely  $\varphi$ , and thus,  $\text{rank}(\varphi^\diamond) = \text{rank}(\varphi)$ .
- If  $\varphi$  is the atomic formula  $\iota = t$ , then  $\varphi^\diamond$  is  $\alpha_t^i$  with  $\text{rank}(\varphi^\diamond) \leq \text{rank}(t)$ .
- If  $\varphi$  is  $P(v_1, \dots, v_k)$ , then  $\varphi^\diamond$  is

$$\bigvee_{(a_1, \dots, a_k) \in P} \left( \bigwedge_{i \in [k]} \alpha_{v_i}^{j_i}[j_i/a_i] \right)$$

with  $\text{rank}(\varphi^\diamond) \leq \max_{i \in [m]} \{\text{rank}(\alpha_{v_i}^{j_i})\} \leq \text{rank}(\varphi)$ .

- If  $\varphi$  is  $\varphi_1 \wedge \varphi_2$ ,  $\varphi_1 \vee \varphi_2$ ,  $\neg \varphi_1$ , then  $\varphi^\diamond$  is  $\varphi_1^\diamond \wedge \varphi_2^\diamond$ ,  $\varphi_1^\diamond \vee \varphi_2^\diamond$ ,  $\neg \varphi_1^\diamond$ , respectively, of with  $\text{rank}(\varphi^\diamond) = \text{rank}(\varphi)$ .
- If  $\varphi$  is  $\exists v \varphi_1$ ,  $\forall v \varphi_1$ , then  $\varphi^\diamond$  is  $\exists v \varphi_1^\diamond$ ,  $\neg \exists v \neg \varphi_1^\diamond$ , respectively, with  $\text{rank}(\varphi^\diamond) = \text{rank}(\varphi)$ .

This completes the translation of  $\Omega$ -terms of type  $\mathbf{n}$  occurring in  $\varphi$ , and of  $\varphi$  itself. It can be verified that  $\varphi^\diamond$  is indeed an  $\mathbf{L}_C$  formula, with  $\text{FV}(\varphi) = \text{FV}(\varphi^\diamond)$ , such that  $\varphi(\bar{u})(D) = \varphi^\diamond(\bar{u})(D)$ , for every database  $D$  of  $\mathbf{S}$ .  $\square$

We now proceed to show that the  $\mathbf{L}_C$  formula  $\varphi^\diamond$  provided by Lemma 34.9 can be converted into an  $\mathbb{L}_C$  formula  $\psi$  such that  $\varphi^\diamond(\bar{u})$  and  $\psi(\bar{u})$  have the same output on every database of  $\mathbf{S}$ , which will prove Proposition 34.8

**Lemma 34.10.** *There exists an  $\mathbb{L}_C$  formula  $\psi$  over  $\mathbf{S}$ , with  $\text{FV}(\varphi^\diamond) = \text{FV}(\psi)$ , such that  $\varphi^\diamond(\bar{u})(D) = \psi(\bar{u})(D)$ , for every database  $D$  of  $\mathbf{S}$ .*

*Proof.* To prove the claim, we need to replace quantifiers of the form  $\exists^{\geq n} \bar{v} \psi'$  in  $\varphi^\diamond$  with  $\exists^{\geq n} x \psi'$ , where  $x \in \text{Var}_o$ , without increasing the rank. We explain how this is done when  $\bar{v}$  is binary; the general proof is then by induction on the arity of  $\bar{v}$ , using the case when  $\bar{v}$  is binary as the base step.

Consider a subformula  $\exists^{\geq n}(v, w) \psi'$  of  $\varphi^\diamond$ , where  $v, w \in \text{FV}(\psi')$ . The idea of replacing this by simpler quantifiers is to say that there are at least  $k_1$   $v$ 's for which there exist exactly  $\ell_1$   $w$ 's satisfying  $\psi'$ , and there are exactly  $k_2$   $v$ 's

for which there exist exactly  $\ell_2$   $w$ 's satisfying  $\psi'$ , and so on, with all the  $\ell_i$ 's being distinct in order to ensure that the same pair of values is never counted twice. Formally, a finite set of pairs of integers  $\{(k_1, \ell_1), \dots, (k_s, \ell_s)\}$  is an  $n$ -witness if  $\sum_{i=1}^s k_i \cdot \ell_i \geq n$ , and  $\ell_i \neq \ell_j$  for each  $i, j \in [s]$  with  $i \neq j$ . Let  $\mathcal{W}_n$  be the set of all  $n$ -witnesses, which is clearly infinite. Then,  $\exists^{\geq n}(v, w) \psi'$  is replaced by the infinitary disjunction

$$\bigvee_{\{(k_1, \ell_1), \dots, (k_s, \ell_s)\} \in \mathcal{W}_n} \left( \bigwedge_{i=1}^s \exists^{\geq k_i} v \exists^{\ell_i} w \psi' \right)$$

whose rank is  $\text{rank}(\exists^{\geq n}(v, w) \psi') = \text{rank}(\psi') + 2$ .

According to the definition of  $\mathbb{L}_{\mathbf{C}}$ , we can only quantify variables of type  $\mathbf{o}$ . Thus, we need to eliminate from the above infinitary disjunction the quantifiers over numerical variables. This is done by using infinitary disjunctions as follows: a formula of the form  $\exists^{\geq n} \iota \psi''$ , where  $\iota \in \text{FV}(\psi'')$ , is written as

$$\bigvee_{A \subseteq \text{Num} : |A| \geq n} \left( \bigwedge_{a \in A} \psi''[\iota/a] \right)$$

whose rank is  $\text{rank}(\exists^{\geq n} \iota \psi'') - 1$ . This completes the proof of Lemma 34.10.  $\square$

By Lemma 34.9 and 34.10, we get an  $\mathbb{L}_{\mathbf{C}}$  formula  $\psi$ , with  $\text{FV}(\varphi) = \text{FV}(\psi)$ , such that  $\varphi(\bar{u})(D) = \psi(\bar{u})(D)$ , for every database  $D$  of  $\mathbf{S}$ , as needed.  $\square$

## Locality of Counting Logic

The last step of the proof of Theorem 34.3 is proving the locality of  $\mathbb{L}_{\mathbf{C}}$ .

### Proposition 34.11

Consider an  $\mathbb{L}_{\mathbf{C}}$  formula  $\varphi$  over a schema  $\mathbf{S} = \{R_1 : \tau_1, \dots, R_n : \tau_n\}$  with  $\tau_i \in \{\mathbf{o}\}^{k_i}$  and  $k_i \geq 0$ , for each  $i \in [n]$ , of rank  $k$  such that  $\text{FV}(\varphi) \subseteq \text{Var}_{\mathbf{o}}$ , and a tuple  $\bar{x}$  over  $\text{FV}(\varphi)$  that mentions all the variables of  $\text{FV}(\varphi)$ . Then  $\varphi(\bar{x})$  is  $(3^k - 1)/2$ -local.

*Proof.* We proceed exactly as in the proof of Theorem 31.3, using rank in place of quantifier rank. The cases of atomic formulae are identical. For infinitary connectives  $\bigvee_{\varphi \in \Phi} \varphi$  and  $\bigwedge_{\varphi \in \Phi} \varphi$ , we note that  $k$  is the bound on the rank of all formulae  $\varphi \in \Phi$ , and thus, the same proof as in Theorem 31.3 applies. The case of negation does not change either.

The last step is to consider formulae of the form  $\varphi(\bar{x}) = \exists^{\geq n} y \psi$ , with  $\text{rank}(\varphi) = k$ . Thus,  $\text{rank}(\psi) = k - 1$ , and by the induction hypothesis,  $\psi(\bar{x}, y)$  is  $r$ -local, where  $r = (3^{k-1} - 1)/2$ . As in Theorem 31.3, we must show that



$\varphi$  is  $(3r + 1)$ -local. The proof is the same as the one given before. Indeed, recall that we established that, given a database  $D$  of  $\mathbf{S}$ , with tuples  $\bar{a}, \bar{b}$  over  $\text{Const}$  such that  $N_{3r+1}^D(\bar{a})$  and  $N_{3r+1}^D(\bar{b})$  are isomorphic, we have a bijection  $f : \text{Dom}(D) \rightarrow \text{Dom}(D)$  such that  $N_r^D(\bar{a}, c)$  and  $N_r^D(\bar{b}, f(c))$  are isomorphic, for each  $c \in \text{Dom}(D)$ . By  $r$ -locality of  $\psi$ , we get that the number of witnesses to  $\psi(\bar{a}, y)$  and  $\psi(\bar{b}, y)$  are the same, and thus  $\varphi(\bar{a})$  and  $\varphi(\bar{b})$  are either simultaneously true, or simultaneously false.  $\square$

We can now finalize the proof of Theorem 34.3. Given an  $\text{RA}_{\text{Aggr}}(\Omega)$  query  $e$  of type  $\{\mathbf{o}\}^k$ , for  $k \geq 0$ , over a schema  $\mathbf{S}$  without numerical types, by Propositions 34.6 and 34.8, there exists an  $\mathbb{L}_{\mathbf{C}}$  formula  $\varphi_e$  over  $\mathbf{S}$  with  $\text{FV}(\varphi_e) \subseteq \text{Var}_{\mathbf{o}}$ , and a tuple  $\bar{x}_e$  over  $\text{FV}(\varphi_e)$  that mentions all the variables of  $\text{FV}(\varphi_e)$ , such that  $e(D) = \varphi_e(\bar{x}_e)(D)$ , for every database  $D$  of  $\mathbf{S}$ . Therefore, by Proposition 34.11,  $e$  is  $(3^k - 1)/2$ -local, where  $k$  is the rank of the formula  $\varphi_e$ .

## Adding Recursion: Datalog

As discussed in Chapter 34, a serious limitation of relational algebra with aggregates, and in fact all of the query languages encountered so far in the previous chapters, is their inability to express recursive queries such as the reachability query. In this chapter, we introduce a rule-based language, called Datalog, that is powerful enough to express such queries. It can be seen as an extension of UCQs with the key feature of recursion.

### Syntax of Datalog

We start by defining the syntax of Datalog rules by using a rule-based syntax similar to that of CQs when seen as rules.

#### Definition 35.1: Syntax of Datalog

A *Datalog rule* over a schema  $\mathbf{S}$  is an expression of the form

$$R_0(\bar{x}) :- R_1(\bar{u}_1), \dots, R_n(\bar{u}_n)$$

for  $n \geq 1$ , where

- $R_i \in \mathbf{S}$ , for each  $i \in [0, n]$ ,
- $R_i(\bar{u}_i)$  is a relational atom, and  $\bar{u}_i$  is a tuple of constants and variables, for each  $i \in [n]$ ,
- $R_0(\bar{x})$  is a relational atom, and  $\bar{x}$  is a tuple of variables, and
- each variable mentioned in  $\bar{x}$  is also mentioned in  $\bar{u}_k$  for some  $k \in [n]$ ; this is known as the *safety condition*.

A *Datalog program* over  $\mathbf{S}$  is a finite set of Datalog rules over  $\mathbf{S}$ .

As we shall see, the key idea underlying Datalog queries is to declaratively specify what the query output should be by means of a Datalog program. The Datalog program that provides the specification of the reachability query over directed graphs follows.

### Example 35.2: Graph Reachability

Consider the following (named) database schema:

```
Edge [ node1, node2 ]
Reachable [ node1, node2 ]
```

The Edge relation stores the edges of the input directed graph  $G$ , and the Reachable relation stores the pairs of nodes  $(v, u)$  of  $G$  such that  $u$  is reachable from  $v$ . We can now inductively compute the Reachable relation via the following Datalog program over the above schema:

```
Reachable( $x, y$ ) :- Edge( $x, y$ )
Reachable( $x, y$ ) :- Reachable( $x, z$ ), Edge( $z, y$ ).
```

The first rule, which is the base step of the inductive definition, simply states that if there is an edge from  $x$  to  $y$ , then  $y$  is reachable from  $x$ . The second rule, which corresponds to the inductive step, states that if  $z$  is reachable from  $x$  and there is an edge from  $z$  to  $y$ , then  $y$  is reachable from  $x$ . Notice that the second rule is recursive in the sense that the definition of the relation Reachable depends on itself.

The relational atom that appears on the left of the  $\text{:-}$  symbol in a Datalog rule is called the *head* of the rule, while the expression that appears on the right of the  $\text{:-}$  symbol is called the *body* of the rule. Given a Datalog program  $\Pi$  over a schema  $\mathbf{S}$ , it is crucial to have a way to distinguish between the relation names of  $\mathbf{S}$  that appear only in the bodies of the rules of  $\Pi$ , and those that appear in the head of at least one rule of  $\Pi$ . In particular, a relation name  $R \in \mathbf{S}$  occurring in the Datalog program  $\Pi$  is called:

- *extensional* if there is no rule of the form  $R(\bar{x}) \text{ :- body}$  in  $\Pi$ , that is,  $R$  occurs only in rule-bodies, and
- *intensional* if there exists at least one rule of the form  $R(\bar{x}) \text{ :- body}$  in  $\Pi$ , that is,  $R$  appears in the head of at least one rule of  $\Pi$ .

Intuitively, extensional relation names correspond to the input relations, while intensional relation names correspond to the relations that are computed by the Datalog program. The *extensional (database) schema* of  $\Pi$ , denoted  $\text{edb}(\Pi)$ , consists of the extensional relation names in  $\Pi$ , while the *intentional schema* of  $\Pi$ , denoted  $\text{idb}(\Pi)$ , is the set of all intensional relation names in  $\Pi$ .

The *schema* of  $\Pi$ , denoted  $\text{sch}(\Pi)$ , is the set of relation names  $\text{edb}(\Pi) \cup \text{idb}(\Pi)$ . Note that  $\text{sch}(\Pi)$  is, in general, a subset of  $\mathbf{S}$  since it might be the case that some relation names of  $\mathbf{S}$  do not appear in  $\Pi$ .

## Semantics of Datalog

An interesting property of Datalog programs is the fact that their semantics can be defined either declaratively by adopting a *model-theoretic* approach, or operationally by following a *fixpoint* approach. In the model-theoretic approach, the Datalog rules are considered as logical sentences asserting a property of the desired result, while in the fixpoint approach the semantics is defined as a particular solution of a fixpoint procedure.

### Model-Theoretic Semantics

Recall that a set  $\Phi$  of first-order sentences over a schema  $\mathbf{S}$  is called a *first-order theory* over  $\mathbf{S}$ , or simply a *theory* over  $\mathbf{S}$ . A database of  $\mathbf{S}$  is a *model* of the theory  $\Phi$  if  $D \models \varphi$ , for every sentence  $\varphi \in \Phi$ .<sup>1</sup> The idea underlying the model-theoretic approach is to consider a Datalog program  $\Pi$  as a first-order theory  $\Phi_\Pi$  over  $\text{sch}(\Pi)$  that describes the desired outcome of the program. In other words, the desired outcome is a particular model of  $\Phi_\Pi$ ; hence the name model-theoretic semantics. However, there might be infinitely many models of the theory  $\Phi_\Pi$ , which means that the theory alone does not uniquely determine the desired outcome of the program. It is therefore crucial to specify which model is the intended outcome. We proceed to formalize the above discussion. In particular, we are going to explain how a Datalog program is converted into a first-order theory, and which model of this theory is the intended one.

#### Definition 35.3: From a Program to a Theory

Given a Datalog rule  $\rho$  of the form  $R_0(\bar{x}) :- R_1(\bar{u}_1), \dots, R_n(\bar{u}_n)$ , we write  $\varphi_\rho$  for the first-order sentence

$$\forall x_1 \dots \forall x_m (R_1(\bar{u}_1) \wedge \dots \wedge R_n(\bar{u}_n) \rightarrow R_0(\bar{x})),$$

where  $x_1, \dots, x_m$  are the variables in  $\rho$ . Given a Datalog program  $\Pi$ , we define the first-order theory  $\Phi_\Pi$  over  $\text{sch}(\Pi)$  as  $\{\varphi_\rho \mid \rho \in \Pi\}$ .

For brevity, we refer to the models of a Datalog program  $\Pi$  meaning the models of the theory  $\Phi_\Pi$ . Interestingly, the notion of satisfaction of a sentence  $\varphi_\rho$ , where  $\rho \in \Pi$ , by a database  $D$  of  $\text{sch}(\Pi)$ , can be characterized by means of

<sup>1</sup> The notion of first-order theory, together with the notion of model of such a theory, have been already used in Chapter 32

homomorphisms. Similarly to CQs, the body of a Datalog rule can be viewed as a set of atoms. More precisely, given a Datalog rule  $\rho$  of the form

$$R_0(\bar{x}) :- R_1(\bar{u}_1), \dots, R_n(\bar{u}_n)$$

we define the set of relational atoms

$$A_\rho = \{R_1(\bar{u}_1), \dots, R_n(\bar{u}_n)\}.$$

We can thus talk about homomorphisms from rule-bodies to databases. It is then easy to verify the following proposition that provides a useful characterization of rule satisfaction that will be used in our later proofs:

**Proposition 35.4**

Consider a Datalog rule  $\rho$  of the form  $R_0(\bar{x}) :- R_1(\bar{u}_1), \dots, R_n(\bar{u}_n)$ , and a database  $D$ . The following are equivalent:

1.  $D \models \varphi_\rho$ .
2. For every homomorphism  $h$  from  $A_\rho$  to  $D$ ,  $R_0(h(\bar{x})) \in D$ .

It should be clear that a Datalog program  $\Pi$  admits infinitely many models. We proceed to explain how we choose the intended one. The idea is that the intended model should not contain more atoms than needed for satisfying  $\Phi_\Pi$ . In other words, from all the models of  $\Phi_\Pi$ , we choose the  $\subseteq$ -minimal ones, i.e., those models  $D$  such that, for every atom  $R(\bar{a}) \in D$ , it is the case that  $D - \{R(\bar{a})\}$  is not a model of  $\Phi_\Pi$ . Based on this simple idea, we proceed to define the semantics of a Datalog program on an input database.

Given a Datalog program  $\Pi$  and a database  $D$  of  $\text{edb}(\Pi)$ , we define

$$\text{MM}(\Pi, D) = \{D' \mid D' \text{ is a } \subseteq\text{-minimal model of } \Pi \text{ and } D \subseteq D'\}.$$

We can show that  $\text{MM}(\Pi, D)$  contains *exactly one* database, which will give rise to the semantics of  $\Pi$  on  $D$ . But first we need to establish an auxiliary result. Let  $\text{B}(\Pi, D)$  be the union of  $D$  with the set of all relational atoms that can be formed using relation names from  $\text{idb}(\Pi)$  and constants from  $\text{Dom}(D)$ :

$$\text{B}(\Pi, D) = D \cup \left\{ R(\bar{a}) \mid R \in \text{idb}(\Pi) \text{ and } \bar{a} \in \text{Dom}(D)^{\text{ar}(R)} \right\}.$$

We can show the following:

**Lemma 35.5.** *Consider a Datalog program  $\Pi$ , and a database  $D$  of  $\text{edb}(\Pi)$ . It holds that  $\text{B}(\Pi, D)$  is a model of  $\Pi$  that contains  $D$ .*

*Proof.* The fact that  $\text{B}(\Pi, D)$  contains  $D$  follows by definition. It remains to show that  $\text{B}(\Pi, D)$  is a model of  $\Pi$ . Consider an arbitrary rule  $\rho \in \Pi$  of the

form  $R_0(\bar{x}) :- R_1(\bar{u}_1), \dots, R_n(\bar{u}_n)$ , and assume that there is a homomorphism  $h$  from  $A_\rho$  to  $D$ . Due to the safety condition, every variable in  $\bar{x}$  occurs in  $A_\rho$ . This implies that  $h(\bar{x})$  is a tuple over  $\text{Dom}(D)$ . Since  $R \in \text{idb}(\Pi)$  we conclude that  $R(h(\bar{x})) \in \mathbf{B}(\Pi, D)$ , and thus  $\mathbf{B}(\Pi, D) \models \varphi_\rho$ . Consequently,  $\mathbf{B}(\Pi, D)$  is a model of  $\Phi_\Pi$ , and thus, a model of  $\Pi$ , as needed.  $\square$

We are now ready to show the claimed statement concerning  $\text{MM}(\Pi, D)$ :

**Proposition 35.6**

Consider a Datalog program  $\Pi$ , and a database  $D$  of  $\text{edb}(\Pi)$ . Then

$$|\text{MM}(\Pi, D)| = 1.$$

*Proof.* By Lemma 35.5,  $\mathbf{B}(\Pi, D)$  is a model of  $\Pi$  that contains  $D$ . Therefore, there exists a subset of  $\mathbf{B}(\Pi, D)$  that belongs to  $\text{MM}(\Pi, D)$ , which implies that  $|\text{MM}(\Pi, D)| \geq 1$ . Assume now that  $|\text{MM}(\Pi, D)| \geq 2$ , and let  $D_1, \dots, D_\ell$ , for  $\ell \geq 2$ , be its members. We proceed to show that the database

$$D_\cap = D_1 \cap \dots \cap D_\ell$$

is a model of  $\Pi$  that contains  $D$ , which contradicts the fact that  $D_1, \dots, D_\ell$  are  $\subseteq$ -minimal. Since  $D \subseteq D_i$ , for each  $i \in [\ell]$ , we get that  $D \subseteq D_\cap$ . Consider now a Datalog rule  $\rho \in \Pi$  of the form  $R_0(\bar{x}) :- R_1(\bar{u}_1), \dots, R_n(\bar{u}_n)$ , and assume there exists a homomorphism  $h$  from  $A_\rho$  to  $D_\cap$ . For each  $i \in [\ell]$ ,  $D_i$  is a model of  $\Pi$ , and thus,  $R_0(h(\bar{x})) \in D_i$ . Therefore,  $R_0(h(\bar{x})) \in D_\cap$ , which means, due to Proposition 35.4, that  $D_\cap \models \varphi_\rho$ . Hence,  $D_\cap$  is a model of  $\Pi$ , as needed.  $\square$

Having the above result in place, we are now ready to define the semantics of a Datalog program on an input database.

**Definition 35.7: Semantics of Datalog**

Given a Datalog program  $\Pi$ , and a database  $D$  of  $\text{edb}(\Pi)$ , the *output of  $\Pi$  on  $D$* , denoted  $\Pi(D)$ , is the  $\subseteq$ -minimal model of  $\Pi$  that contains  $D$ .

By Proposition 35.6, we get that  $\Pi(D)$  is uniquely determined by the program and the database, and thus, Definition 35.7 provides a well-defined semantics for Datalog programs. The crucial question that comes up is whether we can devise an algorithm that computes the semantics of a Datalog program  $\Pi$  on a database  $D$ . The next result provides such an algorithm. Let  $\mathbf{M}(\Pi, D)$  be all the subsets of  $\mathbf{B}(\Pi, D)$  that are models of  $\Pi$  and contain  $D$ . Formally,

$$\mathbf{M}(\Pi, D) = \{D' \mid D' \text{ is a model of } \Pi \text{ and } D \subseteq D' \subseteq \mathbf{B}(\Pi, D)\}.$$

Interestingly, the intersection of the databases occurring in  $\mathbf{M}(\Pi, D)$  coincides with the  $\subseteq$ -minimal model of  $\Pi$  that contains  $D$ . By giving a proof similar to that of Proposition 35.6, we can show that  $\bigcap \mathbf{M}(\Pi, D)$  is a model of  $\Pi$  that contains  $D$ , while the fact that it is a  $\subseteq$ -minimal model follows by construction.

**Theorem 35.8**

Consider a Datalog program  $\Pi$ , and a database  $D$  of  $\text{edb}(\Pi)$ . Then

$$\Pi(D) = \bigcap_{D' \in \mathbf{M}(\Pi, D)} D'.$$

It is clear that Theorem 35.8 suggests the following procedure for computing the semantics of a Datalog program  $\Pi$  on a database  $D$ : construct all the possible subsets of  $\mathbf{B}(\Pi, D)$  that are models of  $\Pi$  and contain  $D$ , and then compute their intersection. However, this is computationally a very expensive procedure. As we shall see in the next section, the fixpoint approach provides a more efficient algorithm for computing the database  $\Pi(D)$ .

*Fixpoint Semantics*

We present an alternative way to define the semantics of Datalog that relies on an operator called the *immediate consequence operator*. This operator is applied on a database in order to produce new relational atoms. The model-theoretic semantics presented above coincides with the smallest solution of a fixpoint equation that involves the immediate consequence operator.

**Definition 35.9: Immediate Consequence Operator**

Consider a Datalog program  $\Pi$ , and a database  $D$  of  $\text{sch}(\Pi)$ . A relational atom  $R(\bar{a})$  is an *immediate consequence* for  $\Pi$  and  $D$  if:

1.  $R(\bar{a}) \in D$ , or
2. There exists a rule  $\rho \in \Pi$  of the form  $R(\bar{x}) :- R_1(\bar{u}_1), \dots, R_n(\bar{u}_n)$  such that  $(A_\rho, \bar{x}) \rightarrow (D, \bar{a})$ .

The *immediate consequence operator* of  $\Pi$  is defined as the function

$$T_\Pi : \text{Inst}(\text{sch}(\Pi)) \rightarrow \text{Inst}(\text{sch}(\Pi))$$

such that

$$T_\Pi(D) = \{R(\bar{a}) \mid R(\bar{a}) \text{ is an immediate consequence for } \Pi \text{ and } D\}.$$

A database  $D$  of  $\text{sch}(\Pi)$  is called a *fixpoint* of  $T_\Pi$  if  $T_\Pi(D) = D$ .

The next lemma, which is easy to prove, collects some useful properties of the  $T_\Pi$  operator that we are going to use below.

**Lemma 35.10.** *Consider a Datalog program  $\Pi$ . The following hold:*

1.  $T_\Pi$  is monotone, i.e., for every two databases  $D$  and  $D'$  of  $\text{sch}(\Pi)$ , if  $D \subseteq D'$  then  $T_\Pi(D) \subseteq T_\Pi(D')$ .
2. A database  $D$  of  $\text{sch}(\Pi)$  is a model of  $\Pi$  if and only if  $T_\Pi(D) \subseteq D$ .
3. Every fixpoint of  $T_\Pi$  is a model of  $\Pi$ .

We are now ready to establish the following crucial result, which states that the model-theoretic semantics of a Datalog program on a database  $D$  coincides with the  $\subseteq$ -minimal fixpoint of  $T_\Pi$  that contains  $D$ .

**Theorem 35.11**

Consider a Datalog program  $\Pi$ , and a database  $D$  of  $\text{edb}(\Pi)$ . It holds that  $\Pi(D)$  is the  $\subseteq$ -minimal fixpoint of  $T_\Pi$  that contains  $D$ .

*Proof.* We first show that  $\Pi(D)$  is a fixpoint of  $T_\Pi$ , i.e.,  $T_\Pi(\Pi(D)) = \Pi(D)$ . Since  $\Pi(D)$  is a model of  $\Pi$ , Lemma 35.10 implies that  $T_\Pi(\Pi(D)) \subseteq \Pi(D)$ . By Lemma 35.10,  $T_\Pi$  is monotone, and thus,  $T_\Pi(T_\Pi(\Pi(D))) \subseteq T_\Pi(\Pi(D))$ . Therefore, by Lemma 35.10,  $T_\Pi(\Pi(D))$  is a model of  $\Pi$  that contains  $D$ . But since  $\Pi(D)$  is the  $\subseteq$ -minimal mode of  $\Pi$  that contains  $D$ , we immediately get that  $\Pi(D) \subseteq T_\Pi(\Pi(D))$ . Consequently,  $T_\Pi(\Pi(D)) = \Pi(D)$ . By Lemma 35.10, each fixpoint of  $T_\Pi$  that contains  $D$  is a model of  $\Pi$  that contains  $D$ . Hence,  $\Pi(D)$  is the  $\subseteq$ -minimal fixpoint of  $T_\Pi$  that contains  $D$ .  $\square$

It remains to explain how the  $\subseteq$ -minimal fixpoint of the  $T_\Pi$  operator that contains the database  $D$  is constructed. This is essentially done by iteratively applying the  $T_\Pi$  operator starting from the database  $D$ .

**Definition 35.12: Application of the  $T_\Pi$  Operator**

Consider a Datalog program  $\Pi$ , and a database  $D$  of  $\text{edb}(\Pi)$ . We define

$$T_\Pi^0(D) = D \quad \text{and} \quad T_\Pi^{i+1}(D) = T_\Pi(T_\Pi^i(D)), \text{ for } i \in \mathbb{N},$$

and we let

$$T_\Pi^\infty(D) = \bigcup_{i \geq 0} T_\Pi^i(D).$$

At first glance, the construction of  $T_\Pi^\infty(D)$  requires infinitely many iterations. However, since  $T_\Pi^\infty(D) \subseteq \mathbf{B}(\Pi, D)$ , it is the case that  $T_\Pi^\infty(D)$  is obtained in at most  $|\mathbf{B}(\Pi, D)|$  iterations. It is easy to verify that

$$T_\Pi^\infty(D) = T_\Pi^{|\mathbf{B}(\Pi, D)|}(D).$$

We now show the following result, which essentially states that the semantics of a Datalog program  $\Pi$  on a database  $D$  can be computed by iteratively applying the operator  $T_\Pi$  starting from  $D$  until a fixpoint is reached.



**Theorem 35.13**

Consider a Datalog program  $\Pi$ , and a database  $D$  of  $\text{edb}(\Pi)$ . It holds that  $T_\Pi^\infty(D)$  is the  $\subseteq$ -minimal fixpoint of  $T_\Pi$  that contains  $D$ .

*Proof.* Recall first that the following hold:

$$T_\Pi^\infty(D) = T_\Pi^{|\mathcal{B}(\Pi, D)|}(D) \quad \text{and} \quad T_\Pi(T_\Pi^{|\mathcal{B}(\Pi, D)|}(D)) = T_\Pi^{|\mathcal{B}(\Pi, D)|}(D).$$

Therefore,  $T_\Pi^\infty(D)$  is a fixpoint of  $T_\Pi$  that contains  $D$ . It remains to show that  $T_\Pi^\infty(D)$  is  $\subseteq$ -minimal, or, equivalently,  $T_\Pi^\infty(D) \subseteq D'$ , for every fixpoint  $D'$  of  $T_\Pi$  that contains  $D$ . Fix such a fixpoint  $D'$ . We can show via an easy inductive argument that  $T_\Pi^i(D) \subseteq D'$ , for every  $i \in \mathbb{N}$ , which implies that  $T_\Pi^\infty(D) \subseteq D'$ . In fact,  $T_\Pi^0(D) \subseteq D'$  since  $T_\Pi^0(D) = D$ . Moreover,  $T_\Pi^i(D) \subseteq D'$  implies  $T_\Pi(T_\Pi^i(D)) = T_\Pi^{i+1}(D) \subseteq T_\Pi(D') = D'$  by monotonicity of  $T_\Pi$ .  $\square$

The next result is an immediate corollary of Theorems 35.11 and 35.13.

**Corollary 35.14**

Consider a Datalog program  $\Pi$ , and a database  $D$  of  $\text{edb}(\Pi)$ . Then

$$\Pi(D) = T_\Pi^\infty(D).$$

## Datalog Queries

Recall that a  $k$ -ary query  $q$  produces a finite set of  $k$ -ary tuples  $q(D) \subseteq \text{Const}^k$ , for every database  $D$ . Datalog programs can be used to define queries. In order to do this, we simply specify together with a Datalog program  $\Pi$  a relation name  $R$  from  $\text{idb}(\Pi)$  that indicates the relation that collects the output of the query. In other words, given a database  $D$  of  $\text{edb}(\Pi)$ , after computing the database  $\Pi(D)$ , the output of the query is the set of tuples  $\bar{a}$  over  $\text{Dom}(D)$  such that  $R(\bar{a}) \in \Pi(D)$ . For example, the Datalog query over  $\{\text{Edge}[2]\}$  that computes the pairs  $(v, u)$  such that  $u$  is reachable from  $v$  is  $(\Pi, \text{Reachable})$ , where  $\Pi$  is the Datalog program over  $\{\text{Edge}[2], \text{Reachability}[2]\}$  given in Example 35.2. The formal definition of Datalog queries follows.

**Definition 35.15: Datalog Queries**

A *Datalog query* over a schema  $\mathbf{S}$  is a pair  $(\Pi, R)$ , where  $\Pi$  is a Datalog program over a schema  $\mathbf{S} \cup \mathbf{S}'$ , with  $\mathbf{S}'$  being a schema disjoint from  $\mathbf{S}$ , such that  $\text{edb}(\Pi) \subseteq \mathbf{S}$ ,  $\text{idb}(\Pi) \subseteq \mathbf{S}'$ , and  $R \in \text{idb}(\Pi)$ .

Having the semantics of a Datalog program  $\Pi$  on a database  $D$  (see Definition 35.7), we can naturally define what is the output of a Datalog query  $(\Pi, R)$  on  $D$ ; simply collect the tuples in the relation  $R$  after computing  $\Pi(D)$ .

**Definition 35.16: Evaluation of Datalog Queries**

Given a database  $D$  of a schema  $\mathbf{S}$ , and a Datalog query  $q = (\Pi, R)$  over  $\mathbf{S}$ , the *output* of  $q$  on  $D$  is defined as the set of tuples

$$q(D) = \left\{ \bar{a} \in \text{Const}^{\text{ar}(R)} \mid R(\bar{a}) \in \Pi(D) \right\}.$$

It is clear that the set  $q(D)$  belongs to  $\mathcal{P}(\text{Const}^{\text{ar}(R)})$ . However, to be able to say that  $q$  defines a query over  $\mathbf{S}$  as in Definition 2.5, we need to ensure that  $q(D) \in \mathcal{P}_{\text{fin}}(\text{Const}^{\text{ar}(R)})$ , i.e., the output of  $q$  on  $D$  is finite. This is guaranteed by the following result, which is an immediate consequence of Theorem 35.8 and the fact that  $\text{Dom}(D') = \text{Dom}(D)$ , for every database  $D' \in \mathbf{M}(\Pi, D)$ .

**Proposition 35.17**

For a database  $D$  of schema  $\mathbf{S}$ , and a Datalog query  $q = (\Pi, R)$  over  $\mathbf{S}$ ,

$$q(D) = \left\{ \bar{a} \in \text{Dom}(D)^{\text{ar}(R)} \mid R(\bar{a}) \in \Pi(D) \right\}.$$

Since  $\text{Dom}(D)$  is finite, Proposition 35.17 implies that  $q(D) \in \mathcal{P}_{\text{fin}}(\text{Const}^k)$ , and thus,  $q$  defines a query over  $\mathbf{S}$  in the sense of Definition 2.5.

At the beginning of the chapter, we claimed that Datalog extends UCQ s with the feature of recursion. We can easily show that indeed Datalog leads to a strictly more expressive language that is able to express recursive queries:

**Theorem 35.18**

The language of Datalog queries is strictly more expressive than the language of UCQs.

*Proof.* From Theorem 34.1, we conclude that the reachability query on directed graphs cannot be expressed as a UCQ, but we have already seen that it can be easily expressed as a Datalog query. On the other hand, it is straightforward to see that every UCQ  $q(\bar{x}) = q_1 \cup \dots \cup q_n$  can be equivalently written as the Datalog query  $(\Pi_q, \text{Answer})$ , where  $\Pi_q$  consists of the CQs  $q_1, \dots, q_n$  seen as rules of the form  $\text{Answer}(\bar{x}) \text{ :- body}$ .  $\square$

Recall from Chapter 28 that UCQ s with variable-constant equality form a strictly more expressive language than UCQ s. It turns out that there exists a UCQ with variable-constant equality that cannot be expressed as a Datalog query, which means that Datalog and UCQ s with variable-constant equality form incomparable languages in terms of expressive power. This is because UCQ s with variable-constant equality may have in their output constants not

from the domain of the database, which is not possible for Datalog queries (Proposition 35.17). Consider, for example, the simple query  $q = \varphi(x)$  with

$$\varphi = (x = a),$$

where  $a$  is a constant. For  $D = \{R(b)\}$ , we get that  $q(D) = \{(a)\}$ .

## Expressiveness of Datalog Queries

We have already seen in the previous chapter that Datalog queries are strictly more expressive than UCQs. We have also seen an easy inexpressibility result, i.e., there are UCQs with variable-constant equality (in fact, the query  $\varphi(x)$  with  $\varphi = (x = a)$ ) that cannot be expressed as a Datalog query. The question that comes up is how Datalog queries compare in terms of expressive power with UCQs with inequality, and more generally, whether Datalog queries can express negation, or at least a restricted form of negation.

Our goal in this chapter is to show that Datalog queries are inherently positive. Note that the easy inexpressibility result that Datalog queries cannot express  $\varphi(x)$  with  $\varphi = (x = a)$  relies on the property of Datalog queries provided by Proposition 35.17, that is, the output of a Datalog query mentions only constants from the domain of the database. However, this property is not powerful enough to show that Datalog queries are inherently positive. We proceed to establish that Datalog queries are preserved under homomorphisms and monotone, and then use those properties to show that indeed Datalog queries cannot express inequality, negative relational atoms, and difference.

### Preservation Under Homomorphisms

The notion of preservation under homomorphisms for Datalog queries is defined in the same way as for FO queries. For a Datalog program  $\Pi$ , we write  $\text{Dom}(\Pi)$  for the set of constants occurring in the rules of  $\Pi$ .

#### Definition 36.1: Preservation Under Homomorphisms

Consider a Datalog query  $q = (\Pi, R)$  over a schema  $\mathbf{S}$ . We say that  $q$  is *preserved under homomorphisms* if, for every two databases  $D$  and  $D'$  of  $\mathbf{S}$ , and tuples  $\bar{a} \in \text{Dom}(D)^{\text{ar}(R)}$  and  $\bar{b} \in \text{Dom}(D')^{\text{ar}(R)}$ , it holds that

if  $(D, \bar{a}) \rightarrow_{\text{Dom}(\Pi)} (D', \bar{b})$  and  $\bar{a} \in q(D)$  then  $\bar{b} \in q(D')$ .

We proceed to show that Datalog queries are preserved under homomorphisms. The key idea underlying this result is that a Datalog query  $q$  over a schema  $\mathbf{S}$  can be converted into an equivalent UCQ  $q'$  over  $\mathbf{S}$  providing that  $q'$  can have infinitely many disjuncts; such UCQs are called *infinitary*. The evaluation of infinitary UCQs is defined in the same way as for UCQs, that is, given a database  $D$  of a schema  $\mathbf{S}$ , and an infinitary UCQ  $q = q_1 \cup q_2 \cup \dots$  over  $\mathbf{S}$ , where each  $q_i$  for  $i \geq 1$  is a CQ over  $\mathbf{S}$ ,  $q(D) = \bigcup_{i=1}^{\infty} q_i(D)$ . It is easy to show that every infinitary UCQ is preserved under homomorphisms; the proof is essentially the same as the one of Proposition 28.10, which establishes that UCQs are preserved under homomorphisms. Therefore, to show that Datalog queries are preserved under homomorphisms it suffices to show that a Datalog query over  $\mathbf{S}$  can be converted into an equivalent infinitary UCQ over  $\mathbf{S}$ .

### Proposition 36.2

Consider a Datalog query  $q = (\Pi, R)$  over a schema  $\mathbf{S}$ . There exists an infinitary UCQ  $q' = \varphi(\bar{x})$  over  $\mathbf{S}$  with  $\text{Dom}(\Pi) = \text{Dom}(\varphi)$  and  $q \equiv q'$ .

For technical clarity, we show the above result only for Boolean queries. Nevertheless, the given proof illustrates the key elements that are used in the proof for non-Boolean queries, which we leave as an exercise. We first need to introduce the basic notions of *unification* and *unfolding*.

We say that two atoms  $R(\bar{u})$  and  $P(\bar{v})$  *unify* if there exists a function  $\gamma : \text{Const} \cup \text{Var} \rightarrow \text{Const} \cup \text{Var}$ , which is the identity on  $\text{Const}$  and the set of variables not mentioned in  $\bar{u}$  and  $\bar{v}$ , such that  $R(\gamma(\bar{u})) = P(\gamma(\bar{v}))$ ; such a function  $\gamma$  is called a *unifier for  $R(\bar{u})$  and  $P(\bar{v})$* . Observe that for  $R(\bar{u})$  and  $P(\bar{v})$  to unify it is a necessary condition that  $R$  and  $P$  are the same relation names, and  $\bar{u}, \bar{v}$  have the same arity. A *most general unifier* for  $R(\bar{u})$  and  $P(\bar{v})$  is a unifier  $\gamma$  for them such that, for every other unifier  $\gamma'$  for  $R(\bar{u})$  and  $P(\bar{v})$ , there exists a function  $\theta : \text{Const} \cup \text{Var} \rightarrow \text{Const} \cup \text{Var}$  such that  $\gamma' = \theta \circ \gamma$ . For example, the atoms  $R(x, y)$  and  $R(z, a)$  unify due to the function  $\gamma : \text{Const} \cup \text{Var} \rightarrow \text{Const} \cup \text{Var}$  that is the identity on  $\text{Const} \cup (\text{Var} - \{x, y, z\})$ , and assigns the constant  $a$  to the variables  $x, y, z$ , i.e.,  $\gamma(x) = \gamma(y) = \gamma(z) = a$ . However, it is easy to verify that  $\gamma$  is not a most general unifier for  $R(x, y)$  and  $R(z, a)$ . In particular, for the unifier  $\gamma'$  for  $R(x, y)$  and  $R(z, a)$  with  $\gamma'(x) = \gamma'(z) = z$  and  $\gamma'(y) = a$ , we can show that there is no  $\theta : \text{Const} \cup \text{Var} \rightarrow \text{Const} \cup \text{Var}$  such that  $\gamma' = \theta \circ \gamma$ . On the other hand, we can show that  $\gamma'$  is a most general unifier for  $R(x, y)$  and  $R(z, a)$ . For instance,  $\gamma = \theta \circ \gamma'$  with  $\theta : \text{Const} \cup \text{Var} \rightarrow \text{Const} \cup \text{Var}$  being such that  $\theta(u) = a$ .

It is easy to show that, for any two atoms  $R(\bar{u})$  and  $P(\bar{v})$ ,

- if  $R(\bar{u})$  and  $P(\bar{v})$  unify, then there is a most general unifier for them, and

- if  $\gamma_1$  and  $\gamma_2$  are most general unifiers for  $R(\bar{u})$  and  $P(\bar{v})$ , then, for every relational atom  $S(\bar{w})$ , it holds that  $S(\gamma_1(\bar{w}))$  and  $S(\gamma_2(\bar{w}))$  are the same up to variable renaming.

These facts allow us to refer to *the* most general unifier for  $R(\bar{u})$  and  $P(\bar{v})$ .

We now proceed to introduce the notion of unfolding of a Boolean CQ with a Datalog program, which relies on unification. Let  $q$  be the Boolean CQ

$$\text{Answer} :- R_1(\bar{u}_1), \dots, R_n(\bar{u}_n),$$

and  $\rho$  be the Datalog rule

$$P_0(\bar{y}) :- P_1(\bar{v}_1), \dots, P_m(\bar{v}_m).$$

We can always assume that  $q$  and  $\rho$  do not share variables since we can simply rename the variables occurring in  $\rho$  without changing its semantic meaning. Assume now that, for  $i \in [n]$ ,  $R_i(\bar{u}_i)$  and  $P_0(\bar{y})$  unify, and let  $\gamma$  be their most general unifier. We denote by  $q_{\rho, i, \gamma}$  the CQ

$$\begin{aligned} \text{Answer} :- R_1(\gamma(\bar{u}_1)), \dots, R_{i-1}(\gamma(\bar{u}_{i-1})), \\ R_{i+1}(\gamma(\bar{u}_{i+1})), \dots, R_n(\gamma(\bar{u}_n)), P_1(\gamma(\bar{v}_1)), \dots, P_m(\gamma(\bar{v}_m)). \end{aligned}$$

Let  $\{R_{i_1}(\bar{u}_{i_1}), \dots, R_{i_\ell}(\bar{u}_{i_\ell})\}$  be the set that collects all the relational atoms in the body of  $q$  that unify with  $P_0(\bar{y})$ , and let  $\gamma_{i_j}$  be the most general unifier for  $R_{i_j}(\bar{u}_{i_j})$  and  $P_0(\bar{y})$ , for each  $j \in [\ell]$ . The *unfolding of  $q$  with  $\rho$* , denoted  $\text{Unfold}_\rho(q)$ , is defined as the set of CQs  $\{q_{\rho, i_1, \gamma_{i_1}}, \dots, q_{\rho, i_\ell, \gamma_{i_\ell}}\}$ . Now, for a Datalog program  $\Pi$ , the *unfolding of  $q$  with  $\Pi$* , denoted  $\text{Unfold}_\Pi(q)$ , is

$$\bigcup_{\rho \in \Pi} \text{Unfold}_\rho(q).$$

Therefore,  $\text{Unfold}_\Pi(\cdot)$  can be seen as an operator that takes as input a CQ and computes all the CQs that can be obtained by unfolding  $q$  with a Datalog rule from  $\Pi$ . We can then define the set of all CQs that can be obtained starting from  $q$  and exhaustively applying the  $\text{Unfold}_\Pi(\cdot)$  operator.

**Definition 36.3: Application of the  $\text{Unfold}_\Pi(\cdot)$  Operator**

Consider a Datalog program  $\Pi$  over  $\mathbf{S}$ , and a CQ  $q$  over  $\mathbf{S}$ . We define

$$\text{Unfold}_\Pi^0(q) = \{q\} \quad \text{and} \quad \text{Unfold}_\Pi^{i+1}(q) = \bigcup_{q' \in \text{Unfold}_\Pi^i(q)} \text{Unfold}_\Pi(q')$$

for  $i \geq 0$ , and let

$$\text{Unfold}_\Pi^\infty(q) = \bigcup_{i \geq 0} \text{Unfold}_\Pi^i(q).$$

We are now ready to show Proposition 36.2

*Proof (of Proposition 36.2).* We prove this for Boolean queries, i.e., we assume that the Datalog query  $q = (\Pi, R)$  is Boolean, that is,  $\text{ar}(R) = 0$ . We further assume, without affecting the generality of the proof, that there is exactly one rule  $\rho_R \in \Pi$  of the form  $R() :- \text{body}$ , and there is no rule in  $\Pi$  that mentions  $R$  in its body, i.e., we assume that the intensional relation name  $R$  occurs only in the head of  $\rho_R$ . We can indeed make this assumption since we can always rewrite  $q$  into an equivalent Datalog query with the above property: construct  $\Pi^*$  by replacing every occurrence of the relation name  $R$  in  $\Pi$  with a new relation name  $R^*$  not occurring in  $\text{sch}(\Pi)$ , and then consider the query

$$q^* = (\Pi^* \cup \{R() :- R^*()\}, R).$$

It is clear that  $\text{edb}(\Pi) = \text{edb}(\Pi^*)$  and  $q(D) = q^*(D)$  for every  $D$  of  $\text{edb}(\Pi)$ . In what follows, let  $q_R$  for the Boolean CQ such that  $A_{q_R} = A_{\rho_R}$ , i.e.,  $q_R$  is the Boolean CQ that has as its body the body of the rule  $\rho_R$ .

We are now ready to define the desired infinitary UCQ over  $\mathbf{S}$  by using the  $\text{Unfold}_\Pi(\cdot)$  operator. We write  $\text{Unfold}_\Pi^\infty(q_R)|_{\mathbf{S}}$  for the subset of  $\text{Unfold}_\Pi^\infty(q_R)$  that keeps only the CQs over  $\mathbf{S}$ , that is, the CQs that use only relation names from  $\mathbf{S}$ . We then define the infinitary UCQ

$$q_R^\Pi = \bigcup_{q' \in \text{Unfold}_\Pi^\infty(q_R)|_{\mathbf{S}}} q'.$$

By construction,  $q_R^\Pi$  is an infinitary UCQ over  $\mathbf{S}$ , and  $\Pi$  and  $q_R^\Pi$  mention exactly the same constants. It remains to show that  $q$  and  $q_R^\Pi$  are equivalent, i.e.,  $q(D) = q_R^\Pi(D)$ , for every database  $D$  of  $\mathbf{S}$ . To this end, it suffices to show that, for a database  $D$  of  $\mathbf{S}$ , the following are equivalent:

1.  $A_{\rho_R} \rightarrow \Pi(D)$ .
2. There exists a sequence of CQs  $(q_i)_{i \in [0, n]}$ , for some  $n \geq 0$ , such that:
  - $q_0 = q_R$ ,
  - $q_i \in \text{Unfold}_\Pi(q_{i-1})$ , for each  $i \in [n]$ , and
  - $A_{q_n} \rightarrow D$ .

*The Direction (1)  $\Rightarrow$  (2)*

By hypothesis, there exists a sequence of databases  $(D_i)_{i \in [0, n]}$ , for some  $n \geq 0$ , such that: (i)  $D_0 = D$ , (ii) for each  $i \in [n]$ ,  $D_i = D_{i-1} \cup \{P(\bar{a})\}$ , where  $P(\bar{a}) \in T_\Pi(D_{i-1})$ , i.e., there is  $\rho_{i-1} \in \Pi$  of the form  $P(\bar{x}) :- \text{body}$  such that  $(A_{\rho_{i-1}}, \bar{x}) \rightarrow (D_{i-1}, \bar{a})$  via a homomorphism  $h_{i-1}$ , and (iii)  $A_{\rho_R} \rightarrow D_n$ . We proceed to show the following auxiliary lemma:

**Lemma 36.4.** *There exists a sequence of CQs  $(q_i)_{i \in [0, n]}$  such that:*

- $q_0 = q_R$ ,
- $q_i = q_{i-1}$  or  $q_i \in \text{Unfold}_{\rho_{n-i}}(q_{i-1})$ , for each  $i \in [n]$ , and
- $A_{q_i} \rightarrow D_{n-i}$ , for each  $i \in [n]$ .

*Proof.* We proceed by induction on the length of  $(D_i)_{i \in [0, n]}$ . For the base case the statement holds trivially since  $A_{\rho_R} \rightarrow D_0$ , which in turn implies that  $A_{q_R} \rightarrow D_0$ . In other words, the desired sequence of CQs consists only of  $q_0$ .

We proceed with the inductive step. By hypothesis,  $A_{q_R} \rightarrow D_n$  via a homomorphism  $\mu$ . We consider two cases:

- Assume first that  $P(h_{n-1}(\bar{x})) \notin \mu(A_{q_R})$ , or  $P(h_{n-1}(\bar{x})) \in \mu(A_{q_R})$  and  $P(h_{n-1}(\bar{x})) \in D_{n-1}$ . This implies that  $A_{q_R} \rightarrow D_{n-1}$ . By induction hypothesis, there exists a sequence of CQs  $(q'_i)_{i \in [0, n-1]}$ , where  $q'_0 = q_R$ , and, for each  $i \in [n-1]$ ,  $q'_i = q'_{i-1}$  or  $q'_i \in \text{Unfold}_{\rho_{n-i}}(q'_{i-1})$ , and  $A_{q'_i} \rightarrow D_{n-i}$ . Therefore, the claim follows due to the sequence of CQs  $q'_0, q'_0, q'_1, \dots, q'_{n-1}$ .
- The interesting case is when  $P(h_{n-1}(\bar{x})) \notin \mu(A_{q_R})$  and  $P(h_{n-1}(\bar{x})) \in D_n - D_{n-1}$ . It is clear that there exists  $P(\bar{u}) \in A_{q_R}$  such that  $P(\mu(\bar{u})) = P(h_{n-1}(\bar{x}))$ , which means that  $\gamma = \mu \cup h_{n-1}$  is a unifier for  $P(\bar{u})$  and  $P(\bar{x})$ . This implies that there exists a most general unifier  $\hat{\gamma}$  for  $P(\bar{u})$  and  $P(\bar{x})$ . Let  $\hat{q}$  be the unfolding of  $q_R$  with  $\rho_{n-1}$  using  $\hat{\gamma}$ . In fact, we can show that  $A_{\hat{q}} \rightarrow D_{n-1}$ . By definition of most general unifiers,  $\gamma = \theta \circ \hat{\gamma}$  for some function  $\theta : \text{Const} \cup \text{Var} \rightarrow \text{Const} \cup \text{Var}$ . It is clear that  $\theta$  is a homomorphism from  $A_{\hat{q}}$  to  $D_{n-1}$  since  $\gamma$  maps  $A_{\rho_{n-1}}$  to  $D_{n-1}$ . Therefore,  $A_{\hat{q}} \rightarrow D_{n-1}$  as claimed above. By induction hypothesis, there is a sequence of CQs  $(q'_i)_{i \in [0, n-1]}$ , where  $q'_0 = \hat{q}$ , and, for each  $i \in [n-1]$ ,  $q'_i = q'_{i-1}$  or  $q'_i \in \text{Unfold}_{\rho_{n-i}}(q'_{i-1})$ , and  $A_{q'_i} \rightarrow D_{n-i}$ . The claim follows due to the sequence of CQs  $q_R, q'_0, \dots, q'_{n-1}$ .

This completes the proof of Lemma 36.4 □

We can now complete the proof of the direction (1)  $\Rightarrow$  (2). Let  $(q_i)_{i \in [0, n]}$  be the sequence of CQs provided by Lemma 36.4. Clearly,  $q_0 = q_R$  and  $A_{q_n} \rightarrow D_n$ . However, it is not the case that  $q_i \in \text{Unfold}_\Pi(q_{i-1})$ , for each  $i \in [n]$ , due to the fact that some CQs in  $(q_i)_{i \in [0, n]}$  are simply repeated. This can be easily fixed by removing the redundant CQs. Let  $Ind$  be the set of indices

$$\{i_j \mid j \in [n] \text{ and } q_{i_j} \notin \text{Unfold}_\Pi(q_{i_j-1})\}.$$

Observe that  $0 \notin Ind$ , and that, for each  $k \in Ind$ ,  $q_k = q_{k-1}$ . Therefore, the sequence of CQs  $(q'_i)_{i \in [0, m]}$ , where  $m \leq n$ , obtained from  $(q_i)_{i \in [0, n]}$  by simply removing the CQs  $\{q_k \mid k \in Ind\}$  is such that  $q'_0 = q_R$ ,  $q'_i \in \text{Unfold}_\Pi(q'_{i-1})$ , for each  $i \in [m]$ , and  $A_{q'_m} \rightarrow D$ . This implies that (2) holds, as needed.



The Direction (2)  $\Rightarrow$  (1)

We first establish the following auxiliary lemma:

**Lemma 36.5.** *For every  $i \in [n]$ ,  $A_{q_i} \rightarrow \Pi(D)$  implies  $A_{q_{i-1}} \rightarrow \Pi(D)$ .*

*Proof.* By hypothesis, there exists a homomorphism  $h$  that maps  $A_{q_i}$  to  $\Pi(D)$ . Since  $q_i \in \text{Unfold}_\Pi(q_{i-1})$ , we conclude that  $q_i \in \text{Unfold}_\rho(q_{i-1})$  for some  $\rho \in \Pi$  of the form  $P(\bar{x}) :- \text{body}$ . This means that there exists an atom  $P(\bar{u}) \in A_{q_{i-1}}$  that unifies with  $P(\bar{x})$ , and  $q_i$  is the unfolding of  $q_{i-1}$  with  $\rho$  using the most general unifier  $\gamma$  for  $P(\bar{x})$  and  $P(\bar{u})$ . We show that the function  $\mu = h \circ \gamma$  is a homomorphism from  $A_{q_{i-1}}$  to  $\Pi(D)$ , which witnesses that  $A_{q_{i-1}} \rightarrow \Pi(D)$ .

Since  $h$  maps  $A_{q_i}$  to  $\Pi(D)$ , we get that  $h$  maps  $\gamma(A_{q_{i-1}} - \{P(\bar{u})\})$  to  $\Pi(D)$ , i.e.,  $\mu$  maps  $A_{q_{i-1}} - \{P(\bar{u})\}$  to  $\Pi(D)$ . It remains to show that  $P(\mu(\bar{u})) \in \Pi(D)$ . Since  $\gamma(A_\rho) \subseteq \gamma(A_{q_i})$ , we conclude that  $h$  maps  $\gamma(A_\rho)$  to  $\Pi(D)$ , i.e.,  $\mu$  is a homomorphism from  $A_\rho$  to  $\Pi(D)$ . This implies that  $P(\mu(\bar{x})) \in \Pi(D)$ . Observe that  $P(\mu(\bar{x})) = P(\mu(\bar{u}))$ . Indeed, since  $\gamma(\bar{x}) = \gamma(\bar{u})$ , we get that

$$P(\mu(\bar{x})) = P(h(\gamma(\bar{x}))) = P(h(\gamma(\bar{u}))) = P(\mu(\bar{u})),$$

which in turn implies that  $P(\mu(\bar{u})) \in \Pi(D)$ , and the claim follows.  $\square$

We can now complete the proof of the direction (2)  $\Rightarrow$  (1). By hypothesis,  $A_{q_n} \rightarrow D$ , and thus,  $A_{q_n} \rightarrow \Pi(D)$ ; the latter holds due to the monotonicity of CQs (Corollary 13.7). By repeatedly applying Lemma 36.5, we get that  $A_{q_0} \rightarrow \Pi(D)$ . Since  $q_0 = q_R$  and  $A_{q_R} = A_{\rho_R}$ , we conclude that  $A_{\rho_R} \rightarrow \Pi(D)$ .  $\square$

By Proposition 36.2 and the fact that infinitary UCQs are preserved under homomorphisms, we immediately get the following result:

**Corollary 36.6**

Every Datalog query is preserved under homomorphisms.

Another key property is that of monotonicity. Recall that a query  $q$  over a schema  $\mathbf{S}$  is *monotone* if, for every two databases  $D, D'$  of  $\mathbf{S}$ , we have that

$$D \subseteq D' \text{ implies } q(D) \subseteq q(D').$$

We can show that homomorphism preservation implies monotonicity of Datalog queries. In fact, the proof is exactly the same as the one of Corollary 13.7 which establishes that every CQ is monotone.

**Corollary 36.7**

Every Datalog query is monotone.

## Datalog Queries and Negation

We now delineate the expressiveness boundaries of Datalog queries. We show that they cannot express inequality, negative relational atoms, and difference.

**Datalog queries cannot express inequality.** This is because already CQs with inequality are not preserved under homomorphisms. Consider

$$q_1 = \exists x \exists y (R(x, y) \wedge x \neq y).$$

For  $D = \{R(a, b)\}$  and  $D' = \{R(c, c)\}$ , we have that  $D \rightarrow_\emptyset D'$ . However,  $D \models q_1$  while  $D' \not\models q_1$ . As a second example, consider the CQ $^\neq$

$$q_2 = \exists x (S(x) \wedge x \neq a),$$

where  $a$  is a constant. Given  $D = \{S(b)\}$  and  $D' = \{S(a)\}$ , we have that  $D \rightarrow_{\{a\}} D'$ . However,  $D \models q_2$  while  $D' \not\models q_2$ .

**Datalog queries cannot express negative relational atoms.** The reason is because such queries are not monotone. Consider the query

$$q = \neg P(a),$$

where  $a$  is a constant. If we take  $D = \emptyset$  and  $D' = \{P(a)\}$ , then  $D \subseteq D'$  but  $D \models q$  while  $D' \not\models q$ .

**Datalog queries cannot express difference.** This is because difference is not monotone. Consider, for example, the FO query

$$q = \exists x (P(x) \wedge \neg Q(x)).$$

For  $D = \{P(a)\} \subseteq D' = \{P(a), Q(a)\}$ , we have that  $D \models q$  while  $D' \not\models q$ .



## Datalog Query Evaluation

In this chapter, we study the complexity of evaluating Datalog queries, that is, **Datalog-Evaluation**. This is the problem of checking whether  $\bar{a} \in q(D)$  for a Datalog query  $q$ , a database  $D$ , and a tuple  $\bar{a}$  over  $\text{Dom}(D)$ .

### Combined Complexity

We first look at the combined complexity of the problem, i.e., when the input consists of a Datalog query  $q = (\Pi, R)$ , a database  $D$  of  $\text{edb}(\Pi)$ , and a tuple  $\bar{a} \in \text{Dom}(D)^{\text{ar}(R)}$ . Recall that the fixpoint approach for defining the semantics of Datalog programs provides an algorithm for computing the database  $\Pi(D)$ . In particular, by Corollary 35.14,  $\Pi(D) = T_\Pi^\infty(D)$ , which in turn implies that

$$\bar{a} \in q(D) \text{ if and only if } R(\bar{a}) \in T_\Pi^\infty(D).$$

We proceed to analyze the time complexity of checking whether the fact  $R(\bar{a})$  belongs to  $T_\Pi^\infty(D)$ . Recall that for computing  $T_\Pi^\infty(D)$  we need to apply the  $T_\Pi$  operator at most  $|\mathbf{B}(\Pi, D)|$  times. We first analyze the time complexity of the  $i$ -th application of the  $T_\Pi$  operator. Let  $\text{maxvar}$  and  $\text{maxbody}$  be the maximum number of variables and body atoms, respectively, in a rule of  $\Pi$ . The  $i$ -th application of  $T_\Pi$ , for  $i \in \{1, \dots, |\mathbf{B}(\Pi, D)|\}$ , takes time

$$O(|\Pi| \cdot |\text{Dom}(D)|^{\text{maxvar} \cdot \text{maxbody}} \cdot |T_\Pi^{i-1}(D)|)$$

since, for each  $\rho \in \Pi$ , we need to consider all the possible functions  $h$ , which are the identity on  $\text{Const}$ , from the variables and constants in  $\rho$  to  $\text{Dom}(D)$ , and then check whether  $h$  is a homomorphism from the set of atoms in the body of  $\rho$  to  $T_\Pi^{i-1}(D)$ . Recall that, for each  $i \in \{1, \dots, |\mathbf{B}(\Pi, D)|\}$ ,  $|T_\Pi^{i-1}(D)| \leq |\mathbf{B}(\Pi, D)|$ . Hence, each application of the  $T_\Pi$  operator takes time

$$O(|\Pi| \cdot |\text{Dom}(D)|^{\text{maxvar} \cdot \text{maxbody}} \cdot |\mathbf{B}(\Pi, D)|).$$

As said above, for computing  $T_H^\infty(D)$  we need to apply the  $T_H$  operator at most  $|\mathbf{B}(H, D)|$  times, which implies that  $T_H^\infty(D)$  can be computed in time

$$O(|H| \cdot |\text{Dom}(D)|^{\maxvar} \cdot \maxbody \cdot |\mathbf{B}(H, D)|^2).$$

It is easy to verify that

$$|\mathbf{B}(H, D)| \leq |\text{sch}(H)| \cdot |\text{Dom}(D)|^{\text{ar}(H)},$$

where  $\text{ar}(H)$  is the maximum arity over all relation names of  $\text{sch}(H)$ . Consequently,  $T_H^\infty(D)$  can be computed in exponential time, which implies that checking whether  $R(\bar{a}) \in T_H^\infty(D)$  is feasible in exponential time.

One may think that there is a more clever procedure than naively computing  $T_H^\infty(D)$  that allows us to show that the complexity of **Datalog-Evaluation** matches the complexity of **UCQ-Evaluation**, that is, NP-complete. However, we can show that exponential time is the best that we can achieve.

**Theorem 37.1**

**Datalog-Evaluation** is EXPTIME-complete.

*Proof.* We have already seen that **Datalog-Evaluation** is in EXPTIME. We proceed to show that **Datalog-Evaluation** is EXPTIME-hard. This is done by showing that an arbitrary language  $L$  in EXPTIME is polynomial time reducible to **Datalog-Evaluation**. Let  $M = (Q, \Sigma, \delta, s)$  be a (deterministic) Turing Machine that decides  $L$  in exponential time; details on Turing Machines can be found in Appendix [B](#). The goal is, on input  $w$ , to construct in polynomial time in  $|w|$  a database  $D$ , and a Boolean Datalog query  $q = (H, \text{Yes})$ , i.e., Yes is a 0-ary relation name, such that

$$M \text{ accepts } w \quad \text{if and only if} \quad q(D) = \text{true}.$$

We first describe the high level idea of the reduction.

Consider a pair  $(p, a) \in (Q - \{\text{"yes"}, \text{"no"}\}) \times \Sigma$ . The transition rule  $\delta(p, a) = (p', b, \text{dir})$  expresses the following if-then statement:

**if** at some time instant  $t$  of the computation of  $M$  on  $w$ , we have that  $M$  is in state  $p$ , the head points to the tape cell  $c$ , and  $c$  contains the symbol  $a$   
**then** at time instant  $t+1$ , we have that  $M$  is in state  $p'$ , the cell  $c$  contains  $b$ , and the head points to the cell  $c'$ , where  $c'$  is the cell right to  $c$  (respectively, the cell left to  $c$ ,  $c$  itself) if  $\text{dir} = \rightarrow$  (respectively,  $\text{dir} = \leftarrow$ ,  $\text{dir} = -$ ).

We can naturally encode such an if-then statement via Datalog rules since a Datalog rule is essentially an if-then statement. This in turn allows us to describe the complete evolution of  $M$  on input  $w$  from its start configuration  $sc(w)$  to configuration  $c$  that can be reached in  $2^m$  steps, where  $m = |w|^k$  for

some  $k \in \mathbb{N}$ . To achieve this, we need a way to refer to the  $i$ -th time instant of the computation of  $M$  on  $w$ , and the  $i$ -th tape cell of  $M$ , where  $0 \leq i \leq 2^m - 1$ . This can be done by representing the time instances and the tape cells from 0 to  $2^m - 1$  by tuples of size  $m$  over  $\{0, 1\}$ , on which the functions “next time instant” and “next tape cell” are realized by means of a successor relation  $\text{Succ}^m$  from a linear order  $\preceq^m$  on  $\{0, 1\}^m$ . We now formalize this description.

#### *The Extensional and Intensional Schema*

We begin by describing the extensional and intensional schema of  $\Pi$ . As we shall see, there will be relations  $\text{Succ}^i$ ,  $\text{First}^i$  and  $\text{Last}^i$ , for each  $i \in [m]$ , which tell the successor, the first, and the last element from a linear order  $\preceq^i$  on  $\{0, 1\}^i$ , respectively, that will be inductively constructed by  $\Pi$  starting from  $\text{Succ}^1$ ,  $\text{First}^1$  and  $\text{Last}^1$ . The extensional schema  $\text{edb}(\Pi)$  is

$$\{\text{Succ}^1, \text{First}^1, \text{Last}^1\},$$

where  $\text{Succ}^1$  is a binary relation name, and  $\text{First}^1$ ,  $\text{Last}^1$  are unary relation names. The intensional schema  $\text{idb}(\Pi)$  is defined as

$$\begin{aligned} & \{\text{Symbol}_a \mid a \in \Sigma\} \cup \{\text{Head}\} \cup \{\text{State}_p \mid p \in Q\} \cup \\ & \{\text{Yes}\} \cup \bigcup_{i \in [2, m]} \{\text{Succ}^i, \text{First}^i, \text{Last}^i\} \cup \{\preceq^m\}, \end{aligned}$$

where the arity of the relations names  $\text{Symbol}_a$ ,  $\text{Head}$ , and  $\preceq^m$  is  $2m$ , of  $\text{State}_p$  is  $m$ , of  $\text{Succ}^i$  is  $2i$ , of  $\text{First}^i$  and  $\text{Last}^i$  is  $i$ , and of  $\text{Yes}$  is 0.

The intuitive meaning of the relation names of  $\text{idb}(\Pi)$ , apart from  $\text{Succ}^i$ ,  $\text{First}^i$ ,  $\text{Last}^i$ , and  $\preceq^m$  that have been discussed above, is as follows:

- $\text{Symbol}_a(t, c)$ : at time instant  $t$ , the tape cell  $c$  contains the symbol  $a$ .
- $\text{Head}(t, c)$ : at time instant  $t$ , the head points at cell  $c$ .
- $\text{State}_p(t)$ : at time instant  $t$ ,  $M$  is in state  $p$ .
- $\text{Yes}()$ :  $M$  has reached an accepting configuration.

Having  $\text{edb}(\Pi)$  and  $\text{idb}(\Pi)$  in place, we can now proceed with the definition of the database  $D$  and the Datalog program  $\Pi$ .

#### *The Database $D$*

We only need to store the relations  $\text{Succ}^1$ ,  $\text{First}^1$ , and  $\text{Last}^1$ , which form the base case of the inductive definition of  $\text{Succ}^i$ ,  $\text{First}^i$ , and  $\text{Last}^i$ . In particular,

$$D = \{\text{Succ}^1(0, 1), \text{First}^1(0), \text{Last}^1(1)\}.$$

*The Program  $\Pi$* 

The program  $\Pi$ , which is responsible for faithfully describing the evolution of  $M$  on  $w$  starting from  $sc(w)$ , is the union of the following five programs:

1.  $\Pi_{\preceq}$  that inductively constructs  $\text{Succ}^i$ ,  $\preceq^i$ , and  $\text{First}^i$ , for each  $i \in [m]$ .
2.  $\Pi_{\text{start}}$  that constructs the start configuration  $sc(w) = (s, \triangleright, w, \sqcup, \dots, \sqcup)$ .
3.  $\Pi_{\delta}$  that simulates the transition function of  $M$ .
4.  $\Pi_{\text{inertia}}$  that ensures that the tape cells that have not been changed at time instant  $t$  keep their values at time instant  $t + 1$ .
5.  $\Pi_{\text{accept}}$  that checks whether  $M$  has reached an accepting configuration.

The definitions of the above Datalog program follow. For notational convenience, we write  $\bar{x}$  for  $x_1, \dots, x_k$ , where  $k \geq 0$  will be clear from the context, and  $\bar{x}_i$  for  $x_{i,1}, \dots, x_{i,m}$ .

**The Program  $\Pi_{\preceq}$ .** For each  $i \in [m - 1]$ , we add the Datalog rules:

$$\begin{aligned}
 \text{Succ}^{i+1}(z, \bar{x}, z, \bar{y}) &:- \text{Succ}^i(\bar{x}, \bar{y}), \text{First}^1(z) \\
 \text{Succ}^{i+1}(z, \bar{x}, z, \bar{y}) &:- \text{Succ}^i(\bar{x}, \bar{y}), \text{Last}^1(z) \\
 \text{Succ}^{i+1}(z, \bar{x}, v, \bar{y}) &:- \text{Succ}^1(z, v), \text{Last}^i(\bar{x}), \text{First}^i(\bar{y}) \\
 \text{First}^{i+1}(x, \bar{y}) &:- \text{First}^1(x), \text{First}^i(\bar{y}) \\
 \text{Last}^{i+1}(x, \bar{y}) &:- \text{Last}^1(x), \text{Last}^i(\bar{y}) \\
 \preceq^m(\bar{x}, \bar{y}) &:- \text{Succ}^m(\bar{x}, \bar{y}) \\
 \preceq^m(\bar{x}, \bar{z}) &:- \preceq^m(\bar{x}, \bar{y}), \text{Succ}^m(\bar{y}, \bar{z}).
 \end{aligned}$$

**The Program  $\Pi_{\text{start}}$ .** Assuming that  $w = a_0, \dots, a_{|w|-1}$ , we add the rules:

$$\begin{aligned}
 \text{State}_s(\bar{x}) &:- \text{First}^m(\bar{x}) \\
 \text{Symbol}_{a_0}(\bar{x}, \bar{x}) &:- \text{First}^m(\bar{x}) \\
 \text{Symbol}_{a_1}(\bar{x}_0, \bar{x}_1) &:- \text{First}^m(\bar{x}_0), \text{Succ}^m(\bar{x}_0, \bar{x}_1) \\
 &\vdots \\
 \text{Symbol}_{a_{|w|-1}}(\bar{x}_0, \bar{x}_i) &:- \text{First}^m(\bar{x}_0), \text{Succ}^m(\bar{x}_0, \bar{x}_1), \dots, \text{Succ}^m(\bar{x}_{|w|-2}, \bar{x}_{|w|-1}) \\
 \text{Symbol}_{\sqcup}(\bar{x}_0, \bar{y}) &:- \text{First}^m(\bar{x}_0), \text{Succ}^m(\bar{x}_0, \bar{x}_1), \dots, \\
 &\quad \text{Succ}^m(\bar{x}_{|w|-2}, \bar{x}_{|w|-1}), \preceq^m(\bar{x}_{|w|-1}, \bar{y}) \\
 \text{Head}(\bar{x}, \bar{x}) &:- \text{First}^m(\bar{x})
 \end{aligned}$$

**The Program  $\Pi_{\delta}$ .** For each pair  $(p, a) \in (Q - \{\text{"yes"}, \text{"no"}\}) \times \Sigma$ , with  $\delta(p, a) = (p', b, \text{dir})$ , we add the following Datalog rules. For brevity, let

$$\Phi_{(p,a)}(\bar{x}, \bar{y}, \bar{z}) = \text{State}_p(\bar{x}), \text{Head}(\bar{x}, \bar{y}), \text{Symbol}_a(\bar{x}, \bar{y}), \text{Succ}^m(\bar{x}, \bar{z}).$$

The following rules change the state from  $p$  to  $p'$ , and the symbol from  $a$  to  $b$  at the next time instant of the computation:

$$\begin{aligned}\text{State}_{p'}(\bar{z}) &:- \Phi_{(p,a)}(\bar{x}, \bar{y}, \bar{z}) \\ \text{Symbol}_b(\bar{z}, \bar{y}) &:- \Phi_{(p,a)}(\bar{x}, \bar{y}, \bar{z}).\end{aligned}$$

The next rule, which is responsible for moving the head, depends on the direction  $\text{dir} \in \{\rightarrow, \leftarrow, -\}$ . In particular, if  $\text{dir} = \rightarrow$ , then we add the rule

$$\text{Head}(\bar{z}, \bar{v}) :- \Phi_{(p,a)}(\bar{x}, \bar{y}, \bar{z}), \text{Succ}^m(\bar{y}, \bar{v}).$$

If  $\text{dir} = \leftarrow$ , then we add the rule

$$\text{Head}(\bar{z}, \bar{v}) :- \Phi_{(p,a)}(\bar{x}, \bar{y}, \bar{z}), \text{Succ}^m(\bar{v}, \bar{y}).$$

Finally, if  $\text{dir} = -$ , then we add the rule

$$\text{Head}(\bar{z}, \bar{y}) :- \Phi_{(p,a)}(\bar{x}, \bar{y}, \bar{z}).$$

**The Program  $\Pi_{\text{inertia}}$ .** Recall that this program is responsible for, essentially, copying the content of the tape cells that have not been affected during the transition from time instant  $t$  to time instant  $t + 1$ . The following rule achieves this for the tape cells coming before the current cell

$$\text{Symbol}_a(\bar{v}, \bar{y}) :- \text{Symbol}_a(\bar{x}, \bar{y}), \text{Head}(\bar{x}, \bar{z}), \preceq^m(\bar{y}, \bar{z}), \text{Succ}^m(\bar{x}, \bar{v}).$$

The next rule does the same for the tape cells coming after the current cell

$$\text{Symbol}_a(\bar{x}, \bar{y}) :- \text{Symbol}_a(\bar{x}, \bar{y}), \text{Head}(\bar{x}, \bar{z}), \preceq^m(\bar{z}, \bar{y}), \text{Succ}^m(\bar{x}, \bar{v}).$$

**The Program  $\Pi_{\text{accept}}$ .** Finally, we check whether  $M$  has reached an accepting configuration via the Datalog rule

$$\text{Yes} :- \text{State}_{\text{yes}}(\bar{x}).$$

It is not difficult to verify that  $D$  and  $\Pi$  can be constructed from  $M$  and  $w$  in polynomial time. It is also not hard to see that  $\Pi$  faithfully describes the computation of  $M$  on input  $w$ . This means that, with  $q = (\Pi, \text{Yes})$ ,  $M$  accepts  $w$  if and only if  $q(D) = \text{true}$  (we leave the proof as an exercise).  $\square$

## Data Complexity

We now concentrate on the data complexity of Datalog-Evaluation. As discussed in Chapter 2, when we study the data complexity of query evaluation, we essentially consider the query to be fixed, and only the database and the candidate output are considered as input. Formally, we are interested in the



complexity of the problem  $q$ -Evaluation for a Datalog query  $q$ , which takes as input a database  $D$  and a tuple  $\bar{a}$  over  $\text{Dom}(D)$ , and asks whether  $\bar{a} \in q(D)$ . As usual, by convention, we say that Datalog-Evaluation is  $\mathcal{C}$ -complete in data complexity, for a complexity class  $\mathcal{C}$ , if  $q$ -Evaluation is in  $\mathcal{C}$  for every Datalog query  $q$ , and there exists a Datalog query  $q$  such that  $q$ -Evaluation is  $\mathcal{C}$ -hard. We show that fixing the query has an impact on the complexity of the problem, that is, Datalog-Evaluation, from provably intractable, becomes tractable.

**Theorem 37.2**

Datalog-Evaluation is PTIME-complete in data complexity.

*Proof.* The upper bound follows from the analysis performed at the beginning of the chapter. Fix a Datalog query  $q = (\Pi, R)$ . Given a database  $D$  of  $\text{edb}(\Pi)$ , and a tuple  $\bar{a}$  over  $\text{Dom}(D)$ , the analysis performed above shows that  $T_\Pi^\infty(D)$  can be computed in time  $O(|\text{Dom}(D)|^k)$  for some  $k \in \mathbb{N}$  that solely depends on  $q$ , which implies that checking whether  $R(\bar{a}) \in T_\Pi^\infty(D)$  is feasible in time  $O(|\text{Dom}(D)|^k)$ . Therefore,  $q$ -Evaluation is in PTIME, as needed.

For the lower bound we provide a reduction from a standard PTIME-hard problem known as *monotone circuit value*. For  $n \in \mathbb{N}$ , an  $n$ -input, single-output monotone Boolean circuit is a directed acyclic graph  $C$  with exactly  $n$  nodes without incoming edges, called the *sources*, and exactly one node without outgoing edges, called the *sink*. All the nodes that are not sources are labeled with either  $\wedge$  or  $\vee$  ( $\neg$  is not allowed, hence the term monotone). We write  $C(v_1, \dots, v_n)$  to indicate that the  $i$ -th source of  $C$  is the node  $v_i$ , for  $i \in [n]$ . An *input* to such a Boolean circuit  $C(v_1, \dots, v_n)$  is a tuple  $(w_1, \dots, w_n) \in \{0, 1\}^n$ . The *output of  $C(v_1, \dots, v_n)$  on  $(w_1, \dots, w_n)$* , denoted  $C(w_1, \dots, w_n)$ , is defined by recursively assigning to every node  $v$  a value  $b_v$  as follows:

- $b_{v_i} = w_i$ , for each  $i \in [n]$ , and
- for every node  $u \notin \{v_1, \dots, v_n\}$ , assuming that the two incoming edges of  $u$  are coming from  $u_1$  and  $u_2$ ,  $b_u = b_{u_1} \diamond b_{u_2}$ , where  $\diamond$  is the label of  $u$ .

The output  $C(w_1, \dots, w_n)$  is defined as  $b_{v_s}$ , where  $v_s$  is the sink of  $C$ . We are now ready to introduce the monotone circuit value problem:

**Problem: MCVP**

**Input:** An  $n$ -input, single-output monotone Boolean circuit  $C(\bar{v})$ , and a tuple  $\bar{w} \in \{0, 1\}^n$ , where  $n \in \mathbb{N}$   
**Output:** true if  $C(\bar{w}) = 1$ , and false otherwise

Our goal is to show that there exists a Datalog query  $q$  such that MCVP can be reduced to  $q$ -Evaluation via a reduction that is computable in deterministic

logarithmic space. Intuitively, the query  $q$  should specify a generic procedure for evaluating monotone Boolean circuits. This can be straightforwardly done via the query  $q = (\Pi, \text{Yes})$ , where  $\Pi$  is the Datalog program

$$\begin{aligned} \text{True}(x) &:- \text{Or}(x, y, z), \text{True}(y) \\ \text{True}(x) &:- \text{Or}(x, y, z), \text{True}(z) \\ \text{True}(x) &:- \text{And}(x, y, z), \text{True}(y), \text{True}(z) \\ \text{Yes} &:- \text{Sink}(x), \text{True}(x). \end{aligned}$$

We now proceed to show that **MCVP** can be reduced to  $q$ -**Evaluation** via a reduction that is computable in deterministic logarithmic space. Consider an instance of **MCVP**, i.e., an  $n$ -input, single-output monotone Boolean circuit  $C(v_1, \dots, v_n)$ , and a tuple  $\bar{w} = (w_1, \dots, w_n) \in \{0, 1\}^n$ , for an integer  $n \in \mathbb{N}$ . For brevity, we write  $u_i = u_j \wedge u_k$  for the fact that the node  $u_i$  is labeled by  $\wedge$ , and its incoming edges are coming from the nodes  $u_j$  and  $u_k$ ; analogously, we write  $u_i = u_j \vee u_k$ . We define the database  $D_{C, \bar{w}}$  of **edb**( $\Pi$ ) as follows:

$$\begin{aligned} &\{\text{True}(v_i) \mid i \in [n] \text{ and } w_i = 1\} \\ &\cup \{\text{And}(u_i, u_j, u_k) \mid u_i, u_j, u_k \text{ are nonsource nodes of } C, \text{ and } u_i = u_j \wedge u_k\} \\ &\cup \{\text{Or}(u_i, u_j, u_k) \mid u_i, u_j, u_k \text{ are nonsource nodes of } C, \text{ and } u_i = u_j \vee u_k\} \\ &\cup \{\text{Sink}(v_s) \mid v_s \text{ is the sink of } C\}. \end{aligned}$$

It is clear that the database  $D_{C, \bar{w}}$  can be computed in deterministic logarithmic space in the size of  $C$  and  $\bar{w}$ . It is also easy to verify that

$$C(\bar{w}) = 1 \quad \text{if and only if} \quad q(D_{C, \bar{w}}) = \mathbf{true}.$$

Therefore,  $q$ -**Evaluation** is **P**TIME-hard, and the claim follows.  $\square$



## Static Analysis of Datalog Queries

In this chapter, we discuss central static analysis tasks for Datalog queries that are important for query optimization purposes. In fact, we consider the three fundamental tasks that we have also studied for first-order and conjunctive queries, namely satisfiability, containment, and equivalence. We also discuss a new static analysis task, known as boundedness, that is relevant for recursive query languages such as Datalog. In simple words, a Datalog query is bounded if it can be equivalently rewritten as a Datalog query without recursion.

As we shall see, we can effectively check whether a Datalog query is satisfiable. On the other hand, the problem of checking whether a Datalog query is contained into (or is equivalent to) another Datalog query, as well as the problem of checking whether a Datalog query is bounded, are undecidable.

### Satisfiability

We start by considering the satisfiability problem: given a Datalog query  $q = (\Pi, R)$ , is there a database  $D$  of  $\text{edb}(\Pi)$  such that  $q(D) \neq \emptyset$ ? We proceed to show that this problem is decidable. We call a Datalog query  $(\Pi, R)$  *constant-free* if the rules of the Datalog program  $\Pi$  do not mention constants.

#### Theorem 38.1

Datalog-Satisfiability is in EXPTIME, and in PTIME for constant-free Datalog queries.

*Proof.* Consider a Datalog query  $q = (\Pi, R)$ , and assume that  $a_1, \dots, a_n \in \text{Const}$  are the constants mentioned by  $\Pi$ . We first characterize the satisfiability of  $q$  via a very simple database. Let

$$D_\Pi = \left\{ P(\bar{c}) \mid P \in \text{edb}(\Pi) \text{ and } \bar{c} \in \{\star, a_1, \dots, a_n\}^{\text{ar}(P)} \right\},$$

where  $\star$  is a value from  $\text{Const} - \{a_1, \dots, a_n\}$ . We can show the following lemma:

**Lemma 38.2.** *It holds that  $q$  is satisfiable if and only if  $q(D_\Pi) \neq \emptyset$ .*

*Proof.* It is clear that if  $q(D_\Pi) \neq \emptyset$ , then  $q$  is satisfiable. Assume now that  $q$  is satisfiable. This implies that there exists a database  $D$  of  $\text{edb}(\Pi)$  such that  $q(D) \neq \emptyset$ , i.e., there exists a tuple  $\bar{c}$  over  $\text{Dom}(D)$  with  $R(\bar{c}) \in \Pi(D)$ . Let  $h : \text{Dom}(D) \rightarrow \{\star, a_1, \dots, a_n\}$  be the function that maps each element in  $\text{Dom}(D) - \{a_1, \dots, a_n\}$  to  $\star$  and each element in  $\text{Dom}(D) \cap \{a_1, \dots, a_n\}$  to itself. Clearly,  $h(D) \subseteq D_\Pi$ , which in turn allows us to show that  $h(\Pi(D)) \subseteq \Pi(D_\Pi)$ ; the latter can be shown via an easy inductive argument. Therefore,  $R(h(\bar{c})) \in \Pi(D_\Pi)$ , which implies that  $h(\bar{c}) \in q(D_\Pi)$ . Thus,  $q(D_\Pi) \neq \emptyset$ .  $\square$

Lemma 38.2 leads to the following simple procedure for checking whether the Datalog query  $q = (\Pi, R)$  is satisfiable:

if  $q(D_\Pi) \neq \emptyset$ , then yes; otherwise, no.

From the analysis performed in Chapter 37 (see the discussion before Theorem 37.1), we conclude that checking whether  $q(D_\Pi) \neq \emptyset$  is feasible in time

$$O\left(|\Pi| \cdot |\text{Dom}(D_\Pi)|^{\text{maxvar}} \cdot \text{maxbody} \cdot |\text{sch}(\Pi)|^2 \cdot |\text{Dom}(D_\Pi)|^{2 \cdot \text{ar}(\Pi)}\right),$$

where *maxvar* and *maxbody* are the maximum number of variables and body atoms, respectively, in a rule of  $\Pi$ . This implies that Datalog-Satisfiability is in EXPTIME, as needed.

Now, observe that in the case where  $\Pi$  is constant-free, we have that

$$D_\Pi = \{P(\star, \dots, \star) \mid P \in \text{edb}(\Pi)\}.$$

Therefore, checking whether  $q(D_\Pi) \neq \emptyset$  is feasible in time

$$O(|\Pi| \cdot \text{maxbody} \cdot |\text{sch}(\Pi)|^2)$$

since  $|\text{Dom}(D_\Pi)| = 1$ . This implies that, for constant-free Datalog queries, Datalog-Satisfiability is in PTIME, and the claim follows.  $\square$

## Containment and Equivalence

We now concentrate on the containment problem for Datalog: given two Datalog queries  $q_1 = (\Pi_1, R_1)$  and  $q_2 = (\Pi_2, R_2)$  over a schema  $\mathbf{S}$  (in particular, with  $\text{edb}(\Pi_1) = \text{edb}(\Pi_2)$ ), is it the case that  $q_1 \subseteq q_2$ . We can show that:

### Theorem 38.3

Datalog-Containment is undecidable.

The above result is shown via a reduction from a known undecidable problem, namely containment for context-free grammars. A context-free grammar is a set of production rules that describe how to produce words over a certain alphabet. Consider, for example, the grammar  $G$  consisting of the rules

$$S \rightarrow AA \quad A \rightarrow a \quad A \rightarrow b.$$

The first rule states that we can replace  $S$  with  $AA$ , while the other two rules state that  $A$  can be replaced with  $a$  or  $b$ . Assuming that  $S$  is the starting point of the production, and  $\{a, b\}$  is the underlying alphabet, the above grammar produces the words  $aa$ ,  $ab$ ,  $ba$ ,  $bb$ . Let us formalize the above discussion.

A *context-free grammar* (CFG) is a tuple  $(N, T, P, S)$ , where

- $N$  is a finite set, the *non-terminal symbols*,
- $T$  is a finite set disjoint from  $N$ , the *terminal symbols*,
- $P$  is a finite subset of  $N \times (N \cup T)^*$ , the *production rules*, and
- $S \in N$ , the *start symbol*.

For any two words  $v, w \in (N \cup T)^*$ , we say that  $v$  *directly yields*  $w$ , written  $v \Rightarrow w$ , if there exists  $(x, y) \in P$  and  $z_1, z_2 \in (N \cup T)^*$  such that  $v = z_1xz_2$  and  $w = z_1yz_2$ . Now, for any two words  $v, w \in (N \cup T)^*$ , we say that  $v$  *yields*  $w$ , written as  $v \Rightarrow^* w$ , if there exists  $k \geq 1$ , and words  $z_1, \dots, z_k \in (N \cup T)^*$  such that  $v = z_1 \Rightarrow z_2 \cdots \Rightarrow z_k = w$ . The *language* of  $G$ , denoted  $L(G)$ , is the set of words  $\{w \in T^* \mid S \Rightarrow^* w\}$ , that is, all the words  $w$  over  $T$  that can be obtained starting from the symbol  $S$  and applying production rules of  $P$ .

Given two context-free grammars  $G_1$  and  $G_2$ , we say that  $G_1$  is *contained* in  $G_2$ , denoted  $G_1 \subseteq G_2$ , if  $L(G_1) \subseteq L(G_2)$ . The containment problem for context-free grammars is defined as expected:

**Problem: CFG-Containment**

**Input:** Two context-free grammars  $G_1$  and  $G_2$

**Output:** **true** if  $G_1 \subseteq G_2$ , and **false** otherwise

*Proof (of Theorem 38.3).* We provide a reduction from CFG-Containment to Datalog-Containment. In other words, given two context-free grammars  $G_1$  and  $G_2$ , the goal is to construct two Datalog queries  $q_1$  and  $q_2$  such that  $G_1 \subseteq G_2$  if and only if  $q_1 \subseteq q_2$ . We first explain how to transform a CFG into a Datalog query. Let us clarify that, in what follows, given a CFG  $G = (N, T, P, S)$ , we assume that  $P$  is a finite subset of  $N \times ((N - \{S\}) \cup T)^* - \{\epsilon\}$ , where  $\epsilon$  denotes the empty string. In other words, there is no rule in  $P$  that produces the empty string, and the start symbol does not occur in the right-hand side of a rule. This does not affect the generality of our proof since CFG-Containment remains undecidable even with the above simplifying assumptions.

We proceed to define the Datalog program  $\Pi_G$ , where  $G = (N, T, P, S)$  is a CFG. The extensional schema  $\text{edb}(\Pi_G)$  and intensional schema  $\text{idb}(\Pi_G)$  are

$$\{\text{Symbol}_A \mid A \in T\} \quad \text{and} \quad \{\text{Symbol}_A \mid A \in N\},$$

respectively, where all the relations are binary. For each production rule in  $P$

$$(A, A_1 \cdots A_n),$$

for  $n \geq 1$ , we add to  $\Pi_G$  the Datalog rule

$$\begin{aligned} \text{Symbol}_A(x_1, x_{n+1}) \quad :- \quad & \text{Symbol}_{A_1}(x_1, x_2), \text{Symbol}_{A_2}(x_2, x_3), \dots, \\ & \text{Symbol}_{A_n}(x_n, x_{n+1}). \end{aligned}$$

We finally define the Datalog query  $q_G = (\Pi_G, \text{Symbol}_S)$ .

#### Example 38.4: From CFG to Datalog

Consider the CFG  $G = (N, T, P, S)$ , where  $N = \{A, S\}$ ,  $T = \{a, b\}$ , and  $P = \{(S, Aa), (A, abA), (A, aa)\}$ . The Datalog program  $\Pi_G$  is

$$\begin{aligned} \text{Symbol}_S(x_1, x_3) \quad & :- \quad \text{Symbol}_A(x_1, x_2), \text{Symbol}_a(x_2, x_3) \\ \text{Symbol}_A(x_1, x_4) \quad & :- \quad \text{Symbol}_a(x_1, x_2), \text{Symbol}_b(x_2, x_3), \text{Symbol}_A(x_3, x_4) \\ \text{Symbol}_A(x_1, x_3) \quad & :- \quad \text{Symbol}_a(x_1, x_2), \text{Symbol}_a(x_2, x_3), \end{aligned}$$

while the query  $q_G = (\Pi_G, \text{Symbol}_S)$ .

To show the correctness of the above construction, we need to introduce the notion of proof tree of an atom from a Datalog program. Roughly speaking, such a proof tree explains how an atom can be derived from a Datalog program, i.e., it provides a proof for that atom. As we shall see below, this notion is closely related to the notion of derivation tree in context-free languages that essentially explains how a word can be derived from a CFG.

Consider an atom  $R(\bar{a})$ , with  $\bar{a} \in \text{Const}^{\text{ar}(R)}$ , and a Datalog program  $\Pi$ . A *proof tree of  $R(\bar{a})$  from  $\Pi$*  is a labeled rooted tree  $T = (V, E, \lambda)$ , where  $\lambda$  is a function from  $V$  to the set of atoms that can be formed using relations from  $\text{sch}(\Pi)$  and constants from  $\text{Const}$ , such that

1. assuming that  $v \in V$  is the root node of  $T$ ,  $\lambda(v) = R(\bar{a})$ , and
2. for each internal node  $v \in V$  with children  $u_1, \dots, u_n$  for  $n \geq 1$ , there exists a rule  $\rho \in \Pi$  of the form  $R_0(\bar{x}_0) :- R_1(\bar{x}_1), \dots, R_n(\bar{x}_n)$ , and a function  $h$  from the constants and variables in  $\rho$  to  $\text{Const}$ , which is the identity on  $\text{Const}$ , such that  $\lambda(v) = R_0(h(\bar{x}_0))$  and  $\lambda(u_i) = R_i(h(\bar{x}_i))$ , for each  $i \in [n]$ .

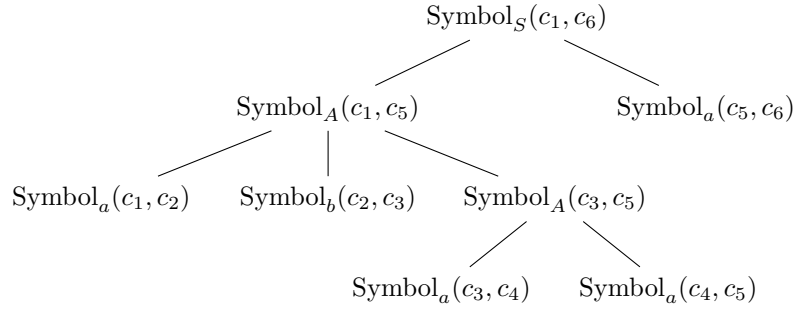
We say that the sequence of atoms  $R_1(\bar{a}_1), \dots, R_n(\bar{a}_n)$  is *induced* by the proof tree  $T$  if, assuming that the leaf nodes of  $T$  are  $v_1, \dots, v_n$  (in this order), then  $\lambda(v_i) = R_i(\bar{a}_i)$ , for each  $i \in [n]$ . Given a database  $D$  of  $\text{edb}(\Pi)$ , a *proof tree of  $R(\bar{a})$  from  $\Pi$  and  $D$*  is a proof tree  $T = (V, E, \lambda)$  of  $R(\bar{a})$  from  $\Pi$  such that, for each  $v \in V$ ,  $\lambda(v) \in \mathbf{B}(\Pi, D)$ , i.e.,  $\lambda(v)$  is an atom with a relation from  $\text{sch}(\Pi)$  and constants from  $\text{Dom}(D)$ , and for each leaf node  $v$  of  $T$ ,  $\lambda(v) \in D$ .

### Example 38.5: Proof Tree

Consider the Datalog program  $\Pi_G$  obtained from the CFG  $G$  as in Example 38.4. A proof tree of the atom  $\text{Symbol}_S(c_1, c_6)$  from  $\Pi_G$  and

$$D \supseteq \{\text{Symbol}_a(c_1, c_2), \text{Symbol}_b(c_2, c_3), \text{Symbol}_a(c_3, c_4), \\ \text{Symbol}_a(c_4, c_5), \text{Symbol}_a(c_5, c_6)\}$$

is the following one



Clearly, the above proof tree induces the sequence of atoms

$$\text{Symbol}_a(c_1, c_2), \text{Symbol}_b(c_2, c_3), \text{Symbol}_a(c_3, c_4), \\ \text{Symbol}_a(c_4, c_5), \text{Symbol}_a(c_5, c_6).$$

It is an easy exercise to show the following lemma:

**Lemma 38.6.** *Consider a Datalog query  $q = (\Pi, R)$ , a database  $D$  of  $\text{edb}(\Pi)$ , and a tuple  $\bar{a} \in \text{Dom}(D)^{\text{ar}(R)}$ . The following are equivalent:*

1.  $\bar{a} \in q(D)$ .
2. There exists a proof tree of  $R(\bar{a})$  from  $\Pi$  and  $D$ .

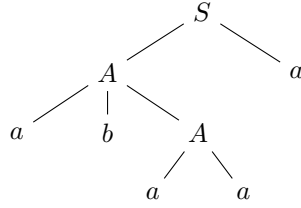
The next technical lemma makes apparent the intention underlying the transformation of a CFG  $G$  to a Datalog program  $\Pi_G$ .



**Lemma 38.7.** Consider a CFG  $G = (N, T, P, S)$ , and two words  $a_1 \cdots a_n \in T^*$  and  $c_1 \cdots c_{n+1} \in \text{Const}^*$ , for  $n \geq 1$ . The following are equivalent:

1.  $a_1 \cdots a_n \in L(G)$ .
2. There exists a proof tree of  $\text{Symbol}_S(c_1, c_{n+1})$  from  $\Pi_G$  that induces the sequence of atoms  $\text{Symbol}_{a_1}(c_1, c_2), \dots, \text{Symbol}_{a_n}(c_n, c_{n+1})$ .

The proof of the above lemma, which is left as an exercise, relies on the correspondence between proof trees and derivation trees in context-free languages. For example, the following tree is a derivation tree of the word  $abaaa$  from the CFG  $G$  given in Example 38.4



The correspondence between the proof tree of the atom  $\text{Symbol}_S(c_1, c_6)$  given above, and the derivation tree of the word  $abaaa$  should be apparent. By exploiting the above technical lemmas, we can now show the following result:

**Proposition 38.8**

Consider the CFGs  $G_i = (N_i, T_i, P_i, S_i)$  for  $i \in \{1, 2\}$ . Then

$$G_1 \subseteq G_2 \quad \text{if and only if} \quad q_{G_1} \subseteq q_{G_2}.$$

*Proof.* We only show the ‘if’ direction; the ‘only if’ direction is shown analogously. Consider a database  $D$  of  $\text{edb}(\Pi_{G_1})$ , and assume that  $(c, d) \in q_{G_1}(D)$ . By Lemma 38.6, there exists a proof tree of  $\text{Symbol}_{S_1}(c, d)$  from  $\Pi_{G_1}$  and  $D$ . Assume that this proof tree induces the sequence of atoms

$$s = \text{Symbol}_{a_1}(c_1, c_2), \text{Symbol}_{a_2}(c_2, c_3), \dots, \text{Symbol}_{a_n}(c_n, c_{n+1}),$$

where  $c = c_1$ ,  $d = c_{n+1}$ , and  $a_1 \cdots a_n \in T_1^*$ . By Lemma 38.7, we conclude that  $a_1 \cdots a_n \in L(G_1)$ . Since, by hypothesis,  $G_1 \subseteq G_2$ , we get that  $a_1 \cdots a_n \in L(G_2)$ . By Lemma 38.7, there exists a proof tree of  $\text{Symbol}_{S_2}(c_1, c_{n+1})$  from  $\Pi_{G_2}$  that induces the sequence of atoms  $s$ . Since the atoms in  $s$  are atoms of  $D$ , Lemma 38.6 implies that  $(c, d) \in q_{G_2}(D)$ , and the claim follows.  $\square$

Since CFG-Containment is undecidable, Proposition 38.8 implies the same for Datalog-Containment. This completes the proof of Theorem 38.3.  $\square$

Let us now turn our attention on the equivalence problem: given two Datalog queries  $q$  and  $q'$ , is it the case that  $q \equiv q'$ . By exploiting the fact that the containment problem is undecidable, we can easily show the following:

**Theorem 38.9**

Datalog-Equivalence is undecidable.

*Proof.* It suffices to reduce Datalog-Containment to Datalog-Equivalence. Consider two Datalog queries  $q_1 = (\Pi_1, R_1)$  and  $q_2 = (\Pi_2, R_2)$ , where  $\text{edb}(\Pi_1) = \text{edb}(\Pi_2)$ . We assume, without loss of generality, that  $\text{idb}(\Pi_1) \cap \text{idb}(\Pi_2) = \emptyset$ . We define the Datalog query

$$q_{12} = (\Pi_1 \cup \Pi_2 \cup \{R_{12}(\bar{x}) :- R_1(\bar{x}), R_{12}(\bar{x}) :- R_2(\bar{x})\}, R_{12}),$$

where  $R_{12}$  is a new relation name not occurring in  $\text{sch}(\Pi_1) \cup \text{sch}(\Pi_2)$ . It is easy to verify that  $q_1 \subseteq q_2$  if and only if  $q_{12} \equiv q_2$ , and the claim follows.  $\square$

**Boundedness**

As discussed in Chapter 35, given a Datalog program  $\Pi$ , and a database  $D$  of  $\text{edb}(\Pi)$ , the semantics of  $\Pi$  on  $D$ , i.e., the database  $\Pi(D)$ , can be computed by repeatedly applying the immediate consequence operator  $T_\Pi$  of  $\Pi$  starting from  $D$  until a fixpoint is reached; in fact, by Corollary 35.14,  $\Pi(D) = T_\Pi^\infty(D)$ . We have also seen that the construction of  $T_\Pi^\infty(D)$  does not require infinitely many iterations. Actually, there exists an integer  $k \leq |\mathbf{B}(\Pi, D)|$  such that  $T_\Pi^\infty(D) = T_\Pi^k(D)$ . The smallest integer  $k$  such that  $T_\Pi^\infty(D) = T_\Pi^k(D)$  is called the *stage of  $\Pi$  and  $D$* , denoted  $\text{stage}(\Pi, D)$ .

Given a Datalog program  $\Pi$ , it is generally the case that, for some arbitrary database  $D$  of  $\text{edb}(\Pi)$ , the integer  $\text{stage}(\Pi, D)$  depends on both  $\Pi$  and  $D$ . This essentially means that the Datalog program  $\Pi$  is inherently recursive, or, in other words, the depth of recursion of  $\Pi$  is unbounded. On the other hand, if there is a uniform upper bound (i.e., a bound that depends only on  $\Pi$ ) for  $\text{stage}(\Pi, D)$ , then the recursion of  $\Pi$  is bounded, which actually means that  $\Pi$  is non-recursive despite the fact that syntactically it may look recursive. The following example illustrates that a bounded (seemingly recursive) Datalog program can be replaced by an equivalent non-recursive Datalog program.

**Example 38.10: Program Boundedness**

Consider the Datalog program  $\Pi$  consisting of the rules

$$P(x, y) :- R(x), P(z, y) \quad P(x, y) :- S(x, y).$$

Notice that  $\Pi$  is syntactically recursive due to the first rule ( $P$  depends on itself). However,  $\Pi$  is bounded, and equivalent to the program

$$P(x, y) :- R(x), S(z, y) \quad P(x, y) :- S(x, y),$$

that is non-recursive. We can safely replace the atom  $P(z, y)$  in the first rule with the atom  $S(z, y)$ , which leads to a non-recursive program, since it does not share any variable with the atom  $R(x)$ .

Boundedness for Datalog programs is defined as expected:

**Definition 38.11: Program Boundedness**

A Datalog program  $\Pi$  is *bounded* if there exists  $k \in \mathbb{N}$  such that  $\text{stage}(\Pi, D) \leq k$ , for every database  $D$  of  $\text{edb}(\Pi)$ .

As explained above, boundedness essentially removes from Datalog programs the feature of recursion. Therefore, it should not come as a surprise the fact that a Datalog query  $(\Pi, R)$ , where  $\Pi$  is bounded, can always be written as a UCQ. An interesting question at this point is whether the opposite holds, namely whenever a Datalog query  $(\Pi, R)$  is equivalent to a UCQ, then  $\Pi$  is bounded. It is easy to see that, in general, this is not the case.

**Example 38.12: Program Boundedness and UCQ s**

Consider the Datalog query  $q = (\Pi, R)$ , where  $\Pi$  consists of the rules

$$\begin{aligned} P(x, y) &:- S(x, y) \\ P(x, y) &:- P(x, z), S(z, y) \\ R(x, y) &:- S(x, y) \\ R(x, y) &:- T(x, y). \end{aligned}$$

It is clear that  $\Pi$  is not bounded due to the first two rules that compute the transitive closure of the binary relation  $P$ . On the other hand,  $q$  is equivalent to the UCQ  $q' = \varphi(x, y)$  with

$$\varphi = S(x, y) \vee T(x, y).$$

Indeed, for every database  $D$  of  $\text{edb}(\Pi) = \{S, T\}$ ,  $q(D) = q'(D)$  since the relation name  $R$  depends only on  $S$  and  $T$ .

Observe that the key reason why the query  $q = (\Pi, R)$  from Example 38.12 can be written as a UCQ, despite the fact that  $\Pi$  is not bounded, is because the relation name  $R$  does not depend on a recursive relation name, but only on non-recursive ones (in this case, on the extensional relation names  $S$  and  $T$ ). This leads to the notion of boundedness of relation names.

Given a Datalog program  $\Pi$ , and a database  $D$  of  $\text{edb}(\Pi)$ , analogously to the stage of  $\Pi$  and  $D$ , for a relation name  $R \in \text{idb}(\Pi)$  we define the *stage of  $R$  with respect to  $\Pi$  and  $D$*  as the smallest integer  $k$  with  $R^{T_\Pi^\infty(D)} =$

$R^{T_H^k(D)}$ , denoted  $\text{stage}_{\Pi,D}(R)$ . Considering again the Example 38.12, although the stage of  $\Pi$  and  $D$  is not bounded,  $\text{stage}_{\Pi,D}(R) = 1$ . This essentially tells that, even though the program  $\Pi$  may be inherently recursive, the part of it that is responsible for computing the relation  $R$  is actually non-recursive. We can now define when a Datalog query (instead of a program) is bounded.

**Definition 38.13: Query Boundedness**

A Datalog query  $q = (\Pi, R)$  is *bounded* if there exists a  $k \in \mathbb{N}$  such that  $\text{stage}_{\Pi,D}(R) \leq k$ , for every database  $D$  of  $\text{edb}(\Pi)$ .

The notion of query boundedness essentially removes from Datalog queries the feature of recursion. Therefore, it should be expected that a bounded Datalog query  $(\Pi, R)$  can always be written as a UCQ, even if  $\Pi$  is not bounded. What is more interesting, though, is the fact that query boundedness, unlike program boundedness, characterises the fragment of Datalog queries that can be written, not only as UCQs, but actually as FO queries.

**Theorem 38.14**

Let  $q = (\Pi, R)$  be a Datalog query over  $\mathbf{S}$ . The following are equivalent:

1.  $q$  is bounded.
2. There exists an FO query  $q'$  over  $\mathbf{S}$  such that  $q \equiv q'$ .
3. There exists a UCQ  $q''$  over  $\mathbf{S}$  such that  $q \equiv q''$ .

*Proof.* We discuss how (1) implies (2) can be shown, and leave the formal proof as an exercise. We know from Proposition 36.2 that there is an infinitary UCQ  $q'$  over  $\mathbf{S}$  such that  $q \equiv q'$ . Recall that  $q'$  is defined by exhaustively applying the  $\text{Unfold}_{\Pi}(\cdot)$  operator, introduced in Chapter 36, starting from a certain CQ  $q_R$  obtained from  $q$ , and then keeping only the infinitary CQs over  $\mathbf{S}$ .<sup>1</sup> Now, in the case of bounded Datalog queries, it can be shown that  $q'$  is an ordinary UCQ over  $\mathbf{S}$ , and thus, an FO query over  $\mathbf{S}$ . This is due to the fact that the  $\text{Unfold}_{\Pi}(\cdot)$  operator, starting from  $q_R$ , it constructs a CQ over  $\mathbf{S}$  after finitely many steps, in fact, in at most  $k$  steps, where  $k \geq 0$  is the integer that bounds  $\text{stage}_{\Pi,D}(R)$ , for every database  $D$  of  $\text{edb}(\Pi)$ .

We now proceed to show the direction (2) implies (3). By hypothesis, there is an FO query  $q'$  over  $\mathbf{S}$  such that  $q \equiv q'$ . By Corollary 36.6, we get that  $q'$  is preserved under homomorphisms. This in turn implies, due to Theorem 28.11, that there is a UCQ with variable-constant equality  $\hat{q}$  over  $\mathbf{S}$  such that  $q' \equiv \hat{q}$ . It remains to explain how the equational atoms in  $\hat{q}$  can be eliminated. Since

<sup>1</sup> Note that the proof of Proposition 36.2 was given for Boolean queries, but it can be extended to non-Boolean queries; this extension was left as an exercise.

$q \equiv \hat{q}$ , by Proposition 35.17 we get that, for every database  $D$  of  $\mathbf{S}$ ,  $\hat{q}(D)$  consists of tuples over  $\text{Dom}(D)$ , i.e., it is not possible to have a value in the output of  $\hat{q}$  on  $D$  that occurs in  $\hat{q}$  but not in  $D$ . Assuming that  $\hat{q} = \varphi(\bar{x})$ , we then conclude the following: if a variable  $y$  in  $\varphi$  occurs in an equational atom, but not in a relational atom, then  $y$  is not among the free variables of  $\varphi$ , i.e.,  $y$  is not mentioned in  $\bar{x}$ . This observation allows us to eliminate an equational atom ( $y = a$ ) from  $\varphi$  by simply replacing each occurrence of  $y$  with  $a$ , and then removing ( $y = a$ ) from  $\varphi$ . This eventually leads to a UCQ  $q''$  over  $\mathbf{S}$  such that  $\hat{q} \equiv q''$ , and thus,  $q \equiv q''$ , as needed.

We finally show that (3) implies (1). Consider an arbitrary database  $D$  of  $\text{edb}(\Pi)$ , and a tuple  $\bar{a}$  over  $\text{Dom}(D)$ . It suffices to show that  $R(\bar{a}) \in \Pi(D)$  implies  $R(\bar{a}) \in T_{\Pi}^k(D)$  for some  $k \in \mathbb{N}$  that does not depend on  $D$  and  $\bar{a}$ . Indeed, if this is the case, then  $\text{stage}_{\Pi,D}(R) \leq k$ , for every database  $D$  of  $\text{edb}(\Pi)$ , which means that  $q$  is bounded. By hypothesis, there exists a UCQ  $q''$ , let say of the form  $q_1 \cup \dots \cup q_n$ , such that  $q \equiv q''$ . Therefore, if  $R(\bar{a}) \in \Pi(D)$ , which means that  $\bar{a} \in q(D)$ , then  $\bar{a} \in q''(D)$ . Let  $m$  be the maximum number of atoms occurring in the CQs  $q_1, \dots, q_n$ , that is,  $m = \max_{i \in [n]} \{|A_{q_i}|\}$ . It is clear that there is a database  $D' \subseteq D$  with  $|D'| \leq m$  such that  $\bar{a} \in q''(D')$ . Since  $q \equiv q''$ , we have that  $R(\bar{a}) \in \Pi(D')$  or  $R(\bar{a}) \in T_{\Pi}^{\infty}(D')$ . Recall that  $T_{\Pi}^{\infty}(D') = T_{\Pi}^{|\text{B}(\Pi, D')|}(D')$  with  $|\text{B}(\Pi, D')| \leq |\text{sch}(\Pi)| \cdot |\text{Dom}(D')|^{\text{ar}(\Pi)}$ , where  $\text{ar}(\Pi)$  is the maximum arity over all relation names of  $\text{sch}(\Pi)$ . Since  $|D'| \leq m$ , we get that  $|\text{Dom}(D')| \leq m \cdot \text{ar}(\Pi)$ . Observe that  $T_{\Pi}^{|\text{B}(\Pi, D')|}(D') \subseteq T_{\Pi}^{|\text{B}(\Pi, D')|}(D)$ . Therefore,  $R(\bar{a}) \in T_{\Pi}^k(D)$  for some  $k \in \mathbb{N}$  that does not depend on  $D$  and  $\bar{a}$  (it only depends on  $q''$  and  $\Pi$ ), and the claim follows.  $\square$

It is clear that checking whether a Datalog program or query is bounded are important static analysis tasks that are relevant for optimization purposes.

#### Problem: Datalog-PBoundedness

**Input:** A Datalog program  $\Pi$

**Output:** **true** if  $\Pi$  is bounded, and **false** otherwise

#### Problem: Datalog-QBoundedness

**Input:** A Datalog query  $q$

**Output:** **true** if  $q$  is bounded, and **false** otherwise

It turns out that both problems are undecidable. It can be shown via a reduction from the *Post Correspondence Problem*, a classical undecidable problem, that checking whether a Datalog program is bounded is undecidable. This can be then easily transferred to query boundedness via an easy reduction.

Consider a Datalog program  $\Pi$ . We define the Datalog query  $q = (\Pi \cup \{\rho\}, R)$ , where, assuming that  $\text{idb}(\Pi) = \{P_1, \dots, P_n\}$ , the Datalog rule  $\rho$  is

$$R(x_1^1, \dots, x_{\text{ar}(P_1)}^1, \dots, x_1^n, \dots, x_{\text{ar}(P_n)}^n) :- P_1(x_1^1, \dots, x_{\text{ar}(P_1)}^1), \dots, \\ P_n(x_1^n, \dots, x_{\text{ar}(P_n)}^n)$$

and  $R$  is a  $(\text{ar}(P_1) + \dots + \text{ar}(P_n))$ -ary relation name not occurring in  $\text{idb}(\Pi)$ . It is easy to see that  $\Pi$  is bounded iff  $q$  is bounded. We then have that:

**Theorem 38.15**

Datalog-PBoundedness and Datalog-QBoundedness are undecidable.



---

## Exercises

**Exercise 4.1.** Prove that query evaluation for  $\exists\text{FO}^+$  and  $\text{RA}^+$  is in NP.

**Exercise 4.2.** Prove that SPJU-Containment containment is  $\Pi_2^p$ -hard.

**Exercise 4.3.** Prove that containment for  $\exists\text{FO}^+$  and  $\text{RA}^+$  is  $\Pi_2^p$ -complete.

**Exercise 4.4.** Prove that  $\text{CQ}^\neq$ -Containment is  $\Pi_2^p$ -hard.

**Exercise 4.5.** Let  $q, q'$  be  $\text{CQ}^\neq$ s. Prove that  $q \subseteq q'$  if, and only if, for every  $i \in [n]$ , there exists  $j \in [m]$  such that  $q_i \subseteq q'_j$ .

**Exercise 4.6.** The language  $\text{CQ}^<$  is defined in the same way as  $\text{CQ}^\neq$  (see Definition 30.1), but instead of  $\neq$  we use  $<$ , assuming that there is an order on the set of constant **Const** from which database entries are drawn. Analogously, we can define the language  $\text{UCQ}^<$ . Prove that the problem of containment remains decidable for  $\text{CQ}^<$  and  $\text{UCQ}^<$ .

**Exercise 4.7.** The class BCCQ consists of Boolean combinations of CQs, i.e., queries obtained by repeatedly applying the operations of union ( $q_1 \cup q_2$ ), intersection ( $q_1 \cap q_2$ ), and difference ( $q_1 - q_2$ ) to CQs of the same arity with the obvious semantics. Prove that containment for BCCQs is decidable.

**Exercise 4.8.** Prove that the sentences  $\psi_k(R\bar{S})$  and  $\psi_k(\bar{R}S)$ , used in the proof of Theorem 32.5 are true in almost all databases of  $\mathbf{S} = \{R[1], S[1]\}$ .

**Exercise 4.9.** In general, Theorem 32.5 (0–1 law) does not hold if we focus on a restricted class  $C$  of databases, i.e., for an FO sentence  $\varphi$ ,  $\mu_n(\varphi)$  is defined by considering only databases from the class  $C$ . Let  $C_\subseteq$  be the class of databases  $D$  of the schema  $\mathbf{S} = \{R[1], S[1]\}$  such that  $S^D \subseteq R^D$ . Adapt the proof of Theorem 32.5, given in Chapter 30, to show that the 0–1 law holds even if we focus on the class of databases  $C_\subseteq$ . Use this result to show that the *parity* query  $q$  over  $\mathbf{S}$ , which checks whether the cardinality of the relation  $B$  is even, is not expressible as an FO query if we focus on the class of databases  $C_\subseteq$ .



**Exercise 4.10.** Use Exercise 4.9 (not just the result but the proof that you produced) to infer the following: for every Boolean FO query  $q$  over the schema  $\mathbf{S} = \{R[1], S[1]\}$ , there exist numbers  $k \geq 0$  and  $m \geq 0$  such that, for every database  $D$  of  $\mathbf{S}$ ,  $D \models q$  iff  $|S^D| \geq k$  and  $|R^D - S^D| \geq m$ .

Use this fact to derive Theorem 32.7 when  $\diamond$  is  $=$ , i.e., to show that the query  $q_=_$  over  $\mathbf{S}$ , which checks whether, for a database  $D$  of  $\mathbf{S}$ ,  $|R^D| = |S^D|$ , cannot be expressed as a constant-free FO query over  $\mathbf{S}$ .

**Exercise 4.11.** Let  $\varphi_{\geq}()$  be the constant-free FO query over  $\mathbf{S} = \{R[1], S[1]\}$  such that, for every database  $D$  of  $\mathbf{S}$ ,  $D \models \varphi_{\geq}()$  iff  $|R^D| \geq |S^D|$ . Analogously, we define the constant-free FO queries  $\varphi_{>}()$  and  $\varphi_{=}$  over  $\mathbf{S}$ . Show that

$$\lim_{n \rightarrow \infty} \mu_n(\varphi_{\geq}) = \frac{1}{2} \quad \text{and} \quad \lim_{n \rightarrow \infty} \mu_n(\varphi_{>}) = \frac{1}{2}.$$

Also show that

$$\lim_{n \rightarrow \infty} \mu_n(\varphi_{=}) = 0.$$

**Exercise 4.12.** Recall the estimate used in the proof of Theorem 32.7

$$F_n^= = \sum_{k \leq \lfloor n/2 \rfloor} \binom{n}{k} \binom{n-k}{k}.$$

Show that

$$\lim_{n \rightarrow \infty} \frac{F_n^=}{3^n} = 0.$$

**Exercise 4.13.** Recall that Theorem 32.5 (0–1 law) was shown for the special case of constant-free FO sentences over a schema with two unary relation names. Prove the result for the schema  $\mathbf{S} = \{E[2]\}$  with a single binary relation name (i.e., for undirected graphs). Recall that you need to construct a theory  $T$  which has a unique, up to isomorphism, countable model, and whose sentences are true in almost all databases of  $\mathbf{S}$ . Such a theory  $T$  has sentences  $\psi_{k,m}$  that express the following: for every two disjoint sets  $X$  and  $Y$  of nodes of cardinalities  $k$  and  $m$ , respectively, there exists a node  $z$  such that there are edges  $(z, x)$  for each  $x \in X$ , and there is no edge  $(z, y)$  for  $y \in Y$ . While the proof of condition 2 of Lemma 32.6 follows the same ideas as those we saw, the proof of condition 1 is more elaborate and requires ideas not seen in this book; the interested reader is advised to consult [12].

**Exercise 4.14.** Give an example of an FO sentence *with constants* for which Theorem 32.5 (0–1 law) does not hold. A much more difficult task is to show that, for every FO sentence  $\varphi$  with constants,  $\lim_{n \rightarrow \infty} \mu_n(\varphi)$  exists, and is a rational number with the denominator being of the form  $2^k$  for some  $k \geq 0$ .

**Exercise 4.15.** Recall that Theorem 34.1 was shown for  $\text{RA}_{\text{Aggr}}(\Omega)$  without constants. Extend the proof to the version of  $\text{RA}_{\text{Aggr}}(\Omega)$  that allows the use of constants from  $\text{Const}$ . This is done in the following four steps:

1. First extend  $\mathbf{L}_C$  and  $\mathbb{L}_C$  with constants from  $\text{Const}$ , and show that a translation of an  $\text{RA}_{\text{Aggr}}(\Omega)$  expression that uses constants from  $\{c_1, \dots, c_n\} \subsetneq \text{Const}$  into  $\mathbf{L}_C$ , and then into  $\mathbb{L}_C$ , can be carried out in such a way that only constants from  $\{c_1, \dots, c_n\}$  are present in formulae.
2. Next, prove a locality result for the extension of  $\mathbb{L}_C$  with constants. Consider an expression  $\varphi(\bar{x})$ , where  $\varphi$  is an  $\mathbb{L}_C$  formula using constants from  $\{c_1, \dots, c_n\}$  over a schema  $\mathbf{S} = \{R_1 : \tau_1, \dots, R_n : \tau_n\}$  with  $\tau_i \in \{\mathbf{o}\}^{k_i}$  and  $k_i \geq 0$ , for each  $i \in [n]$ , such that  $\text{FV}(\varphi) \subseteq \text{Var}_{\mathbf{o}}$ , and  $\bar{x}$  is a tuple over  $\text{FV}(\varphi)$  that mentions all the variables of  $\text{FV}(\varphi)$ . There exists  $r \geq 0$  such that, for a database  $D$  of  $\mathbf{S}$ , and tuples  $\bar{a}, \bar{b}$  over  $\text{Dom}(D)$ , if  $N_r^D(\bar{a}, \bar{c})$  is isomorphic to  $N_r^D(\bar{b}, \bar{c})$ , for  $\bar{c} = (c_1, \dots, c_n)$ , then either both  $\bar{a}, \bar{b}$  belong to  $\varphi(\bar{x})(D)$ , or none of them belongs to  $\varphi(\bar{x})(D)$ .
3. Finally, using the above locality result for the extension of  $\mathbb{L}_C$  with constants, show that a modified version of the reachability query (that incorporates constants from  $\{c_1, \dots, c_n\}$ ) still violates locality.

**Exercise 4.16.** Use the translation of Theorem 34.6 to find syntactic restrictions on  $\text{FO}_{\text{Aggr}}(\Omega)$  that lead to a logical formalism that can be used to define a query language, i.e., to ensure that the output of an expression  $\varphi(\bar{u})$ , where  $\varphi$  is a formula from the restricted formalism, and  $\bar{u}$  a tuple of variables over  $\text{FV}(\varphi)$  that mentions all the variables of  $\text{FV}(\varphi)$ , on a database is always finite.

**Exercise 4.17.** Recall that Proposition 34.11 was shown for the schema  $\mathbf{S} = \{R : (\mathbf{o}, \mathbf{o})\}$ . Generalize the proof to arbitrary schemas  $\{R_1 : \tau_1, \dots, R_n : \tau_n\}$  with  $\tau_i \in \{\mathbf{o}\}^{k_i}$  and  $k_i \geq 0$ , for each  $i \in [n]$ .

**Exercise 4.18.** There is a different notion of locality, known in the literature as *Hanf-locality* (as opposed to *Gaifman-locality* used in Chapter 34). Given two databases  $D$  and  $D'$  of a two-sorted schema  $\mathbf{S}$ , and tuples  $\bar{a}$  and  $\bar{a}'$  of the same arity over  $\text{Dom}(D)$  and  $\text{Dom}(D')$ , respectively, we write  $(D, \bar{a}) \sim_r (D', \bar{a}')$  if there is a bijection  $f : \text{Dom}(D) \rightarrow \text{Dom}(D')$  such that, for every  $b \in \text{Dom}(D)$ ,  $N_r^D(\bar{a}, b)$  is isomorphic to  $N_r^{D'}(\bar{a}', f(b))$ . Note that, in particular,  $(D, \bar{a}) \sim_r (D', \bar{a}')$  implies  $|\text{Dom}(D)| = |\text{Dom}(D')|$ . A query  $q$  of type  $\tau \in \{\mathbf{o}, \mathbf{n}\}^k$ , for  $k \geq 0$ , over  $\mathbf{S}$  is *Hanf-local* if there exists  $r \geq 0$  such that, for every  $\bar{a}, \bar{a}' \in (\text{Const} \cup \text{Num})^k$ ,  $(D, \bar{a}) \sim_r (D', \bar{a}')$  implies  $\bar{a} \in q(D)$  iff  $\bar{a}' \in q(D')$ .

Prove Theorem 34.3 for Hanf-locality instead of Gaifman locality. To do so, show, by extending the argument in the proof of Proposition 34.11, the following: for an  $\mathbb{L}_C$  formula  $\varphi$  over a schema  $\mathbf{S} = \{R_1 : \tau_1, \dots, R_n : \tau_n\}$  with  $\tau_i \in \{\mathbf{o}\}^{k_i}$  and  $k_i \geq 0$ , for each  $i \in [n]$ , such that  $\text{FV}(\varphi) \subseteq \text{Var}_{\mathbf{o}}$ , and a tuple  $\bar{x}$  over  $\text{FV}(\varphi)$  that mentions all the variables of  $\text{FV}(\varphi)$ ,  $\varphi(\bar{x})$  is Hanf-local.

**Exercise 4.19.** Consider a restriction of  $\mathbf{L}_C$  where infinitary connectives are not allowed, we have quantification over variables of both types, and we have counting terms  $\sharp x \varphi(x, \bar{y})$  counting the number of elements of the input database satisfying  $\varphi$ . Concerning the semantics, is defined in the expected

way with the crucial restriction that numerical variables can only range over the values  $[0, n - 1]$ , where  $n$  is the number of elements in the input database.

Prove that, for every expression  $\varphi()$ , where  $\varphi$  is a sentence from this restricted logic, there exists a Boolean  $\text{RA}_{\text{Aggr}}(\Omega)$  query  $e$ , where  $\Omega$  contains  $<$ ,  $+$ , and  $\cdot$ , and the summation aggregate  $\sum$ , such that, for every database  $D$  of a schema  $\mathbf{S}$ ,  $\varphi()(D) = e(D)$  for the following two cases:

1.  $\mathbf{S} = \{R_1 : \tau_1, \dots, R_n : \tau_n\}$  such that, for each  $i \in [n]$ ,  $\tau_i \in \{\mathbf{n}\}^{k_i}$  for  $k_i \geq 0$ , i.e., we focus on databases where all elements are of type  $\mathbf{n}$ .
2.  $\mathbf{S} = \{R_1 : \tau_1, \dots, R_n : \tau_n\}$  such that, for each  $i \in [n]$ ,  $\tau_i \in \{\mathbf{o}\}^{k_i}$  for  $k_i \geq 0$ , i.e., we focus on databases where all elements are of type  $\mathbf{o}$ , and we have access to an order relation  $<$  over the constants of  $\text{Const}$ .

The importance of these results stems from the fact the restricted logic we defined captures a uniform version of a complexity class called  $\text{TC}^0$  (this stands for threshold circuits of constant depth).  $\text{TC}^0$  has not been separated from others above it such as  $\text{PTIME}$  or  $\text{NLOGSPACE}$ . In particular, this means that bounds on the expressivity of  $\text{RA}_{\text{Aggr}}(\Omega)$  cannot be proved either over ordered non-numerical domains, or numerical domains, without resolving deep problems in complexity theory.

**Exercise 4.20.** A *path system* is a tuple  $P = (V, R, S, T)$ , where  $V$  is a finite set of nodes,  $R \subseteq V \times V \times V$ ,  $S \subseteq V$  and  $T \subseteq V$ . A node  $v \in V$  is *admissible* if  $v \in T$ , or there are admissible nodes  $u, w \in V$  such that  $(v, u, w) \in R$ .

Show that the set of admissible nodes for  $P$  can be computed via a Datalog query. In other words, there exists a Datalog query  $q = (\Pi, A)$ , where  $A$  is unary relation, such that, for every path system  $P$ ,  $q(D_P)$  is the set of admissible nodes for  $P$ , where  $D_P$  stores  $P$  in the obvious way, i.e.,  $V, S$  and  $T$  via unary relations, and  $R$  via a ternary relation.

**Exercise 4.21.** An undirected graph  $G = (V, E)$  is *2-colorable* if there exists a function  $f : V \rightarrow \{0, 1\}$  such that  $(v, u) \in E$  implies  $f(v) \neq f(u)$ .

Show that *non-2-colorability* is expressible via a Datalog query, i.e., there exists a Datalog query  $q = (\Pi, \text{Yes})$ , where  $\text{Yes}$  is a 0-ary relation, such that, for every undirected graph  $G$ ,  $q(D_G) = \text{true}$  if and only if  $G$  is *not* 2-colorable, where  $D_G$  stores  $G$  via the binary relation  $\text{Edge}(\cdot, \cdot)$ .

**Exercise 4.22.** Prove Theorem [35.8](#)

**Exercise 4.23.** Prove Lemma [35.10](#)

**Exercise 4.24.** Show that the query that asks whether an undirected graph is 2-colorable is not expressible via a Datalog query. To do so, exploit the fact that Datalog queries are monotone.

**Exercise 4.25.** Prove that there exists a monotone query that is not expressible via a Datalog query. The proof should not rely on any complexity-theoretic assumption. (It is very easy to show that this holds under the assumption that  $\text{PTIME} \neq \text{NP}$ .)

**Exercise 4.26.** A Datalog program  $\Pi$  is called *linear* if, for every rule  $\rho \in \Pi$ , the body of  $\rho$  mentions at most one relation from  $\text{idb}(\Pi)$ . A Datalog query  $(\Pi, R)$  is linear if  $\Pi$  is linear.

Show that there is no linear Datalog query that computes the set of admissible nodes for a path system  $P$ . The proof should not rely on any complexity-theoretic assumption. (It is easier to show this statement if we assume that  $\text{NLOGSPACE} \neq \text{PTIME}$ .)

**Exercise 4.27.** The *predicate graph* of a program  $\Pi$  is the directed graph  $G_\Pi = (V, E)$ , where  $V$  consists of the relations of  $\text{sch}(\Pi)$ , and  $(P, R) \in E$  if and only if there exists a rule  $\rho \in \Pi$  of the form

$$R(\bar{x}) :- \dots, P(\bar{y}), \dots$$

We call  $\Pi$  *non-recursive* if  $G_\Pi$  is acyclic. A Datalog query  $(\Pi, R)$  is non-recursive if  $\Pi$  is non-recursive.

Show that, for every non-recursive Datalog query  $q = (\Pi, R)$ , there exists a finite UCQ  $q'$  such that, for every database  $D$  of  $\text{edb}(\Pi)$ ,  $q(D) = q'(D)$ .

**Exercise 4.28.** Let  $D$  and  $q$  be the database and the Datalog query, respectively, constructed from the Turing machine  $M$  and the input word  $w$  in the proof of Theorem 37.1. Show that  $M$  accepts  $w$  if and only if  $q(D) = \text{true}$ .

**Exercise 4.29.** A Datalog program  $\Pi$  is called *guarded* if, for every rule  $\rho \in \Pi$ , the body of  $\rho$  has an atom that contains (or “guards”) all the variables occurring in  $\rho$ . A Datalog query  $(\Pi, R)$  is guarded if  $\Pi$  is guarded.

Prove that the EXPTIME-hardness shown in Theorem 37.1 holds even if we focus on guarded Datalog queries.

**Exercise 4.30.** Consider a Datalog query  $q = (\Pi, R)$ , where the arity of the relations of  $\text{sch}(\Pi)$  is bounded by some integer constant, a database  $D$  of  $\text{edb}(\Pi)$ , and a tuple  $\bar{a}$  of arity  $\text{ar}(R)$  over  $\text{Dom}(D)$ . Show that the problem of deciding whether  $\bar{a} \in q(D)$  is NP-complete.

**Exercise 4.31.** Consider a Datalog query  $q = (\Pi, R)$ , where  $\text{sch}(\Pi)$  consists only of 0-ary relations, and a database  $D$  of  $\text{edb}(\Pi)$ . Show that the problem of deciding whether  $q(D) = \text{true}$  is PTIME-complete.

**Exercise 4.32.** Show that the evaluation problem for linear Datalog queries is PSPACE-complete in combined complexity, and NLOGSPACE-complete in data complexity.

**Exercise 4.33.** Show that the evaluation problem for non-recursive Datalog queries is PSPACE-complete in combined complexity, and in DLOGSPACE in data complexity.



---

## Bibliographic Comments

To be done.



## Part V

---

### Uncertainty





## Incomplete Databases

---

So far we assumed that all information in databases was complete: we have a set **Const** from which all entries in databases come, and for every fact  $R(\bar{a})$  in a database, every element of the tuple  $\bar{a}$  comes from **Const**. The reality of course is different: information is often missing, and databases are full of entries shown as null. These null elements can have different semantics: sometimes they mean that an attribute is not known (e.g., date of birth), sometimes they mean that an attribute is not applicable (e.g., date of death for a living person), and sometimes we simply do not have any information why no data value was provided and a null replaced it.

In the next two chapters we study the most commonly considered model of nulls representing missing, and currently unknown data. We define this model using the framework of sets of atoms, and explain what such incomplete instances mean. In fact, they can represent multiple complete databases, and thus query answering must be consistent with answers obtained in all of them.

In this chapter, we present the model of incomplete data and the notion of *certain* query answers. We then analyze its complexity (which is very high), look at what real-world databases do (obviously they keep the complexity low), thus arriving at the inevitable mismatch between the correct way of answering queries and what happens in reality. In the following chapter, we look at the ways of bridging this gap.

### A Formal Model of Incomplete Data and Nulls

The idea of the model is simple: some entries in databases are replaced by variables, which are interpreted as values that exist but are currently unknown. To separate them from variables used in queries, we assume that the set **Var** of variables has a countably infinite subset **Null** of nulls whose elements will be denoted by  $\perp_1, \perp_2, \perp_3, \dots$ . Unlike other elements of **Var**, however, they will not be explicitly referenced in queries. An example of such a database is shown in Figure [39.1](#). Some entries in the **Profession** table contain nulls,

Person			Profession	
pid	pname	cid	pid	prname
1	Aretha	MPH	1	singer
2	Billie	BLT	1	$\perp_1$
3	Bob	DLT	1	$\perp_2$
			2	singer
			3	$\perp_1$
			3	author

Fig. 39.1: A incomplete database of the schema in Example 3.2.

denoted by  $\perp_1$  and  $\perp_2$ . They stand for two currently unknown professions of the person with id 1, namely **Aretha**. Moreover, one of the nulls,  $\perp_1$ , is also a profession of the person with id 3, i.e., **Bob**. Thus, while we do not know what these professions are, we do know that Aretha and Bob share a profession.

### Definition 39.1: Incomplete Database Instances

An *incomplete database instance*  $\mathcal{I}$  of a schema  $\mathbf{S}$  is a finite set of atoms over  $\mathbf{S}$ , where all variables come from the set **Null**.

Since an incomplete instance  $\mathcal{I}$  may contain both constants and nulls, so can its domain  $\text{Dom}(\mathcal{I})$ . We thus use the notation  $\text{Dom}_{\text{Const}}(\mathcal{I})$  for  $\text{Dom}(\mathcal{I}) \cap \text{Const}$  and  $\text{Dom}_{\text{Null}}(\mathcal{I})$  for  $\text{Dom}(\mathcal{I}) \cap \text{Null}$ .

An incomplete instance  $\mathcal{I}$  may represent multiple complete instances. To start with, nulls in it may be replaced by elements of **Const**, i.e., by their real, but currently unknown values. If this is the only way to obtain further information about  $\mathcal{I}$ , one refers to the semantics under the *Closed World Assumption*, or CWA-semantics. Another, more permissive interpretation, permits not only replacing nulls with constant values but also adding tuples that may have been missing in  $\mathcal{I}$ . This is known as the semantics under the *Open World Assumption*, or OWA.

To define them formally, we use the machinery of homomorphisms. Recall that a homomorphism  $h$  on a set of atoms (i.e., an incomplete instance)  $\mathcal{I}$  is a mapping  $h$  defined on  $\text{Dom}(\mathcal{I})$  such that  $h(a) = a$  whenever  $a \in \text{Dom}_{\text{Const}}(\mathcal{I})$ , and otherwise  $h(a)$  is in  $\text{Const} \cup \text{Var}$ . In providing semantics of incompleteness, we want to replace nulls with constant values rather than other nulls. Thus, we look at homomorphisms whose image is a subset of **Const**. These are called *valuations*.

### Definition 39.2: Closed and Open-World Semantics

Given an incomplete instance  $\mathcal{I}$  of schema  $\mathbf{S}$ , its CWA-semantics is defined as

$$\llbracket \mathcal{I} \rrbracket_{\text{CWA}} = \{D \in \text{Inst}(\mathbf{S}) \mid D = h(\mathcal{I}) \text{ and } h \text{ is a valuation}\}.$$

Its OWA-semantics is defined as

$$\llbracket D \rrbracket_{\text{OWA}} = \{D \in \text{Inst}(\mathbf{S}) \mid \text{there is a valuation } h : \mathcal{I} \rightarrow D\}.$$

It is easy to see that  $D \in \llbracket \mathcal{I} \rrbracket_{\text{OWA}}$  iff  $D \supseteq D'$  for some  $D' \in \llbracket \mathcal{I} \rrbracket_{\text{CWA}}$ .

The database in Figure 3.1 of Example 3.2 belongs to  $\llbracket \mathcal{I} \rrbracket_{\text{OWA}}$  for the instance  $\mathcal{I}$  in Figure 39.1. It does not belong to  $\llbracket \mathcal{I} \rrbracket_{\text{CWA}}$  due to the presence in it of relation **City**, and information about the person with **pid** 4, namely Freddie. However, if all the tuples about that person are eliminated from the instance in Figure 3.1 together with the relation **City**, then the resulting database will be in  $\llbracket \mathcal{I} \rrbracket_{\text{CWA}}$ .

## Certain Answers

Given an incomplete instance  $\mathcal{I}$  and a query  $q$  written in some query language, what is  $q(\mathcal{I})$ ? Since the instance  $\mathcal{I}$  represents multiple complete instances (those in  $\llbracket \mathcal{I} \rrbracket_{\text{CWA}}$  or  $\llbracket \mathcal{I} \rrbracket_{\text{OWA}}$ , depending on the semantics),  $q(\mathcal{I})$  must take into account query answers  $q(D)$  for all  $D$  that  $\mathcal{I}$  can represent.

To arrive at the definition of answers to  $q$ , consider initially Boolean queries. That is,  $q(D)$  is either **true** or **false**. When can we say that the answer to  $q$  in  $\mathcal{I}$  is **true**? This can be asserted with certainty only if  $q(D) = \text{true}$  for all  $D \in \llbracket \mathcal{I} \rrbracket$  (where  $\llbracket \cdot \rrbracket$  is one of the semantics we use). In that case, we say that the *certain answer* to  $q$  over  $\mathcal{I}$  is **true**; otherwise it is **false**. From now on, we shall use the notation  $\text{cert}_*(q, \mathcal{I})$  for certain answers, where  $*$  is OWA or CWA. Recall that **true** is viewed as the set  $\{()\}$  containing the empty tuple while **false** is viewed as the empty set. Thus,  $\text{cert}_*(q, \mathcal{I}) = \bigcap \{q(D) \mid D \in \llbracket \mathcal{I} \rrbracket_*\}$ .

Next, we move to queries of arbitrary arity. Suppose we have an incomplete instance  $\mathcal{I}$  and a tuple  $\bar{a}$  that consists of elements of  $\text{Dom}_{\text{Const}}(\mathcal{I})$ . When can we state with certainty that this tuple is in the answer to  $q$  on  $\mathcal{I}$ ? Again, this happens if it is in all the answers  $q(D)$ , for  $D \in \llbracket \mathcal{I} \rrbracket_*$ . So again we can define  $\text{cert}_*(q, \mathcal{I})$  as  $\bigcap \{q(D) \mid D \in \llbracket \mathcal{I} \rrbracket_*\}$ .

The definition gets more complicated if a tuple  $\bar{a}$  contains both nulls and constants. Then the previous definition does not work: if  $\bar{a}$  contains nulls, it cannot belong to  $q(D)$  for any complete  $D$ , as  $D$  itself contains no nulls. But how can a query return, with certainty, a tuple with nulls? To see this, consider a simple query that returns the relation **Profession** in Figure 39.1. We know with certainty that the tuple  $\bar{a} = (1, \perp_1)$  is in it, and thus we expect it to be in the answer. The reason we know this with certainty is that for every valuation  $h$ , we have  $h(\bar{a}) \in h(\text{Profession})$ . We now extend this to arbitrary queries.

**Definition 39.3: Certain Answers**

Given a query  $q$ , and an incomplete instance  $\mathcal{I}$ , its certain answers under CWA and under OWA are defined by:

$$\bar{a} \in \text{cert}_{\text{CWA}}(q, \mathcal{I}) \Leftrightarrow h(\bar{a}) \in q(h(\mathcal{I})) \text{ for every valuation } h$$

$$\bar{a} \in \text{cert}_{\text{OWA}}(q, \mathcal{I}) \Leftrightarrow h(\bar{a}) \in q(D) \text{ for every valuation } h : \mathcal{I} \rightarrow D$$

If  $\mathcal{I}$  has no nulls, then  $\llbracket \mathcal{I} \rrbracket_{\text{CWA}} = \{\mathcal{I}\}$  and thus in this case  $\text{cert}_{\text{CWA}}(q, \mathcal{I}) = q(\mathcal{I})$ . Hence, certain query answering under CWA properly generalizes query answering in complete databases.

For tuples with no nulls in them, this definition of certain answers coincides with the intersection of query answers in all databases represented by  $\mathcal{I}$ .

**Proposition 39.4**

Given a query  $q$ , an incomplete instance  $\mathcal{I}$ , and a tuple  $\bar{a}$  without nulls,  $\bar{a} \in \text{cert}_*(q, \mathcal{I})$  iff  $\bar{a} \in \bigcap \{q(D) \mid D \in \llbracket \mathcal{I} \rrbracket_*\}$ , where  $*$  is OWA or CWA.

*Proof.* We show this for CWA; for OWA the proof is very similar. Assume  $\bar{a} \in \text{cert}_{\text{CWA}}(q, \mathcal{I})$  and take an arbitrary valuation  $h$ . Then  $h(\bar{a}) = \bar{a}$  since  $\bar{a}$  has no nulls and thus  $\bar{a} \in q(h(\mathcal{I}))$ , proving  $\bar{a} \in \bigcap \{q(D) \mid D \in \llbracket \mathcal{I} \rrbracket_{\text{CWA}}\}$ . Conversely, assume  $\bar{a} \in \bigcap \{q(D) \mid D \in \llbracket \mathcal{I} \rrbracket_{\text{CWA}}\}$ . Take any valuation  $h$ ; then  $h(\mathcal{I}) \in \llbracket \mathcal{I} \rrbracket_{\text{CWA}}$  and hence  $\bar{a} \in q(h(\mathcal{I}))$ . Since  $\bar{a}$  has no nulls, we have  $h(\bar{a}) = \bar{a} \in q(h(\mathcal{I}))$ , thus proving  $\bar{a} \in \text{cert}_{\text{CWA}}(q, \mathcal{I})$ .

**Example 39.5: Certain Answers to Queries**

We consider two queries similar to those in Example 3.2 now asked on the incomplete instance  $\mathcal{I}$  shown in Figure 39.1. As the language we use FO, with one clarification: as mentioned earlier, variables in its formulae use come from **Var** – **Null**. That is, elements of **Null** cannot be used as explicit constants in formulae.

The first query asks for names of persons with two different professions:

$$\begin{aligned} \exists x \exists z \exists u_1 \exists u_2 \big( \text{Person}(x, y, z) \wedge \\ \text{Profession}(x, u_1) \wedge \text{Profession}(x, u_2) \wedge \neg(u_1 = u_2) \big). \end{aligned} \quad (39.1)$$

The certain answer to this query is empty, under both OWA and CWA interpretation. Indeed, Billie has only one profession, and even though it may seem like Aretha and Bob have multiple professions, there are possible completions of  $\mathcal{I}$  where this is not the case. For example, a valuation where  $h(\perp_1) = h(\perp_2) = \text{'singer'}$  results in a complete database

where Aretha has only one profession, while a valuation with  $h(\perp_1) =$  author results in a complete database where the same is true of Bob.

The second asks for pairs of persons who share a profession:

$$\exists x \exists x' \exists z \exists z' \exists u (\text{Person}(x, y, z) \wedge \text{Person}(x', y', z') \wedge (\text{Profession}(x, u) \wedge \text{Profession}(x', u))).$$

This query will return pairs of names, including all pairs of the form  $(n, n)$ , but also (Aretha, Billie) and (Billie, Aretha) who share the profession ‘singer’, as well as (Aretha, Bob) and (Bob, Aretha), who too share a profession,  $\perp_1$ , though we do not know what it is. Still, in every database  $D$  that is equal or contains  $h(\mathcal{I})$ , this common profession will be represented by  $h(\perp_1)$ , and thus these pairs will be output by the query, making them certain answers.

It might seem from these examples that certain answers are the same under OWA and CWA. This is not the case however, and we shall see soon that they differ dramatically in terms of their complexity. To give a simple example where certain answers are different, consider a Boolean query  $q$  given by

$$\forall x \forall y \forall z (\text{Person}(x, y, z) \rightarrow \exists u \text{Profession}(x, u))$$

asking whether every person has a recorded profession. Then  $\text{cert}_{\text{CWA}}(q, \mathcal{I}) = \text{true}$  while  $\text{cert}_{\text{OWA}}(q, \mathcal{I}) = \text{false}$ . For example, a database  $D \in \llbracket \mathcal{I} \rrbracket_{\text{OWA}}$  obtained by giving values to nulls and by adding a tuple to table **Person** with an id that is not present in table **Profession** violates the condition expressed by  $q$ .

## Complexity of Certain Answers

Next, we analyze the complexity of certain answers, and see that it is significantly higher than the complexity of query evaluation.

### Theorem 39.6

- Under the CWA semantics, data complexity of certain answers to FO queries is CONP-complete and their combined complexity is PSPACE-complete.
- Under the OWA semantics, combined complexity of certain answers to FO queries is undecidable.

*Proof. Upper bounds.* We look at the complement problem, namely whether  $\bar{a} \notin \text{cert}_{\text{CWA}}(q, \mathcal{I})$  for an incomplete instance  $\mathcal{I}$ , a query  $q$ , and a tuple  $\bar{a}$ . We will show that this problem is in PSPACE if all of  $\bar{a}$ ,  $\mathcal{I}$ , and  $q$  are inputs, and in NP,

if  $q$  is fixed. This will suffice, since PSPACE is closed under complementation. Note that  $\bar{a} \notin \text{cert}_{\text{CWA}}(q, \mathcal{I})$  iff  $h(\bar{a}) \notin q(h(\mathcal{I}))$  for some valuation  $h$ . This looks like the proof: just guess that valuation  $h$ . However, to prove membership in NP, a guess must be of polynomial size, but the above  $h$  could be an arbitrary valuation, i.e., its range could be an arbitrary subset of the infinite set  $\text{Const}$ . Instead, we show that if  $\bar{a} \notin \text{cert}_{\text{CWA}}(q, \mathcal{I})$ , then such  $h$  can be chosen among a finite set of valuations with a fixed range that is linear in the size of  $\text{Dom}(\mathcal{I})$ . This will prove membership in NP for a fixed  $q$ , as well as membership in PSPACE for combined complexity. Indeed, PSPACE is closed under nondeterminism (hence one can guess  $h$ ) and checking  $h(\bar{a}) \notin q(h(\mathcal{I}))$  is in PSPACE, since  $q$  is an FO query.

Assume that  $\text{Dom}_{\text{Null}}(\mathcal{I}) = \{\perp_1, \dots, \perp_m\}$ . Let  $C = \{c_1, \dots, c_m\}$  be a set of  $m$  elements of  $\text{Const}$  disjoint from  $\text{Dom}(\mathcal{I})$  and  $\text{Dom}(q)$ , the set of constants mentioned in  $q$ . We claim that if  $h(\bar{a}) \notin q(h(\mathcal{I}))$  for some  $h$ , then this also happens for some  $h$  whose range is contained in  $\text{Dom}_{\text{Const}}(\mathcal{I}) \cup \text{Dom}(q) \cup C$ ; clearly this will suffice. Consider  $h$  such that  $h(\bar{a}) \notin q(h(\mathcal{I}))$  and let  $i_1, \dots, i_k$  be such that  $h(\perp_{i_j}) \notin \text{Dom}_{\text{Const}}(\mathcal{I}) \cup \text{Dom}(q)$  for  $j \in [k]$ . We then define  $h'$  that coincides with  $h$  everywhere except that it sets  $h'(\perp_{i_j}) = c_{i_j}$  for  $j \in [k]$ . Let  $\pi : \text{Const} \rightarrow U$  be a bijection that is the identity on  $\text{Dom}_{\text{Const}}(\mathcal{I}) \cup \text{Dom}(q)$ , sends each  $c_{i_j}$  to  $h(\perp_{i_j})$ , and is arbitrarily defined elsewhere. Note then that  $\pi(h'(a)) = h(a)$  for every  $a \in \text{Dom}(\mathcal{I}) \cup \text{Dom}(q)$ . If we now assume that  $h'(\bar{a}) \in q(h'(\mathcal{I}))$ , then by genericity of FO queries (see Exercise 1.2) that would mean that  $\pi(h'(\bar{a})) \in q(\pi(h'(\mathcal{I})))$ , i.e.,  $h(\bar{a}) \in q(h(\mathcal{I}))$ , which is a contradiction. Hence, we can find a witness valuation with the range contained in  $\text{Dom}_{\text{Const}}(\mathcal{I}) \cup \text{Dom}(q) \cup C$ . This finishes the proof of upper bounds.

*Lower bounds.* The PSPACE lower bound for combined complexity is already known for queries on databases with no nulls. Indeed, for such databases  $\text{cert}_{\text{CWA}}(q, \mathcal{I}) = q(\mathcal{I})$ , and the latter is known to be PSPACE-hard.

We next show coNP-hardness of data complexity. Actually, we will look at the complement problem, which for Boolean queries is whether  $\text{cert}_{\text{CWA}}(q, \mathcal{I}) = \text{false}$  and prove it to be NP-hard, for a fixed query  $q$ . For this we reduce from graph 3-colorability, a well-known NP-complete problem. Given a graph  $G = \langle V, E \rangle$ , it is 3-colorable if there a coloring map  $c : V \rightarrow \{r, g, b\}$  such that for every edge  $(v, v') \in E$ , the colors  $c(v)$  and  $c(v')$  are different. Now given such a graph with  $V = \{v_1, \dots, v_n\}$ , we create an incomplete instance  $\mathcal{I}_G$  of the schema containing relations  $E[2]$  and  $C[2]$  such that  $E$  is the edge relation of the graph, and  $C$  contains tuples  $(v_1, \perp_1), \dots, (v_n, \perp_n)$ , where  $\perp_1, \dots, \perp_n$  are distinct nulls. The intuition is that these are yet unknown colors assigned to vertices; a database  $D \in \llbracket \mathcal{I}_G \rrbracket_{\text{CWA}}$  then assigns colors to vertices. We next define a Boolean FO query  $q$  stating that such an assignment is not a 3-coloring. Specifically, this happens if

1. there are 4 or more different colors (values of the second column of  $C$ ); or
2. there is an edge whose endpoints are given the same color.

The first condition is expressed by an FO query

$$\exists x_1 \exists x_2 \exists x_3 \exists x_4 \exists y_1 \exists y_2 \exists y_3 \exists y_4 \left( \bigwedge_{i=1}^4 C(x_i, y_i) \wedge \text{diff}(y_1, y_2, y_3, y_4) \right)$$

where the formula  $\text{diff}(y_1, y_2, y_3, y_4)$  says that all of  $y_1, y_2, y_3, y_4$  are pairwise distinct; the second condition is expressed by

$$\exists x \exists y \exists u \left( C(x, u) \wedge C(y, u) \wedge \neg(x = y) \right).$$

Thus, their disjunction is an FO query  $q$  such that  $q$  is true in  $h(\mathcal{I}_G)$  iff  $h$ , that assigns values to nulls, is not a 3-coloring of  $G$ . Hence  $\text{cert}_{\text{CWA}}(q, \mathcal{I}_G) = \text{false}$  iff there is at least one  $h$  that is a 3-coloring, i.e., iff the graph is 3-colorable. Hence, checking  $\text{cert}_{\text{CWA}}(q, \mathcal{I}_G) = \text{false}$  is NP-hard in data complexity, and thus checking  $\text{cert}_{\text{CWA}}(q, \mathcal{I}_G) = \text{true}$  is coNP-hard.

We finally show undecidability of the combined complexity of certain answers under OWA. Consider a Boolean FO query  $q$  given by a sentence  $\varphi$ , and let  $\mathcal{I}_\emptyset$  be the empty instance, where every relation has no tuples. Then an arbitrary database  $D$  is in  $\llbracket \mathcal{I}_\emptyset \rrbracket_{\text{OWA}}$ . Hence  $\text{cert}_{\text{OWA}}(q, \mathcal{I}_\emptyset) = \text{true}$  iff  $\varphi$  is a valid sentence, i.e., true in every database. But we know from Chapter 8 that this is an undecidable condition.  $\square$

Note that Theorem 39.6 says nothing about data complexity of finding certain answers under OWA. Indeed, to prove undecidability, we use the fixed empty database and it was the query that varied. This alone rules out certain answers under OWA as a viable query answering strategy, but we can show an even stronger undecidability of data complexity. That is, there is a fixed FO query  $q$  such that  $\text{cert}_{\text{OWA}}(q, \mathcal{I})$  is undecidable when  $\mathcal{I}$  is the only input; see Exercise 5.1.

Note also that in the coNP-hardness reduction in the proof, we used an example of an incomplete instance in which every nulls occurs exactly once. Furthermore, the fixed query used in the reduction is a union of CQs with inequalities, meaning that already for them the problem of finding certain answers is intractable.

## SQL Nulls and Three-Valued Logic

If finding certain answers is intractable even under the CWA, what do real-life databases do? Clearly they cannot implement intractable exponential-time algorithms. They do not indeed, and instead they use a very different solution: SQL resorts to a *3-valued logic* with an extra truth value **unknown**, which is reserved for comparisons involving nulls. In fact in relational DBMSs there is just one single null value, i.e., we cannot indicate that two values currently unknown are equal. But this restriction, as we just saw, does not lower the complexity.



The 3-valued logic, used by SQL, and also known as Kleene's logic, is as follows:

- On **true** and **false**, operations  $\wedge$ ,  $\vee$ , and  $\neg$  have their standard Boolean interpretation;
- $\neg \text{unknown} = \text{unknown}$  (if we do not know a truth value, we do not know its negation);
- $\text{true} \vee \text{unknown} = \text{true}$  (disjunction of **true** with anything is **true**) while  $\text{false} \vee \text{unknown} = \text{unknown} \vee \text{unknown} = \text{unknown}$  (if one truth value is unknown its disjunction with anything other than **true** is unknown);
- $\text{false} \wedge \text{unknown} = \text{false}$  (conjunction of **false** with anything is **false**) while  $\text{true} \wedge \text{unknown} = \text{unknown} \wedge \text{unknown} = \text{unknown}$  (if one truth value is unknown its conjunction with anything other than **false** is unknown).

For the definition of satisfaction, it is no longer sufficient to use only the notion of  $(\mathcal{I}, \eta) \models \varphi$ , which says that  $\varphi$  evaluates to **true** under the assignment  $\eta$  in  $\mathcal{I}$ . In Boolean logic, the alternative to it is that  $\varphi$  evaluates to **false**, but now we need to account for the possibility of  $\varphi$  evaluating to **unknown**. The key rule that SQL uses is the following: *If at least one argument of a comparison is null, then the result of the comparison is unknown.*

**Definition 39.7: Eval( $\varphi, \mathcal{I}, \eta$ )**

$\text{Eval}(\varphi, \mathcal{I}, \eta)$ , the truth value to which  $\varphi$  evaluates in an incomplete instance  $\mathcal{I}$  under assignment  $\eta$ , is defined as follows.

- $\text{Eval}(x = y, \mathcal{I}, \eta)$  is

$$\begin{cases} \text{true} & \text{if } \eta(x), \eta(y) \in \text{Dom}_{\text{Const}}(\mathcal{I}) \text{ and } \eta(x) = \eta(y) \\ \text{false} & \text{if } \eta(x), \eta(y) \in \text{Dom}_{\text{Const}}(\mathcal{I}) \text{ and } \eta(x) \neq \eta(y) \\ \text{unknown} & \text{if } \eta(x) \in \text{Dom}_{\text{Null}}(\mathcal{I}) \text{ or } \eta(y) \in \text{Dom}_{\text{Null}}(\mathcal{I}) \end{cases}$$

- If  $a$  is a constant, then  $\text{Eval}(x = a, \mathcal{I}, \eta)$  is

$$\begin{cases} \text{true} & \text{if } \eta(x) = a \\ \text{false} & \text{if } \eta(x) \in \text{Dom}_{\text{Const}}(\mathcal{I}) \text{ and } \eta(x) \neq a \\ \text{unknown} & \text{if } \eta(x) \in \text{Dom}_{\text{Null}}(\mathcal{I}) \end{cases}$$

- $\text{Eval}(R(u_1, \dots, u_k), \mathcal{I}, \eta)$  is

$$\begin{cases} \text{true} & \text{if } R(\eta(u_1), \dots, \eta(u_k)) \in \mathcal{I} \\ \text{false} & \text{otherwise} \end{cases}$$

- $\text{Eval}(\varphi_1 \wedge \varphi_2, \mathcal{I}, \eta) = \text{Eval}(\varphi_1, \mathcal{I}, \eta) \wedge \text{Eval}(\varphi_2, \mathcal{I}, \eta)$
- $\text{Eval}(\varphi_1 \vee \varphi_2, \mathcal{I}, \eta) = \text{Eval}(\varphi_1, \mathcal{I}, \eta) \vee \text{Eval}(\varphi_2, \mathcal{I}, \eta)$
- $\text{Eval}(\neg \varphi_1, \mathcal{I}, \eta) = \neg \text{Eval}(\varphi_1, \mathcal{I}, \eta)$
- $\text{Eval}(\exists x \psi, \mathcal{I}, \eta) = \bigvee_{a \in \text{Dom}(\mathcal{I})} \text{Eval}(\psi, \mathcal{I}, \eta[x/a])$ .
- $\text{Eval}(\exists x \psi, \mathcal{I}, \eta) = \bigvee_{a \in \text{Dom}(\mathcal{I})} \text{Eval}(\psi, \mathcal{I}, \eta[x/a])$ .

A formula then defines a query  $q = (\varphi, \bar{x})$  whose result on an incomplete instance  $\mathcal{I}$  under the 3-valued logic interpretation is defined as

$$q_{3\text{vl}}(\mathcal{I}) = \{\eta(\bar{x}) \mid \text{Eval}(\varphi, \mathcal{I}, \eta) = \text{true}\}.$$

Thus, in computing query answers, we are only interested in tuples that are **true**, dismissing not only **false** but also **unknown**; in other words, **unknown** is a truth value used for evaluating conditions but not for producing query output. Notice that a query evaluated on an incomplete instance may return tuples with nulls.

To give an example, consider query (39.1) from Example 39.5, evaluated on the incomplete instance  $\mathcal{I}$  in Figure 39.1. It asks for names of people having two different professions. This query will return nothing under the 3-valued interpretation as all comparisons  $\neg(u_1 = u_2)$  in it will evaluate to either **false**, when  $u_1$  and  $u_2$  are the same, or **unknown**, when one of them is a null. Thus, when the free variable  $y$  is instantiated as Aretha or Bob (pid is 1 or 3), the formula evaluates to **unknown**, and when it is instantiated as Billie, with pid 2 and a single profession recorded, the formula is **false**.

By a simple extension of the query evaluation algorithm for FO queries, we can show

#### Proposition 39.8

Query evaluation of FO queries under the 3-valued semantics is in DLOGSPACE.

Thus, the 3-valued semantics cannot compute certain answers. The question is then how close its answers are to certain answers. In general, they can deviate from each other in two possible ways:

- SQL evaluation can produce false negatives, i.e., miss some certain answers: that is,  $q_{3\text{vl}}(\mathcal{I}) \subsetneq \text{cert}_{\text{CWA}}(q, \mathcal{I})$ ;
- SQL evaluation can produce false positives, i.e., give answers which are now certain:  $\text{cert}_{\text{CWA}}(q, \mathcal{I}) \subsetneq q_{3\text{vl}}(\mathcal{I})$ .

We now show that the evaluation strategy based on 3-valued logic that is chosen by SQL can produce *both*. In the next chapter we study alternative query evaluation strategies that provide some guarantees on the quality of the answers.

**Proposition 39.9**

Query evaluation of FO queries under the 3-valued semantics can produce both false positives and false negatives.

*Proof.* Let  $\mathcal{I}$  contain a binary relation  $R$  with tuples  $(1, \perp)$  and  $(2, \perp)$ . Consider a Boolean query  $q$  given by  $\exists x \exists y (R(x, y) \wedge \neg(x = y))$ . Then  $\text{cert}_{\text{CWA}}(q, \mathcal{I}) = \text{true}$ . Indeed,  $h(\perp)$  cannot be both 1 and 2 and thus in  $h(\mathcal{I})$  at least one tuple will have different values. On the other hand,  $q_{3\text{val}}(\mathcal{I}) = \emptyset$ , since  $\text{Eval}(q, \mathcal{I}) = \text{unknown}$  due to all the conditions  $x = y$  evaluating to **unknown**. Since  $\text{true} = \{()\}$ , we have an example of a false negative.

To give an example of a false positive, consider an incomplete instance  $\mathcal{I}$  with two unary relations  $R$  and  $S$  containing  $(1)$  and  $(\perp)$  respectively. Let  $q$  be given by  $\varphi(x) = R(x) \wedge \neg(R(x) \wedge \neg S(y))$ . Then its certain answer is empty, as witnessed by  $h(\mathcal{I})$  with  $h(\perp) = 2$ . However, a routine inspection shows  $q_{3\text{val}}(\mathcal{I}) = \{(1)\}$ , thus producing a false positive.  $\square$

## Computing Certain Answers

To summarize the previous chapter, we have seen that:

- finding certain answers for FO queries is computationally intractable;
- SQL resorts to a 3-valued logic to achieve low complexity of query evaluation; however, this leads to the presence of both false positives and false negatives in query answers.

In this chapter we look at three different query evaluation strategies over incomplete instances that combine efficiency and correctness.

**Naïve evaluation.** One can try to evaluate a query naïvely, as if nulls were actual values. This cannot always give us certain answers – otherwise they would have had tractable data complexity – but naïve evaluation avoids false negatives and for some queries actually computes certain answers.

**Approximation: absolute guarantees.** We show to evaluate all FO queries tractably, over incomplete instances, without producing false positives. This is achieved by means of query translation: a query  $q$  is translated into a query  $q^t$  whose evaluation is guaranteed to produce a subset of certain answers.

**Approximation: probabilistic guarantees.** A different approach is to compute the probability that a given tuple is in the answer for a randomly picked assignment of values to nulls. It turns out that for all FO queries, their naïve evaluation gives us correct answers with probability 1.

### Naïve Evaluation and Correctness of Answers

We first need to formalize the idea of answering a query on an incomplete instance treating nulls as new constants. For this, we use the standard notion of satisfaction of FO formulae, but extended to incomplete instances. More

precisely, to define  $(\mathcal{I}, \eta) \models \varphi$ , we follow, word by word, Definition 3.3 of satisfaction of FO formulae. Since **Const** and **Null** are disjoint sets, it means that equalities  $\perp = a$  for a constant  $a$  and  $\perp = \perp'$  for two distinct nulls evaluate to **false**, while  $\perp = \perp$  evaluates to **true**.

**Definition 40.1: Naïve Evaluation**

The naïve evaluation of an FO query  $q = (\varphi, \bar{x})$  on an incomplete instance  $\mathcal{I}$  is defined as

$$\text{naïve}(q, \mathcal{I}) = \{\eta(\bar{x}) \mid (\mathcal{I}, \eta) \models \varphi\}$$

where  $\eta$  ranges over all assignments  $\bar{x} \mapsto \text{Dom}(\mathcal{I}) \cup \text{Dom}(q)$ .

In other words, the definition is really unchanged compared to the usual definition of FO query evaluation except that free variables can now be assigned nulls, and not only constants. For example, given an incomplete instance with atoms  $R(1, \perp), R(\perp, \perp), R(\perp, \perp')$ , the naïve evaluation of the query  $R(x, y) \wedge x = y$  produces  $(\perp, \perp)$ , while the naïve evaluation of  $R(x, y) \wedge \neg(x = y)$  produces  $(1, \perp)$  and  $(\perp, \perp')$ .

What is the relationship between naïve evaluation and certain answers? We saw that SQL's 3-valued logic-based evaluation can produce both false positives and false negatives. Naïve evaluation, on the other hand, avoids false negatives.

**Proposition 40.2**

For every FO query  $q$ , we have  $\text{cert}_{\text{OWA}}(q, \mathcal{I}) \subseteq \text{cert}_{\text{CWA}}(q, \mathcal{I}) \subseteq \text{naïve}(q, \mathcal{I})$ .

*Proof.* The first inclusion is an immediate consequence of  $\llbracket \mathcal{I} \rrbracket_{\text{CWA}} \subseteq \llbracket \mathcal{I} \rrbracket_{\text{OWA}}$ . Assume next that  $q$  is given as  $(\varphi, \bar{x})$  and  $\bar{a} \in \text{cert}_{\text{CWA}}(q, \mathcal{I})$ . We say that a valuation  $h$  on  $\text{Dom}_{\text{Null}}(\mathcal{I})$  is *simple* if it is a bijection and furthermore  $h(\perp) \notin \text{Dom}_{\text{Const}}(\mathcal{I}) \cup \text{Dom}(\varphi)$  for every  $\perp \in \text{Dom}_{\text{Null}}(\mathcal{I})$ . If  $\bar{a} \in \text{cert}_{\text{OWA}}(q, \mathcal{I})$  then  $h(\bar{a}) \in q(h(\mathcal{I}))$  for every valuation  $h$ , and in particular for every simple valuation  $h$ . The latter means  $(h(\mathcal{I}), h \circ \eta) \models \varphi$ , where  $\eta$  assigns  $\bar{a}$  to  $\bar{x}$ . Using  $\text{Dom}(\varphi)$ -genericity of  $\varphi$  (see Exercise 1.2) and applying it to the bijection  $h^{-1}$ , we conclude that  $(h^{-1}(h(\mathcal{I})), h^{-1} \circ h \circ \eta) \models \varphi$ , i.e.,  $(\mathcal{I}, \eta) \models \varphi$ , and thus  $\bar{a} \in \text{naïve}(q, \mathcal{I})$ .  $\square$

At the same time, naïve evaluation can produce false positives. Consider an FO query  $\varphi(x) = R(x) \wedge \neg S(x)$  evaluated on a database with facts  $R(1), S(\perp)$ . Its naïve evaluation outputs 1, while certain answers are empty, as witnessed by the valuation  $h(\perp) = 1$ .

Nonetheless, for a large class of queries, naïve evaluation produces exactly the certain answers. The class depends on the semantics of incompleteness, and is larger for CWA. Indeed, under CWA, there are fewer instances in  $\llbracket \mathcal{I} \rrbracket_{\text{CWA}}$  than in  $\llbracket \mathcal{I} \rrbracket_{\text{OWA}}$  and thus it is “easier” for a tuple to be a certain answer.

**Definition 40.3: Positive Formulae with Atomic Guards (PAG)**

The class PAG of *positive atomic formulae with guards* is given by the following rules:

- every atomic formula is in PAG;
- if  $\varphi$  and  $\psi$  are in PAG, then so are  $\varphi \wedge \psi$  and  $\varphi \vee \psi$ ;
- if  $\varphi$  is in PAG, then so are  $\exists x \varphi$  and  $\forall x \varphi$ ;
- if  $\varphi$  is in PAG, and  $\alpha(\bar{x})$  is an atomic formula in which no variable repeats, then  $\forall \bar{x} (\alpha(\bar{x}) \rightarrow \varphi)$  is in PAG.

A PAG query is an FO query given by a PAG formula.

This class of formulae excludes the negation from FO, while keeping everything else, and replaces negation by a weak form that has a guarded implication  $\alpha(\bar{x}) \rightarrow \varphi$  under the scope of a universal quantifier. These are rather common queries of the shape “objects that participate in every relationship of a certain kind”, e.g., customers who buy every product. Using the schema of our example, such a query would be

$$\forall x_1 \forall x_2 (\text{Profession}(x_1, x_2) \rightarrow \exists w \text{Person}(x, y, w) \wedge \text{Profession}(x, x_2))$$

asking asking for ids ( $x$ ) and names ( $y$ ) of people who have every profession listed in the database. These queries correspond to the *division* operation of relational algebra; see Exercise 5.4 for the explicit connection. Notice that the guard  $\alpha$  is not allowed to repeat variables: for example, a guard  $R(x, x)$  is not allowed, while  $R(x, y)$  is.

**Theorem 40.4**

1. If  $q$  is a UCQ, then  $\text{naïve}(q, \mathcal{I}) = \text{cert}_{\text{OWA}}(q, \mathcal{I}) = \text{cert}_{\text{CWA}}(q, \mathcal{I})$ .
2. If  $q$  is a PAG query, then  $\text{naïve}(q, \mathcal{I}) = \text{cert}_{\text{CWA}}(q, \mathcal{I})$ .

The key idea behind the proof is that the capturing of certain answers by naïve evaluation is very closely connected with preservation under homomorphisms, seen in Chapter 13. The intuition behind this connection is as follows:  $\llbracket \mathcal{I} \rrbracket_{\text{OWA}}$  consists of all database instances  $D$  such that we have a homomorphism  $\mathcal{I} \rightarrow D$ . Thus, if we have a Boolean query  $q$  and it is naïve ly true in  $\mathcal{I}$ , and furthermore preserved under all homomorphisms, then  $q$  is true in all  $D \in \llbracket \mathcal{I} \rrbracket_{\text{OWA}}$ . This reasoning extends to queries with free variables and to the CWA semantics. For the latter, we need the notion of *strong onto homomorphisms*. These are homomorphisms  $h : S \rightarrow S'$  between sets of atoms such that  $h(S) = S'$ .

We now recall the definition of preservation under homomorphisms. A query  $q$  is preserved under homomorphisms if  $(D, \bar{a}) \rightarrow_{\text{Dom}(q)} (D', \bar{b})$  and

$\bar{a} \in q(D)$  imply  $\bar{b} \in q(D')$ , where  $(D, \bar{a}) \rightarrow_C (D', \bar{b})$  means that there exists a homomorphism from  $(V_C(D), V_C(\bar{a}))$  to  $(D', \bar{b})$ . Here  $V_C$  turns every constant in  $D$  other than those in  $C$  into a new null, see Chapter 13. If in the above definition we replace homomorphisms by strong onto homomorphisms, then we have the definition of preservation under strong onto homomorphisms.

**Proposition 40.5**

Let  $q$  be an FO query.

1. If  $q$  is preserved under homomorphisms, then  $\text{naïve}(q, \mathcal{I}) = \text{cert}_{\text{OWA}}(q, \mathcal{I})$ .
2. If  $q$  is preserved under strong onto homomorphisms, then  $\text{naïve}(q, \mathcal{I}) = \text{cert}_{\text{CWA}}(q, \mathcal{I})$ .

*Proof.* We prove the first item for the OWA semantics; the proof for CWA is very similar. Assume  $q = (\varphi, \bar{x})$  is preserved under homomorphisms. We already know that  $\text{cert}_{\text{OWA}}(q, \mathcal{I}) \subseteq \text{naïve}(q, \mathcal{I})$ , see Proposition 40.2. For the converse, assume  $\bar{a} \in \text{naïve}(q, \mathcal{I})$ . Then  $(\mathcal{I}, \eta) \models \varphi$ , where  $\eta$  maps  $\bar{x}$  to  $\bar{a}$ . Suppose  $h : \mathcal{I} \rightarrow D$  is an arbitrary valuation. Consider now  $V_{\text{Dom}(\varphi)}$ ; since it is a bijection, by  $\text{Dom}(\varphi)$ -genericity of  $\varphi$  we have  $(V_{\text{Dom}(\varphi)}(\mathcal{I}), V_{\text{Dom}(\varphi)} \circ \eta) \models \varphi$ . Next notice that  $h \circ V_{\text{Dom}(\varphi)}^{-1}$  is a homomorphism, as a composition of two homomorphisms. Thus, by homomorphism preservation, we obtain

$$(h \circ V_{\text{Dom}(\varphi)}^{-1} \circ V_{\text{Dom}(\varphi)}(\mathcal{I}), h \circ V_{\text{Dom}(\varphi)}^{-1} \circ V_{\text{Dom}(\varphi)} \circ \eta) \models \varphi,$$

or in other words  $(D, h \circ \eta) \models \varphi$ , and thus  $h(\bar{a}) \in q(D)$ , proving that  $\bar{a} \in \text{cert}_{\text{OWA}}(q, \mathcal{I})$ . This proves Proposition 40.5.  $\square$

To complete the proof of Theorem 40.4, we simply use the fact that every UCQ is preserved under homomorphisms, see Proposition 28.10. The last remaining bit is to show that formulae in PAG are preserved under strong onto homomorphism. This is done by a routine induction on the structure of PAG formulae, see Exercise 5.5.

The connection between homomorphism preservation and coincidence of naïve evaluation and certain answers extends beyond FO queries (Exercise 5.9), but for FO queries the classes of formulae in Theorem 40.4 are optimal (Exercise 5.10).

## Absolute Approximation of Certain Answers

Let us now summarize what we know about tractably evaluating queries on incomplete instances:

- we could use 3-valued evaluation,  $q_{3vl}$ , as SQL does (see Chapter 39) but then  $q_{3vl}(\mathcal{I})$  can produce both false positives and false negatives (see Proposition 39.9); or
- we could use naïve evaluation which always produces a superset of certain answers (Proposition 40.2) and thus no false negatives, but it can produce false positives.

As false positives may be viewed as the worse of the two, since they give completely false results, as opposed to omitting some true results, we now turn to *approximations* of FO queries and demonstrate that FO queries could be easily modified so that false positives could be avoided.

To achieve this, we transform an FO query  $q = (\varphi, \bar{x})$  into another FO query  $q^t = (\varphi^t, \bar{x})$  such that  $\text{naïve}(q^t, \mathcal{I}) \subseteq \text{cert}_{\text{CWA}}(q, \mathcal{I})$ . Thus, evaluating  $q^t$  with its usual DLOGSPACE complexity is guaranteed to avoid false positives. Of course there is a simple and totally uninteresting way to define such  $q^t$  simply by setting  $\varphi^t = \text{false}$ . We obviously want a better approximation; at the very least we need to ensure that  $q^t(\mathcal{I}) = q(\mathcal{I})$  if  $\mathcal{I}$  has no nulls, i.e., is a database rather than an incomplete instance.

To produce such a transformation we need two additional ingredients. One is an extra atomic formula  $\text{null}(x)$  testing if  $x$  is a null. Its semantics is

- $(\mathcal{I}, \eta) \models \text{null}(x)$  iff  $\eta(x) \in \text{Dom}_{\text{Null}}(\mathcal{I})$ .

The second ingredient we need is a notion of *unification* of tuples. Given two tuples  $\bar{u}$  and  $\bar{w}$  of the same length over  $\text{Const} \cup \text{Null}$ , we say they are *unifiable*, and write  $\bar{u} \uparrow \bar{w}$ , if there is a valuation  $h$  so that  $h(\bar{u}) = h(\bar{w})$ . For example, tuples  $(1, \perp_1, 2)$  and  $(\perp_1, 1, \perp_2)$  are unifiable via  $h(\perp_1) = 1$  and  $h(\perp_2) = 2$ , which maps both into the same tuple  $(1, 1, 2)$ . On the other hand,  $(1, \perp_1, 2)$  and  $(\perp_1, 2, \perp_2)$  are not unifiable.

We make two observations about unifiability. First, the condition  $\bar{x} \uparrow \bar{y}$  can be expressed in FO with  $\text{null}(\cdot)$ . For example, if  $\bar{x} = (x)$  and  $\bar{y} = (y)$ , it is expressed by  $\text{null}(x) \vee \text{null}(y) \vee (x = y)$ . For tuples of higher arity, it is a more complicated case analysis, checking which elements among  $(x_1, \dots, x_n)$  and  $(y_1, \dots, y_n)$  are the same, which are different, and which are nulls. In other words, the explicit listing of all equalities and inequalities between  $x_i$ s and  $y_j$ s, as well as explicit listing of all statements  $\text{null}(x_i)$  and  $\text{null}(y_j)$  or their negations, which determine whether  $\bar{x}$  and  $\bar{y}$  unify (see Exercise 5.11 for more details).

An FO formula for checking whether  $\bar{x} \uparrow \bar{y}$  is of exponential size in  $|\bar{x}|$ . At the same time, checking whether  $\bar{u} \uparrow \bar{w}$  can be done in linear time (see Exercise 5.12). Thus we can also assume that  $\bar{x} \uparrow \bar{y}$  are new atomic formulae added to FO, and they have linear time evaluation complexity.

With these formulae in place, we can now produce a translation that eliminates false positives.



**Definition 40.6: Translation  $\varphi \mapsto (\varphi^t, \varphi^f)$  of FO queries**

- $R(\bar{x})^t = R(\bar{x})$
  - $(x = y)^t = (x = y)$
  - $(\varphi \wedge \psi)^t = \varphi^t \wedge \psi^t$
  - $(\neg \varphi)^t = \varphi^f$
  - $(\exists x \varphi)^t = \exists x \varphi^t$
- 
- $(R(\bar{x}))^f = \neg \exists \bar{y} (R(\bar{y}) \wedge \bar{x} \uparrow \bar{y})$
  - $(x = y)^f = \neg(x = y) \wedge \neg \text{null}(x) \wedge \neg \text{null}(y)$
  - $(\varphi \wedge \psi)^f = \varphi^f \vee \psi^f$
  - $(\neg \varphi)^f = \varphi^t$
  - $(\exists x \varphi)^f = \forall x \varphi^f$

We now show that these translations give us a way to avoid one of false positives or false negatives. Specifically,  $\varphi^t$  avoids false positive answers while  $\neg \varphi^f$  avoids false negatives.

**Theorem 40.7**

For every incomplete instance  $\mathcal{I}$ , a FO query  $q = (\varphi, \bar{x})$ , its negation  $\neg q = (\neg \varphi, \bar{x})$ , and every valuation  $h$  on  $\mathcal{I}$ ,

$$h(\text{naïve}(q^t, \mathcal{I})) \subseteq q(h(\mathcal{I})) \text{ and } h(\text{naïve}(q^f, \mathcal{I})) \subseteq \neg q(h(\mathcal{I})).$$

*Proof.* Given  $q$  and  $\mathcal{I}$  as in the statement of the theorem, we will show the following two implications, for every valuation  $h$ :

1.  $(\mathcal{I}, \eta) \models \varphi^t \Rightarrow (h(\mathcal{I}), h \circ \eta) \models \varphi$
2.  $(\mathcal{I}, \eta) \models \varphi^f \Rightarrow (h(\mathcal{I}), h \circ \eta) \models \neg \varphi$

These clearly imply the statement: if  $\bar{a} \in \text{naïve}(q^t, \mathcal{I})$ , then  $(\mathcal{I}, \eta) \models \varphi^t$  for  $\eta$  such that  $\eta(\bar{x}) = \bar{a}$ , and thus  $(h(\mathcal{I}), h \circ \eta) \models \varphi$ , implying that  $h(\bar{a}) \in q(h(\mathcal{I}))$ , and likewise for  $q^f$ . We now prove these two statements by induction on the structure of  $\varphi$ .

Proof for  $\varphi^t$ .

- Assume  $\varphi$  is  $R(\bar{x})$ ; then  $\varphi^t = R(\bar{x})$ . If  $(\mathcal{I}, \eta) \models \varphi^t$  then  $R(\eta(\bar{x}))$  is an atom in  $\mathcal{I}$  and hence  $R(h(\eta(\bar{x})))$  is a fact in  $h(\mathcal{I})$  and hence  $(h(\mathcal{I}), h \circ \eta) \models \varphi$ .
- Assume  $\varphi$  is  $x = y$ , and  $\varphi^t = \varphi$ . If  $(\mathcal{I}, \eta) \models \varphi^t$  then  $\eta(x) = \eta(y)$  and hence  $h(\eta(x)) = h(\eta(y))$  implying  $(h(\mathcal{I}), h \circ \eta) \models \varphi$ .

- Assume  $\varphi$  is  $\psi \wedge \chi$ , and  $\varphi^t = \psi^t \wedge \chi^t$ . If  $(\mathcal{I}, \eta) \models \varphi^t$  then  $(\mathcal{I}, \eta) \models \psi^t$  and  $(\mathcal{I}, \eta) \models \chi^t$  and by the induction hypothesis  $(h(\mathcal{I}), h \circ \eta) \models \psi$  and  $(h(\mathcal{I}), h \circ \eta) \models \chi$  and therefore  $(h(\mathcal{I}), h \circ \eta) \models \varphi$ .
- Assume  $\varphi$  is  $\neg\psi$ ; then  $\varphi^t = \psi^f$ . If  $(\mathcal{I}, \eta) \models \varphi^t$  then  $(\mathcal{I}, \eta) \models \psi^f$  and by the induction hypothesis  $(h(\mathcal{I}), h \circ \eta) \models \neg\psi$  and therefore  $(h(\mathcal{I}), h \circ \eta) \models \varphi$ .
- Assume  $\varphi$  is  $\exists x \psi$ ; then  $\varphi^t = \exists x \psi^t$ . If  $(\mathcal{I}, \eta) \models \varphi^t$  then  $(\mathcal{I}, \eta[a/x]) \models \psi^t$  for some  $a \in \text{Dom}(\mathcal{I})$  and by the induction hypothesis  $(h(\mathcal{I}), h \circ \eta[a/x]) \models \psi$  and therefore  $(h(\mathcal{I}), h \circ \eta) \models \varphi$  as  $h(a)$  is a witness for  $x$ .

Proof for  $\varphi^f$ .

- Assume  $\varphi$  is  $R(\bar{x})$ . Towards contradiction, suppose that  $(\mathcal{I}, \eta) \models \varphi^f$  but  $(h(\mathcal{I}), h \circ \eta) \not\models \varphi$ . The latter means that  $R(h(\eta(\bar{x})))$  is a fact of  $h(\mathcal{I})$ . This happens if there an atom  $R(\bar{u})$  in  $\mathcal{I}$  such that  $h(\bar{u}) = h(\eta(\bar{x}))$ , i.e.,  $\bar{u} \uparrow \eta(\bar{x})$ . Thus,  $(\mathcal{I}, \eta) \models \exists \bar{y} (R(\bar{y}) \wedge \bar{y} \uparrow \bar{x})$ , i.e.,  $(\mathcal{I}, \eta) \models \varphi^f$ , which is a contradiction.
- Assume  $\varphi$  is  $x = y$ , and  $\varphi^f = (x \neq y) \wedge \neg \text{null}(x) \wedge \neg \text{null}(y)$ . If  $(\mathcal{I}, \eta) \models \varphi^f$  then  $\eta(x)$  and  $\eta(y)$  are different constants, and thus so are  $h(\eta(x))$  and  $h(\eta(y))$ . Hence  $(h(\mathcal{I}), h \circ \eta) \models \neg(x = y)$ , i.e.,  $(h(\mathcal{I}), h \circ \eta) \models \neg\varphi$ .
- Assume  $\varphi$  is  $\psi \wedge \chi$ , and  $\varphi^f = \psi^f \vee \chi^f$ . If  $(\mathcal{I}, \eta) \models \varphi^f$  then by the induction hypothesis  $(h(\mathcal{I}), h \circ \eta) \models \neg\psi \vee \neg\chi$  which is of course  $\neg(\psi \wedge \chi)$ .
- Assume  $\varphi$  is  $\neg\psi$ ; then  $\varphi^f = \psi^t$ . If  $(\mathcal{I}, \eta) \models \varphi^f$  then by the induction hypothesis  $(h(\mathcal{I}), h \circ \eta) \models \psi$  and therefore  $(h(\mathcal{I}), h \circ \eta) \models \neg\varphi$ .
- Assume  $\varphi$  is  $\exists x \psi$ ; then  $\varphi^f = \forall x \psi^f$ . If  $(\mathcal{I}, \eta) \models \varphi^f$ , then for every  $a \in \text{Dom}(\mathcal{I})$  we have  $(\mathcal{I}, \eta[a/x]) \models \psi^f$  and thus  $(h(\mathcal{I}), h \circ \eta[a/x]) \models \neg\psi$  and hence  $(h(\mathcal{I}), h \circ \eta) \models \neg\exists x \psi$  which is  $\neg\varphi$ . This concludes the proof.  $\square$

#### Corollary 40.8: Avoiding False Positives and Negatives

For every FO query  $q$ , the query  $q^t$  returns no false positives, and the query  $\neg q^f$  returns no false negatives. Moreover, on databases without nulls,  $q(D) = q^t(D) = \neg q^f(D)$ .

*Proof.* The first statement follows immediately from the theorem and the definition of certain answers. The second is by inspection of the translation, taking into account that without nulls,  $\bar{x} \uparrow \bar{y}$  can be replaced by  $\bar{x} = \bar{y}$  and  $\neg \text{null}(x)$  by **true**.  $\square$

As an illustration, we return to the query  $\varphi(x) = R(x) \wedge \neg S(x)$  seen earlier, for which naïve evaluation produces a false positive answer. Applying the transformation, we have

$$\begin{aligned}
\varphi^{\mathbf{t}}(x) &= (R(x) \wedge \neg S(x))^{\mathbf{t}} = (R(x))^{\mathbf{t}} \wedge (\neg S(x))^{\mathbf{t}} \\
&= R(x) \wedge (S(x))^{\mathbf{f}} = R(x) \wedge \neg \exists y (S(y) \wedge x \uparrow y) \\
&= R(x) \wedge \neg \exists y (S(y) \wedge (\text{null}(x) \vee \text{null}(y) \vee (x = y)))
\end{aligned}$$

It is then easy to check that  $\varphi^{\mathbf{t}}$  returns the empty set when evaluated naïvely on an incomplete instance with atoms  $R(1)$  and  $S(\perp)$ .

We remark that while simple to present, the translation of Definition 40.6 is not going to be very efficient when implemented in relational algebra for queries involving negation. A couple of examples illustrating this are presented in Exercise 5.13, and an alternative more efficient translation is given in Exercise 5.14.

We conclude by showing that the above translation also helps us eliminate false positives or false negatives under the 3-valued evaluation used in SQL. Recall that at the end of Chapter 39, we defined  $q_{3\text{vl}}(\mathcal{I})$  as  $\{\eta(\bar{x}) \mid \text{Eval}(\varphi, \mathcal{I}, \eta) = \text{true}\}$ , for a query  $q = (\varphi, \bar{x})$ . We now define  $q_{3\text{vl}}^{\mathbf{t}}$  and  $q_{3\text{vl}}^{\mathbf{f}}$  in the same way, but with  $\varphi$  replaced by  $\varphi^{\mathbf{t}}$  and  $\varphi^{\mathbf{f}}$ . Of course we also need to extend  $\text{Eval}$  to new formulae  $\text{null}(x)$  and  $\bar{x} \uparrow \bar{y}$ . These formulae can only take values **true** and **false**. We have  $\text{Eval}(\text{null}(x), \mathcal{I}, \eta) = \text{true}$  iff  $\eta(x)$  is a null, and  $\text{Eval}((\bar{x} \uparrow \bar{y}), \mathcal{I}, \eta) = \text{true}$  iff  $\eta(\bar{x})$  and  $\eta(\bar{y})$  unify.

Then we have an analog of Corollary 40.8

**Corollary 40.9: Avoiding False Positives and Negatives**

For every FO query  $q$ , the query  $q_{3\text{vl}}^{\mathbf{t}}$  returns no false positives, and the complement of the query  $q_{3\text{vl}}^{\mathbf{f}}$  returns no false negatives. Moreover, on databases without nulls, both of these queries coincide with  $q$ .

The proof is by simple adaptation of the inductive proof of Theorem 40.7.

## Probabilistic Approximation of Certain Answers

We now address the following question: how close to certainty is the query answer produced by the naïve evaluation? It turns out that the answer is “very close”. Recall that  $\bar{a} \in \text{cert}_{\text{CWA}}(q, \mathcal{I})$  if  $h(\bar{a}) \in q(h(\mathcal{I}))$  for every valuation  $h$ . The framework in which we pose our question is the following:

- Pick a valuation  $h$  uniformly at random. What is the probability that  $h(\bar{a}) \in q(h(\mathcal{I}))$  holds?

We then show – again, describing the intuition behind the result – that for *all* FO queries  $q$ ,

- this probability is either 0 or 1;

- it is 1 if and only if  $\bar{a} \in \text{naïve}(q, \mathcal{I})$ .

Thus, naïve evaluation, while only providing absolute guarantees of certainty of query answers for a restricted class of queries, gives high probability guarantees for many more queries.

To prove such a result, there is still a technical problem: how do we choose a valuation  $h : \text{Dom}_{\text{Null}}(\mathcal{I}) \rightarrow \text{Const}$  uniformly at random? Since the set  $\text{Const}$  is infinite, there are infinitely many such valuations, and thus no uniform distribution on them. Instead, we can use a trick inspired by the 0–1 law seen in Chapter 32. Namely, we consider finite subsets of  $\text{Const}$ , restricted to which there are finitely many valuations. We then look at the asymptotic behavior of the probability that a tuple is in the answer.

To do this, fix some enumeration of elements of  $\text{Const}$  as  $\{c_1, c_2, c_3, \dots\}$ . The exact enumeration is not important, as we will see shortly. Let  $\text{Const}_n$  be the set  $\{c_1, \dots, c_n\}$ . Let  $\mathcal{H}_n(\mathcal{I})$  be the finite set of all valuations  $h : \text{Dom}_{\text{Null}}(\mathcal{I}) \rightarrow \text{Const}_n$ . For a  $k$ -ary query  $q$ , an incomplete instance  $\mathcal{I}$ , and a  $k$ -ary tuple over  $\text{Dom}(\mathcal{I})$ , we now define

$$\mu_n(q, \mathcal{I}, \bar{a}) = \frac{|\{h \in \mathcal{H}_n \mid h(\bar{a}) \in q(h(\mathcal{I}))\}|}{|\mathcal{H}_n(\mathcal{I})|}$$

as the probability that a valuation  $h$ , picked uniformly at random from the finite set  $\mathcal{H}_n(\mathcal{I})$ , witnesses that  $\bar{a}$  is in the answer to  $q$  on  $\mathcal{I}$ , i.e., the condition  $h(\bar{a}) \in q(h(\mathcal{I}))$ . Then we define its asymptotic behavior:

$$\mu(q, \mathcal{I}, \bar{a}) = \lim_{n \rightarrow \infty} \mu_n(q, \mathcal{I}, \bar{a}).$$

Notice that this value does not depend on the exact enumeration of  $\text{Const}$ : as long as  $n$  is big enough so that it contains  $\text{Dom}_{\text{Const}}(\mathcal{I})$  and  $\text{Dom}(q)$ , the value of  $\mu_n(q, \mathcal{I}, \bar{a})$  does not depend on how the rest of  $\text{Const}$  is enumerated.

Intuitively,  $\mu(q, \mathcal{I}, \bar{a})$  gives us the probability that a valuation – now without a restriction on its range – picked randomly witnesses that  $\bar{a}$  is an answer. The closer this is to 1, the better  $\bar{a}$  is as the answer to  $q$ . We now show that naïve evaluation gives us precisely the tuples that are good (although not absolutely certain) answers.

#### Theorem 40.10

For every FO  $k$ -ary query  $q$ , every incomplete instance  $\mathcal{I}$ , and every  $k$ -tuple  $\bar{a}$ ,

$$\mu(q, \mathcal{I}, \bar{a}) = \begin{cases} 1 & \text{if } \bar{a} \in \text{naïve}(q, \mathcal{I}) \\ 0 & \text{if } \bar{a} \notin \text{naïve}(q, \mathcal{I}). \end{cases}$$

*Proof.* Consider again simple valuations  $h : \text{Dom}_{\text{Null}}(\mathcal{I}) \rightarrow \text{Const}_n$  that are bijections whose range contains no elements of  $\text{Dom}_{\text{Const}}(\mathcal{I})$  or  $\text{Dom}(q)$ ,

and let  $\mathcal{H}_n^s(\mathcal{I})$  stand for the set of such valuations. We next show that  $\lim_{n \rightarrow \infty} |\mathcal{H}_n^s(\mathcal{I})|/|\mathcal{H}_n(\mathcal{I})| = 1$ .

Assume  $|\text{Dom}_{\text{Null}}(\mathcal{I})| = m$  and  $|\text{Dom}_{\text{Const}}(\mathcal{I}) \cup \text{Dom}(q)| = k$ . A valuation  $h$  fails to be simple if one of the following is true.

- The range of  $h$  contains an element of  $\text{Dom}_{\text{Const}}(\mathcal{I}) \cup \text{Dom}(q)$ , i.e.,  $h(\perp)$  is in that set for some  $\perp \in \text{Dom}_{\text{Null}}(\mathcal{I})$ . There are at most  $mk \cdot n^{m-1}$  such valuations. Indeed, we pick one of  $m$  nulls  $\perp$  and one of  $k$  elements of  $\text{Dom}_{\text{Const}}(\mathcal{I}) \cup \text{Dom}(q)$  it is mapped into, and on the remaining  $m-1$  nulls this valuation can be arbitrary.
- The valuation  $h$  is not a bijection: there are two nulls  $\perp$  and  $\perp'$  such that  $h(\perp) = h(\perp')$ . The number of such valuations is at most  $\binom{m}{2} \cdot n \cdot n^{m-2}$ . Indeed, there are  $\binom{m}{2}$  ways to pick  $\perp$  and  $\perp'$ , then  $n$  possible assignments to an element of  $\text{Const}_n$ , and on the remaining  $m-2$  elements of  $\text{Dom}_{\text{Null}}(\mathcal{I})$ , the valuation can be chosen arbitrarily.

Summing up, the number of valuations not in  $\mathcal{H}_n^s(\mathcal{I})$  is at most  $(mk + \binom{m}{2}) \cdot n^{m-1}$ . Since  $|\mathcal{H}_n(\mathcal{I})| = n^m$ , we have

$$\lim_{n \rightarrow \infty} \frac{|\mathcal{H}_n^s(\mathcal{I})|}{|\mathcal{H}_n(\mathcal{I})|} \geq \lim_{n \rightarrow \infty} \frac{n^m - (mk + \binom{m}{2}) \cdot n^{m-1}}{n^m} = 1 - \lim_{n \rightarrow \infty} \frac{mk + \binom{m}{2}}{n} = 1.$$

Next assume  $\bar{a} \in \text{naïve}(q, \mathcal{I})$ . By  $\text{Dom}(q)$ -genericity of FO queries (Exercise 1.2), we have  $h(\bar{a}) \in q(h(\mathcal{I}))$  for each  $h \in \mathcal{H}_n^s(\mathcal{I})$ . Hence,  $\mu(q, \mathcal{I}, \bar{a}) \geq \lim_{n \rightarrow \infty} |\mathcal{H}_n^s(\mathcal{I})|/|\mathcal{H}_n(\mathcal{I})| = 1$  which implies  $\mu(q, \mathcal{I}, \bar{a}) = 1$ . If  $\bar{a} \notin \text{naïve}(q, \mathcal{I})$ , we apply the previous argument to the complement of  $q$ , obtaining  $\mu(\neg q, \mathcal{I}, \bar{a}) = 1$  from which  $\mu(q, \mathcal{I}, \bar{a}) = 0$  follows.  $\square$

---

## Exercises

**Exercise 5.1.** . We have seen in Theorem 39.6 that combined complexity of certain answers under OWA is undecidable. Prove that their data complexity is undecidable too. (Hint: use the fixed formula to encode the universal Turing machine, and the database to encode a Turing machine whose run is simulated.)

**Exercise 5.2.** Prove that the following problem is NP-complete: given an incomplete instance  $\mathcal{I}$  and a database  $D$  of the same schema, is  $D \in \llbracket \mathcal{I} \rrbracket_{\text{CWA}}$ ? Prove that the same problem can be solved in polynomial time if no null repeats in  $\mathcal{I}$  and give a  $O(n^2\sqrt{n})$  algorithm for it. (Hint: use bipartite matching.) What about the same problems under OWA?

**Exercise 5.3.** Prove that every query under the 3-valued interpretation can be expressed as an FO query under the standard two-valued interpretation with the addition of an extra unary predicate  $\text{isnull}(\cdot)$  that checks whether an element is in  $\text{Dom}_{\text{Null}}(\mathcal{I})$ .

**Exercise 5.4.** The class PAG has a natural relational algebra counterpart. Recall the division operation  $R \div S$ , where  $R$  has  $n$  attributes,  $S$  has  $m$  attributes, and  $n > m$ . Its result contains all tuples  $\bar{a}$  in the projection of  $R$  on its first  $n - m$  attributes such that  $(\bar{a}, \bar{c}) \in R$  for every  $\bar{c} \in S$ .

Prove that PAG queries are exactly the same as those in the fragment of relational algebra that contains selection, projection, join, union, and division by a relation. That is, if  $q$  is a query in this fragment and  $S$  is a relation in the database, then  $q \div S$  is in the fragment.

**Exercise 5.5.** Prove by induction on the structure of the formulae that all PAG queries are preserved under strong onto homomorphisms.

**Exercise 5.6.** Prove that there are FO queries preserved under strong onto homomorphisms on all finite databases that are not expressible in PAG.

**Exercise 5.7.** We say that a incomplete instance  $\mathcal{I}$  is at least as informative as  $\mathcal{I}'$ , written  $\mathcal{I}' \preceq_* \mathcal{I}$ , if  $\llbracket \mathcal{I} \rrbracket_* \subseteq \llbracket \mathcal{I}' \rrbracket_*$ , where  $*$  is OWA or CWA. The intuition is as follows: the less we know about an instance the more databases it can possibly denote. Prove the following:

- Both  $\preceq_{\text{OWA}}$  and  $\preceq_{\text{CWA}}$  are *pre-orders*, i.e., reflexive and transitive.
- $\mathcal{I}' \preceq_{\text{OWA}} \mathcal{I}$  iff there is a homomorphism  $h : \mathcal{I}' \rightarrow \mathcal{I}$ ;
- $\mathcal{I}' \preceq_{\text{CWA}} \mathcal{I}$  iff there is a homomorphism  $h : \mathcal{I}' \rightarrow \mathcal{I}$  such that  $h(\mathcal{I}') = \mathcal{I}$  (recall that such homomorphisms are called strong surjective homomorphisms).

**Exercise 5.8.** Define the *most informative answer* to a query  $q$  on an incomplete instance  $\mathcal{I}$ , denoted by  $\text{mia}_*(q, \mathcal{I})$ , as the greatest lower bound of  $\{q(D) \mid D \in \llbracket \mathcal{I} \rrbracket_*\}$ , in the pre-order  $\preceq_*$ , where  $*$  is OWA or CWA. That is,

1.  $\text{mia}_*(q, \mathcal{I}) \preceq_* q(D)$  for each  $D \in \llbracket \mathcal{I} \rrbracket_*$ ; and
2. for each  $\mathcal{I}_0$  satisfying the above condition,  $\mathcal{I}_0 \preceq_* \text{mia}_*(q, \mathcal{I})$ .

Note that since  $\preceq_*$  are preorders,  $\text{mia}_*(q, \mathcal{I})$  is not uniquely defined, but any two instances satisfying 1 and 2 are homomorphically equivalent with respect to arbitrary or strong onto homomorphisms.

Prove the following:

- $\text{cert}_*(q, \mathcal{I}) \preceq_* \text{mia}_*(q, \mathcal{I})$ ;
- $\text{mia}_{\text{OWA}}(q, \mathcal{I})$  may not exist for FO queries;
- $\text{mia}_{\text{CWA}}(q, \mathcal{I})$  exists for all FO queries, but there are queries for which the size of the smallest instance satisfying the conditions for  $\text{mia}_{\text{CWA}}(q, \mathcal{I})$  is exponential in the size of  $\mathcal{I}$ ;
- if  $q$  is a UCQ, then  $\text{mia}_{\text{CWA}}(q, \mathcal{I})$  exists and its size is polynomial in the size of  $\mathcal{I}$ .

**Exercise 5.9.** The connection between homomorphism preservation, naïve evaluation and certain answers works beyond FO queries. Recall that a query  $q$  is  $C$ -generic if, for every bijection  $\rho : \text{Const} \rightarrow \text{Const}$  that is the identity on  $C$ , we have  $q(\rho(D)) = \rho(q(D))$  for each database  $D$ . For a  $C$ -generic query, we define  $\text{naïve}(q, \mathcal{I})$  by  $h^{-1}(q(h(D)))$ , where  $h$  is a bijection that is the identity on  $C \cup \text{Dom}_{\text{Const}}(\mathcal{I})$ . As a first simple exercise, verify that this definition does not depend on the choice of such  $h$ .

A  $C$ -generic query is preserved under (strong onto) homomorphisms if  $(D, \bar{a}) \rightarrow_C (D', \bar{b})$  and  $\bar{a} \in q(D)$  imply  $\bar{b} \in q(D')$ . Prove the following:

1. If  $q$  is  $C$ -generic and preserved under homomorphisms, then  $\text{naïve}(q, \mathcal{I}) = \text{cert}_{\text{OWA}}(q, \mathcal{I})$

2. If  $q$  is  $C$ -generic and preserved under strong onto homomorphisms, then  $\text{naïve}(q, \mathcal{I}) = \text{cert}_{\text{CWA}}(q, \mathcal{I})$

**Exercise 5.10.** Prove that if  $q$  is a Boolean FO query and  $\text{cert}_{\text{OWA}}(q, \mathcal{I}) = \text{naïve}(q, \mathcal{I})$  for every incomplete instance, then  $q$  is equivalent to a UCQ (hint: use Theorem 28.11).

In a certain sentence, the result about the class PAG is the best one can get. The following exercise is for the reader with a strong background in logic and model theory. Consider an FO sentence that is preserved under strong onto homomorphisms over both finite and infinite databases. Prove that such a sentence is equivalent to a PAG sentence.

**Exercise 5.11.** Given two tuples  $\bar{x} = (x_1, \dots, x_n)$  and  $\bar{y} = (y_1, \dots, y_n)$ , consider conjunctions of the form

$$\bigwedge_{i,j \in [n]} (\neg)(x_i = x_j) \wedge \bigwedge_{i,j \in [n]} (\neg)(y_i = y_j) \wedge \bigwedge_{i,j \in [n]} (\neg)(x_i = y_j) \\ \wedge \bigwedge_{i \in [n]} (\neg)\text{null}(x_i) \wedge \bigwedge_{j \in [n]} (\neg)\text{null}(y_j) \quad (43.1)$$

where the atomic formulae are optionally preceded by negation. Since there are  $3n^2 + 2n$  possible atomic subformulae of these formulae, there are  $2^{3n^2+2n}$  formulae (43.1).

Prove that  $\bar{x} \uparrow \bar{y}$  can be expressed as a disjunction of formulae of the form (43.1).

**Exercise 5.12.** Given two tuples  $\bar{u} = (u_1, \dots, u_n)$  and  $\bar{w} = (w_1, \dots, w_n)$  containing constants and nulls, prove that checking whether  $\bar{u}$  and  $\bar{w}$  unify can be done in time  $O(n)$ .

**Exercise 5.13.** In this exercise we see two examples of the high cost of translation  $\varphi \mapsto (\varphi^t, \varphi^f)$  when queries are written explicitly in relational algebra. Let  $\text{Dom}_{\mathbf{S},A}$  be a query that returns  $\text{Dom}(D)$  as a single unary relation of attribute  $A$ . Show that for a unary relation  $R$  of attribute  $B$ , the query  $R(x)^f$  is expressed as

$$R^f = \text{Dom}_{\mathbf{S},A} - \pi_A \left( \sigma_{(A=B) \vee \text{null}(A) \vee \text{null}(B)} (\text{Dom}_{\mathbf{S},A} \times R) \right).$$

Show that if we also have a unary relation  $S$  with attribute  $C$ , the query  $(R(x) \vee S(y))^f$  is expressed in RA as

$$R^f \times \text{Dom}_{\mathbf{S},C} \cup \text{Dom}_{\mathbf{S},B} \times S^f.$$

**Exercise 5.14.** Consider a different translation of FO formulae  $\varphi \mapsto (\varphi^\square, \varphi^\diamond)$  given by



- $(R(\bar{x}))^\square = R(\bar{x})$
- $(x = y)^\square = (x = y)$
- $(x \neq y)^\square = (x \neq y) \wedge \neg \text{null}(x) \wedge \neg \text{null}(y)$
- $(\varphi \wedge \psi)^\square = \varphi^\square \wedge \psi^\square$
- $(\neg \varphi)^\square = \neg \exists \bar{y} (\varphi^\diamond(\bar{y}) \wedge \bar{x} \uparrow \bar{y})$
- $(\exists x \varphi)^\square = \exists x \varphi^\square$

and

- $R(\bar{x})^\diamond = R(\bar{x})$
- $(x = y)^\diamond = (x = y) \vee \text{null}(x) \vee \text{null}(y)$
- $(x \neq y)^\diamond = (x \neq y)$
- $(\varphi \wedge \psi)^\diamond = \varphi^\diamond \wedge \exists \bar{y} (\bar{x} \uparrow \bar{y} \wedge \psi^\diamond[\bar{y}/\bar{x}])$
- $(\neg \varphi)^\diamond = \neg \varphi^\square$
- $(\exists x \varphi)^\diamond = \exists x \varphi^\diamond$

Prove that for every incomplete instance  $\mathcal{I}$ , a FO query  $q = (\varphi, \bar{x})$ , and valuation  $h$  on  $\mathcal{I}$ , we have  $h(q^\square(\mathcal{I})) \subseteq q(h(\mathcal{I})) \subseteq h(q^\diamond(D))$ . Conclude that  $(\varphi^\square, \bar{x})$  returns no false positives and  $(\varphi^\diamond, \bar{x})$  returns no false negatives. Work out this translation in RA, and show that the subqueries of the form  $\text{Dom}_{\mathbf{S}, A}$  from Exercise 5.13 are not present in the translation.

**Exercise 5.15.** Extend the framework of measuring the quality of answers of Theorem 40.10 to the case of integrity constraints. Let  $\Sigma$  be a set of integrity constraints; define  $\mu_n(q|\Sigma, \mathcal{I}, \bar{a})$  as the ratio

$$\frac{|\{h \in \mathcal{H}_n(\mathcal{I}) \mid h(\mathcal{I}) \models \Sigma \text{ and } h(\bar{a}) \in q(h(\mathcal{I}))\}|}{|\{h \in \mathcal{H}_n(\mathcal{I}) \mid h(\mathcal{I}) \models \Sigma\}|}$$

and then, as before, set  $\mu(q|\Sigma, \mathcal{I}, \bar{a})$  to be  $\lim_{n \rightarrow \infty} \mu_n(q|\Sigma, \mathcal{I}, \bar{a})$ .

- Give an example of a Boolean CQ  $q$ , a single inclusion constraint  $\Sigma$ , and an incomplete instance  $\mathcal{I}$  such that  $\mu(q|\Sigma, \mathcal{I}) \notin \{0, 1\}$ .
- Prove that for every FO query  $q$ , every set of FO integrity constraints  $\Sigma$ , every incomplete instance  $\mathcal{I}$  and tuple  $\bar{a}$ , the value  $\mu(q|\Sigma, \mathcal{I}, \bar{a})$  exists, and is a rational number. Provide an exponential time algorithm (in data complexity) to compute it.
- Prove that if  $\Sigma$  contains only functional dependencies, then  $\mu(q|\Sigma, \mathcal{I}, \bar{a})$  is either 0 or 1. Prove furthermore that it equals 1 iff  $\mu(q, \mathcal{I}', \bar{a}) = 1$ , where  $\mathcal{I}'$  is the result of chasing  $\mathcal{I}$  with  $\Sigma$ .

---

## Bibliographic Comments



# Background for Tree- and Graph-Structured Data



## Background For Tree and Graph Structured Data

Conceptually, tree- and graph-structured data is fundamentally different from relational data. The reason is that tree- and graph-structured data does not necessarily need to adhere to a schema. The reader does not need to be familiar with the entire background chapter (Chapter 2) to start with Parts VIII and X. We assume here that the reader is familiar with the *Basic Notions and Notation* from Chapter 2 and introduce everything else here or give the reader an explicit pointer to Chapter 2.

### Graphs

Throughout Parts VIII and X we assume that we have a countably infinite set

Nodes of *nodes*

available, from which we will take the nodes in trees and graphs. We also assume that we have a countably infinite set

Lab of *labels*.

Furthermore, we assume the existence of a special symbol  $*$  that does not occur in any of the aforementioned sets and which we will use as a “wildcard” in query languages over trees.

We next define node- and edge-labeled directed graphs, which will be the basis of different data models that we will use throughout Parts VIII and X. In particular, we will use special cases of the following definition. The main reason why we introduce this general definition here is because we want to define *queries* and some of their computational problems in this chapter.

**Definition 55.1: Node- and Edge-Labeled, Directed Graph**

Let  $k \in \mathbb{N}$ . A  $k$ -sorted node- and edge-labeled, directed graph with wildcards is a tuple

$$G = (V, E_1, \dots, E_k, \text{lab}) ,$$

where

- $V \subseteq \mathbf{Nodes}$  is a finite, nonempty set of *nodes*,
- $E_i \subseteq V \times V$  is a set of directed *edges* for every  $i \in [k]$ , and
- $\text{lab}: V \cup E_1 \cup \dots \cup E_k \rightarrow \mathbf{Lab} \cup \{*\}$  is a partial function that assigns to every node and edge its label from  $\mathbf{Lab}$  or the wildcard symbol.

We say that  $G$  is *node-labeled* (resp., *edge-labeled*) if  $\text{Dom}(\text{lab}) = V$  (resp.,  $\text{Dom}(\text{lab}) = E_1 \cup \dots \cup E_k$ ).

We will use the  $k$  different sorts of edges to be able to assign different roles to edges in the graph. For example, in Part [VIII](#) we will use this feature to define *child* edges and *next-sibling* edges in trees.

Usually, we will use a restricted labeling function that does not assign wildcards, that is,  $\text{lab}: V \cup E_1 \cup \dots \cup E_k \rightarrow \mathbf{Lab}$ . If this is the case, we will omit “with wildcards” when we refer to  $G$ , that is, we will talk about a  $k$ -sorted node-labeled, directed graph, etc. Furthermore, we will often leave “ $k$ -sorted” implicit, since the number  $k$  can be inferred from the tuple  $(V, E_1, \dots, E_k, \text{lab})$ .

**Example 55.2**

Figure [55.1](#) has a graphical representation of a 2-sorted node- and edge-labeled directed graph  $G = (V, E_1, E_2, \text{lab})$  where

- $V = \{v_1, \dots, v_{10}\}$ ,
- $E_1 = \{(v_1, v_2), (v_2, v_3), (v_2, v_4), (v_1, v_5), (v_5, v_6), (v_6, v_7), (v_6, v_8), (v_5, v_9), (v_5, v_{10})\}$ ,
- $E_2 = \{(v_2, v_5), (v_3, v_4), (v_6, v_9), (v_7, v_8), (v_9, v_{10})\}$ , and
- $\text{lab}$  is defined in the Figure.

The two edge sorts are represented with solid and dashed arrows, respectively.

Since presenting graphs as in Figure [55.1](#) is very verbose, we will leave the elements of  $V$  implicit whenever we can. For example, the graph of Figure [55.1](#) can be represented more compactly as in Figure [55.2](#). Notice that, for the sake of brevity, we omitted some of the elements of  $V$ . (We only mention  $v_4$ ,  $v_5$ ,  $v_8$ , and  $v_{10}$ .) We may do this in cases where only some of the elements of  $V$  are relevant. For instance, if we view

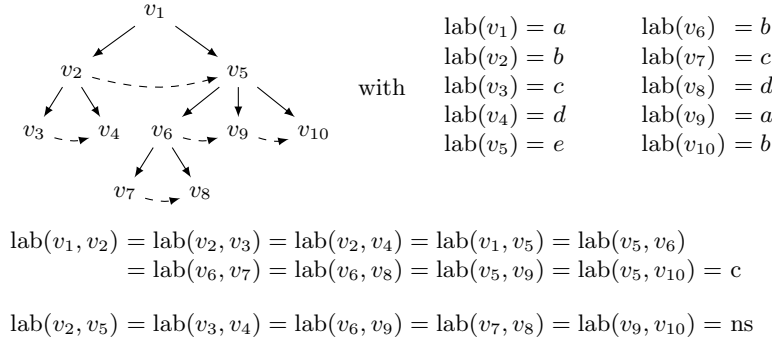


Fig. 55.1: Example of a Node- and Edge-Labeled Graph

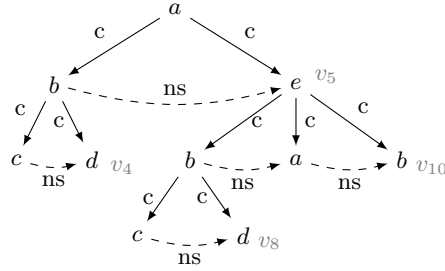


Fig. 55.2: Less verbose representation of the graph in Figure 55.1, leaving some of the elements of Nodes implicit

Figure 55.2 as a tree-like structure, we may want to say that  $v_4$ ,  $v_5$ ,  $v_8$ , and  $v_{10}$  are “rightmost children of their respective parents”. The nodes we omitted are irrelevant for this statement.

In Parts VIII-X, we will use the letters  $u$  and  $v$  for referring to nodes. Let  $G = (V, E_1, \dots, E_k, \text{lab})$  and  $E = E_1 \cup \dots \cup E_k$ . An *undirected path* in  $G$  is a non-empty sequence of nodes

$$\pi = v_0 v_1 \cdots v_n$$

where  $(v_{i-1}, v_i) \in E$  or  $(v_i, v_{i-1}) \in E$  for every  $i \in [n]$ . We say that  $\pi$  is *from*  $v_0$  *to*  $v_n$  and has *length*  $n$ . (The path of length zero is from  $v_0$  to  $v_0$ .) A *directed path* in  $G$  is an undirected path where  $(v_{i-1}, v_i) \in E$  for every  $i \in [n]$ . We assume in Parts VIII-X that paths are directed unless stated otherwise. A graph is *connected* if, for every pair of nodes  $u, v \in V$ , there is an undirected path from  $u$  to  $v$ . Notice that these definitions are consistent with the corresponding definitions for directed graphs in Chapter 2.



## Data Graphs

For the following definition, we recall from Chapter 2 that we use  $\text{Const}$  throughout the book to denote the set of *values* in databases.

### Definition 55.3: $k$ -Sorted Data Graph

Let  $k \in \mathbb{N}$ . A  $k$ -sorted data graph is a tuple

$$G = (V, E_1, \dots, E_k, \text{lab}, d),$$

where:

- $(V, E_1, \dots, E_k, \text{lab})$  is a  $k$ -sorted node- and edge-labeled directed graph, and
- $d: V \cup E_1 \cup \dots \cup E_k \rightarrow \text{Const}$  is a partial function that assigns to every node and edge its *data value*.

We write  $\mathcal{G}$  for the set of all  $k$ -sorted data graphs for some  $k \in \mathbb{N} - \{0\}$ .

When the number  $k$  of sorts is not important, we may use the term *data graph* to refer to a  $k$ -sorted data graph for some  $k \in \mathbb{N} - \{0\}$ .

Notice that we assume a single data value per node. This is just for notational simplicity and is not a restriction at all, since a node  $u$  with  $k > 1$  attributes can be modeled by a node  $u$  with  $k$  outgoing edges, leading to  $k$  new nodes containing the data values.

## Queries

Next we define queries over data graphs. Since Parts VIII–X will sometimes consider strict subsets of the data graphs (e.g., node-labeled trees, edge-labeled graphs), we define queries for subsets  $\mathcal{G}'$  of  $\mathcal{G}$ .

### Definition 55.4: Queries and Query Languages over Graphs

Let  $\mathcal{G}' \subseteq \mathcal{G}$ . A  $k$ -ary data graph query over  $\mathcal{G}'$  is a function of the form

$$q : \mathcal{G}' \rightarrow (\text{Nodes} \cup \text{Const})^k.$$

A *data graph query language* is a set of data graph queries.

## Key Problems: Query Evaluation and Query Analysis

We now define some key problems concerning tree and graph-structured data. In their most common form, these problems are parameterized by a data

graph query language  $\mathcal{L}$ . Such query languages in Parts VIII-X will always be associated to a *data model*, i.e., a subset  $\mathcal{G}_{\mathcal{L}} \subseteq \mathcal{G}$  onto which the semantics of  $\mathcal{L}$  is defined. The intention is that queries from  $\mathcal{L}$  will always be evaluated over elements of  $\mathcal{G}_{\mathcal{L}}$ .

#### Query Evaluation

##### Problem: $\mathcal{L}$ -Evaluation

**Input:** A query  $q$  from  $\mathcal{L}$ , a data graph  $G \in \mathcal{G}_{\mathcal{L}}$ , a tuple  $\bar{a}$  over  $\text{Nodes} \cup \text{Const}$   
**Output:** **true** if  $\bar{a} \in q(G)$ , and **false** otherwise

#### Query Containment and Equivalence

Let  $\mathcal{L}$  be a data graph query language. We say that a query  $q \in \mathcal{L}$  is *contained* in a query  $q' \in \mathcal{L}$ , written as  $q \subseteq q'$ , if  $q(G) \subseteq q'(G)$  for every data graph  $G \in \mathcal{G}_{\mathcal{L}}$ . Query  $q$  is *equivalent* to  $q'$ , written as  $q \equiv q'$ , if  $q(G) = q'(G)$  for every data graph  $G \in \mathcal{G}_{\mathcal{L}}$ . Since, by definition, queries return sets of tuples, the notions of set inclusion and equality can be applied.

In relation to containment and equivalence, we consider the following decision problems, again parameterized by a query language  $\mathcal{L}$ .

##### Problem: $\mathcal{L}$ -Containment

**Input:** Two queries  $q$  and  $q'$  from  $\mathcal{L}$   
**Output:** **true** if  $q \subseteq q'$ , and **false** otherwise

##### Problem: $\mathcal{L}$ -Equivalence

**Input:** Two queries  $q$  and  $q'$  from  $\mathcal{L}$   
**Output:** **true** if  $q \equiv q'$ , and **false** otherwise



## Part VIII

---

### Tree-Structured Data



Tree-structured data came to the attention of the data management community with the introduction of the Extensible Markup Language (XML) in 1998, which has since then become one of the most widespread formats for exchanging data on the Web. The other popular data interchange format, the JavaScript Data Interchange Standard (JSON), was defined in 1999 and is similar to XML in the sense that it also treats data in a tree structured manner. Tree-structured data is therefore an important part of data management and we devote this part to discuss some of its fundamental aspects.

Part [VIII](#) is organized as follows. In Chapter [56](#), we discuss the data models that we will use throughout this part. [TODO add explaining sentence.] We then define *tree pattern queries* in Chapter [61](#). Tree pattern queries are a fragment of the language XPath, and can navigate through trees using the *child* and *descendant* relations. As such, they are among the most fundamental languages one can devise for querying trees. They arguably play a similar role for tree-structured data as conjunctive queries do for relational data. In Chapters [62](#) and [63](#) we treat the basic static analysis and optimization problems for tree pattern queries, namely *containment* and *minimization*.

Chapters [61](#) and [63](#) are presented in terms of unordered trees, that is, trees in which siblings are unordered. However, since tree pattern queries are agnostic of the sibling ordering in the data, all the results in Chapters [61](#) and [63](#) remain to hold in the setting where sibling ordering is present in the data.



## Data Model

In this chapter we introduce several abstractions of tree-structured data. We focus on abstractions that are simple and elegant, yet powerful enough to prove results that help us understand the nature of querying and reasoning about tree-structured data in practice. Before we dive into our mathematical models, we take a look at how tree-structured data is represented in practice.

### Tree-Structured Data in Practice

The most widespread formats for storing tree-structured data are XML (eXtensible Markup Language) and JSON (JavaScript Object Notation). We briefly discuss the nuts and bolts of these formats by example.

#### XML

Consider Figure 56.1(a), which describes a simple XML document that contains IDs, names, and birthplaces of three persons. The first line states that the document uses version 1.0 of the XML standard and is encoded in UTF-8. The actual data starts in the second line and is hierarchically structured, much like HTML.<sup>1</sup> In XML, the tags such as `persons` and `person` are called *element names*, whereas `pers_id`, `name`, and `birthplace` are *attribute names*. In the document, the top-level element `person` contains three elements, labeled `person`, each of which has attributes with names `pers_id`, `name`, and `birthplace`. The values of the respective attributes are given as strings inside double quotes. For instance, the birthplace of Jimi is Seattle.

Figure 56.1(a), however, does not significantly use the capabilities of XML to structure data as a tree. So let us consider Figure 56.1(b), which contains similar information than Figure 56.1(a), but with more detail. Compared to Figure 56.1(a), we made two significant changes. First, instead of writing the

---

<sup>1</sup> For brevity, we assume that the reader has basic familiarity with HTML, but at the same time encourage her/him to continue reading even if this is not the case.



```
<?xml version="1.0" encoding="UTF-8"?>
<persons>
  <person pers_id="1" name="Jimi" birthplace="Seattle"/>
  <person pers_id="2" name="Saul" birthplace="Stoke-on-Trent"/>
  <person pers_id="3" name="Mark" birthplace="Glasgow"/>
</persons>
```

(a)

```
<?xml version="1.0" encoding="UTF-8"?>
<persons>
  <person pers_id="1">
    <name> Jimi </name>
    <birthplace>
      <city> Seattle </city>
      <state> Washington </state>
      <country> United States </country>
    </birthplace>
  </person>
  <person pers_id="2">
    <name> Saul </name>
    <birthplace>
      <city> Stoke-on-Trent </city>
      <country> United Kingdom </country>
    </birthplace>
  </person>
  <person pers_id="3">
    <name> Mark </name>
    <birthplace>
      <city> Glasgow </city>
      <country> United Kingdom </country>
    </birthplace>
  </person>
</persons>
```

(b)

Fig. 56.1: Two XML documents

names of persons as XML attributes, we wrote them inside **name**-tags which are nested inside **person**-tags. (Notice that the values of names are no longer enclosed in double quotes.) Second, we did the same with birthplaces, but additionally provided more detail in terms of the state and country of the city of birth. Furthermore, we did not provide state information for every city of birth.

We now discuss some principled differences between the two documents. First of all, the ordering between attributes in XML is irrelevant. That is, the element

```
<person pers_id="1" name="Jimi" birthplace="Seattle"/>
```

conveys the same information than

```
<person name="Jimi" birthplace="Seattle" pers_id="1"/> .
```

Intuitively, XML considers both fragments as encoding an element with name **person** for which the attribute with name **pers\_id** has value “1”, **name** has value “Jimi”, and **birthplace** has value “Seattle”.

Ordering between elements, however is important. That is, the fragment

```
<person> <pers_id> 3 </pers_id> <name> Mark </name> </person>
```

is considered to be different from

```
<person> <name> Mark </name> <pers_id> 3 </pers_id> </person>
```

in XML. The reason is that XML treats the hierarchical nesting of elements as a tree where, in the first fragment, the *first child* of **person** is a **pers\_id**-element, whereas in the second fragment, the *first child* of **person** is a **name**-element, which is different.

### JSON

Let us now consider how JSON represents tree-structured data. In JSON, a natural way to represent the information similar to that in Figure 56.1(a) would be the document in Figure 56.2(a). Another way to represent this information would be in Figure 56.2(b), which is closer to our XML representation in Figure 56.1(a).

Before we dive into detail on these examples of JSON documents, let us give a quick overview of how JSON documents are constructed, focusing on aspects that are relevant to this book. *JSON documents* can be inductively defined as follows. The simplest JSON documents are *string values*, which are strings enclosed in double quotes, such as “**persons**”, “**person**”, “**pers\_id**”, “1”, etc. If  $k_1, \dots, k_n$  are pairwise distinct string values and  $v_1, \dots, v_n$  are JSON documents, then

$$\{ k_1 : v_1, \dots, k_n : v_n \}$$

is a JSON document, called a *JSON object*. For each  $i \in [n]$ , we refer to  $k_i : v_i$  as a *key-value* pair of the object. Furthermore,

$$[ v_1, \dots, v_n ]$$

is a JSON document, called a *JSON array*. It is important to note that, inside JSON objects and arrays, the  $v_i$  can again be objects and arrays, which gives rise to the hierarchical or tree-structured nature of JSON documents.

Objects and arrays are used in JSON to represent ordered and unordered information, respectively. That is,

```
{"pers_id":"3", "name":"Mark", "birthplace":"Glasgow"}
```

and

```
{"name":"Mark", "birthplace":"Glasgow", "pers_id":"3"}
```

represent the same information, but

```
["pers_id":"3", "name":"Mark", "birthplace":"Glasgow"]
```

and

```
["name":"Mark", "birthplace":"Glasgow", "pers_id":"3"]
```

do not, since the order inside arrays is important.

Let us now discuss the two JSON documents in Figure 56.2. We can now see that the document in Figure 56.2(b) indeed represents a similar structure than the XML document in Figure 56.1(a). We have a **persons** entity, which contains an ordered list of three **person** entities. Each of these contains a JSON object that describes their “attribute names” **pers\_id**, **name**, and **birthplace**, and their respective values. The data in Figure 56.1(a) could be described in exactly the same way. Notice that, in this example, the requirement that keys inside JSON objects are pairwise disjoint corresponds to the requirement that attribute names in XML elements are pairwise disjoint.

The JSON document in Figure 56.2(a) represents the information a bit differently, since it omits the **person**-keys and immediately packs the information on the three persons in an array that contains three objects.

Finally, we show in Figure 56.3 how tree-structured information as in Figure 56.1b can be modeled in JSON. Notice that we needed to add a string **data** to give a name to the array that contains the ordered name and birthplace information of persons. The ordering between the person ID and this array is irrelevant.

## Labeled Unordered Trees

We now turn to mathematical abstractions of tree-structured data. The simplest such abstraction is the one of *labeled unordered trees*, which is based on node-labeled directed graphs.

```
{
  "persons": [
    {
      "pers_id": "1",
      "name": "Jimi",
      "birthplace": "Seattle"
    },
    {
      "pers_id": "2",
      "name": "Saul",
      "birthplace": "Stoke-on-Trent"
    },
    {
      "pers_id": "3",
      "name": "Mark",
      "birthplace": "Glasgow"
    }
  ]
}
```

(a)

```
{
  "persons": [
    {
      "person": {
        "pers_id": "1",
        "name": "Jimi",
        "birthplace": "Seattle"
      }
    },
    {
      "person": {
        "pers_id": "2",
        "name": "Saul",
        "birthplace": "Stoke-on-Trent"
      }
    },
    {
      "person": {
        "pers_id": "3",
        "name": "Mark",
        "birthplace": "Glasgow"
      }
    }
  ]
}
```

(b)

Fig. 56.2: Two similar JSON documents

**Definition 56.1: Labeled Unordered Tree**

A connected, node-labeled, directed graph  $T = (V, E, \text{lab})$  is a *labeled unordered tree* if,

- for every node  $v$ , there is at most one node  $u$  with  $(u, v) \in E$  and
- there is exactly one node  $v$  (called the *root* of  $T$ ) without an incoming edge  $(u, v)$ .

If  $T$  is a labeled unordered tree, we denote its root by  $\text{Root}(T)$ . Furthermore, if  $(u, v) \in E$  then we call  $v$  a *child* of  $u$  and  $u$  the *parent* of  $v$ . Likewise, we refer to  $E$  as the set of *child edges* of  $T$ . We call  $v$  a *descendant* of  $u$  if there exists a non-empty path from  $u$  to  $v$  in  $T$ . In this case, we also call  $u$  an *ancestor* of  $v$ . A *leaf* is a node without a child, that is, a node  $u$  such that  $\{v \mid (u, v) \in E\}$  is empty. Two nodes are *siblings* if they have the same parent. We extend this terminology to sets of nodes and say that  $S$  is a set of siblings if all pairs of elements from  $S$  are siblings. A set of siblings  $S$  is *maximal* if there does not exist a node  $u \in V - S$  such that  $S \cup \{u\}$  is a set of siblings. By  $T|_u$  we denote the subtree of  $T$  rooted at node  $u$ , formally defined as  $(V|_u, E|_u, \text{lab}|_u)$ , where  $V|_u$  consists of  $u$  and all its descendants in  $T$ , the relation  $E|_u$  is  $E \cap V|_u \times V|_u$ , and  $\text{lab}|_u$  is the restriction of  $\text{lab}$  to  $V|_u$ .

We denote the empty labeled unordered tree by  $\varepsilon$ . For  $a \in \mathbf{Lab}$ , we use the notation  $a(T_1, \dots, T_n)$  to denote the labeled unordered tree in which the root has label  $a$  and has  $n$  children  $u_1, \dots, u_n$ . Furthermore,  $T_i$  is its subtree at node  $u_i$  for each  $i \in [n]$ . That is, for each  $i \in [n]$ , we have that  $T_i$  is a labeled unordered tree and  $a(T_1, \dots, T_n)|_{u_i} = T_i$ . Notice that this definition implies that  $T_i$  is never empty.

```

{"persons": [
  {"person":
    {"pers_id": "1",
     "data": [
       {"name": "Jimi"},
       {"birthplace": [
         {"city": "Seattle" },
         {"state": "Washington"},
         {"country": "United States"}]}
     ]}
  },
  {"person":
    {"pers_id": "2",
     "data": [
       {"name": "Saul"},
       {"birthplace": [
         {"city": "Stoke-on-Trent"},
         {"country": "United Kingdom"}]}
     ]}
  },
  {"person":
    {"pers_id": "3",
     "data": [
       {"name": "Mark"},
       {"birthplace": [
         {"city": "Glasgow"},
         {"country": "United Kingdom"}]}
     ]}
  }
]}

```

Fig. 56.3: A JSON document that describes ordered and unordered information in a similar way as the XML document in Figure 56.1b

The *degree* of a node is the number of its children. The degree of a tree is the maximum degree of any of its nodes. The *depth* of the empty tree  $\varepsilon$  is zero, the depth of a tree  $T = a$  is one, and the depth of  $T = a(T_1, \dots, T_n)$  is one plus the maximal depth of  $T_1, \dots, T_n$ .

We depict trees in the usual way as in Figure 56.4(b), with the root node on top and the leaves at the bottom. Since all edges are directed downward (or away from the root) in such pictures, we omit arrows. The trees we defined until now are *unordered*, which means that we consider the children of a node as an unordered set. Therefore, the trees in Figures 56.4(b) and 56.4(c)

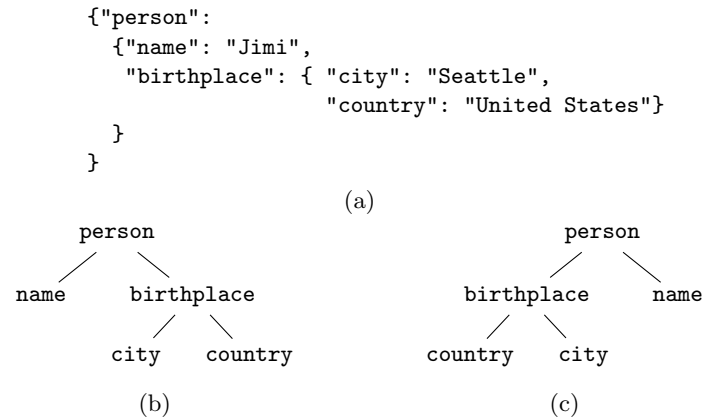


Fig. 56.4: A JSON document and two illustrations of the unordered tree that captures the nesting of its keys

formally represent the same unordered tree, which corresponds to the nesting structure of the keys in the JSON document in Figure 56.4(a).

Recall, however, that there is a fundamental difference between the nesting of keys in JSON objects and labeled unordered trees. Labeled unordered trees can have siblings with the same label, whereas key values on the same level of a JSON object need to be unique. We choose not to impose this restriction on labeled unordered trees, because

- (1) this allows us to present results that are relevant to both JSON and XML in Chapters 61–63 and
- (2) the results in Chapters 61–63 are not affected by this choice.

## Labeled Ordered Trees

In many situations, the ordering between siblings is important. For instance, if we want to describe an address, we may want to start with a street name and number, followed by a city, followed by a country. We will model such behavior by labeled trees in which the siblings are ordered. To this end, let  $S$  be a finite set. We say that a binary relation  $R$  is a *successor relation on  $S$*  if the restriction of  $R$  to  $S$ , that is,  $\{(a, b) \in R \mid a \in S, b \in S\}$  is isomorphic to  $\{(1, 2), (2, 3), \dots, (|S| - 1, |S|)\}$ .

### Definition 56.2: Labeled Ordered Tree

A *labeled ordered tree* is a tuple

$$T = (V, E^c, E^{\text{ns}}, \text{lab}) ,$$

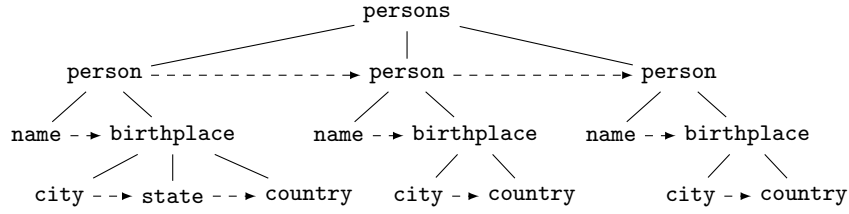


Fig. 56.5: A labeled ordered tree, capturing the nesting of the element names of the XML document in Figure 56.1(b).

where  $(V, E^c, \text{lab})$  is a labeled unordered tree and  $E^{\text{ns}}$  is a successor relation on every maximal set of siblings of  $(V, E^c, \text{lab})$ .

Notice that a labeled ordered tree is a two-sorted node-labeled directed graph according to Definition 55.1. We will use all the terminology that we introduced for labeled unordered trees (such as *root*, *child*, *degree*, etc.) in the context of labeled ordered trees. We refer to  $E^{\text{ns}}$  as the *next sibling* relation. The *first child* (resp., *last child*) of  $u$  in  $T$  is only defined if  $u$  has at least one child and is the smallest (resp., largest) element among all the children of  $u$  according to  $E^{\text{ns}}$ . We say that  $v$  is a *following sibling* of  $u$  if there exists a non-empty path from  $u$  to  $v$  in  $T$  whose edges are all in  $E^{\text{ns}}$ . Whenever  $T$  is clear from the context, we may omit “in  $T$ ” when we discuss nodes and the relationships between them. A set of labeled ordered trees is also called a *tree language*.

Similarly to unordered labeled trees, we denote by  $T|_u$  the subtree of  $T$  rooted at node  $u$ , formally defined as  $(V|_u, E|_u^c, E|_u^{\text{ns}}, \text{lab}|_u)$ , where  $V|_u$  consists of  $u$  and all its descendants,  $\text{lab}|_u$  is the restriction of  $\text{lab}$  to  $V|_u$  and the relations  $E|_u^c$  and  $E|_u^{\text{ns}}$  are  $E^c \cap V|_u \times V|_u$ , and  $E^{\text{ns}} \cap V|_u \times V|_u$ , respectively.

We denote the empty labeled ordered tree by  $\varepsilon$ . For  $a \in \text{Lab}$ , we use the notation  $a(T_1, \dots, T_n)$  to denote the labeled ordered tree in which the root has label  $a$  and has  $n$  children  $u_1, \dots, u_n$ . Furthermore,  $T_i$  is its subtree at node  $u_i$  for each  $i \in [n]$ . That is, for each  $i \in [n]$ , we have that  $T_i$  is a labeled ordered tree and  $a(T_1, \dots, T_n)|_{u_i} = T_i$ .

Figure 56.5 depicts the ordered tree that captures the nested structure of the element names in the XML document of Figure 56.1(b). The edges in  $E^c$  are solid, whereas the edges in  $E^{\text{ns}}$  are dashed. As in Figure 56.4(c), we omit the arrows on the edges in  $E^c$ , since they all point downward. Furthermore, notice that all the edges in  $E^{\text{ns}}$  can be inferred from the figure, once we know that the tree is supposed to be a labeled ordered tree. (They just go from left to right among the children of nodes.) For this reason, we will therefore also omit the arrows in  $E^{\text{ns}}$  in illustrations of labeled ordered trees. To avoid confusion, we will always make clear in the context whether a figure illustrates a labeled unordered tree or a labeled ordered tree.

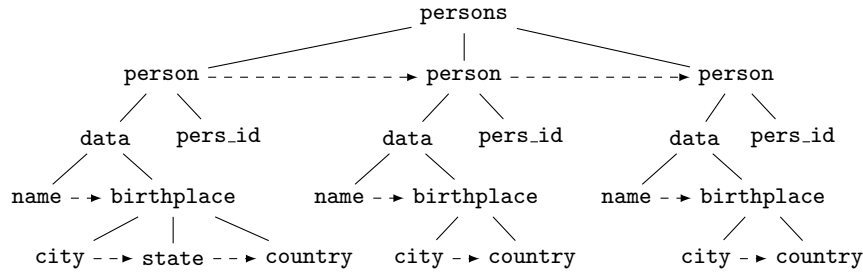


Fig. 56.6: A two-sorted node-labeled graph, capturing the nesting and ordering of the element keys of the JSON document in Figure 56.3

## Theory versus Practice

We have now seen two fundamental theoretical abstractions of tree-structured data, namely *labeled unordered trees* and *labeled ordered trees*, together with a few examples of how tree-structured data can be represented in practice, using JSON and XML. As the careful reader undoubtedly realizes, our two abstractions are not even powerful enough to precisely capture our practical examples. For instance, for representing the nesting of keys in the JSON document in Figure 56.3, we would arguably require a structure as in Figure 56.6, which is neither a labeled ordered tree, nor a labeled unordered tree.

So why do we work with labeled ordered or unordered trees? This is not because we are lazy or uninterested in the gritty details of reality. Our aim is rather to work with elegant and simple abstractions that still allow us to study practically relevant issues. In fact, simple abstractions are often very helpful to distill practical challenges to their mathematical core. They enable us to discover the mathematical questions whose solution would imply a solution to the practical challenge. Good abstractions help us to focus our attention on the fundamental underlying issues.

In this spirit, all results in Part VIII are presented in terms of labeled ordered or unordered trees. As such, they are immediately applicable to those parts of XML or JSON that can be understood as labeled ordered or unordered trees, respectively. In many cases, the results can even be trivially generalized to abstractions that are closer to reality, e.g., abstractions that consider a mixture of ordered and unordered siblings as in Figure 56.6.

We finally argue that there is no “best” way to abstract tree-structured data as a tree. For example, the JSON document in Figure 56.4(a) can also be represented as the labeled unordered tree in Figure 56.7, in which we treat both the JSON keys and the *data values* ‘Jimi’, ‘Seattle’, and ‘United States’ as node labels. This aspect is important for understanding many results in Part VIII. Although these results are about labeled ordered/unordered trees, this does not mean that the results are only about the nesting structure of element names in XML document or keys in JSON documents and therefore



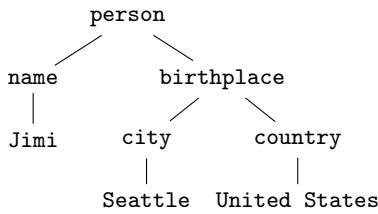


Fig. 56.7: Another way of seeing the JSON document in Figure 56.4(a) as a labeled unordered tree

do not concern the actual data. In fact, when it comes to interpreting the results in Part VIII for managing XML or JSON data, it is more instructive to consider the interaction between *data* and *queries*. In Part VIII, we will consider query languages that allow to compare labels to fixed values that we specify in the query. An example of such a query is

*Select all nodes with label 'Seattle',*

which compares node labels to the fixed value 'Seattle'. Therefore, such results also generalize to query languages that can compare data values to constants. Later in the book, we will consider query languages that can compare data values to each other, which is fundamentally different. An example of a query that compares values to each other is

*Select all node pairs  $(x, y)$  such that  $x$  and  $y$  have different labels.*

## First-Order Logic over Trees

In this chapter we introduce first-order logic over labeled ordered trees, to which we refer as “first order logic over trees” for the sake of brevity. This logic is similar to first-order logic (Chapter 3), but is tuned towards expressing properties over labeled ordered trees. Just like its relational counterpart, first-order logic over trees is an important yardstick for the expressiveness of query languages over trees.

Formulae in first-order logic over trees will use relation names from the set

$$\{\text{Lab}_a \mid a \in \mathbf{Lab}\} \cup \{E^{\text{fc}}, E^{\text{ns}}, E^{\text{d}}, E^{\text{fs}}\}.$$

Here,  $\text{Lab}_a$  is called the *a-label relation*,  $E^{\text{fc}}$  is called the *first-child relation*,  $E^{\text{ns}}$  the *next-sibling relation*,  $E^{\text{d}}$  the *descendant relation*, and  $E^{\text{fs}}$  the *following-sibling relation*, respectively.<sup>1</sup>

We will follow the Background chapter (Chapter 2) in the sense that we assume that we have a countably infinite set  $\mathbf{Var}$  of *variables*, disjoint from  $\mathbf{Nodes}$  and  $\mathbf{Lab}$ . As in Chapter 3 variables will be typically denoted by  $x, y, z, \dots$  (possibly with subscripts and superscripts). Differently from Chapter 3, the symbols  $a, b, c, \dots$  are now used to denote *labels* from  $\mathbf{Lab}$ . Formulae of first-order logic are inductively defined using terms, conjunction ( $\wedge$ ), disjunction ( $\vee$ ), negation ( $\neg$ ), existential quantification ( $\exists$ ), and universal quantification ( $\forall$ ).

### Definition 57.1: Syntax of First-Order Logic over Trees

We define *formulae of first-order logic* (FO) over trees as follows:

- If  $x$  is a variable, then  $\text{Lab}_a(x)$ , for every  $a \in \mathbf{Lab}$  is an atomic formula.
- If  $x$  and  $y$  are variables from  $\mathbf{Var}$ , then  $x = y$  is an atomic formula.

<sup>1</sup> The reader may wonder why we don't use the child relation or the last-child relation, but we will see that these can be expressed using the available relations.

- If  $x$  and  $y$  are variables from  $\text{Var}$ , then  $E^{\text{fc}}(x, y)$ ,  $E^{\text{ns}}(x, y)$ ,  $E^{\text{d}}(x, y)$ ,  $E^{\text{fs}}(x, y)$ , and are atomic formulae.
- If  $\varphi_1$  and  $\varphi_2$  are formulae, then  $(\varphi_1 \wedge \varphi_2)$ ,  $(\varphi_1 \vee \varphi_2)$ , and  $(\neg \varphi_1)$  are formulae.
- If  $\varphi$  is a formula and  $x \in \text{Var}$ , then  $(\exists x \varphi)$  and  $(\forall x \varphi)$  are formulae.

Analogously to the first-order formulae for relational databases in Chapter 3, the *size*  $\|\varphi\|$  of  $\varphi$  is defined to be the total number of labels, variables, and symbols from  $\{\wedge, \vee, \neg, =, \exists, \forall\}$  occurring in  $\varphi$ . The *free variables*  $\text{FV}(\varphi)$  of  $\varphi$  are defined analogously as in first-order logic over relational databases. We also omit brackets in the same manner as for first order formulae for relational databases to avoid notional clutter.

### Example 57.2: First-Order Formulae

Consider the following FO-formula over trees:

$$\text{Lab}_{\text{'person'}}(x) \wedge \text{Lab}_{\text{'name'}}(y) \wedge E^{\text{fc}}(x, y) . \quad (57.1)$$

The free variables of this formula are  $x$  and  $y$ . The three formulae

$$\begin{aligned} &\text{Lab}_{\text{'person'}}(x) \wedge \neg \exists y \left( \text{Lab}_{\text{'state'}}(y) \wedge (\exists z \exists z_1 \exists z_2 \text{Lab}_{\text{'birthplace'}}(z) \right. \\ &\quad \left. \wedge E^{\text{fc}}(x, z_1) \wedge E^{\text{ns}}(z_1, z) \wedge E^{\text{fc}}(z, z_2) \wedge E^{\text{ns}}(z_2, y)) \right) , \end{aligned} \quad (57.2)$$

$$\neg \exists y E^{\text{ns}}(x, y) , \quad (57.3)$$

and

$$\neg \exists y (E^{\text{fc}}(y, x) \vee E^{\text{ns}}(y, x)) \quad (57.4)$$

have only one free variable, namely  $x$ . The formula

$$E^{\text{fc}}(x, y) \vee (\exists z E^{\text{fc}}(x, z) \wedge E^{\text{fs}}(z, y)) \quad (57.5)$$

has two free variables, namely  $x$  and  $y$ .

Given a labeled ordered tree  $T = (V, E^c, E^{\text{ns}}, \text{lab})$ , we define the notion of satisfaction of a formula  $\varphi$  with respect to an *assignment*  $\eta$  for  $\varphi$  over  $T$ . Such an assignment is a function from  $\text{FV}(\varphi)$  to  $V$ . If  $x$  is a variable and  $v \in V$  a node of  $T$ , we use  $\eta[x/v]$  to denote the assignment that modifies  $\eta$  by setting  $\eta(x) = v$ . We are now ready for defining the semantics of first-order formulae over trees.

**Definition 57.3: Semantics of First-Order Logic over Trees**

Given a tree  $T = (V, E^c, E^{ns}, \text{lab})$ , a formula  $\varphi$ , and an assignment  $\eta$  for  $\varphi$  over  $T$ , we inductively define when  $\varphi$  is *satisfied* in  $T$  under  $\eta$ , written  $(T, \eta) \models \varphi$ , as follows:

- If  $\varphi$  is  $\text{Lab}_a(x)$ , then  $(T, \eta) \models \varphi$  if  $\text{lab}(\eta(x)) = a$ .
- If  $\varphi$  is  $x = y$ , then  $(T, \eta) \models \varphi$  if  $\eta(x) = \eta(y)$ .
- If  $\varphi$  is  $E^{\text{fc}}(x, y)$ , then  $(T, \eta) \models \varphi$  if  $\eta(y)$  is the first child of  $\eta(x)$  in  $T$ .
- If  $\varphi$  is  $E^{ns}(x, y)$ , then  $(T, \eta) \models \varphi$  if  $\eta(y)$  is a next sibling of  $\eta(x)$  in  $T$ .
- If  $\varphi$  is  $E^d(x, y)$ , then  $(T, \eta) \models \varphi$  if  $\eta(y)$  is a descendant of  $\eta(x)$  in  $T$ .
- If  $\varphi$  is  $E^{\text{fs}}(x, y)$ , then  $(T, \eta) \models \varphi$  if  $\eta(y)$  is a following sibling of  $\eta(x)$  in  $T$ .
- If  $\varphi = \varphi_1 \wedge \varphi_2$ , then  $(T, \eta) \models \varphi$  if  $(T, \eta) \models \varphi_1$  and  $(T, \eta) \models \varphi_2$ .
- If  $\varphi = \varphi_1 \vee \varphi_2$ , then  $(T, \eta) \models \varphi$  if  $(T, \eta) \models \varphi_1$  or  $(T, \eta) \models \varphi_2$ .
- If  $\varphi = \neg\psi$ , then  $(T, \eta) \models \varphi$  if  $(T, \eta) \models \psi$  does not hold.
- If  $\varphi = \exists x \psi$ , then  $(T, \eta) \models \varphi$  if  $(T, \eta[x/v]) \models \psi$  for *some* node  $v \in V$ .
- If  $\varphi = \forall x \psi$ , then  $(T, \eta) \models \varphi$  if  $(T, \eta[x/v]) \models \psi$  for *every* node  $v \in V$ .

Formulae of FO over trees are therefore very similar to formulae of FO over relations (Chapter 3) “over  $\{\text{Lab}_a \mid a \in \text{Lab}\} \cup \{E^{\text{fc}}, E^{ns}, E^d, E^{\text{fs}}\}$ ”, with the following differences:

- (1) Formally, assignments for formulae of FO over relations map variables to **Const**, whereas they map variables to **Nodes** for formulae of FO over trees. This means that formulae in FO over relations are interpreted over **Const**, whereas formulae in FO over trees are interpreted over **Nodes**.
- (2) Formulae in FO over relations have atomic formulae of the form “ $x = a$ ”, which allow equality tests between variables and constants. We do not have a similar construct here (which would be an equality test between a variable and a node).
- (3) Whereas a schema **S** in Chapter 3 is always finite, the set  $\{\text{Lab}_a \mid a \in \text{Lab}\} \cup \{E^{\text{fc}}, E^{ns}, E^d, E^{\text{fs}}\}$  is infinite, because **Lab** is infinite. This is only a minor difference in the sense that it is similar to “omitting” the schema in the relational world. (Each FO formula “over Rel” is an FO formula over some schema **S**.)

**Example 57.4: Semantics of First-Order Formulae**

We give an intuitive description of the semantic meaning of the formulae in Example 57.2:

- Formula (57.1) is satisfied by all nodes  $x$  and  $y$  such that  $x$  has label ‘person’,  $y$  has label ‘name’, and  $y$  is the first child of  $x$ .
- Formula (57.2) is satisfied by all nodes  $x$  with label ‘person’, such that there does not exist a node  $y$  which is labeled ‘state’, and  $y$  is the second child of a node labeled ‘birthplace’, which is in turn the second child of  $x$ . In Figure 56.5, such nodes  $x$  would be the second and third children of the root, i.e., nodes that do not have a ‘state’ child of their ‘birthplace’ child.
- Formula (57.3) is satisfied by all nodes  $x$  that do not have next siblings. That is,  $x$  is a last sibling (or the root).
- Formula (57.4) is satisfied by all nodes  $x$  such that there does not exist a node  $y$  from which  $x$  is the first child or the next sibling. That is,  $x$  is the root of the tree.
- Formula (57.5) is satisfied by all nodes  $x$  and  $y$  such that  $y$  is a child of  $x$ .

**First-Order Queries over Trees**

Recall that a  $k$ -ary data graph query is a function  $q$  that takes a data graph  $G$  as input and produces a set  $q(G) \subseteq (\text{Nodes} \cup \text{Const})^k$ . First-order queries over trees are a special case of such queries, since their inputs are labeled ordered trees and they produce a set of  $k$ -tuples of nodes, that is, a set in  $\text{Nodes}^k$ . We first define the syntax of first-order queries over trees.

**Definition 57.5: First-Order Queries over Trees**

A *first-order query over trees* is an expression of the form  $\varphi(\bar{x})$ , where  $\varphi$  is an FO formula over trees, and  $\bar{x}$  is a tuple of free variables of  $\varphi$  such that each free variable of  $\varphi$  occurs in  $\bar{x}$  at least once.

We define *size*  $\|\varphi(\bar{x})\|$  of a first-order query  $\varphi(\bar{x})$  as  $\|\varphi\| + \|\bar{x}\|$ .

We now define the semantics of first-order queries over trees. Let  $\varphi(\bar{x})$  be an FO query over trees. Given a labeled ordered tree  $T = (V, E^c, E^{\text{ns}}, \text{lab})$  and a tuple  $\bar{a}$  of elements from  $\text{Nodes}$ , we say that  $T$  *satisfies the query*  $\varphi(\bar{x})$  *using the values*  $\bar{a}$ , denoted by  $T \models \varphi(\bar{a})$ , if there exists an assignment  $\eta$  for  $\varphi$  over  $T$  such that  $\eta(\bar{x}) = \bar{a}$  and  $(T, \eta) \models \varphi$ . We can now define the data model and the output of FO queries over trees. The data model defines the possible

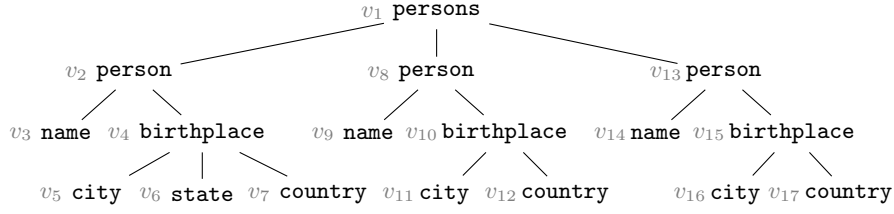


Fig. 57.1: A labeled ordered tree

inputs of the queries. For each data graph in the data model, we define the output of the query.

#### Definition 57.6: Evaluation of First-Order Queries over Trees

Given a labeled ordered tree  $T = (V, E^c, E^{ns}, \text{lab})$  and an FO query over trees  $q = \varphi(x_1, \dots, x_k)$ , where  $k \geq 0$ , the *output* of  $q$  on  $T$  is defined as the set of tuples

$$q(T) = \{\bar{a} \in \text{Nodes}^k \mid T \models \varphi(\bar{a})\}.$$

The *data model* of FO queries over trees is the set of labeled ordered trees.

We conclude the chapter with a few examples of first-order queries over trees.

#### Example 57.7: First-Order Queries over Trees

Let  $T$  be the tree in Figure 57.1, where the nodes are depicted in grey.

- Let  $q_1$  be the query  $\varphi_1(x, y)$ , where  $\varphi_1$  is the formula (57.1). Then

$$q_1(T) = \{(v_1, v_3), (v_8, v_9), (v_{13}, v_{14})\}.$$

- Let  $q_2$  be the query  $\varphi_2(x)$ , where  $\varphi_2$  is the formula (57.2). Then

$$q_2(T) = \{v_8, v_9\}.$$

- Let  $q_3$  be the query  $\varphi_3(x)$ , where  $\varphi_3$  is the formula (57.3). Then

$$q_3(T) = \{v_1, v_4, v_7, v_{10}, v_{12}, v_{13}, v_{15}, v_{17}\}.$$

- Let  $q_4$  be the query  $\varphi_4(x)$ , where  $\varphi_4$  is the formula (57.4). Then

$$q_4(T) = \{v_1\}.$$

- Let  $q_5$  be the query  $\varphi_5(y, x)$ , where  $\varphi_5$  is the formula (57.5). Then

$$q_5(T) = \{(v_2, v_1), (v_3, v_2), (v_4, v_2), (v_5, v_4), (v_6, v_4), (v_7, v_4), \\ (v_8, v_1), (v_9, v_8), (v_{10}, v_8), (v_{11}, v_{10}), \\ (v_{12}, v_{10}), (v_{13}, v_1), (v_{14}, v_{13}), (v_{15}, v_{13}), (v_{16}, v_{15}), (v_{17}, v_{15})\}$$

## XPath

XPath is a powerful language that is designed for navigation and node-selection in tree-structured data. In this chapter, we offer a principled view on XPath, which means that we describe and formally define an elegant and mathematically clean language for navigating through labeled ordered trees. We will point out differences with the official XPath standard whenever relevant.

### XPath Axes

XPath uses so-called *axes* as primitive operations for navigating in trees. We use the following axes in this chapter:

<code>self</code>	<code>descendant</code>	<code>next-sibling</code>	<code>following</code>
<code>child</code>	<code>descendant-or-self</code>	<code>following-sibling</code>	<code>preceding</code>
<code>parent</code>	<code>ancestor</code>	<code>previous-sibling</code>	
	<code>ancestor-or-self</code>	<code>preceding-sibling</code>	

We note that the XPath standard does not consider the axes `next-sibling` and `previous-sibling`, but we consider them here for completeness. We now explain the meaning of each axis  $x$  by associating it to a binary relation  $R_x$  over **Nodes**. To this end, if  $R$  and  $R'$  are two binary relations, then

- we denote by  $R \circ R'$  the *composition* of  $R$  and  $R'$ , which is defined as  $\{(a, c) \mid (a, b) \in R, (b, c) \in R'\}$ ,
- we denote by  $R^+$  the *transitive closure* of  $R$ , which is defined as  $\{(a_1, a_n) \mid \exists a_2, \dots, a_{n-1} \text{ such that } (a_i, a_{i+1}) \in R \text{ for every } i \in [n-1]\}$ ,
- we denote by  $R^*$  the *transitive reflexive closure* of  $R$ , which is defined as  $R^+ \cup \{(a, a) \mid \exists b \text{ such that } (a, b) \in R \text{ or } (b, a) \in R\}$ , and
- we denote by  $R^{-1}$  the *inverse* of  $R$ , which is defined as  $\{(b, a) \mid (a, b) \in R\}$ .



Now consider a labeled ordered tree  $T = (V, E^c, E^{ns}, \text{lab})$ . Then, we define

$$R_{\text{self}} = \{(u, u) \mid u \in V\}, \quad R_{\text{child}} = E^c, \quad \text{and} \quad R_{\text{parent}} = (R_{\text{child}})^{-1}.$$

Furthermore, we define

$$\begin{aligned} R_{\text{descendant}} &= (R_{\text{child}})^+, & R_{\text{descendant-or-self}} &= (R_{\text{child}})^*, \\ R_{\text{ancestor}} &= (R_{\text{parent}})^+, & R_{\text{ancestor-or-self}} &= (R_{\text{parent}})^*, \\ R_{\text{next-sibling}} &= E^{ns}, & R_{\text{following-sibling}} &= (R_{\text{next-sibling}})^+, \\ R_{\text{previous-sibling}} &= (R_{\text{next-sibling}})^{-1}, & R_{\text{preceding-sibling}} &= (R_{\text{previous-sibling}})^+, \\ R_{\text{following}} &= R_{\text{ancestor-or-self}} \circ R_{\text{following-sibling}} \circ R_{\text{descendant-or-self}}, \text{ and} \\ R_{\text{preceding}} &= (R_{\text{following}})^{-1}. \end{aligned}$$

## Core XPath and Conditional XPath

We are now ready to define the syntax of *Core XPath* and *Conditional XPath*, which are two languages that have received considerable attention in the research literature.

### Definition 58.1: Syntax of Core- and Conditional XPath

The syntax of Core XPath expressions is defined by the grammar

$$\begin{aligned} \text{pe} &::= \text{step} \mid (\text{pe}/\text{pe}) \mid (\text{pe} \cup \text{pe}) \\ \text{step} &::= \text{axis} \mid \text{step}[\text{ne}] \\ \text{ne} &::= ?\text{pe} \mid (\text{lab} = a) \mid (\text{lab} = *) \mid (\text{ne} \wedge \text{ne}) \mid (\text{ne} \vee \text{ne}) \mid \neg \text{ne} \end{aligned}$$

Here, **axis** stands for one of the aforementioned axes and  $a$  is a label from **Lab**. The non-terminal **pe** defines the syntax of Core XPath *path expressions* and **ne** defines the syntax of Core XPath *node expressions*.

The syntax of *Conditional XPath expressions* is obtained from the above definition by replacing the rule for **step** with

$$\text{step} ::= \text{axis} \mid \text{step}[\text{ne}] \mid (\text{step}[\text{ne}])^*$$

In the grammar for Conditional XPath expressions, the non-terminal **pe** defines the syntax of Conditional XPath *path expressions* and **ne** defines the syntax of Conditional XPath *node expressions*.

In the remainder of the chapter, we abbreviate Core XPath with CoreXPath and Conditional XPath with CondXPath. We now define the semantics of these languages.

**Definition 58.2: Semantics of Core- and Conditional XPath**

Let  $T = (V, E^c, E^{ns}, \text{lab})$  be a labeled ordered tree. We define the semantics of CoreXPath and CondXPath expressions on  $T$  using a mutual induction that defines the semantics of path expressions and node expressions simultaneously. For every path expression  $e_p$ , its semantics  $\llbracket e_p \rrbracket$  is a subset of  $V^2$  and, for every node expression  $e_n$ , its semantics  $\llbracket e_n \rrbracket$  is a subset of  $V$ . More precisely, we have

$$\begin{aligned} \llbracket \text{axis} \rrbracket_T &:= R_{\text{axis}} \\ \llbracket \text{step}[e_n] \rrbracket_T &:= \{(u, v) \in \llbracket \text{step} \rrbracket_T \mid v \in \llbracket e_n \rrbracket_T\} \\ \llbracket (\text{step}[e_n])^* \rrbracket_T &:= \llbracket \text{step}[e_n] \rrbracket_T^* \\ \llbracket (e_p/e'_p) \rrbracket_T &:= \llbracket e_p \rrbracket_T \circ \llbracket e'_p \rrbracket_T \\ \llbracket (e_p \cup e'_p) \rrbracket_T &:= \llbracket e_p \rrbracket_T \cup \llbracket e'_p \rrbracket_T \end{aligned}$$

and

$$\begin{aligned} \llbracket (\text{lab} = a) \rrbracket_T &:= \{v \mid \text{lab}(v) = a\} \\ \llbracket (\text{lab} = *) \rrbracket_T &:= V \\ \llbracket ?e_p \rrbracket_T &:= \{u \mid \exists v \text{ with } (u, v) \in \llbracket e_p \rrbracket_T\} \\ \llbracket (e_n \wedge e'_n) \rrbracket_T &:= \llbracket e_n \rrbracket_T \cap \llbracket e'_n \rrbracket_T \\ \llbracket (e_n \vee e'_n) \rrbracket_T &:= \llbracket e_n \rrbracket_T \cup \llbracket e'_n \rrbracket_T \\ \llbracket \neg e_n \rrbracket_T &:= V - \llbracket e_n \rrbracket_T \end{aligned}$$

We note that subexpressions of the form  $(\text{lab} = *)$ , i.e., wildcard tests, are syntactic sugar, since they can be written as  $((\text{lab} = a) \vee \neg(\text{lab} = a))$  for an arbitrary  $a \in \text{Lab}$ . Our main reasons for having  $(\text{lab} = *)$  in our definition are that the XPath standard has this construct and it allows us to have a clean correspondence between XPath and tree pattern queries in Chapter 61. On the other hand, since it is syntactic sugar, we do not need to consider the case  $(\text{lab} = *)$  in several proofs.

**Example 58.3**

Consider the labeled ordered tree in Figure 57.1. The CoreXPath path expression

$$(\text{self}[(\text{lab} = \text{person})]/\text{child}[(\text{lab} = \text{name})]) \quad (58.1)$$

selects pairs of nodes  $(x, y)$  where  $x$  is labeled ‘person’,  $y$  is labeled ‘name’, and  $y$  is a child of  $x$ .

The CoreXPath node expression

$$\begin{aligned} &?(\text{descendant}[\llbracket (\text{lab} = \text{person}) \wedge \neg ?\text{child}[(\text{lab} = \text{birthplace})] \rrbracket_T] / \\ &\quad \text{child}[(\text{lab} = \text{name})]) \end{aligned}$$

selects the nodes  $x$  that are labeled ‘person’ and have a child with label ‘name’, but do not have a child with label ‘birthplace’.

The CondXPath path expression

$$(\text{next-sibling}[(\text{lab} = \text{person}))]^* \quad (58.2)$$

selects node pairs  $(x, y)$  such that  $x$  and  $y$  are siblings with label ‘person’ and every sibling between them also has label ‘person’.

Notice that our definition of the semantics of node and path expressions allows us to interpret them as queries over the set  $\mathcal{T}$  of labeled ordered trees. Indeed, we can associate to each node expression  $e_n$  a query  $q_{e_n}$  and to each path expression  $e_p$  a query  $q_{e_p}$  such that, on each labeled ordered tree  $T$ ,

- the *output of  $q_{e_n}$  on  $T$*  is the set of nodes

$$q_{e_n}(T) = \llbracket e_n \rrbracket_T$$

and

- the *output of  $q_{e_p}$  on  $T$*  is the set of node pairs

$$q_{e_p}(T) = \llbracket e_p \rrbracket_T .$$

From now on, we will therefore also treat  $e_n$  and  $e_p$  as unary and binary queries over  $\mathcal{T}$ , respectively.

## Complexity of Evaluation

We show that unary CoreXPath and CondXPath queries can be evaluated in linear time combined complexity.

### Theorem 58.4

Let  $T$  be a labeled ordered tree and  $q$  be a unary CoreXPath or CondXPath query. Then we can compute  $q(T)$  in time  $O(\|q\|\|T\|)$ .

*Proof.* Let  $T = (V, E^c, E^{\text{ns}}, \text{lab})$  be a labeled ordered tree. We assume without loss of generality that  $V = [k]$ , where  $k = |V|$ . Since  $q$  is a unary query, we have that  $q = e_n$  for some node expression  $e_n$ . Since node and path expressions are defined by mutual induction, we will prove the following by induction on subexpressions  $e$  of  $q$ :

- if  $e = e_n$  for some node expression  $e_n$ , then the set  $\llbracket e_n \rrbracket_T \subseteq V$  can be computed in time  $O(\|e_n\|\|T\|)$  and

- (b) if  $S \subseteq V$  and  $e = e_p$  for some path expression  $e_p$ , then the set  $\llbracket e_p \rrbracket_{T,S} := \{u \mid \exists v \in S \text{ with } (u, v) \in \llbracket e_p \rrbracket_T\} \subseteq V$  can be computed in time  $O(\|e_p\| \|T\|)$ .

Notice that the theorem immediately follows from the first bullet. The algorithm recursively computes and stores the sets  $\llbracket e_n \rrbracket_T$  and  $\llbracket e_p \rrbracket_{T,S}$ , for a relevant set of nodes  $S \subseteq V$ , for every subexpression  $e_n$  and  $e_p$  of  $q$ . Each such set can be stored as an array of  $|V|$  bits.

Let  $e$  be a subexpression of  $q$ . We first consider case (a) where  $e$  is a node expression. The first base case is when  $e$  is of the form  $(\mathbf{lab} = a)$ . Since we can indeed compute the set of  $a$ -labeled nodes in  $T$  in time  $O(\|T\|)$ , the induction hypothesis holds. The second base case, where  $e$  is of the form  $(\mathbf{lab} = *)$  is analogous.

Moving to the induction, if  $e$  is one of  $(e_n \wedge e'_n)$  or  $(e_n \vee e'_n)$ , then  $\llbracket e_n \rrbracket_T$  can be computed by iterating over every node  $v$  of  $T$  and testing if  $v \in \llbracket e_n \rrbracket_T \cap \llbracket e'_n \rrbracket_T$ . The total time for doing this, including the computation of  $\llbracket e_n \rrbracket_T$  and  $\llbracket e'_n \rrbracket_T$  is  $2|V| + O(\|e_n\| \|T\|) + O(\|e'_n\| \|T\|)$ , which is in  $O(\|e\| \|T\|)$ . The proof where  $e = \neg e_n$  is analogous. The final case, where  $e = ?e_p$ , is immediate from the induction hypothesis (b), taking  $S = V$ .

We now consider case (b). To this end, let  $S \subseteq V$  and  $e$  be a path expression. If  $e = \mathbf{axis}$ , then  $\llbracket e \rrbracket_{T,S}$  can be computed for each possibility of  $\mathbf{axis}$  in time  $O(\|T\|)$ . E.g., when  $\mathbf{axis} = \mathbf{descendant}$ , then  $\llbracket e \rrbracket_{T,S}$  is the set of nodes  $\{u \mid \exists v \in S \text{ with } (u, v) \in R_{\mathbf{descendant}}\}$ , which are the ancestors of nodes in  $S$ . (We leave the other cases as an exercise.) If  $e = (e_p / e'_p)$ , then we first compute  $S' = \llbracket e'_p \rrbracket_{T,S}$  and then  $\llbracket e_p \rrbracket_{T,S'}$ . For the case  $e = \mathbf{step}[e_n]$ , we use the induction hypothesis to compute  $S' = \llbracket e_n \rrbracket_T$  and then  $\llbracket \mathbf{step} \rrbracket_{T,S'}$ . The final case  $e = (e_p \cup e'_p)$  is also immediate from the induction, by computing  $\llbracket e \rrbracket_{T,S}$  as  $\llbracket e_p \rrbracket_{T,S} \cup \llbracket e'_p \rrbracket_{T,S}$ . This concludes the proof for CoreXPath.

In the case of CondXPath, the only extra case we need to deal with is  $e = (\mathbf{step}[e_n])^*$ . According to the definition of CondXPath,  $e$  can either be of the form  $(\mathbf{axis}[e_n])^*$ ,  $(\mathbf{step}'[e'_n][e_n])^*$ , or  $((\mathbf{step}'[e'_n])^*[e_n])^*$ . We will show how  $\llbracket e \rrbracket_{T,S}$  can be computed in time  $O(\|q\| \|T\|)$  in all these cases.

In the first case, if  $\llbracket \mathbf{axis} \rrbracket_T$  is a transitive relation (e.g., for  $\mathbf{axis} = \mathbf{descendant}$ ), then  $e$  is equivalent to  $\mathbf{self} \cup \mathbf{axis}[e_n]$ , which we have already solved. If  $\mathbf{axis} = \mathbf{child}$ , then we need to compute all ancestors  $u$  of nodes in  $S$  such that all nodes on the path from  $u$  to  $S$  are in  $\llbracket e_n \rrbracket$ , which can obviously be done in the required time. The other cases where  $\llbracket \mathbf{axis} \rrbracket_T$  is not transitive are analogous.

In the second case,  $e$  is of the form  $(\mathbf{step}'[e'_n][e_n])^*$  which is equivalent to  $(\mathbf{step}'[e'_n \wedge e_n])^*$  and therefore reduces to the first or third case.

In the third case,  $e$  is equivalent to  $\mathbf{self} \cup (\mathbf{step}'[e'_n])^*[e_n]$ , which also reduces to cases that we already dealt with (union and  $\mathbf{step}[e_n]$ ). This concludes the proof.  $\square$



## Expressiveness of XPath

In this chapter we compare the expressiveness of CoreXPath and CondXPath queries to first-order queries over trees (Definition 57.5). In a nutshell, we will see that unary CoreXPath queries are equally expressive as the unary FO queries over trees that can be defined using formulae that use two variables. CondXPath on the other hand can express precisely the unary and binary FO queries over trees. In this chapter we present some of the main ideas behind these correspondences. We first focus on unary queries and discuss binary queries at the end of the chapter.

### Unary CoreXPath and FO<sup>2</sup>

For our first result, we consider the *two-variable fragment* of FO over trees, which is the set of all FO formulae over trees that use at most two variables (repetitions are allowed). We denote this subset of formulae as FO<sup>2</sup>. We define *FO<sup>2</sup> queries over trees* analogously to FO queries over trees. We illustrate such a query in the next example.

#### Example 59.1: FO<sup>2</sup> Query over Trees

Consider the query  $\varphi(x)$ , where  $\varphi$  is the formula

$$\exists y \left( E^{\text{ns}}(x, y) \wedge (\exists x E^{\text{ns}}(y, x) \wedge \exists y E^{\text{ns}}(x, y)) \right).$$

This query selects all nodes  $x$  that have at least three siblings to the right. Notice that the existence of the second and third sibling to the right is expressed by re-using the variables  $x$  and  $y$  respectively.

It turns out that the unary CoreXPath queries are precisely the same as the unary FO<sup>2</sup> queries over trees.

**Theorem 59.2**

Unary CoreXPath is equally expressive as unary  $\text{FO}^2$  over trees.

*Proof (Sketch).* For simplicity, we will only prove the result in *one dimension*, that is, we only consider the CoreXPath queries for which the set of axes is restricted to **self**, **next-sibling**, **previous-sibling**, **following-sibling**, and **preceding-sibling**. We will prove that these queries are equally expressive as those  $\text{FO}^2$  queries in which formulae use the predicates  $\text{Lab}_a(x)$ ,  $E^{\text{ns}}(x_1, x_2)$ , and  $E^{\text{fs}}(x_1, x_2)$ . The general result can be proved using the same technique.

For the translation from CoreXPath to  $\text{FO}^2$ , we only need to show that CoreXPath queries defined by node expressions can be translated to  $\text{FO}^2$  queries. To this end, we will define

- (a) for every CoreXPath node expression  $e$  an  $\text{FO}^2$  query  $\varphi_e(x)$  that is equivalent to the query  $e$  and
- (b) for every CoreXPath path expression  $e$  an  $\text{FO}^2$  query  $\varphi_e(x)$  that is equivalent to the query  $?e$ .

If  $e$  is a node expression, the definition of  $\varphi_e$  is as follows:

- If  $e$  is of the form  $(\text{lab} = a)$ , then  $\varphi_e = \text{Lab}_a(x)$ .
- If  $e$  is of the form  $(e_1 \wedge e_2)$ ,  $(e_1 \vee e_2)$ , or  $\neg e_1$ , then  $\varphi_e = (\varphi_{e_1} \wedge \varphi_{e_2})$ ,  $\varphi_e = (\varphi_{e_1} \vee \varphi_{e_2})$ , or  $\varphi_e = (\neg \varphi_{e_1})$ , respectively.

If  $e$  is a path expression, the definition of  $\varphi_e$  is as follows:

- If  $e$  is of the form **self**, **next-sibling** or **preceding-sibling**, then  $\varphi_e = x$ ,  $\varphi_e = \exists y E^{\text{ns}}(x, y)$  or  $\varphi_e(x) = \exists y E^{\text{ns}}(y, x)$ , respectively.
- If  $e$  is of the form **following-sibling** or **preceding-sibling**, then  $\varphi_e = \exists y E^{\text{fs}}(x, y)$  or  $\varphi_e = \exists y E^{\text{fs}}(y, x)$ , respectively.
- If  $e$  is of the form  $(e_1/e_2)$ , then  $\varphi_e = \varphi_{e_1} \wedge (\varphi_{e_2}[x/y, y/x])$ .
- If  $e$  is of the form  $(e_1 \cup e_2)$ , then  $\varphi_e = \varphi_{e_1} \vee \varphi_{e_2}$ .
- If  $e$  is of the form  $e_1[e_2]$ , where  $e_2$  is a node expression, then  $\varphi_e = \varphi_{e_1} \wedge \varphi_{e_2}[x/y, y/x]$ .

Recall that  $\varphi_{e_2}[x/y, y/x]$  is the formula obtained from  $\varphi_{e_2}$  by swapping  $x$  and  $y$ , that is, simultaneously replacing every occurrence of  $x$  with  $y$  and vice versa. We leave the correctness of the translation as an exercise.

Conversely, let  $q = \varphi(x)$  be a unary  $\text{FO}^2$  query. We show a recursive translation procedure to turn  $\varphi(x)$  into an equivalent CoreXPath node expression  $e_\varphi$ . We can assume without loss of generality that  $\varphi(x)$  only uses the Boolean connectives  $\vee$  and  $\neg$  and the quantifier  $\exists$ .

If the formula  $\varphi$  is atomic, that is, of the form  $\text{Lab}_a(x)$ , then  $e_\varphi = (\text{lab} = a)$ . If  $\varphi$  is of the form  $\psi_1 \vee \psi_2$  or  $\neg\psi$ , then we recursively compute  $e_{\psi_1} \vee e_{\psi_2}$  or  $\neg e_\psi$ , respectively. Now, there is one remaining case, which is that  $\varphi$  is of the form  $\exists y \psi$ . Since  $\varphi(x)$  is an  $\text{FO}^2$  query, the variable  $x$  must be free in  $\varphi$ , which means that either  $\text{FV}(\psi) = \{x\}$  or  $\text{FV}(\psi) = \{x, y\}$ .

In the first case, the formula  $\varphi$  is equivalent to

$$\exists y (\psi \wedge (\text{lab}(y) = a \vee \neg \text{lab}(y) = a)),$$

which therefore reduces to the case where  $\text{FV}(\psi) = \{x, y\}$ . So the only remaining case is  $\varphi = \exists y \psi$  with  $\text{FV}(\psi) = \{x, y\}$ . In the remainder of the proof, for a Boolean formula  $\beta$  over variables  $\{x_1, \dots, x_k\}$ , and for  $\text{FO}^2$  formulae  $\psi_1, \dots, \psi_k$ , we write  $\beta[\psi_1, \dots, \psi_k]$  for the formula obtained from  $\beta$  by replacing each  $x_i$  with  $\psi_i$ . Using this notation, and since  $\psi$  only mentions variables  $x$  and  $y$ , we can write  $\psi$  as

$$\beta[\chi_1, \dots, \chi_r, \xi_1, \dots, \xi_s, \zeta_1, \dots, \zeta_t],$$

where

- $\beta$  is a Boolean formula,
- each  $\chi_i$  is an atomic  $\text{FO}^2$  formula with  $\text{FV}(\chi_i) = \{x, y\}$ ,
- each  $\xi_i$  is an atomic or existential  $\text{FO}^2$  formula with  $\text{FV}(\xi_i) = \{x\}$ , and
- each  $\zeta_i$  is an atomic or existential  $\text{FO}^2$  formula with  $\text{FV}(\zeta_i) = \{y\}$ .

In order to be able to recurse on subformulas of  $\varphi(x)$ , we have to separate the  $\xi_i$ 's from the  $\zeta_i$ 's. We first introduce a case distinction on which of the subformulas  $\xi_i$ 's hold or not and obtain that  $\varphi$  is equivalent to

$$\bigvee_{\bar{\gamma} \in \{\text{true}, \text{false}\}^s} \left( \bigwedge_{i \in [s]} (\xi_i \leftrightarrow \gamma_i) \wedge \exists y \beta(\chi_1, \dots, \chi_r, \gamma_1, \dots, \gamma_s, \zeta_1, \dots, \zeta_t) \right).$$

We proceed with a case distinction on which order relation holds between  $x$  and  $y$ . These are five mutually exclusive cases, determined by the following formulas, which we call *order types*:  $x = y$ ,  $E^{\text{ns}}(x, y)$ ,  $E^{\text{ns}}(y, x)$ ,  $E^{\text{fs}}(x, y) \wedge \neg E^{\text{ns}}(x, y)$ , or  $E^{\text{fs}}(y, x) \wedge \neg E^{\text{ns}}(y, x)$ . When we assume that one of these order types  $\tau$  is true, each atomic order formula evaluates to either **true** or **false**. In particular, each of the  $\chi_i$ 's evaluates to either **true** or **false**, and we will denote this truth value by  $\chi_i^\tau$ . Taking  $\mathcal{T}$  as the set containing the five order types, we can now obtain that  $\varphi$  is equivalent to

$$\bigvee_{\bar{\gamma} \in \{\text{true}, \text{false}\}^s} \left( \bigwedge_{i \in [s]} (\xi_i \leftrightarrow \gamma_i) \wedge \bigvee_{\tau \in \mathcal{T}} \exists y (\tau \wedge \beta(\chi_1^\tau, \dots, \chi_r^\tau, \bar{\gamma}, \bar{\zeta})) \right).$$

If  $\tau$  is an order type,  $\eta(x)$  an  $\text{FO}^2$  query, and  $e_\eta$  an equivalent CoreXPath formula, there is an obvious way to obtain a CoreXPath expression  $e\langle \tau, \eta \rangle$ , as shown in the following table.



$\tau$	$e\langle\tau, \eta\rangle$
$x = y$	$e_\eta$
$E^{\text{ns}}(x, y)$	<b>next-sibling</b> $[e_\eta]$
$E^{\text{ns}}(y, x)$	<b>previous-sibling</b> $[e_\eta]$
$E^{\text{fs}}(x, y) \wedge \neg E^{\text{ns}}(x, y)$	<b>next-sibling/following-sibling</b> $[e_\eta]$
$E^{\text{fs}}(y, x) \wedge \neg E^{\text{ns}}(y, x)$	<b>previous-sibling/preceding-sibling</b> $[e_\eta]$

Let us denote by  $\zeta_i[x]$  the formula  $\zeta_i$  in which we substituted the free occurrences of  $y$  with  $x$ . Our recursive procedure will then recursively compute  $e_{\zeta_i}$  for every  $i \in [s]$  and  $e_{\zeta_i[x]}$  for every  $i \in [t]$  and output

$$\bigvee_{\bar{\gamma} \in \{\text{true}, \text{false}\}^s} \left( \bigwedge_{i \in [s]} (e_{\zeta_i} \leftrightarrow \gamma_i) \wedge \bigvee_{\tau \in \mathcal{T}} e\langle\tau, \beta(\chi_1^\tau, \dots, \chi_r^\tau, \bar{\gamma}, e_{\zeta_1[x]}, \dots, e_{\zeta_t[x]})\rangle \right).$$

This concludes the translation. We leave the proof of its correctness as an exercise.  $\square$

## Unary CondXPath and FO

We now turn to CondXPath and FO. First, we shed some light on the difference between  $\text{FO}^2$  and FO. By definition, every  $\text{FO}^2$  query is also an FO query, but the converse does not hold. Indeed, one can show that the unary FO query  $\varphi(x)$  with

$$\begin{aligned} \varphi = & \text{following-sibling}(x, y) \wedge \text{Lab}_a(x) \wedge \text{Lab}_b(y) \\ & \wedge \forall z ((\text{following-sibling}(x, z) \wedge \text{following-sibling}(z, y)) \\ & \quad \rightarrow \text{Lab}_c(z)) \end{aligned}$$

cannot be expressed as an  $\text{FO}^2$  query over trees. The query returns nodes  $x$  that have a following-sibling  $y$  such that  $x$  is labeled  $a$ ,  $y$  is labeled  $b$  and every node *between*  $x$  and  $y$  is labeled  $c$ . Intuitively, the two-variable fragment cannot express this condition, because one needs one variable for  $x$  and  $y$  each, and then one lacks the third variable to express the condition between  $x$  and  $y$ .

Therefore, concerning unary queries, FO over trees is strictly more expressive than CoreXPath and  $\text{FO}^2$  over trees. However, it can be shown that the addition of  $(\text{step}[\text{ne}])^*$ , which we have in CondXPath, makes the language equally expressive as FO. In fact, this is one of the reasons why such a construct has been added to the XPath standard. We state the following result without proof.

**Theorem 59.3**

Unary CondXPath is equally expressive as unary FO over trees.

**Binary Queries**

It is natural to ask if Theorem 59.2 can be generalized to binary queries. It turns out that this is not the case, since the binary CoreXPath query

(next-sibling/next-sibling)

cannot be defined as a binary  $\text{FO}^2$  query (Exercise 8.3).

The connection between CondXPath and FO is stronger, though, since the binary FO queries over trees are indeed the same as the binary CondXPath queries.

**Theorem 59.4**

Binary CondXPath is equally expressive as binary FO over trees.



## Static Analysis of XPath

Static analysis tasks such as satisfiability, containment, and equivalence for queries are useful for query optimization. Indeed, if a (sub)query is not satisfiable, then the empty result can be returned without looking at the data. In turn, if a query is satisfiable, it makes sense to test if it can be rewritten into an equivalent query that can be more efficiently evaluated on the data. In this chapter we study satisfiability, containment, and equivalence for CondXPath.

### Satisfiability

We first investigate **Satisfiability**. From a complexity theory point of view, satisfiability for CoreXPath and CondXPath are rather complex, since both problems are EXPTIME-complete (Exercise 8.5). Here, we present the main argument why the problems are in EXPTIME.

#### Theorem 60.1

CondXPath-Satisfiability is in EXPTIME.

*Proof (Sketch).* We prove the result only for CondXPath expressions that

- do not use unions of path expressions,
- only use  $(\text{step}[\text{ne}])^*$  in the form  $(\text{child}[\text{ne}])^*$ , and
- only use the axes **child**, **descendant**, and **descendant-or-self**.

To this end, let  $q$  be a CondXPath query that satisfies these conditions. We will first rewrite  $q$  into an equivalent query that allows us to reduce the number of cases we need to consider later in the proof. Using De Morgan's law, we can rewrite  $q$  such that it only uses  $\wedge$  and  $\neg$ . Using associativity, we can rewrite all subexpressions of the form  $((\text{pe}/\text{pe})/\text{pe})$  as  $(\text{pe}/(\text{pe}/\text{pe}))$ . We

can replace subexpressions of the form  $\text{step}[e_1][e_2]$  by  $\text{step}[e_1 \wedge e_2]$  until no such expression occurs anymore. Finally, we can rewrite **descendant** into **child/descendant-or-self** and then **descendant-or-self** into  $(\text{child}[(\text{lab} = *)])^*$ , so that  $q$  only uses the **child** axis.

We will prove how we can iteratively compute sets  $S_n(T)$  and  $S_p(T)$  of subexpressions of  $q$  such that  $T$  is a tree and

- $S_n(T) = \{e_n \mid e_n \text{ is a node subexpression of } q \text{ with } \text{Root}(T) \in \llbracket e_n \rrbracket_T\}$  and
- $S_p(T) = \{e_p \mid e_p \text{ is a path subexpression of } q \text{ with } \text{Root}(T) \in \llbracket ?e_p \rrbracket_T\}$ .

Let  $\mathcal{S}$  be the set of all such sets. Notice that  $|\mathcal{S}|$ , the number of sets in  $\mathcal{S}$ , is at most  $2^{\|q\|}$ .

The computation of  $\mathcal{S}$  will be a *fixpoint computation* in a process that iterates over the increasing depth of trees  $T$ . More precisely, in a first iteration, we compute  $S_n(T)$  and  $S_p(T)$  for the empty tree and single-node trees. Later iterations consider trees  $T = a(T_1, \dots, T_n)$  of depth  $k$ , where we already computed  $S_n(T_i)$  and  $S_p(T_i)$  for all  $i \in [n]$  in a previous iteration. We stop once we find an iteration that has not found any new set in  $\mathcal{S}$ . This would mean that, for every tree of  $T$  depth  $k$ , the set  $S_n(T)$  equals the set  $S_n(T')$  for some  $T'$  of depth smaller than  $k$  (similarly for  $S_p(T)$ ).

During the entire proof, we will use a set  $\Sigma = \{a \mid (\text{lab} = a) \text{ is a subexpression of } q\} \cup \{\#\}$ . (We assume without loss of generality that  $(\text{lab} = \#)$  is not a subexpression of  $q$ . The symbol  $\#$  plays the role of representing “every label that does not appear in  $q$ ”.) The first iteration of the fixpoint algorithm computes all the sets  $S_n(T)$  and  $S_p(T)$  for  $T = \varepsilon$  and for  $T = \sigma$  for every  $\sigma \in \Sigma$ . According to Theorem 58.4, these sets can be computed in polynomial time.

We now assume that we have computed  $S_n(T)$  and  $S_p(T)$  for all trees of depth  $k - 1$ , degree at most  $\|q\|$ , and with labels in  $\Sigma$ . The next iteration now considers every  $\sigma \in \Sigma$  and  $\ell \leq \|q\|$ . We then consider every tree  $T = \sigma(T_1, \dots, T_\ell)$ , such that we have computed  $S_n(T_i)$  and  $S_p(T_i)$  for all  $i \in [\ell]$ . We will now decide, for every node subexpression  $e_n$  and every path subexpression  $e_p$  of  $q$ , if  $e_n \in S_n(T)$  or not and if  $e_p \in S_p(T)$  or not. We make these decisions inductively on the structure of the subexpressions.

- If  $e_n$  is of the form  $(\text{lab} = a)$ , then we decide  $e_n \in S_n(T)$  if and only if  $\sigma = a$ .
- If  $e_n$  is of the form  $(\text{lab} = *)$ , then we decide  $e_n \in S_n(T)$ .
- If  $e_n$  is of the form  $(e_1 \wedge e_2)$ , then we decide  $e_n \in S_n(T)$  if and only if  $e_1 \in S_n(T)$  and  $e_2 \in S_n(T)$ .
- If  $e_n$  is of the form  $\neg e$ , then we decide  $e_n \in S_n(T)$  if and only if  $e \notin S_n(T)$ .
- If  $e_n$  is of the form  $?e$  for a path expression  $e$ , then we decide  $e_n \in S_n(T)$  if and only if  $e \in S_p(T)$ .

This concludes all cases for node subexpressions. We now turn to path subexpressions.

- If  $e_p$  is of the form `child`, then we decide that  $e_p \in S_p(T)$ .
- If  $e_p$  is of the form `child` $[e]$ , then we decide that  $e_p \in S_p(T)$  if and only if  $e \in S_n(T_i)$  for some  $i \in [\ell]$ .
- If  $e_p$  is of the form `(child` $[e])^*[f]$ , then we decide that  $e_p \in S_p(T)$  if and only if  $f \in S_n(T)$  or there exists some  $i \in [\ell]$  such that  $e \in S_n(T_i)$  and  $e_p \in S_p(T_i)$ .
- If  $e_p$  is of the form  $e_1/e_2$ , then we do a case distinction on the structure of  $e_1$ .
  - If  $e_1$  is of the form `child`, then we decide that  $e_p \in S_p(T)$  if and only if  $e_2 \in S_p(T_i)$  for some  $i \in [\ell]$ .
  - If  $e_1$  is of the form `(child` $[e])^*[f]$ , then we decide that  $e_p \in S_p(T)$  if and only if either (1)  $f \in S_n(T)$  and  $e_2 \in S_p(T)$  or (2) there exists some  $i \in [\ell]$  such that  $e \in S_n(T_i)$  and  $e_p \in S_p(T_i)$ .

This concludes the computation of  $\mathcal{S}$ .

We argue why  $\mathcal{S}$  can be computed in exponential time. This holds because  $|\mathcal{S}| \leq 2^{\|q\|}$  at all times of the algorithm. Therefore, we can do at most  $|\mathcal{S}| \leq 2^{\|q\|}$  iterations before we reach a fixpoint. Furthermore, in every iteration, we consider at most  $|\Sigma| \cdot \|q\|^{\|q\|}$  trees.

We leave the proof of correctness as an exercise. (This requires proving that considering the types of trees that we have considered in this proof, i.e., trees that use labels from  $\Sigma$  and have degree at most  $\|q\|$  is sufficient.)  $\square$

Since every CoreXPath expression is also a CondXPath expression, we have the following corollary.

### Corollary 60.2

CoreXPath-Satisfiability is in EXPTIME.

## Containment and Equivalence

Since CoreXPath and CondXPath node expressions are closed under the Boolean operations, we can reduce equivalence and containment to satisfiability.

**Theorem 60.3**

1. CoreXPath-Containment is in EXPTIME.
2. CondXPath-Containment is in EXPTIME.

*Proof (Sketch).* For unary queries, we have that  $e_n$  is contained in  $e'_n$  if and only if  $(e_n \wedge \neg e'_n)$  is not satisfiable. Concerning binary queries, there is a construction that reduces the containment problem of binary queries to the one of unary queries (Exercise 8.7).

Since equivalence can be decided by testing containment in both directions, we also have the following.

**Corollary 60.4**

1. CoreXPath-Equivalence is in EXPTIME.
2. CondXPath-Equivalence is in EXPTIME.

## Tree Pattern Queries

*Tree pattern queries* are a simple query language that uses the child and descendant relation in trees in order to select tuples of nodes. They are important for querying tree-structured data because of two reasons. The first is that they closely correspond to a natural part of XPath and CoreXPath, namely the fragment that navigates with **child** and **descendant**, and uses conjunction, label tests, wildcards, and filter expressions (denoted  $[\cdot]$ ). The second reason is that this XPath fragment is widely used in practice, which makes tree pattern queries important from a practical point of view.

In this sense, tree pattern queries play a similar fundamental role for tree-structured data as conjunctive queries for relational data. An important difference between tree pattern queries and conjunctive queries though is that tree pattern queries are always acyclic. This is due to the syntax of CoreXPath, which cannot specify cyclic queries.

### Definition and Semantics

In order to define tree pattern queries, we first generalize the definition of labeled unordered trees to incorporate wildcards. We say that a connected, node-labeled, directed graph with wildcards  $T = (V, E, \text{lab})$  is a *labeled unordered tree with wildcards* if,

- for every node  $v$ , there is at most one node  $u$  with  $(u, v) \in E$  and
- there is exactly one node  $v$  (called the *root* of  $T$ ) without an incoming edge  $(u, v)$ .

We are now ready to define tree pattern queries.

#### Definition 61.1: Tree Pattern Query

A  $k$ -ary tree pattern query (TPQ) is a tuple



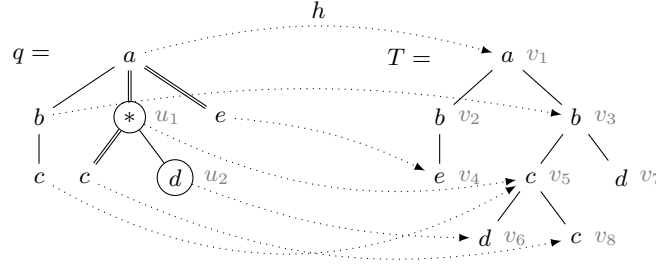


Fig. 61.1: A tree pattern query  $q$  (left), a tree  $T$  (right), and a tree pattern homomorphism from  $q$  to  $T$ .

$$q = (V, E^c, E^d, \text{lab}, \bar{v})$$

where

- $E^c$  and  $E^d$  are disjoint,
- $(V, E^c \cup E^d, \text{lab})$  is a labeled unordered tree with wildcards, and
- $\bar{v} = (v_1, \dots, v_k)$  is a  $k$ -tuple of *output nodes*

Similar to first-order queries over trees, we will sometimes denote a tree pattern query  $q$  as  $q(\bar{v})$  to emphasize that  $\bar{v}$  is the tuple of output nodes of  $q$ . We refer to  $E^c$  and  $E^d$  as the *child edges* and *descendant edges* of  $q$ , respectively. If a node is labeled “\*”, we call it a *wildcard node*. When we represent TPQs graphically, we draw child edges using single lines and descendant edges using double lines, see Figure 61.1.

We define  $\|q\|$  to be  $|V| + k$ , i.e., the number of nodes of  $q$  plus the arity  $k$  of the tuple of output nodes. Notice that we do not require the  $k$  output nodes to be distinct. We call  $q$  a  $k$ -ary TPQ. If  $k = 0$ , we call the query *Boolean*. We sometimes write Boolean TPQs as  $q = (V, E^c, E^d, \text{lab})$  to simplify notation. We define BoolTPQ to be the set of Boolean TPQs.

Intuitively, a TPQ can be matched in a tree if there exists a function from the nodes of the TPQ to the nodes of tree that satisfies all constraints imposed by the query. We define this next.

### Definition 61.2: Semantics of Tree Pattern Queries

Let  $q = (V, E^c, E^d, \text{lab}, \bar{v})$  be a TPQ and  $T = (V_T, E_T^c, \text{lab}_T)$  be a labeled unordered tree. We say that a label  $a \in \text{Lab}$  *matches* a node  $v \in V$  if  $\text{lab}(v) = a$  or  $\text{lab}(v) = *$ . A function  $h: V \rightarrow V_T$  is a *tree pattern homomorphism* from  $q$  to  $T$  if it fulfills all the following conditions:

- for every  $v \in V$ , the label  $\text{lab}_T(h(v))$  matches  $v$ ;

- if  $(u, v) \in E^c$ , then  $h(v)$  is a child of  $h(u)$  in  $T$ ; and
- if  $(u, v) \in E^d$ , then  $h(v)$  is a descendant of  $h(u)$  in  $T$ .

By  $q \rightarrow T$  we denote that there exists a tree pattern homomorphism from  $q$  to  $T$ . Conversely,  $q \nrightarrow T$  denotes that there does not exist such a tree pattern homomorphism. The *output* of  $q(\bar{v})$  on  $T$  is defined as

$$q(T) = \{h(\bar{v}) \mid h \text{ is a homomorphism from } q \text{ to } T\}.$$

As such, the *data model* of TPQs is the set of labeled unordered trees.

Notice that we do not require a homomorphism to be injective.

Notice that, even though we define the data model of TPQs to be the set of labeled unordered trees, their semantics can also be defined on ordered trees  $T = (V, E^c, E^{ns}, \text{lab})$ . (In fact, the definition is exactly the same.)

### Example 61.3: Semantics of Tree Pattern Queries

Figure 61.1 depicts a TPQ  $q$ , a tree  $T$ , and a tree pattern homomorphism  $h$  from  $q$  to  $T$ . Assume that the tuple of output nodes of  $q$  is  $(u_1, u_2)$  — we circled these two nodes in the Figure. The homomorphism  $h$  produces the answer  $(v_5, v_6)$  in  $T$ . Another homomorphism can be obtained from  $h$  by mapping  $u_1$  and  $u_2$  to nodes  $v_3$  and  $v_7$ , respectively. Furthermore,  $q(T) = \{(v_3, v_7), (v_5, v_6)\}$ . We note that there exist different tree pattern homomorphisms that produce  $(v_3, v_7)$  (by mapping the  $c$ -labeled sibling of node  $u_2$  in  $q$  to node  $v_5$  or to node  $v_8$ , respectively).

## Relationship to XPath

Unary and binary tree pattern queries naturally correspond to a fragment of CoreXPath that uses the **child** and **descendant** axes, label tests, wildcards, conjunction, and the  $[\cdot]$ -operator. This fragment of CoreXPath turns out to be very widely used in practice. We will not prove this correspondence formally, but provide an illustrating example.

### Example 61.4: XPath versus Tree Pattern Queries

The pattern  $q$  in Figure 61.1 can be defined in CoreXPath as

$$\text{self} \left[ \begin{array}{l} (\text{lab} = *) \wedge ?\text{descendant}[(\text{lab} = c)] \\ \wedge ?\text{ancestor} \left[ \begin{array}{l} (\text{lab} = a) \wedge ?\text{child}[(\text{lab} = b) \wedge ?\text{child}[(\text{lab} = c)]] \\ \wedge ?\text{descendant}[(\text{lab} = e)] \end{array} \right] \end{array} \right] / \text{child}[(\text{lab} = d)]$$

More generally, one can prove the following.

**Proposition 61.5**

For every unary or binary tree pattern query  $q$ , there exists an equivalent CoreXPath query  $q'$ . Furthermore,

1.  $q'$  can be constructed in linear time and
2.  $q'$  only uses `child`, `descendant`,  $\wedge$ ,  $(\text{lab} = a)$ ,  $(\text{lab} = *)$ , and  $[\cdot]$ .

## Evaluation

Since TPQs have labeled unordered trees as their associated data model, the problem TPQ-Evaluation is defined as follows.

**Problem: TPQ-Evaluation**

**Input:** A TPQ  $q$ , a labeled unordered tree  $T = (V, E^c, \text{lab})$ , and a tuple  $\bar{u} \in V^k$

**Output:** `true` if  $\bar{u} \in q(T)$  and `false` otherwise

We will show that this problem can be solved in polynomial time. First we reduce it to a simpler, Boolean variant of the problem. We present a logarithmic space reduction, but it is easy to see that the reduction can also be carried out in linear time on a random-access machine.

**Lemma 61.1.** *Let  $q$  be a  $k$ -ary TPQ,  $T$  a tree, and  $\bar{u}$  a  $k$ -tuple of nodes of  $T$ . Then there exists a Boolean TPQ  $q_b$  and a tree  $T_b$  such that  $\bar{u} \in q(T)$  if and only if  $q_b \rightarrow T_b$ . Furthermore,  $q_b$  and  $T_b$  can be computed in logarithmic space.*

*Proof.* Let  $q = (V_q, E_q^c, E_q^d, \text{lab}_q, \bar{v})$  with  $\bar{v} = (v_1, \dots, v_k)$ , let  $T = (V, E^c, \text{lab})$ , and let  $\bar{u} = (u_1, \dots, u_k)$ . Let  $z$  and  $o_1, \dots, o_k$  be elements of **Lab** not appearing in  $T$  or  $q$ . The tree  $T_b$  is obtained from  $T$  by attaching to each node  $u_i$  a new child  $u'_i$  with label  $o_i$  and to each leaf node  $u \notin \{u_1, \dots, u_k\}$  a new child  $u'$

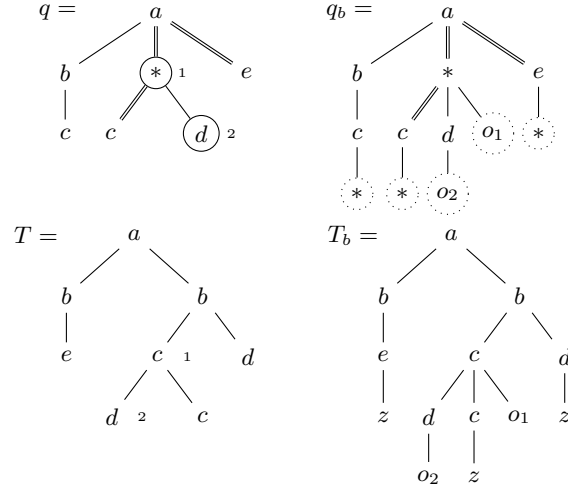


Fig. 61.2: Illustration of the reduction from TPQ-Evaluation to BoolTPQ-Evaluation



Fig. 61.3: Illustration why the reduction in Lemma 61.1 needs to attach the leaf nodes labeled  $z$  and  $*$ .

with label  $z$ . Similarly, pattern  $q_b$  is obtained from  $q$  by attaching to each node  $v_i$  a new child  $v'_i$  with label  $o_i$  and to each leaf node  $v \notin \{v_1, \dots, v_k\}$  a new child  $v'$  with label  $*$ . It is easy to show that  $\bar{u} \in q(T)$  if and only if  $q_b \rightarrow T_b$ .  $\square$

#### Example 61.6: Construction in Lemma 61.1

We illustrate the construction in Lemma 61.1 on an example TPQ  $q$  and labeled unordered tree  $T$  in Figure 61.2. The construction essentially adds nodes with special labels to the output nodes of  $q$  and the nodes in  $T$  that are to be tested. Leaf nodes different from these special nodes receive an additional child.

One may wonder why it is necessary to attach new nodes (labeled  $*$  and  $z$ , respectively) to the leaf nodes of  $q$  and  $T$  in the proof of Lemma 61.1. The

reason is that this is necessary for the correctness of the construction, as we explain next. Consider the pattern  $q$  and tree  $T$  in Figure 61.3. Furthermore,  $\tilde{q}_b$  and  $\tilde{T}_b$  are obtained from  $q$  and  $T$  by the construction as in Lemma 61.1 but *without attaching the new nodes to the leaf nodes*. Assume that  $u_1$  is the  $b$ -labeled node in  $T$ . Then  $u_1 \notin q(T)$  because  $q$  requires that its output node is not a child. However,  $\tilde{q}_b \rightarrow \tilde{T}_b$ , which means that attaching new leaf nodes is indeed necessary.

**Theorem 61.7**

TPQ-Evaluation is in PTIME.

*Proof.* By Lemma 61.1, it suffices to show how to test  $q \rightarrow T$  in PTIME for a Boolean pattern  $q$ . So, assume that  $T = (V, E^c, \text{lab})$  and  $q = (V_q, E_q^c, E_q^d, \text{lab}_q)$ . The idea is that we compute two sets of nodes of  $q$  for each node  $u$  of  $T$ , namely

- $\text{match}(u) = \{v \in V_q \mid q|_v \rightarrow T|_u\}$  and
- $\text{match}'(u) = \{v \in V_q \mid \text{there exists a descendant } u' \text{ of } u \text{ such that } q|_v \rightarrow T|_{u'}\}$

We compute these sets in a bottom-up fashion. First, for each leaf  $u$  of  $T$ , we define  $\text{match}'(u) = \emptyset$  and  $\text{match}(u) = \{v \in V_q \mid v \text{ is a leaf in } q \text{ and } \text{lab}(u) \text{ matches } v\}$ .

Now, assume that  $u$  is a node in  $T$  with children  $\{u_1, \dots, u_n\}$  such that we know all sets  $\text{match}(u_i)$  and  $\text{match}'(u_i)$ . Then,  $\text{match}(u)$  is the set of nodes  $v \in V_q$  that are matched by  $\text{lab}(u)$  and such that, for every edge  $(v, v')$  of  $q$ , there exists a child  $u_i$  of  $u$  for which one of the following holds:

- If  $(v, v')$  is a child edge, then  $v' \in \text{match}(u_i)$ .
- If  $(v, v')$  is a descendant edge,  $v' \in \text{match}(u_i) \cup \text{match}'(u_i)$ .

The set  $\text{match}'(u)$  is simply the union of all sets  $\text{match}(u_i)$  and  $\text{match}'(u_i)$ . Finally, the algorithm accepts if there exists a node  $u \in V$  such that  $\text{Root}(q) \in \text{match}(u)$ . It is easy to see that the algorithm runs in polynomial time.  $\square$

The complexity in Theorem 61.7 can be improved to  $O(\|q\| \cdot \|T\|)$  by a direct algorithm that does not use the reduction to Boolean patterns, see Exercise 8.8.

## Tree Pattern Query Containment and Equivalence

Recall that the main static analysis problems for query languages are satisfiability, containment, equivalence, and minimization. We will study these problems in the present and the next chapter.

First, observe that satisfiability for tree pattern queries is trivial, since every tree pattern query is satisfiable by the tree obtained from the pattern by

1. replacing each descendant edge with a child edge and
2. replacing each wildcard by an arbitrary label.

Therefore, the first nontrivial problems we study are containment and equivalence.

### Containment and Equivalence

Recall that tree pattern query  $q_1$  is *contained* in tree pattern query  $q_2$  (denoted  $q_1 \subseteq q_2$ ), if  $q_1(T) \subseteq q_2(T)$  for every labeled unordered tree  $T$ . Tree pattern query  $q_1$  is *equivalent* to  $q_2$  (denoted  $q_1 \equiv q_2$ ), if  $q_1 \subseteq q_2$  and  $q_2 \subseteq q_1$ .

We first prove that it suffices to focus on Boolean TPQs to study the complexity of containment and equivalence.

**Lemma 62.1.** *Let  $q_1$  and  $q_2$  be TPQs. Then there exist Boolean TPQs  $q_1^b$  and  $q_2^b$  such that  $q_1 \subseteq q_2$  if and only if  $q_1^b \subseteq q_2^b$ . Furthermore,  $q_1^b$  and  $q_2^b$  can be computed in logarithmic space.*

*Proof.* The construction of  $q_i^b$  from  $q_i$  is the same as the construction of  $q_b$  from  $q$  in the proof of Lemma [61.1](#).  $\square$

Next, we prove that it suffices to focus on the containment problem, i.e., for Boolean TPQs, containment and equivalence are interreducible.

Assume that  $q_1$  and  $q_2$  are Boolean TPQs. On the one hand,  $q_1 \equiv q_2$  if and only if  $q_1 \subseteq q_2$  and  $q_2 \subseteq q_1$ . On the other hand, suppose that we want to test if  $q_1 \subseteq q_2$ . Let  $q'_1$  be the TPQ obtained from  $q_1$  by adding a new root labeled  $*$  and connecting it with a child edge to the root of  $q_1$ . Likewise, let  $q_1 \cap q_2$  be the TPQ obtained from  $q_1$  and  $q_2$  by adding a new root labeled  $*$  and attaching two child edges, one to the root of  $q_1$  one to the root of  $q_2$ , respectively. Then, we have that  $q_1 \subseteq q_2$  if and only if  $q'_1 \equiv q_1 \cap q_2$ . Therefore, and by Lemma 62.1, we have now obtained the following.

**Proposition 62.1**

There exists logarithmic-space reductions between all of the following problems:

- TPQ-Containment
- TPQ-Equivalence
- BoolTPQ-Containment
- BoolTPQ-Equivalence

For this reason, our focus for studying containment and equivalence of tree pattern queries will be on the BoolTPQ-Containment problem.

## Containment Via Homomorphisms

We proved in Theorem 15.4, the Homomorphism Theorem, that a conjunctive query  $q_1$  is contained in  $q_2$  if and only if there exists a homomorphism from  $q_2$  to  $q_1$ . We now prove a similar result for TPQs. This time, however, the equivalence between containment and the existence of a homomorphism only holds for a fragment of TPQs. To this end, let  $\text{BoolTPQ}_{\text{Lab}}$  be the set of Boolean TPQs without wildcard nodes, i.e., every node has some label from  $\text{Lab}$ .

We generalize the definition of tree pattern homomorphisms (Definition 61.2) to include functions from tree pattern queries to tree pattern queries.

**Definition 62.2: Homomorphism between Tree Pattern Queries**

Let  $q_1 = (V_1, E_1^c, E_1^d, \text{lab}_1)$  and  $q_2 = (V_2, E_2^c, E_2^d, \text{lab}_2)$  be Boolean TPQs. A function  $h : V_1 \rightarrow V_2$  is a *tree pattern homomorphism from  $q_1$  to  $q_2$* , if

- every  $v \in V_{q_1}$  matches  $\text{lab}_2(h(v))$ ,
- if  $(u, v) \in E_1^c$  then  $h(u, v) \in E_2^c$ , and
- if  $(u, v) \in E_1^d$  then  $h(u)$  is an ancestor of  $h(v)$  in  $q_2$ .

We write  $q_1 \rightarrow q_2$  if there exists a tree pattern homomorphism from  $q_1$  to  $q_2$ . Similarly, we write  $q_1 \nrightarrow q_2$  if no such tree pattern homomorphism exists.

We are now ready to state a result for TPQs that is similar to the Homomorphism Theorem for conjunctive queries.

**Theorem 62.3**

If  $q_1$  and  $q_2$  are in  $\text{BoolTPQ}_{\text{Lab}}$ , then

$$q_1 \subseteq q_2 \iff q_2 \rightarrow q_1$$

*Proof.* Assume that  $h$  is a tree pattern homomorphism from  $q_2$  to  $q_1$ . Then if  $h_1$  is a tree pattern homomorphism from  $q_1$  to a tree  $T$ , we have that  $h_1 \circ h$  is a tree pattern homomorphism from  $q_2$  to  $T$ .

Conversely, assume that  $q_1 \subseteq q_2$ . Let  $z$  be a label from  $\text{Lab}$  not appearing in  $q_2$ . Notice that  $z$  exists because  $\text{Lab}$  is infinite. Let  $T$  be the tree obtained from  $q_1$  by replacing each descendant edge  $(u, v)$  with two child edges  $(u, u_{uv})$  and  $(u_{uv}, v)$ , where  $u_{uv}$  is a new node, labeled  $z$ . Since  $q_1(T) = \text{true}$  and  $q_1 \subseteq q_2$ , we also have that  $q_2(T) = \text{true}$ . Therefore, there is a tree pattern homomorphism  $h$  from  $q_2$  to  $T$ . Since  $q_2$  does not have wildcard nodes, the image of  $h$  does not contain any new nodes of the form  $n_{uv}$  and, therefore,  $h$  is a tree pattern homomorphism from  $q_2$  to  $q_1$ .  $\square$

We now turn to complexity. The following result is analogous to Theorem 61.7

**Proposition 62.4**

If  $q_1$  and  $q_2$  are Boolean TPQs, then it can be tested in PTIME whether  $q_1 \rightarrow q_2$ .

The following corollary is immediate from Theorem 62.3 and Proposition 62.4

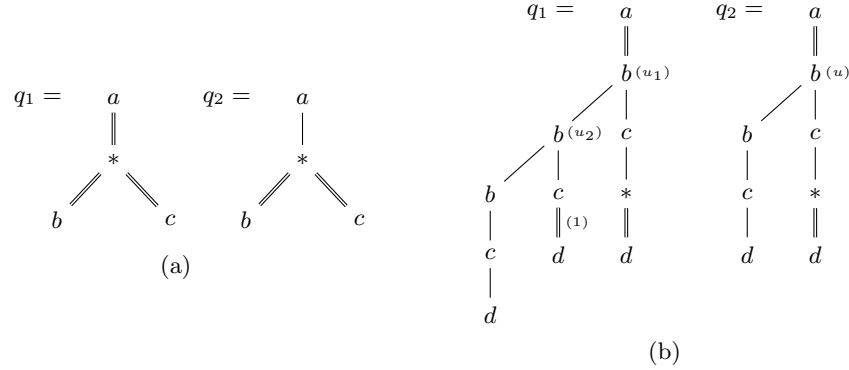
**Corollary 62.5**

$\text{BoolTPQ}_{\text{Lab}}$ -Containment and  $\text{BoolTPQ}_{\text{Lab}}$ -Equivalence are in PTIME.

We note that Corollary 62.5 also holds in the case where the queries are not Boolean, see Exercise 8.9

We now discuss why we restricted the aforementioned results to tree pattern queries without wildcards. If the queries use wildcard nodes, there are obvious examples that show that the existence of a homomorphism is not necessary for containment, see Figure 62.1(a). Here, the left query is contained



Fig. 62.1: Two examples where  $q_1 \subseteq q_2$  even though  $q_2 \not\sim q_1$ 

in the right query (and is even equivalent), even though no homomorphism exists from right to left. Figure 62.1(b) illustrates a more complex example. Here one can see that the left query is contained in the right query by case distinction on the edge marked (1). Assume that  $T \models q_1$  with match  $m$ . If (1) is mapped to a single edge in  $T$ , then node  $u$  of  $q_2$  can be mapped to  $m(u_1)$ . Otherwise, it can be mapped to  $m(u_2)$ .

## Star Extensions and Canonical Trees

Since testing the existence of a homomorphism is not sufficient to test containment of TPQs, the question arises if there exists an alternative method. Here we make a first step towards such a method, by showing that it is sufficient to consider *canonical trees*, which we define next.

### Definition 62.6: Star Extensions and Canonical Trees

Let  $q$  be a Boolean TPQ.

- A *star extension* of  $q$  is a TPQ without descendant edges, obtained from  $q$  by replacing each descendant edge  $(u, v)$  by a path  $uu_1^{uv} \cdots u_k^{uv}v$  where, for every  $i \in [k]$  the node  $u_i^{uv}$  is new and labeled  $*$ .
- A tree  $T$  is a *canonical tree* of  $q$  if it can be obtained from a star extension  $q^*$  of  $q$  by labeling each wildcard node by some fixed label  $z$  that does not appear in  $q$ . By  $T_z[q^*]$  we denote the canonical tree that was obtained in this manner. Notice that canonical trees  $T$  of  $q$  always match  $q$ , since the identity function on the nodes of  $q$  is always a match of  $q$  in  $T$ .

The importance of canonical trees is captured in the following Lemma, which shows that it suffices to restrict our attention to canonical trees for deciding whether  $q_1 \subseteq q_2$ .

**Lemma 62.2.** *If  $q_1$  and  $q_2$  are Boolean TPQs. Then  $q_1 \subseteq q_2$  if and only if  $q_2(T) = \mathbf{true}$  for every canonical tree  $T$  of  $q_1$ .*

*Proof.* The direction from left to right is trivial. We prove the other direction by contraposition. Assume that  $q_1 \not\subseteq q_2$  and let  $z$  be a label that does not occur in  $q_2$ . Let  $T$  be a (not necessarily canonical) tree such that  $q_1(T) = \mathbf{true}$  and  $q_2(T) = \mathbf{false}$ . Let  $h_1$  be a homomorphism from  $q_1$  to  $T$ . For each descendant edge  $(u, v)$  of  $q_1$ , let  $k_{uv}$  be the number of nodes between  $h_1(u)$  and  $h_1(v)$  in  $T$ . Let  $q_1^*$  be the star extension of  $q_1$  where each descendant edge  $(u, v)$  is replaced by the path  $uu_1^{uv} \cdots u_{k_{uv}}^{uv}v$ . Let  $h^*$  be a homomorphism from  $q_1^*$  to  $T$ .

We claim that  $q_2(T_z[q_1^*]) = \mathbf{false}$ . Towards a contradiction, assume that  $h_2$  is a tree pattern homomorphism from  $q_2$  to  $T_z[q_1^*]$ . Notice that  $h_2$  can only map wildcard nodes of  $q_2$  to  $z$ -labeled nodes of  $T_z[q_1^*]$ , since  $z$  does not appear in  $q_2$ . But this means that  $h = h^* \circ h_2$  is a tree pattern homomorphism<sup>1</sup> of  $q_2$  in  $T$ , which contradicts that  $q_2(T) = \mathbf{false}$ .  $\square$

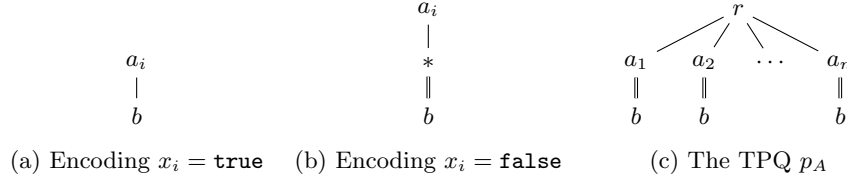
Next, we establish that it is even sufficient to consider canonical trees of polynomial size.

**Lemma 62.3.** *If  $q_1$  and  $q_2$  are Boolean TPQs such that  $q_1 \not\subseteq q_2$ , then there exists a tree  $T$  with  $|T| \leq |q_1|(|q_2| + 1)$  such that  $q_1(T) = \mathbf{true}$  and  $q_2(T) = \mathbf{false}$ .*

*Proof.* Let  $z$  be a label that does not occur in  $q_1$  or  $q_2$ . By Lemma 62.2 we know that there exists a canonical tree  $T_{\text{long}} = T_z[q_1^*]$  such that  $q_1(T_{\text{long}}) = \mathbf{true}$  and  $q_2(T_{\text{long}}) = \mathbf{false}$ .

We will now prove that long paths of  $z$ -labeled nodes in  $T_{\text{long}}$  can be shortened. Assume that  $T_{\text{long}}$  has a path  $\pi = uu_1^{uv} \cdots u_k^{uv}v$  with  $(u, v)$  a descendant edge in  $q_1$  and  $k > |q_2| + 1$ , so  $u_1^{uv}, \dots, u_k^{uv}$  are new nodes in  $q_1^*$ . Let  $T_{\text{short}}$  be obtained from  $T_{\text{long}}$  by replacing  $\pi$  with  $\pi' = uu_1^{uv} \cdots u_{|q_2|+1}^{uv}$ . Then clearly  $q_1(T_{\text{short}}) = \mathbf{true}$ , because it is a canonical tree. We will show that  $q_2(T_{\text{short}}) = \mathbf{false}$ . Assume the contrary. Then there is a tree pattern homomorphism  $h : q_2 \rightarrow T_{\text{short}}$ . By the pigeon hole principle, there is at least one node  $v_j$  in  $\{v_1^{uv}, \dots, v_{|q_2|+1}^{uv}\}$  that is not in the image of  $h$ . But then  $h$  is also a homomorphism from  $q_2$  to the tree  $T'$ , obtained from  $T_{\text{short}}$  by replacing  $\pi'$  with  $\pi'' = uv_1^{uv} \cdots v_{j-1}^{uv}v_{|q_2|+2}^{uv} \cdots v_k^{uv}v_j^{uv} \cdots v_{|q_2|+1}^{uv}$ . However,  $T'$  is isomorphic with  $T_{\text{long}}$ , which would mean that  $q_2 \rightarrow T_{\text{long}}$  and therefore  $q_2(T_{\text{long}}) = \mathbf{true}$ . This is a contradiction.  $\square$

<sup>1</sup> We note that  $h^*$  and  $h_2$  can indeed be composed, since  $T_z[q_1^*]$  has the same set of nodes as  $q_1^*$ .

Fig. 62.2: Encoding **true**, **false**, and truth assignments

## The Complexity of TPQ Containment

We are now ready to settle the complexity of TPQ-Equivalence and TPQ-Containment.

### Theorem 62.7

TPQ-Containment and TPQ-Equivalence are **coNP**-complete.

*Proof.* Lemma 62.3 implies that **BoolTPQ-Containment** is in **coNP**, since we can guess  $T$  and check deterministically (using Theorem 61.7) that  $q_1 \rightarrow T$  and  $q_2 \not\rightarrow T$ . Using Proposition 62.1, this shows that **TPQ-Containment** and **TPQ-Equivalence** are in **coNP**.

We now prove that the problem is **coNP**-hard. In particular, we prove that **BoolTPQ-Containment** is **coNP** hard. Hardness for equivalence follows again from Proposition 62.1.

We reduce from the validity problem of 3DNF formulae, which is well known to be **coNP**-complete. To this end, let  $\varphi$  be a Boolean formula in 3DNF with variables  $\{x_1, \dots, x_n\}$ . We will construct Boolean TPQs  $p_1$  and  $p_2$  such that  $p_1 \subseteq p_2$  if and only if  $\varphi$  is valid. Let

$$\varphi = c_1 \vee \dots \vee c_k,$$

where  $c_i$  is a clause of the form  $(\ell_{i,1} \wedge \ell_{i,2} \wedge \ell_{i,3})$  for each  $i \in [k]$ . Here, each  $\ell_{i,j}$  is called a *literal* and is either a variable or the negation thereof. We will first construct TPQs  $p_A$  and  $p_{C_1}, \dots, p_{C_k}$  such that  $\varphi$  is valid if and only if, for every tree  $T$ ,

$$p_A \rightarrow T \implies \exists i \in [k] \text{ such that } p_{C_i} \rightarrow T$$

The task of  $p_A$  is to “generate all the truth assignments”. To this end, we interpret the tree in Figure 62.2a as “ $x_i = \text{true}$ ” and every tree that matches the TPQ in Figure 62.2b as “ $x_i = \text{false}$ ”. As such, every tree that matches the TPQ  $p_A$  in Figure 62.2c can be understood as a truth assignment for  $\{x_1, \dots, x_n\}$ .

We now construct the patterns  $p_{C_1}, \dots, p_{C_k}$ . We first illustrate how to construct  $p_{C_i}$  by example. Assume that  $c_i = (x_1 \wedge \neg x_3 \wedge x_7)$ . Then  $p_{C_i}$  is

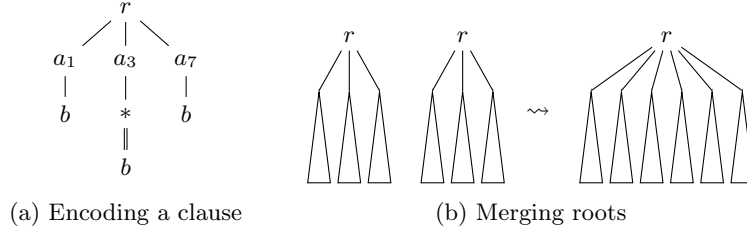
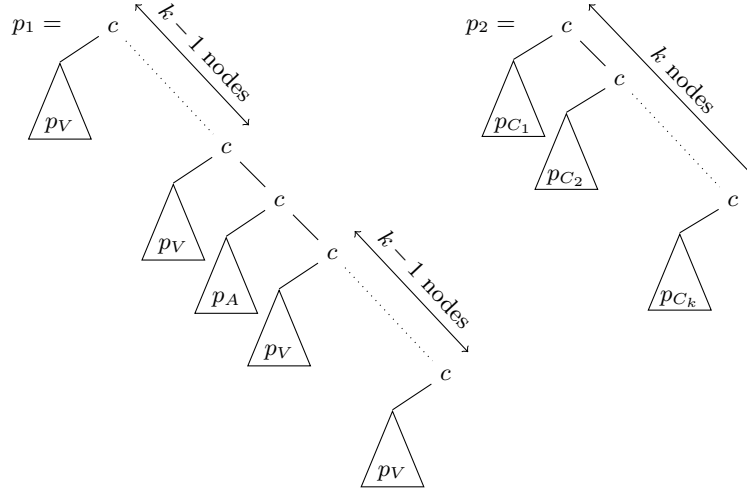


Fig. 62.3: Constructing clause patterns and all-match patterns

Fig. 62.4: Patterns  $p_1$  and  $p_2$  in Theorem 62.7

the TPQ in Figure 62.3a. More generally, assume that  $c_i = (\ell_{i,1} \wedge \ell_{i,2} \wedge \ell_{i,3})$ . Let  $f(i, j) \in [n]$  be the number such that  $x_{f(i,j)}$  is the variable mentioned in  $\ell_{i,j}$  for every  $j \in [3]$ . Let  $b_{i,j}$  be a Boolean value indicating whether  $\ell_{i,j}$  is a positive or negative literal. That is,  $b_{i,j} = \text{true}$  if  $\ell_{i,j} = x_{i,j}$  and  $b_{i,j} = \text{false}$  if  $\ell_{i,j} = \neg x_{i,j}$ . Then,  $p_{C_i}$  consists of a root, labeled  $r$ , below which we attach the three subpatterns encoding  $x_{f(i,j)} = b_{i,j}$ , for  $j \in [3]$ . It is easy to prove that  $\varphi$  is valid if and only if every tree that matches  $p_A$  matches at least one  $p_{C_i}$  for some  $i \in [k]$ .

In the next part of the proof, we turn  $p_A$  and  $p_{C_1}, \dots, p_{C_k}$  into two patterns  $p_1$  and  $p_2$  such that  $\varphi$  is valid if and only if  $p_1 \subseteq p_2$ . To this end, let  $V$  be the pattern obtained by merging the roots of  $p_{C_1}, \dots, p_{C_k}$  (see Figure 62.3b). We then define  $p_1$  and  $p_2$  as illustrated in Figure 62.4. We leave the proof that  $\varphi$  is valid if and only if  $p_1 \subseteq p_2$  as an exercise.  $\square$



## Tree Pattern Query Minimization

Tree pattern query minimization amounts to transforming a given TPQ into an equivalent one that has as few nodes as possible. We introduce minimizations of TPQs in a similar way as we did for conjunctive queries.

### Definition 63.1: Minimization of TPQs

Given a TPQ  $q$ , a TPQ  $q'$  is called a *minimization* of  $q$  if it is equivalent to  $q$  and has the smallest number of nodes among all the TPQs that are equivalent to  $q$ . A TPQ  $q$  is *minimal* if it is a minimization of itself.

In Chapter 16 we saw that, in the case of CQs, minimizations can be obtained by removing atoms from the query. We can now ask ourselves if something similar is true for TPQs. More precisely, we ask if TPQs can be minimized by *deleting nodes*. To this end, for a TPQ  $q = (V_q, E_q, \text{lab}_q, \bar{v})$  and a leaf  $u \in V_q$  that does not appear in  $\bar{v}$ , denote by  $(q - u)$  the query obtained from  $q$  by removing  $u$  from  $V_q$  and every edge of the form  $(u', u)$  from  $E_q$ . Notice that every subpattern of  $q$  that contains the root and output nodes of  $q$  can be obtained by repeatedly deleting leaves.

In Algorithm 17 we describe a procedure TRIMTPQ, which is similar to the COMPUTECORE algorithm for CQs, i.e., Algorithm 4 from Chapter 16. It repeatedly removes leaves from the pattern as long as the pattern stays equivalent. In some important cases, it is not very difficult to show that TRIMTPQ indeed computes a minimization. We leave the proof as Exercise 8.13.

### Theorem 63.2

Let  $q$  be a TPQ. Then TRIMTPQ( $q$ ) computes a minimization of  $q$  if either

**Algorithm 17** TRIMTPQ( $q$ )**Input:** A tree pattern query  $q$ **Output:** A subpattern  $q^*$  of  $q$  that is equivalent to  $q$ 

- 1:  $q^* := q$
  - 2: **while** there is a leaf  $u$  of  $q^*$  such that  $(q^* - u) \equiv q^*$  **do**
  - 3:    $q^* := (q^* - u)$
- return**
- $q^*$

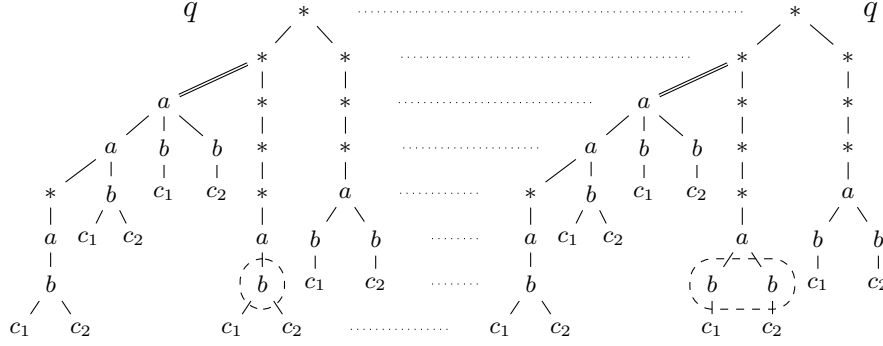
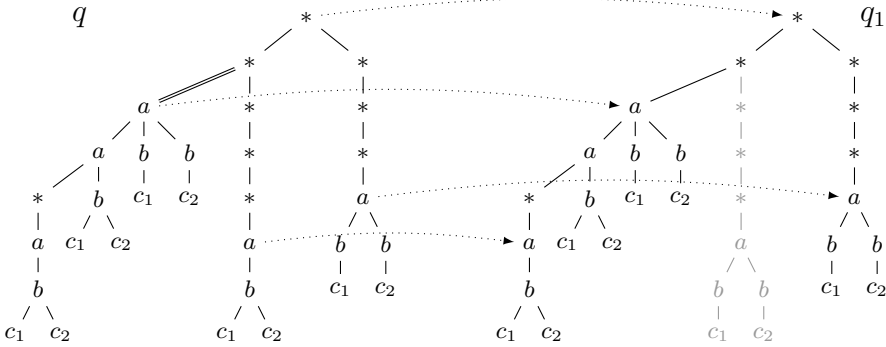


Fig. 63.1: A tree pattern query  $q$  that has no equivalent proper subpattern (right) and a tree pattern query  $q$  that is equivalent and smaller (left). The only difference between the two patterns is in the dashed lines.

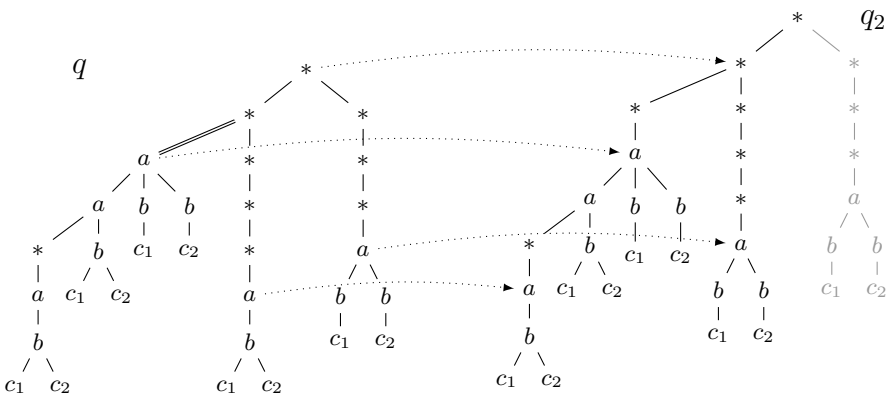
- (1)  $q$  does not contain descendant edges, or
- (2)  $q$  does not have wildcard nodes.

Perhaps surprisingly, Theorem 63.2 does not hold for TPQs in general. Figure 63.1 contains on the right a tree pattern query  $q$  that has no equivalent proper subpattern, that is, there is no node  $v$  with  $q \equiv (q - v)$ . However, it is not minimal, because the pattern on the left is equivalent and smaller. We prove that  $q$  and  $q$  are equivalent and leave the proof that  $q$  has no equivalent proper subpattern as an exercise (Exercise 8.14).

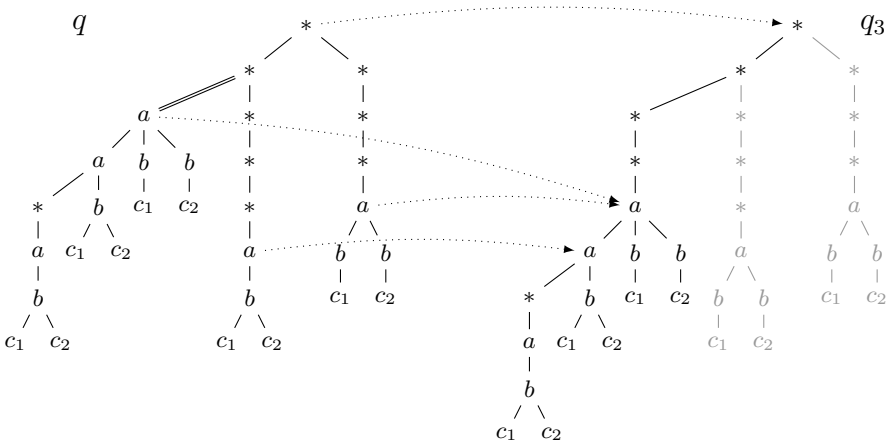
Consider the tree pattern queries  $q_1, \dots, q_5$ , depicted on the right hand sides of Figures 63.2 and 63.3. Observe that  $q$  is equivalent to  $q_1 \cup q_2 \cup q_3 \cup q_4 \cup q_5$ , since the only difference between these patterns and  $q$  is that they replace the lowermost descendant edge in  $q$  by paths of increasing length, where  $q_5$  has a descendant edge to deal with paths of length at least five. Figures 63.2 and 63.3 then show how homomorphisms from  $q$  to  $q_1, \dots, q_5$  can be constructed. This shows that  $q \subseteq q$ . The inclusion  $q \subseteq q$  is easy to see, because there exists a homomorphism from  $q$  to  $q$ . The homomorphism maps the two nodes in the dashed lines in  $q$  to the single such node in  $q$ .



(a) Homomorphism from  $q$  to  $q_1$



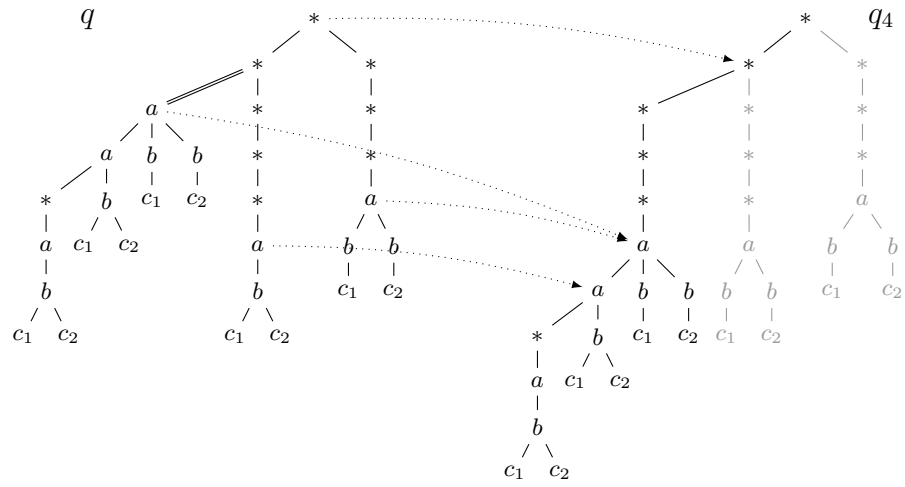
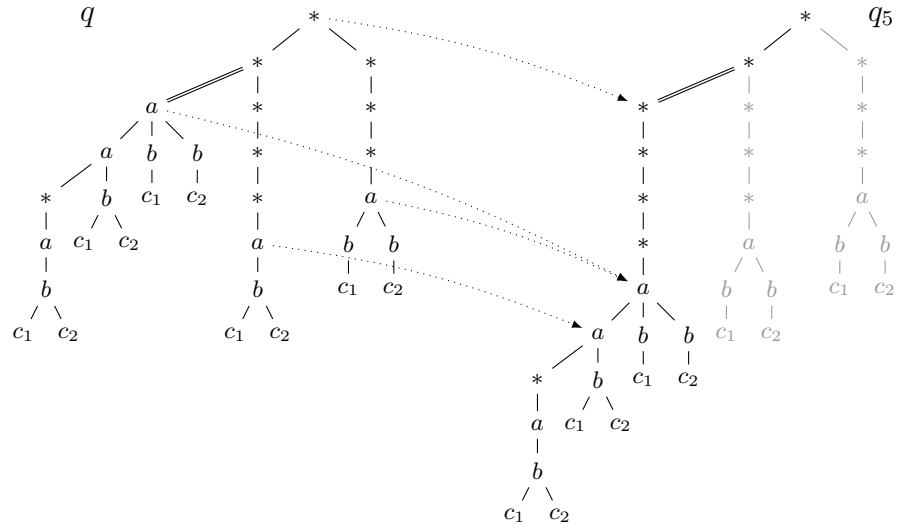
(b) Homomorphism from  $q$  to  $q_2$



(c) Homomorphism from  $q$  to  $q_3$

Fig. 63.2: Showing that patterns  $q$  and  $q$  in Figure 63.1 are equivalent



(a) Homomorphism from  $q$  to  $q_4$ (b) Homomorphism from  $q$  to  $q_5$ Fig. 63.3: Showing that patterns  $q$  and  $q$  in Figure 63.1 are equivalent

We conclude this chapter with a note on the complexity of computing minimizations of TPQs. Minimizations of a TPQ  $q$  can be computed by a naive algorithm that iterates through all TPQs that are smaller than  $q$  and tests whether they are equivalent, but such an algorithm hardly seems satisfactory. Surprisingly, we do not know an algorithm that is significantly better! From the example in Figure 63.1, we know that deleting nodes is not sufficient for minimizing TPQs since, in this case, nodes also need to be merged. But there are examples that show that deleting and merging nodes still is not sufficient: in some cases, nodes also need to be split. As such minimization is a more complex problem for TPQs than for CQs, for example. Indeed, one can show that, whereas it is  $\text{CONP}$ -complete to test if a given CQ is minimal (Exercise 2.17), this problem is  $\Pi_2^P$ -complete for TPQs (Exercise 8.15). The latter means that the “obvious algorithm” that tests if, for a given TPQ  $q$ , every smaller TPQ is not equivalent to  $q$  is worst-case optimal.



---

## Exercises

**Exercise 8.1.** Let  $T = (V, E^c, E^{ns}, \text{lab})$  be a labeled ordered tree. In the proof of Theorem 58.4, show that, for every **axis** in the set

`{self, child, parent, descendant, descendant-or-self, ancestor,  
ancestor-or-self, next-sibling, following-sibling,  
previous-sibling, preceding-sibling, following, preceding}`

and every set  $S \subseteq V$ , it is possible to compute  $\llbracket \text{axis} \rrbracket_{T,S}$  in time  $O(\|T\|)$ .

**Exercise 8.2 (\*).** What is the precise complexity of evaluation for downward XPath (child, descendant, filter, and, or, not)? For which complexity class is the evaluation problem complete under logspace reductions? (Open problem.)

**Exercise 8.3.** Show that the binary CoreXPath query

$$(\text{next-sibling}/\text{next-sibling}) \tag{63.1}$$

cannot be defined as a binary  $\text{FO}^2$  query.

*Hint:* It will be convenient to use a pebble characterization for  $\text{FO}^2$ .

**Exercise 8.4.** Show that every FO query over trees can be defined as an  $\text{FO}^3$  query, that is, an FO query that only uses three variables.

*Hint:* Show that every binary CondXPath query can be defined in  $\text{FO}^3$  over trees.

*Hint:* Look into the paper Immerman & Kozen. Definability by a bounded number of variables.

**Exercise 8.5.** Show that CoreXPath-Satisfiability and CondXPath-Satisfiability are EXPTIME-hard.

*Hint:* Reduce from two-player corridor tiling. For checking vertical constraints, an idea of encoding their violations in unary as in Bjorklund Martens Schwentick MFCS 13 will be useful.

**Exercise 8.6.** Show that the algorithm in the proof of Theorem 60.1 is correct.

**Exercise 8.7.** Show that CoreXPath-Containment and CondXPath-Containment for binary queries can be reduced in polynomial time to CoreXPath-Containment and CondXPath-Containment for unary queries, respectively

*Hint:* The construction is the proof of Theorem 6 in [?]. It consists of rewriting the queries in such a way that some axes are reversed.

**Exercise 8.8.** We note that the complexity in Theorem 61.7 can be improved to  $O(|p||T|)$ . *Hint:* Let the input to TPQ Evaluation be  $p$ ,  $T$ , and  $(u_1, \dots, u_k)$ . Let  $(v_1, \dots, v_k)$  be the output nodes of  $p$ . Then the overall idea of the direct algorithm is the same as in Theorem 61.7 but it treats the output nodes of  $p$  differently. More precisely, it only allows an output node  $v_i$  to be in a set  $\text{match}(u)$  if  $u = u_i$ .

**Exercise 8.9.** Prove that the restriction to Boolean queries in Corollary 62.5 can be lifted. *Hint:* Notice that one cannot just apply Proposition 62.1 since the reduction to Boolean queries introduces wildcard nodes. Instead, generalize Definition 62.2 and Theorem 62.3.

**Exercise 8.10.** Prove that  $\varphi$  is valid if and only if  $p_1 \subseteq p_2$  in the proof of Theorem 62.7.

**Exercise 8.11.** Consider *conjunctive queries over trees*, i.e., conjunctive queries evaluated over tree structures, where the built-in relations are *child* and *descendant*. What is the complexity of their evaluation problem?

**Exercise 8.12.** Show that the output of  $\text{COMPUTECORE}(p)$  (Algorithm 17 on page 510) is not unique up to isomorphism.

**Exercise 8.13.** Prove Theorem 63.2.

**Exercise 8.14.** Show that the pattern  $p$  in Figure 63.1 has no equivalent proper subpattern.

**Exercise 8.15.** (a) Let TPQ-Minimization be the problem where, given a Boolean TPQ  $q$  and integer  $k \in \mathbb{N}$ , the question is if there exists a TPQ  $q'$  such that  $q' \equiv q$  and  $|q'| \leq k$ . Prove that TPQ-Minimization is  $\Sigma_2^P$ -complete.

(b) Let TPQ-Minimality be the problem where, given a Boolean TPQ  $q$ , the question is to answer **true** if  $q$  is minimal and **false** otherwise. Prove that TPQ-Minimality is  $\Pi_2^P$ -complete.

---

## Bibliographic Comments

To be done.

**Expressive Languages for Tree-Structured Data**





In this part we move to more expressive languages for querying and specifying tree-structured data. In Chapter 64, we study *monadic second-order logic (MSO)* on trees. Boolean MSO formulas can define precisely the *regular tree languages* which, similarly to word languages, are also characterized by finite automata. We will see a third characterization, namely through *monadic Datalog*.

The idea of storing data as trees has led to the development of schema languages for trees, based on ideas from extended context-free grammars and regular tree languages. In Chapter 67 we define the structural core of the three most widespread schema languages for XML, namely DTD, XML Schema, and Relax NG and give some insights in their expressiveness.

We study queries and schema together in Chapter 68, where we look into query optimization *under schema information*. Queries can indeed be optimized more aggressively when schema information is taken into account, but the computational complexity for problems such as satisfiability and containment with respect to schemas also increases.



## Monadic Second Order Logic

In this chapter, we consider a significant extension of first-order logic over trees with quantification over sets. This means that we will be able to write formulae that say “there exists a set of nodes  $X$  such that ...” or “for every set  $X$  of nodes ...”, which is not possible in first-order logic. The logic we obtain in this manner is called *monadic second-order logic* or simply *MSO*.

In this chapter, we will define the syntax and semantics of MSO queries over trees. In the next chapters, we study their connections to *finite automata over trees* and *monadic Datalog*.

### Monadic Second Order Logic over Trees

We assume a countably infinite set

$$\text{SVar}$$

of *set variables*, disjoint from  $\text{Var}$  and  $\text{Nodes}$ . Set variables will be denoted by  $X, Y, Z$ , etc. Analogously to first-order logic over trees, formulae of monadic second-order logic over trees will use relation symbols from the set

$$\{\text{Lab}_a \mid a \in \text{Lab}\} \cup \{E^{\text{fc}}, E^{\text{ns}}, E^{\text{d}}, E^{\text{fs}}\},$$

where  $\text{Lab}_a$  is the *a-label relation*. In the following, we abbreviate *monadic second-order logic* as MSO.

#### Definition 64.1: Syntax of MSO over Trees

We define *formulae of MSO over trees* as follows:

- Every first-order formula over trees is an MSO formula over trees.
- If  $X \in \text{SVar}$  and  $y \in \text{Var}$ , i.e.,  $X$  is a set variable and  $y$  is a first-order variable, then  $X(y)$  is an MSO formula over trees.

- If  $\varphi_1$  and  $\varphi_2$  are MSO formulae over trees, then  $(\varphi_1 \wedge \varphi_2)$ ,  $(\varphi_1 \vee \varphi_2)$ , and  $(\neg \varphi_1)$  are MSO formulae over trees.
- If  $\varphi$  is an MSO formula over trees and  $X$  is a set variable, then  $(\exists X \varphi)$  and  $(\forall X \varphi)$  are MSO formulae over trees.

To avoid notional clutter, we will omit brackets in the same manner as for FO formulae over relational databases (Chapter 3). The set of *free* variables of a formula  $\varphi$ , denoted  $FV(\varphi)$ , is defined as follows:

- $FV(\text{Lab}_a(x)) = \{x\}$  for every  $a \in \text{Lab}$ .
- $FV(x = y) = FV(R(x, y)) = \{x, y\}$  for every  $R \in \{E^{\text{fc}}, E^{\text{ns}}, E^{\text{d}}, E^{\text{fs}}\}$ .
- $FV(X(y)) = \{X, y\}$ .
- $FV(\neg \varphi) = FV(\varphi)$ .
- $FV(\varphi_1 \vee \varphi_2) = FV(\varphi_1 \wedge \varphi_2) = FV(\varphi_1) \cup FV(\varphi_2)$ .
- $FV(\exists x \varphi) = FV(\forall x \varphi) = FV(\varphi) - \{x\}$ .
- $FV(\exists X \varphi) = FV(\forall X \varphi) = FV(\varphi) - \{X\}$ .

An MSO formula without free variables is called an MSO *sentence*.

#### Example 64.2: Monadic Second-Order Formulae

Consider the MSO formula over trees

$$\exists X (\forall z_1 \forall z_2 (E^{\text{ns}}(z_1, z_2) \rightarrow X(z_1) \leftrightarrow X(z_2))) \wedge X(x) \wedge \neg X(y) . \quad (64.1)$$

The free variables of this formula are  $x$  and  $y$ .

Consider the FO-formula over trees

$$\varphi_c(x, y) = E^{\text{fc}}(x, y) \vee (\exists z E^{\text{fc}}(x, z) \wedge E^{\text{fs}}(z, y)) ,$$

which is satisfied by nodes  $x$  and  $y$  if  $y$  is a child of  $x$ . The formulae

$$\varphi_{\text{root}}(x) = \neg \exists y E^{\text{fc}}(y, x) \text{ and}$$

$$\varphi_{\text{leaf}}(x) = \neg \exists y E^{\text{fc}}(x, y)$$

are satisfied by  $x$  if  $x$  is the root or a leaf of a tree, respectively. Consider the following MSO-formula over trees:

$$\begin{aligned} & \exists x \exists y \left( \varphi_{\text{root}}(x) \wedge \varphi_{\text{leaf}}(y) \wedge \right. \\ & \left. \exists X (\forall z_1 \forall z_2 (\varphi_c(z_1, z_2) \rightarrow (X(z_1) \leftrightarrow \neg X(z_2)))) \wedge X(x) \wedge X(y) \right) \quad (64.2) \end{aligned}$$

This formula has no free variables.

Given a tree  $T = (V, E^c, E^{ns}, \text{lab})$ , we define the notion of satisfaction of a formula  $\varphi$  with respect to an *assignment*  $\eta$  for  $\varphi$  over  $T$ . Such an assignment is a function from  $\text{FV}(\varphi)$  to  $V \cup 2^V$  such that  $\eta(x) \in V$  for every  $x \in \text{Var}$  and  $\eta(X) \subseteq V$  for every  $X \in \text{SVar}$ . If  $X$  is a set variable and  $S \subseteq V$  a set of nodes of  $T$ , we use  $\eta[X/S]$  to denote the assignment that modifies  $\eta$  by setting  $\eta(X) = S$ . We are now ready for defining the semantics of MSO formulae over trees.

#### Definition 64.3: Semantics of MSO over Trees

Given a tree  $T = (V, E^c, E^{ns}, \text{lab})$ , an MSO formula  $\varphi$ , and an assignment  $\eta$  for  $\varphi$  over  $T$ , we inductively define when  $\varphi$  is *satisfied* in  $T$  under  $\eta$ , written  $(T, \eta) \models \varphi$ , as follows. We only highlight the parts of the definition that are different from the definition of the semantics of FO over trees:

- If  $\varphi = X(y)$  then  $(T, \eta) \models \varphi$  if and only if  $\eta(y) \in \eta(X)$ .
- If  $\varphi = (\exists X \psi)$ , then  $(T, \eta) \models \varphi$  if and only if  $(T, \eta[S/X]) \models \psi$  for some  $S \subseteq V$ .
- If  $\varphi = (\forall X \psi)$ , then  $(T, \eta) \models \varphi$  if and only if  $(T, \eta[S/X]) \models \psi$  for every  $S \subseteq V$ .

#### Example 64.4: Monadic Second-Order Formulae

We give an intuitive description of the semantic meaning of the formulae in Example 64.2:

- Formula (64.1) is satisfied by all nodes  $x$  and  $y$  that are siblings and have an even number of nodes between them.
- Formula (64.2) is satisfied in every nonempty tree that has a path from root to leaf of even length.

### MSO Queries over Trees

Just like first-order queries over trees, monadic second-order queries over trees are  $k$ -ary data graph queries, that is, they take a data graph  $G$  as input and produce a set  $q(G) \subseteq (\text{Nodes} \cup \text{Const})^k$ . We first define the syntax of first-order queries over trees.

**Definition 64.5: Monadic Second-Order Queries over Trees**

A *monadic second-order query over trees* or *MSO query over trees* for short is an expression of the form  $\varphi(\bar{x})$ , where  $\varphi$  is an MSO formula over trees without free set variables, and  $\bar{x}$  is a tuple of free first-order variables of  $\varphi$  such that each free variable of  $\varphi$  occurs in  $\bar{x}$  at least once.

We define *size*  $\|\varphi(\bar{x})\|$  of an MSO query over trees  $\varphi(\bar{x})$  as  $\|\varphi\| + \|\bar{x}\|$ .

We now define the semantics of MSO queries over trees. Let  $\varphi(\bar{x})$  be an MSO query over trees. Given a labeled ordered tree  $T = (V, E^c, E^{\text{ns}}, \text{lab})$  and a tuple  $\bar{a}$  of elements from **Nodes**, we say that  $T$  *satisfies the query*  $\varphi(\bar{x})$  *using the values*  $\bar{a}$ , denoted by  $T \models \varphi(\bar{a})$ , if there exists an assignment  $\eta$  for  $\varphi$  over  $T$  such that  $\eta(\bar{x}) = \bar{a}$  and  $(D, \eta) \models \varphi$ . We can now define the data model and the output of MSO queries over trees. The data model defines the possible inputs of the queries. For each data graph in the data model, we define the output of the query.

**Definition 64.6: Evaluation of MSO Queries over Trees**

Given a labeled ordered tree  $T = (V, E^c, E^{\text{ns}}, \text{lab})$  and an MSO query over trees  $q = \varphi(x_1, \dots, x_k)$ , where  $k \geq 0$ , the *output* of  $q$  on  $T$  is defined as the set of tuples

$$q(T) = \{\bar{a} \in \text{Nodes}^k \mid T \models \varphi(\bar{a})\}.$$

The *data model* of FO queries over trees is the set of labeled ordered trees.

## Tree Automata

We now define *tree automata*, which are an extension of finite state automata, geared towards trees. We assume that readers are familiar with the basic theory of finite automata, regular expressions, and regular word languages. We provide an introduction and a reminder of the standard notation in Appendix [E](#).

### Definition 65.1: Nondeterministic Unranked Tree Automaton

A *nondeterministic unranked tree automaton (NTA)* is a quadruple

$$A = (Q, \Sigma, \delta, F) ,$$

where

- $Q$  is a finite set of *states*,
- $\Sigma \subseteq \mathbf{Lab}$  is a finite, non-empty set of labels,
- $F \subseteq Q$  is a set of *accepting states*, and
- $\delta : Q \times \Sigma \rightarrow 2^{Q^*}$  is its *transition function*, which is such that  $\delta(q, a)$  is a regular word language over  $Q$  for every  $a \in \Sigma$  and  $q \in Q$ .

Unless we say otherwise, we always assume that  $\delta(q, a)$  is represented by a nondeterministic finite automaton (NFA) over words.

The term *unranked* comes from *unranked trees*, which are trees in which every node may have an arbitrary number of children. This is opposed to *ranked trees* in which the number of children of a node is determined by the node's label. This number of children is called the *rank* of the label in the literature.

The *size* of an NTA  $A$ , denoted  $\|A\|$  is defined as  $|Q| + \sum_{(q,a) \in Q \times \Sigma} \|N_{q,a}\|$ , where  $\|N_{q,a}\|$  is the size of the NFA representing  $\delta(q, a)$  (we discuss  $\|N_{q,a}\|$  in Appendix [E](#)).

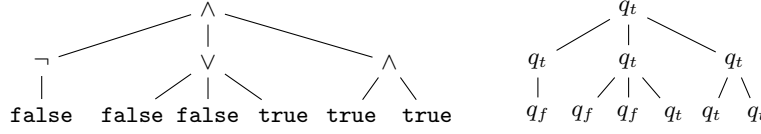


Fig. 65.1: A tree (left) and a run of the tree automaton of Example 65.2 over the tree.

A *run* of  $A$  on a tree  $T = (V, E^c, E^{ns}, \text{lab})$  is a mapping  $\lambda : V \rightarrow Q$  such that, for every  $u \in V$  with ordered sequence of children  $u_1, \dots, u_n$ , we have that

$$\lambda(u_1) \cdots \lambda(u_n) \in \delta(\lambda(u), \text{lab}(u)) .$$

Notice that, when  $u$  has no children, this condition reduces to  $\varepsilon \in \delta(\lambda(u), \text{lab}(u))$ . A run is *accepting* if  $\lambda(\text{Root}(T)) \in F$ , that is, it maps the root to an accepting state. A tree  $T$  is *accepted by*  $A$  if there exists an accepting run of  $A$  on  $T$ . The set of all accepted trees is denoted by  $L(A)$  and is called the *language of*  $A$ .

We extend the definition of  $\delta$  to trees by defining a function  $\delta^*$  as follows:

- $\delta^*(a) := \{q \mid \varepsilon \in \delta(q, a)\}$ ; and
- $\delta^*(a(T_1, \dots, T_n)) := \{q \mid \exists q_1 \in \delta^*(T_1), \dots, \exists q_n \in \delta^*(T_n) \text{ with } q_1 \cdots q_n \in \delta(q, a)\}$ .

Notice that a tree  $T$  is accepted by  $A$  if and only if  $\delta^*(T) \cap F \neq \emptyset$ .

#### Example 65.2: Unranked Tree Automaton

We give an example of a tree automaton  $A$  that recognizes tree representations of Boolean formulas that evaluate to **true** (as in Figure 65.1 left). The automaton has state set  $Q = \{q_t, q_f\}$  and accepting states  $F = \{q_t\}$ . The transition function is defined as follows (we use regular expressions to denote the regular languages):

$$\begin{aligned} \delta(q_t, \wedge) &= q_t^* & \delta(q_t, \vee) &= q_f^* q_t (q_t + q_f)^* & \delta(q_t, \neg) &= q_f & \delta(q_t, \text{true}) &= \varepsilon \\ \delta(q_f, \vee) &= q_f^* & \delta(q_f, \wedge) &= q_t^* q_f (q_t + q_f)^* & \delta(q_f, \neg) &= q_t & \delta(q_f, \text{false}) &= \varepsilon \end{aligned}$$

A run of  $A$  on a tree  $T$  in Figure 65.1 (left) is depicted in Figure 65.1 (right). For every node  $u$  of  $T$ , the state  $\lambda(u)$  appears on the position that corresponds to  $u$  in the tree on the right. For instance, if  $u$  is the rightmost leaf of  $T$ , then  $\lambda(u) = q_t$ .

The languages that can be accepted by unranked tree automata form a very robust class, known as the *regular tree languages*.



**Definition 65.3: Regular Tree Language**

A set  $S$  of trees is called a *tree language*. If there exists an unranked tree automaton  $A$  such that  $S = L(A)$ , then  $S$  is also called a *regular tree language*.

**Algorithms on Tree Automata**

We review some standard algorithms on tree automata. First we show that it can be tested in linear time combined complexity if a given tree is accepted by a given NTA.

**Problem: NTA-Acceptance**

**Input:** A labeled ordered tree  $T$  and an NTA  $A$   
**Output:** `true` if  $T \in L(A)$  and `false` otherwise

**Theorem 65.4**

NTA-Acceptance can be decided in PTIME.

*Proof.* The algorithm computes the sets  $\delta^*(T')$  for all subtrees  $T'$  of  $T$  in a bottom-up manner. For each leaf labeled  $a$  and each state  $q$  of  $A$ , it can be tested in constant time if  $q \in \delta^*(a)$  by testing if  $\varepsilon \in L(A_{q,a})$ , where  $A_{q,a}$  is the NFA for  $\delta(q, a)$ .

For each internal node, rooting a subtree  $a(T_1, \dots, T_n)$ , we can test if  $q \in \delta^*(a(T_1, \dots, T_n))$  by testing if there is a word in  $(\delta^*(T_1) \cdots \delta^*(T_n)) \cap L(A_{q,a})$ .  
 $\square$

Next, we show that unions and intersections of tree automata can be constructed in time proportional to the product of the sizes of the two respective automata.

**Theorem 65.5**

Let  $A_1$  and  $A_2$  be NTAs. Then NTAs for  $L(A_1) \cap L(A_2)$  and  $L(A_1) \cup L(A_2)$  can be constructed in time  $O(\|A_1\| \|A_2\|)$ .

*Proof.* The result follows from a “product construction” for the involved automata. Let  $A_1 = (Q_1, \Sigma, \delta_1, F_1)$  and  $A_2 = (Q_2, \Sigma, \delta_2, F_2)$ . We first explain how to construct an automaton  $A$  for  $L(A_1) \cap L(A_2)$ . Here,  $A = (Q_1 \times Q_2, \Sigma, \delta, F_1 \times F_2)$ , where  $\delta((q, q'), a) = \{(q_1, q'_1) \cdots (q_n, q'_n) \mid q_1 \cdots q_n \in \delta(q, a)\}$

**Algorithm 18** Computing the set  $R$  of reachable states of an NTA

---

```

1:  $R_1 := \{q \in Q \mid \exists a \in \Sigma \text{ such that } \varepsilon \in \delta(q, a)\};$ 
2: for  $i := 2$  to  $|Q|$  do
3:    $R_i := \{q \in Q \mid \text{there exists an } a \in \Sigma \text{ such that } \delta(q, a) \cap R_{i-1}^* \neq \emptyset\};$ 
4:  $R := R_{|Q|};$ 

```

---

and  $q'_1 \cdots q'_n \in \delta(q', a)\}$ . If the NFAs for  $\delta(q, a)$  and  $\delta(q', a)$  use  $k_1$  and  $k_2$  states respectively, then an NFA for  $\delta((q, q'), a)$  with  $k_1 k_2$  states can be constructed. The construction for  $L(A_1) \cup L(A_2)$  is analogous but uses  $(F_1 \times Q_2) \cup (Q_1 \times F_2)$  as accepting states. We leave the proof of correctness as an exercise.  $\square$

We now discuss the problem of testing whether a tree automaton accepts at least one tree. This problem is an automata-theoretic version of *satisfiability* in logic and is typically called *non-emptiness*, since it asks if the language of the automaton is non-empty.

**Problem: NTA-Non-Emptiness**

**Input:** A non-deterministic tree automaton  $A$

**Output:** **true** if  $L(A) \neq \emptyset$  and **false** otherwise

**Theorem 65.6**

NTA-Non-Emptiness is in PTIME.

*Proof.* Solving NTA-Non-Emptiness amounts to deciding if there exists a tree  $T$  such that  $\delta^*(T)$  contains an accepting state. Let  $A = (Q, \Sigma, \delta, F)$  be an NTA. Algorithm 18 computes the set of states  $R := \{q \mid \exists \text{ tree } T \text{ such that } q \in \delta^*(T)\}$  in a bottom-up manner. Clearly,  $L(A) \neq \emptyset$  if and only if  $R \cap F \neq \emptyset$ . Note that  $R_i \subseteq R_{i+1}$  and  $R_1 = \{\delta^*(a) \mid a \in \Sigma\}$ . We argue that the algorithm is in PTIME. Clearly,  $R_1$  can be computed in linear time. Further, the for-loop makes a linear number of iterations. Every iteration is a linear number of non-emptiness tests of the intersection of an NFA with  $R_{i-1}^*$  where  $R_{i-1} \subseteq Q$ . As emptiness for NFAs is in linear time, the latter is also in linear time. We leave the proof of correctness as an exercise.  $\square$

**Equivalence of MSO and Tree Automata**

We will prove that MSO sentences and tree automata are equally expressive over *binary labeled ordered trees*, that is, labeled ordered trees in which every node has zero or two children. The result also holds for general labeled ordered

trees, but the proof for the binary case is less technical. The difference between the proofs for the binary and the general case mainly lies in the following lemma, which shows that NTAs can be complemented. We denote the set of binary trees with labels in  $\Sigma$  by  $\mathcal{T}_\Sigma^b$ .

**Lemma 65.7.** *Let  $A$  be an NTA such that  $L(A) \subseteq \mathcal{T}_\Sigma^b$ . Then an NTA for  $\mathcal{T}_\Sigma^b - L(A)$  can be constructed in time  $O(2^{\|A\|})$ .*

*Proof.* Let  $A = (Q_A, \Sigma, \delta_A, F_A)$ . The main idea of the proof is to compute a powerset automaton for  $A$  and to complement it. More precisely, we first construct an NTA  $B = (Q_B, \Sigma, \delta_B, F_B)$  for  $L(A)$  as follows. Define  $Q_B = 2^{Q_A}$  and  $F_B = \{S \mid S \cap F_A \neq \emptyset\}$ . The transition function  $\delta_B$  is defined such that, for each tree  $T$ ,  $\delta_B^*(T) = \{S\}$ , where  $S = \delta_A^*(T)$ . That is,  $\delta_B(S, a) = \{\varepsilon \mid \varepsilon \in \delta(q, a) \text{ for every } q \in S\} \cup \{S_1 S_2 \mid \forall q \in S, \exists q_1 \in S_1, \exists q_2 \in S_2 \text{ such that } q_1 q_2 \in \delta_A(q, a)\}$ . Finally, the NTA for  $\mathcal{T}_\Sigma^b - L(A)$  is  $\bar{B} = (Q_B, \Sigma, \delta_B, Q_B - F_B)$ . We leave the proof of correctness as an exercise.  $\square$

### Definition 65.8: MSO-definable tree languages

Let  $S$  be a tree language. We say that  $S$  is *MSO-definable* if there exists an MSO sentence  $\varphi$  such that  $S = \{T \mid T \models \varphi\}$ .

### Theorem 65.9

Let  $L$  be a set of trees using labels from a nonempty finite set  $\Sigma$ . Then  $L$  is MSO-definable if and only if it is regular.

*Proof.* We will prove the theorem for binary trees only. Therefore, let  $L$  be a regular language of binary trees, which means that  $L$  is the language of an NTA  $A = (Q, \Sigma, \delta, F)$ . We can assume w.l.o.g. that  $Q = [n]$ . We construct an MSO formula  $\varphi$  such that  $L = \{T \mid T \models \varphi\}$ . Intuitively,  $\varphi$  expresses the existence of an accepting run. For every  $i \in [n]$ , it uses a set variable  $X_i$  for the set of nodes visited in state  $i$ . Consider the formulae

$$\begin{aligned} \varphi_1 &= (\forall x (\bigvee_{i=1}^n X_i(x))) \wedge (\bigvee_{i \neq j} \forall x (X_i(x) \rightarrow \neg X_j(x))), \\ \varphi_2 &= \forall x (\neg(\exists y E^{\text{fc}}(y, x)) \rightarrow \bigvee_{i \in F} X_i(x)), \\ \varphi_3 &= \bigwedge_{a \in \Sigma} \forall x (\text{Lab}_a(x) \wedge \neg(\exists y E^{\text{fc}}(x, y))) \rightarrow \bigvee_{\varepsilon \in \delta(i, a)} X_i(x), \text{ and} \\ \varphi_4 &= \bigwedge_{a \in \Sigma} \forall x \forall x_1 \forall x_2 ((\text{Lab}_a(x) \wedge E^{\text{fc}}(x, x_1) \wedge E^{\text{ns}}(x_1, x_2)) \\ &\quad \rightarrow (\bigvee_{(k, \ell) \in \delta(i, a)} (X_i(x) \wedge X_k(x_1) \wedge X_\ell(x_2))))). \end{aligned}$$

Formula  $\varphi_1$  expresses that every node is assigned exactly one state,  $\varphi_2$  expresses that the root is assigned an accepting state,  $\varphi_3$  expresses that the leaves are labeled in accordance with  $\delta$ , and  $\varphi_4$  expresses the internal nodes are labeled in accordance with  $\delta$ . It is easy to show that  $A$  accepts a tree  $T$

if and only if there exist sets  $X_1, \dots, X_n$  that satisfy  $\varphi_1, \varphi_2, \varphi_3$ , and  $\varphi_4$ . We conclude this direction by defining

$$\varphi = \exists X_1 \cdots \exists X_n (\varphi_1 \wedge \varphi_2 \wedge \varphi_3 \wedge \varphi_4).$$

For the other direction, let  $\varphi$  be an MSO-formula. Since the formulae  $E^{\text{fs}}(x, y)$  and  $E^{\text{d}}(x, y)$  are expressible in MSO using the relations  $E^{\text{fc}}$  and  $E^{\text{ns}}$ , we assume that  $\varphi$  only uses the relations  $E^{\text{fc}}$  and  $E^{\text{ns}}$ . In fact, we can assume w.l.o.g. that  $\varphi$  is generated by the grammar

$$\varphi := \text{Lab}_a(x) \mid E^{\text{fc}}(x, y) \mid E^{\text{ns}}(x, y) \mid x = y \mid \varphi \wedge \varphi \mid \neg \varphi \mid \exists x \varphi \mid \exists X \varphi,$$

where  $a \in \Sigma$ .<sup>1</sup> We will further simplify the syntax of MSO-formulae by eliminating the first-order variables. That is, we introduce another logic with the same expressive power than MSO over trees that we call  $MSO_0$ . The idea is that  $MSO_0$  introduces some syntactic sugar, but constrains formulae elsewhere.  $MSO_0$ -formulae  $\varphi$  over  $\sigma_{\text{fcns}}$  are generated over the grammar

$$\varphi := \text{sing}(X) \mid \text{Lab}_a(X) \mid X \subseteq Y \mid E^{\text{fc}}(X, Y) \mid E^{\text{ns}}(X, Y) \mid \varphi \wedge \varphi \mid \neg \varphi \mid \exists X \varphi,$$

meaning respectively that  $X$  is a singleton; all variables in  $X$  are labeled  $a$ ; the set  $X$  is a subset of  $Y$ ; the sets  $X$  and  $Y$  are singletons  $\{x\}, \{y\}$  with  $E^{\text{fc}}(x, y)$ ; the sets  $X$  and  $Y$  are singletons  $\{x\}, \{y\}$  with  $E^{\text{ns}}(x, y)$ ; and  $\exists X \varphi$ . The semantics of these formulae can be defined in MSO over  $\sigma_{\text{fcns}}$ . For instance,  $\text{sing}(X)$  can be written as  $\forall x \forall y (X(x) \wedge X(y)) \rightarrow (x = y)$  and  $\text{Lab}_a(X)$  as  $\forall x (X(x) \rightarrow \text{Lab}_a(x))$ . Conversely, it is easy to show by structural induction that each MSO-formula can be translated into an equivalent  $MSO_0$ -formula.

We can assume w.l.o.g. that  $\varphi$  is an  $MSO_0$  formula, that it does not re-use variables, and that its variables are  $X_1, \dots, X_n$ . For a subformula  $\psi$  of  $\varphi$  with free variables  $S$ , we will construct an NTA that accepts trees for which the labels are pairs  $(a, f)$ , where  $a \in \Sigma$  and  $f \in \{0, 1\}^S$ , i.e.,  $f$  is a function from  $S$  to  $\{0, 1\}$ . Intuitively, such a tree  $T$  represents a tree  $T'$  over  $\Sigma$ , together with an assignment of the free variables, and the task of the NTA is to accept precisely those that satisfy  $\psi$ .

More precisely, let  $T$  be a tree over  $\Sigma \times \{0, 1\}^S$ . We denote by  $T_\Sigma$  the tree obtained from  $T$  by replacing each label  $(a, f)$  by  $a$ . We define  $\eta_T$  to be the assignment obtained from  $T$  by mapping each free variable  $X_i$  to  $\eta(X_i) = \{u \mid \text{lab}(u) = (a, f) \text{ and } f(X_i) = 1\}$ .

We now inductively define NTAs  $A_\psi$  that accept precisely the trees  $T$  such that  $(T_\Sigma, \eta_T) \models \psi$ . An NTA for  $\text{Lab}_a(X_i)$  tests that the tree has no node labeled  $(b, f)$  with  $f(X_i) = 1$  and  $b \neq a$ . An NTA for  $X \subseteq Y$  is similar. An NTA for  $E^{\text{fc}}(X_i, X_j)$  tests that there is exactly one node  $u_1$  labeled  $(a, f_1)$  with  $f_1(X_i) = 1$ , exactly one node  $u_2$  labeled  $(b, f_2)$  with  $f_2(X_j) = 1$ , and that  $u_2$  is a child of  $u_1$ . The NTAs for  $\text{sing}(X)$  and  $E^{\text{ns}}(x_i, x_j)$  are similar.

<sup>1</sup> For instance, subformulae  $\text{Lab}_b(x)$  with  $b \notin \Sigma$  can be replaced by  $\text{Lab}_a(x) \wedge \neg \text{Lab}_a(x)$  for some  $a \in \Sigma$ .

For the inductive step, in the cases  $\psi = \neg\psi_1$  and  $\psi = \psi_1 \wedge \psi_2$ , we can use Theorem 65.5 and Lemma 65.7. In the case  $\psi = \exists X_i \psi'$ , assume that we have an NTA for  $\psi'$ . The NTA for  $\psi$  reads labels  $(a, f)$  for which  $f(X_i)$  is undefined. It guesses labels of the form  $(a, f')$  for which  $f$  and  $f'$  agree on every free variable from  $\psi$ , and then runs as if it would be the NTA for  $\psi'$ .  $\square$



## Monadic Datalog

We now establish another characterization of MSO-definable queries, namely in terms of Datalog programs (see Chapter 35). We call a Datalog program  $\Pi$  *monadic* if all intensional relations are unary. In the following, we will consider monadic Datalog programs over the relations

$$\{\text{Lab}_a \mid a \in \text{Lab}\} \cup \{E^{\text{fc}}, E^{\text{ns}}, E^{\text{d}}, E^{\text{fs}}\} \cup \{\text{Leaf}, \text{LastSib}\}.$$

Given a tree  $T = (V, E^c, E^{\text{ns}}, \text{lab})$ , we have that

- $\text{Lab}_a = \{v \mid \text{lab}(v) = a\}$ ,
- $E^{\text{fc}} = \{(u, v) \mid v \text{ is the first child of } u\}$ ,
- $E^{\text{ns}} = \{(u, v) \mid v \text{ is the next sibling of } u\}$ ,
- $E^{\text{d}} = \{(u, v) \mid v \text{ is a descendant of } u\}$ ,
- $E^{\text{fs}} = \{(u, v) \mid v \text{ is a following sibling of } u\}$ ,
- $\text{Leaf} = \{v \mid v \text{ is a leaf}\}$ , and
- $\text{LastSib} = \{v \mid v \text{ is a last sibling}\}$ .

A Boolean monadic Datalog program  $\Pi$  *recognizes* a set  $L$  of trees if  $L = \{T \mid \Pi(T) = \text{true}\}$ . A set of trees  $L$  is *recognizable* by Boolean monadic Datalog if there exists a Boolean monadic Datalog program  $\Pi$  that recognizes  $L$ .

We will prove that the Boolean monadic Datalog queries are precisely the Boolean MSO-definable queries. The equivalence between unary MSO-definable queries and monadic Datalog queries is known to hold for unary queries as well, but the result presented here is easier to prove.

### Theorem 66.1

A set  $L$  of binary trees is MSO-definable if and only if it is recognizable by a Boolean monadic Datalog program.

*Proof.* Let  $(\Pi, R_1)$  be a monadic Datalog query, where  $\text{edb}(\Pi) = \sigma_{\text{tree}}$  and  $\text{idb}(\Pi) = \{R_1, \dots, R_n\}$ . Then, the MSO query can be defined as

$$\varphi := \forall R_1 \dots \forall R_n \bigwedge_{\rho \in \Pi} \varphi_\rho,$$

where  $\varphi_\rho$  is the first-order sentence associated to Datalog rule  $\rho$ , that was used for defining the model-theoretic semantics of Datalog in Chapter 35.

Conversely, let  $A = (Q, \Sigma, \delta, F)$  be an NTA. We define a Boolean monadic Datalog program  $(\Pi, R)$  that tests if there exists an accepting run of  $A$  on a given tree  $T$ . To this end,  $(\Pi, R)$  will have a unary predicate  $q$  for each  $q \in Q$  that will be satisfied in node  $u$  of  $T$  if and only if  $q \in \delta^*(T|_u)$ .

For every state  $q \in Q$  and  $a \in \Sigma$ , let  $N_{q,a} = (Q_{q,a}, Q, \delta_{q,a}, I_{q,a}, F_{q,a})$  be an NFA for the language  $\delta(q, a)$ . We can assume w.l.o.g. that  $Q$  and all state sets  $Q_{q,a}$  are pairwise disjoint. In the following, we will always use  $q, q'$  to denote states from  $Q$  and  $p, p'$  to denote states from the sets  $Q_{q,a}$ .

Notice that we can view all  $N_{q,a}$  as one automaton  $(Q_N, Q, \delta_N, I_N, F_N)$ , where  $Q_N = \cup_{q \in Q} \cup_{a \in \Sigma} Q_{q,a}$ ,  $I_N = \cup_{q,a} I_{q,a}$ ,  $F_N = \cup_{q,a} F_{q,a}$ , and  $\delta_N$  is defined as follows. For a given  $p \in \cup_{q,a} Q_{q,a}$ , we define  $\delta_N(p, q)$  to be the set  $\delta_{q',a}(p, q)$  for the unique  $q'$  and  $a$  such that  $p \in Q_{q',a}$ . We define  $I_q$  to be the set of states reachable in  $N$  by reading  $q$  from an initial state, that is,  $I_q := \{p \mid \exists p' \in I_N \text{ such that } p \in \delta_N(p', q)\}$ .

For each transition such that  $\varepsilon \in \delta(q, a)$ , the program  $\Pi$  has the rule

$$q(x) :- \text{Leaf}(x) \wedge \text{Lab}_a(x).$$

Next, we define rules that capture the computation of the automata  $N_{q,a}$ . Therefore, the rules

$$\begin{array}{ll} p(x) :- E^{\text{fc}}(y, x), \text{Leaf}(x), q(x). & \text{for every } p \in I_q \\ p(x) :- E^{\text{fc}}(y, x), E^{\text{fc}}(x, y'), q(x). & \text{for every } p \in I_q \\ p(x) :- E^{\text{ns}}(y, x), p'(y), \text{Leaf}(x), q(x). & \text{for every } p \in \delta'(p', q) \\ p(x) :- E^{\text{ns}}(y, x), p'(y), E^{\text{fc}}(x, y'), q(x). & \text{for every } p \in \delta'(p', q) \end{array}$$

are added to  $\Pi$ . The rules

$$q(x) :- \text{Lab}_a(x), E^{\text{fc}}(x, y), \text{LastSib}(y, y'), p(y') \quad \text{for every } p \in F_{i,a}$$

test if the children of  $x$  can be visited in a sequence of states that is accepted by  $N_{q,a}$ . Finally, the rules  $R() :- \text{Root}(x), q(x)$  for every  $q \in F$  tests if the root can be labeled with an accepting state.  $\square$



## Schemas for XML

Database schemas for tree-structured data follow a different approach from those for relational data. The reason is that their task is more complex. Such schemas in practice do not only describe how attributes are associated to values, or which data types (integer, date, etc.) are allowed where, but they also describe *what is the admissible tree structure of the data*. Since one can follow different philosophies in how to specify all this information, there are different formalisms for specifying schemas for XML data — these are called *schema languages*. The most prominent schema languages are *Document Type Definition (DTD)*, *XML Schema*, and *Relax NG*. In this chapter, we explain how these three languages describe the structure of labeled ordered trees. The underlying ideas for all these languages come from *context-free grammars* (see Appendix E for a quick refresher) and from *tree automata*, which we defined in Chapter 65.

The first schema language for XML data is also the simplest and is called *Document Type Definition (DTD)*.

### Definition 67.1: Document Type Definition

A *Document Type Definition (DTD)* over alphabet  $\Sigma$  is a triple

$$d = (\Sigma, \rho, S) ,$$

where

- $\Sigma \subseteq \mathbf{Lab}$  is a finite set of *labels*,
- $\rho$  is a function from  $\Sigma$  to the set of regular expression over  $\Sigma$ , and
- $S \subseteq \Sigma$  is a set of *start labels*.

In practice, the labels in  $\Sigma$  are also called *element names*, which is a convention that we will also follow.

```

<?xml version="1.0" encoding="UTF-8"?>
<inventory>
  <category name="concert guitars">
    <item>
      <maker> Tandler </maker>
      <model> Advanced Student </model>
      <description> Spruce top,
                    Indian rosewood back and sides </description>
      <price> Please ask </price>
    </item>
    <item>
      <maker> Hanika </maker>
      <model> Grand Concert </model>
      <description> Spruce top,
                    Palo Escrito back and sides </description>
      <price> 4299 </price>
      <discount> 10 </discount>
    </item>
    [...]
  </category>
  [...]
</inventory>

```

Fig. 67.1: Fragment of an XML document representing the inventory of a store

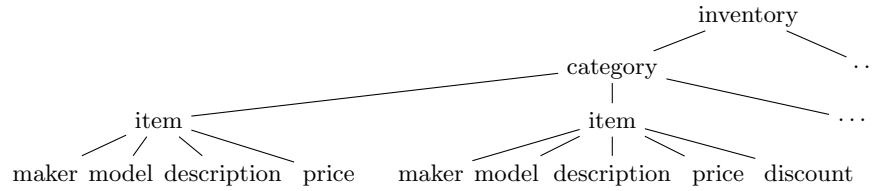


Fig. 67.2: The tag structure of the data in Figure 67.1 as a sibling-ordered tree

We now define the semantics of DTDs. Let  $T = (V, E^c, E^{ns}, \text{lab})$  be a labeled ordered tree. We say that  $T$  *satisfies*  $d$  if  $\text{lab}(\text{Root}(T)) \in S$  and, for every  $v \in V$  with ordered sequence of children  $v_1, \dots, v_n$ , the word  $\text{lab}(v_1) \cdots \text{lab}(v_n)$  is in  $L(r)$ , where  $r = \rho(\text{lab}(v))$ . By  $L(d)$  we denote the set of trees that satisfy  $d$ .

In real-world DTDs,  $\rho$  is usually written as a set of rules rather than a function. We will follow this convention and write  $a \rightarrow r$  if  $\rho(a) = r$ . The *size* of a DTD  $d$  is  $\sum_{a \in \Sigma} \|\rho(a)\|$ .

**Example 67.2**

Throughout the chapter, we will illustrate schemas using a running example as depicted in Figure 67.1. The data is an XML document describing the inventory of a store, which sells items that are classified into categories. Each item has a maker, model, price, and optionally a description and a discount. For space reasons, we only list one category, in which we only list two items. (The reader can imagine further items and categories where we have put the placeholder [...].)

Figure 67.2 depicts the tree structure induced by the nesting of the tags in the XML document. It is at this level of abstraction that DTDs operate.

A DTD describing such data can be defined with  $\Sigma = \{\text{inventory, category, item, maker, model, description, price, discount}\}$  and  $S = \{\text{inventory}\}$ . The function  $\rho$  is defined as follows:

$$\begin{aligned} \text{inventory} &\rightarrow \text{category}^* \\ \text{category} &\rightarrow \text{item}^* \\ \text{item} &\rightarrow \text{maker model description? price discount?} \end{aligned}$$

It describes that inventory should have zero or more category children which, in turn, have zero or more item children. Each item has children labeled maker, model, price and, optionally description and discount, from left to right. The tree in Figure 67.2 satisfies this DTD.

Now assume that we want to express in our running example that each category contains at least one item with a discount child. DTDs cannot express this property<sup>1</sup> since they only associate a single regular expression to item. Schema languages such as XML Schema and Relax NG solve this limitation by extending DTDs with *types*.

**Definition 67.3: Extended DTD**

An *extended DTD (EDTD)* is a tuple

$$D = (\Sigma, \Gamma, \rho, S, \mu),$$

where

- $\Sigma \subseteq \text{Lab}$  is a finite set of *labels*,
- $\Gamma \subseteq \text{Lab}$  is a finite set of *types*,

<sup>1</sup> In Theorem 67.9 we will learn a method to formally prove this.

- $(\Gamma, \rho, S)$  is a DTD over alphabet  $\Gamma$ , and
- $\mu : \Gamma \rightarrow \Sigma$ .

Intuitively, a labeled ordered tree  $T = (V, E^c, E^{ns}, \text{lab})$  satisfies an EDTD if there exists an assignment of types to the nodes in  $V$  such that the typed tree is a derivation tree of the underlying DTD  $(\Gamma, \rho, S)$ .

Formally, for a tree  $T^\Gamma$  using labels in  $\Gamma$ , let us denote by  $\mu(T^\Gamma)$  the tree obtained from  $T^\Gamma$  by replacing each label  $a \in \Gamma$  with  $\mu(a) \in \Sigma$ . Hence,  $\mu(T^\Gamma)$  only uses labels in  $\Sigma$ . We now say that  $T$  *satisfies*  $D$  if there exists a tree  $T^\Gamma$  such that  $T^\Gamma \in L((\Gamma, \rho, S))$  and  $\mu(T^\Gamma) = T$ . We call  $T^\Gamma$  a *witness* for  $T$ . Again, we denote the set of trees satisfying  $D$  by  $L(D)$ .

The *size* of an EDTD  $D$  is the size of the underlying DTD  $(\Gamma, \rho, S)$ . In examples, we denote EDTDs as sets of rules, just as we did before with DTDs.

#### Example 67.4

Continuing Example 67.2 we will define an EDTD that describes inventories in which each category has at least one item with a discount child. We take  $\Sigma$  as before and use  $\Gamma = \{\text{inventory, category, normal, discounted, maker, model, description, price, discount}\}$ . For defining  $\rho$ , we will use a simplified notation, much like how schemas are written in XML Schema or Relax NG: we use  $a[b]$  to denote  $a \in \Sigma$  and  $b \in \Gamma$  with  $\mu(b) = a$ , and we abbreviate  $a[a]$  by  $a$ . The mapping  $\rho$  can now be defined as follows.

```

inventory      → category*
category       → item[normal]* item[discounted]
                  (item[normal] + item[discounted])*
item[normal]    → maker model description? price
item[discounted] → maker model description? price discount

```

The rule for category requires that at least one child is assigned the type discounted. As we can see in the rule for item[discounted], this means that the item must have a child with label discount. The (partial) tree in Figure 67.3 shows how a witness for the tree in Figure 67.2 w.r.t. the present EDTD can be found.

Figure 67.4a shows how the rule for category in Example 67.4 can be defined in Relax NG (with a partial definition of the rule for the type “normal”). In Relax NG, **element** blocks are used to define the elements from  $\Sigma$  (and their content), whereas **define** blocks are used to define types. Since we don’t need a special type for category, we can define it immediately using an **element** block.

In Figure 67.4b, we try to mimic this behavior using XML Schema syntax. However, the code fragment is not syntactically correct because it violates

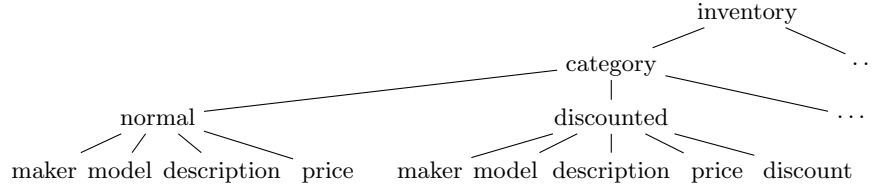


Fig. 67.3: A (partial) witness of the tree in Figure 67.2 for the EDTD in Example 67.4

XML Schema's *Element Declarations Consistent (EDC)* constraint, which forbids the occurrence of different types associated to the same element name in a regular expression  $\rho(t)$ . So, the occurrence of both types normal and discount associated to the same element name item is a violation of this constraint. We formalize this constraint next.

#### Definition 67.5: Single-Type EDTDs

Let  $D = (\Sigma, \Gamma, \rho, S, \mu)$  be an EDTD. A regular expression  $r$  over  $\Gamma$  is *single-type* if it does not contain distinct types  $t_1$  and  $t_2$  with  $\mu(t_1) = \mu(t_2)$ . We say that  $D$  is a *single-type EDTD (stEDTD)* when  $S$  does not contain distinct types  $t_1$  and  $t_2$  with  $\mu(t_1) = \mu(t_2)$  and the regular expression  $\rho(t)$  is single-type for every  $t \in \Gamma$ .

Although single-type EDTDs do not precisely capture the power of XML Schema on the level of labeled ordered trees,<sup>2</sup> they are a very useful theoretical abstraction.

## EDTDs are Equivalent to Tree Automata

We prove that EDTDs have the same expressiveness of tree automata. Furthermore, EDTDs and tree automata can even be converted back and forth in polynomial time. To this end, if  $D$  and  $D'$  are tree automata, DTDs, or EDTDs, we say that  $D$  is *equivalent to*  $D'$  if  $L(D) = L(D')$ .

#### Theorem 67.6

Each EDTD  $D$  can be translated in time  $O(\|D\|)$  to an equivalent NTA. Conversely, each NTA  $A$  over alphabet  $\Sigma$  can be translated in  $O(\|A\| \|\Sigma\|^2)$  time to an equivalent EDTD.

<sup>2</sup> XML Schema has an additional restriction on the regular expressions in  $\rho$ , which we briefly discuss in the bibliographic notes.

```

<element name="category">
  <zeroOrMore><ref name="normal"/></zeroOrMore>
  <ref name="discounted"/>
  <zeroOrMore>
    <choice><ref name="normal"/><ref name="discounted"/></choice>
  </zeroOrMore>
</element>

<define name="normal">
  <element name="item"> ... </element>
</define>

```

(a) Relax NG code fragment corresponding to the EDTD rule for category

```

<element name="category">
  <complexType>
    <sequence>
      <element name="item" type="normal"
        minOccurs="0" maxOccurs="unbounded"/>
      <element name="item" type="discounted"/>
      <choice minOccurs="0" maxOccurs="unbounded">
        <element name="item" type="normal"/>
        <element name="item" type="discounted"/>
      </choice>
    </sequence>
  </complexType>
</element>

<complexType name="normal">
  <element name="item"> ... </element>
</complexType>

```

(b) Code fragment in XML Schema syntax (violating EDC), corresponding to the EDTD of Example 67.4

Fig. 67.4: Fragments of Relax NG and XML Schema code

*Proof.* Let  $D = (\Sigma, \Gamma, \rho, S, \mu)$  be an EDTD. An equivalent NTA  $A = (Q, \Sigma, \delta, F)$  is obtained by taking  $Q = \Gamma$ ,  $F = S$ , and  $\delta(t, \mu(t)) = \rho(t)$  for every type  $t \in \Gamma$ .

Conversely, let  $A = (Q, \Sigma, \delta, F)$  be an NTA. An equivalent EDTD  $D = (\Sigma, \Gamma, \rho, S, \mu)$  is obtained by taking  $\Gamma = Q \times \Sigma$ ,  $S = F \times \Sigma$ , and  $\mu((q, a)) = a$  for each type  $(q, a)$ . Let  $L_{(q,a)}$  denote the language obtained from  $\delta(q, a)$  by replacing, in each word, every symbol  $p \in Q$  by the disjunction  $\sum_{b \in \Sigma} (p, b)$ . Since  $\delta(q, a)$  is a regular language over  $Q$ , the language  $L_{(q,a)}$  is also regular.

For every  $(q, a) \in \Gamma$ , we define  $\rho((q, a))$  to be a regular expression for  $L_{(q,a)}$ . We leave the proofs that  $D$  and  $A$  are equivalent as exercises to the reader.  $\square$

## Complexity of Validation

We now consider the problem of testing if a tree satisfies a given schema. Since we have three types of schema languages (DTD, EDTD, and stEDTD), we define the problem in a general form.

### Problem: $\mathcal{L}$ -Validation

**Input:** A tree  $T$  and a schema  $D$  from the class  $\mathcal{L}$

**Output:** true if  $t \in L(D)$  and false otherwise

### Theorem 67.7

EDTD-Validation is in PTIME.

*Proof.* This is an immediate consequence of Theorems 65.4 and 67.6  $\square$

We have the following corollary since DTDs and stEDTDs are special cases of EDTDs.

**Corollary 67.1.** DTD-Validation and stEDTD-Validation are in PTIME.

## Expressiveness

We now investigate the expressiveness of DTDs, stEDTDs, and EDTDs. The following is immediate from Theorem 67.6

### Theorem 67.8

EDTDs recognize precisely the regular tree languages.

If we want to characterize the languages that are definable by DTDs and stEDTDs, we need the notion of *subtree exchange*. Let  $T = (V, E^c, E^{ns}, \text{lab})$  and  $T' = (V', E'_\downarrow, E'_\rightarrow, \text{lab}')$  be two trees, and let  $u \in V$ . We assume w.l.o.g. that  $V$  and  $V'$  are disjoint (otherwise, the nodes in  $V'$  can be renamed). We denote by  $T[u \leftarrow T']$  the tree obtained by replacing the subtree  $T|_u$  of  $T$  rooted at  $u$  by the tree  $T'$ . (We omit a formal definition, but  $T[u \leftarrow T']$  has none of the nodes in  $T_u$ , all the nodes of  $T'$  and has the edge  $(v, \text{Root}(T'))$  instead of the edge  $(v, u)$ .)

We say that a tree language  $L$  is *closed under node-guarded subtree exchange*, if the following holds. Whenever two trees  $T_1, T_2 \in L$  have nodes  $u_1$  and  $u_2$  respectively, with  $\text{lab}(u_1) = \text{lab}(u_2)$ , then  $T_1[u_1 \leftarrow T_2|_{u_2}] \in L$ .

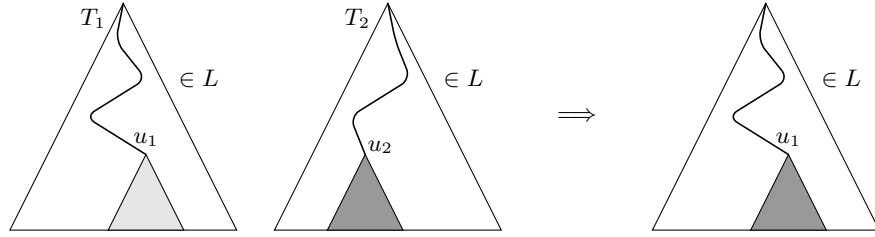


Fig. 67.5: Ancestor-guarded subtree exchange

**Theorem 67.9**

DTDs recognize precisely the regular tree languages that are closed under node-guarded subtree exchange.

*Proof.* Observe that, by definition, every language recognized by a DTD is closed under node-guarded subtree exchange. Conversely, let  $D = (\Sigma, \Gamma, \rho, S, \mu)$  be an EDTD such that  $L(D)$  is closed under node-guarded subtree exchange. We construct a DTD  $d = (\Sigma, \rho_d, S_d)$  as follows. Define  $S_d = \{\mu(t) \mid t \in S\}$ . For a regular expression  $r$  over  $\Gamma$ , denote by  $\mu(r)$  the expression obtained from  $r$  by replacing each symbol  $t$  with  $\mu(t)$ . Then we define  $\rho_d(a)$  as the disjunction of all  $\mu(\rho(t))$  over  $\{t \mid \mu(t) = a\}$ . It is now easy to show that  $L(d) = L(D)$ .  $\square$

We now provide a similar characterization for stEDTDs. For a tree  $T = (V, E^c, E^{ns}, \text{lab})$  and node  $u \in V$ , denote by  $\text{ancestorlab}^T(u)$  the word  $\text{lab}(u_1) \cdots \text{lab}(u_n)$ , where  $u_1 \cdots u_n$  is the path from  $\text{Root}(T)$  to  $u$  in  $T$ . A tree language  $L$  is *closed under ancestor-guarded subtree exchange*, if the following holds. Whenever two trees  $T_1, T_2 \in L$  have nodes  $u_1$  and  $u_2$  respectively, with  $\text{ancestorlab}^{T_1}(u_1) = \text{ancestorlab}^{T_2}(u_2)$ , then  $T_1[u_1 \leftarrow T_2|_{u_2}] \in L$ . We illustrate the concept in Figure 67.5.

**Theorem 67.10**

stEDTDs recognize precisely the regular tree languages that are closed under ancestor-guarded subtree exchange.

*Proof.* It is not difficult to prove that every language definable by an stEDTD is a regular language closed under ancestor-guarded subtree exchange, see Exercise 9.5.

Conversely, let  $D = (\Sigma, \Gamma, \rho, S, \mu)$  be an EDTD such that  $L(D)$  is closed under ancestor-guarded subtree exchange. We will construct a single-type EDTD  $E$  such that  $L(E) = L(D)$ . We will assume w.l.o.g. that  $D$  only uses *useful* types, that is, for every type  $t \in \Gamma$ , there exists a fixed tree



$T_t \in L((\Gamma, \rho, \{t\}))$ . We will make use of the following general property, which immediately follows from the definition of EDTDs:

- (†) If  $T_1, T_2$  are trees in  $L(D)$  with witnesses  $T'_1, T'_2$ , respectively, such that  $u_1$  in  $T_1$  and  $u_2$  in  $T_2$  have the same type in  $T'_1$  and  $T'_2$ , respectively, then  $T_1[u_1 \leftarrow T_2|_{u_2}] \in L(D)$ .

For a word  $w \in \Sigma^*$  and  $a \in \Sigma$ , let  $\text{types}(wa)$  be the set of all types  $t_a$  for which there is a tree  $T \in L(D)$  with witness  $T' \in L((\Gamma, \rho, S))$ , and a node  $v$  such that  $\text{ancestorlab}^T(v) = wa$  and  $\text{lab}^{T'}(v) = t_a$ . For each  $a \in \Sigma$ , let  $\Gamma(D, a)$  be the set of all nonempty sets  $\text{types}(wa)$ , with  $w \in \Sigma^*$ . Clearly, each  $\Gamma(D, a)$  is finite.

We next define  $E = (\Sigma, \Gamma_E, \rho_E, S_E, \mu_E)$ . Its set of types  $\Gamma_E$  is  $\bigcup_{a \in \Sigma} \Gamma(D, a)$  and its set of start types  $S_E$  is  $\{\text{types}(a) \mid a \in \Sigma\}$ . For every type  $\tau \in \Gamma(D, a)$ , set  $\mu_E(\tau) = a$ . The mapping  $\rho_E$  maps each type  $\text{types}(wa)$  to the disjunction of all  $\rho(t_a)$  for  $t_a \in \text{types}(wa)$ , with each  $t_b$  in  $\rho(t_a)$  with  $\mu(t_b) = b$  replaced by  $\text{types}(wab)$ . Notice that, by (†), the definition of the rules of  $\rho_E$  does not depend on the actual choice of  $wa$ .

Clearly,  $E$  is single-type and  $L(D) \subseteq L(E)$ . We show that  $L(E) \subseteq L(D)$ . To this end, let  $T \in L(E)$  and let  $T'$  be a witness. We call a set  $N$  of nodes of  $T$  *well-formed*, if (1) for each node  $v \in N$ , all its ancestors are in  $N$  and (2) if a child  $u$  of a node  $v$  is in  $N$  then all children of  $v$  are in  $N$ . The singleton set  $N_0$  containing the root is well-formed. We say that a tree  $T_1$  *agrees* with  $T$  on a well-formed set  $N_1$  of nodes of  $T_1$ , if  $N_1$  can be mapped to a well-formed set of nodes  $N$  of  $T$  by a mapping  $m$  which respects the child-relationship, the order of siblings, and the labels of nodes. By definition,  $\text{ancestorlab}(v) = \text{ancestorlab}(m(v))$  for every  $v \in N$ .

As the trees in  $L(D)$  and  $L(E)$  have the same possible root labels, there exists a tree  $T_1 \in L(D)$  which agrees with  $T$  on  $N_0$ . To complete the proof, it is sufficient to prove the following.

**Claim 67.2.** *If there exists a tree  $T_1 \in L(D)$  which agrees with  $T = (V, E^c, E^{ns}, \text{lab})$  on a well-formed set  $N$  with  $m(N) \subsetneq V$  then there exists  $T_2 \in L(D)$  which agrees with  $T$  on a well-formed, strict superset of  $N$ .*

For the proof of this claim, let  $T_1$  be as stated and let  $T'_1$  be its witness. Let  $v \in N$  be such that the children of  $m(v)$  are not in  $m(N)$  and let  $wa = \text{ancestorlab}^{T_1}(v) = \text{ancestorlab}^T(m(v))$ .

By construction of  $E$ , the regular expression  $\rho_E(\text{types}(wa))$  is a disjunction  $\sum_{t \in \text{types}(wa)} \rho'(t)$ , where  $\rho'$  is defined as  $\rho$  but, for every  $b \in \Sigma$ , has the type  $\text{types}(wab) \in \Gamma_E$  instead of  $t_b \in \Gamma$  if  $\mu(t_b) = b$ . Let  $t_a$  be such that the children of  $m(v)$  are typed in  $T'$  according to a disjunct  $\rho'(t_a)$ , with  $t_a \in \text{types}(wa)$ . Thus, there is a tree  $T_3 \in L(D)$  with a node  $u$  such that  $\text{ancestorlab}^{T_3}(u) = wa$  and the type of  $u$  is  $t_a$  in the witness  $T'_3$  for  $T_3$ .

Let  $v_1, \dots, v_n$  be the ordered sequence of children of  $m(v)$  in  $T$  and choose, for each  $i \in [1, n]$ , a type  $f(v_i)$  such that  $f(v_1) \cdots f(v_n)$  matches  $\rho(t_a)$ . Let  $T_4$

be obtained from  $T_3$  by (1) removing everything below  $u$ , (2) adding  $v_1, \dots, v_n$  below  $u$ , and (3) adding for each child  $v_i$  the subtree  $T_{f(v_i)}$  which exists because  $f(v_i)$  is a useful type. Clearly,  $T_4 \in L(D)$  by (†). Furthermore, by the ancestor-closed subtree exchange property, the tree  $T_2$  resulting from  $T_1$  by replacing the subtree rooted at  $v$  by the subtree of  $T_4$  rooted at  $u$  is in  $L(D)$ , too.  $\square$

---

## Exercises

- Exercise 9.1.** (a) Prove that the construction in Theorem 65.5 is correct.  
(b) Prove that the algorithm in Theorem 65.6 is correct.  
(c) Prove that the construction in Lemma 65.7 is correct.

**Exercise 9.2.** Prove that Lemma 65.7 also holds for general labeled ordered trees. That is, let us denote by  $\mathcal{T}_\Sigma$  the set of all labeled ordered trees with labels in  $\Sigma$ . Then, for every NTA  $A$ , we can construct an NTA for  $\mathcal{T}_\Sigma - L(A)$  in time  $O(2^{\|A\|})$ .

**Exercise 9.3.** The cores of XML Schema and Relax NG, seen as grammars, are a bit more involved than what we described in Chapter ??, even if we allow arbitrary regular expressions. For XML Schema, it can be more accurately defined as tuples  $(\Sigma, \tau, R, S)$  where  $S \subseteq \Sigma \cup (\Sigma \times \tau)$  and  $R$  maps  $\Sigma \cup \tau$  to regular expressions over  $\Sigma \cup (\Sigma \times \tau)$ . For Relax NG,  $R$  maps  $\Sigma \cup \tau$  to regular expressions over  $\Sigma \cup \tau \cup (\Sigma \times \tau)$ . Prove that this is equally expressive.

**Exercise 9.4.** Prove that the constructions in Theorem 67.6 are correct, that is, prove that the EDTD  $D$  and automaton  $A$  are indeed equivalent in both constructions.

**Exercise 9.5.** Prove that every language definable by an stEDTD is a regular language closed under ancestor-guarded subtree exchange.

*Hint:* Use a partial function  $f : \Sigma^+ \rightarrow \Gamma$  that maps each “ancestor-string” to the relevant type, i.e.,

- $f(a) = t$ , where  $t \in S$  is unique such that  $\mu(t) = a$  and
- $f(wa) = t$ , where  $f(w) = t'$  and  $t$  is the unique type occurring in  $\rho(t')$  with  $\mu(t) = a$ .

Some ideas for exercises:

- Can we postpone canonical models to here?

- Prove that minimal patterns are not unique in a strong sense?
- Prove that you cannot define “trees for which the deepest node is on even depth” with a tree automaton.
- Prove that you cannot define “the tree has a unique  $b$ -labeled leaf” with XSD (single-type EDTD).
- Turn a regular expression into a deterministic one?

**Exercise 9.6.** Prove that, given a TPQ  $p$  and EDTD  $D$ , you can construct in polynomial time a Boolean TPQ  $p_b$  and EDTD  $D_b$  such that the following are equivalent:

- (1) There exists a tree  $T \in L(D)$  such that  $p(T)$  is not empty.
- (2) There exists a tree  $T \in L(D_b)$  such that  $T \models p_b$ .

*Hint: add children with unique labels to the output nodes of  $p$ .*

**Exercise 9.7.** When reducing a DTD, it is important that steps 1 and 2 in the proof of Lemma 68.1 are performed in that order. Show that, if the order is reversed, the result is not always a reduced DTD.

**Exercise 9.8.** Prove that the NP-hardness of Theorem 68.2 also holds for DTDs and tree patterns without descendant edges or wildcards. *Hint: Try to “flatten” the trees defined by the DTD.*

---

## Bibliographic Comments

To be done.

---

## References

1. S. Abiteboul, R. Hull, and V. Vianu. *Foundations of Databases*. Addison-Wesley, 1995.
2. S. Arora and B. Barak. *Computational Complexity - A Modern Approach*. Cambridge University Press, 2009.
3. P. Bourhis, J. L. Reutter, and D. Vrgoc. JSON: data model and query languages. *Inf. Syst.*, 89:101478, 2020.
4. D. Calvanese, G. De Giacomo, M. Lenzerini, and M. Y. Vardi. Containment of conjunctive regular path queries with inverse. In *KR 2000, Principles of Knowledge Representation and Reasoning Proceedings of the Seventh International Conference, Breckenridge, Colorado, USA, April 11-15, 2000*, pages 176–185, 2000.
5. K. Etessami, M. Y. Vardi, and T. Wilke. First-order logic with two variables and unary temporal logic. *Inf. Comput.*, 179(2):279–295, 2002.
6. H. Garcia-Molina, J. D. Ullman, and J. Widom. *Database Systems - the Complete Book*. Pearson Education, 2009.
7. E. Grädel, P. Kolaitis, L. Libkin, M. Marx, J. Spencer, M. Vardi, and S. Weinstein. *Finite Model Theory and its Applications*. Springer, 2008.
8. P. R. Halmos. *Measure Theory*. Springer, 1974.
9. J. E. Hopcroft, R. Motwani, and J. D. Ullman. *Introduction to Automata Theory, Languages, and Computation*. Addison-Wesley, 2003.
10. N. Immerman. *Descriptive Complexity*. Springer, 1999.
11. D. Kozen. *Automata and Computability*. Springer, 1997.
12. L. Libkin. *Elements of Finite Model Theory*. Springer, 2004.
13. C. H. Papadimitriou. *Computational Complexity*. Addison-Wesley, 1994.
14. R. Ramakrishnan and J. Gehrke. *Database Management Systems*. McGraw-Hill, 2003.
15. K. H. Rosen. *Discrete Mathematics and its Applications*. McGraw-Hill, 2006.
16. A. Silberschatz, H. F. Korth, and S. Sudarshan. *Database System Concepts*. McGraw-Hill Book Company, 2005.

17. M. Sipser. *Introduction to the Theory of Computation*. PWS Publishing Company, 1997.
18. R. E. Tarjan and M. Yannakakis. Simple linear-time algorithms to test chordality of graphs, test acyclicity of hypergraphs, and selectively reduce acyclic hypergraphs. *SIAM J. Comput.*, 13(3):566–579, 1984.
19. D. Vrgoč. *Querying graphs with data*. PhD thesis, University of Edinburgh, 2014.
20. I. Wegener. *Complexity Theory*. Springer, 2005.

**Appendix: Theory of Computation**





## A

---

### Big-O Notation

We write  $\mathbb{R}_0^+$  for the set of non-negative real numbers, and  $\mathbb{R}^+$  for the set of positive real numbers. We typically measure the performance of an algorithm, that is, the number of basic operations it performs, as a function of its input length. In other words, the performance of an algorithm can be captured by a function  $f : \mathbb{N} \rightarrow \mathbb{R}_0^+$  such that  $f(n)$  is the maximum number of basic operations that the algorithm performs on inputs of length  $n$ . However, since  $f$  may heavily depend on the details of the definition of basic operations, we usually concentrate on the overall and asymptotic behaviour of the algorithm. This is achieved via the well-known notion of *big-O notation*.

The big-O notation is typically defined for single variable functions such as  $f$  above. However, in the database setting, where the input to key problems usually consists of several different components, we generally have to deal with multiple variable functions. For example, the performance of a query evaluation algorithm, where the input consists of two distinct components, the *database* and the *query*, can be captured by a function  $f : \mathbb{N}^2 \rightarrow \mathbb{R}_0^+$  such that  $f(n, m)$  is the maximum number of basic operations that the algorithm performs on databases of size  $n$  and queries of size  $m$ . The notion of big-O notation for multiple variable functions follows:

#### Definition 1.1: Big-O Notation

Let  $f, g : \mathbb{N}^\ell \rightarrow \mathbb{R}_0^+$ , where  $\ell \geq 1$ . We say that

$$f(x_1, \dots, x_\ell) \text{ is in } O(g(x_1, \dots, x_\ell))$$

if there exist  $k \in \mathbb{R}^+$  and  $n_0 \in \mathbb{N}$  such that, for every  $(x_1, \dots, x_\ell)$  with  $x_i \geq n_0$  for some  $i \in [\ell]$ , we have  $f(x_1, \dots, x_\ell) \leq k \cdot g(x_1, \dots, x_\ell)$ .

Notice that when  $\ell = 1$ , i.e.,  $f, g$  are single variables function, Definition 1.1 coincides with the standard big-O notation for single variable functions.

Sometimes we will use a variant of the big-O notation that hides logarithmic factors. This variant will be denoted with  $\tilde{O}$  instead of the usual letter  $O$ . We define this notation as follows.

**Definition 1.2: Big-O Notation Without Logarithmic Factors**

Let  $f, g : \mathbb{N}^\ell \rightarrow \mathbb{R}_0^+$ , where  $\ell \geq 1$ . We say that

$$f(x_1, \dots, x_\ell) \text{ is in } \tilde{O}(g(x_1, \dots, x_\ell))$$

if  $f(x_1, \dots, x_\ell)$  is in  $O\left(g(x_1, \dots, x_\ell) \cdot \log(g(x_1, \dots, x_\ell))\right)$ .

## B

---

### Turing Machines and Complexity Classes

Many results in this book will provide bounds on computational resources (time and space), or key database problems such as query evaluation. These are often formulated in terms of membership in, or completeness for, complexity classes. Those, in turn, are defined using the basic model of computation, that is, Turing Machines. We now briefly recall basic concepts related to Turing Machines and complexity classes. For more details, the reader can consult standard textbooks on computability theory and computational complexity.

#### Turing Machines

Turing Machines can work in two modes: either as *acceptors*, for deciding whether an input string belongs to a given language (in which case we speak of *decision problems*), or as computational devices that compute the value of a function applied to its input. When a Turing Machine works as an acceptor, it typically contains a read-write tape, a model of computation that is convenient for defining time complexity classes, or a read-only input tape and a read-write working tape, a model that is convenient for defining space complexity classes. When a Turing Machine works as a computational device, it typically contains a read-only tape where the input is placed, a read-write working tape, and a write-only tape where the output computed by the Turing Machine is placed.

##### *Turing Machines as Acceptors*

We start with the definition of deterministic Turing Machines.

**Definition 2.1: Deterministic Turing Machine**

A (*deterministic*) *Turing Machine* (TM) is defined as a tuple

$$M = (Q, \Sigma, \delta, s),$$

where

- $Q$  is a finite set of states, including the *accepting state* “yes”, and the *rejecting state* “no”,
- $\Sigma$  is a finite set of input symbols, called the *alphabet* of  $M$ , including the symbols  $\sqcup$  (the *blank symbol*) and  $\triangleright$  (the *left marker*),
- $\delta : (Q - \{\text{“yes”}, \text{“no”}\}) \times \Sigma \rightarrow Q \times \Sigma \times \{\rightarrow, \leftarrow, -\}$  is the *transition function* of  $M$ , and
- $s \in Q$  is the *start state* of  $M$ .

Accepting and rejecting states are needed for decision problems: they determine whether the input belongs to the language or not. Notice that, according to  $\delta$ , the accepting and rejecting states do not have outgoing transitions.

A *configuration* of a TM  $M = (Q, \Sigma, \delta, s)$  is a tuple

$$c = (q, u, v),$$

where  $q \in Q$ , and  $u, v$  are words in  $\Sigma^*$  with  $u$  being always non-empty. If  $M$  is in configuration  $c$ , then the tape has content  $uv$  and the head is reading the last symbol of  $u$ . We use left markers, which means that  $u$  always starts with  $\triangleright$ . Moreover, the transition function  $\delta$  is restricted in such a way that  $\triangleright$  occurs exactly once in  $uv$ , and always as the first symbol of  $u$ .

Assume now that  $M$  is in a configuration  $c = (q, ua, v)$ , where  $q \in Q - \{\text{“yes”}, \text{“no”}\}$ ,  $a \in \Sigma$  and  $u, v \in \Sigma^*$ , and assume that  $\delta(q, a) = (q', b, \text{dir})$ , where  $\text{dir} \in \{\rightarrow, \leftarrow, -\}$ . Then, in one step,  $M$  enters the configuration  $c' = (q', u', v')$ , where  $u', v'$  is obtained from  $ua, v$  by replacing  $a$  with  $b$  and moving the head one step in the direction  $\text{dir}$ . By moving the head in the direction “-” we mean that the head stays in its place. Furthermore, the head cannot move left of the  $\triangleright$  symbol (the transition function  $\delta$  is restricted in such a way that this cannot happen: if  $\delta(q, \triangleright) = (q', a, \text{dir})$ , then  $a = \triangleright$  and  $\text{dir} \neq \leftarrow$ ). For example, if  $c = (q, \triangleright 01, 100)$  and  $\delta(q, 1) = (q', 0, \leftarrow)$ , then  $c' = (q', \triangleright 0, 0100)$ . In this case, we write  $c \rightarrow_M c'$ , and we also write  $c \rightarrow_M^m c'$  if  $c'$  can be reached from  $c$  in  $m$  steps, and  $c \rightarrow_M^* c'$  if  $c \rightarrow_M^m c'$  for some  $m \geq 0$  (we assume that  $c \rightarrow_M^0 c$ ). Finally, if  $v = \varepsilon$  and  $\text{dir} = \rightarrow$ , then we insert an additional  $\sqcup$ -symbol in our configuration, that is  $u' = ub\sqcup$  and  $v' = \varepsilon$ .

A TM  $M$  receives an input word  $w = a_1 \cdots a_n$ , where  $n \geq 0$  and  $a_i \in \Sigma - \{\sqcup, \triangleright\}$  for each  $i \in [n]$ . The *start configuration* of  $M$  on input  $w$  is  $sc(w) = (s, \triangleright, w)$ . We call a configuration  $c$  *accepting* if its state is “yes”, and *rejecting* if its state is “no”. The TM  $M$  *accepts* (respectively, *rejects*) input  $w$  if  $sc(w) \rightarrow_M^* c$  for some accepting (respectively, rejecting) configuration  $c$ .

### Nondeterministic Turing Machines as Acceptors

We also use nondeterministic Turing Machines as acceptors, which are defined similarly to deterministic ones, but with the key difference that the a state-symbol pair has more than one outgoing transitions.

**Definition 2.2: Nondeterministic Turing Machine**

A *nondeterministic Turing Machine* (NTM) is defined as a tuple

$$M = (Q, \Sigma, \delta, s),$$

where

- $Q$  is a finite set of states, including the *accepting state* “yes”, and the *rejecting state* “no”,
- $\Sigma$  is a finite set of input symbols, called the *alphabet* of  $M$ , including the symbols  $\sqcup$  (the *blank symbol*) and  $\triangleright$  (the *left marker*),
- $\delta : (Q - \{\text{“yes”}, \text{“no”}\}) \times \Sigma \rightarrow \mathcal{P}(Q \times \Sigma \times \{\rightarrow, \leftarrow, -\})$  is the *transition function* of  $M$ , and
- $s \in Q$  is the *start state* of  $M$ .

Observe that for a given configuration  $c = (q, ua, v)$ , where  $q \in Q - \{\text{“yes”}, \text{“no”}\}$ ,  $a \in \Sigma$  and  $u \in \Sigma^*$ , several alternatives  $(q', b, \text{dir})$  can belong to  $\delta(q, a)$ , each one of which generates a successor configuration  $c'$  as in the case of (deterministic) TMs. If  $c'$  is a possible successor configuration of  $c$ , then we write  $c \rightarrow_M c'$ . Moreover, we write  $c \xrightarrow{m}_M c'$  if there exists a sequence of configurations  $c_1, \dots, c_{m-1}$  such that  $c \rightarrow_M c_1, c_1 \rightarrow_M c_2, \dots, c_{m-1} \rightarrow_M c'$ . In this case, notice that it is possible that  $c \xrightarrow{m}_M c'$  and  $c \xrightarrow{n}_M c'$  with  $m \neq n$ . Moreover, we write  $c \xrightarrow{*}_M c'$  if there exists  $m \geq 0$  such that  $c \xrightarrow{m}_M c'$  (again, we assume that  $c \xrightarrow{0}_M c$ ).

Given an input word  $w$  for a NTM  $M$ , the start configuration  $sc(w)$  of  $M$ , and accepting and rejecting configurations of  $M$ , are defined as in the deterministic case. Moreover,  $M$  accepts input  $w$  if there exists an accepting configuration  $c$  such that  $sc(w) \xrightarrow{*}_M c$ , and  $M$  rejects  $w$  otherwise (i.e.,  $M$  rejects  $w$  if there is no accepting configuration  $c$  such that  $sc(w) \xrightarrow{*}_M c$ ).

*2-Tape Turing Machines as Acceptors*

We now define Turing Machines that, apart from a read-write working tape, they also have a read-only input tape.

**Definition 2.3: 2-Tape Deterministic Turing Machine**

A 2-tape (*deterministic*) Turing Machine (2-TM) is defined as a tuple

$$M = (Q, \Sigma, \delta, s),$$

where

- $Q$  is a finite set of states, including the *accepting state* “yes”, and the *rejecting state* “no”,
- $\Sigma$  is a finite set of input symbols, called the *alphabet* of  $M$ , including the symbols  $\sqcup$  (the *blank symbol*) and  $\triangleright$  (the *left marker*),
- $\delta : (Q - \{\text{“yes”}, \text{“no”}\}) \times \Sigma \times \Sigma \rightarrow Q \times \{\rightarrow, \leftarrow, -\} \times \Sigma \times \{\rightarrow, \leftarrow, -\}$  is the *transition function* of  $M$ , and
- $s \in Q$  is the *start state* of  $M$ .

A *configuration* of a 2-TM is a tuple

$$c = (q, u_1, v_1, u_2, v_2),$$

where  $q \in Q$  and, for every  $i \in \{1, 2\}$ , we have that  $u_i, v_i \in \Sigma^*$  and  $u_i$  is not empty. If  $M$  is in configuration  $c$ , then the input tape has content  $u_1v_1$  and the head of this tape is reading the last symbol of  $u_1$ , while the working tape has content  $u_2v_2$  and the head of this tape is reading the last symbol of  $u_2$ . We use left markers, which means that  $u_i$  always starts with  $\triangleright$ . Besides, the transition function  $\delta$  is restricted in such a way that  $\triangleright$  occurs exactly once in  $u_iv_i$ , and always as the first symbol of  $u_i$ .

Assume that  $M$  is in a configuration  $c = (q, u_1a_1, v_1, u_2a_2, v_2)$ , where  $q \in Q - \{\text{“yes”}, \text{“no”}\}$ ,  $a_1, a_2 \in \Sigma$  and  $u_1, v_1, u_2, v_2 \in \Sigma^*$ , and assume that  $\delta(q, a_1, a_2) = (q', \text{dir}_1, b, \text{dir}_2)$ , where  $\text{dir}_i$  is a direction, i.e., one of  $\{\rightarrow, \leftarrow, -\}$ . Then in one step  $M$  enters configuration  $c' = (q', u'_1, v'_1, u'_2, v'_2)$ , where  $u'_1, v'_1$  is obtained from  $u_1a_1, v_1$  by moving the head one step in the direction  $\text{dir}_1$ , and  $u'_2, v'_2$  is obtained from  $u_2a_2, v_2$  by replacing  $a_2$  with  $b$  and moving the head one step in the direction  $\text{dir}_2$ . Recall that by moving the head in the direction “-” we mean that the head stays in its place. Furthermore, the head cannot move left of the  $\triangleright$  symbol (again, the transition function  $\delta$  is restricted in such a way that this cannot happen). For example, if  $c = (q, \triangleright 01, 100, \triangleright, \varepsilon)$  and  $\delta(q, 1, \triangleright) = (q', \leftarrow, \triangleright, \rightarrow)$ , then  $c' = (q', \triangleright 0, 1100, \triangleright, \varepsilon)$ . In this case, we write  $c \rightarrow_M c'$ . We also write  $c \xrightarrow{m}_M c'$  if  $c'$  can be reached from  $c$  in  $m$  steps, and  $c \xrightarrow{*}_M c'$  if  $c \xrightarrow{m}_M c'$  for some  $m \geq 0$  (we assume that  $c \xrightarrow{0}_M c$ ).

A 2-TM  $M$  receives an input word  $w = a_1 \cdots a_n$ , where  $n \geq 0$  and  $a_i \in \Sigma - \{\sqcup, \triangleright\}$  for each  $i \in [n]$ . The *start configuration* of  $M$  on input  $w$  is  $sc(w) = (s, \triangleright, w, \triangleright, \sqcup)$ . We call a configuration  $c$  *accepting* if its state is “yes”, and

*rejecting* if its state is “no”. The TM  $M$  *accepts* (respectively, *rejects*) input  $w$  if  $sc(w) \rightarrow_M^* c$  for some accepting (respectively, rejecting) configuration  $c$ .

### 2-Tape Nondeterministic Turing Machines as Acceptors

As for TMs, we also have the nondeterministic version of 2-TMs.

#### Definition 2.4: 2-Tape Nondeterministic Turing Machine

A 2-Tape Nondeterministic Turing Machine (2-NTM) is a tuple

$$M = (Q, \Sigma, \delta, s),$$

where

- $Q$  is a finite set of states, including the *accepting state* “yes”, and the *rejecting state* “no”,
- $\Sigma$  is a finite set of input symbols, called the *alphabet* of  $M$ , including the symbols  $\sqcup$  (the *blank symbol*) and  $\triangleright$  (the *left marker*),
- $\delta : (Q - \{\text{“yes”}, \text{“no”}\}) \times \Sigma \times \Sigma \rightarrow \mathcal{P}(Q \times \{\rightarrow, \leftarrow, -\} \times \Sigma \times \{\rightarrow, \leftarrow, -\})$  is the *transition function* of  $M$ , and
- $s \in Q$  is the *start state* of  $M$ .

It is clear that, for a configuration  $c = (q, u_1, a_1v_1, u_2, a_2v_2)$ , where  $q \in Q - \{\text{“yes”}, \text{“no”}\}$ ,  $a_1, a_2 \in \Sigma$  and  $v_1, v_2 \in \Sigma^*$ , several alternatives  $(q', \text{dir}_1, b, \text{dir}_2)$  can belong to  $\delta(q, a_1, a_2)$ , each one of which generates a successor configuration  $c'$  as in the case of 2-TMs. If  $c'$  is a possible successor configuration of  $c$ , then we write  $c \rightarrow_M c'$ . Moreover, we write  $c \rightarrow_M^m c'$  if there exists a sequence of configurations  $c_1, \dots, c_{m-1}$  such that  $c \rightarrow_M c_1, c_1 \rightarrow_M c_2, \dots, c_{m-1} \rightarrow_M c'$ . In this case, notice that it is possible that  $c \rightarrow_M^m c'$  and  $c \rightarrow_M^n c'$  with  $m \neq n$ . Moreover, we write  $c \rightarrow_M^* c'$  if there exists  $m \geq 0$  such that  $c \rightarrow_M^m c'$  (again, we assume that  $c \rightarrow_M^0 c$ ).

Given an input word  $w$  for a 2-NTM  $M$ , the start configuration  $sc(w)$  of  $M$ , and accepting and rejecting configurations of  $M$ , are defined as in the deterministic case. Moreover,  $M$  accepts input  $w$  if there exists an accepting configuration  $c$  such that  $sc(w) \rightarrow_M^* c$ , and  $M$  rejects  $w$  otherwise (i.e.,  $M$  rejects  $w$  if there is no accepting configuration  $c$  such that  $sc(w) \rightarrow_M^* c$ ).

### Turing Machines as Computational Devices

If a 2-TM acts not as a language acceptor but rather as a device for computing a function  $f$ , then a write-only output tape is added and the states “yes” and “no” are replaced with a *halting state* “halt”; once the computation enters the halting state, the output tape contains the value  $f(w)$  for the input  $w$ .



**Definition 2.5: Turing Machine with Output**

A *Turing Machine with output* (TMO) is a tuple

$$M = (Q, \Sigma, \delta, s),$$

where

- $Q$  is a finite set of states, including the *halting state* “halt”,
- $\Sigma$  is a finite set of input symbols, called the *alphabet* of  $M$ , including the symbols  $\sqcup$  (the *blank symbol*) and  $\triangleright$  (the *left marker*),
- $\delta : (Q - \{\text{“halt”}\}) \times \Sigma \times \Sigma \rightarrow Q \times \{\rightarrow, \leftarrow, -\} \times \Sigma \times \{\rightarrow, \leftarrow, -\} \times \Sigma$  is the *transition function* of  $M$ , and
- $s \in Q$  is the *start state* of  $M$ .

If  $\delta(q, a_1, a_2) = (q', \text{dir}_1, b, \text{dir}_2, c)$ , then  $q', \text{dir}_1, b, \text{dir}_2$  are used exactly as in the case of a 2-TM accepting a language. Moreover, if  $c \neq \sqcup$ , then  $c$  is written on the output tape and the head of this tape is moved one position to the right; otherwise, no changes are made on this tape. The start configuration of a TMO  $M$  on input  $w$  is  $sc(w) = (s, \triangleright, w, \triangleright, \varepsilon, \triangleright, \varepsilon)$ . The output of  $M$  on input  $w$  is the word  $u$  such that  $sc(w) \rightarrow_M^* (\text{“halt”}, u_1, v_1, u_2, v_2, \triangleright u, \varepsilon)$ .

## Complexity Classes

We proceed to introduce some central complexity classes that are used in this book. Recall that  $\mathbb{R}_0^+$  is the set of non-negative real numbers. Given a function  $f : \mathbb{N} \rightarrow \mathbb{R}_0^+$ , a TM (respectively, NTM)  $M$  is said to run in *time*  $f(n)$  if, for every input  $w$  and configuration  $c$ ,  $sc(w) \rightarrow_M^m c$  implies  $m \leq f(|w|)$ .<sup>1</sup> We further say that  $M$  *decides* a language  $L$  if  $M$  accepts every word in  $L$  and rejects every word not in  $L$ . Notice that this implies that  $M$ ’s computation is finite on every input. We define the classes of decision problems

$$\begin{aligned} \text{TIME}(f(n)) = \{L \mid \text{there exists a TM that decides } L \\ \text{and runs in time } f(n)\} \end{aligned}$$

and

$$\begin{aligned} \text{NTIME}(f(n)) = \{L \mid \text{there exists an NTM that decides } L \\ \text{and runs in time } f(n)\}. \end{aligned}$$

We use  $\text{TIME}(O(f(n)))$  for the union of all  $\text{TIME}(g(n))$  where  $g(n) \in O(f(n))$ . Furthermore, we use the following time complexity classes in this book:

<sup>1</sup> The running time of a TMO is defined in the same way.

**Definition 2.6: Time Complexity Classes**

$$\begin{aligned}
\text{PTIME} &= \bigcup_{k \in \mathbb{N}} \text{TIME}(n^k) & \text{NP} &= \bigcup_{k \in \mathbb{N}} \text{NTIME}(n^k) \\
\text{EXPTIME} &= \bigcup_{k \in \mathbb{N}} \text{TIME}(2^{n^k}) & \text{NEXPTIME} &= \bigcup_{k \in \mathbb{N}} \text{NTIME}(2^{n^k}) \\
2\text{EXPTIME} &= \bigcup_{k \in \mathbb{N}} \text{TIME}(2^{2^{n^k}})
\end{aligned}$$

Given a function  $f : \mathbb{N} \rightarrow \mathbb{R}_0^+$ , a 2-TM (respectively, 2-NTM)  $M$  is said to run in *space*  $f(n)$  if, for every input  $w$  and configuration  $c = (q, u_1, v_1, u_2, v_2)$ ,  $sc(w) \rightarrow_M^* c$  implies  $|u_2 v_2| \leq f(|w|)$ .<sup>2</sup> We say that  $M$  *decides* a language  $L$  if  $M$  accepts every word in  $L$  and rejects every word not in  $L$ . We define the classes of decision problems

$$\begin{aligned}
\text{SPACE}(f(n)) &= \{L \mid \text{there exists a 2-TM that decides } L \\
&\quad \text{and runs in space } f(n)\}
\end{aligned}$$

and

$$\begin{aligned}
\text{NSPACE}(f(n)) &= \{L \mid \text{there exists a 2-NTM that decides } L \\
&\quad \text{and runs in space } f(n)\}.
\end{aligned}$$

We write  $\text{SPACE}(O(f(n)))$  for the union of all  $\text{SPACE}(g(n))$ , where  $g(n) \in O(f(n))$ . We further use the following space complexity classes in this book:

**Definition 2.7: Space Complexity Classes**

$$\begin{aligned}
\text{DLOGSPACE} &= \text{SPACE}(\log n) & \text{NLOGSPACE} &= \text{NSPACE}(\log n) \\
\text{PSPACE} &= \bigcup_{k \in \mathbb{N}} \text{SPACE}(n^k) & \text{NPSPACE} &= \bigcup_{k \in \mathbb{N}} \text{NSPACE}(n^k) \\
\text{EXPSPACE} &= \bigcup_{k \in \mathbb{N}} \text{SPACE}(2^{n^k}) & \text{NEXPSPACE} &= \bigcup_{k \in \mathbb{N}} \text{NSPACE}(2^{n^k})
\end{aligned}$$

At this point, let us stress that we can always assume that the computation of a space-bounded 2-TM  $M$  is finite on every input word. Intuitively, since the space that  $M$  uses is bounded, the number of different configurations in which  $M$  can be is also bounded. Therefore, by maintaining a counter that “counts” the steps of  $M$ , we can guarantee that  $M$  will never fall in an unnecessarily long computation, which in turn allows us to assume that the computation of  $M$  is finite. Further details on this assumption can be found in any standard textbook on computational complexity.

For a complexity class  $\mathcal{C}$ , the class  $\text{co}\mathcal{C}$  is defined as the set of complements of the problems in  $\mathcal{C}$ , that is,  $\text{co}\mathcal{C} = \{\Sigma^* - L \mid L \in \mathcal{C}\}$ . It is known that

<sup>2</sup> The running space of a TMO is defined without considering the output tape. More precisely, for every input  $w$  and configuration  $c = (q, u_1, v_1, u_2, v_2, u_3, v_3)$ ,  $sc(w) \rightarrow_M^* c$  implies  $|u_2 v_2| \leq f(|w|)$ .

$$\begin{aligned} \text{DLOGSPACE} \subseteq \text{NLOGSPACE} \subseteq \text{PTIME} \subseteq \text{NP} \subseteq \text{PSPACE} = \text{NPSPACE} \\ \subseteq \text{EXPTIME} \subseteq \text{NEXPTIME} \subseteq \text{EXPSPACE} = \text{NEXPSPACE} \subseteq 2\text{EXPTIME} \end{aligned}$$

$$\text{PTIME} \subsetneq \text{EXPTIME} \subsetneq 2\text{EXPTIME}$$

$$\text{NP} \subsetneq \text{NEXPTIME}$$

and that

$$\text{NLOGSPACE} = \text{CONLOGSPACE} \subsetneq \text{PSPACE} \subsetneq \text{EXPSPACE}$$

However, it is still not known whether PTIME (and in fact DLOGSPACE) is properly contained in NP, whether PTIME is properly contained in PSPACE, and whether NP equals CONP.

Key concepts related to complexity classes are *reductions* between problems, and *hardness* and *completeness* of problems. For precise definitions the reader can consult any complexity theory textbook. A reduction between languages  $L$  and  $L'$  over an alphabet  $\Sigma$  is a function  $f : \Sigma^* \rightarrow \Sigma^*$  such that  $w \in L$  if and only if  $f(w) \in L'$ , for every  $w \in \Sigma^*$ . Let  $\mathcal{C}$  be one of the complexity classes introduced above such that  $\text{NP} \subseteq \mathcal{C}$  or  $\text{CONP} \subseteq \mathcal{C}$ . A *problem*, i.e., a language  $L$ , is *hard* for  $\mathcal{C}$ , or  $\mathcal{C}$ -hard, if every problem  $L' \in \mathcal{C}$  is reducible to  $L$  via a reduction that is computable in polynomial time. If  $L$  is also in  $\mathcal{C}$ , then it is *complete* for  $\mathcal{C}$ , or  $\mathcal{C}$ -complete. For the complexity classes NLOGSPACE and PTIME, the notions of hardness and completeness are defined in the same way, but with the crucial difference that we rely on reductions that are computable in deterministic logarithmic space. This is because a reduction is meaningful only within a class that is computationally stronger than the reduction.<sup>3</sup>

We say that a decision problem is *tractable* if it is in PTIME. As such, problems that are hard for EXPTIME are *provably intractable*. We call problems that are hard for NP or CONP *presumably intractable* (if we cannot make a stronger case and prove that they are not in PTIME).

The most fundamental problem that is presumably intractable is the satisfiability problem of Boolean formulae. A *Boolean formula* is defined as follows:

- a variable  $x \in \text{Var}$  is a Boolean formula, and
- if  $\varphi_1$  and  $\varphi_2$  are Boolean formulae, then  $(\varphi_1 \wedge \varphi_2)$ ,  $(\varphi_1 \vee \varphi_2)$ , and  $(\neg \varphi_1)$  are Boolean formulae.

To define the semantics of such Boolean formulae, we need the notion of truth assignment. A *truth assignment* for a set of variables  $V$  is a function  $f : V \rightarrow \{\text{true}, \text{false}\}$ . Consider a Boolean formula  $\varphi$ , and a truth assignment  $f$  for the set of variables in  $\varphi$ . We define when  $f$  *satisfies*  $\varphi$ , written  $f \models \varphi$ :

<sup>3</sup> We could also define hardness for DLOGSPACE by using reductions that can be computed via a computation even more restrictive than deterministic logarithmic space, but this is not needed for the purposes of this book.

- If  $\varphi$  is a variable  $x$ , then  $f \models \varphi$  if  $f(x) = \mathbf{true}$ .
- If  $\varphi = (\varphi_1 \wedge \varphi_2)$ , then  $f \models \varphi$  if  $f \models \varphi_1$  and  $f \models \varphi_2$ .
- If  $\varphi = (\varphi_1 \vee \varphi_2)$ , then  $f \models \varphi$  if  $f \models \varphi_1$  or  $f \models \varphi_2$ .
- If  $\varphi = (\neg\psi)$ , then  $f \models \varphi$  if  $f \models \psi$  does not hold.

We say that  $\varphi$  is *satisfiable* if there exists a truth assignment  $f$  for the set of variables in  $\varphi$  such that  $f \models \varphi$ . The *Boolean satisfiability* problem or SAT, which is known to be an NP-complete problem, is defined as follows.

**Problem: SAT**

**Input:** A Boolean formula  $\varphi$

**Output:** **true** if  $\varphi$  is satisfiable, and **false** otherwise

It is actually the first problem that was proven to be NP-complete, a result known as Cook-Levin Theorem that goes back to the 1970s.

A generalization of SAT is the satisfiability problem of quantified Boolean formulae. For a Boolean formula  $\varphi$  and a tuple of variables  $\bar{x}$ , we denote by  $\varphi\langle\bar{x}\rangle$  the fact that  $\varphi$  uses precisely the variables in  $\bar{x}$ . A *quantified Boolean formula*  $\psi$  is an expression of the form

$$Q_1\bar{x}_1 Q_2\bar{x}_2 \cdots Q_n\bar{x}_n \varphi\langle\bar{x}_1, \dots, \bar{x}_n\rangle,$$

where, for each  $i \in [n]$ ,  $Q_i$  is either  $\exists$  or  $\forall$ , and, for each  $i \in [n-1]$ ,  $Q_i = \exists$  implies  $Q_{i+1} = \forall$  and  $Q_i = \forall$  implies  $Q_{i+1} = \exists$ . Assuming that  $Q_1 = \exists$ , we say that  $\psi$  is *satisfiable* if there exists a truth assignment for  $\bar{x}_1$  such that for every truth assignment for  $\bar{x}_2$  there exists a truth assignment for  $\bar{x}_3$ , and so on up to  $\bar{x}_n$ , such that the overall truth assignment satisfies  $\psi$ . Analogously, we can define when  $\psi$  is satisfiable in the case  $Q_1 = \forall$ . The *quantified satisfiability* problem or QSAT, also known under the name *quantified Boolean formula* or QBF, which is the canonical PSPACE-complete problem, is defined as follows:

**Problem: QSAT**

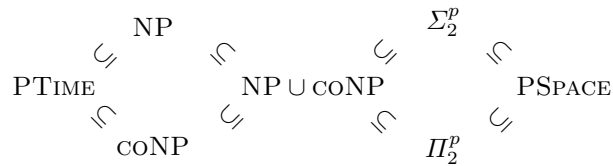
**Input:** A quantified Boolean formula  $\psi$

**Output:** **true** if  $\psi$  is satisfiable, and **false** otherwise

Notice that SAT is the special case of QSAT where  $\psi$  is of the form  $\exists\bar{x}\varphi\langle\bar{x}\rangle$ . Two special cases of QSAT will be particularly important for this book, namely the ones with exactly one quantifier alternation:

**Problem:  $\exists\forall\text{QSAT}$** **Input:** A quantified Boolean formula  $\psi = \exists \bar{x}_1 \forall \bar{x}_2 \varphi(\bar{x}_1, \bar{x}_2)$ **Output:** **true** if  $\psi$  is satisfiable, and **false** otherwise**Problem:  $\forall\exists\text{QSAT}$** **Input:** A quantified Boolean formula  $\psi = \forall \bar{x}_1 \exists \bar{x}_2 \varphi(\bar{x}_1, \bar{x}_2)$ **Output:** **true** if  $\psi$  is satisfiable, and **false** otherwise

We define  $\Sigma_2^P$  as the class of decision problems reducible to  $\exists\forall\text{QSAT}$  in polynomial time. Similarly,  $\Pi_2^P$  is the class of decision problems reducible to  $\forall\exists\text{QSAT}$  in polynomial time. Recall that  $\text{QSAT}$  is  $\text{PSPACE}$ -complete. We know that



We finally remark that the smallest complexity class we consider here is  $\text{DLOGSPACE}$ . In database theory, and especially in its logical counterpart, that is, finite model theory, it is very common to consider parallel complexity classes, of which the smallest one is  $\text{AC}^0$ . These are circuit complexity classes, and the machinery needed to define them is not TMs but rather circuits, parameterized by their fan-in (the number of inputs to their gates), their size, and their depth. Due to the notational overhead this incurs, we shall not be using circuit-based classes in this book. The interested reader can consult books on finite model theory and descriptive complexity to understand the differences between  $\text{DLOGSPACE}$  and classes such as  $\text{AC}^0$ .

## Input Encodings

---

To reason about the computational complexity of problems, we need to represent their inputs (such as databases, queries, and constraints) as inputs to Turing Machines, that is, as words over some finite alphabet.

### *Encoding of Databases*

For databases, the idea is that each value in the active domain can be encoded as a number in binary, and then use further separator symbols that allows us to faithfully encode the facts occurring in the database.

We assume a strict total order  $<_{\text{Rel}}$  on the elements of  $\text{Rel}$ , and a strict total order  $<_{\text{Const}}$  on the elements of  $\text{Const}$ . Consider a schema  $\mathbf{S} = \{R_1, \dots, R_n\}$  with  $R_i <_{\text{Rel}} R_{i+1}$  for each  $i \in [n-1]$ , and a database  $D$  of  $\mathbf{S}$  with  $\text{Dom}(D) = \{a_1, \dots, a_k\}$  and  $a_i <_{\text{Const}} a_{i+1}$  for each  $i \in [k-1]$ . We proceed to explain how  $D$  is encoded as a word over the alphabet

$$\Sigma = \{0, 1, \triangle, \#, \$, \square\}.$$

We first explain how constants, tuples, and relations are encoded:

- The constant  $a_i \in \text{Dom}(D)$ , for  $i \in [k]$ , is encoded as the number  $i$  in binary, and we write  $\text{enc}(a_i)$  for the obtained word over  $\{0, 1\}$ .
- A tuple  $\bar{t} = (a_1, \dots, a_\ell)$  over  $\text{Dom}(D)$ , for  $\ell \geq 0$ , is encoded as the word

$$\text{enc}(\bar{t}) = \begin{cases} \square \text{enc}(a_1) \square \dots \square \text{enc}(a_\ell) \square & \text{if } \ell > 0, \\ \square \square & \text{if } \ell = 0. \end{cases}$$

- A relation  $R_i^D = \{\bar{t}_1, \dots, \bar{t}_m\}$ , for  $i \in [n]$  and  $m \geq 0$ , is encoded as

$$\text{enc}(R_i^D) = \begin{cases} \$ \text{enc}(\bar{t}_1) \$ \dots \$ \text{enc}(\bar{t}_m) \$ & \text{if } m > 0, \\ \$ \$ & \text{if } m = 0. \end{cases}$$

We can now encode the database  $D$  as a word over  $\Sigma$  as follows:

$$\text{enc}(D) = \triangle \text{enc}(a_1) \triangle \cdots \triangle \text{enc}(a_k) \triangle \# \text{enc}(R_1^D) \# \cdots \# \text{enc}(R_n^D) \#.$$

The key property of the above encoding is that, for a database  $D$  of a schema  $\mathbf{S}$ , and a tuple  $\bar{t}$ , given as their encodings  $\text{enc}(D)$  and  $\text{enc}(\bar{t})$ , respectively, we can check via a deterministic computation, which uses logarithmic space in the size of  $\text{enc}(D)$ , whether  $\bar{t} \in R^D$  for some  $R \in \mathbf{S}$ . In what follows, we write  $\text{enc}(i)$  for the binary representation of an integer  $i > 0$ .

**Lemma C.1.** *Let  $\mathbf{S}$  be the schema  $\{R_1, \dots, R_n\}$  with  $R_1 <_{\text{Rel}} \cdots <_{\text{Rel}} R_n$ . Consider a database  $D$  of  $\mathbf{S}$ , a tuple  $\bar{t}$ , and an integer  $i \in [n]$ , and let  $w$  be the word  $\triangleright \text{enc}(D) \flat \text{enc}(\bar{t}) \flat \text{enc}(i)$  over  $\Sigma \cup \{\triangleright, \sqcup, \flat\}$ . There exists a 2-TM  $M$  with alphabet  $\Sigma$  such that the following hold:*

1.  $M$  accepts  $w$  if and only if  $\bar{t} \in R_i^D$ , and
2.  $M$  runs in space  $O(\text{ar}(R_i) \cdot \log |\text{enc}(D)|)$  if  $\text{ar}(R_i) > 0$ , and  $O(\log |\text{enc}(D)|)$  if  $\text{ar}(R_i) = 0$ .

*Proof.* We first give a high-level description of the 2-TM  $M$ ; for brevity, we write  $it$  for the symbol read by the head of the input tape:

1. Let  $ctr = 0$  – this is a counter maintained on the work tape in binary.
2. While  $ctr \neq i$  do the following:
  - a) If  $it = \#$ , then  $ctr := ctr + 1$ .
  - b) Move the head of the input tape to the right so that it reads the first  $\$$  symbol of  $\text{enc}(R_i^D)$ .
3. Move the head of the input tape to the right so that it reads the first  $\sqcup$  symbol of  $\text{enc}(\bar{u})$ , where  $\bar{u}$  is the first tuple of  $R_i^D$  (i.e.,  $\text{enc}(R_i^D) = \$ \text{enc}(\bar{u}) \$ \cdots \$$ ), or the second  $\$$  symbol of  $\text{enc}(R_i^D)$  in case  $R_i^D$  is empty (which means that  $\text{enc}(R_i^D) = \$\$$ ).
4. Erase the content of the work tape by replacing every symbol different than  $\sqcup$  with  $\sqcup$  (since  $ctr$  is not needed further), and move its head after the left marker  $\triangleright$ .
5. Repeat the following steps until  $it = \#$  (which means that the relation  $R_i^D$  has been fully explored):
  - a) While  $it \neq \$$  do the following:
    - (i) Copy  $it$  to the work tape.
    - (ii) Move the head of both tapes to the right.
  - b) Assuming that  $\triangleright u \sqcup$  is the content of the work tape, if  $u = \text{enc}(\bar{t})$ , then halt and accept; otherwise:

- (i) Move the head of the input tape to the right so that it reads the first  $\sqcup$  symbol of the encoding of the next tuple of  $R_i^D$ , or the symbol  $\#$  if the last tuple of  $R_i^D$  has just been explored. In other words, the head of the input tape reads the symbol to the right of the last  $\$$  symbol read during the while loop of step (a).
- (ii) Erase the content of the work tape by replacing every symbol different than  $\sqcup$  with  $\sqcup$  (since the copied tuple is not needed further), and move its head after the left marker  $\triangleright$ .

6. Halt and reject.

It is easy to verify that  $M$  accepts  $w$  if and only if  $enc(R_i^D)$  is of the form  $\$ \cdots \$ enc(\bar{t}) \$ \cdots \$$ , or, equivalently,  $\bar{t} \in R_i^D$ . It remains to argue that  $M$  runs in the claimed space. At each step of the computation of  $M$ , the work tape holds either  $ctr$ , or the word  $enc(\bar{t})$  for some  $\bar{t} \in R_i^D$ . The value of  $ctr$  (represented in binary) can be maintained using  $O(|enc(i)|)$  bits. The encoding of a tuple of  $R_i^D$  takes space  $O(\text{ar}(R_i) \cdot \log |\text{Dom}(D)|)$ . Therefore, the space used is

$$O(\log |enc(i)| + \text{ar}(R_i) \cdot \log |\text{Dom}(D)|).$$

Since  $|enc(i)| \leq |enc(D)|$  and  $|\text{Dom}(D)| \leq |enc(D)|$ , we can conclude that the above 2-TM on input  $w$  runs in space  $O(\text{ar}(R_i) \cdot \log |enc(D)|)$  if  $\text{ar}(R_i) > 0$ , and  $O(\log |enc(D)|)$  if  $\text{ar}(R_i) = 0$ , and the claim follows.  $\square$

Note that the encoding described above is not the only way of encoding a database as a word over a finite alphabet. We could employ any other encoding as long as it enjoys the property established in Lemma [C.1](#), without affecting the complexity results presented in this book.

#### *Encoding of Queries and Constraints*

Queries  $q$  and constraints  $\sigma$  will most commonly be coming from a query language and a class of constraints, respectively, for which we define the *size*  $\|q\|$  or  $\|\sigma\|$  throughout the book. Since queries and constraints of size  $n$  can be encoded as words over a finite alphabet of length  $n \log n$ , we define the length of their Turing Machine encoding to be  $\|q\| \log \|q\|$  and  $\|\sigma\| \log \|\sigma\|$ , respectively.



