



SPLIT-CORRECTNESS IN INFORMATION EXTRACTION

Johannes Doleschal - Benny Kimelfeld - Wim Martens
Yoav Nahshon - Frank Neven

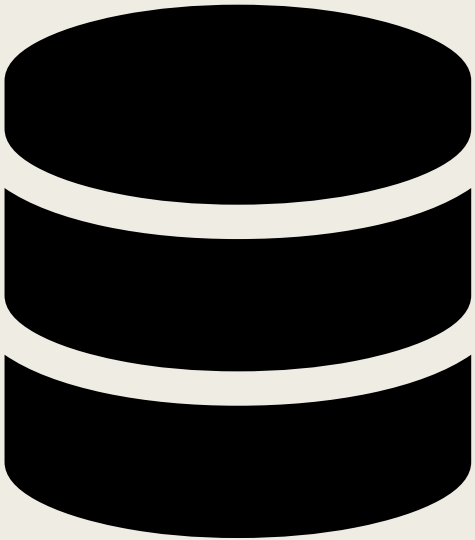
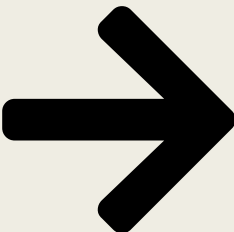


Motivation

Information Extraction

- Input: Huge text documents, e.g. from server logs or financial reports
- Task: Convert the data into a structured relational format

```
{\"timestamp\":\"2017-06-03T18:42:18.018\", \"class\":\"com.orgmanager.handlers.RequestHandler\", \"method\":\"handle\", \"durationMillis\":7}, {\"timestamp\":\"2017-06-03T18:42:18.018\", \"class\":\"com.orgmanager.handlers.RequestHandler\", \"method\":\"handle\", \"durationMillis\":7}, {\"timestamp\":\"2017-06-03T18:42:18.018\", \"class\":\"com.orgmanager.handlers.RequestHandler\", \"method\":\"handle\", \"durationMillis\":7}, {\"timestamp\":\"2017-06-03T18:42:18.018\", \"class\":\"com.orgmanager.handlers.RequestHandler\", \"method\":\"handle\", \"durationMillis\":7}, {\"timestamp\":\"2017-06-03T18:42:18.018\", \"class\":\"com.orgmanager.handlers.RequestHandler\", \"method\":\"handle\", \"durationMillis\":7}, {\"timestamp\":\"2017-06-03T18:42:18.018\", \"class\":\"com.orgmanager.handlers.RequestHandler\", \"method\":\"handle\", \"durationMillis\":7}, {\"timestamp\":\"2017-06-03T18:42:18.018\", \"class\":\"com.orgmanager.handlers.RequestHandler\", \"method\":\"handle\", \"durationMillis\":7}, {\"timestamp\":\"2017-06-03T18:42:18.018\", \"class\":\"com.orgmanager.handlers.RequestHandler\", \"method\":\"handle\", \"durationMillis\":7}, {\"timestamp\":\"2017-06-03T18:42:18.018\", \"class\":\"com.orgmanager.handlers.RequestHandler\", \"method\":\"handle\", \"durationMillis\":7}, {\"timestamp\":\"2017-06-03T18:42:18.018\", \"class\":\"com.orgmanager.handlers.RequestHandler\", \"method\":\"handle\", \"durationMillis\":7}
```



Information Extraction

Naive Extraction

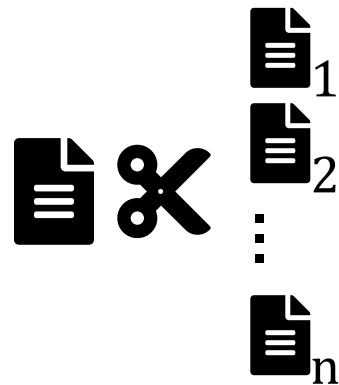


Information Extraction

Naive Extraction



Parallel Extraction



Possible splits ✂:

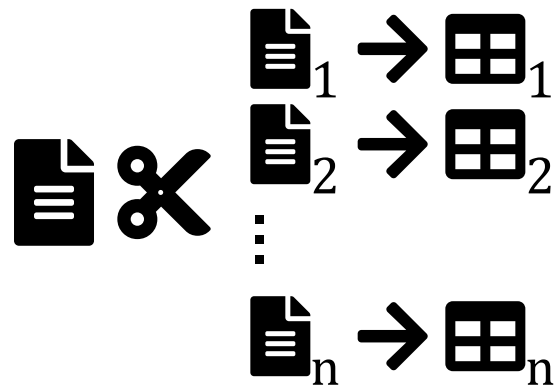
- Text documents into paragraphs or sentences,
- error logs into exceptions,
- server logs into HTTP messages, ...

Information Extraction

Naive Extraction



Parallel Extraction



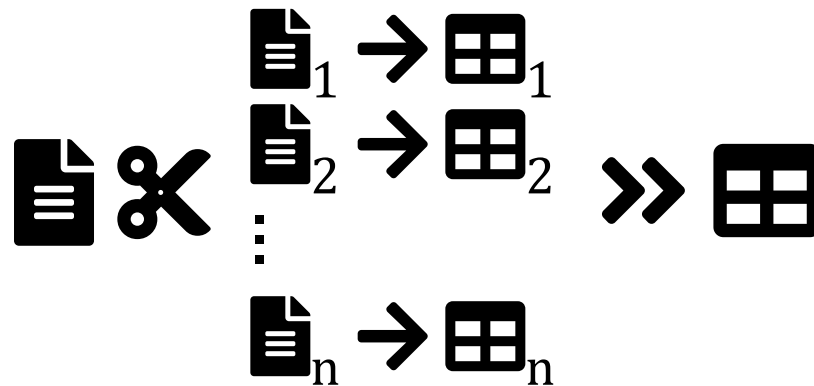
Information Extraction

Naive Extraction



?

Parallel Extraction



More Formally

Information Extractor \rightarrow

$$\rightarrow(\text{document}) = \text{table}$$

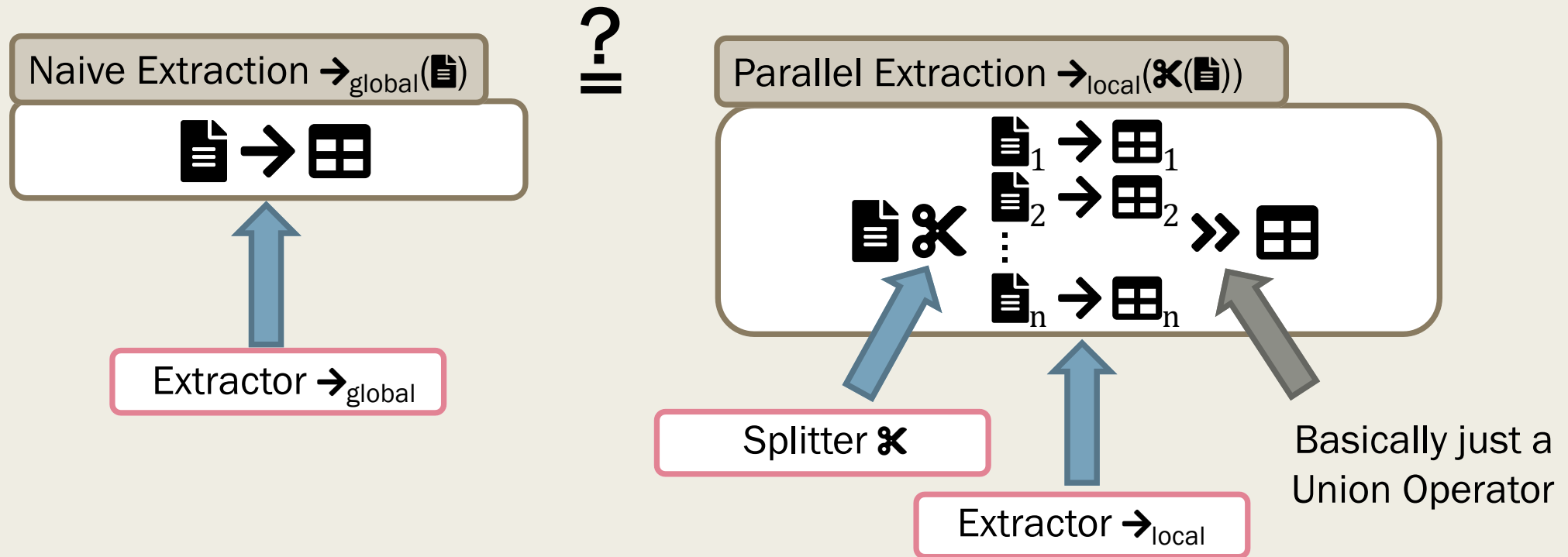
Splitter \bowtie

$$\bowtie(\text{document}) = \text{document}_1, \dots, \text{document}_n$$

Evaluating Extractor \rightarrow on Splitter \bowtie

$$\rightarrow(\bowtie(\text{document})) = \bigcup_{1 \leq i \leq n} \rightarrow(\text{document}_i) = \text{table}$$

Central Question



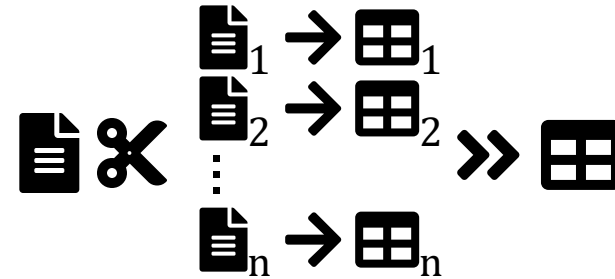
Main Problems

Naive Extraction $\rightarrow_{\text{global}}(\text{📄})$



?

Parallel Extraction $\rightarrow_{\text{local}}(\text{✂}(\text{📄}))$



Split-Correctness

Input: Information Extractors $\rightarrow_{\text{global}}, \rightarrow_{\text{local}}$, Splitter ✂
Question: Is $\rightarrow_{\text{global}}(\text{📄}) = \rightarrow_{\text{local}}(\text{✂}(\text{📄}))$, for all documents 📄 ?

Self-Splittability

Input: Information Extractor \rightarrow , Splitter ✂
Question: Is $\rightarrow(\text{📄}) = \rightarrow(\text{✂}(\text{📄}))$, for all documents 📄 ?

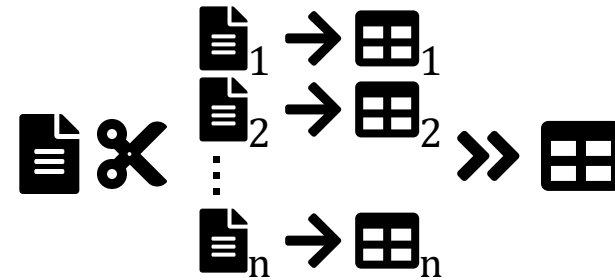
Main Problems

Naive Extraction $\rightarrow_{\text{global}}(\text{📄})$



?

Parallel Extraction $\rightarrow_{\text{local}}(\text{✂️}(\text{📄}))$



Splittability

Input: Information Extractor $\rightarrow_{\text{global}}$, Splitter ✂️

Question: Is there an Information Extractor $\rightarrow_{\text{local}}$, such that $\rightarrow_{\text{global}}(\text{📄}) = \rightarrow_{\text{local}}(\text{✂️}(\text{📄}))$,
for all documents 📄 ?

Regular Document Spanners/Splitters

Information Extractors: Regular document spanners

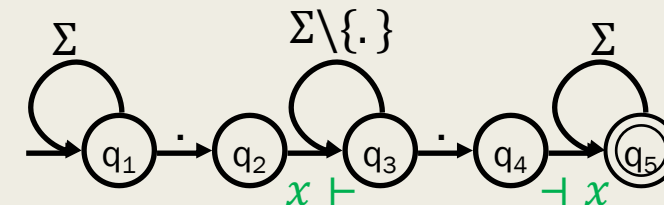
Finite state automata with variable operations.

Essentially, regular expressions with capture variables, closed under relational algebra (\cup, π, \bowtie) .

E.g. sentence extractor

Splitters: Regular document Splitters

Unary regular document spanners



Splitter Disjointness

Splitter Disjointness

A splitter is disjoint, if every part of a document is in at most one output document.

Disjoint Splitter

e.g.

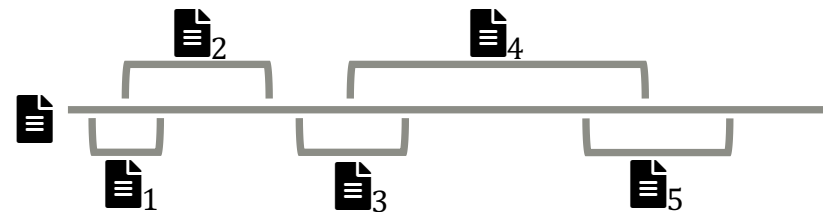
- Sentence segmentation
- Split into paragraphs



Non-Disjoint Splitter

e.g.

- N-Grams
- Pairs of consecutive sentences



Main Results

Self-Splittability

Input: Information Extractor \rightarrow , Splitter \bowtie
Question: Is $\rightarrow(\text{doc}) = \rightarrow(\bowtie(\text{doc}))$, for all documents doc ?

Split-Correctness

Input: Information Extractors $\rightarrow_{\text{global}}, \rightarrow_{\text{local}}$, Splitter \bowtie
Question: Is $\rightarrow_{\text{global}}(\text{doc}) = \rightarrow_{\text{local}}(\bowtie(\text{doc}))$, for all documents doc ?

Theorem 5.1 and Theorem 5.16

Self-Splittability and Split-Correctness for regular document spanners
are PSPACE-complete

Theorem 5.7 and Theorem 5.17

Self-Splittability and Split-Correctness for disjoint splitters and regular document spanners
in normal form* is in polynomial time

* deterministic + functional

Main Results

Splittability




Input: Information Extractor $\rightarrow_{\text{global}}$, Splitter \bowtie

Question: Is there an Information Extractor $\rightarrow_{\text{local}}$, such that $\rightarrow_{\text{global}}(\text{doc}) = \rightarrow_{\text{local}}(\bowtie(\text{doc}))$ for any doc?

Theorem 5.15

Splittability for disjoint splitters and regular document spanners is PSPACE-complete

Split-Constrained Black Boxes

- Information Extraction in practice is often done with other algorithms, e.g.
 - *coreference resolvers*  Splittable by paragraphs
 - *sentiment extractors*  Splittable by paragraphs or sentences
 - *named entity recognition*  Splittable by sentences or 5-grams

Black Box Split-Correctness

Input: Black Box Extractors $\rightarrow_1, \rightarrow_2$, a Splitter \mathfrak{X} , and a set C of Split-Constraints
Question: Is $\rightarrow_1 \bowtie \rightarrow_2$ splittable by \mathfrak{X} under C ?

Split-Constrained Black Boxes

Black Box Split-Correctness

Input: Black Box Extractors $\rightarrow_1, \rightarrow_2$, a Splitter \bowtie , and a set C of Split-Constraints
Question: Is $\rightarrow_1 \bowtie \rightarrow_2$ splittable by \bowtie under C ?

There are cases, where $\rightarrow_1 \bowtie \rightarrow_2$ is splittable by sentences, even though \rightarrow_1 is not!
All 'lost' tuples of \rightarrow_1 are not needed for the join

Theorem 7.3

There are extractors \rightarrow_1 and \rightarrow_2 that are self-splittable by the same splitter \bowtie
but $\rightarrow_1 \bowtie \rightarrow_2$ is not splittable by \bowtie under C

Why should I care about this?

Main Motivation

Parallelized information extraction for huge documents

Further Motivation

1. Reduce the skew, if the a big number of smaller documents must be analyzed
2. Debugging: notify users if extractors reach beyond natural boundaries
3. Handling updates to the input document

Future Work

- Complexity of splittability for non-disjoint splitters
- Try to better understand splittability of split-constrained black boxes
- Empirical analysis of the framework

Thank you for your attention 😊

Paper: doi.org/10.1145/3294052.3319684

