

# Automata and Logic on Trees: Unranked Tree Automata

Wim Martens<sup>1</sup>   Stijn Vansummen<sup>2</sup>

<sup>1</sup>University of Dortmund, Germany

<sup>2</sup>Hasselt University, Belgium

# Outline

- 1 Unranked Tree Automata
- 2 Connection to Ranked Tree Automata
- 3 Minimization
- 4 MSO on Unranked Trees

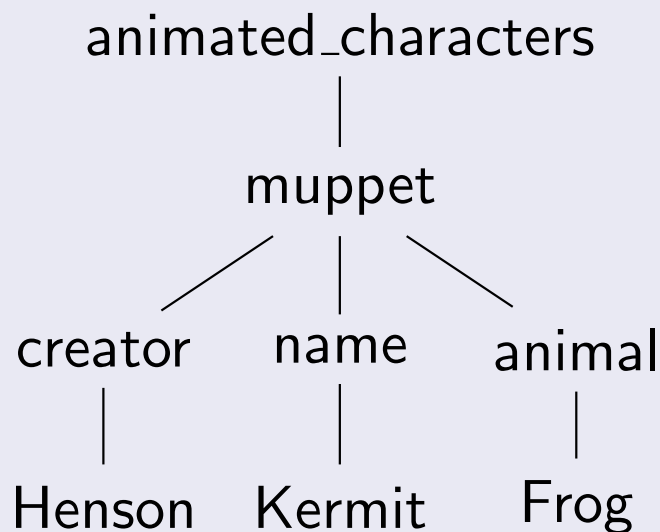
# Outline

- 1 Unranked Tree Automata
- 2 Connection to Ranked Tree Automata
- 3 Minimization
- 4 MSO on Unranked Trees

# Why Unranked Tree Automata?

## Data on the Web

```
<animated_characters>
  <muppet creator="Henson">
    <name> Kermit </name>
    <animal> Frog </animal>
  </muppet>
</animated_characters>
```



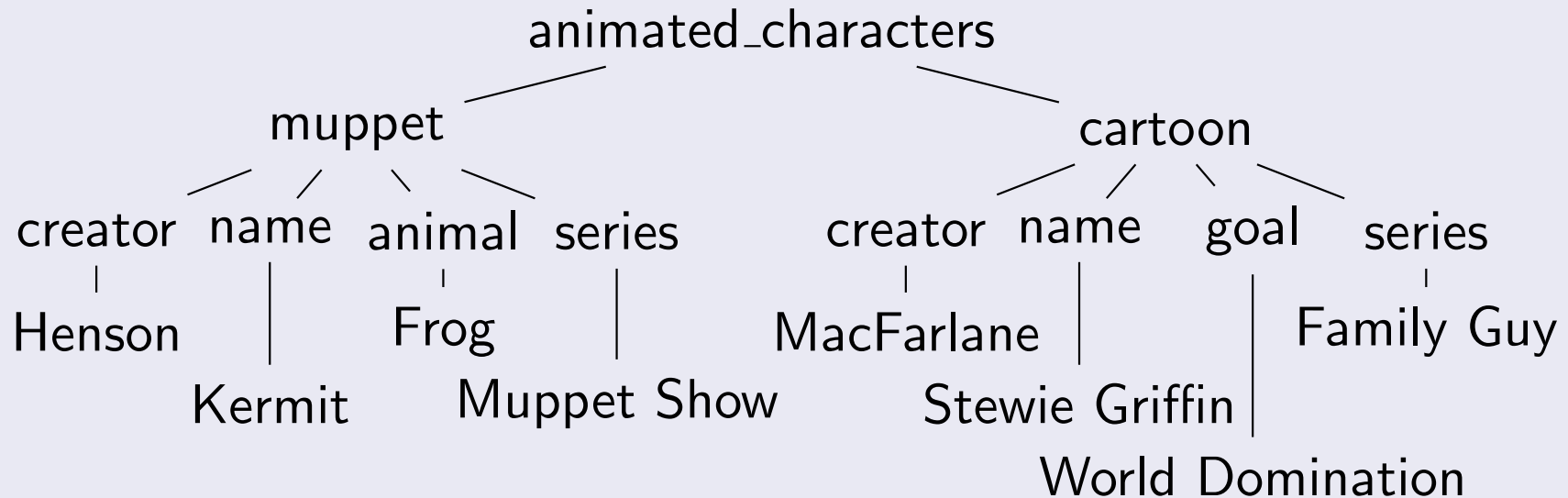
# Why Unranked Tree Automata?

## Data on the Web

```
<animated_characters>
  <muppet creator="Henson">
    <name> Kermit </name>
    <animal> Frog </animal>
    <series> Muppet Show </series>
  </muppet>
  <cartoon creator="MacFarlane">
    <name> Stewie Griffin </name>
    <goal> World Domination </goal>
    <series> Family Guy </series>
  </cartoon>
</animated_characters>
```

# Why Unranked Tree Automata?

## Data on the Web



The number of children is **not** predetermined by a node label  
...and can be arbitrarily (though finitely) many

Hence, the name **unranked** trees

# What are Unranked Tree Automata?

So, let's extend our automata the straightforward way

(Let's take the Boolean Circuit example. That one's easy.)

## Boolean Circuit Evaluation for Binary Trees

$\varepsilon$	$\xrightarrow{\text{true}}$	$t$																		
$\varepsilon$	$\xrightarrow{\text{false}}$	$f$																		
$(t, t)$	$\xrightarrow{\wedge}$	$t$	$(t, f)$	$\xrightarrow{\wedge}$	$f$	$(f, t)$	$\xrightarrow{\wedge}$	$f$	$(f, f)$	$\xrightarrow{\wedge}$	$f$									
$(t, t)$	$\xrightarrow{\vee}$	$t$	$(t, f)$	$\xrightarrow{\vee}$	$t$	$(f, t)$	$\xrightarrow{\vee}$	$t$	$(f, f)$	$\xrightarrow{\vee}$	$f$									

# What are Unranked Tree Automata?

## Boolean Circuit Evaluation for Binary Trees

$\varepsilon \xrightarrow{\text{true}} t$

$\varepsilon \xrightarrow{\text{false}} f$

$(t, t) \xrightarrow{\wedge} t \quad (t, f) \xrightarrow{\wedge} f \quad (f, t) \xrightarrow{\wedge} f \quad (f, f) \xrightarrow{\wedge} f$

$(t, t) \xrightarrow{\vee} t \quad (t, f) \xrightarrow{\vee} t \quad (f, t) \xrightarrow{\vee} t \quad (f, f) \xrightarrow{\vee} f$

## Boolean Circuit Evaluation for Unranked Trees

$\varepsilon \xrightarrow{\text{true}} t$

$\varepsilon \xrightarrow{\text{false}} f$

$(t) \xrightarrow{\wedge} t \quad (t, t) \xrightarrow{\wedge} t$

$(t, t, t) \xrightarrow{\wedge} t \quad (t, t, t, t) \xrightarrow{\wedge} t$



# What are Unranked Tree Automata?

An Unranked Tree Automaton  $A$  consists of:

- $\text{States}(A)$ : its set of states
- $\text{Alphabet}(A)$ : the (non-ranked) alphabet
- $\text{Rules}(A)$ : the transition rules
- $\text{Final}(A)$ : the final states

Here the rules in  $\text{Rules}(A)$  are of the form

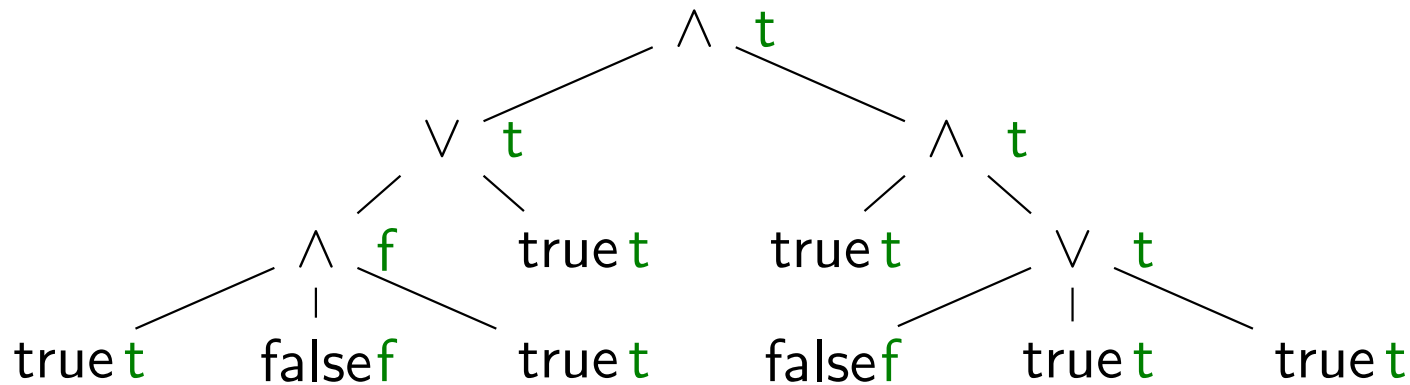
$$L \xrightarrow{a} q,$$

where  $L$  is a regular language over  $\text{States}(A)$

We assume that no two rules  $L_1 \xrightarrow{a} q, L_2 \xrightarrow{a} q$  occur

(Instead,  $L_1 \cup L_2 \xrightarrow{a} q$ )

# Example



## Boolean Circuit Evaluation

$\epsilon$	$\xrightarrow{\text{true}}$	$t$	$t^+$	$\xrightarrow{\wedge}$	$t$	$(t + f)^* f (t + f)^*$	$\xrightarrow{\wedge}$	$f$
$\epsilon$	$\xrightarrow{\text{false}}$	$f$	$f^+$	$\xrightarrow{\vee}$	$f$	$(t + f)^* t (t + f)^*$	$\xrightarrow{\vee}$	$t$

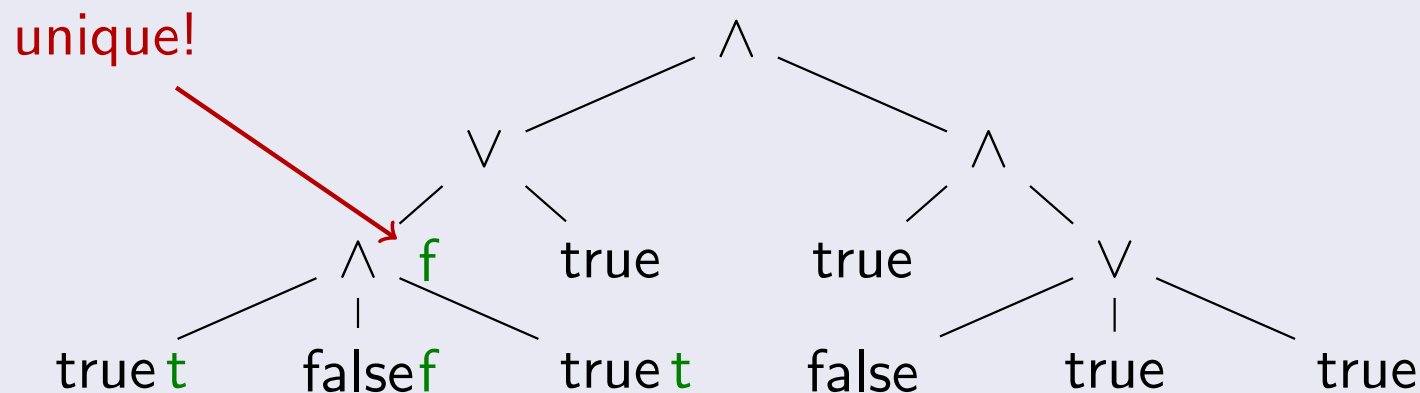
# (Blockwise) Deterministic Unranked Tree Automata

## Definition ((Blockwise) Deterministic Unranked Tree Automaton)

An unranked tree automaton  $A$  is **blockwise deterministic** if

for all rules  $L_1 \xrightarrow{a} q_1$ ,  $L_2 \xrightarrow{a} q_2$ ,  $L_1 \cap L_2 = \emptyset$

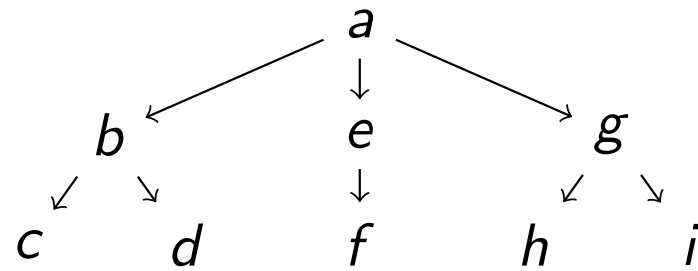
Why can this be seen as bottom-up determinism?



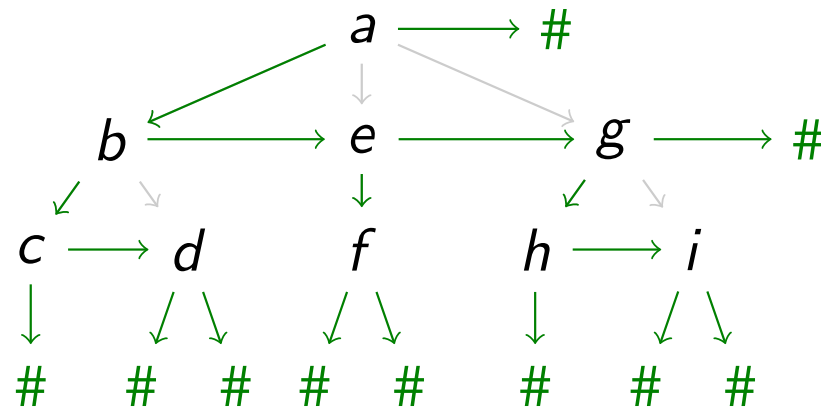
# Outline

- 1 Unranked Tree Automata
- 2 Connection to Ranked Tree Automata
- 3 Minimization
- 4 MSO on Unranked Trees

# First-Child Next-Sibling Encoding



# First-Child Next-Sibling Encoding



## Notation

- $\text{fcns}(L)$ : first-child next-sibling encoding of (unranked) language  $L$
- $\text{fcns}^{-1}(L)$ : first-child next-sibling decoding of (ranked) language  $L$

$\text{fcns}$  is a **bijection** between

- unranked trees over alphabet  $\Sigma$  and
- binary trees over  $\Sigma \uplus \{\#^{(0)}\}$  (every  $a \in \Sigma$  is binary)

# First-Child Next-Sibling Encoding: Automata

## Lemma (Tree Automata Encoding Lemma)

- (1) *For each UTA  $A$ , there is a*  
$$\text{BTA } fcns(A) \text{ accepting } fcns(\text{Language}(A))$$
  - (2) *For each BTA  $A$  there is a*  
$$\text{UTA } fcns^{-1}(A) \text{ accepting } fcns^{-1}(\text{Language}(A))$$
- $fcns(A)$  and  $fcns^{-1}(A)$  are constructible in linear time.*

# The Result Transfer

Corollaries of the FCNS Lemma:

Unranked Regular Tree Languages are Closed Under

- Union
- Intersection
- Complement
- ...



# The Result Transfer

## Corollaries of the FCNS Lemma:

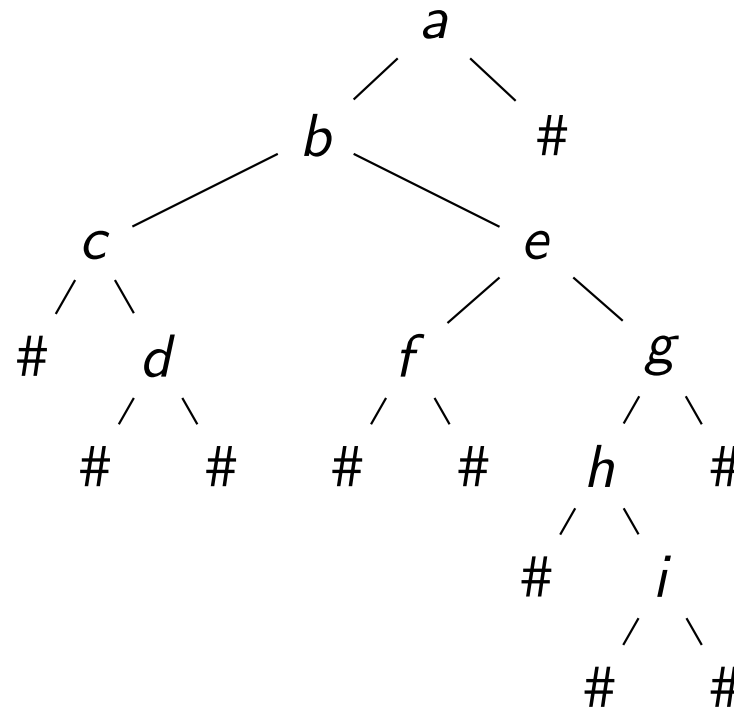
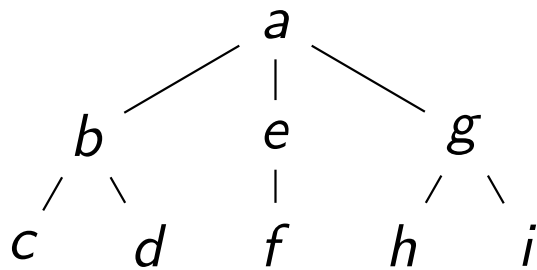
### Complexities for Unranked Tree Automata

- Membership: **PTIME**
- (non)-Emptiness: **PTIME**
- Finiteness: **PTIME**
- Universality: **EXPTIME**
- Containment / Inclusion: **EXPTIME**
- Equivalence: **EXPTIME**
- Intersection (non)-Emptiness: **EXPTIME**

# Outline

- 1 Unranked Tree Automata
- 2 Connection to Ranked Tree Automata
- 3 Minimization**
- 4 MSO on Unranked Trees

# FCNS encoding revisited



# FCNS encoding and determinism

## Observation

Determinism is not preserved by FCNS encoding!

# Determinism Revisited

An Unranked Tree Automaton  $A$  consists of:

- $\text{States}(A)$ : its set of states
- $\text{Alphabet}(A)$ : the (non-ranked) alphabet
- $\text{Rules}(A)$ : the transition rules
- $\text{Final}(A)$ : the final states

Here  $\text{Rules}(A)$  becomes something completely different

## New Rules

$\text{Rules}(A)$  are of the form

$$a \rightarrow D,$$

where  $a \in \text{Alphabet}(A)$  and  $D$  is a finite string automaton

# Determinism Revisited

## New Rules

Rules( $A$ ) are of the form

$$a \rightarrow D,$$

where  $a \in \text{Alphabet}(A)$  and  $D$  is a finite string automaton

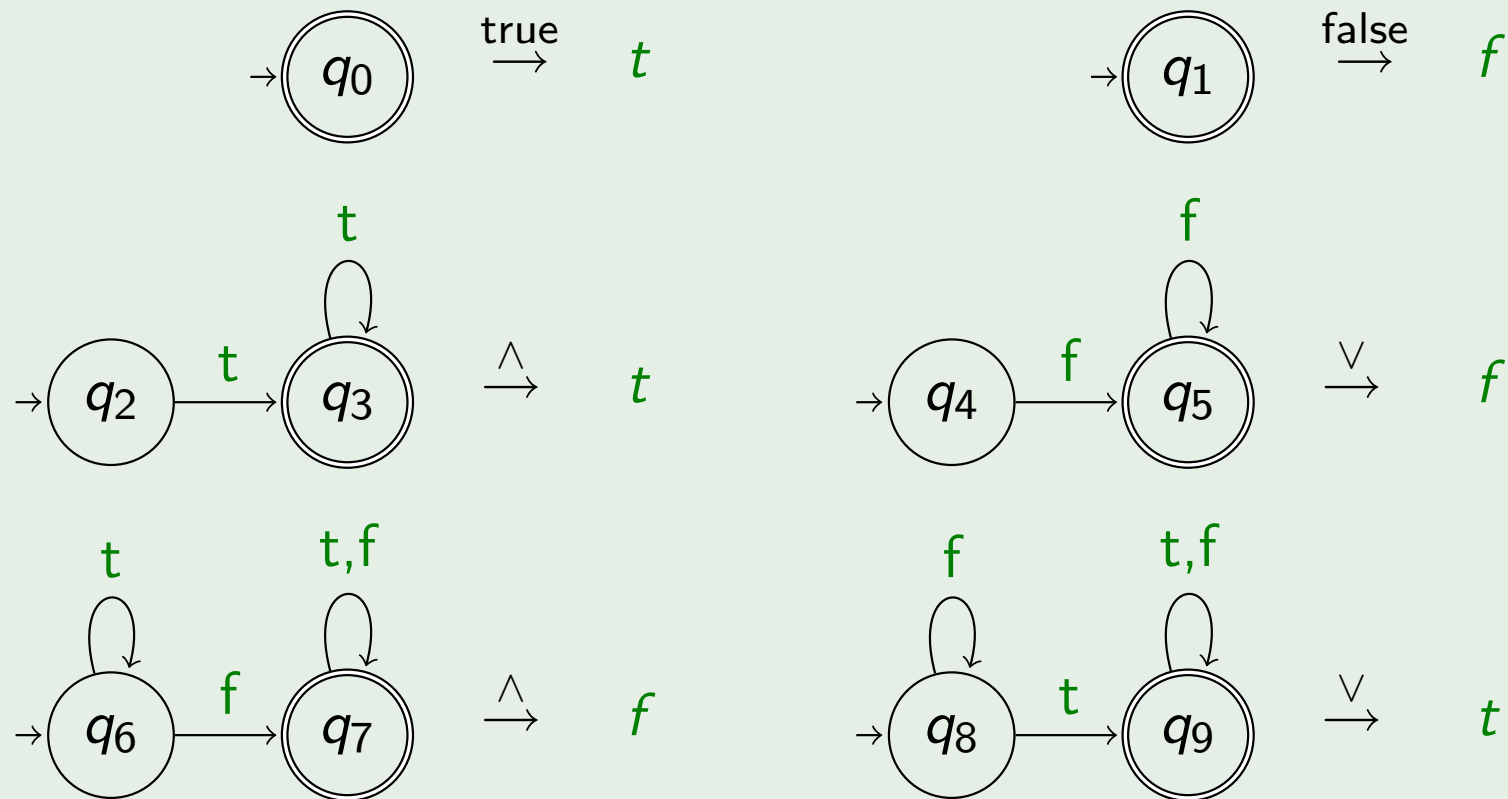
$D$  is a finite string automaton with

- $\text{Alphabet}(D) = \text{States}(A)$
- $\text{States}(D) = \text{States}(A)$

# Intermezzo: Example

## Example (Boolean Circuit Evaluation)

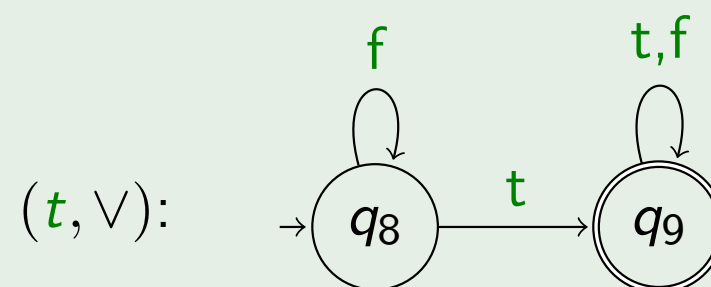
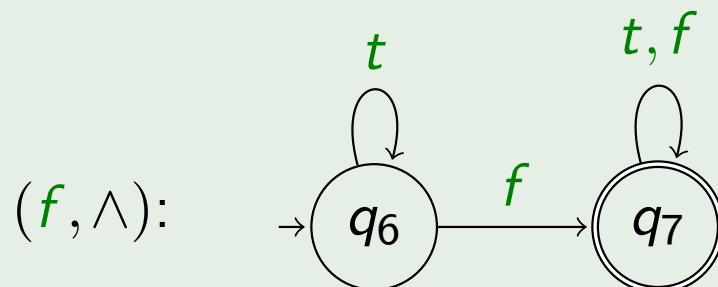
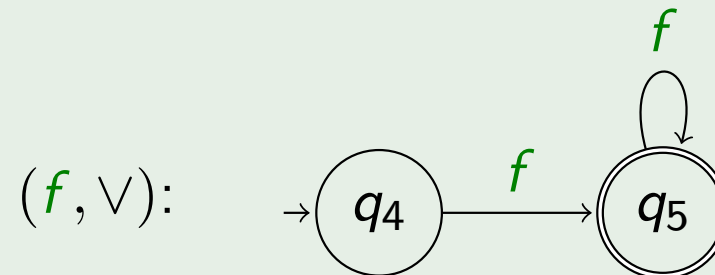
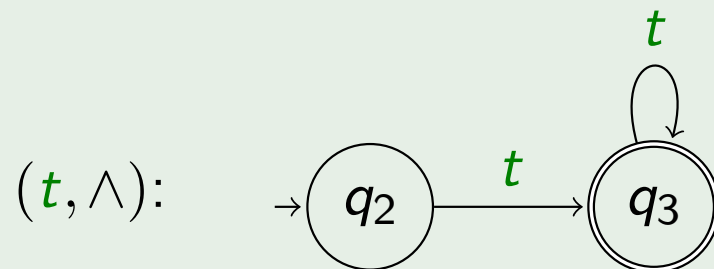
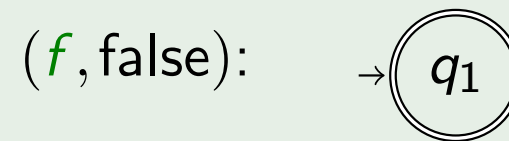
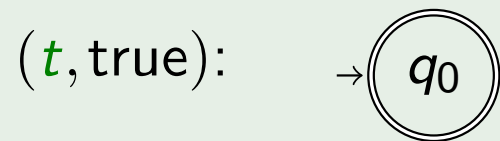
$\varepsilon$	$\xrightarrow{\text{true}}$	$t$	$t^+$	$\xrightarrow{\wedge}$	$t$	$(t+f)^*f(t+f)^*$	$\xrightarrow{\wedge}$	$f$
$\varepsilon$	$\xrightarrow{\text{false}}$	$f$	$f^+$	$\xrightarrow{\vee}$	$f$	$(t+f)^*t(t+f)^*$	$\xrightarrow{\vee}$	$t$



# Intermezzo: Example

## Example (Boolean Circuit Evaluation)

$\varepsilon$	$\xrightarrow{\text{true}}$	$t$	$t^+$	$\xrightarrow{\wedge}$	$t$	$(t+f)^*f(t+f)^*$	$\xrightarrow{\wedge}$	$f$
$\varepsilon$	$\xrightarrow{\text{false}}$	$f$	$f^+$	$\xrightarrow{\vee}$	$f$	$(t+f)^*t(t+f)^*$	$\xrightarrow{\vee}$	$t$

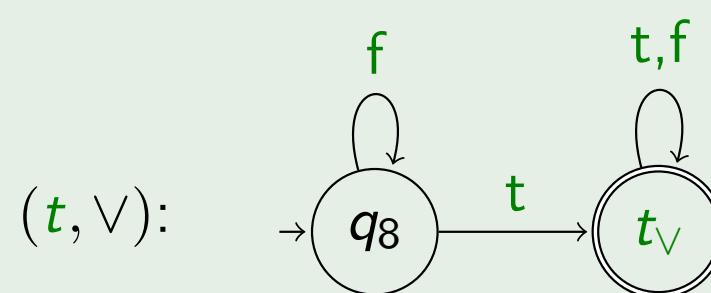
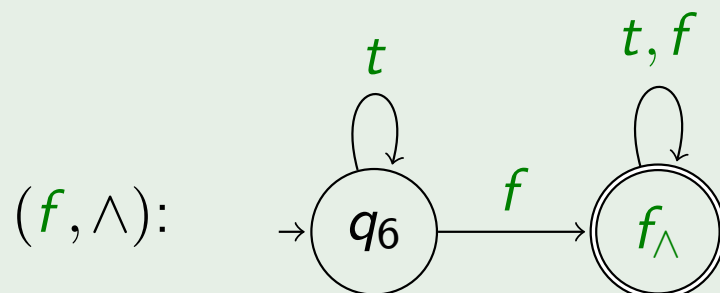
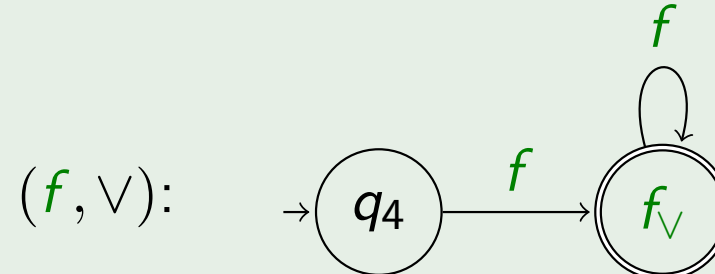
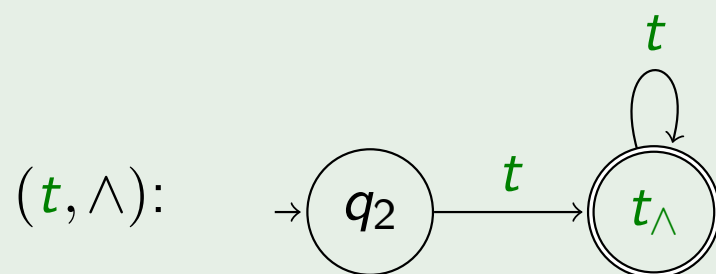
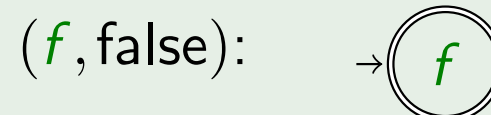
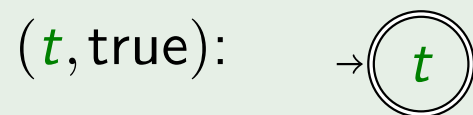




# Intermezzo: Example

## Example (Boolean Circuit Evaluation)

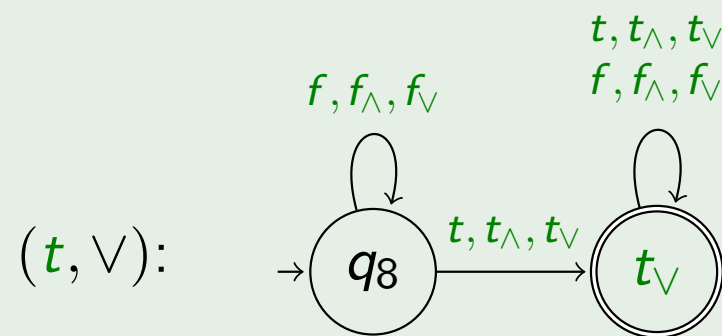
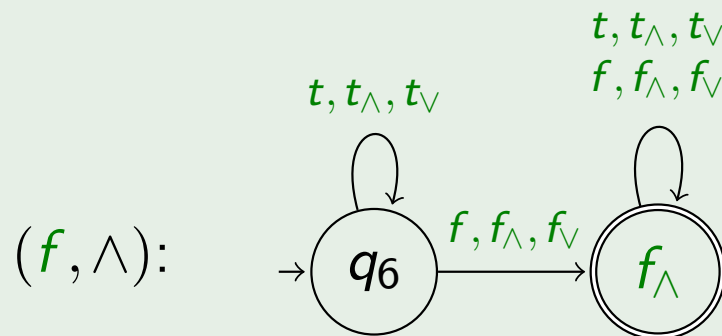
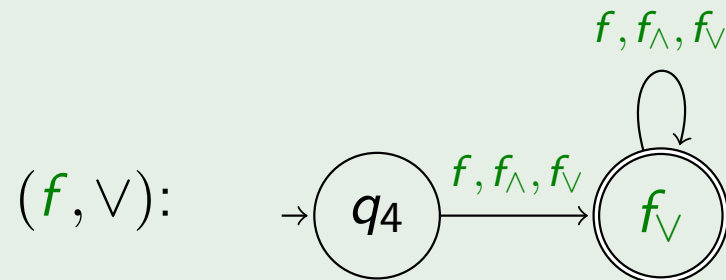
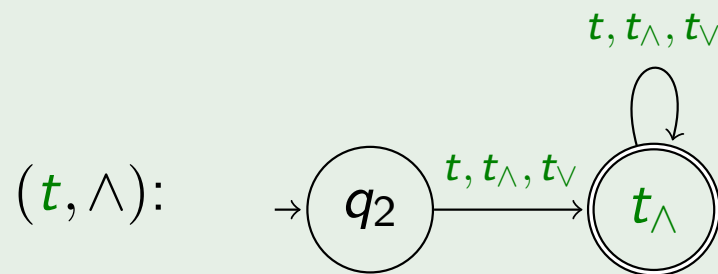
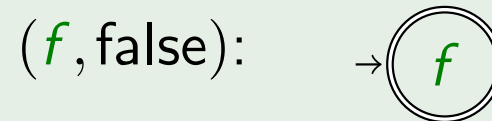
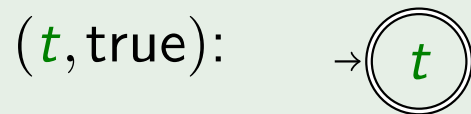
$\varepsilon$	$\xrightarrow{\text{true}}$	$t$	$t^+$	$\xrightarrow{\wedge}$	$t$	$(t+f)^*f(t+f)^*$	$\xrightarrow{\wedge}$	$f$
$\varepsilon$	$\xrightarrow{\text{false}}$	$f$	$f^+$	$\xrightarrow{\vee}$	$f$	$(t+f)^*t(t+f)^*$	$\xrightarrow{\vee}$	$t$



# Intermezzo: Example

## Example (Boolean Circuit Evaluation)

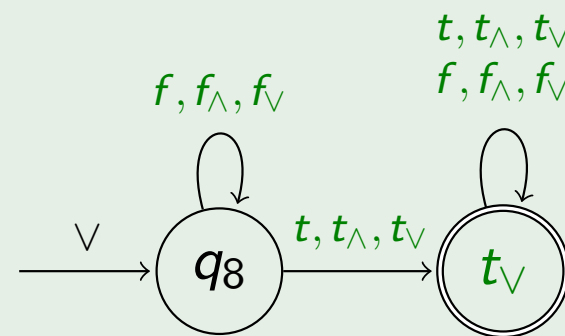
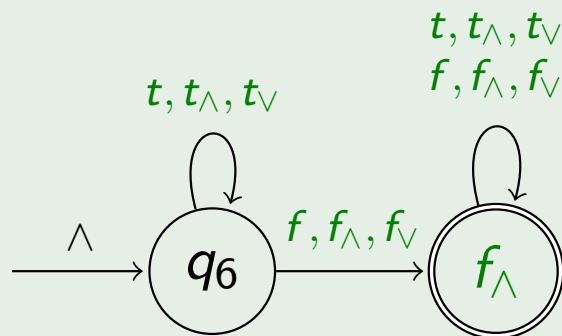
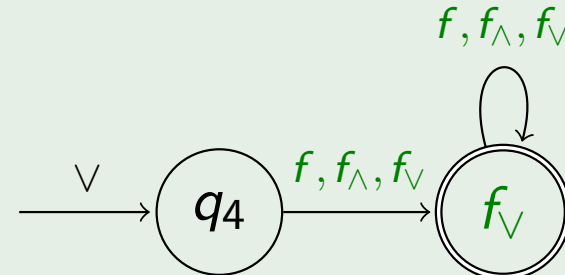
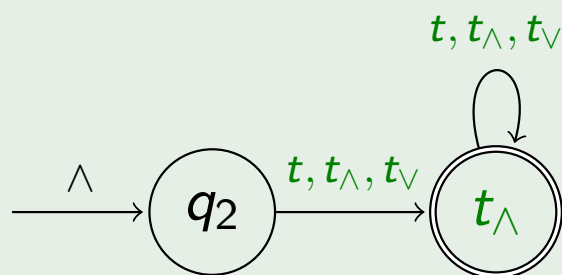
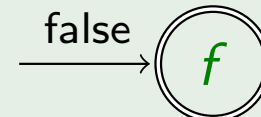
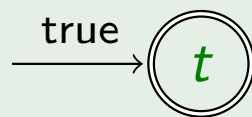
$\varepsilon$	$\xrightarrow{\text{true}}$	$t$	$t^+$	$\xrightarrow{\wedge}$	$t$	$(t + f)^* f (t + f)^*$	$\xrightarrow{\wedge}$	$f$
$\varepsilon$	$\xrightarrow{\text{false}}$	$f$	$f^+$	$\xrightarrow{\vee}$	$f$	$(t + f)^* t (t + f)^*$	$\xrightarrow{\vee}$	$t$



# Intermezzo: Example

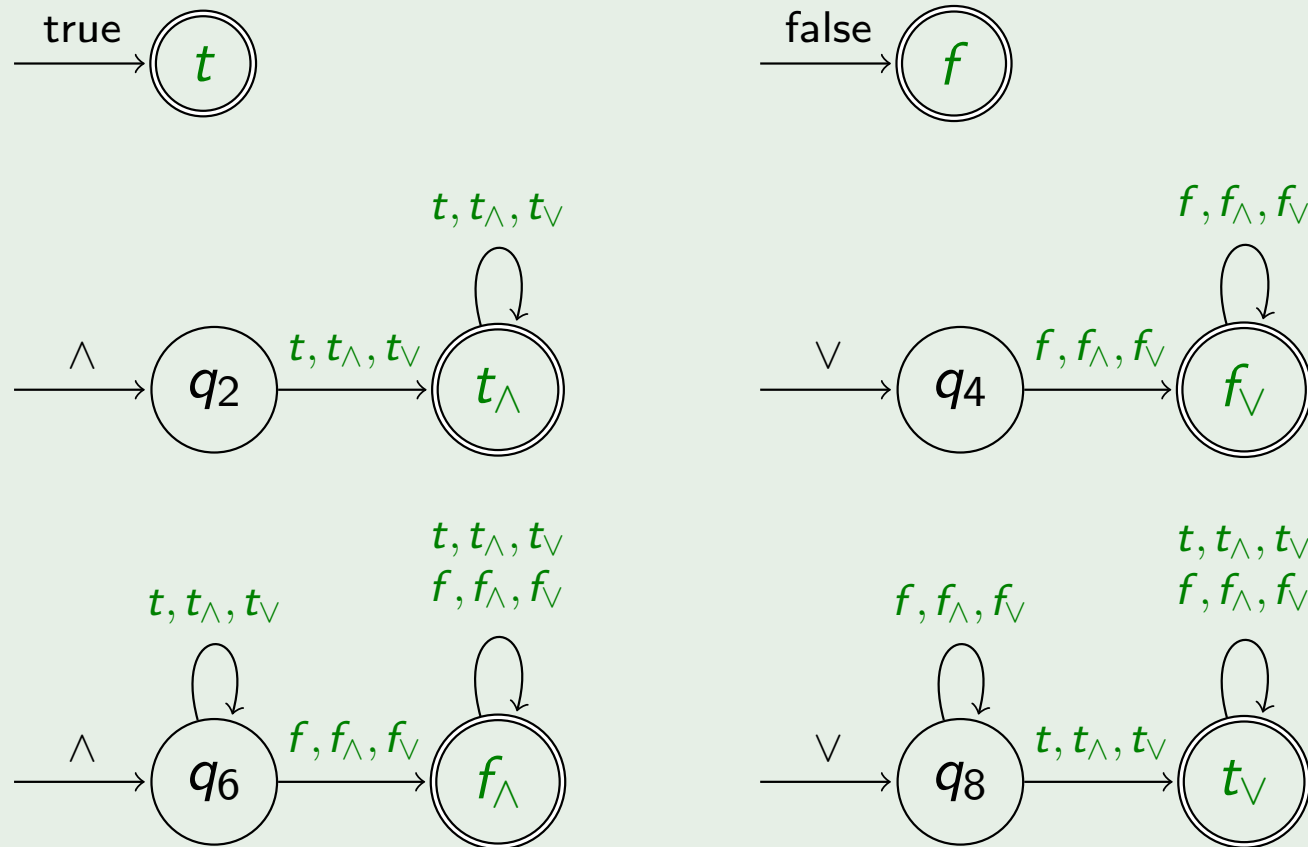
## Example (Boolean Circuit Evaluation)

$\varepsilon$	$\xrightarrow{\text{true}}$	$t$	$t^+$	$\xrightarrow{\wedge}$	$t$	$(t+f)^*f(t+f)^*$	$\xrightarrow{\wedge}$	$f$
$\varepsilon$	$\xrightarrow{\text{false}}$	$f$	$f^+$	$\xrightarrow{\vee}$	$f$	$(t+f)^*t(t+f)^*$	$\xrightarrow{\vee}$	$t$



# Intermezzo: Example

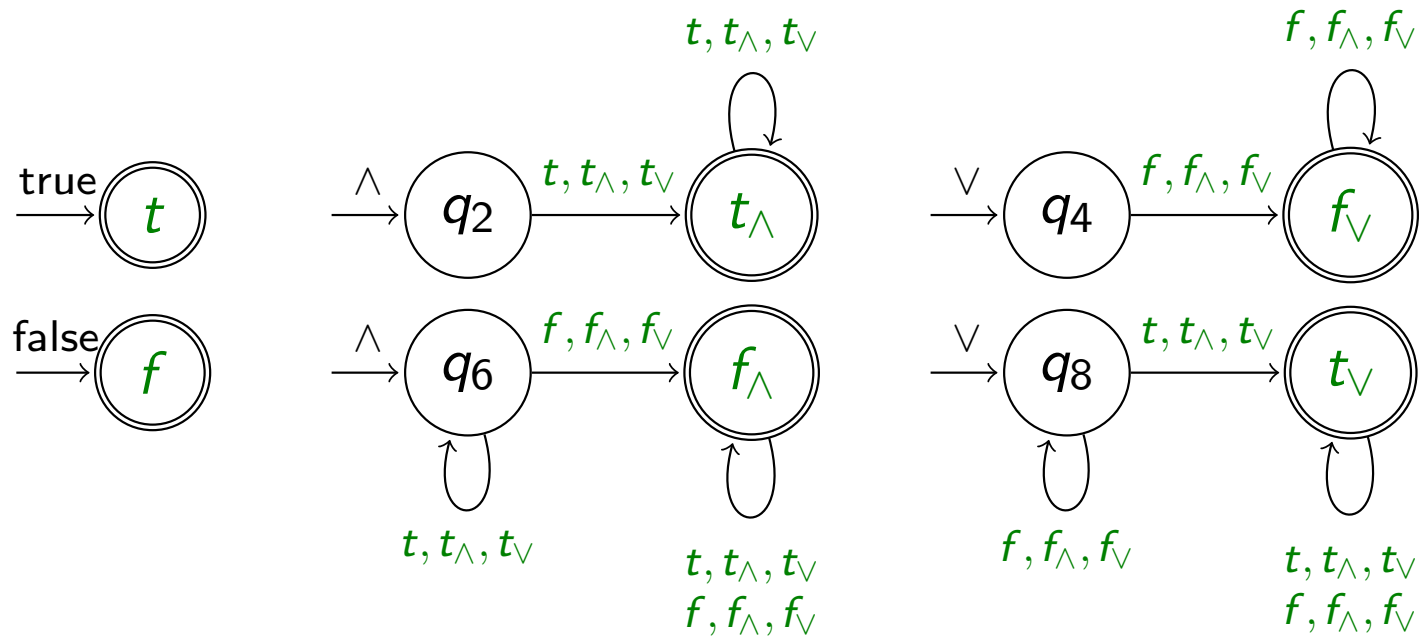
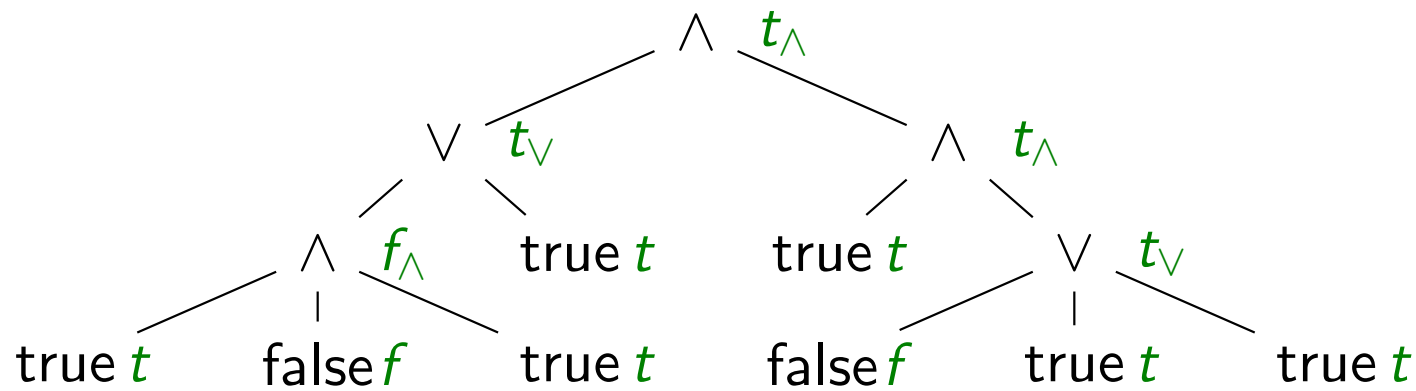
## Example (Boolean Circuit Evaluation)



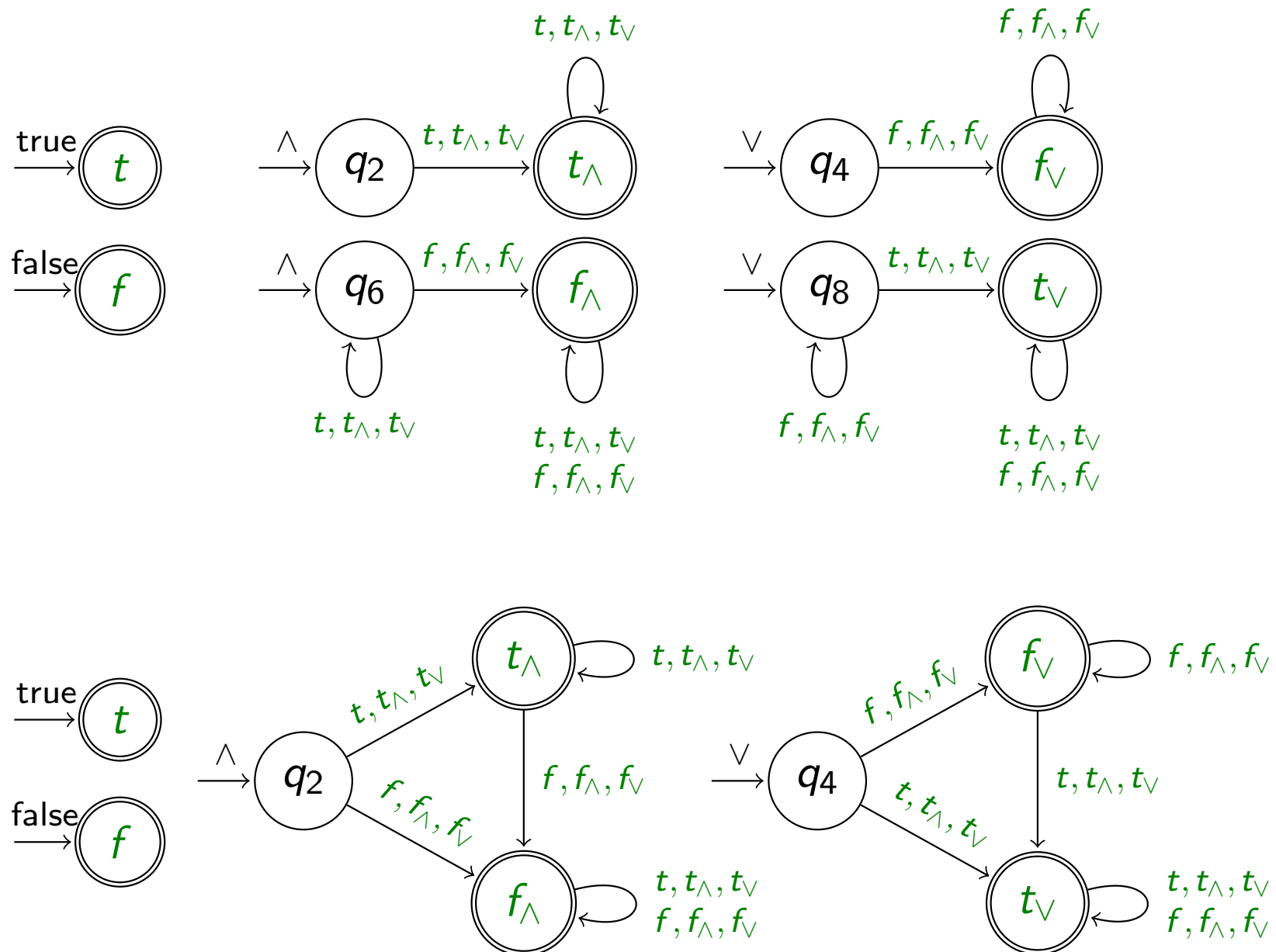
## Definition (Stepwise Determinism (Intuitive))

An unranked tree automaton is (stepwise) deterministic if all the above transitions are deterministic

# Let's run the new guy



# Let's make him Deterministic

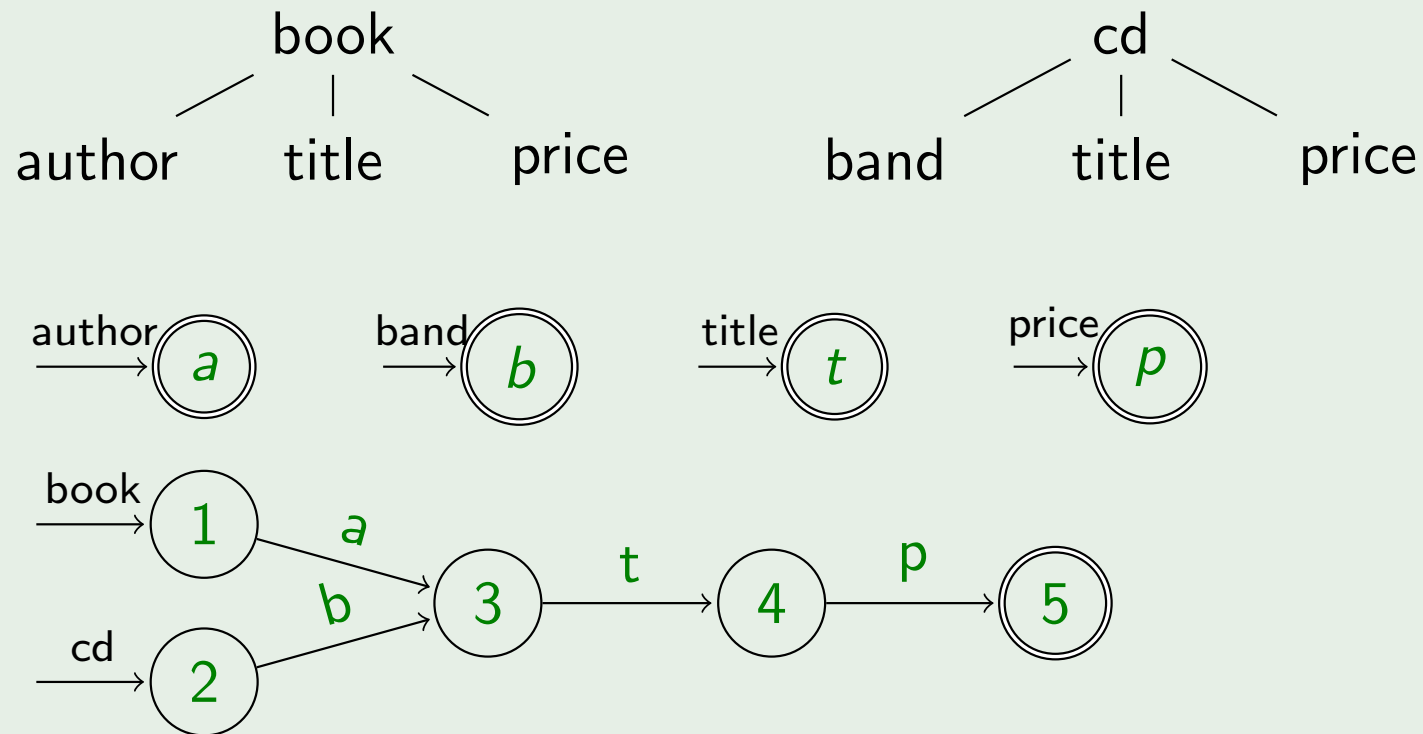


# Another Example (State Sharing)

Recall our rules:  $a \rightarrow D$ , in which  $\text{States}(D) = \text{States}(A)$

This has advantages!

## Example (State Sharing with $A$ )

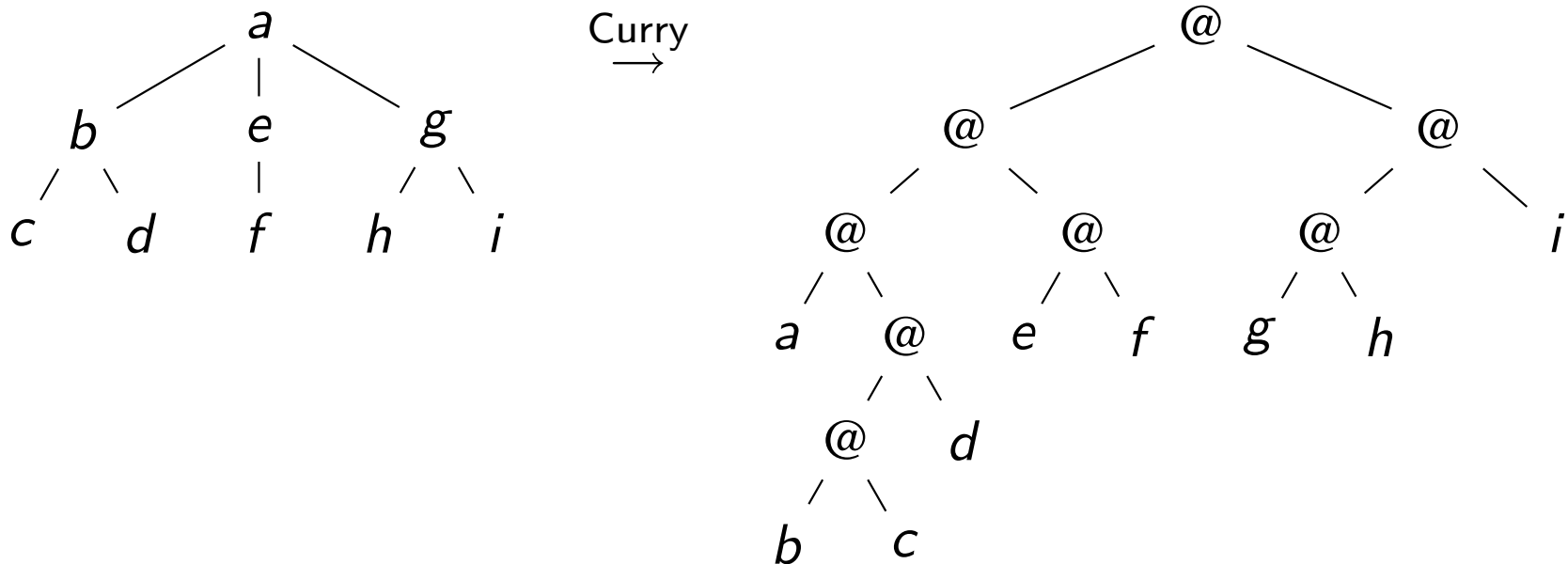


$\text{Final}(A) = \{5\}$

# Currying

## Currying

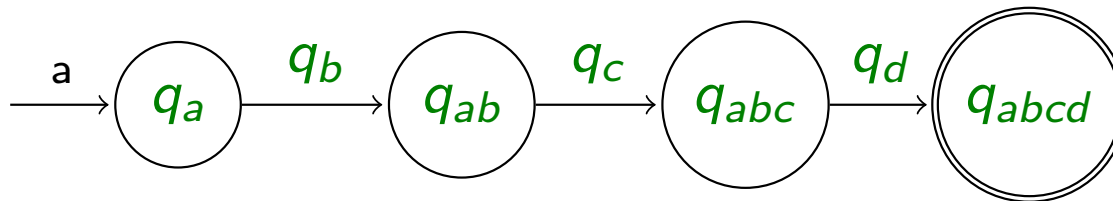
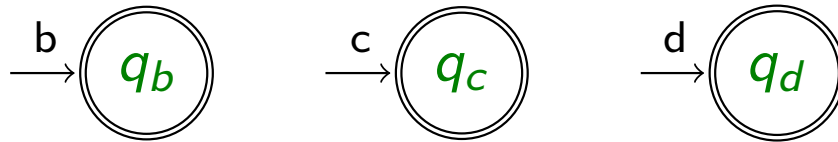
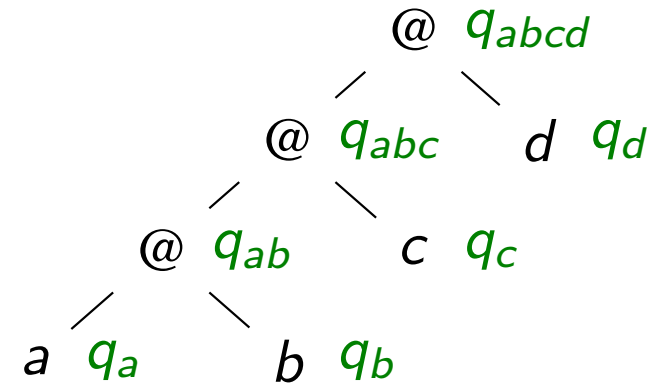
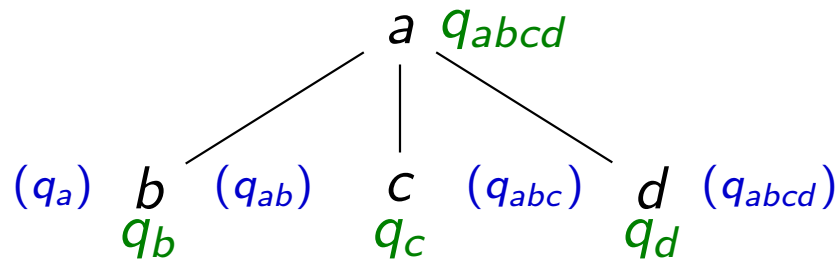
$\text{Curry}(f(a, b, c)) = ((f @ a) @ b) @ c$   
(@: pronounced “applied to”)



Now, a bottom-up deterministic automaton on the right tree...  
does precisely what we want on the left one!

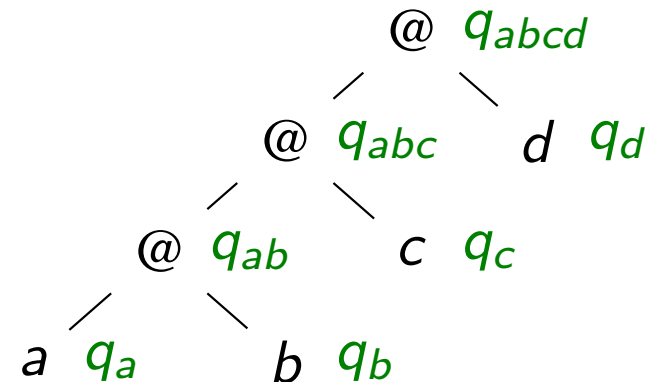
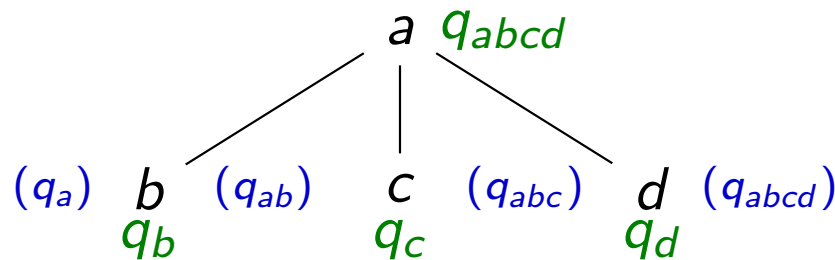


# Currying Relation between Automata



$\text{Final}(A) = \{q_{abcd}\}$

# Currying Relation between Automata



## Currying: transition relation between automata

- $\xrightarrow{a} q$  in unranked automaton  
 $= \varepsilon \xrightarrow{a} q$  in binary automaton
- $q_1 \xrightarrow{p} q_2$  in unranked automaton  
 $= (q_1, p) \xrightarrow{@} q_2$  in binary automaton

One-to-one correspondence in transition rules,  
 which preserves determinism!

# Currying Tree Automata Encoding Lemma

So we actually proved:

## Lemma (Currying Tree Automata Encoding Lemma)

- (1) *For each UTA  $A$ , there is a  
BTA  $\text{Curry}(A)$  accepting  $\text{Curry}(\text{Language}(A))$*
- (2) *For each BTA  $A$  there is a  
UTA  $\text{Curry}^{-1}(A)$  accepting  $\text{Curry}^{-1}(\text{Language}(A))$*

Moreover,

- $\text{Curry}(A)$  and  $\text{Curry}^{-1}(A)$  are constructible in linear time, and
- $\text{Curry}$  and  $\text{Curry}^{-1}$  *preserve (stepwise) determinism!*

# The Result Transfer

- ... all the results from the previous transfer also follow here
- Minimization is in **PTIME** for stepwise deterministic automata!

(Just use the minimization algorithm for the ranked automaton and merge states / transitions in the unranked one iff they should be merged in the ranked one)

## Remark

Minimization is **NP**-hard for the other form of bottom-up determinism, as it's **not entirely deterministic**

# The Result Transfer

- ... all the results from the previous transfer also follow here
- Minimization is in **PTIME** for stepwise deterministic automata!
- Universality: in **PTIME** for stepwise deterministic automata
- Equivalence: in **PTIME** for stepwise deterministic automata
- Containment / Inclusion: in **PTIME** for stepwise deterministic automata

## Remark

Minimization is **NP**-hard for the other form of bottom-up determinism, as it's **not entirely deterministic**

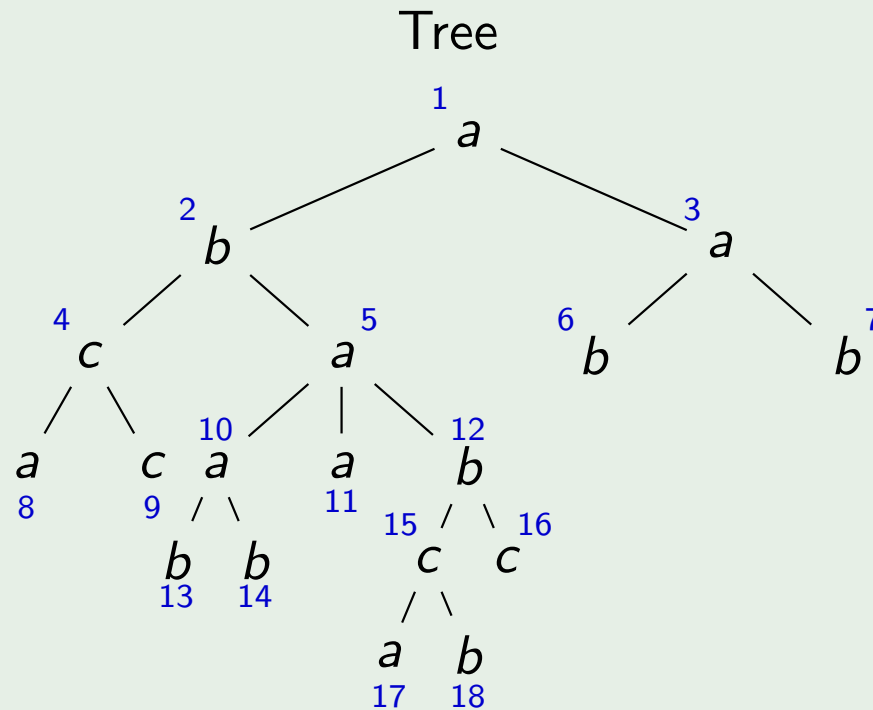
# Outline

- 1 Unranked Tree Automata
- 2 Connection to Ranked Tree Automata
- 3 Minimization
- 4 MSO on Unranked Trees

# Unranked Trees as Structures

An **unranked** tree  $t$  over an alphabet  $\Sigma = \{\sigma, \dots, \sigma'\}$  also naturally corresponds to a structure  $\underline{t}$  over  $V_\Sigma = (E, <, L_\sigma, \dots, L_{\sigma'})$

## Example



## Structure

- $D^{\underline{t}} = \{1, 2, 3, \dots\}$
- $E^{\underline{t}}(1, 2), E^{\underline{t}}(1, 3), E^{\underline{t}}(2, 4), \dots$
- $2 <^{\underline{t}} 3, 4 <^{\underline{t}} 5, \dots$
- $L_a = \{1, 3, 5, \dots\}$
- $L_b = \{2, 6, 7, \dots\}$
- $L_c = \{4, 9, 15, 16\}$

# MSO on Unranked Trees

## Syntax:

$$\begin{aligned} \phi \quad ::= \quad & x = y \mid E(x, y) \mid x < y \mid L_a(x) \mid \cdots \mid L_b(y) \\ & \mid \phi \wedge \phi \mid \neg \phi \mid \exists x \phi \mid X \mid X(x) \mid \exists X \phi \end{aligned}$$

with the usual abbreviations  $\phi \vee \phi$ ,  $\phi \rightarrow \phi$ ,  $\forall x \phi$ ,  $\forall X \phi$ , ...

## Familiar looking example

Every  $a$ -labeled node in  $t$  has a  $b$ -labeled **descendant** if, and only if

$$\underline{t} \models \forall n L_a(n) \rightarrow \exists X (\phi \wedge \psi)$$

where  $\phi := X(n) \wedge \forall m \forall m' (E(m, m') \wedge X(m) \rightarrow X(m'))$

$$\psi := \exists m (X(m) \wedge L_b(m))$$



# MSO on Unranked Trees

## Syntax:

$$\begin{aligned} \phi \quad ::= \quad & x = y \mid E(x, y) \mid x < y \mid L_a(x) \mid \cdots \mid L_b(y) \\ & \mid \phi \wedge \phi \mid \neg \phi \mid \exists x \phi \mid X \mid X(x) \mid \exists X \phi \end{aligned}$$

with the usual abbreviations  $\phi \vee \phi$ ,  $\phi \rightarrow \phi$ ,  $\forall x \phi$ ,  $\forall X \phi$ , ...

## Notation and terminology

- If  $\phi$  is an MSO sentence over  $V_\Sigma$  with  $\Sigma$  an **unranked** alphabet then

$$\text{Language}(\phi) := \{t \mid t \text{ an unranked tree over } \Sigma \text{ such that } t \models \phi\}$$

- An unranked tree language  $S$  is **MSO-definable** if there is some MSO formula  $\phi$  with  $S = \text{Language}(\phi)$

# MSO on Unranked Trees

## Theorem

An **unranked** tree language is MSO-definable if, and only if, it is regular.

The proof is exactly the same as in the ranked case.