

Probabilistic Query Evaluation on Bounded- Treewidth Instances



SIGMOD/PODS PH.D. SYMPOSIUM
JUNE 26, 2016, SAN FRANCISCO

Mikaël Monet
Supervised by Pierre Senellart

Context

- ▶ Boolean queries (yes/no) on relational instances

Context

- ▶ Boolean queries (yes/no) on relational instances
- ▶ We want the answer to contain more information than just « yes/no »:
 - ▶ Add uncertainty
 - ▶ Obtain provenance information

Context

- ▶ Boolean queries (yes/no) on relational instances
- ▶ We want the answer to contain more information than just « yes/no »:
 - ▶ Add uncertainty
 - ▶ Obtain provenance information
- ▶ We need restrictions for all of this to be tractable

A probabilistic database

3

Monet Mikael

R	
a	d
f	e
d	a
b	e

S	
d	e
f	c
a	e
c	e

Q		
c	e	f

A probabilistic database

3

Monet Mikael

R	
a	d
f	e
d	a
b	e

S	
d	e
f	c
a	e
c	e



R		
a	d	0.2
f	e	0.7
d	a	0.13
b	e	0.81

S		
d	e	0.005
f	c	0.9
a	e	0.7
c	e	0.23

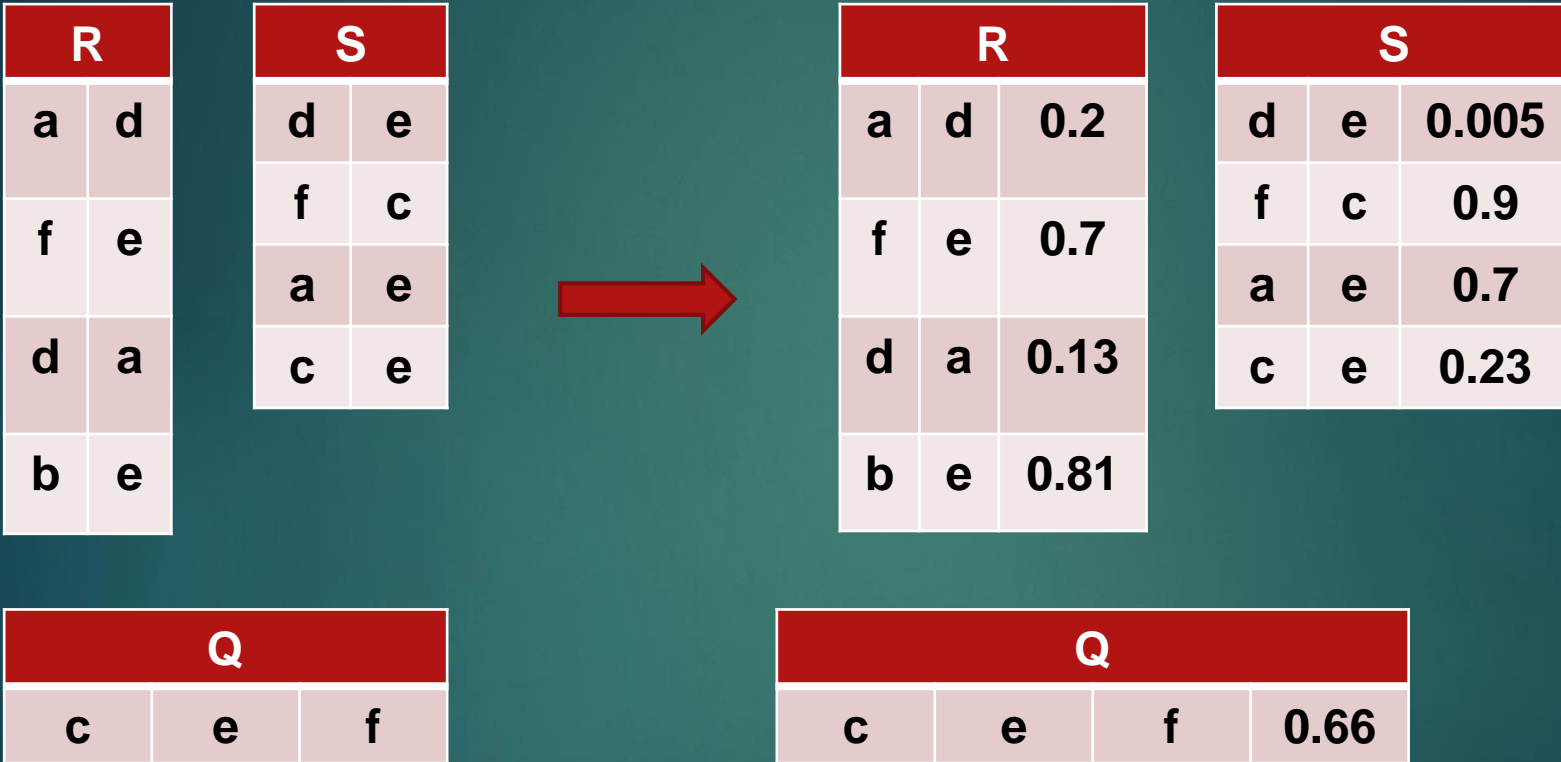
Q		
c	e	f

Q			
c	e	f	0.66

A probabilistic database

3

Monet Mikael




TID model

Probability of a possible world

R		
a	d	0.2
f	e	0.7
d	a	0.13
b	e	0.81

S		
d	e	0.005
f	c	0.9
a	e	0.7
c	e	0.23

A possible
world I



R	
f	e
d	a

S	
c	e

Monet Mikael

Q			
c	e	f	0.66

Q		
c	e	f

Probability of a possible world

R		
a	d	0.2
f	e	0.7
d	a	0.13
b	e	0.81

S		
d	e	0.005
f	c	0.9
a	e	0.7
c	e	0.23

A possible
world I

R	
f	e
d	a

S	
c	e

Monet Mikael

Q			
c	e	f	0.66

Q		
c	e	f

Probability $\Pr(I)$ of this possible world =
 $0.7 * 0.13 * 0.23 * 0.66$

Probability of a possible world

R		
a	d	0.2
f	e	0.7
d	a	0.13
b	e	0.81

S		
d	e	0.005
f	c	0.9
a	e	0.7
c	e	0.23

A possible



world I

R	
f	e
d	a

S	
c	e

Monet Mikael

Q			
c	e	f	0.66

Q		
c	e	f

Probability $\Pr(I)$ of this possible world =

$$0.7 * 0.13 * 0.23 * 0.66$$

$$* (1 - 0.2) * (1 - 0.81) * (1 - 0.005) * (1 - 0.9) * (1 - 0.7)$$

Probabilistic query evaluation (PQE)

- ▶ Focus on Boolean queries (yes/no)

Probabilistic query evaluation (PQE)

- ▶ Focus on Boolean queries (yes/no)
- ▶ Probability of a query Q on probabilistic instance \mathfrak{I} :

$$P(Q) = \sum_{I \subseteq \mathfrak{I}, I \models Q} \Pr(I)$$

Probabilistic query evaluation (PQE)

- ▶ Focus on Boolean queries (yes/no)
- ▶ Probability of a query Q on probabilistic instance \mathfrak{I} :

$$P(Q) = \sum_{I \subseteq \mathfrak{I}, I \models Q} \Pr(I)$$

- ▶ Problem: in general #P-hard

3 possible directions

- ▶ Approximate
- ▶ Restrict queries
- ▶ Restrict instances

1) Approximate probability computation

1) Approximate probability computation

- ▶ Monte-Carlo sampling

1) Approximate probability computation

- ▶ Monte-Carlo sampling
- ▶ Inconvenient: running time quadratic in desired precision

1) Approximate probability computation

- ▶ Monte-Carlo sampling
- ▶ Inconvenient: running time quadratic in desired precision
 - ⇒ Not adequate for low probabilities.

2) Restricting the class of queries

- ▶ [Dalvi and Suciu 2012] shows the following dichotomy for any UCQ Q :

2) Restricting the class of queries

- ▶ [Dalvi and Suciu 2012] shows the following dichotomy for any UCQ Q :
- ▶ Either PQE is PTIME on all instances

2) Restricting the class of queries

- ▶ [Dalvi and Suciu 2012] shows the following dichotomy for any UCQ Q :
- ▶ Either PQE is PTIME on all instances
- ▶ Or PQE is #P-hard on all instances

2) Restricting the class of queries

- ▶ [Dalvi and Suciu 2012] shows the following dichotomy for any UCQ Q :
 - ▶ Either PQE is PTIME on all instances
 - ▶ Or PQE is #P-hard on all instances

- ▶ Simple conjunctive query $\exists x,y R(x),S(x,y),T(y)$ is already #P-hard!

2) Restricting the class of queries

- ▶ [Dalvi and Suciu 2012] shows the following dichotomy for any UCQ Q :
 - ▶ Either PQE is PTIME on all instances
 - ▶ Or PQE is #P-hard on all instances
- ▶ Simple conjunctive query $\exists x,y R(x),S(x,y),T(y)$ is already #P-hard!
- ▶ Criterion is too crisp

3) Restricting the shape of the instances

- ▶ Bound the *treewidth* of instances by a constant
- ▶ Treewidth: measure used to tell how far a graph is from being a tree

Treewidth

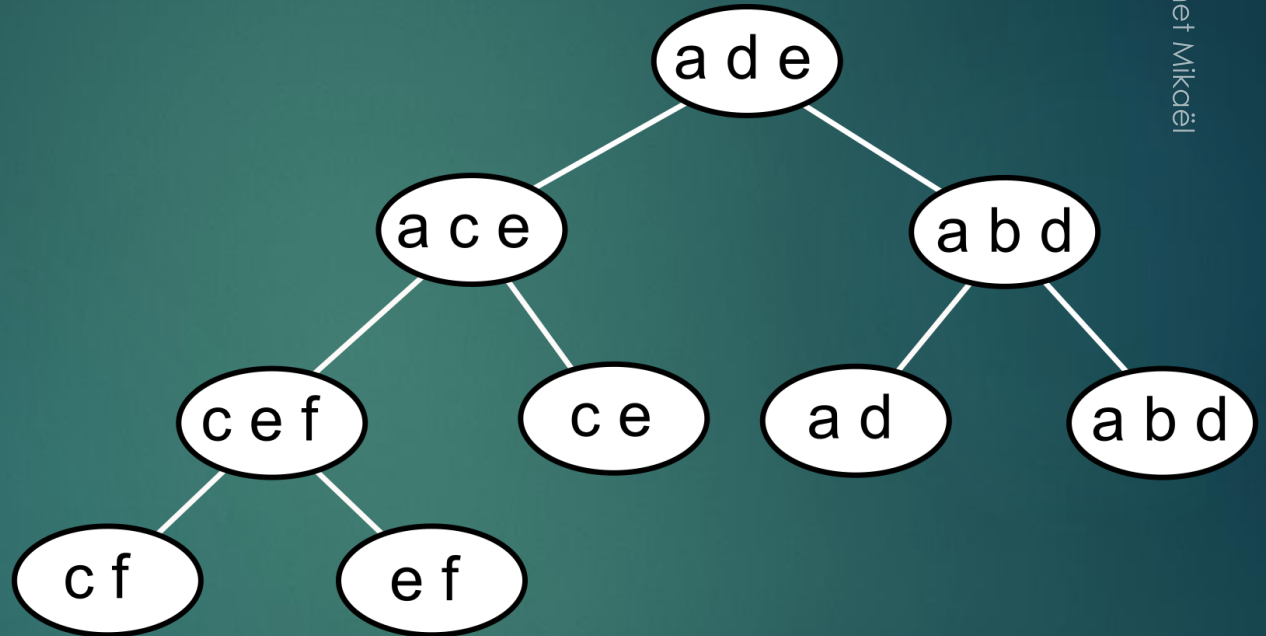
10

Monet Mikael

R	
a	d
f	e
d	a
b	e

S	
d	e
f	c
a	e
c	e

Q		
c	e	f



Treewidth

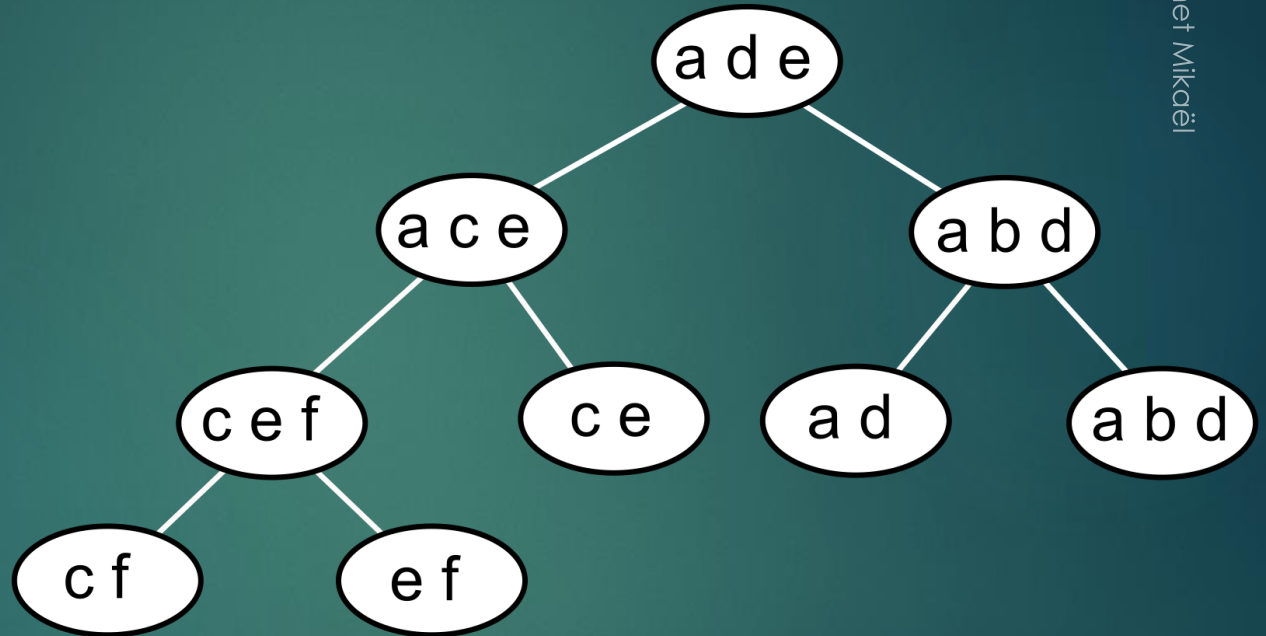
10

Monet Mikael

R	
a	d
f	e
d	a
b	e

S	
d	e
f	c
a	e
c	e

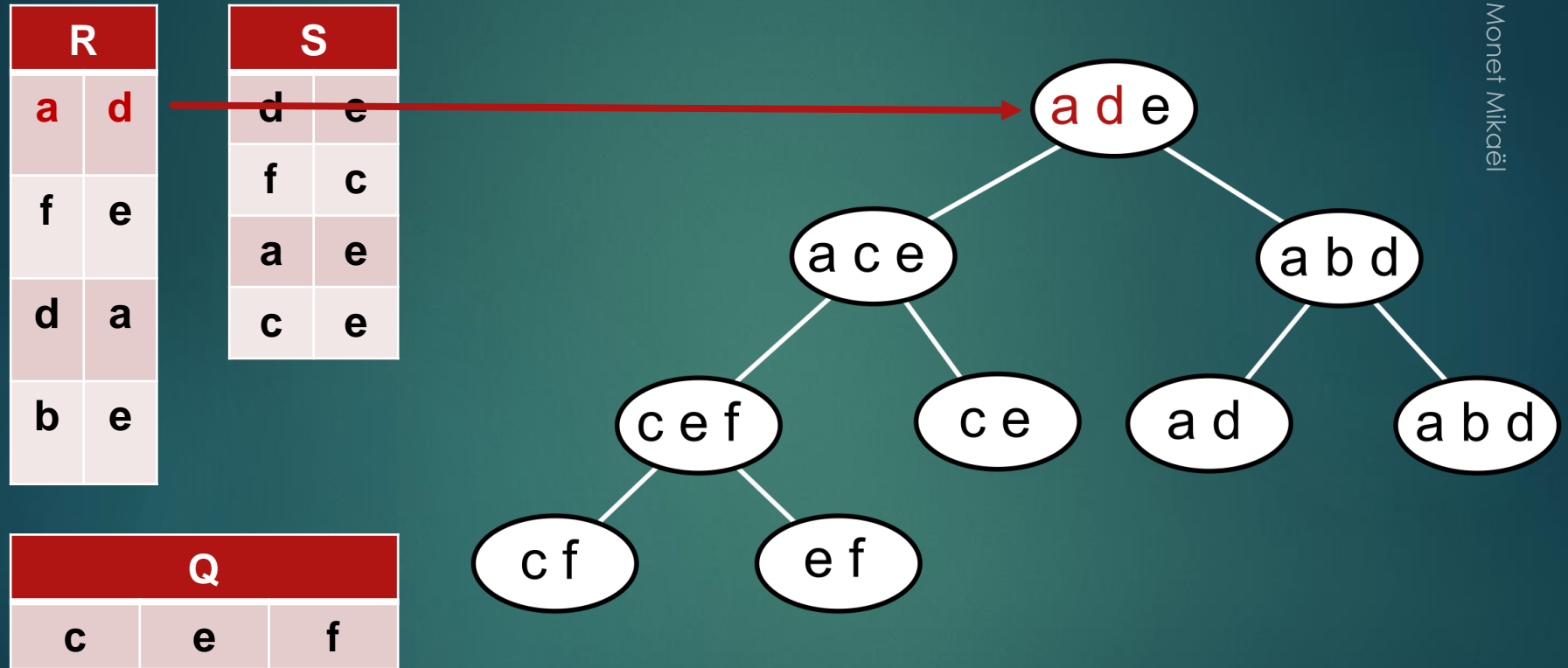
Q		
c	e	f



Treewidth

10

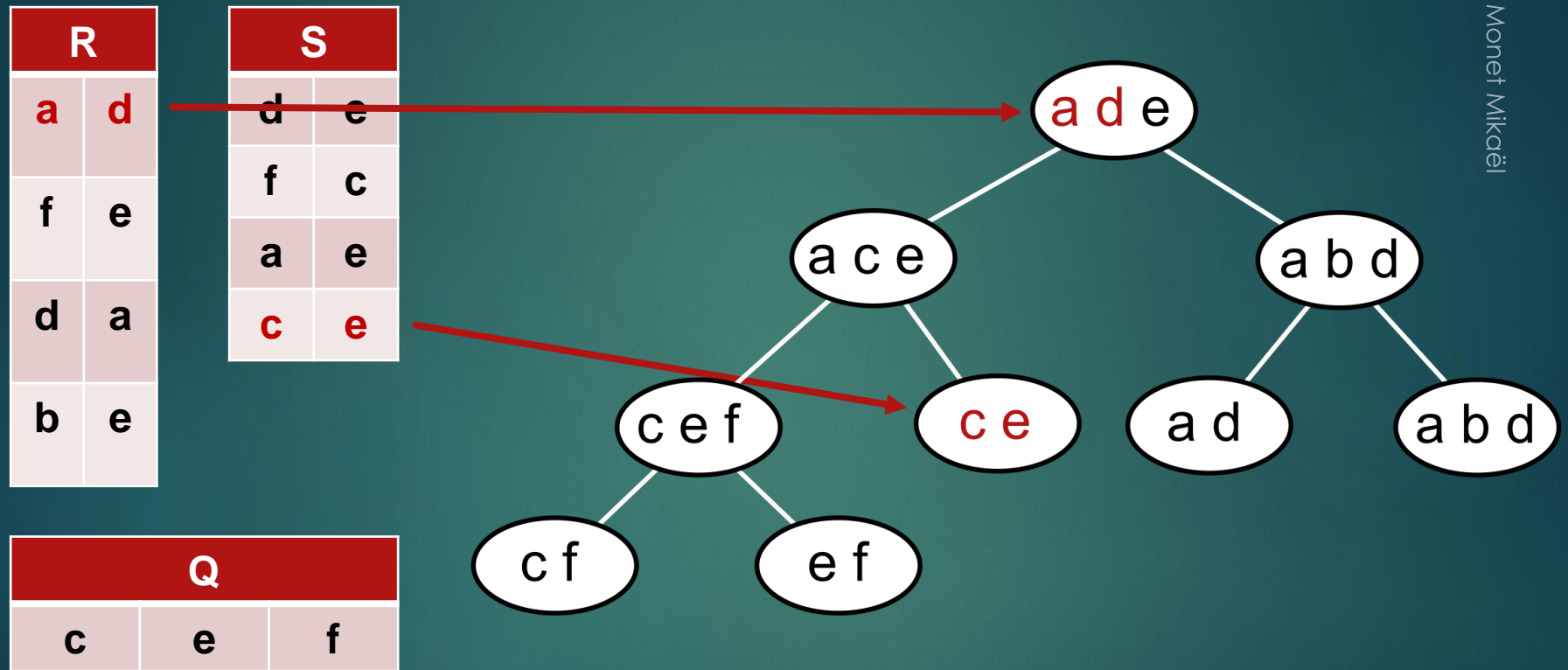
Monet Mikael



Treewidth

10

Monet Mikael



Treewidth

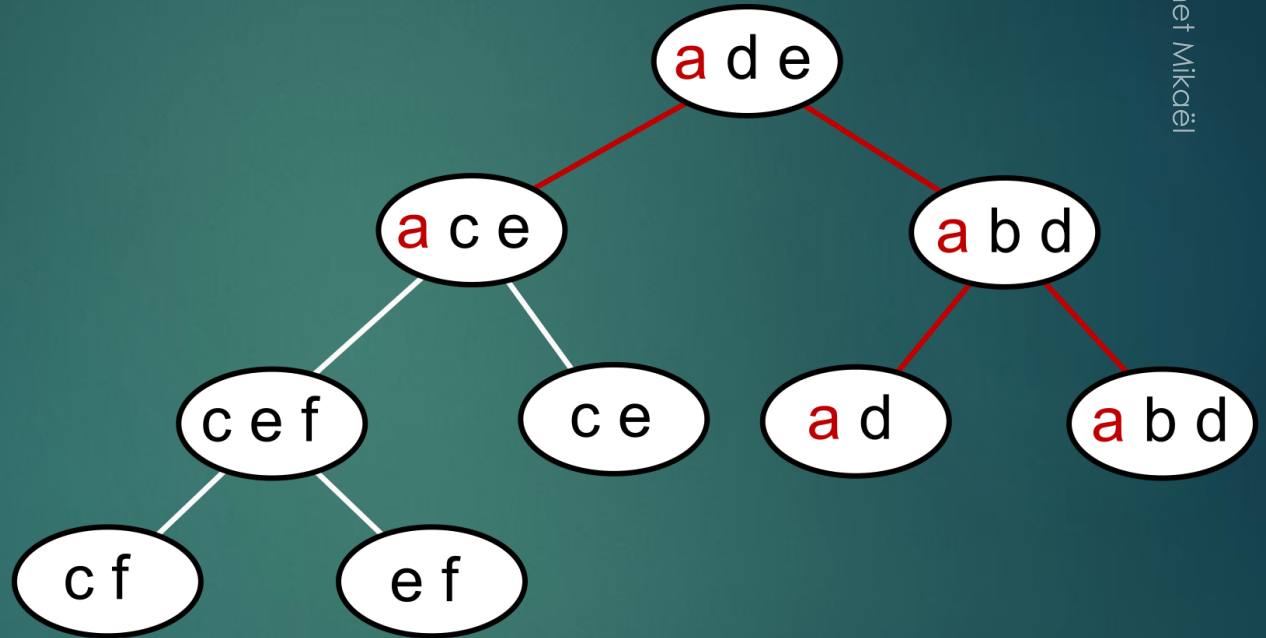
10

Monet Mikael

R	
a	d
f	e
d	a
b	e

S	
d	e
f	c
a	e
c	e

Q		
c	e	f



Treewidth

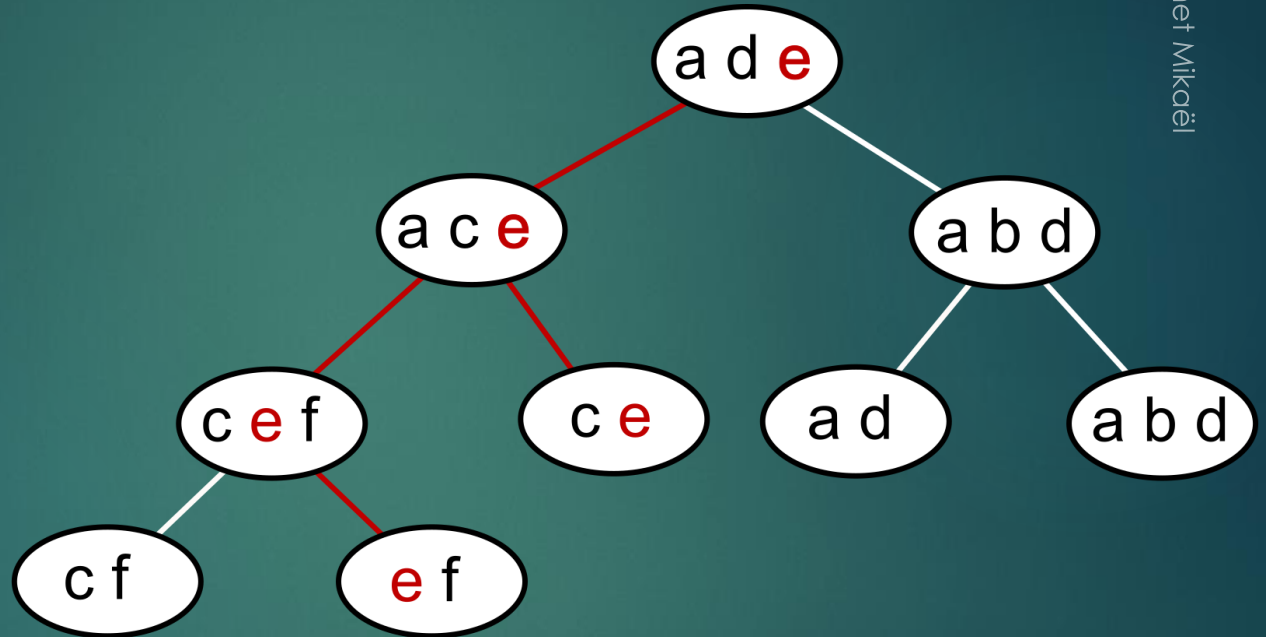
10

Monet Mikael

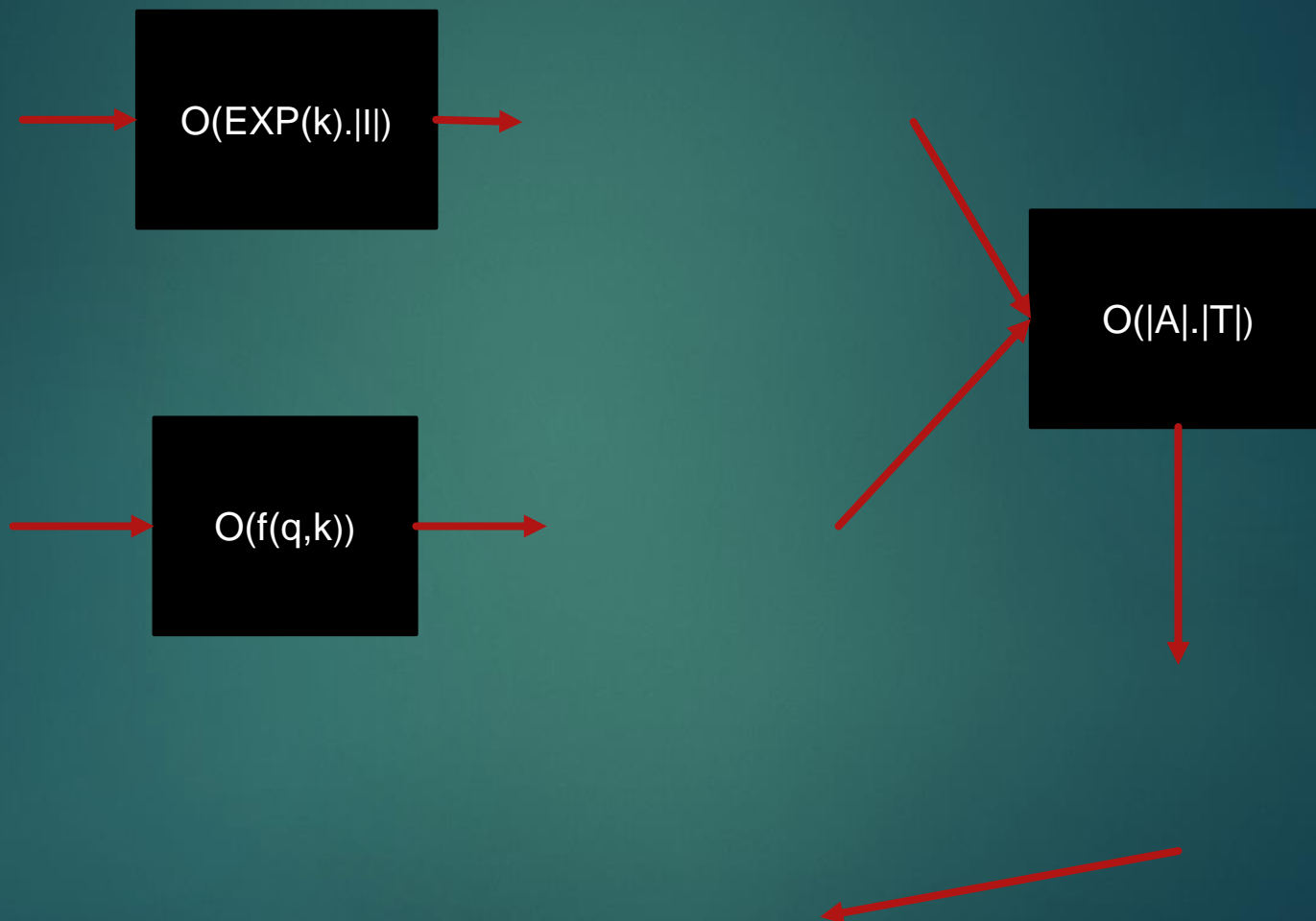
R	
a	d
f	e
d	a
b	e

S	
d	e
f	c
a	e
c	e

Q		
c	e	f



Divide and conquer !



Instance I
of treewidth
 k

$O(\text{EXP}(k) \cdot |I|)$

$O(f(q,k))$

$O(|A| \cdot |T|)$

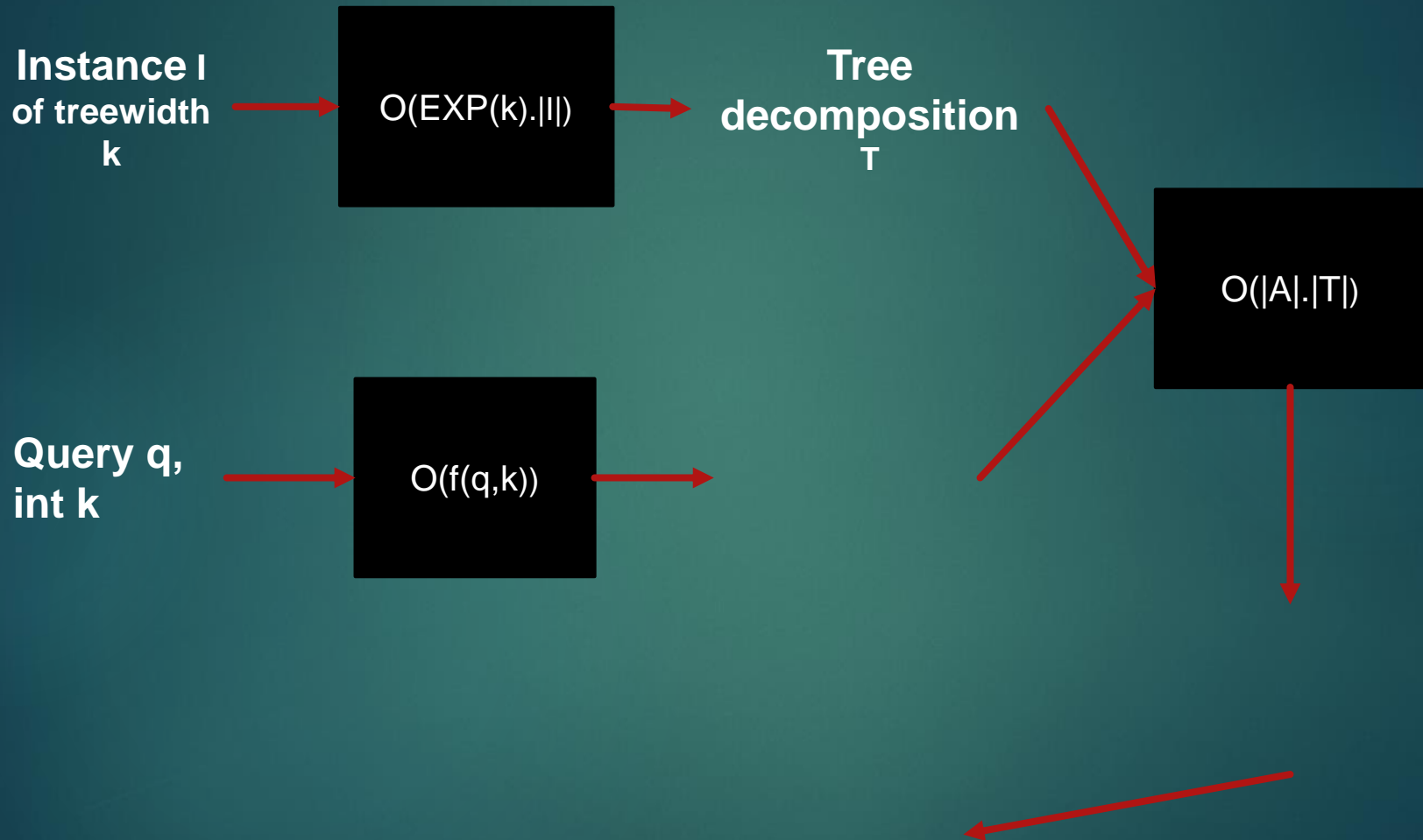
Instance I
of treewidth
 k

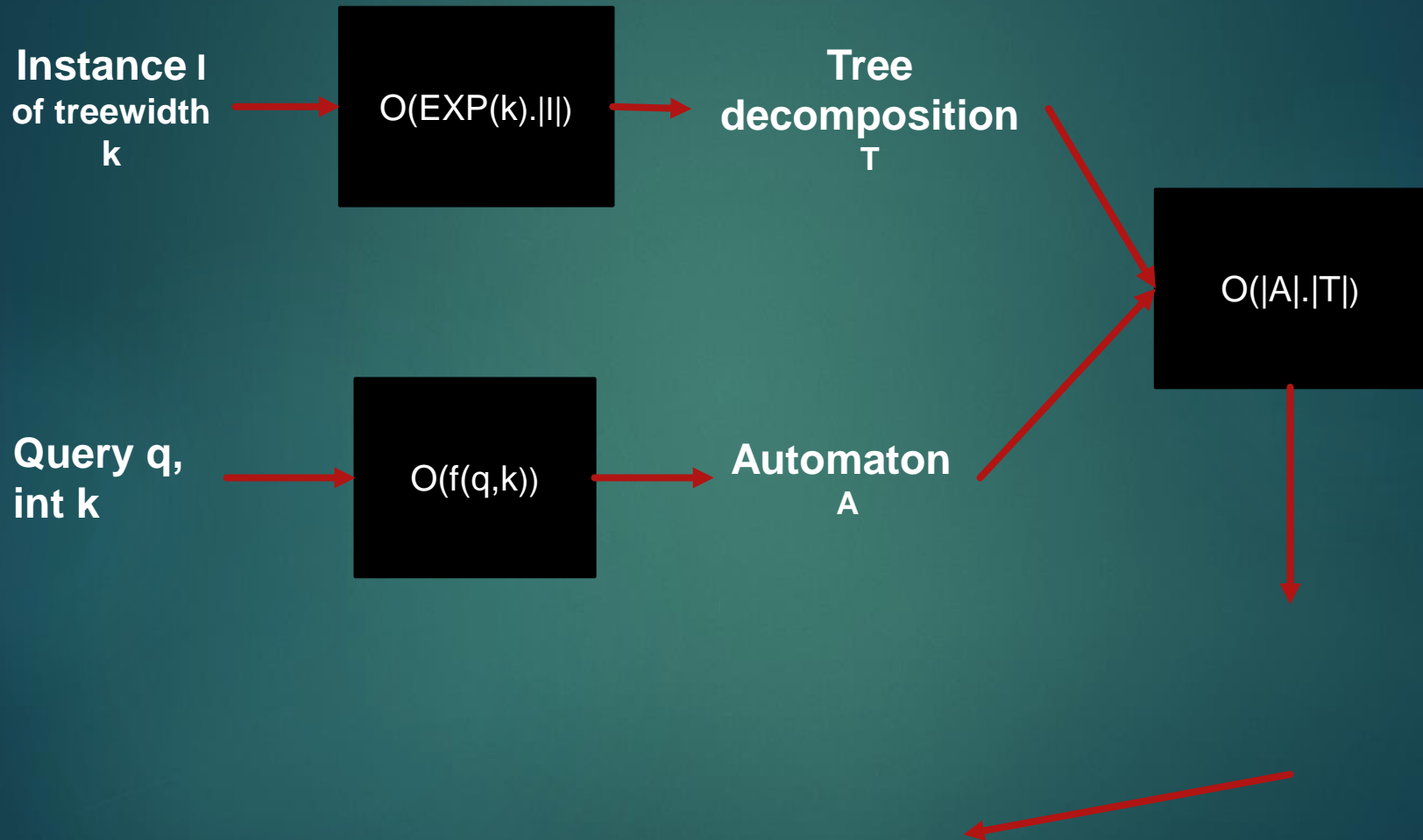
$O(\text{EXP}(k) \cdot |I|)$

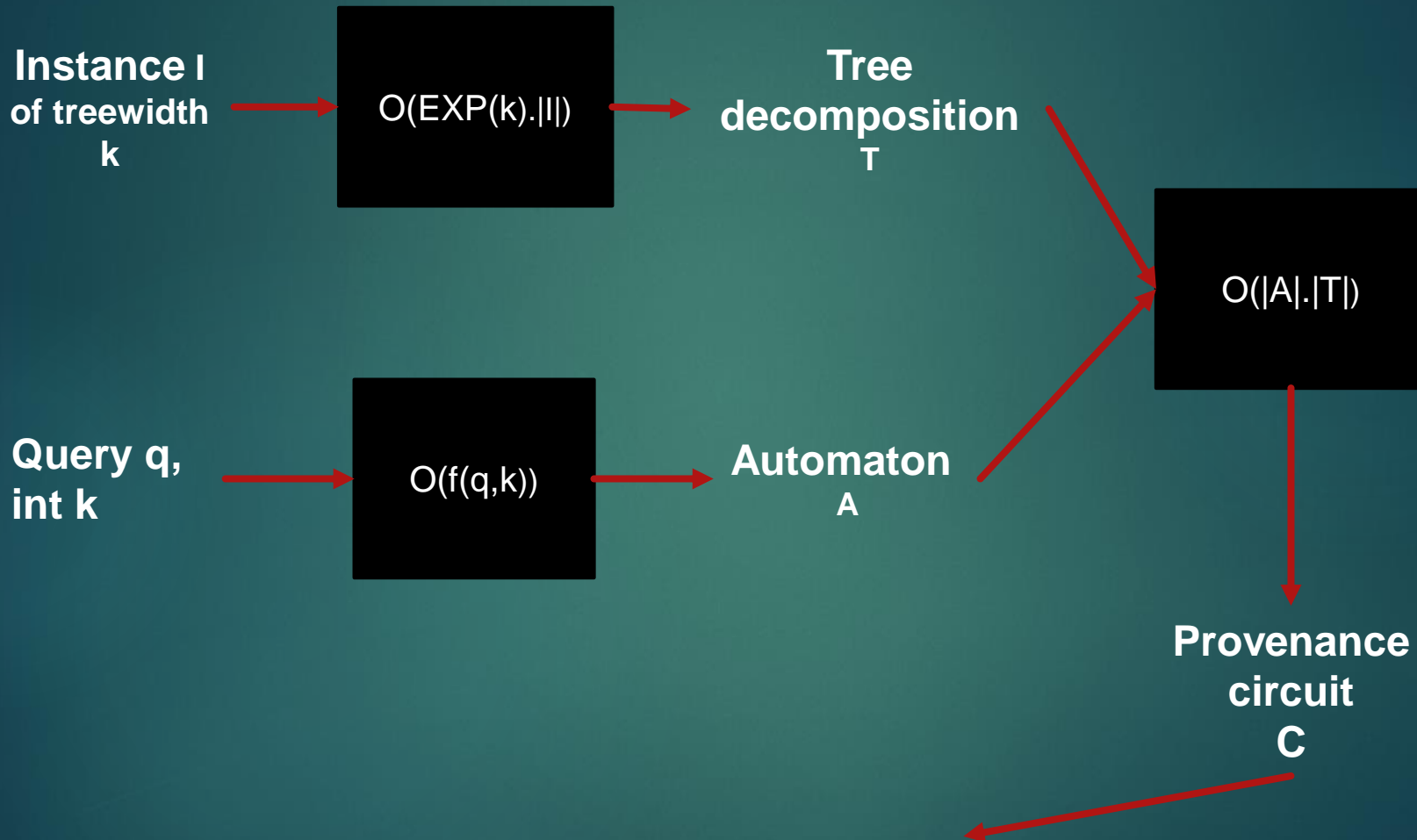
Tree
decomposition
 T

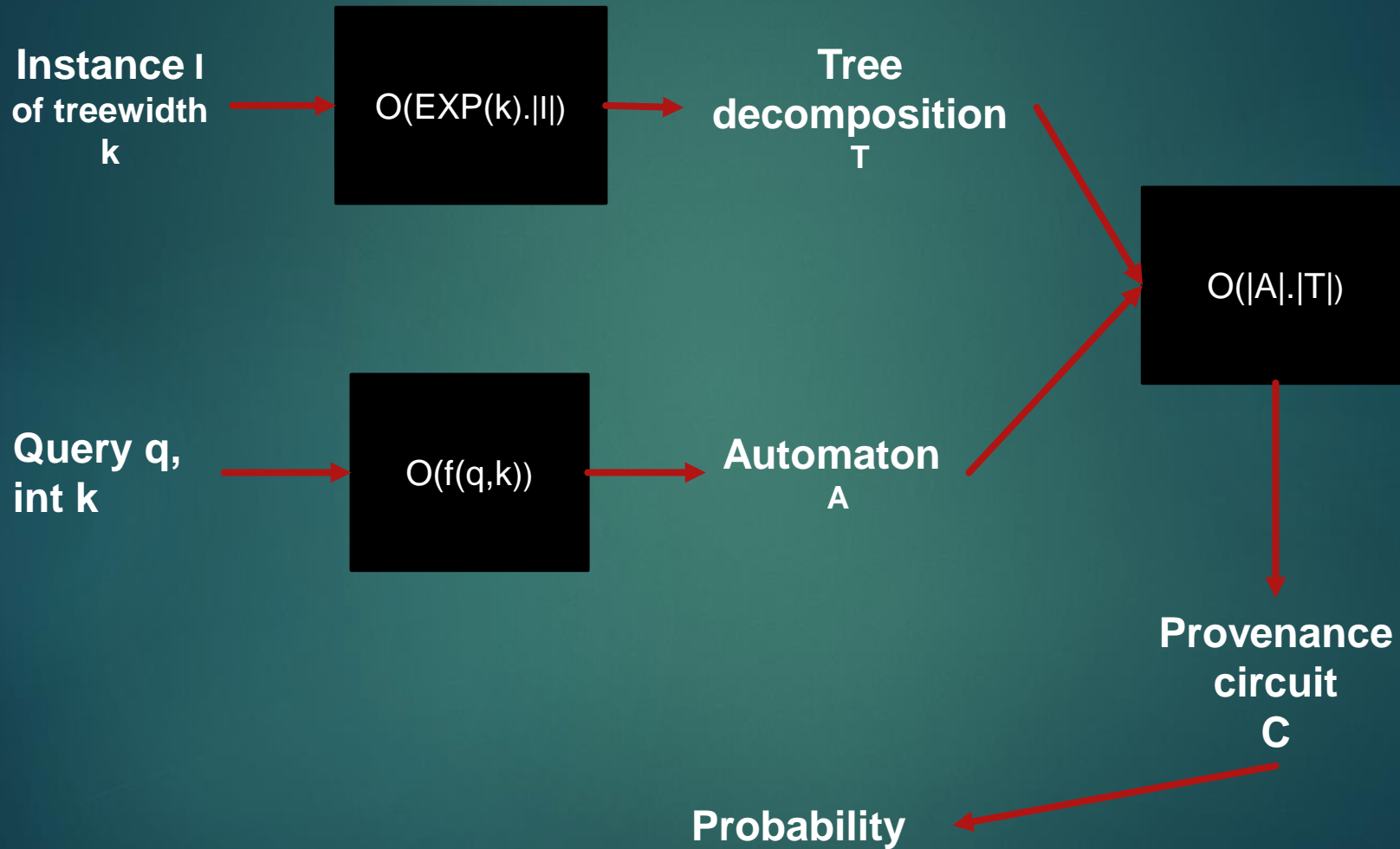
$O(f(q,k))$

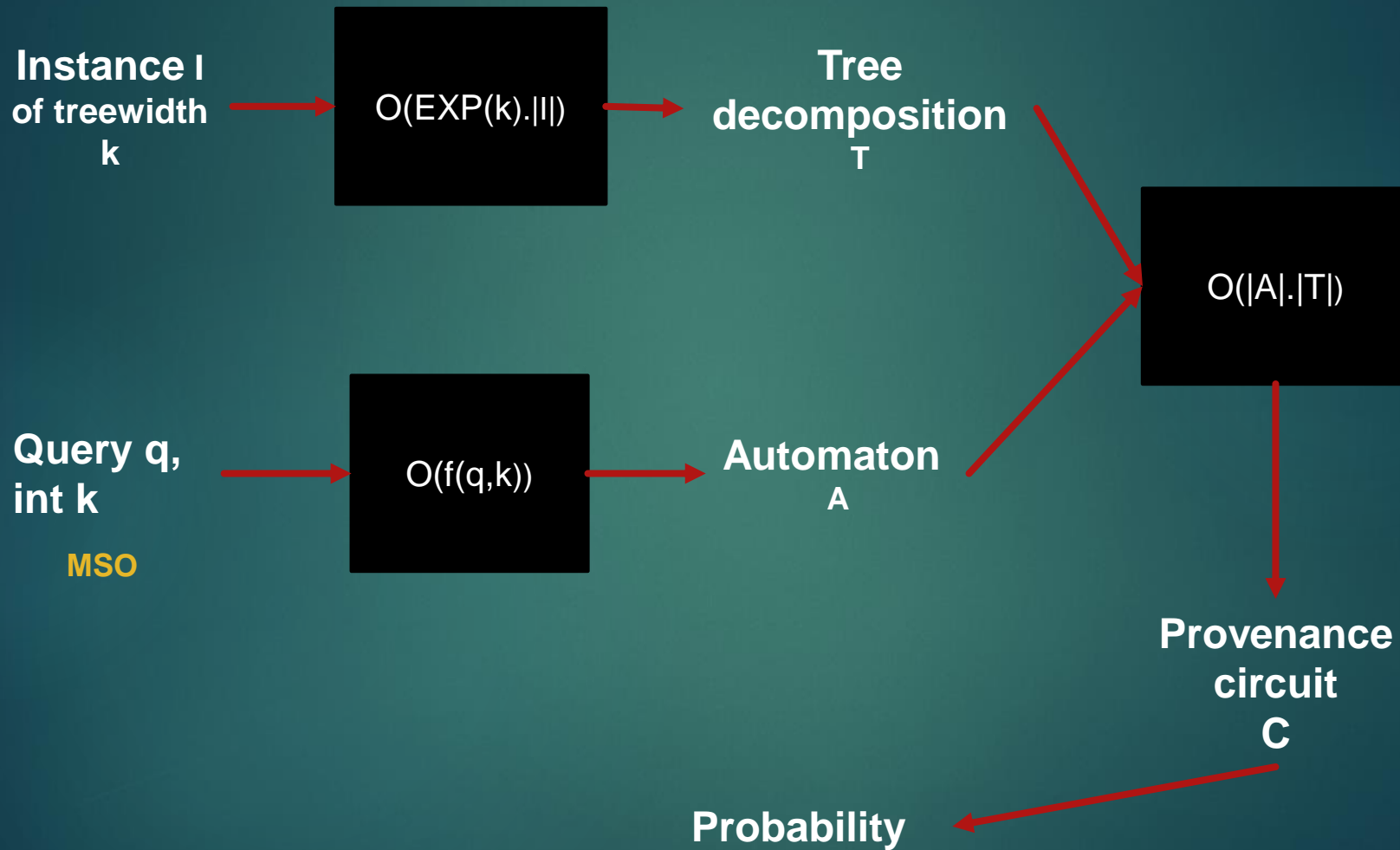
$O(|A| \cdot |T|)$

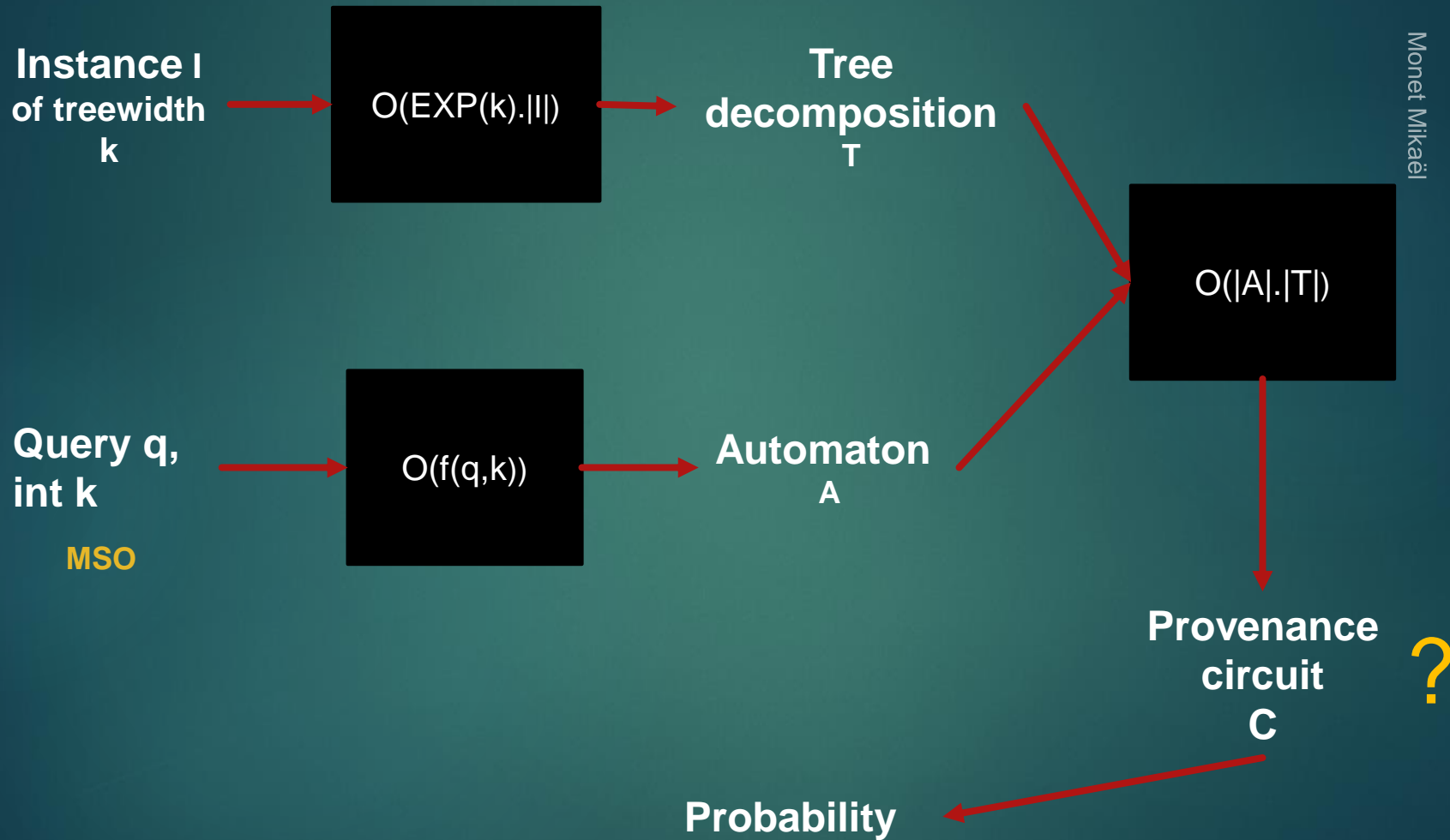


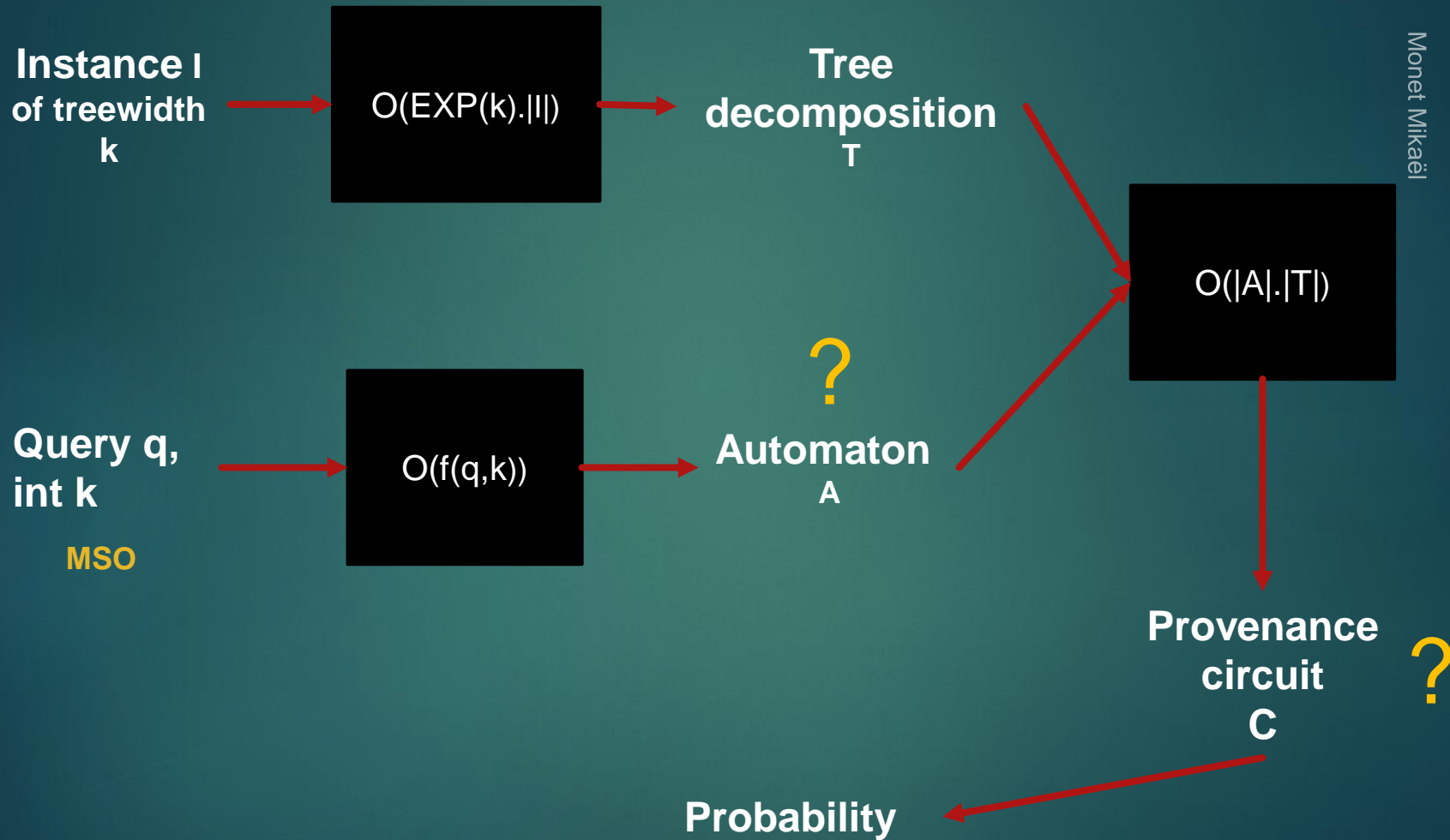












Provenance circuit C of query Q on instance I

12

Provenance circuit C of query Q on instance I

- ▶ Boolean circuit (AND, OR, NOT gates)

Provenance circuit C of query Q on instance I

- ▶ Boolean circuit (AND, OR, NOT gates)
- ▶ Inputs = the facts of I

Provenance circuit C of query Q on instance I

- ▶ Boolean circuit (AND, OR, NOT gates)
- ▶ Inputs = the facts of I

- ▶ For every $v : I \rightarrow \{\text{true}, \text{false}\}$

$$v(I) \models Q \text{ iff } v(C) = 1$$

Tree automata

13

Monet Mikael

- ▶ A bottom-up deterministic tree automaton on $\{a, b\}$ -trees is a tuple $A = (Q, F, \iota, \delta)$ where :
- ▶ Q : finite set of states
- ▶ $F \subseteq Q$: accepting states
- ▶ $\iota : \{a, b\} \rightarrow Q$, determining state for the leaves
- ▶ $\delta : \{a, b\} \times Q^2 \rightarrow Q$, determining the state for internal nodes

Run of an automaton on a tree

14

▶ $Q = \{\text{red}, \text{yellow}, \text{green}\}$

Run of an automaton on a tree

14

▶ $Q = \{\text{red}, \text{yellow}, \text{green}\}$

▶ $F = \{\text{green}\}$

Run of an automaton on a tree

14

- ▶ $Q = \{\text{red}, \text{yellow}, \text{green}\}$
- ▶ $F = \{\text{green}\}$
- ▶ $\iota = \{(a, \text{red}), (b, \text{yellow})\}$

Run of an automaton on a tree

14

Monet Mikael

- ▶ $Q = \{\text{red}, \text{yellow}, \text{green}\}$
- ▶ $F = \{\text{green}\}$
- ▶ $\iota = \{(a, \text{red}), (b, \text{yellow})\}$

$\delta =$

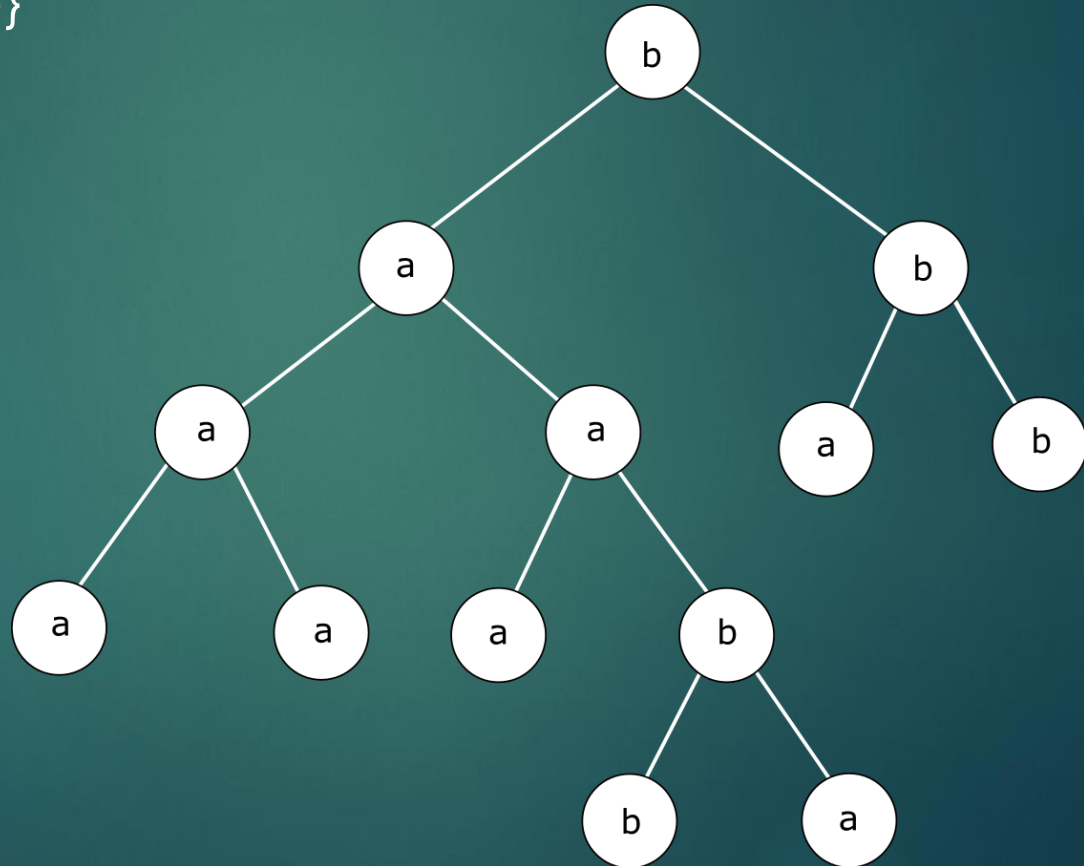
lab	q1	q2	out
a	red	red	red
a	yellow	?	green
a	?	yellow	green
a	green	?	green
a	?	green	green
b	red	red	yellow
b	red	yellow	yellow
b	yellow	red	yellow
b	yellow	yellow	yellow
b	green	?	green
b	?	green	green

Run of an automaton on a tree

- ▶ $Q = \{\text{red}, \text{yellow}, \text{green}\}$
- ▶ $F = \{\text{green}\}$
- ▶ $\iota = \{(a, \text{red}), (b, \text{yellow})\}$

$\delta =$

lab	q1	q2	out
a	red	red	red
a	yellow	?	green
a	?	yellow	green
a	green	?	green
a	?	green	green
b	red	red	yellow
b	red	yellow	yellow
b	yellow	red	yellow
b	yellow	yellow	yellow
b	green	?	green
b	?	green	green

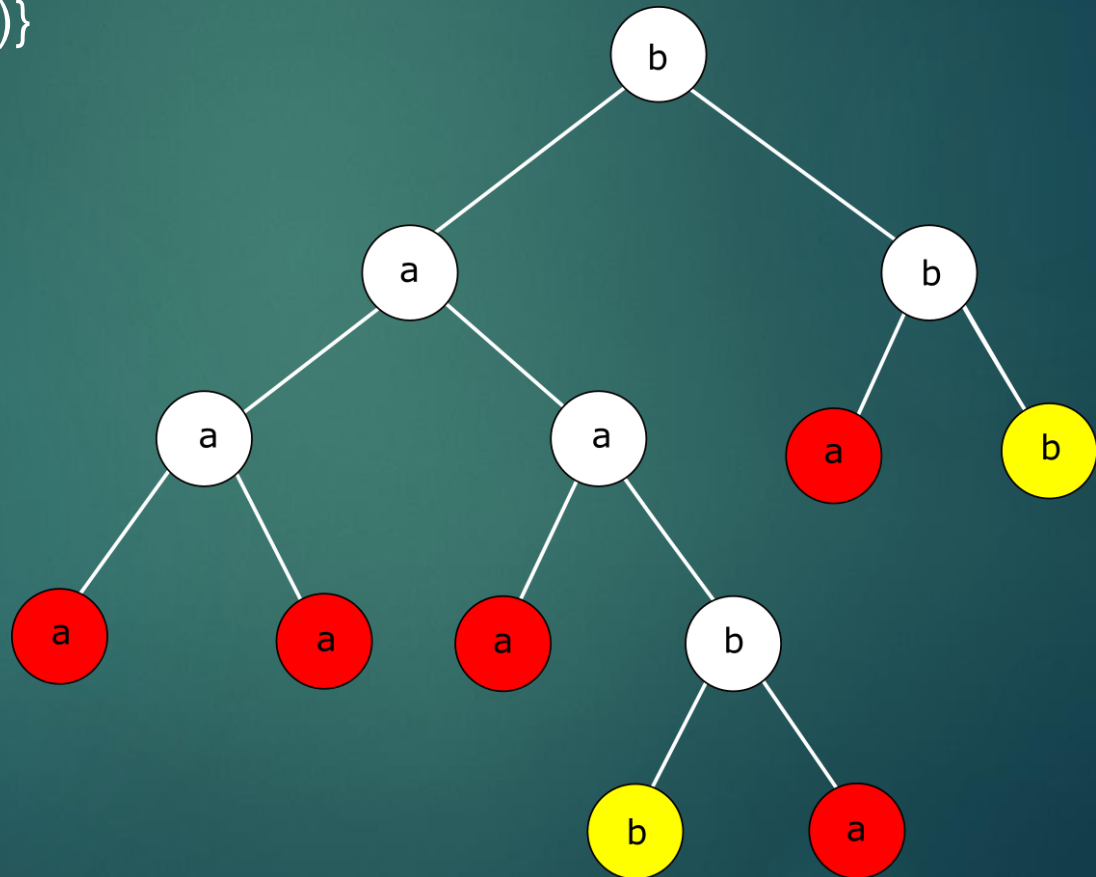


Initialization of the leaves

- ▶ $Q = \{\text{red}, \text{yellow}, \text{green}\}$
- ▶ $F = \{\text{green}\}$
- ▶ $\iota = \{(a, \text{red}), (b, \text{yellow})\}$

 $\delta =$

lab	q1	q2	out
a	red	red	red
a	yellow	?	green
a	?	yellow	green
a	green	?	green
a	?	green	green
b	red	red	yellow
b	red	yellow	yellow
b	yellow	red	yellow
b	yellow	yellow	yellow
b	green	?	green
b	?	green	green

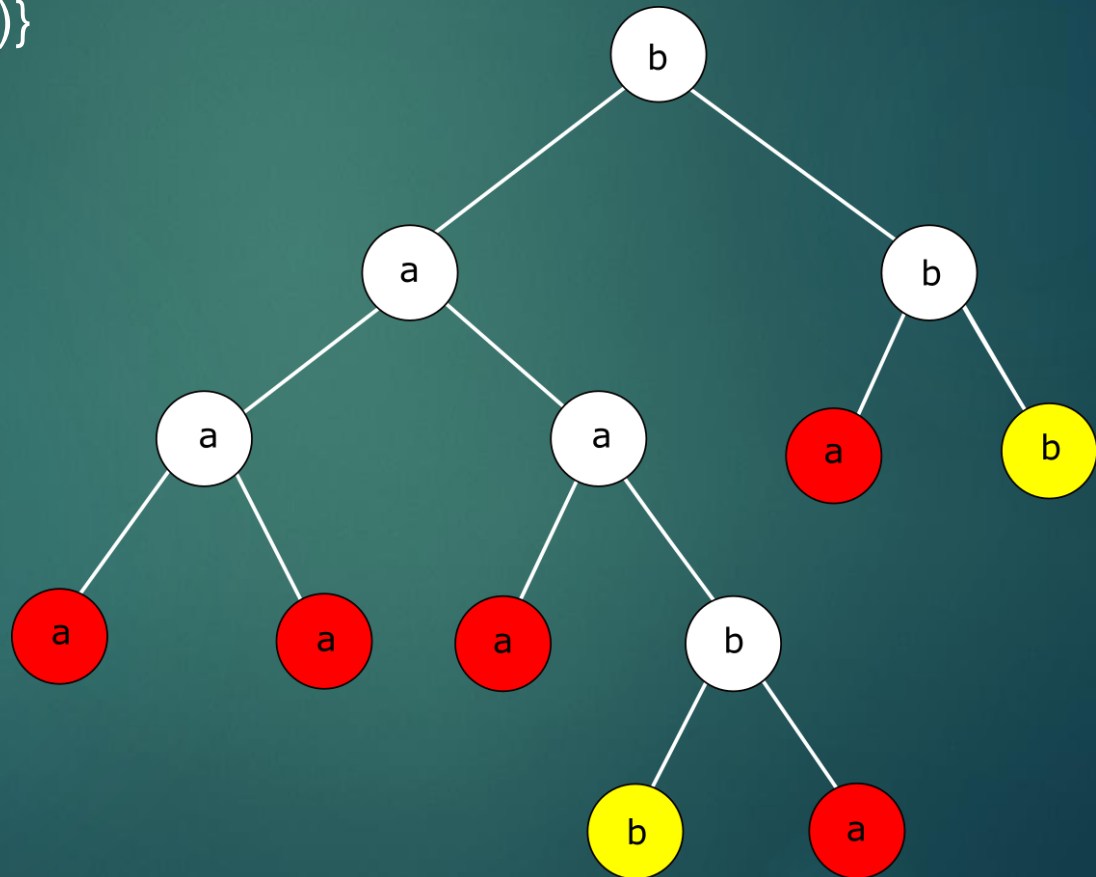


Initialization of the leaves

- ▶ $Q = \{\text{red}, \text{yellow}, \text{green}\}$
- ▶ $F = \{\text{green}\}$
- ▶ $\iota = \{(a, \text{red}), (b, \text{yellow})\}$

$\delta =$

lab	q1	q2	out
a	red	red	red
a	yellow	?	green
a	?	yellow	green
a	green	?	green
a	?	green	green
b	red	red	yellow
b	red	yellow	yellow
b	yellow	red	yellow
b	yellow	yellow	yellow
b	green	?	green
b	?	green	green

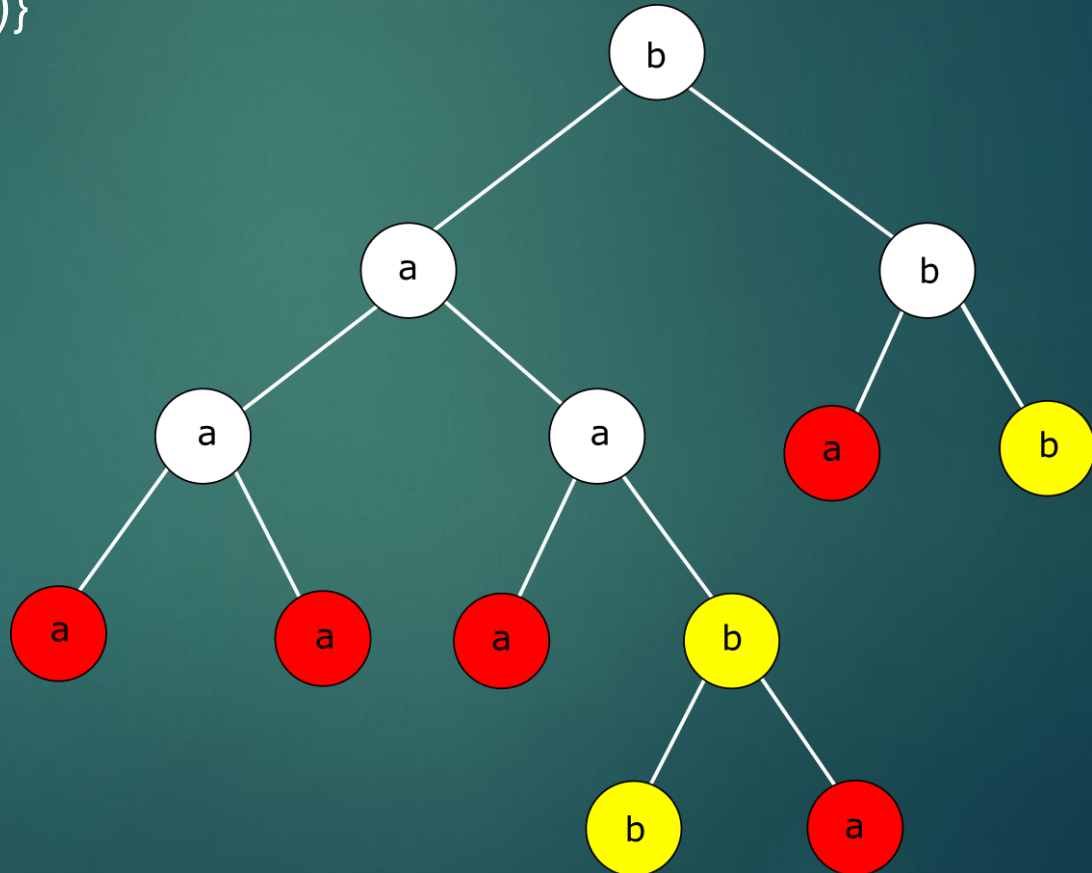


Internal nodes

- ▶ $Q = \{\text{red}, \text{yellow}, \text{green}\}$
- ▶ $F = \{\text{green}\}$
- ▶ $\iota = \{(a, \text{red}), (b, \text{yellow})\}$

$\delta =$

lab	q1	q2	out
a	red	red	red
a	yellow	?	green
a	?	yellow	green
a	green	?	green
a	?	green	green
b	red	red	yellow
b	red	yellow	yellow
b	yellow	red	yellow
b	yellow	yellow	yellow
b	green	?	green
b	?	green	green

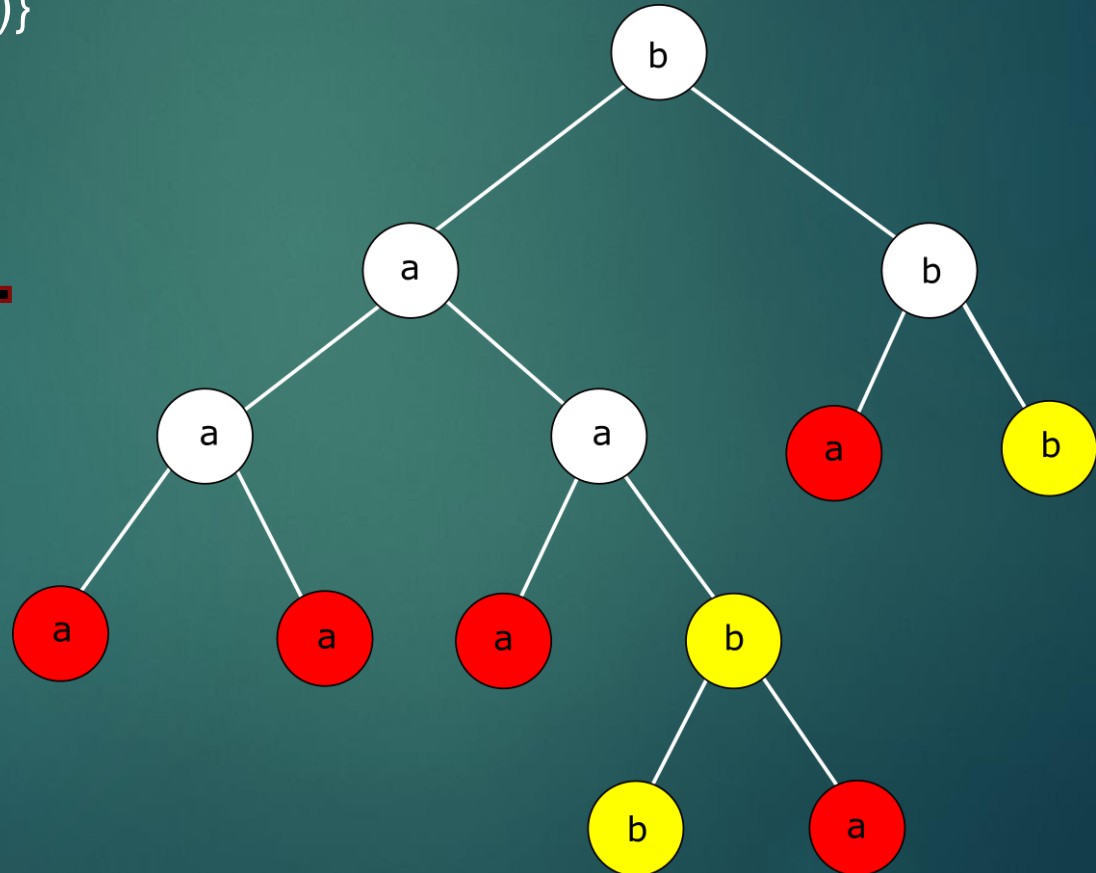


Internal nodes

- ▶ $Q = \{\text{red}, \text{yellow}, \text{green}\}$
- ▶ $F = \{\text{green}\}$
- ▶ $\iota = \{(a, \text{red}), (b, \text{yellow})\}$

$\delta =$

lab	q1	q2	out
a	red	red	red
a	yellow	?	green
a	?	yellow	green
a	green	?	green
a	?	green	green
b	red	red	yellow
b	red	yellow	yellow
b	yellow	red	yellow
b	yellow	yellow	yellow
b	green	?	green
b	?	green	green

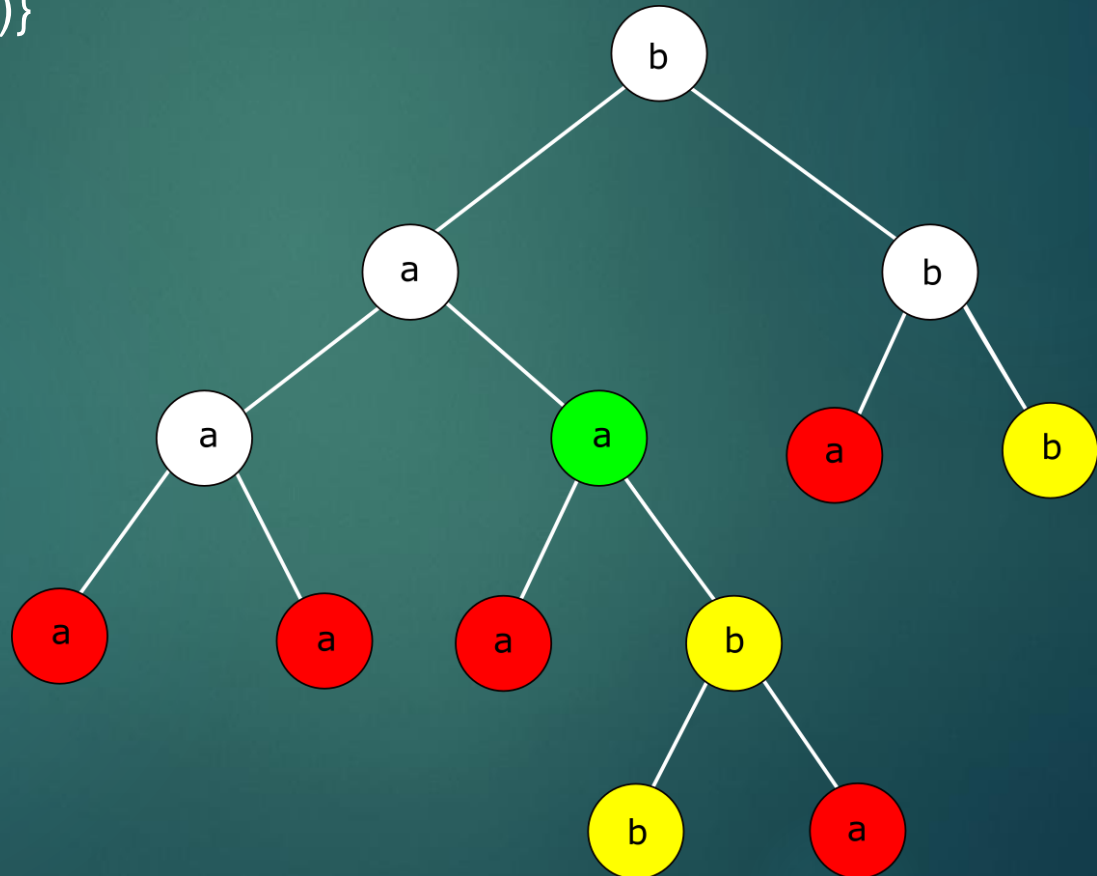


And so on...

- ▶ $Q = \{\text{red}, \text{yellow}, \text{green}\}$
- ▶ $F = \{\text{green}\}$
- ▶ $\iota = \{(a, \text{red}), (b, \text{yellow})\}$

$\delta =$

lab	q1	q2	out
a	red	red	red
a	yellow	?	green
a	?	yellow	green
a	green	?	green
a	?	green	green
b	red	red	yellow
b	red	yellow	yellow
b	yellow	red	yellow
b	yellow	yellow	yellow
b	green	?	green
b	?	green	green



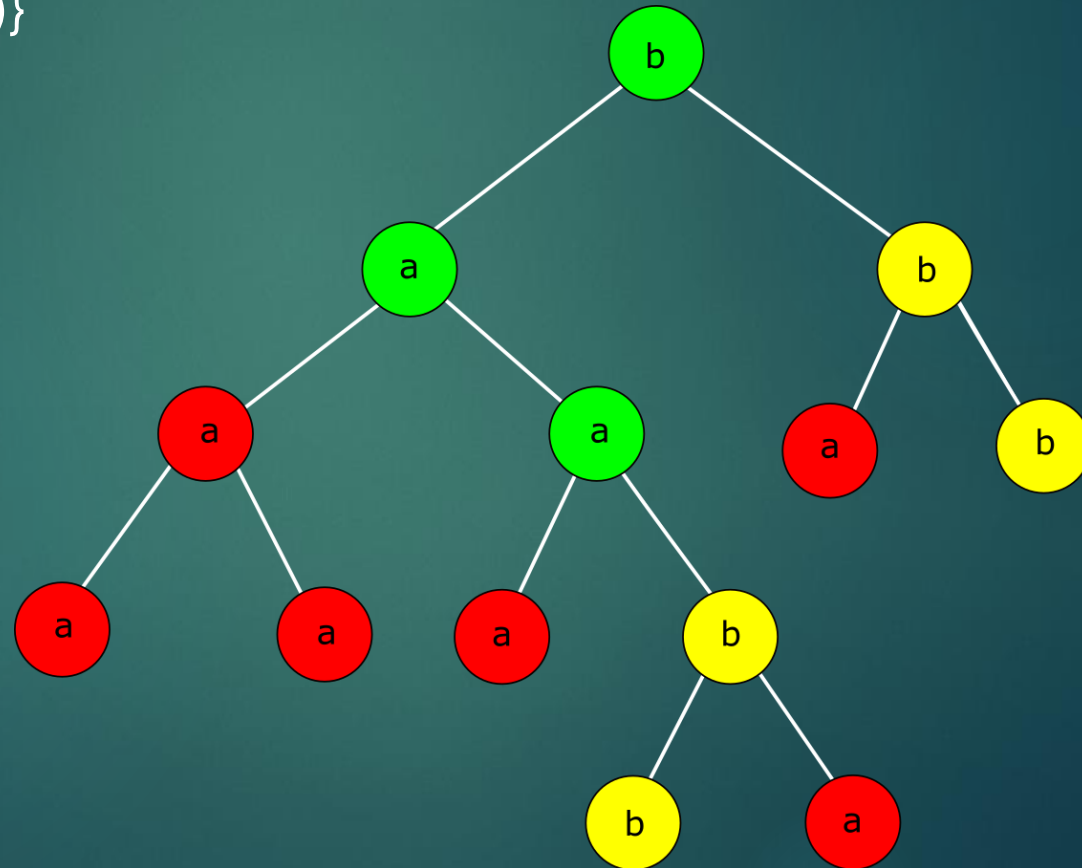
This tree is in the language of A

14

- ▶ $Q = \{\text{red}, \text{yellow}, \text{green}\}$
- ▶ $F = \{\text{green}\}$
- ▶ $\iota = \{(a, \text{red}), (b, \text{yellow})\}$

$\delta =$

lab	q1	q2	out
a	red	red	red
a	yellow	?	green
a	?	yellow	green
a	green	?	green
a	?	green	green
b	red	red	yellow
b	red	yellow	yellow
b	yellow	red	yellow
b	yellow	yellow	yellow
b	green	?	green
b	?	green	green



Major drawbacks

15

Monet Mikael

- ▶ In general, computing the automaton has non-elementary complexity in the query

Major drawbacks

15

Monet Mikael

- ▶ In general, computing the automaton has non-elementary complexity in the query
- ▶ Exponential dependence in the instance treewidth

Major drawbacks

15

Monet Mikael

- ▶ In general, computing the automaton has non-elementary complexity in the query
- ▶ Exponential dependence in the instance treewidth

- ▶ Natural question: restrict queries to obtain tractable combined complexity of PQE on bounded treewidth instances?

Bad news...

- ▶ We proved that:

Path queries on tree instances (treewidth = 1) is already #P-hard. (reduction from #MONOTONE-2-SAT)

Bad news...

- ▶ We proved that:

Path queries on tree instances (treewidth = 1) is already #P-hard. (reduction from #MONOTONE-2-SAT)

- ▶ What to do now?

Lower ambitions

17

Monet Mikael

- ▶ Restrict queries to obtain tractable combined complexity of probabilistic query evaluation on bounded treewidth instances?

Lower ambitions

17

Monet Mikael

- ▶ Restrict queries to obtain tractable combined complexity of probabilistic query evaluation on bounded treewidth instances?
- ▶ We now aim at a tractable combined complexity for deterministic query evaluation: which queries, which automata, which provenance representation?

Two-way Alternating Tree Automata

18

Monet Mikael

- ▶ Can navigate the tree in every direction, can launch simultaneous runs

Two-way Alternating Tree Automata

18

Monet Mikael

- ▶ Can navigate the tree in every direction, can launch simultaneous runs

$$\delta(a,q) = (q,l) \vee [\quad \wedge \quad]$$

Two-way Alternating Tree Automata

18

Monet Mikael

- ▶ Can navigate the tree in every direction, can launch simultaneous runs

$$\delta(a,q) = (q,l) \vee [(q,p) \wedge \quad]$$

Two-way Alternating Tree Automata

18

Monet Mikael

- ▶ Can navigate the tree in every direction, can launch simultaneous runs

$$\delta(a,q) = (q,l) \vee [(q,p) \wedge (q',r)]$$

Two-way Alternating Tree Automata

18

Monet Mikael

- ▶ Can navigate the tree in every direction, can launch simultaneous runs

$$\delta(a,q) = (q,l) \vee [(q,p) \wedge (q',r)]$$

- ▶ Intuition: less things to remember, more parallelizable

Cycluits

- ▶ Boolean circuits with cycles

Cycluits

19

Monet Mikael

- ▶ Boolean circuits with cycles
- ▶ Least fixed-point semantics

Cycluits

- ▶ Boolean circuits with cycles
- ▶ Least fixed-point semantics
- ▶ Beware of negations!

Cycluits

- ▶ Boolean circuits with cycles
- ▶ Least fixed-point semantics
- ▶ Beware of negations!
- ▶ Linear time evaluation

Cycluits

- ▶ Boolean circuits with cycles
- ▶ Least fixed-point semantics
- ▶ Beware of negations!
- ▶ Linear time evaluation
- ▶ Can be acyclified in quadratic time

Cycluits

- ▶ Boolean circuits with cycles
- ▶ Least fixed-point semantics
- ▶ Beware of negations!
- ▶ Linear time evaluation
- ▶ Can be acyclified in quadratic time
- ▶ Are they more concise?

