

Challenges for Efficient Query Evaluation on Structured Probabilistic Data



SUM2016

SEPTEMBER 23, 2016, NICE

Antoine Amarilli, Silviu Maniu, **Mikaël Monet**

A probabilistic database

2

Amarilli, Antoine, Maniu Silviu, 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

2

Amarilli, Antoine, Maniu Sîrbu, 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

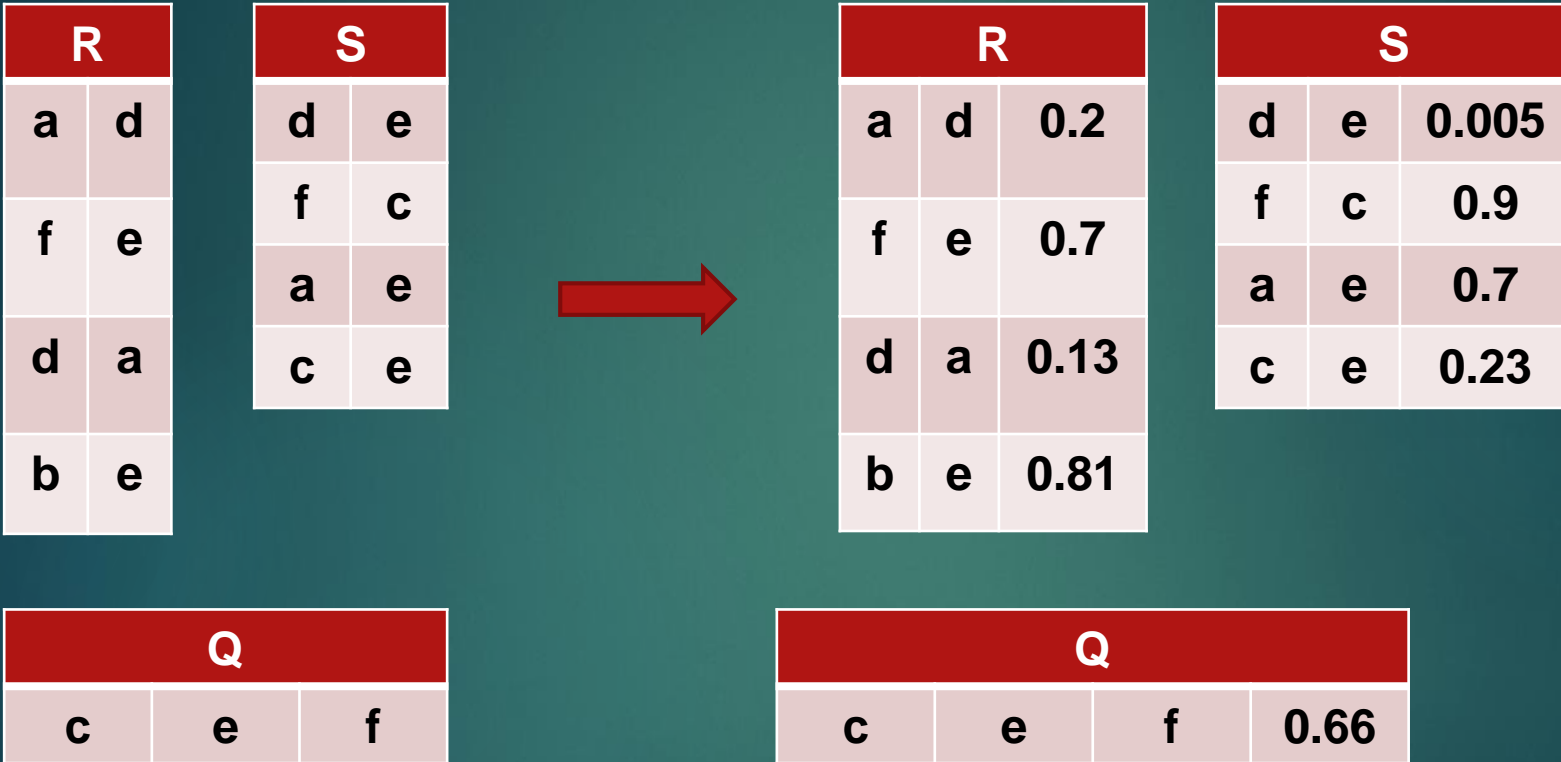
Q		
c	e	f

Q			
c	e	f	0.66

A probabilistic database

2

Amarilli, Antoine, Maniu Sîrbu, 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

Amarilli Antoine, Maniu Silviu, 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

Amarilli Antoine, Maniu Silviu, 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

Amarilli, Antoine, Maniù Sîvîu, Monet Mikaeîl

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)

4

Amarilli, Antoine, Maniu Silviu, Monet Mikael

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} :

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

Probabilistic query evaluation (PQE)

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

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

- ▶ Problem: in general #P-hard

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] show the following dichotomy for any UCQ Q :

2) Restricting the class of queries

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

2) Restricting the class of queries

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

2) Restricting the class of queries

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

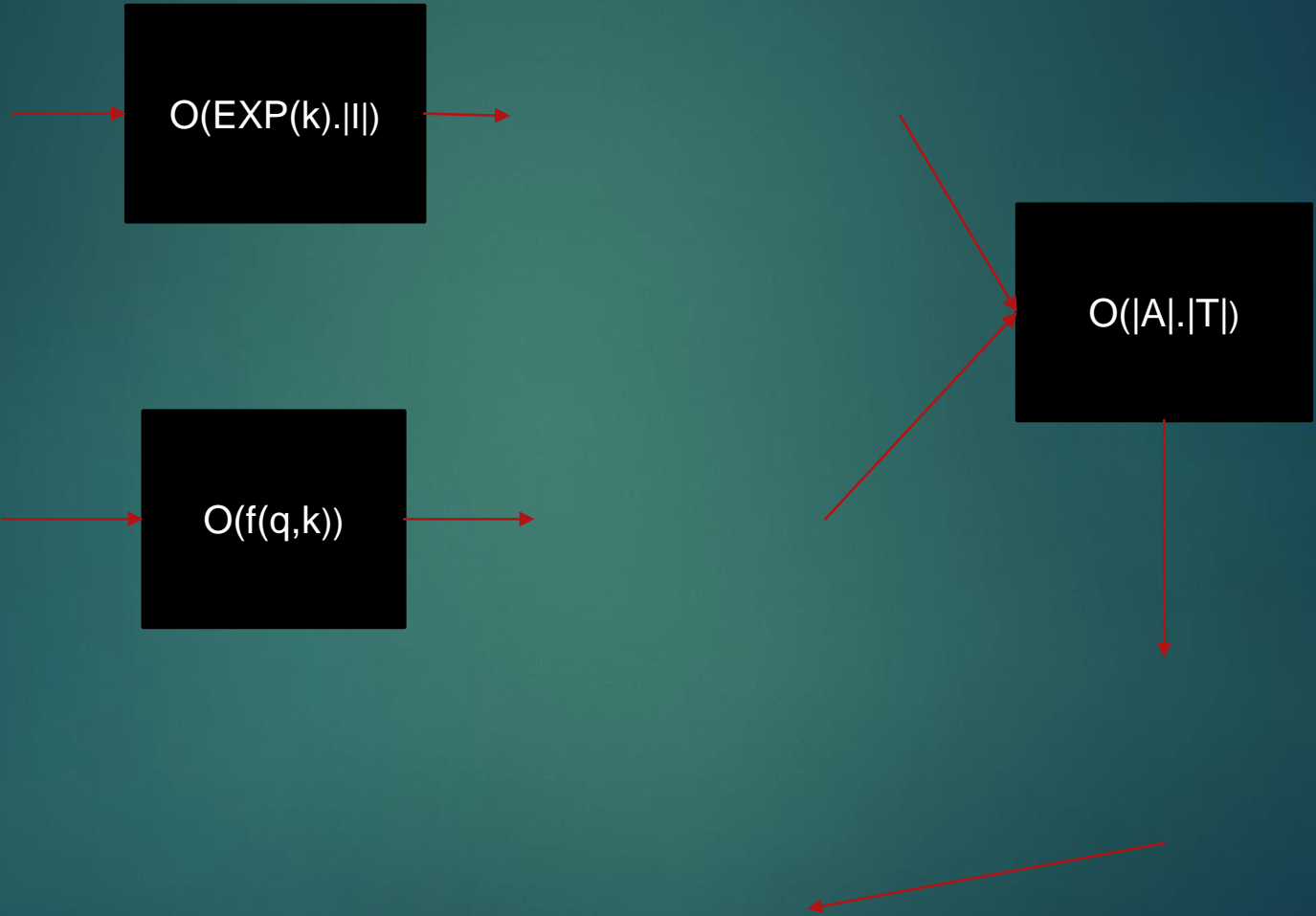
- ▶ Simple conjunctive query $\exists x,y R(x),T(x,y),S(y)$ is already #P-hard !

2) Restricting the class of queries

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

3) Restricting the shape of the instances

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



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

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

$$O(f(q,k))$$

Instance I
of treewidth
 k

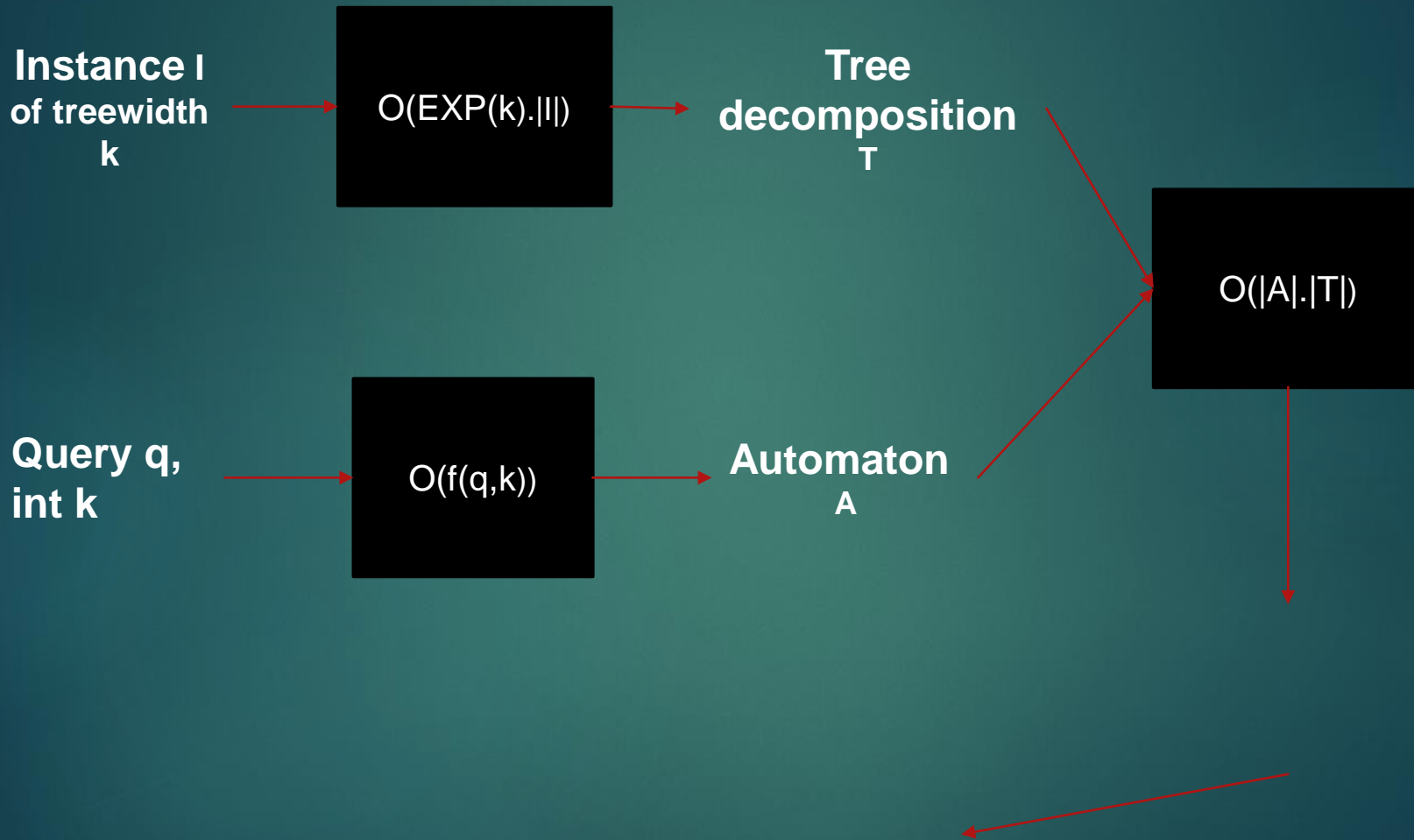
$O(\text{EXP}(k) \cdot |I|)$

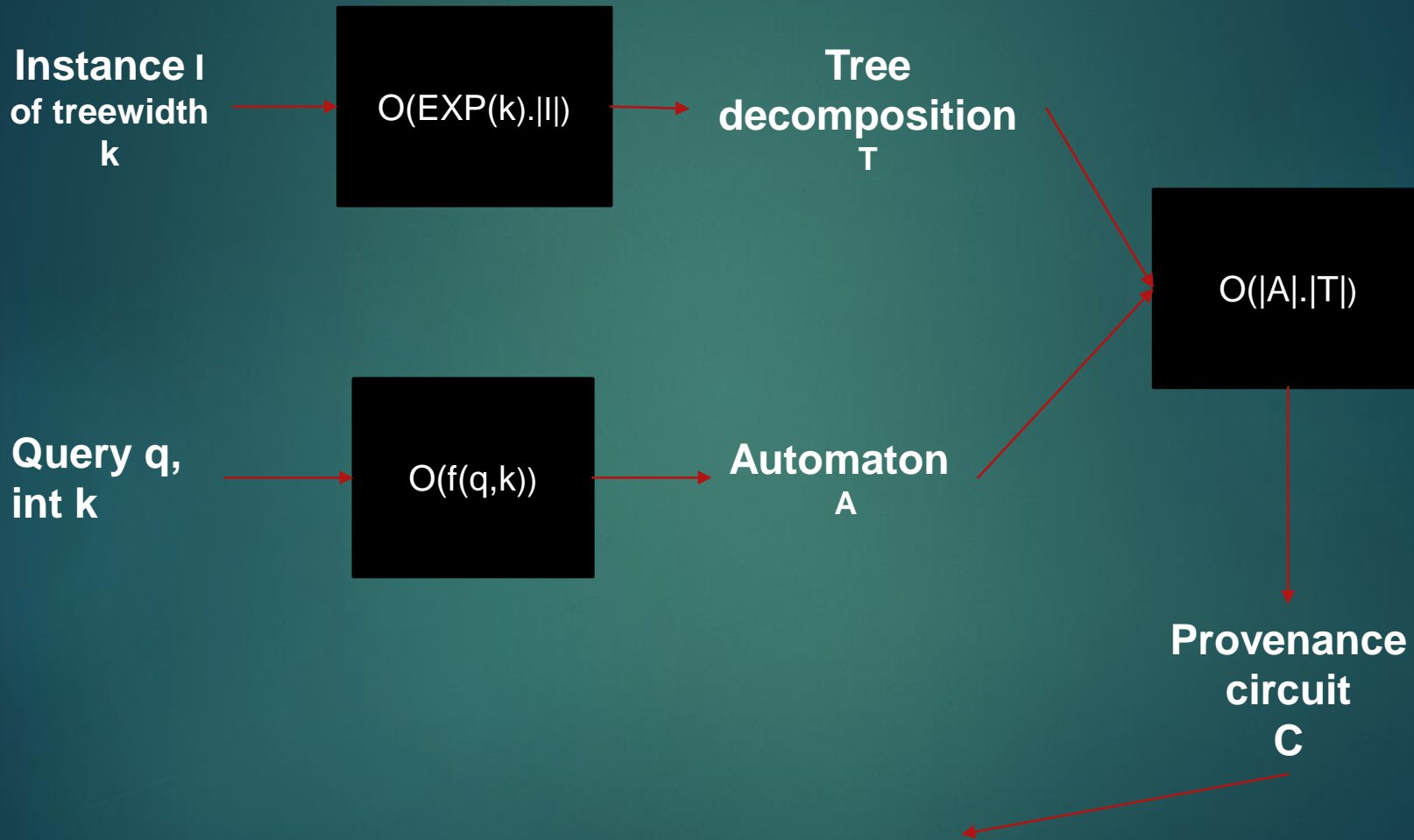
Tree
decomposition
 T

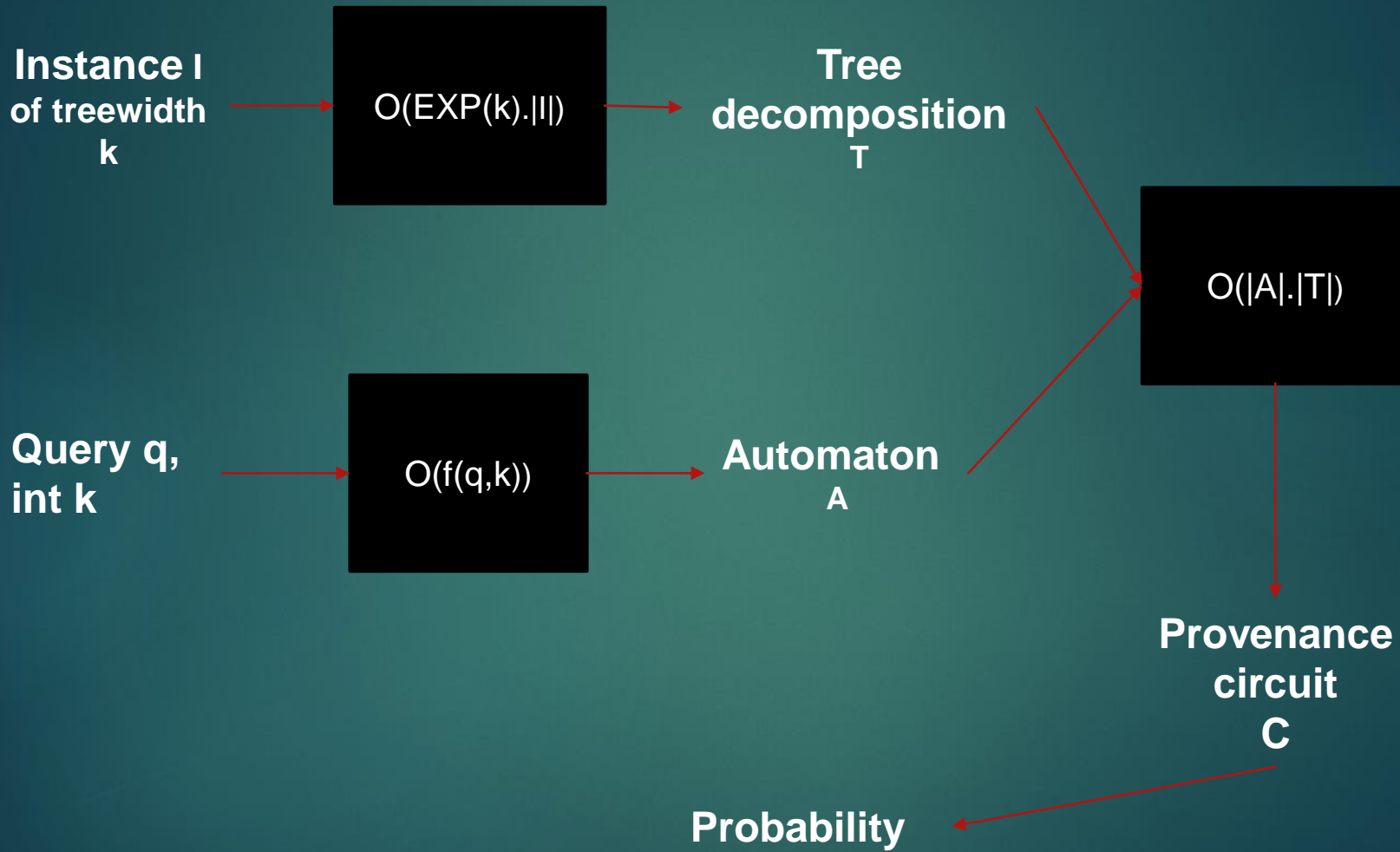
$O(|A| \cdot |T|)$

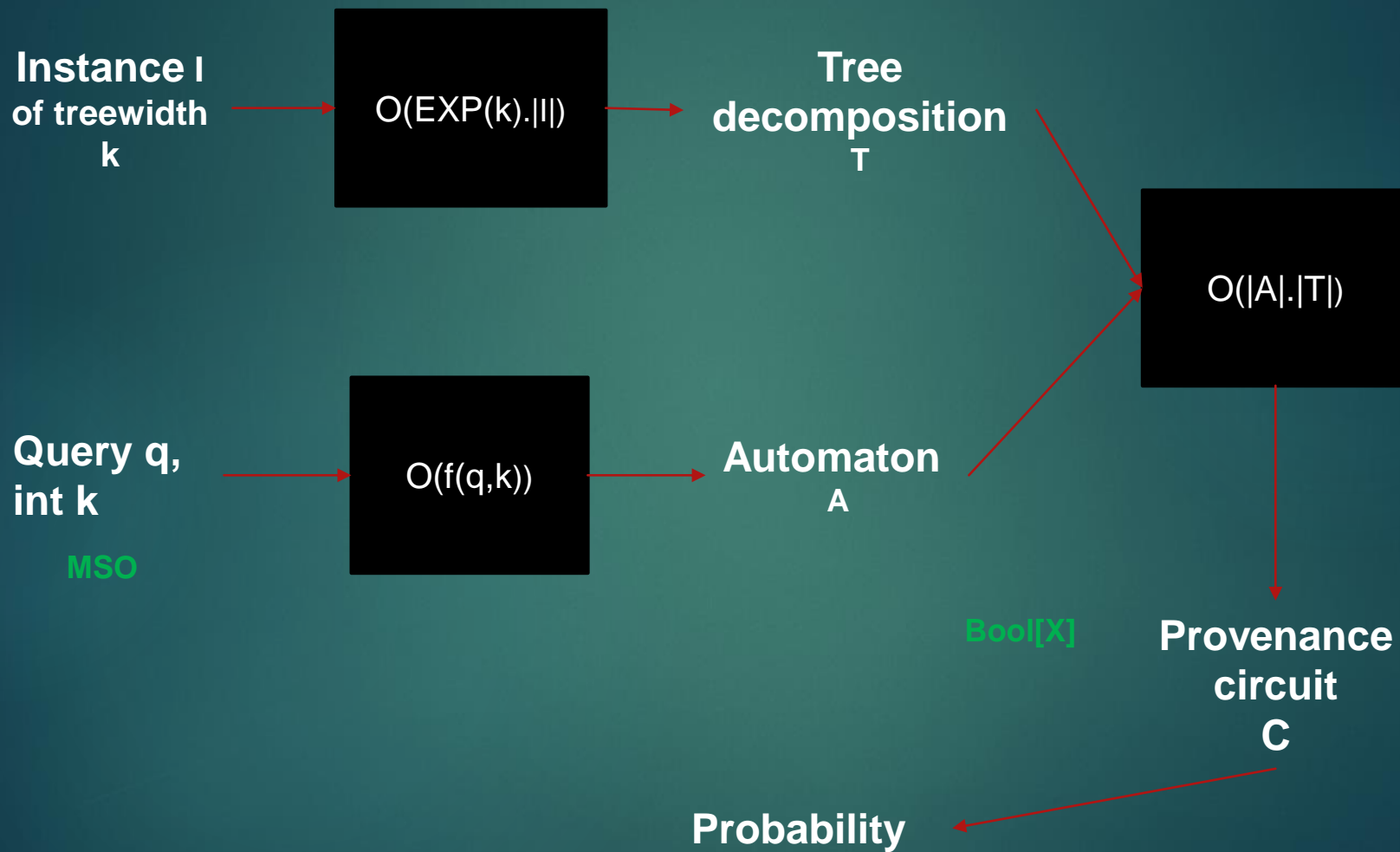
Query q ,
int k

$O(f(q, k))$





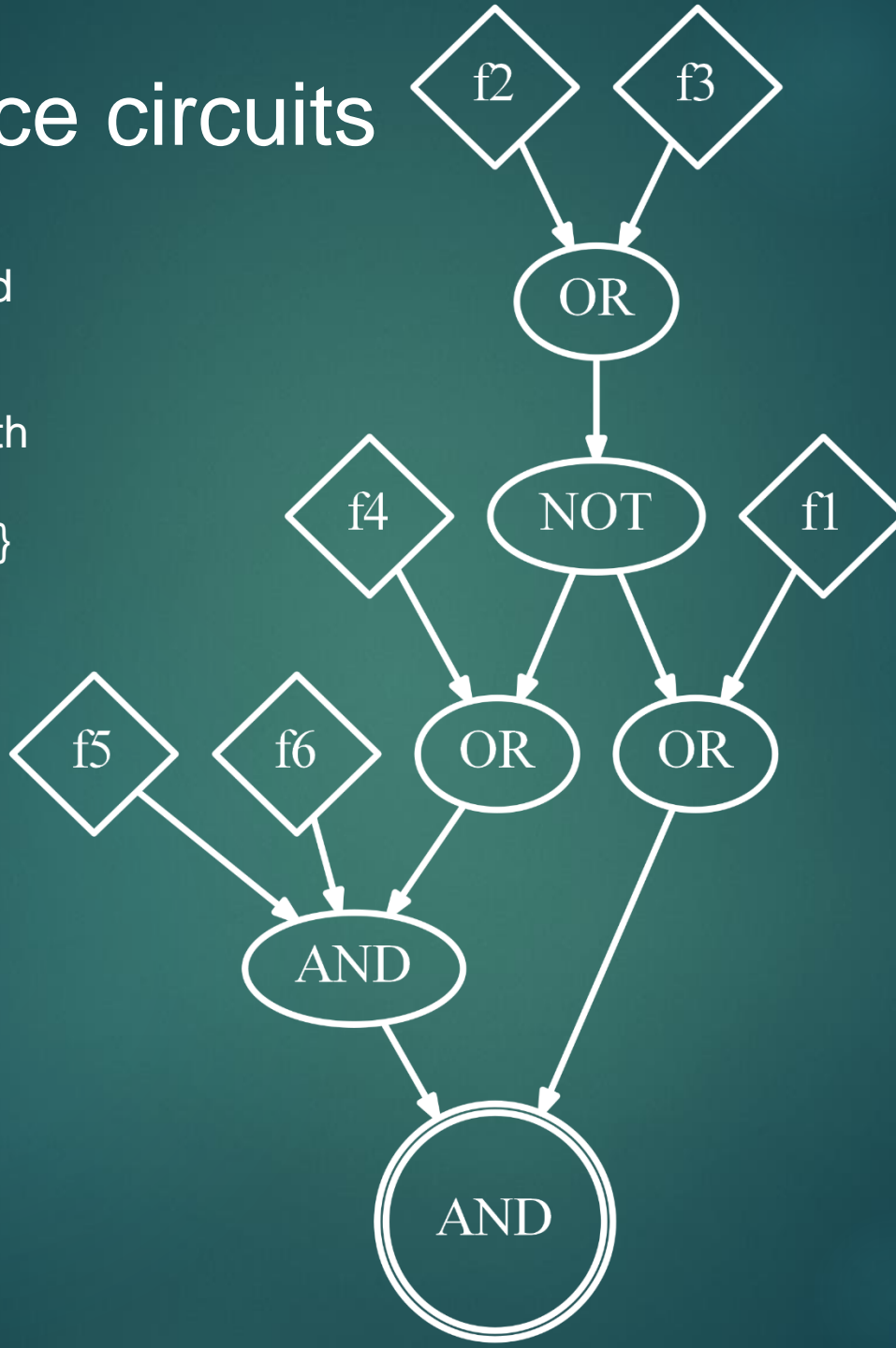




Provenance circuits

Some query q , fixed

Some instance I with facts
 $\{f1, f2, f3, f4, f5, f6\}$

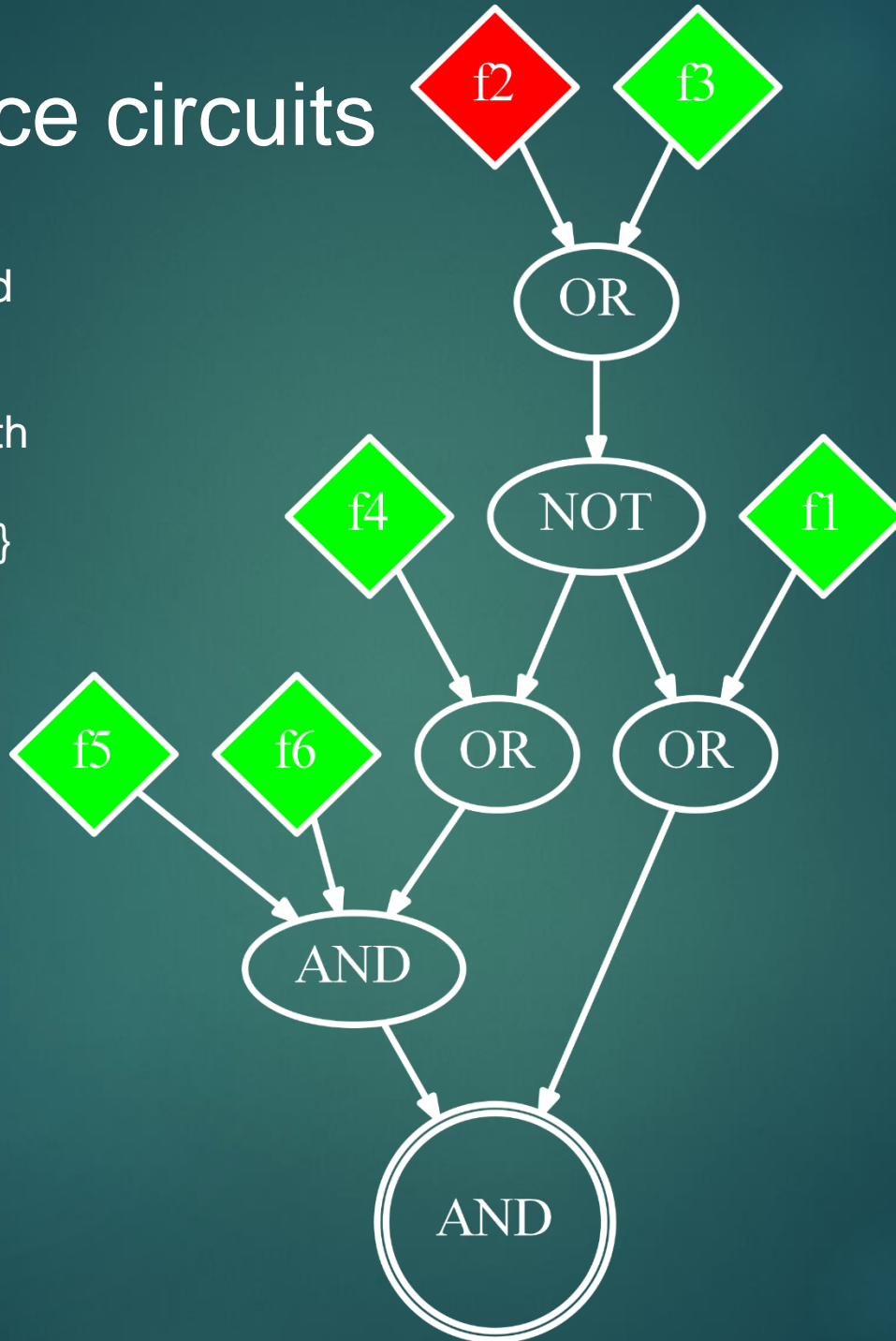


Provenance circuits

Some query q , fixed

Some instance I with facts

$\{f_1, f_2, f_3, f_4, f_5, f_6\}$

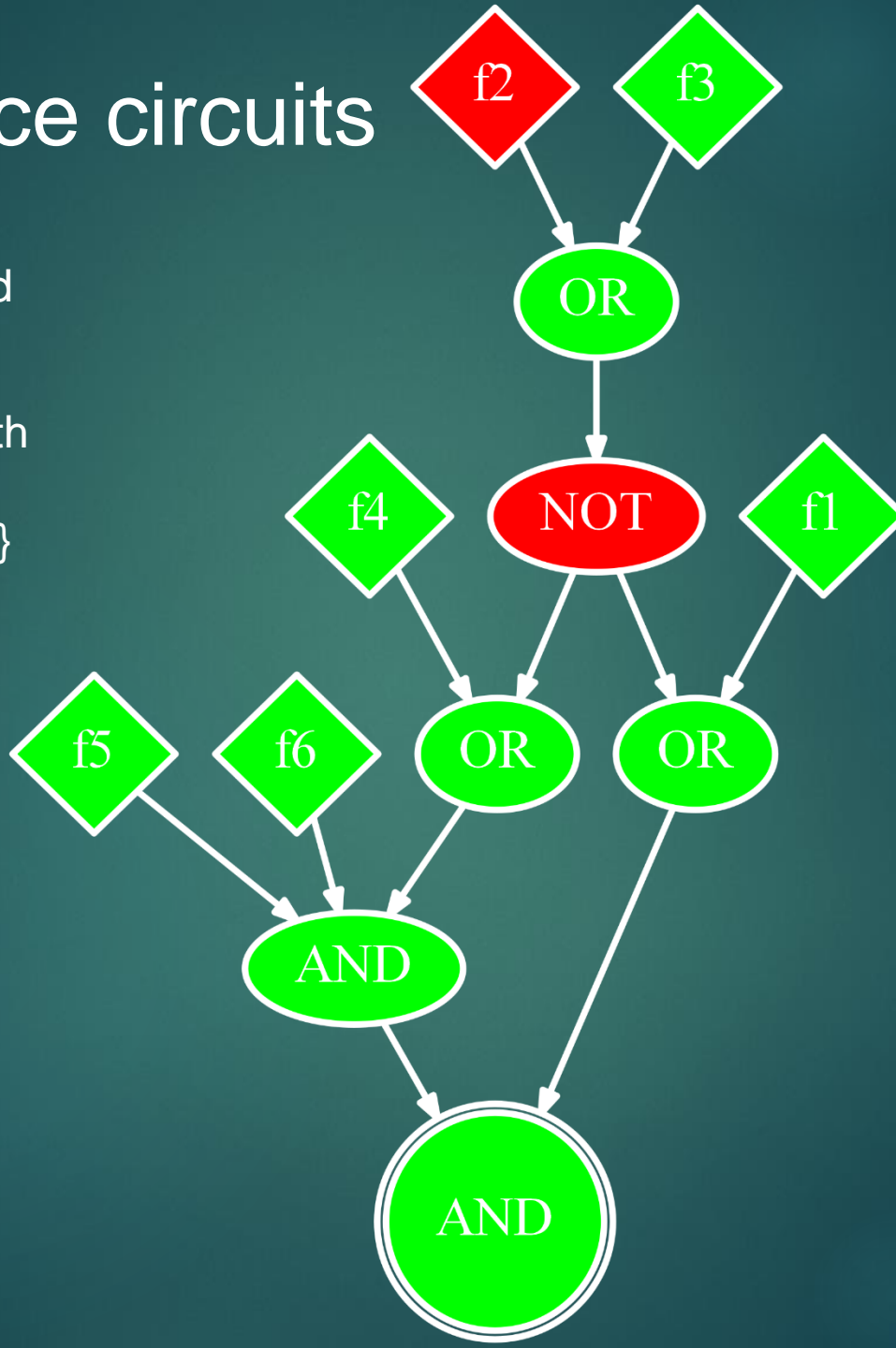


Provenance circuits

Some query q , fixed

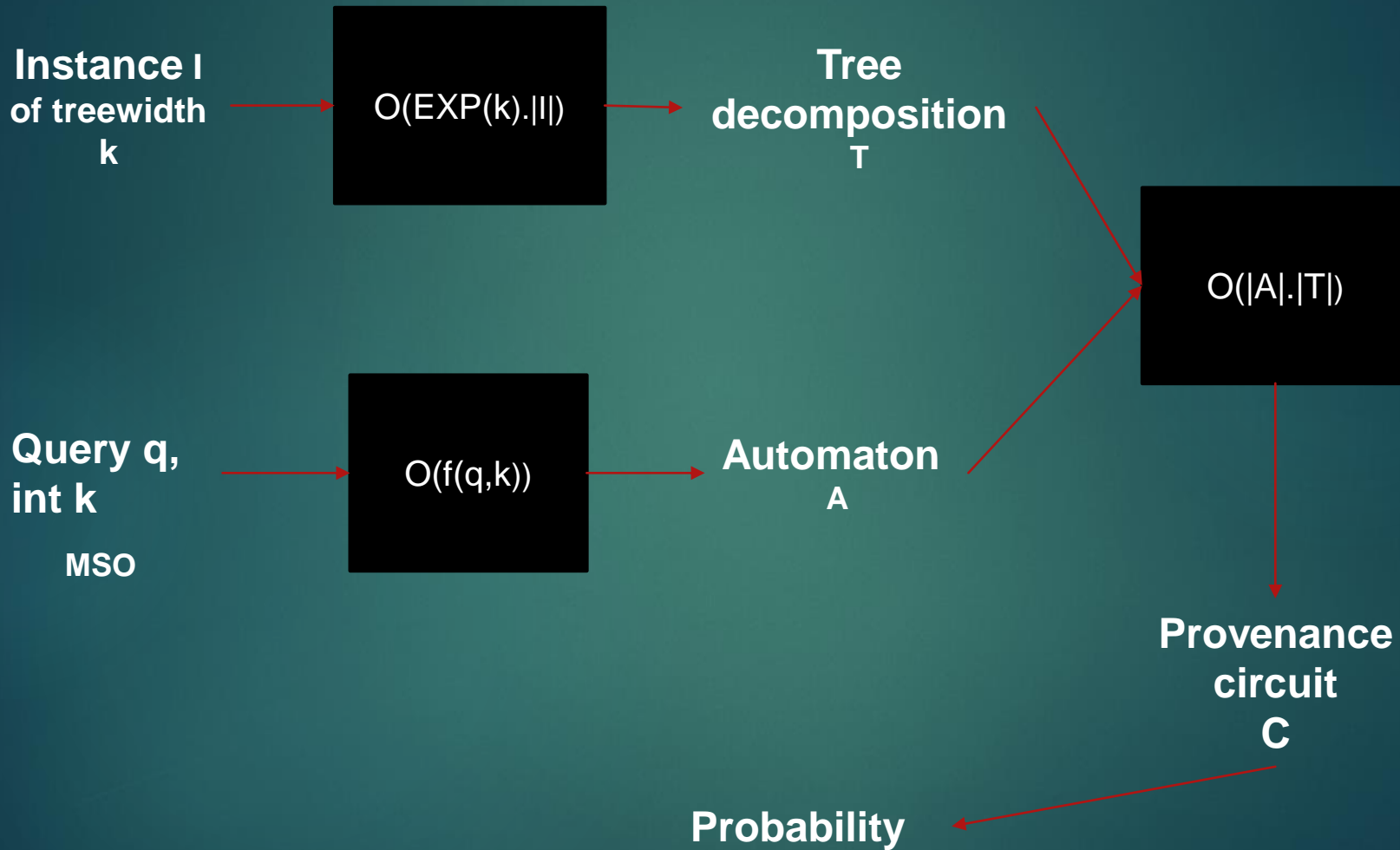
Some instance I with facts

$\{f_1, f_2, f_3, f_4, f_5, f_6\}$



Problems

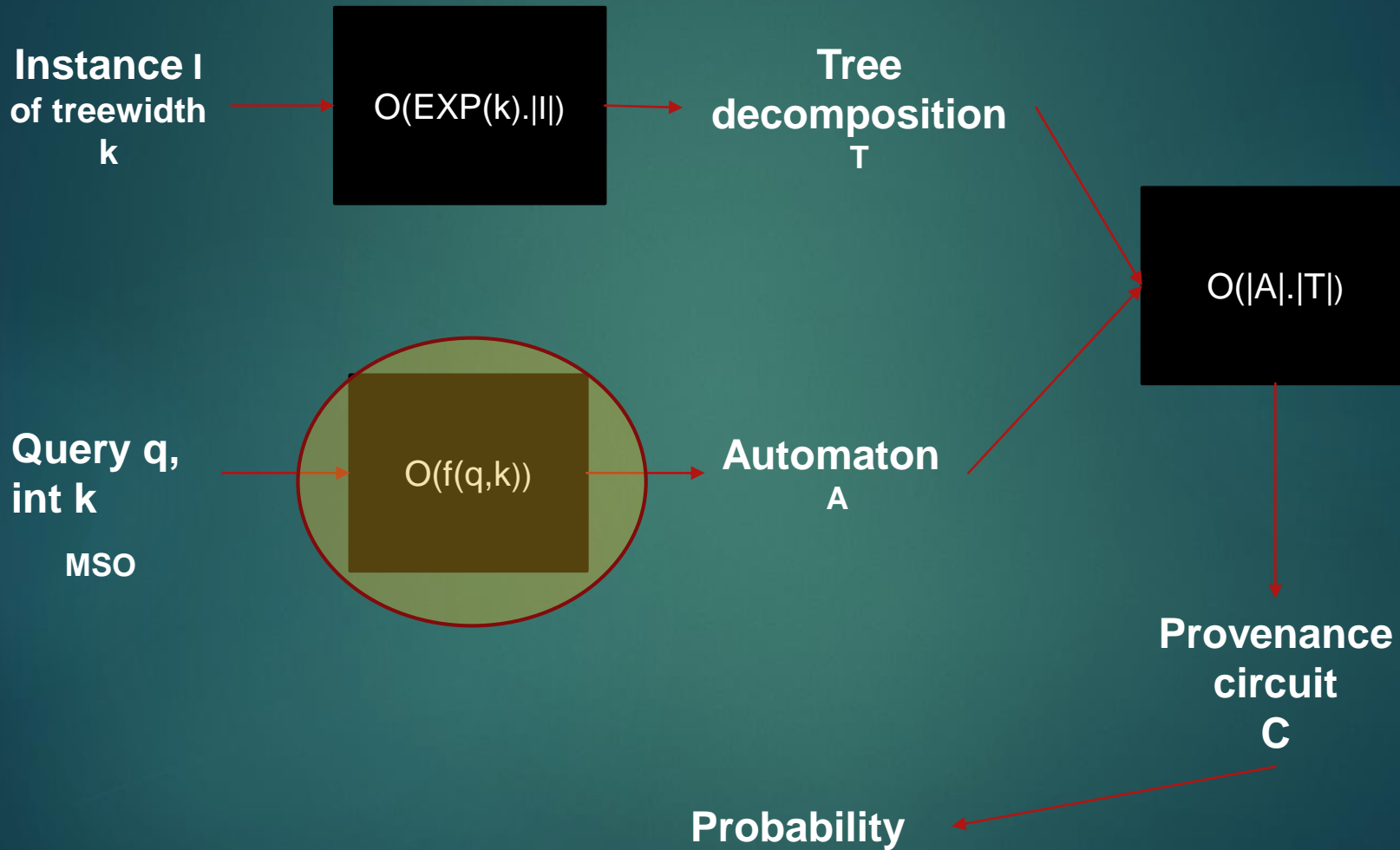
10



Amarilli, Antoine, Maniu Silviu, Monnet Mikael

Problems

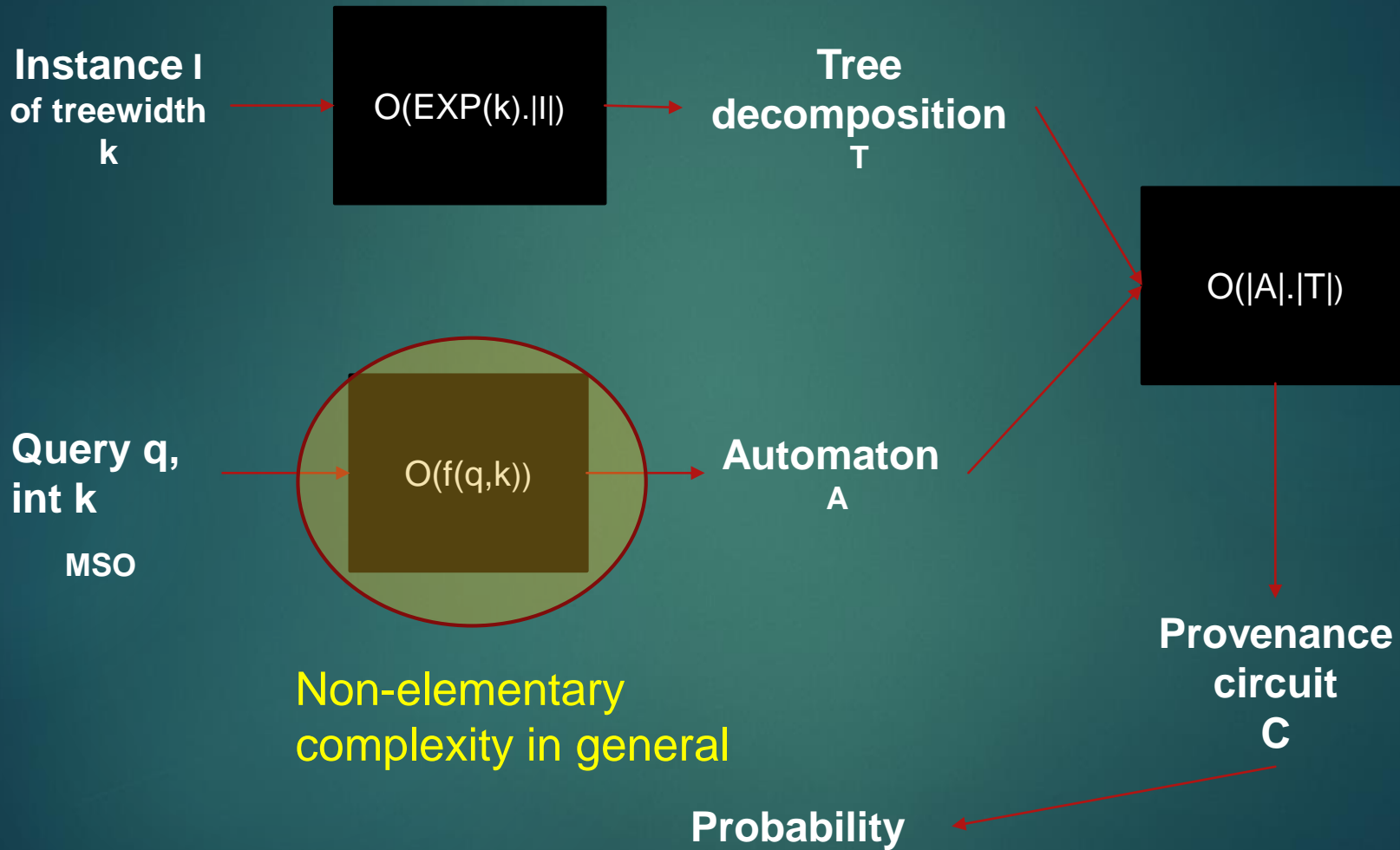
10



Amarilli, Antoine, Maniu Silviu, Monnet Mikael

Problems

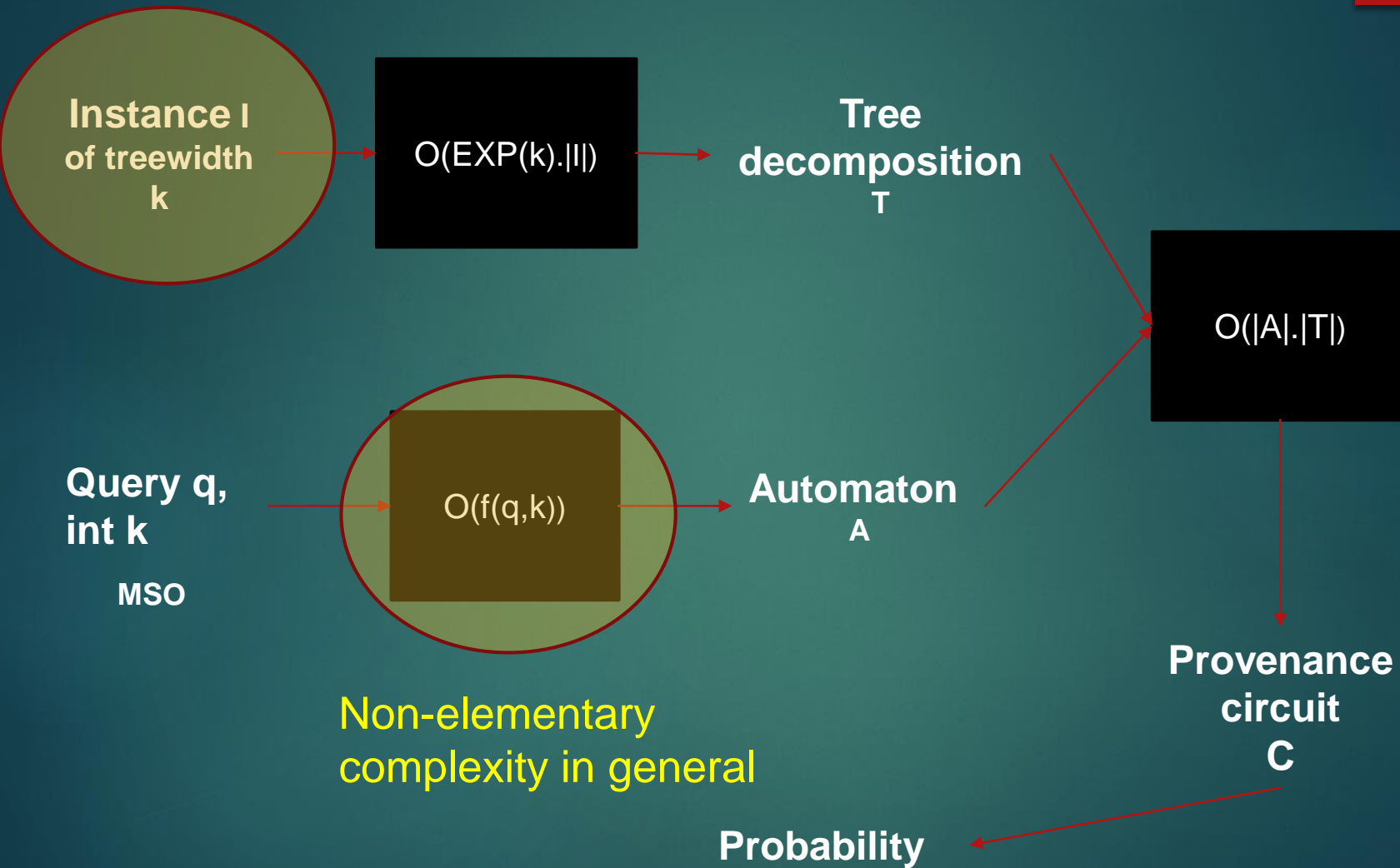
10



Amarilli, Antoine, Maniu Silviu, Monnet Mikael

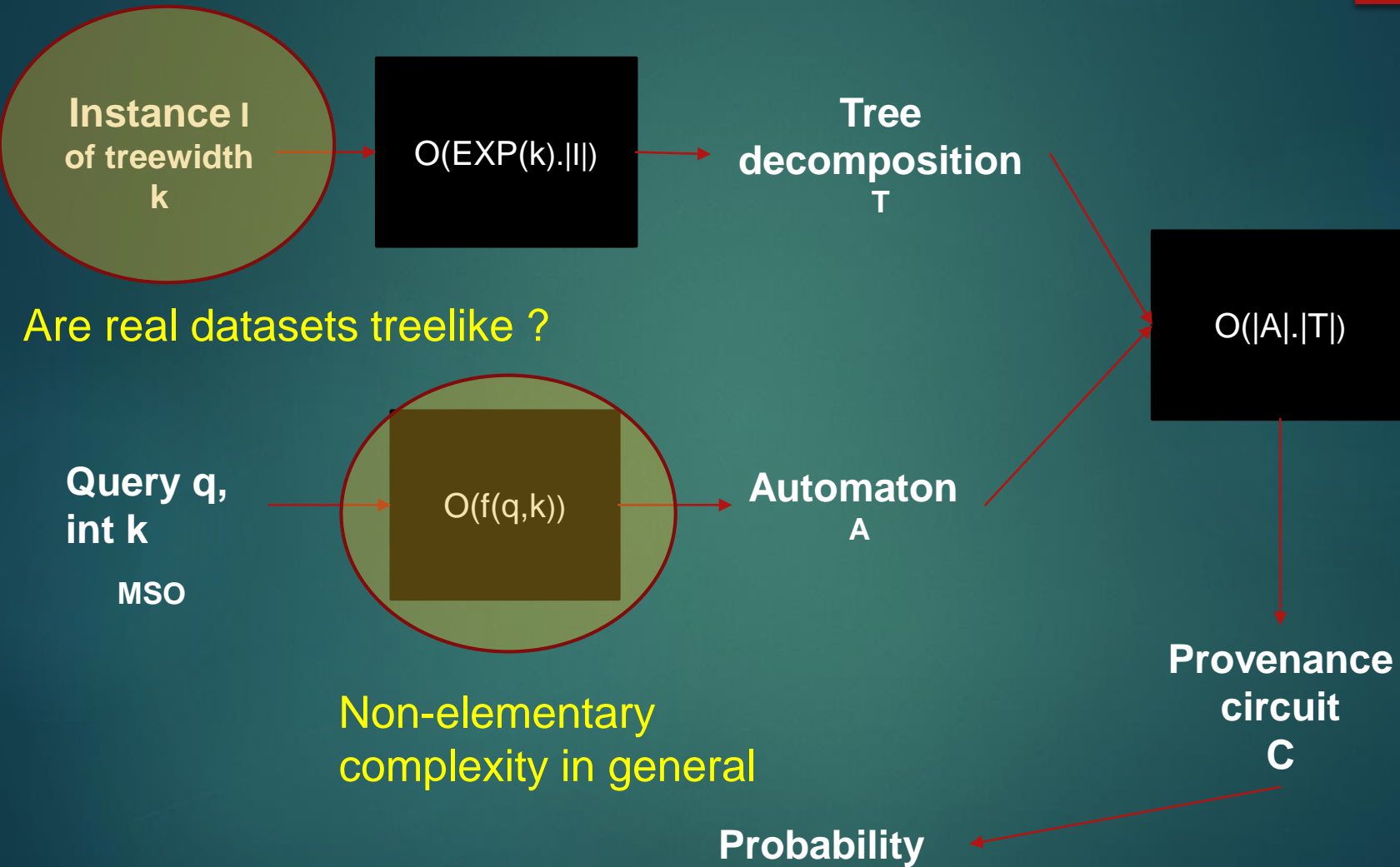
Problems

10



Amarilli, Antoine, Maniu Silviu, Monnet Mikael

Problems



Amarilli Antoine, Maniu Silviu, Monnet Mikael

Current work

11

Amarilli Antoine, Maniu Silviu, Monet Mikael

Current work

11

- ▶ From bottom-up tree automata to alternating two-way automata

Current work

11

- ▶ From bottom-up tree automata to alternating two-way automata
- ▶ Introduce Intensionally-Clique-Guarded Datalog (ICG-Datalog) parameterized by body-size

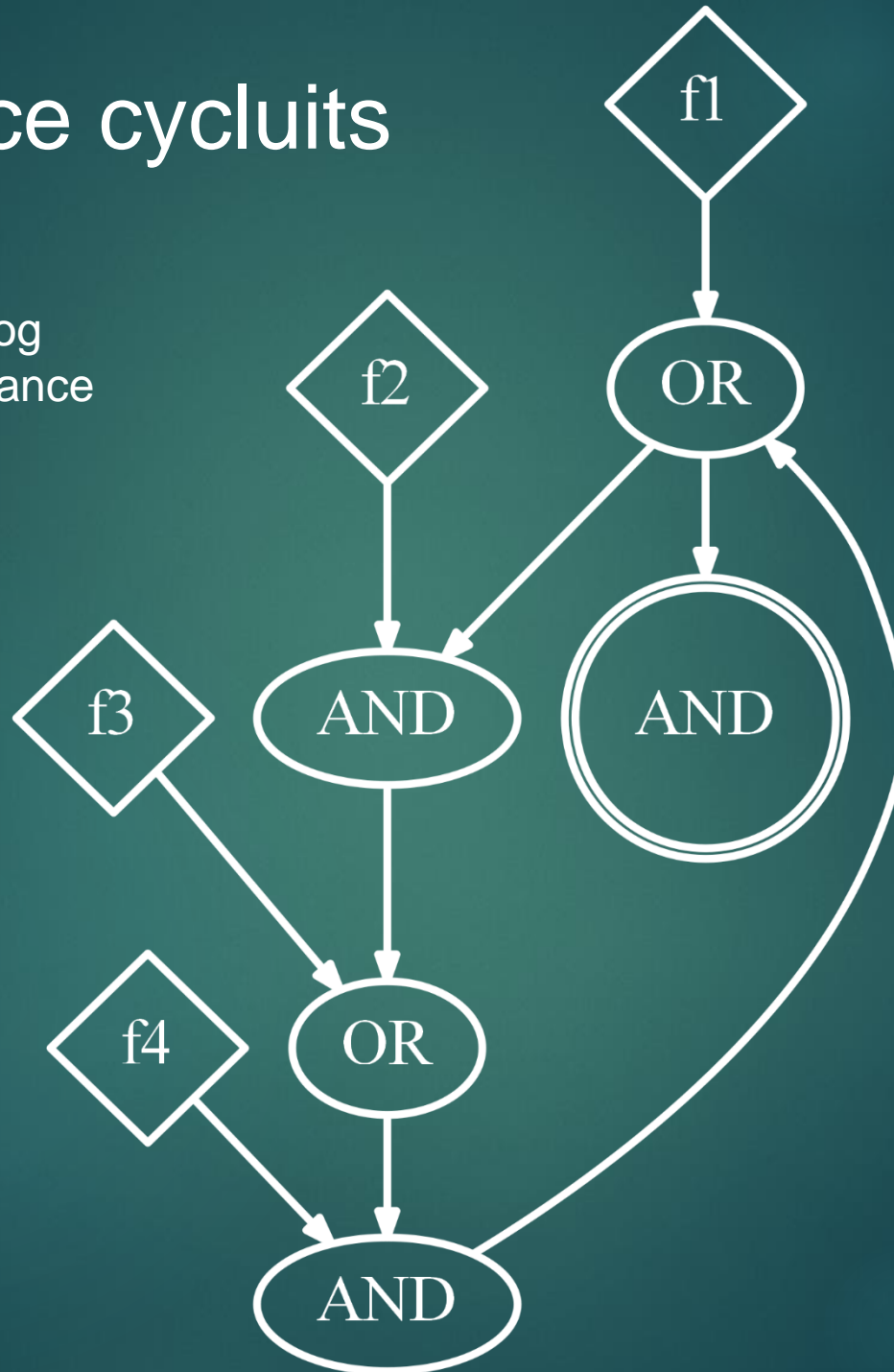
Current work

11

- ▶ From bottom-up tree automata to alternating two-way automata
- ▶ Introduce Intensionally-Clique-Guarded Datalog (ICG-Datalog) parameterized by body-size
- ▶ Provenance as a cyclic circuit ! (cycluit)

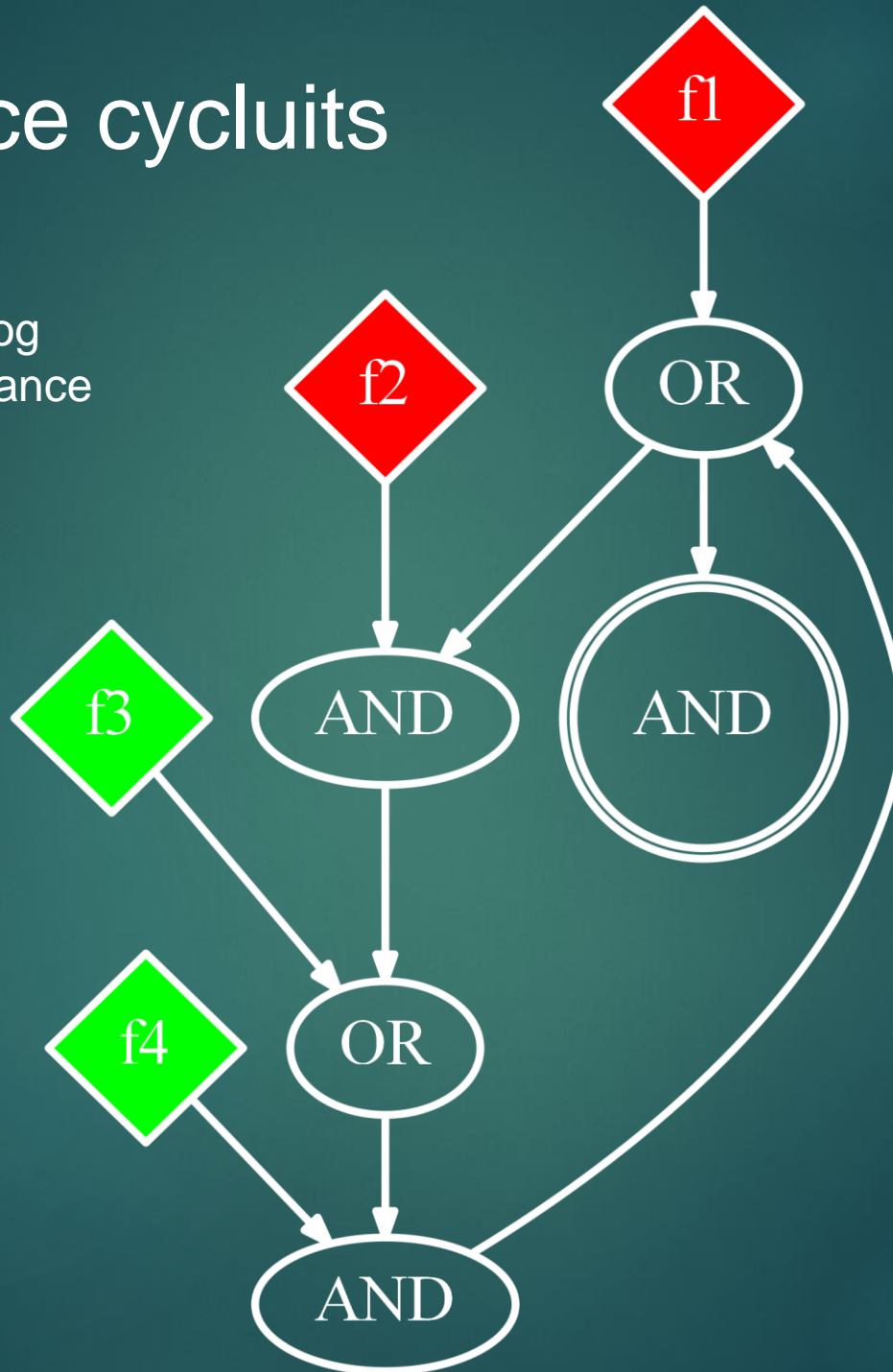
Provenance cycluits

Some ICG-Datalog program, some instance I with facts {f1, f2, f3, f4}



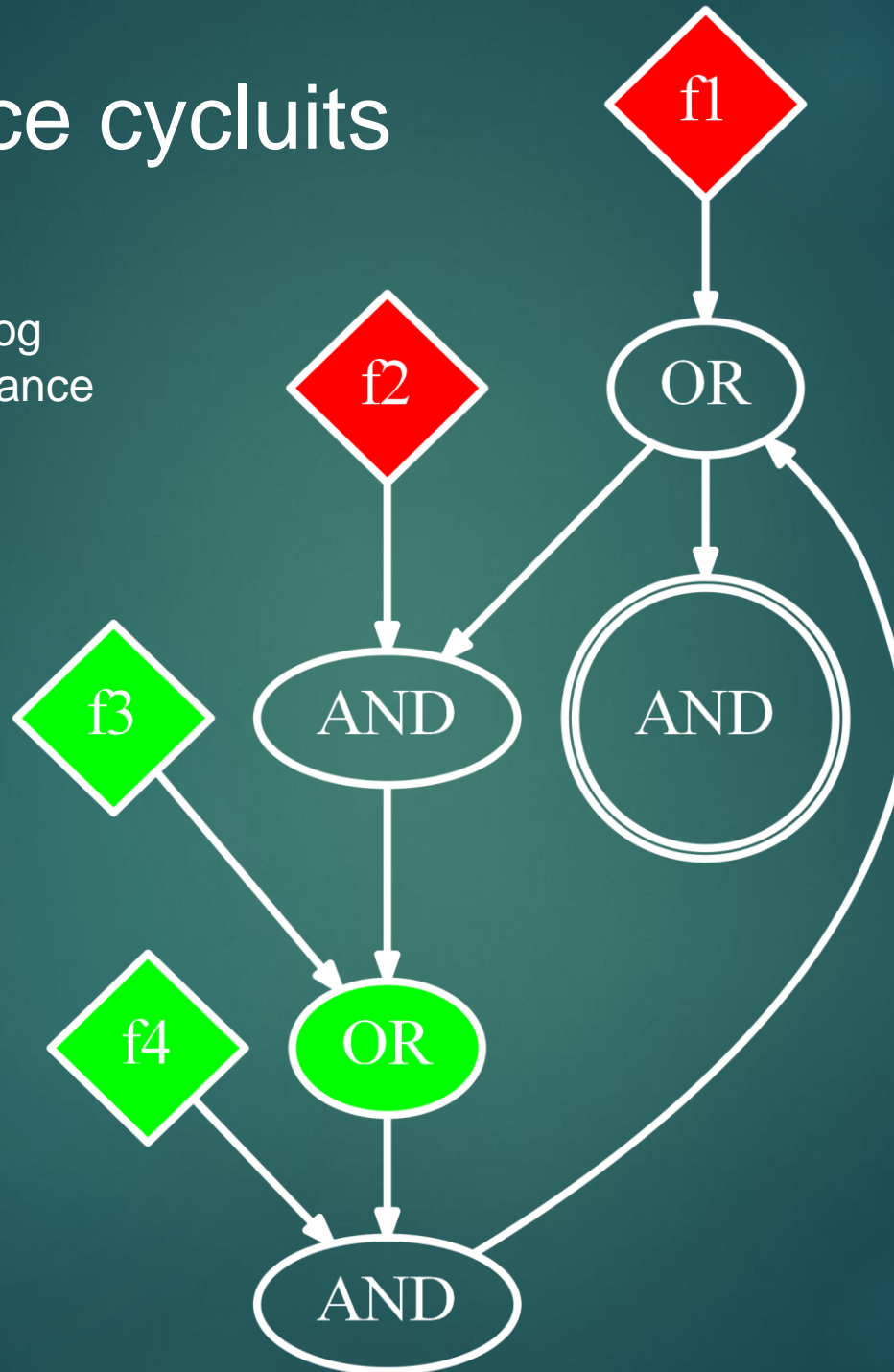
Provenance cycluits

Some ICG-Datalog program, some instance I with facts {f1, f2, f3, f4}



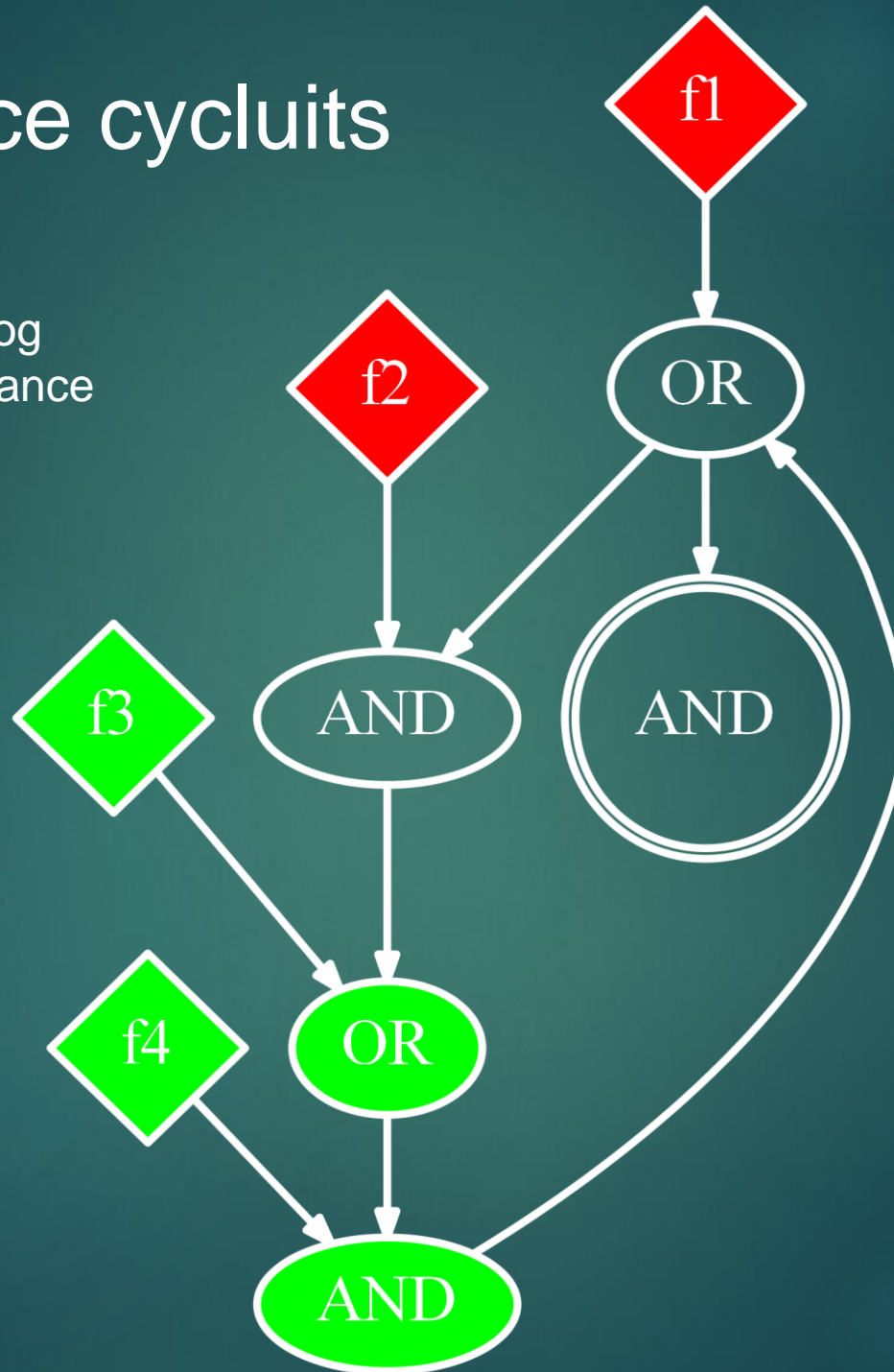
Provenance cycluits

Some ICG-Datalog program, some instance I with facts {f1, f2, f3, f4}



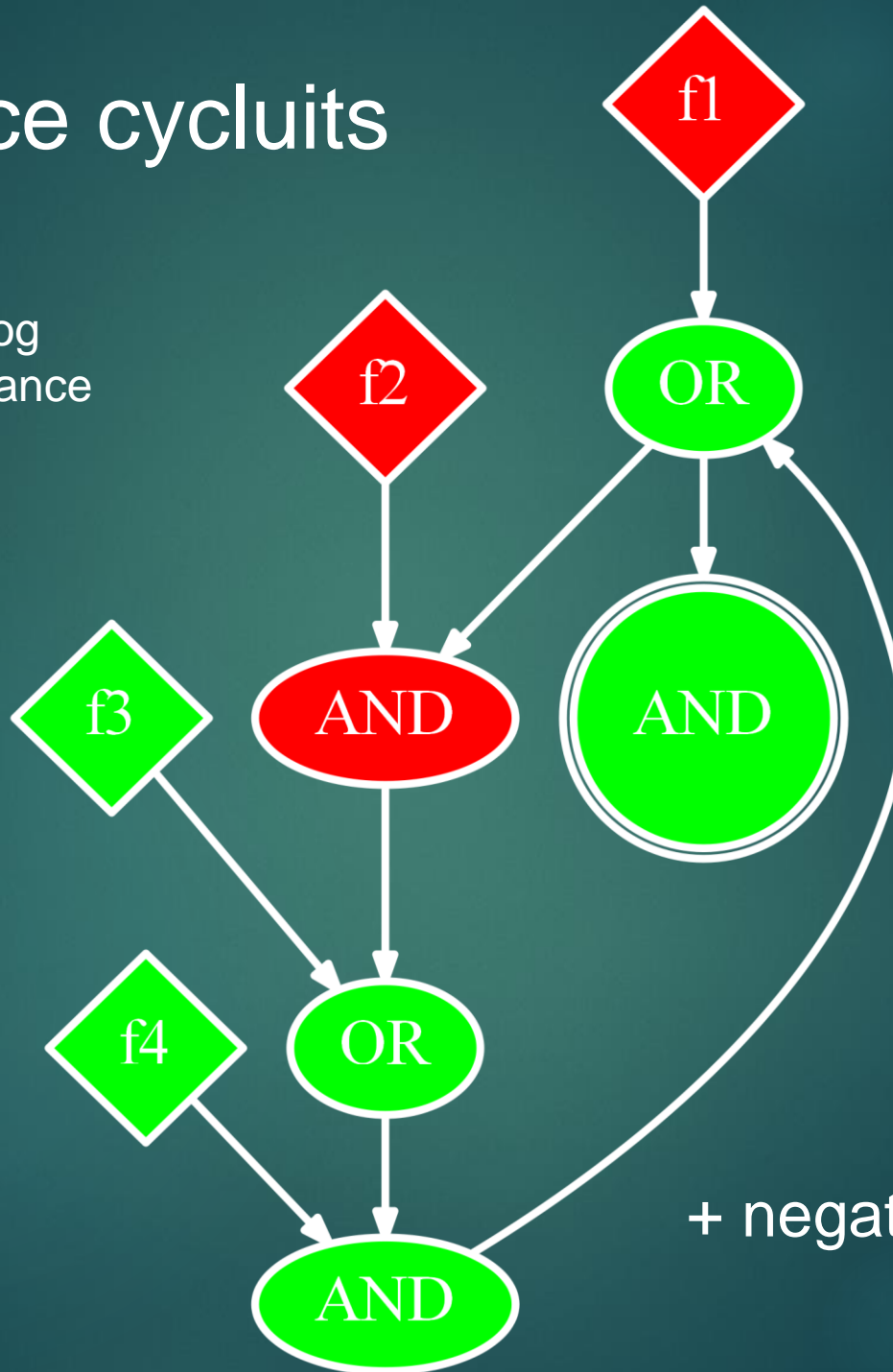
Provenance cycluits

Some ICG-Datalog program, some instance I with facts {f1, f2, f3, f4}



Provenance cycluits

Some ICG-Datalog program, some instance I with facts {f1, f2, f3, f4}



+ negations (stratified)

Instance I
of treewidth
 k

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

Tree
decomposition
 T

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

ICG
program P
of body
size c , int k

$O(f(c, k) \cdot |P|)$

2-way
Automaton
 A

Provenance
CYCLUIT
 C

Probability

Instance I
of treewidth
 k

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

Tree
decomposition
 T

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

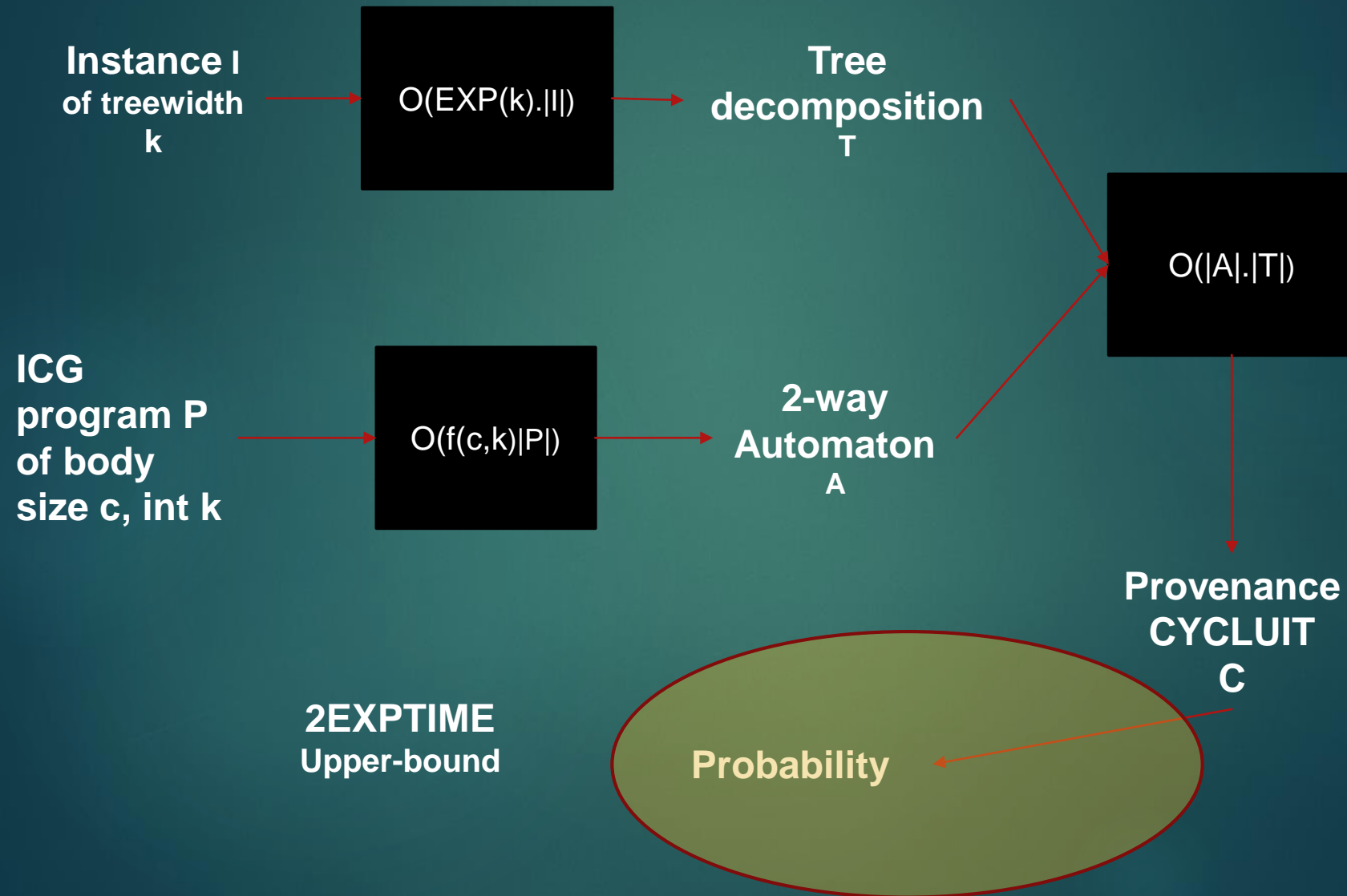
ICG
program P
of body
size c , int k

$O(f(c, k) \cdot |P|)$

2-way
Automaton
 A

Provenance
CYCLUIT
 C

Probability



Bad news...

- ▶ We proved that:

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

- ▶ Still, we obtain a $2EXPTIME$ combined complexity upperbound

Treelike datasets

15

Amarilli, Antoine, Maniu Silviu, Monet Mikael

Treelike datasets

15

- ▶ Transportation networks

Treelike datasets

15

- ▶ Transportation networks
- ▶ Partial decompositions

Treelike datasets

15

- ▶ Transportation networks
- ▶ Partial decompositions
- ▶ Query-specific decompositions

Instance I
of treewidth
 k

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

Tree
decomposition
 T

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

ICG
program P
of body
size c , int k

$O(f(c, k) \cdot |P|)$

2-way
Automaton
 A

Provenance
CYCLUIT
 C