

Leveraging the structure of uncertain data

Antoine Amarilli May 16, 2018











Database theory and query evaluation



- (Hyper)graph
- Collection of ground facts

G(aa₁, ab₂), G(ab₂, ac₃), S(aa₁, m₄), S(ab₂, r_B), ...

Database theory and query evaluation



• (Hyper)graph

Collection of ground facts

G(aa₁, ab₂), G(ab₂, ac₃), S(aa₁, m₄), S(ab₂, r_B), ...

uery

Rechercher mon itinéraire			
De	Rue Monticelli	0	.e
Α	Gare du Nord	J	
Aujo	urd'hui		•
Arriv	rée à 🗸 🗸	21	h ¥

- Regular path (Metro|RER)* |(Bus|Tram)*
- Logic formula

 $\begin{array}{l} \forall X(rm \in X \land \forall xy \\ (x \in X \land G(x, y) \rightarrow \\ y \in X)) \rightarrow gn \in X \end{array}$

Database theory and query evaluation



• (Hyper)graph

Collection of ground facts

G(aa₁, ab₂), G(ab₂, ac₃), S(aa₁, m₄), S(ab₂, r_B), ...



Rechercher mon itineraire			
De	Rue Monticelli	0	
			ປ*
A	Gare du Nord		
Aujou	urd'hui		•
Arriv	ée à 🗸 🗸	21	h •



Result

- Regular path (Metro|RER)* |(Bus|Tram)*
- Logic formula

 $\begin{array}{l} \forall X(rm \in X \land \forall xy \\ (x \in X \land G(x, y) \rightarrow \\ y \in X)) \rightarrow gn \in X \end{array}$

• TRUE/FALSE ↔ Model checking



Q, :	Départ 20h17 - 5 rue Monticelli , Paris	
*** 齐 **	1.1 km 13 min	
B	20h30 - CITE UNIVERSITAIRE, Paris	
	RER B - EPAU Vers Aéroport CDG Terminal 2 TGV	0 0 0
	✓ 6 arrêts 14 min	
•	Arrivée 20h44 - GARE DU NORD RER , Paris	

Panne du RER B : trafic interrompu entre Paris et Roissy, des TGV en renfort

🏠 > Transports | 06 décembre 2016, 9h56 | MAJ: 06 décembre 2016, 17h03 | 🛉 🈏 🗨





Panne du RER B : trafic interrompu entre Paris : pourquoi il y a autant de perturbations sur INCIDENT SUR LE RER B : QUE S'EST-IL PASSÉ CE MATIN ? Malaise voyageur et application des mesures de sécurité : pour quelles raisons le trafic a-t-il été perturbé ce matin sur la ligne B? Pour beaucoup, le voyage a été difficile ce matin. Au fil de vos réactions sur Twitter notamment, je constate que les raisons de ces perturbations ne paraissent pas cohérentes. Je tiens donc à vous apporter des premiers éléments d'explication, que nous pourrons développer

Panne du RER B : trafic interrompu entre Paris : pourquoi il y a autant de perturbations sur le INCIDENT SUB I E REB B : OUE

ACTUALITÉS

Le RER B en panne, les voyageurs n'ont pas eu d'autre choix que de descendre sur les voies

Alors que la circulation alternée a augmenté le nombre de voyageurs dans les transports en commun, le RER B s'est retrouvé à l'arrêt.

O 06/12/2016 11:57 CET | Actualisé 06/12/2016 20:14 CET



Pour beaucoup, le voyage a été difficile ce matin. Au fil de vos réactions sur Twitter notamment, je constate que les raisons de ces perturbations ne paraissent pas cohérentes. Je tiens donc à vous apporter des premiers éléments d'explication, que pous pourrons dévelopmer











- (Hyper)graph
- Collection of ground facts
 + independent probabilities



- (Hyper)graph
- Collection of ground facts + independent probabilities

?	Query
---	-------

Rechercher mon itinéraire			
De	Rue Monticelli	0	۰ŀ
A	Gare du Nord		
Aujo	urd'hui		•
Arriv	réeà v	21	h •

- Regular path (Metro|RER)* |(Bus|Tram)*
- Logic formula

 $\begin{array}{l} \forall X(rm \in X \land \forall xy \\ (x \in X \land G(x, y) \rightarrow \\ y \in X)) \rightarrow gn \in X \end{array}$



- (Hyper)graph
- Collection of ground facts
 + independent
 probabilities

?	Query
لنا	Query

Rechercher mon itinéraire			
De	Rue Monticelli	0	
			ଏ
Α	Gare du Nord		
Aujo	urd'hui		~
Arriv	rée à 🗸 🗸	21	h ¥

- Regular path (Metro|RER)* |(Bus|Tram)*
- Logic formula

 $\begin{array}{l} \forall X(rm \in X \land \forall xy \\ (x \in X \land G(x, y) \rightarrow \\ y \in X)) \rightarrow gn \in X \end{array}$



 Probability according to the input distribution



- Computing paths on a large graph:
 - ightarrow Well-studied problem, efficient algorithms



Computing paths on a large probabilistic graph:
→ ???



- Computing paths on a large **probabilistic** graph:
 - \rightarrow **Exponential** number of possibilities



- Computing paths on a large **probabilistic** graph:
 - \rightarrow Exponential number of possibilities
 - \rightarrow **#P-hard** computational complexity in the **database**















 \rightarrow Shortest path: very easy on a large tree

In this talk:

• Existing results on **non-probabilistic data**:

In this talk:

- Existing results on non-probabilistic data:
 - Tree automata, to evaluate queries on trees

In this talk:

- Existing results on **non-probabilistic data**:
 - Tree automata, to evaluate queries on trees
 - Treewidth, formalizes the notion of being "close to a tree"
- Existing results on **non-probabilistic data**:
 - Tree automata, to evaluate queries on trees
 - Treewidth, formalizes the notion of being "close to a tree"
 - · Courcelle's theorem

- Existing results on non-probabilistic data:
 - Tree automata, to evaluate queries on trees
 - Treewidth, formalizes the notion of being "close to a tree"
 - Courcelle's theorem
- Introduce **new tools** and **results**:

- Existing results on non-probabilistic data:
 - Tree automata, to evaluate queries on trees
 - Treewidth, formalizes the notion of being "close to a tree"
 - Courcelle's theorem
- Introduce **new tools** and **results**:
 - Provenance circuits of tree automata on uncertain trees

- Existing results on non-probabilistic data:
 - Tree automata, to evaluate queries on trees
 - Treewidth, formalizes the notion of being "close to a tree"
 - · Courcelle's theorem
- Introduce **new tools** and **results**:
 - Provenance circuits of tree automata on uncertain trees
 - Applications to **probabilistic query evaluation**

- Existing results on non-probabilistic data:
 - Tree automata, to evaluate queries on trees
 - Treewidth, formalizes the notion of being "close to a tree"
 - · Courcelle's theorem
- Introduce **new tools** and **results**:
 - Provenance circuits of tree automata on uncertain trees
 - Applications to **probabilistic query evaluation**
- Other applications: Counting, enumeration, provenance...

Introduction

Existing tools

Provenance circuits and probabilistic query evaluation

Other applications

Database: a **word w** where nodes have a color from an alphabet $\bigcirc \bigcirc \bigcirc$



Database: a word w where nodes have a color from an alphabet OOO



Query Q: a sentence (yes/no question) in monadic second-order logic (MSO) "Is there both a pink and a blue node?" Database: a word w where nodes have a color from an alphabet OOO



Query Q: a sentence (yes/no question) in monadic second-order logic (MSO) "Is there both a pink and a blue node?"

 \mathbf{i} **Result**: TRUE/FALSE indicating if the word **w** satisfies the query **Q**

Database: a word w where nodes have a color from an alphabet OOO



Query Q: a sentence (yes/no question) in monadic second-order logic (MSO) "Is there both a pink and a blue node?"

 \mathbf{i} **Result**: TRUE/FALSE indicating if the word **w** satisfies the query **Q**

Computational complexity as a function of **w** (the query **Q** is **fixed**)



- $P_{\odot}(x)$ means "x is blue"; also $P_{\odot}(x)$, $P_{\bigcirc}(x)$
- $x \rightarrow y$ means "x is the predecessor of y"



- $x \rightarrow y$ means "x is the predecessor of y"
- Propositional logic: formulas with AND $\wedge,$ OR $\vee,$ NOT \neg

• $P_{\bigcirc}(x) \land P_{\bigcirc}(y)$ means "Node x is pink and node y is blue"



- $x \rightarrow y$ means "x is the predecessor of y"
- Propositional logic: formulas with AND $\wedge,$ OR $\vee,$ NOT \neg
 - $P_{\bigcirc}(x) \land P_{\bigcirc}(y)$ means "Node x is pink and node y is blue"
- First-order logic: adds existential quantifier ∃ and universal quantifier ∀
 - · ∃x y $P_{\bigcirc}(x) \land P_{\bigcirc}(y)$ means "There is both a pink and a blue node"



- $x \rightarrow y$ means "x is the predecessor of y"
- Propositional logic: formulas with AND $\wedge,$ OR $\vee,$ NOT \neg
 - $P_{\bigcirc}(x) \land P_{\bigcirc}(y)$ means "Node x is pink and node y is blue"
- First-order logic: adds existential quantifier ∃ and universal quantifier ∀
 - $\exists x \ y \ P_{\bigcirc}(x) \land P_{\bigcirc}(y)$ means "There is both a pink and a blue node"
- Monadic second-order logic (MSO): adds quantifiers over sets
 - $\exists S \forall x S(x)$ means "there is a set S containing every element x"
 - Can express transitive closure $x \rightarrow^* y$, i.e., "x is before y"
 - $\forall x P_{\bigcirc}(x) \Rightarrow \exists y P_{\bigcirc}(y) \land x \rightarrow^{*} y$ means "There is a blue node after every pink node"

• States: $\{\bot, B, P, \top\}$

- States: $\{\bot, B, P, \top\}$
- Final states: $\{\top\}$

- **States:** {⊥, *B*, *P*, ⊤}
- Final states: $\{\top\}$
- Initial function: $\bigcirc \bot \quad \bigcirc P \quad \bigcirc B$

- **States:** {⊥, *B*, *P*, ⊤}
- Final states: $\{\top\}$
- Initial function: $\bigcirc \bot \quad \bigcirc P \quad \bigcirc B$

- **States:** {⊥, *B*, *P*, ⊤}
- Final states: $\{\top\}$
- Initial function: $\bigcirc \bot \quad \bigcirc P \quad \bigcirc B$
- Transitions (examples): $\perp \bigcirc_{P} P \bigcirc_{T} \top \bigcirc_{T}$

Translate the query **Q** to a **deterministic word automaton**

- **States:** {⊥, *B*, *P*, ⊤}
- Final states: $\{\top\}$
- Initial function: $\bigcirc \bot \quad \bigcirc P \quad \bigcirc B$
- Transitions (examples): $\perp \bigcirc_{P} P \bigcirc_{\top} \top \bigcirc_{\top}$

Translate the query **Q** to a **deterministic word automaton**

Alphabet: $\bigcirc \bigcirc \bigcirc \qquad w: \bigcirc - \bigcirc - \bigcirc - \bigcirc \qquad Q: \exists x \ y \ P_{\bigcirc}(x) \land P_{\bigcirc}(y)$

- **States:** {⊥, *B*, *P*, ⊤}
- Final states: $\{\top\}$
- Initial function: $\bigcirc \bot \quad \bigcirc P \quad \bigcirc B$
- Transitions (examples): $\perp \bigcirc_{P} P \bigcirc_{\top} \top \bigcirc_{\top}$

Translate the query **Q** to a **deterministic word automaton**

- States: $\{\perp, B, P, \top\}$
- Final states: $\{\top\}$
- Initial function: $\bigcirc \bot \quad \bigcirc P \quad \bigcirc B$
- Transitions (examples): $\perp \bigcirc_{P} P \bigcirc_{\top} \top \bigcirc_{\top}$

Translate the query **Q** to a **deterministic word automaton**

- **States:** {⊥, *B*, *P*, ⊤}
- Final states: $\{\top\}$
- Initial function: $\bigcirc \bot \quad \bigcirc P \quad \bigcirc B$
- Transitions (examples): $\perp \bigcirc_{P} P \bigcirc_{\top} \top \bigcirc_{\top}$

Translate the query **Q** to a **deterministic word automaton**

- **States:** {⊥, *B*, *P*, ⊤}
- Final states: $\{\top\}$
- Initial function: $\bigcirc \bot \quad \bigcirc P \quad \bigcirc B$
- Transitions (examples): $\perp \bigcirc_{\mathbf{p}} \mathbf{P} \bigcirc_{\mathbf{T}} \top \bigcirc_{\mathbf{T}}$

Theorem (Büchi, 1960)

MSO and word automata have the same expressive power on words

- States: $\{\perp, B, P, \top\}$
- Final states: $\{\top\}$
- Initial function: $\bigcirc \bot \quad \bigcirc P \quad \bigcirc B$
- Transitions (examples): $\perp \bigcirc_{P} P \bigcirc_{T} \top \bigcirc_{T}$

Theorem (Büchi, 1960)

MSO and word automata have the same expressive power on words

Corollary

Query evaluation of MSO on words is in **linear time**.







Database: a **tree** T where nodes have a color from an alphabet $\bigcirc \bigcirc \bigcirc$





- $\cdot P_{\odot}(x)$ means "x is blue"
- $\cdot x
 ightarrow y$ means "x is the parent of y"

"Is there both a pink and a blue node?" ∃x y P_⊙(x) ∧ P_⊙(y)



Database: a **tree** T where nodes have a color from an alphabet $\bigcirc \bigcirc \bigcirc$





Query Q: a **sentence** in monadic second-order logic (MSO)

- $\cdot P_{\odot}(x)$ means "x is blue"
- $\cdot x
 ightarrow y$ means "x is the parent of y"

"Is there both a pink and a blue node?" ∃x y P_⊙(x) ∧ P_⊙(y)

 \bigcirc **Result**: TRUE/FALSE indicating if the tree *T* satisfies the query *Q*



Database: a **tree** T where nodes have a color from an alphabet $\bigcirc \bigcirc \bigcirc$





Query Q: a sentence in monadic second-order logic (MSO)
 Po(x) means "x is blue"

 $\cdot x \rightarrow y$ means "x is the parent of y"

"Is there both a pink and a blue node?" ∃x y P_⊙(x) ∧ P_⊙(y)

 \bigcirc **Result**: TRUE/FALSE indicating if the tree *T* satisfies the query *Q*

Computational complexity as a function of **T** (the query **Q** is **fixed**)





- Bottom-up deterministic tree automaton
- "Is there both a pink and a blue node?"

Tree alphabet:

- Bottom-up deterministic tree automaton
- "Is there both a pink and a blue node?"
- States: $\{\bot, B, P, \top\}$



- Bottom-up deterministic tree automaton
- "Is there both a pink and a blue node?"
- States: $\{\bot, B, P, \top\}$
- Final states: $\{\top\}$



- Bottom-up deterministic tree automaton
- "Is there both a pink and a blue node?"
- States: $\{\bot, B, P, \top\}$
- Final states: $\{\top\}$
- Initial function: $\bigcirc \bot \quad \bigcirc P \quad \bigcirc B$



- Bottom-up deterministic tree automaton
- "Is there both a pink and a blue node?"
- States: $\{\bot, B, P, \top\}$
- Final states: $\{\top\}$
- Initial function: $\bigcirc \bot \quad \bigcirc P \quad \bigcirc B$


- Bottom-up deterministic tree automaton
- "Is there both a pink and a blue node?"
- States: $\{\bot, B, P, \top\}$
- Final states: $\{\top\}$

В

- Initial function: $\bigcirc \bot \quad \bigcirc P \quad \bigcirc B$
- Transitions (examples):



• "Is there both a pink and a blue node?"

• Initial function: $\bigcirc \bot \bigcirc P$ В



В



Theorem [Thatcher and Wright, 1968]

MSO and tree automata have the same expressive power on trees



Theorem [Thatcher and Wright, 1968]

MSO and tree automata have the same expressive power on trees

Corollary

Query evaluation of MSO on trees is in linear time.



Database: a **tree** T where nodes have a color from an alphabet $\bigcirc \bigcirc \bigcirc$





Query Q: a sentence in monadic second-order logic (MSO)
 Po(x) means "x is blue"

 $\cdot x
ightarrow y$ means "x is the parent of y"

"Is there both a pink and a blue node?" ∃x y P_⊙(x) ∧ P_⊙(y)

 \bigcirc **Result**: TRUE/FALSE indicating if *T* satisfies the query *Q*

Computational complexity as a function of the **tree** *T* (the query *Q* is **fixed**)







(Metro|RER)* | (Bus|Tram)*

Result: TRUE/FALSE indicating if **T** satisfies the query **Q**

Computational complexity as a function of the **tree** *T* (the query *Q* is **fixed**)

















































- Trees have treewidth 1
- Cycles have treewidth 2
- k-cliques and (k 1)-grids have treewidth k 1





- Trees have treewidth 1
- Cycles have treewidth 2
- k-cliques and (k 1)-grids have treewidth k 1
- $\rightarrow~\mbox{Treelike}:$ the $\mbox{treewidth}$ is bounded by a $\mbox{constant}$

Treelike data



MSO query (RER|metro)*

|(bus|tram)*

Treelike data











Theorem [Courcelle, 1990]

For any fixed Boolean MSO query **Q** and $k \in \mathbb{N}$, given a database **D** of treewidth $\leq k$, we can compute in **linear time** in **D** whether **D** satisfies **Q** Introduction

Existing tools

Provenance circuits and probabilistic query evaluation

Other applications

- Database D with treewidth ≤ k for some constant k
- **Probability** of each fact of **D** to be actually present in the data (independently from other facts)



- Database D with treewidth ≤ k for some constant k
- Probability of each fact of D to be actually present in the data (independently from other facts)





(Metro|RER)* | (Bus|Tram)*

- Database D with treewidth ≤ k for some constant k
- Probability of each fact of D to be actually present in the data (independently from other facts)





(Metro|RER)* | (Bus|Tram)*

(i) **Result**: **Probability** that the database **D** satisfies query **Q**

- Database D with treewidth ≤ k for some constant k
- Probability of each fact of D to be actually present in the data (independently from other facts)





(Metro|RER)* | (Bus|Tram)*

1 Result: **Probability** that the database **D** satisfies query **Q**

Computational complexity as a function of the **database** *D* (the query **Q** is **fixed**)

Roadmap



Roadmap



Roadmap








Valuation: $\{2, 3, 7 \mapsto 1, * \mapsto 0\}$



Valuation: $\{2 \mapsto 1, * \mapsto 0\}$



Valuation: $\{2, 7 \mapsto 1, * \mapsto 0\}$



Valuation: $\{2, 7 \mapsto 1, * \mapsto 0\}$

A: "Is there both a pink and a blue node?"



Valuation: $\{2, 3, 7 \mapsto 1, * \mapsto 0\}$

A: "Is there both a pink and a blue node?"

The tree automaton A accepts



Valuation: $\{2 \mapsto 1, * \mapsto 0\}$

A: "Is there both a pink and a blue node?"

The tree automaton A rejects



Valuation: $\{2, 7 \mapsto 1, * \mapsto 0\}$

A: "Is there both a pink and a blue node?"

The tree automaton A accepts



• Directed acyclic graph of gates



- Directed acyclic graph of gates
- Output gate:





- Directed acyclic graph of gates
- Output gate:
- Variable gates:





• Directed acyclic graph of gates

Х

Λ

- Output gate:
- Variable gates:
- Internal gates:



- Directed acyclic graph of gates
- Output gate:
- Variable gates: (
- Internal gates: (V) (A) (¬)
- Valuation: function from variables to $\{0, 1\}$ Example: $\nu = \{x \mapsto 0, y \mapsto 1\}$...



- Directed acyclic graph of gates
- Output gate:
- Variable gates: (
- Internal gates: (\vee) (\wedge) (\neg)
- Valuation: function from variables to $\{0, 1\}$ Example: $\nu = \{x \mapsto 0, y \mapsto 1\}$...



- Directed acyclic graph of gates
- Output gate:
- Variable gates: (
- Internal gates: (\vee) (\wedge) (\neg)
- Valuation: function from variables to $\{0, 1\}$ Example: $\nu = \{x \mapsto 0, y \mapsto 1\}$...



- Directed acyclic graph of gates
- Output gate:
- Variable gates: (
- Internal gates: (V) (^) (¬)
- Valuation: function from variables to $\{0, 1\}$ Example: $\nu = \{x \mapsto 0, y \mapsto 1\}$... mapped to 1

Provenance circuit



Query: Is there both a pink and a blue node?

Provenance circuit



Provenance circuit



Formally:

- Tree automaton A, uncertain tree T, circuit C
- Variable gates of C: nodes of T
- Condition: Let ν be a valuation of T, then $\nu(C)$ iff A accepts $\nu(T)$

- Alphabet: 🔿 🔵 🔵
- Automaton: "Is there both a pink and a blue node?"
- States:
- $\{\perp, B, P, \top\}$
- Final: {⊤}

- Transitions:
- $\mathbf{P}^{\top} \mathbf{P}^{\top}$

- Alphabet: 🔿 🔵 🔵
- Automaton: "Is there both a pink and a blue node?" • Final: $\{\top\}$ P \perp
- States:
 - $\{\perp, B, P, \top\}$

- Transitions:



- Alphabet: 🔿 🔵 🔵
- Automaton: "Is there both a pink and a blue node?"
- States:
 - $\{\perp, \textit{B},\textit{P}, \top\}$
- Final: $\{\top\}$ P \perp
- Transitions:
- $P^{\top} P^{\top}$



- Alphabet: 🔿 🔵 🔵
- Automaton: "Is there both a pink and a blue node?"
- States:
 - $\{\perp, \textit{B},\textit{P}, \top\}$
- Final: $\{\top\}$

- Transitions:
- $\mathbb{R}^{\top} \mathbb{R}^{P}$



- Alphabet: 🔿 🔵 🔵
- Automaton: "Is there both a pink and a blue node?"
- States:
 - $\{\bot, B, P, \top\}$
- Final: $\{\top\}$

- Transitions:
- $\begin{array}{c} \mathbb{R}^{\top} & \mathbb{R}^{P} \\ \mathbb{R}^{\perp} & \mathbb{R}^{P} \end{array}$



- Alphabet: 🔿 🔵 🔵
- Automaton: "Is there both a pink and a blue node?"
- States:
 - $\{\bot, B, P, \top\}$
- Final: $\{\top\}$

- Transitions:



- Alphabet: 🔿 🔵 🔵
- Automaton: "Is there both a pink and a blue node?"
- States:
 - $\{\perp, \textit{B},\textit{P},\top\}$
- Final: $\{\top\}$

- Transitions:



Probabilistic query evaluation



Probabilistic treelike **data**



Each **fact** can **disappear** with some probability

Probabilistic treelike **data**







Each **fact** can **disappear** with some probability Each **node label** can **disappear** with the probability of the coded **fact**



Each fact can disappear with some probability

Each node label can **disappear** with the probability of the coded **fact**

Each variable can **be true** with the probability of the coded **fact**

٨

́у



Each **fact** can **disappear** with some probability Each **node label** can **disappear** with the probability of the coded **fact** Each **variable** can **be true** with the probability of the coded **fact** Probability that the **circuit** evaluates to **true**



Each **fact** can **disappear** with some probability Each **node label** can **disappear** with the probability of the coded **fact** Each **variable** can **be true** with the probability of the coded **fact** Probability that the **circuit** evaluates to **true**

 \rightarrow How to compute **efficiently** the probability of the circuit?

Computing the probability of a circuit

• We are given a circuit and a probability P for each variable



Computing the probability of a circuit

• We are given a **circuit** and a **probability P** for each variable



- P(x) = 40%
 P(y) = 50%
- We are given a circuit and a probability P for each variable
- Each variable x is true **independently** with probability P(x)



- P(x) = 40%
- P(y) = 50%

- We are given a circuit and a probability P for each variable
- Each variable x is true **independently** with probability P(x)
- What is the probability that the circuit evaluates to true?



- P(x) = 40%
- P(y) = 50%

- We are given a circuit and a probability P for each variable
- Each variable x is true **independently** with probability P(x)
- What is the probability that the circuit evaluates to true?
- In general, **#P-hard** (harder than SAT)

- P(x) = 40%
- P(y) = 50%

- We are given a circuit and a probability P for each variable
- Each variable x is true **independently** with probability P(x)
- What is the probability that the circuit evaluates to true?



- In general, **#P-hard** (harder than SAT)
- Here it's easy:

- P(x) = 40%
- P(y) = 50%

- We are given a circuit and a probability P for each variable
- Each variable x is true **independently** with probability P(x)
- What is the probability that the circuit evaluates to true?



- In general, **#P-hard** (harder than SAT)
- Here it's easy:
 - $\cdot\,$ The inputs to the $\wedge\text{-gate}$ are independent

- P(x) = 40%
- P(y) = 50%

- We are given a circuit and a probability P for each variable
- Each variable x is true **independently** with probability P(x)
- What is the probability that the circuit evaluates to true?

- In general, **#P-hard** (harder than SAT)
- Here it's easy:
 - $\cdot\,$ The inputs to the $\wedge\text{-gate}$ are independent

• P(x) = 40%

Х

20%

- We are given a circuit and a probability P for each variable
- Each variable x is true **independently** with probability P(x)
- What is the probability that the circuit evaluates to true?

- In general, **#P-hard** (harder than SAT)
- Here it's easy:
 - $\cdot\,$ The inputs to the $\wedge\text{-gate}$ are independent
 - The ¬-gate has probability 1 − P(input)

• P(x) = 40%

Х

20%

- We are given a circuit and a probability P for each variable
- Each variable x is true **independently** with probability P(x)
- What is the probability that the circuit evaluates to true?

- In general, **#P-hard** (harder than SAT)
- Here it's easy:
 - The inputs to the $\wedge\text{-gate}$ are independent
 - The \neg -gate has probability 1 P(input)

• P(x) = 40%

60%

Х

20%

v

- We are given a circuit and a probability P for each variable
- Each variable x is true **independently** with probability P(x)
- What is the probability that the circuit evaluates to true?

- In general, **#P-hard** (harder than SAT)
- Here it's easy:
 - $\cdot\,$ The inputs to the $\wedge\text{-gate}$ are independent
 - The \neg -gate has probability 1 P(input)
 - The V-gate has mutually exclusive inputs

• P(x) = 40%

60%

Х

20%

v

- We are given a circuit and a probability P for each variable
- Each variable x is true **independently** with probability P(x)
- What is the probability that the circuit evaluates to true?

- In general, **#P-hard** (harder than SAT)
- Here it's easy:
 - $\cdot\,$ The inputs to the $\wedge\text{-gate}$ are independent
 - The ¬-gate has probability 1 − P(input)
 - The V-gate has mutually exclusive inputs

• P(x) = 40%

60%

Х

80%

20%

У

- We are given a circuit and a probability P for each variable
- Each variable x is true **independently** with probability P(x)
- What is the probability that the circuit evaluates to true?

- In general, **#P-hard** (harder than SAT)
- Here it's easy:
 - $\cdot\,$ The inputs to the $\wedge\text{-gate}$ are independent
 - The \neg -gate has probability 1 P(input)
 - The V-gate has mutually exclusive inputs
- Let's focus on a **restricted class** of circuits that satisfies these conditions

• P(x) = 40%

60%

Х

80%

The circuit is a **d-DNNF**...

The circuit is a **d-DNNF**...



The circuit is a **d-DNNF**...

• V gates always have mutually exclusive inputs

The circuit is a **d-DNNF**...

• V gates always have mutually exclusive inputs



The circuit is a **d-DNNF**...

... so probability computation is **easy**!

• V gates always have mutually exclusive inputs

A gates are all on independent inputs

The circuit is a **d-DNNF**...

... so probability computation is **easy**!





• V gates always have mutually exclusive inputs



The circuit is a **d-DNNF**...

... so probability computation is **easy**!





$$P(g) := 1 - P(g')$$

• V gates always have mutually exclusive inputs



The circuit is a **d-DNNF**...

... so probability computation is **easy**!



$$P(g) := 1 - P(g')$$

• V gates always have mutually exclusive inputs

A gates are all on independent inputs

The circuit is a **d-DNNF**...

... so probability computation is **easy**!



• V gates always have mutually exclusive inputs



The circuit is a **d-DNNF**...

... so probability computation is **easy**!

• V gates always have mutually exclusive inputs

independent inputs

gates are all on

g P(q) := 1 - P(q')g $P(g) := P(g_1') + P(q_2')$ g'_1 g_2' g q_2' g'_1

The circuit is a **d-DNNF**...

... so probability computation is **easy**!

• V gates always have mutually exclusive inputs

• (A) gates are all on independent inputs

g P(q) := 1 - P(q')g $P(g) := P(g_1') + P(q_2')$ g'_1 g_2' g $P(g) := P(g_1') \times P(q_2')$ q'_1 q_2'

The circuit is a **d-DNNF**...

... so probability computation is **easy**!

• V gates always have **mutually exclusive** inputs

(A) gates are all on independent inputs

g P(q) := 1 - P(q')g $P(g) := P(g_1') + P(q_2')$ g'_1 g_2' g $P(g) := P(g_1') \times P(q_2')$ q' q_2'

Lemma

The **provenance circuit** computed in our construction is a **d-DNNF**









Theorem [Amarilli et al., 2015]

For any **fixed** Boolean MSO query **Q** and $k \in \mathbb{N}$, given a database **D** of **treewidth** $\leq k$ with **independent probabilities**, we can compute in **linear time** the probability that **D** satisfies **Q** 28/33 Introduction

Existing tools

Provenance circuits and probabilistic query evaluation

Other applications

"Is there both a pink and a blue node?"

 $Q(): \exists x \ y \ P_{\odot}(x) \land P_{\odot}(y)$

"Is there both a pink and a blue node?"

 $Q(): \exists x \ y \ P_{\odot}(x) \land P_{\odot}(y)$

• In practice, queries often return some results:

"Find all pairs of a pink and a blue node?"

 $Q(x,y): P_{\odot}(x) \wedge P_{\odot}(y)$

"Is there both a pink and a blue node?"

 $Q(): \exists x \ y \ P_{\odot}(x) \land P_{\odot}(y)$

• In practice, queries often return some results:

"Find all pairs of a pink and a blue node?"

 $Q(x,y): P_{\odot}(x) \wedge P_{\odot}(y)$

• We can consider each pair (a, b) and test if Q(a, b) is true

"Is there both a pink and a blue node?"

 $Q(): \exists x \ y \ P_{\odot}(x) \land P_{\odot}(y)$

• In practice, queries often return some results:

"Find all pairs of a pink and a blue node?"

 $Q(x,y): P_{\odot}(x) \wedge P_{\odot}(y)$

- We can consider each pair (a, b) and test if Q(a, b) is true
- Can we do **better**?

• Query: $Q(X_1, \ldots, X_n)$ with free variables X_1, \ldots, X_n

- Query: $Q(X_1, \ldots, X_n)$ with free variables X_1, \ldots, X_n
- Goal: find all tuples a_1, \ldots, a_n such that $Q(a_1, \ldots, a_n)$ holds

- Query: $Q(X_1, \ldots, X_n)$ with free variables X_1, \ldots, X_n
- Goal: find all tuples a_1, \ldots, a_n such that $Q(a_1, \ldots, a_n)$ holds
- \rightarrow Add **special facts** to materialize all possible assignments
 - e.g., $X_i(a_j)$ means element a_i is mapped to variable X_j

- Query: $Q(X_1, \ldots, X_n)$ with free variables X_1, \ldots, X_n
- Goal: find all tuples a_1, \ldots, a_n such that $Q(a_1, \ldots, a_n)$ holds
- \rightarrow Add **special facts** to materialize all possible assignments
 - e.g., $X_i(a_j)$ means element a_i is mapped to variable X_j
- \rightarrow The provenance circuit of Q is now a factorized representation which describes all the tuples that make Q true
- Query: $Q(X_1, \ldots, X_n)$ with free variables X_1, \ldots, X_n
- Goal: find all tuples a_1, \ldots, a_n such that $Q(a_1, \ldots, a_n)$ holds
- \rightarrow Add **special facts** to materialize all possible assignments
 - · e.g., $X_i(a_j)$ means element a_i is mapped to variable X_j
- \rightarrow The provenance circuit of Q is now a factorized representation which describes all the tuples that make Q true

Example query:

 $Q(X_1,X_2):P_{\bigcirc}(x)\wedge P_{\bigcirc}(y)$

- Query: $Q(X_1, \ldots, X_n)$ with free variables X_1, \ldots, X_n
- Goal: find all tuples a_1, \ldots, a_n such that $Q(a_1, \ldots, a_n)$ holds
- ightarrow Add **special facts** to materialize all possible assignments
 - · e.g., $X_i(a_j)$ means element a_i is mapped to variable X_j
- \rightarrow The provenance circuit of Q is now a factorized representation which describes all the tuples that make Q true

Example query:

 $Q(X_1, X_2) : P_{\bigcirc}(x) \land P_{\bigcirc}(y)$ Database: (1)-(2)-(3)-(4)-(5)

- Query: $Q(X_1, \ldots, X_n)$ with free variables X_1, \ldots, X_n
- Goal: find all tuples a_1, \ldots, a_n such that $Q(a_1, \ldots, a_n)$ holds
- ightarrow Add **special facts** to materialize all possible assignments
 - · e.g., $X_i(a_j)$ means element a_i is mapped to variable X_j
- \rightarrow The provenance circuit of Q is now a factorized representation which describes all the tuples that make Q true

Example query:

 $Q(X_1, X_2) : P_{\odot}(x) \land P_{\odot}(y)$ Database: 1-2-3-4-5

Results:

 $\begin{array}{ccc} X_1 & X_2 \\ 1 & 3 \\ 1 & 5 \end{array}$

- Query: $Q(X_1, \ldots, X_n)$ with free variables X_1, \ldots, X_n
- Goal: find all tuples a_1, \ldots, a_n such that $Q(a_1, \ldots, a_n)$ holds
- ightarrow Add **special facts** to materialize all possible assignments
 - · e.g., $X_i(a_j)$ means element a_i is mapped to variable X_j
- \rightarrow The provenance circuit of Q is now a factorized representation which describes all the tuples that make Q true

Example query:

 $Q(X_{1}, X_{2}) : P_{\odot}(x) \land P_{\odot}(y)$ Database:
1-2-3-4-5 <u>Results:</u> $X_{1} X_{2}$ 1 3
1 5

Provenance circuit:



- Query: $Q(X_1, \ldots, X_n)$ with free variables X_1, \ldots, X_n
- Goal: find all tuples a_1, \ldots, a_n such that $Q(a_1, \ldots, a_n)$ holds
- ightarrow Add **special facts** to materialize all possible assignments
 - · e.g., $X_i(a_j)$ means element a_i is mapped to variable X_j
- \rightarrow The provenance circuit of Q is now a factorized representation which describes all the tuples that make Q true

Example query:

$Q(X_{1}, X_{2}) : P_{\odot}(x) \land P_{\odot}(y)$ Database: 1-2-3-4-5 <u>Results:</u> $X_{1} X_{2}$ 1 3 1 5

Provenance circuit:



- Query: $Q(X_1, \ldots, X_n)$ with free variables X_1, \ldots, X_n
- Goal: find all tuples a_1, \ldots, a_n such that $Q(a_1, \ldots, a_n)$ holds
- ightarrow Add **special facts** to materialize all possible assignments
 - e.g., $X_i(a_j)$ means element a_i is mapped to variable X_j
- \rightarrow The provenance circuit of Q is now a factorized representation which describes all the tuples that make Q true

Example query: $Q(X_1, X_2) : P_{\odot}(x) \land P_{\odot}(y)$ Database: 1-2-3-4-5 <u>Results:</u> $X_1 X_2$ 1 3 1 5

Provenance circuit:



This factorized representation of the results of the query can be computed in **linear time** in the data

• Application: Counting query results [Arnborg et al., 1991]

- Application: Counting query results [Arnborg et al., 1991]
 - Exclusive \lor means +, independent \land means \times
 - Reproves existing result: [Arnborg et al., 1991]

- Application: Counting query results [Arnborg et al., 1991]
 - Exclusive \lor means +, independent \land means \times
 - Reproves existing result: [Arnborg et al., 1991]
- Application: Constant-delay enumeration of query results

- Application: Counting query results [Arnborg et al., 1991]
 - + Exclusive \lor means +, independent \land means \times
 - Reproves existing result: [Arnborg et al., 1991]
- Application: Constant-delay enumeration of query results
 - Requires some linear-time preprocessing of the input circuit
 - + Exclusive \lor means <code>disjoint</code> \cup , <code>independent</code> \land means <code>relational</code> \times
 - New **modular proof** of existing enumeration result [Bagan, 2006, Kazana and Segoufin, 2013, Amarilli et al., 2017a]
 - Extensions to support **updates** on the database [Amarilli et al., 2018]

- Application: Counting query results [Arnborg et al., 1991]
 - + Exclusive \lor means +, independent \land means \times
 - Reproves existing result: [Arnborg et al., 1991]
- Application: Constant-delay enumeration of query results
 - Requires some linear-time preprocessing of the input circuit
 - + Exclusive \lor means <code>disjoint</code> \cup , <code>independent</code> \land means <code>relational</code> \times
 - New **modular proof** of existing enumeration result [Bagan, 2006, Kazana and Segoufin, 2013, Amarilli et al., 2017a]
 - Extensions to support **updates** on the database [Amarilli et al., 2018]
- Application: Semiring provenance [Green et al., 2007]

Conclusion and perspectives

- Other results:
 - Lower bounds: probabilistic query evaluation is hard unless treewidth is bounded (modulo assumptions) [Amarilli et al., 2016]
 - **Complexity in the query:** generally nonelementary but can be improved [Amarilli et al., 2017b, Amarilli et al., 2017c]

Conclusion and perspectives

- Other results:
 - Lower bounds: probabilistic query evaluation is hard unless treewidth is bounded (modulo assumptions) [Amarilli et al., 2016]
 - **Complexity in the query:** generally nonelementary but can be improved [Amarilli et al., 2017b, Amarilli et al., 2017c]
- Ongoing work (with my wonderful co-authors):
 - More efficient enumeration algorithms on words
 - More lower bounds results, connections to knowledge compilation
 - More expressive provenance: cycluits (circuits with cycles)
 - Combined tractability for probabilistic query evaluation



Pierre



Louis



Stefan



Mikaël



Matthias



Pierre

Conclusion and perspectives

- Other results:
 - Lower bounds: probabilistic query evaluation is hard unless treewidth is bounded (modulo assumptions) [Amarilli et al., 2016]
 - **Complexity in the query:** generally nonelementary but can be improved [Amarilli et al., 2017b, Amarilli et al., 2017c]
- Ongoing work (with my wonderful co-authors):
 - More efficient enumeration algorithms on words
 - More lower bounds results, connections to knowledge compilation
 - More expressive provenance: cycluits (circuits with cycles)
 - Combined tractability for probabilistic query evaluation



Pierre



Louis







Pierre

Stefan Mikaël Matthias Thanks for your attention!

References i

Amarilli, A., Bourhis, P., Jachiet, L., and Mengel, S. (2017a). **A Circuit-Based Approach to Efficient Enumeration.** In *ICALP*.

- Amarilli, A., Bourhis, P., and Mengel, S. (2018). Enumeration on Trees under Relabelings. In ICDT.
- Amarilli, A., Bourhis, P., Monet, M., and Senellart, P. (2017b).
 Combined Tractability of Query Evaluation via Tree Automata and Cycluits.

In ICDT.

Amarilli, A., Bourhis, P., and Senellart, P. (2015).
 Provenance Circuits for Trees and Treelike Instances.
 In ICALP.



Amarilli, A., Bourhis, P., and Senellart, P. (2016).

Tractable Lineages on Treelike Instances: Limits and Extensions. In *PODS*.

 Amarilli, A., Monet, M., and Senellart, P. (2017c).
 Conjunctive Queries on Probabilistic Graphs: Combined Complexity.

In PODS.

Arnborg, S., Lagergren, J., and Seese, D. (1991). **Easy problems for tree-decomposable graphs.** J. Algorithms, 12(2):308–340.

References iii



Bagan, G. (2006).

MSO queries on tree decomposable structures are computable with linear delay.

In CSL.

Courcelle, B. (1990).

The monadic second-order logic of graphs. I. Recognizable sets of finite graphs.

Inf. Comput., 85(1).

Green, T. J., Karvounarakis, G., and Tannen, V. (2007). **Provenance semirings.**

In PODS.

Kazana, W. and Segoufin, L. (2013). Enumeration of monadic second-order queries on trees. TOCL, 14(4).

Thatcher, J. W. and Wright, J. B. (1968).
 Generalized finite automata theory with an application to a decision problem of second-order logic.

Mathematical systems theory, 2(1):57–81.

Image credits

- Slides 2 and 5-6:
 - Subway map: https://commons.wikimedia.org/wiki/File:Paris_Metro_map.svg (edited), by user Umx on Wikimedia Commons, public domain
 - Ticket t+: http://www.parisvoyage.com/images/cartoon18.jpg, ParisVoyage, fair use
 - Terms and conditions: http://www.vianavigo.com/fileadmin/galerie/pdf/CGU_t_.pdf (cropped), RATP, fair use
- Slides 3-4: screenshots from http://lab.vianavigo.com, Stif, fair use
- Slide 4: newpaper articles (fair use) :
 - http://www.leparisien.fr/transports/

http:

//www.rtl.fr/actu/societe-faits-divers/paris-le-trafic-totalement-interrompu-gare-du-nord-7786171150

- https://www.rerb-leblog.fr/incident-rer-b-sest-passe-matin/
- http://www.huffingtonpost.fr/2016/12/06/le-rer-b-en-panne-les-voyageurs-nont-pas-eu-dautres-choix-que/
- http://www.lexpress.fr/actualite/societe/trafic/
 rer-b-en-panne-retards-du-d-circulation-alternee-deuxieme-journee-de-galere_1857905.html
- http://www.lemonde.fr/entreprises/article/2016/12/07/ ile-de-france-le-trafic-toujours-interrompu-sur-le-rer-b-en-direction-de-roissy_5044717_1656994.html
- Slides 6, 16, 19, 24–25, 28: Train map https://commons.wikimedia.org/wiki/File:Carte_TGV.svg?uselang=fr (edited), by
 users Jack ma, Muselaar, Benjism89, Pic-Sou, Uwe Dedering, Madcap on Wikimedia Commons, license CC-BY-SA 3.0
- Slide 33: Photos http://www.lifl.fr/~bourhis/pb.png, http://tyrex.inria.fr/people/img/jachiet.png, http://www.cril.univ-artois.fr/~mengel/snap.jpeg, http://mikael-monet.net/images/moi.jpg, https://sigmodrecord.org/wp-content/uploads/2017/05/Matthias-Niewerth-matthias.niewerth.jpg, http://pierre.senellart.com/bubu.jpg, fair use