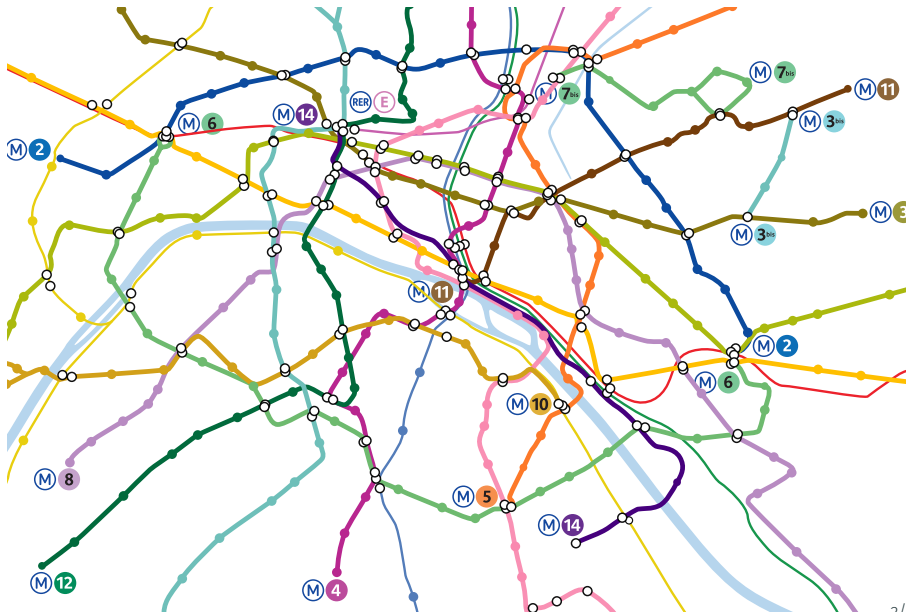


Leveraging the structure of uncertain data

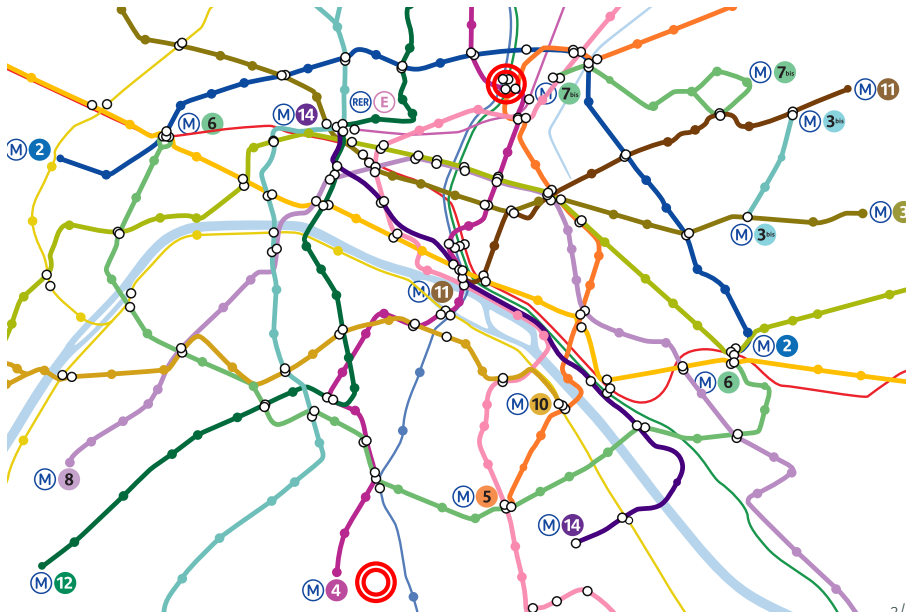
Antoine Amarilli

March 5, 2019

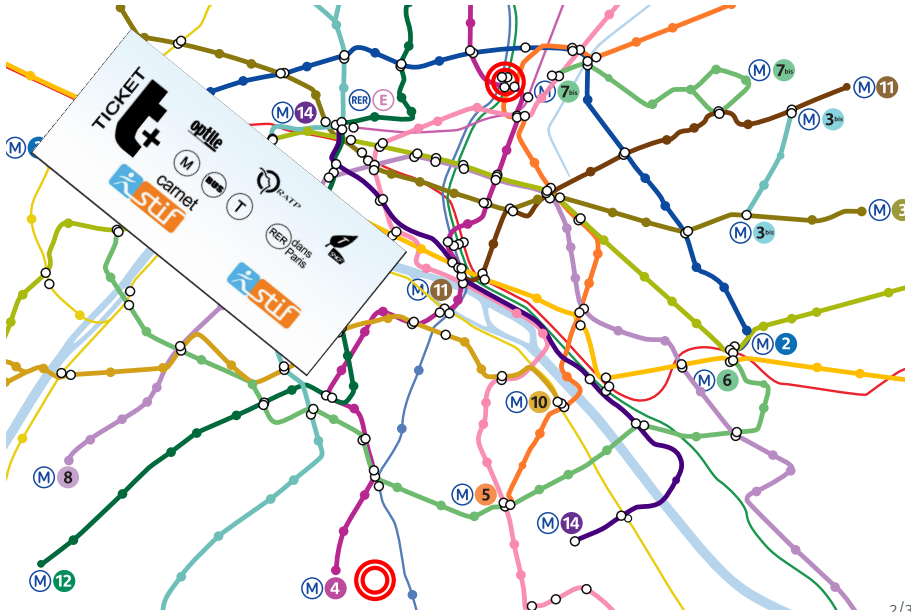
Example application: Subway routing



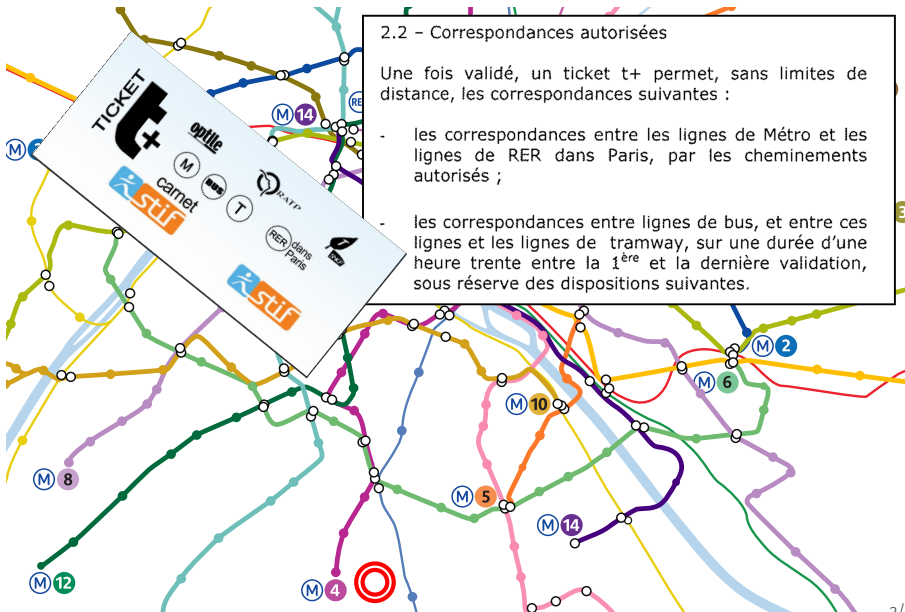
Example application: Subway routing



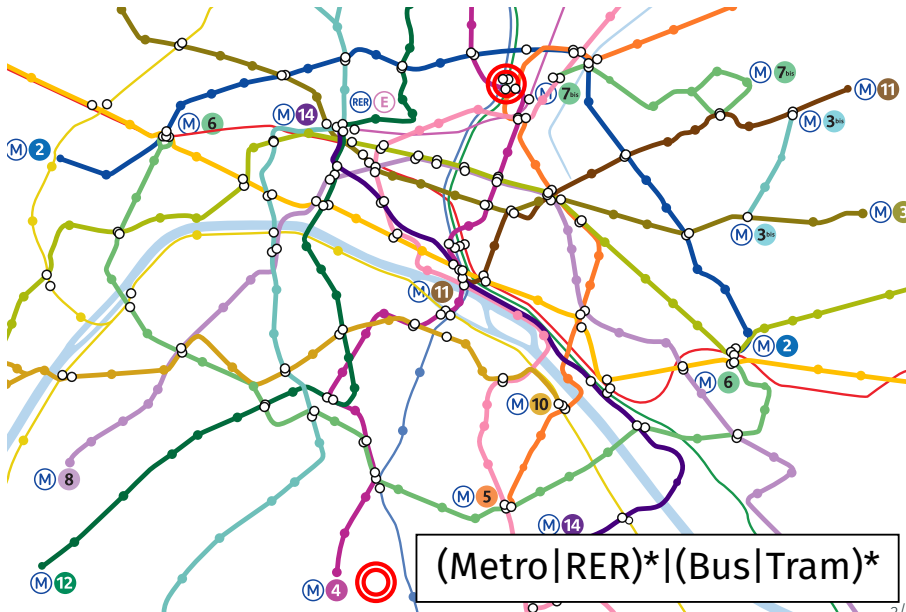
Example application: Subway routing



Example application: Subway routing



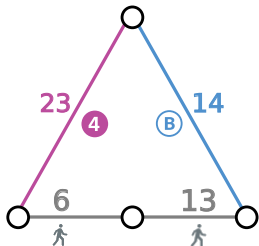
Example application: Subway routing



Database theory and query evaluation



Database



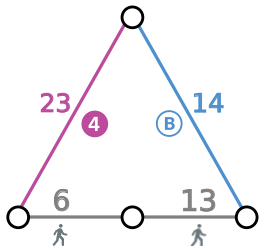
- (Hyper)graph
- Collection of ground facts

$G(aa_1, ab_2), G(ab_2, ac_3),$
 $S(aa_1, m_4), S(ab_2, r_B), \dots$

Database theory and query evaluation



Database



Query

Rechercher mon itinéraire

De

A

Aujourd'hui

Arrivée à 21 h

- (Hyper)graph

- Collection of ground facts

$G(aa_1, ab_2), G(ab_2, ac_3),$
 $S(aa_1, m_4), S(ab_2, r_B), \dots$

- Regular path

$(\text{Metro}|\text{RER})^*$
 $|(\text{Bus}|\text{Tram})^*$

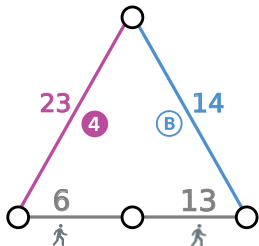
- Logic formula

$\forall X (rm \in X \wedge \forall xy$
 $(x \in X \wedge G(x, y) \rightarrow$
 $y \in X)) \rightarrow gn \in X$

Database theory and query evaluation



Database



Query

Rechercher mon itinéraire

De

A

Arrivée à



Result

Départ
20h17 - 5 rue Monticelli, Paris

1,1 km | 13 min

20h30 - CITE UNIVERSITAIRE, Paris

RER B - EPAU
Vers Aéroport CDG Terminal 2 TGV

6 arrêts | 14 min

Arrivée
20h44 - GARE DU NORD RER, Paris

- (Hyper)graph

- Collection of ground facts

$G(aa_1, ab_2), G(ab_2, ac_3),$
 $S(aa_1, m_4), S(ab_2, r_B), \dots$

- Regular path

$(\text{Metro}|\text{RER})^*$
 $|(\text{Bus}|\text{Tram})^*$

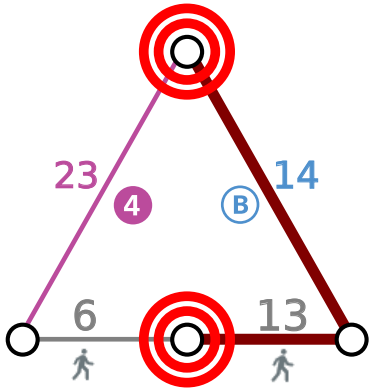
- Logic formula

$\forall X(rm \in X \wedge \forall xy$
 $(x \in X \wedge G(x, y) \rightarrow$
 $y \in X)) \rightarrow gn \in X$

- TRUE/FALSE

\leftrightarrow Model checking

Probabilistic query evaluation



Départ

20h17 - 5 rue Monticelli, Paris

1.1 km | 13 min

20h30 - CITE UNIVERSITAIRE, Paris

RER B - EPAU

Vers Aéroport CDG Terminal 2 TGV

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 | [f](#) [t](#) [m](#)



Panne du RER B : trafic interrompu entre

Paris : pourquoi il y a autant de perturbations sur le RER B et à Gare du Nord

La circulation de l'ensemble des trains au départ de gare du Nord est totalement interrompue à la suite d'une panne électrique.



Panne du RER B : trafic interrompu entre

Paris : pourquoi il y a autant de perturbations sur le RER B ?

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 pourrions développer

Panne du RER B : trafic interrompu entre

Paris : pourquoi il y a autant de perturbations sur le RER B

INCIDENT SUR LE RER B · OUF

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.

© 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 nous pourrions développer

Panne du RER B : trafic interrompu entre

Paris : pourquoi il y a autant de perturbations sur le

INCIDENT SUR LE RER B · OUF

ACTUALITÉS

Le RER B en panne, les voyageurs n'ont pas eu

d'autre choix que de descendre sur les voies

RER B et D en panne, gare du Nord paralysée, pollution: deuxième jour de galère

Actualité / Société / Trafic / Par Iris Péron, publié le 07/12/2016 à 13:40 , mis à jour à 16:07

partages



Partager



Tweeter



Partager



réaction

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 pourrions développer

Panne du RER B : trafic interrompu entre

Paris : pourquoi il y a autant de perturbations sur le

INCIDENT SUR LE RER B • OUF

ACTUALITÉS

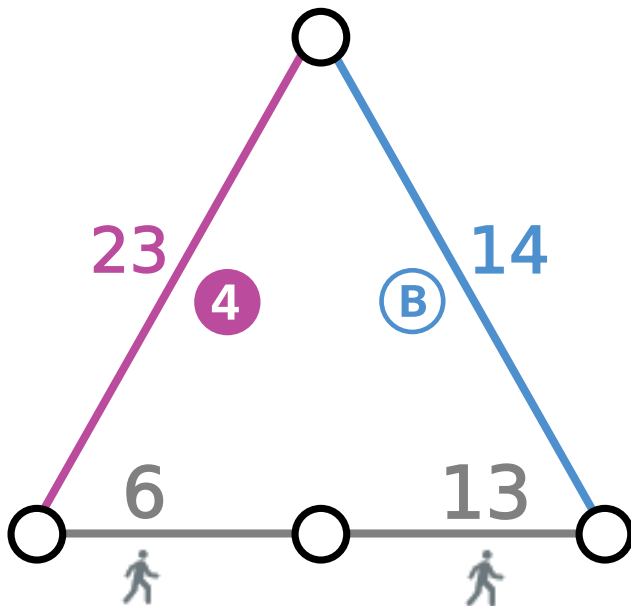
Le RER B en panne, les voyageurs n'ont pas eu

Ile-de-France : le trafic toujours interrompu sur le RER B entre Aulnay-sous-Bois et Roissy

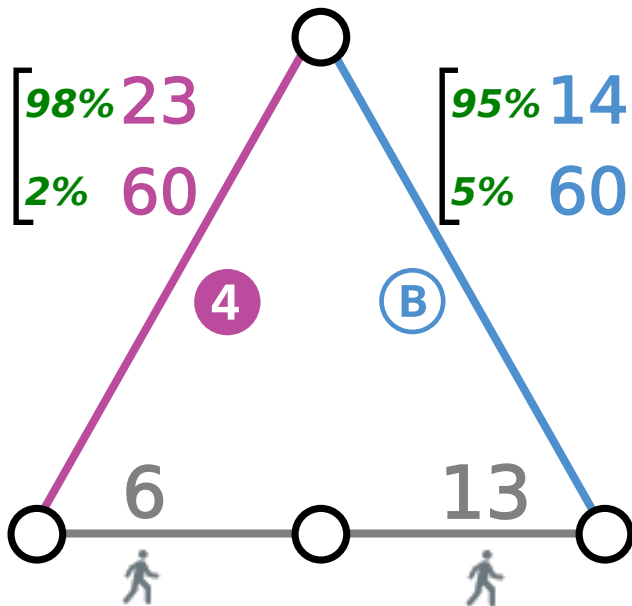
La circulation est arrêtée depuis mardi matin en raison d'une panne de caténaire. Le retour à la normale a été plusieurs fois retardé mais devrait avoir lieu mercredi vers 16 heures, selon la SNCF.

Le Monde | 07/12/2016 à 10h48 • Mis à jour le 07/12/2016 à 16h09

Probabilistic query evaluation



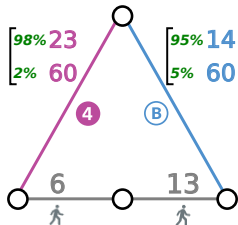
Probabilistic query evaluation



Probabilistic query evaluation



Probabilistic database



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

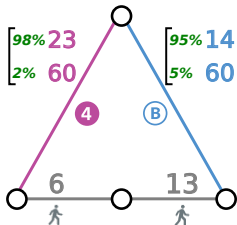
Probabilistic query evaluation



**Probabilistic
database**



Query



Rechercher mon itinéraire

De

A

Arrivée à

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

- Regular path
(Metro|RER)*
|(Bus|Tram)*
- Logic formula
$$\forall X (rm \in X \wedge \forall xy$$
$$(x \in X \wedge G(x, y) \rightarrow$$
$$y \in X)) \rightarrow gn \in X$$

Probabilistic query evaluation

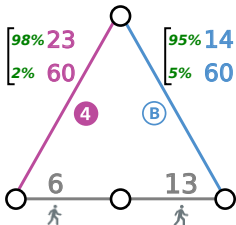


**Probabilistic
database**



Query

**Probabilistic
Result**



+

Rechercher mon itinéraire

De Rue Monticelli

A Gare du Nord

Aujourd'hui

Arrivée à 21 h



Départ
20h14 - 5 rue Monticelli, Paris

425 m | 6 min

20h20 - Porte d'Orléans (Général Leclerc), Paris

Métro 4
Vers Porte de Clignancourt

20 arrêts | 23 min

Arrivée
20h43 - Gare du Nord, Paris

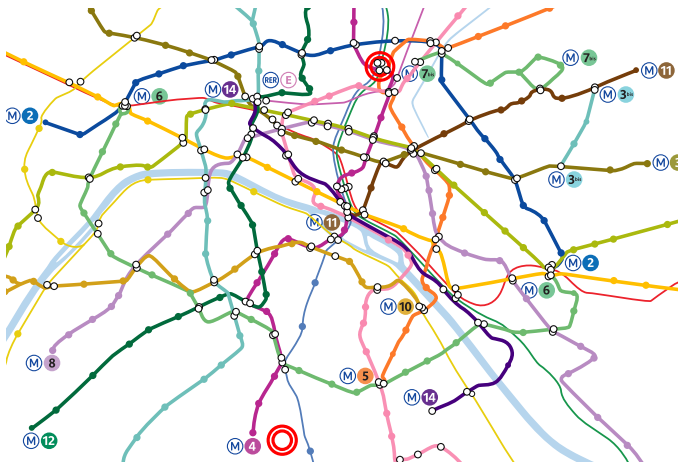
proba to be on time: **98%**

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

- Regular path
(Metro|RER)*
|(Bus|Tram)*
- Logic formula
$$\forall X (rm \in X \wedge \forall xy$$
$$(x \in X \wedge G(x, y) \rightarrow$$
$$y \in X)) \rightarrow gn \in X$$

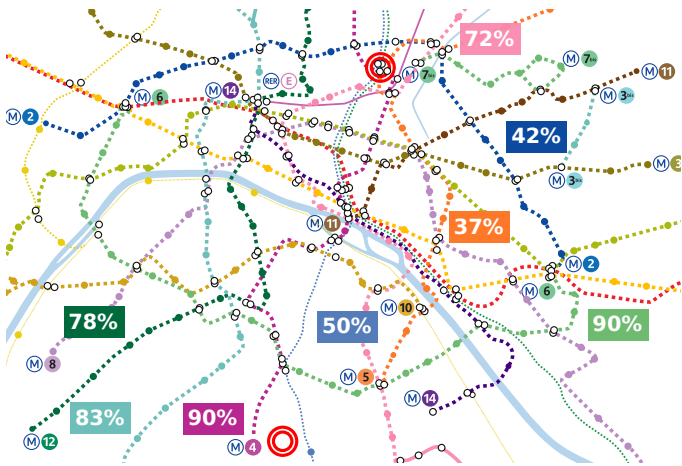
- Probability according to the input distribution

Computational complexity



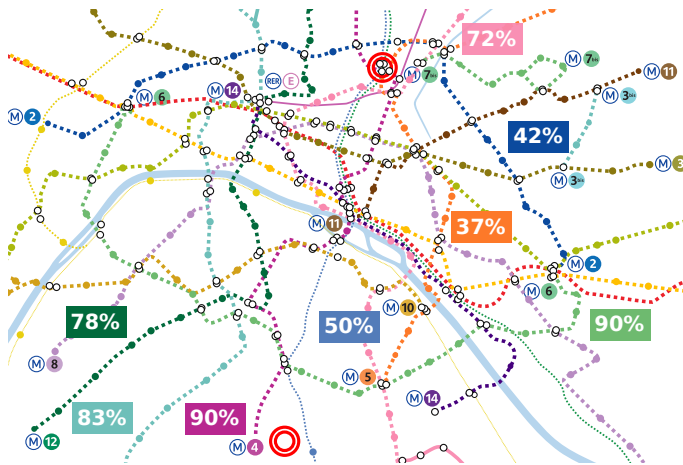
- Computing paths on a large graph:
 - Well-studied problem, efficient algorithms

Computational complexity



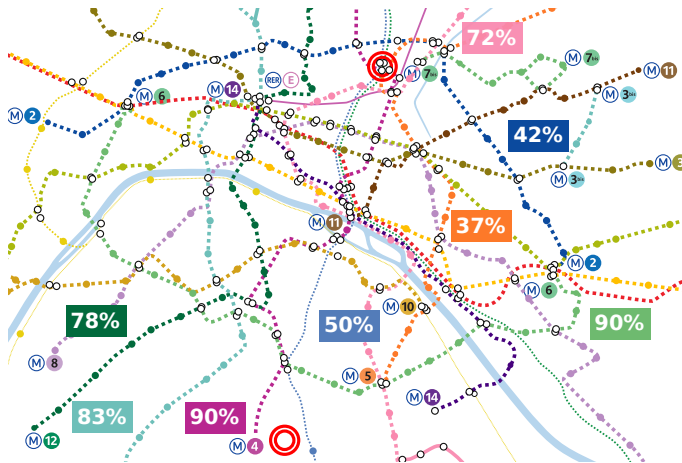
- Computing paths on a large **probabilistic** graph:
→ ???

Computational complexity



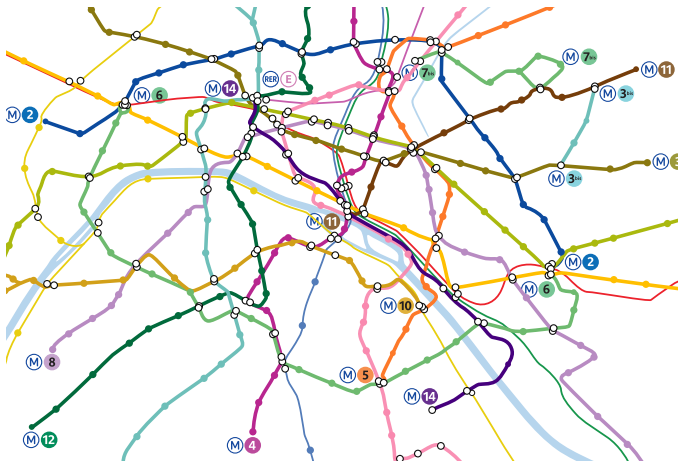
- Computing paths on a large **probabilistic** graph:
→ **Exponential** number of possibilities

Computational complexity

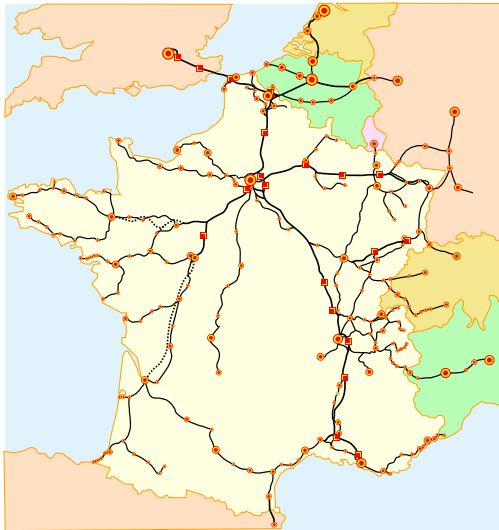


- Computing paths on a large **probabilistic** graph:
 - **Exponential** number of possibilities
 - **#P-hard** computational complexity in the **database**

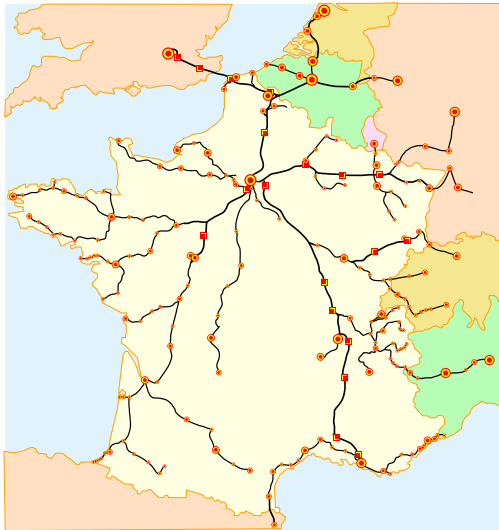
Idea: use the structure of data



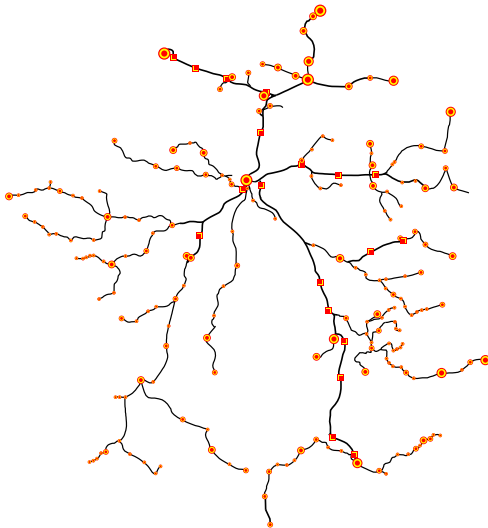
Idea: use the structure of data



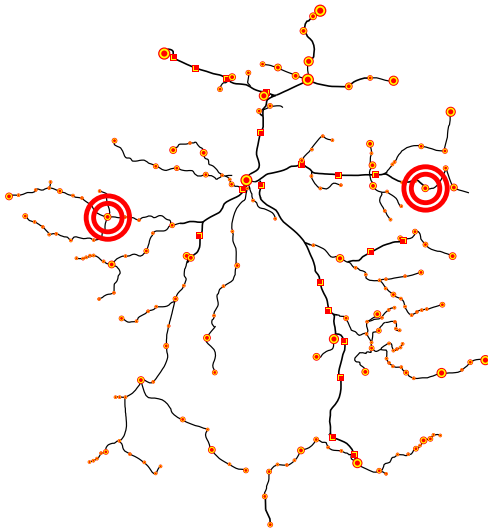
Idea: use the structure of data



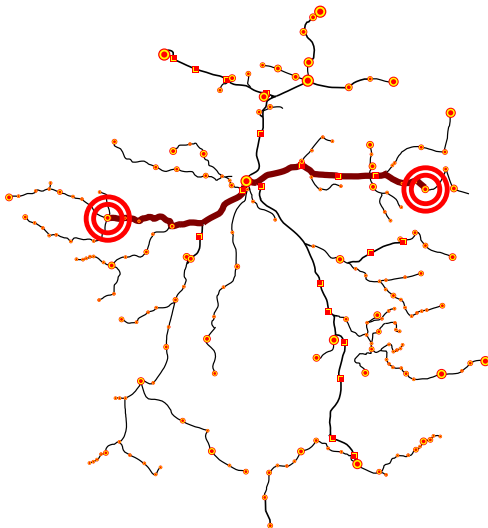
Idea: use the structure of data



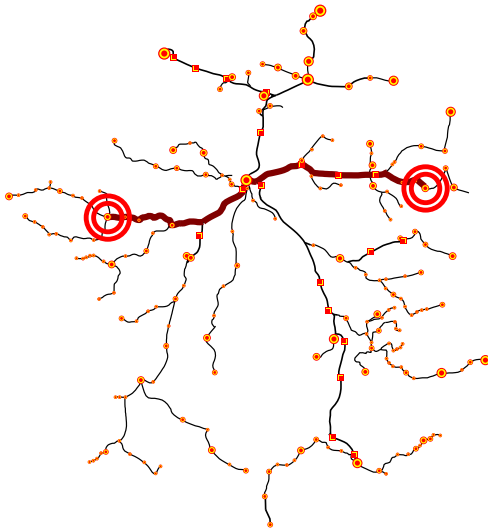
Idea: use the structure of data



Idea: use the structure of data



Idea: use the structure of data



→ Shortest path: **very easy** on a **large tree**

Leveraging the structure of uncertain data

*Does query evaluation on probabilistic data have **lower complexity** when the **structure** of the data is **close to a tree**?*

Leveraging the structure of uncertain data

*Does query evaluation on probabilistic data have **lower complexity** when the **structure** of the data is **close to a tree**?*

In this talk:

- Existing tools for **non-probabilistic data**:

Leveraging the structure of uncertain data

*Does query evaluation on probabilistic data have **lower complexity** when the **structure** of the data is **close to a tree**?*

In this talk:

- Existing tools for **non-probabilistic data**:
 - **Tree automata**, to evaluate queries on trees

Leveraging the structure of uncertain data

*Does query evaluation on probabilistic data have **lower complexity** when the **structure** of the data is **close to a tree**?*

In this talk:

- Existing tools for **non-probabilistic data**:
 - **Tree automata**, to evaluate queries on trees
 - **Treewidth**, formalizes the notion of being “close to a tree”

Leveraging the structure of uncertain data

*Does query evaluation on probabilistic data have **lower complexity** when the **structure** of the data is **close to a tree**?*

In this talk:

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

Leveraging the structure of uncertain data

*Does query evaluation on probabilistic data have **lower complexity** when the **structure** of the data is **close to a tree**?*

In this talk:

- Existing tools for **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**:

Leveraging the structure of uncertain data

*Does query evaluation on probabilistic data have **lower complexity** when the **structure** of the data is **close to a tree**?*

In this talk:

- Existing tools for **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

Leveraging the structure of uncertain data

*Does query evaluation on probabilistic data have **lower complexity** when the **structure** of the data is **close to a tree**?*

In this talk:

- Existing tools for **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
 - Application to **probabilistic query evaluation**

Leveraging the structure of uncertain data

*Does query evaluation on probabilistic data have **lower complexity** when the **structure** of the data is **close to a tree**?*

In this talk:

- Existing tools for **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
 - Application to **probabilistic query evaluation**
- **Other application**: enumeration of query results

Table of contents

Introduction



Existing tools for non-probabilistic data

Provenance circuits and probabilistic query evaluation

Application to enumeration

Query evaluation on words



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



Query evaluation on words



Database: a **word** w where nodes have a color from an alphabet $\bigcirc \text{ (white)} \text{ (pink)} \text{ (blue)}$



Query Q : a **sentence** (yes/no question) in **monadic second-order logic** (MSO)

“Is there both a pink and a blue node?”

Query evaluation on words



Database: a **word** w where nodes have a color from an alphabet $\bigcirc \text{ (white)} \text{ } \textcolor{red}{\bigcirc} \text{ (pink)} \text{ } \textcolor{blue}{\bigcirc} \text{ (blue)}$



Query Q : a **sentence** (yes/no question) in **monadic second-order logic** (MSO)

“Is there both a pink and a blue node?”



Result: TRUE/FALSE indicating if the word w satisfies the query Q

Monadic second-order logic (MSO)



- $P_{\text{blue}}(x)$ means “ x is blue”; also $P_{\text{red}}(x)$, $P_{\text{white}}(x)$
- $x \rightarrow y$ means “ x is the predecessor of y ”

Monadic second-order logic (MSO)



- $P_{\text{blue}}(x)$ means “ x is blue”; also $P_{\text{pink}}(x)$, $P_{\text{white}}(x)$
- $x \rightarrow y$ means “ x is the predecessor of y ”
- **Propositional logic:** formulas with **AND** \wedge , **OR** \vee , **NOT** \neg
 - $P_{\text{pink}}(x) \wedge P_{\text{blue}}(y)$ means “Node x is pink and node y is blue”

Monadic second-order logic (MSO)



- $P_{\text{blue}}(x)$ means “ x is blue”; also $P_{\text{pink}}(x)$, $P_{\text{white}}(x)$
- $x \rightarrow y$ means “ x is the predecessor of y ”
- **Propositional logic:** formulas with **AND** \wedge , **OR** \vee , **NOT** \neg
 - $P_{\text{pink}}(x) \wedge P_{\text{blue}}(y)$ means “Node x is pink and node y is blue”
- **First-order logic:** adds **existential quantifier** \exists and **universal quantifier** \forall
 - $\exists x y P_{\text{pink}}(x) \wedge P_{\text{blue}}(y)$ means “There is both a pink and a blue node”

Monadic second-order logic (MSO)



- $P_{\bullet}(x)$ means “ x is blue”; also $P_{\bullet}(x)$, $P_{\circ}(x)$
- $x \rightarrow y$ means “ x is the predecessor of y ”
- **Propositional logic:** formulas with **AND** \wedge , **OR** \vee , **NOT** \neg
 - $P_{\circ}(x) \wedge P_{\bullet}(y)$ means “Node x is pink and node y is blue”
- **First-order logic:** adds **existential quantifier** \exists and **universal quantifier** \forall
 - $\exists x y P_{\circ}(x) \wedge P_{\bullet}(y)$ 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 ”
 - $\forall x P_{\circ}(x) \Rightarrow \exists y P_{\bullet}(y) \wedge x \rightarrow^* y$
means “There is a blue node after every pink node”

Word automata

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

Alphabet:  **w:**  **Q:** $\exists x y P_{\text{red}}(x) \wedge P_{\text{blue}}(y)$

Word automata

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

Alphabet:  **w:**  **Q:** $\exists x y P_{\text{red}}(x) \wedge P_{\text{blue}}(y)$

- States:** $\{\perp, B, P, \top\}$

Word automata

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




Alphabet:  **w:**  **Q:** $\exists x y P_{\text{red}}(x) \wedge P_{\text{blue}}(y)$

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

Word automata



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




Alphabet:  **w:**  **Q:** $\exists x y P_{\text{red}}(x) \wedge P_{\text{blue}}(y)$

- **States:** $\{\perp, B, P, \top\}$
- **Final states:** $\{\top\}$
- **Initial function:**  \perp  P  B

Word automata



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




Alphabet:  **w:**  **Q:** $\exists x y P_{\text{red}}(x) \wedge P_{\text{blue}}(y)$

- **States:** $\{\perp, B, P, \top\}$
- **Final states:** $\{\top\}$
- **Initial function:**  \perp  P  B

Word automata



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




Alphabet:  **w:**  **Q:** $\exists x y P_{\text{red}}(x) \wedge P_{\text{blue}}(y)$

- **States:** $\{\perp, B, P, \top\}$
- **Final states:** $\{\top\}$
- **Initial function:**  \perp  P  B
- **Transitions** (examples): $\perp \xrightarrow{P} P$ $P \xrightarrow{B} \top$ $\top \xrightarrow{P} \top$

Word automata



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




Alphabet:  **w:**  **Q:** $\exists x y P_{\text{red}}(x) \wedge P_{\text{blue}}(y)$

- **States:** $\{\perp, B, P, \top\}$
- **Final states:** $\{\top\}$
- **Initial function:**  \perp  P  B
- **Transitions** (examples): $\perp \xrightarrow{\quad} \text{red circle } P$ $P \xrightarrow{\quad} \text{blue circle } \top$ $\top \xrightarrow{\quad} \text{red circle } \top$

Word automata



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




Alphabet:  **w:**  **Q:** $\exists x y P_{\text{red}}(x) \wedge P_{\text{blue}}(y)$

- **States:** $\{\perp, B, P, \top\}$
- **Final states:** $\{\top\}$
- **Initial function:**  \perp  P  B
- **Transitions** (examples): $\perp \xrightarrow{P} P$ $P \xrightarrow{B} \top$ $\top \xrightarrow{P} \top$

Word automata



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




Alphabet:  **w:**  **Q:** $\exists x y P_{\text{red}}(x) \wedge P_{\text{blue}}(y)$

- **States:** $\{\perp, B, P, \top\}$
- **Final states:** $\{\top\}$
- **Initial function:**  \perp  P  B
- **Transitions** (examples): $\perp \xrightarrow{P} P$ $P \xrightarrow{\top} B$ $B \xrightarrow{\top} P$

Word automata



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




Alphabet:  **w:**  **Q:** $\exists x y P_{\text{red}}(x) \wedge P_{\text{blue}}(y)$

- **States:** $\{\perp, B, P, \top\}$
- **Final states:** $\{\top\}$
- **Initial function:**  \perp  P  B
- **Transitions** (examples): $\perp \xrightarrow{P} P$ $P \xrightarrow{\top} \top$ $\top \xrightarrow{P} \top$

Word automata

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

Alphabet:  **w:**  **Q:** $\exists x y P_{\text{red}}(x) \wedge P_{\text{blue}}(y)$



- **States:** $\{\perp, B, P, \top\}$
- **Final states:** $\{\top\}$
- **Initial function:**  \perp  P  B
- **Transitions** (examples): $\perp \xrightarrow{P} P$ $P \xrightarrow{\top} \top$ $\top \xrightarrow{P} \top$




Theorem (Büchi, 1960)

MSO and word automata have the same expressive power on words

Word automata

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

Alphabet:  **w:**  **Q:** $\exists x y P_{\text{red}}(x) \wedge P_{\text{blue}}(y)$

- **States:** $\{\perp, B, P, \top\}$
- **Final states:** $\{\top\}$
- **Initial function:**  \perp  P  B
- **Transitions** (examples): $\perp \xrightarrow{P} P$ $P \xrightarrow{\top} \top$ $\top \xrightarrow{P} \top$

Theorem (Büchi, 1960)

MSO and word automata have the same expressive power on words

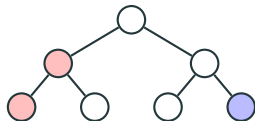
Corollary

Query evaluation of MSO on words is in linear time.

Query evaluation on trees



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



Query evaluation on trees

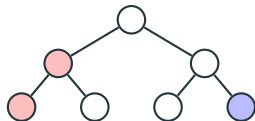


Database: a **tree** T where nodes have a color from an alphabet $\bigcirc \bigcirc \bigcirc$



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

- $P_{\bigcirc}(x)$ means “ x is blue”
- $x \rightarrow y$ means “ x is the parent of y ”



“Is there both a pink and a blue node?”

$$\exists x y P_{\bigcirc}(x) \wedge P_{\bigcirc}(y)$$

Query evaluation on trees



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

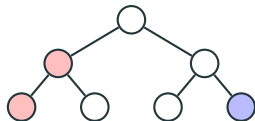


Query Q : a **sentence** in monadic second-order logic (MSO)

- $P_{\bigcirc}(x)$ means “ x is blue”
- $x \rightarrow y$ means “ x is the parent of y ”



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



“Is there both a pink and a blue node?”

$$\exists x y P_{\bigcirc}(x) \wedge P_{\bigcirc}(y)$$

Query evaluation on trees



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



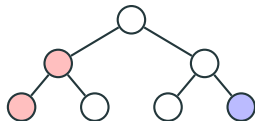
Query Q : a **sentence** in monadic second-order logic (MSO)

- $P_{\bigcirc}(x)$ means “ x is blue”
- $x \rightarrow y$ means “ x is the parent of y ”



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

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

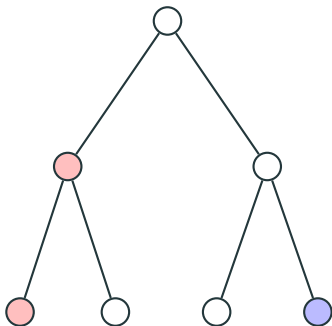


“Is there both a pink and a blue node?”

$$\exists x y P_{\bigcirc}(x) \wedge P_{\bigcirc}(y)$$

Tree automata

Tree alphabet:

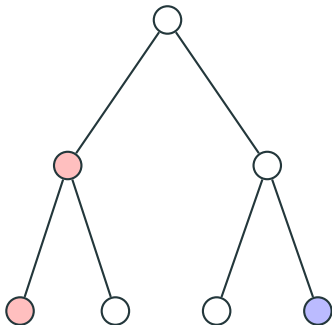


Tree automata

Tree alphabet:

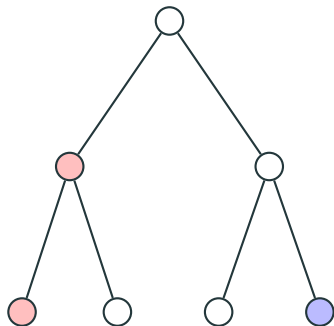


- Bottom-up deterministic **tree automaton**
- *“Is there both a pink and a blue node?”*



Tree automata

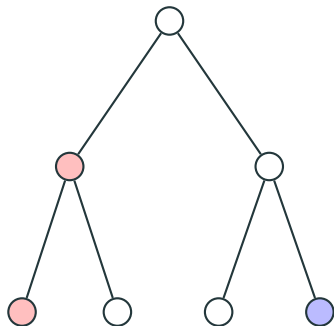
Tree alphabet:



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

Tree automata

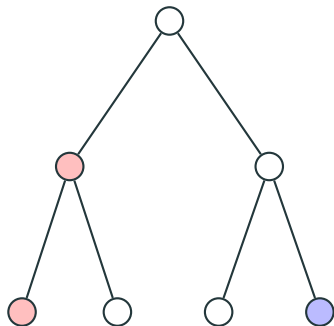
Tree alphabet:



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

Tree automata

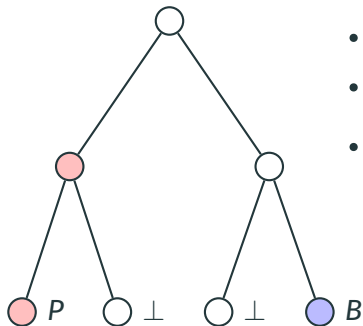
Tree alphabet:



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

Tree automata

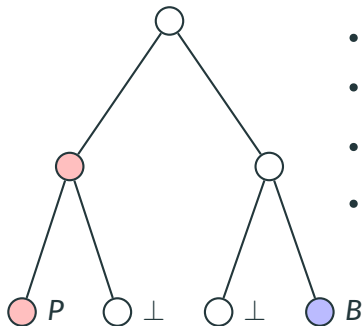
Tree alphabet:



- Bottom-up deterministic **tree automaton**
- “Is there both a pink and a blue node?”
- **States:** $\{\perp, B, P, \top\}$
- **Final states:** $\{\top\}$
- **Initial function:** $\bigcirc \perp \quad \textcolor{red}{\bigcirc} P \quad \textcolor{blue}{\bigcirc} B$

Tree automata

Tree alphabet:

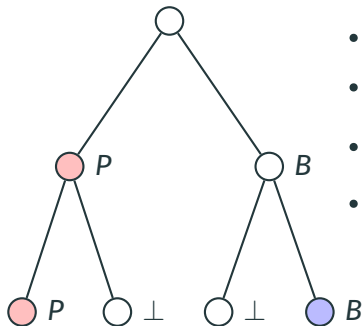


- Bottom-up deterministic **tree automaton**
- *"Is there both a pink and a blue node?"*
- **States:** $\{\perp, B, P, \top\}$
- **Final states:** $\{\top\}$
- **Initial function:** $\bigcirc \perp \quad \textcolor{red}{\bigcirc} P \quad \textcolor{blue}{\bigcirc} B$
- **Transitions** (examples):



Tree automata

Tree alphabet:

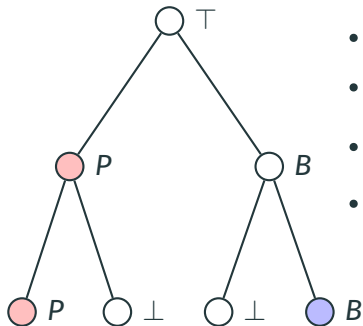


- Bottom-up deterministic **tree automaton**
- *"Is there both a pink and a blue node?"*
- **States:** $\{\perp, B, P, \top\}$
- **Final states:** $\{\top\}$
- **Initial function:** $\bigcirc \perp \quad \textcolor{red}{\bigcirc} P \quad \textcolor{blue}{\bigcirc} B$
- **Transitions** (examples):



Tree automata

Tree alphabet:

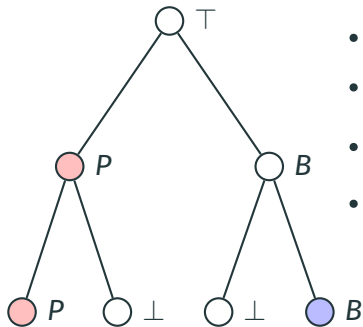


- Bottom-up deterministic **tree automaton**
- “Is there both a pink and a blue node?”
- **States:** $\{\perp, B, P, T\}$
- **Final states:** $\{T\}$
- **Initial function:** $\bigcirc \perp \quad \text{pink} P \quad \text{blue} B$
- **Transitions** (examples):



Tree automata

Tree alphabet:



- Bottom-up deterministic **tree automaton**
- “Is there both a pink and a blue node?”
- **States:** $\{\perp, B, P, T\}$
- **Final states:** $\{T\}$
- **Initial function:** $\bigcirc \perp \quad \textcolor{red}{\bigcirc} P \quad \textcolor{blue}{\bigcirc} B$
- **Transitions** (examples):

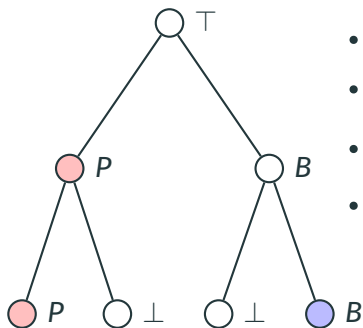


Theorem [Thatcher and Wright, 1968]

MSO and **tree automata** have the same **expressive power** on trees

Tree automata

Tree alphabet:



- Bottom-up deterministic **tree automaton**
- “Is there both a pink and a blue node?”
- **States:** $\{\perp, B, P, T\}$
- **Final states:** $\{T\}$
- **Initial function:** $\bigcirc \perp \quad \text{pink } P \quad \text{blue } B$
- **Transitions** (examples):



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**.

Query evaluation on trees



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



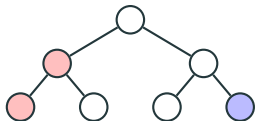
Query Q : a **sentence** in monadic second-order logic (MSO)

- $P_{\bigcirc}(x)$ means “ x is blue”
- $x \rightarrow y$ means “ x is the parent of y ”



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

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



*“Is there both a pink
and a blue node?”*

$$\exists x y P_{\bigcirc}(x) \wedge P_{\bigcirc}(y)$$

Query evaluation on treelike data



Database: a **treelike database** T

???



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

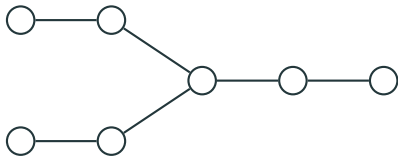
$(\text{Metro}|\text{RER})^*$
 $|\ (\text{Bus}|\text{Tram})^*$



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

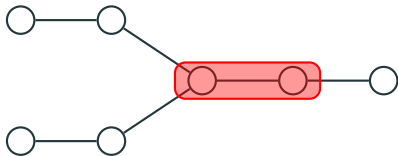
Computational complexity as a function of T
(the query Q is **fixed**)

Treewidth by example:



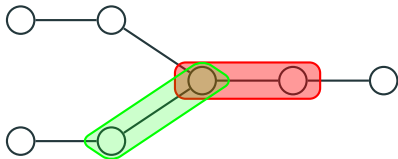
Treewidth

Treewidth by example:



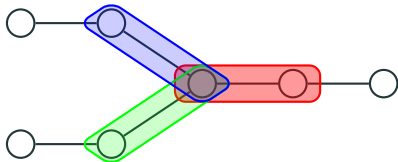
Treewidth

Treewidth by example:



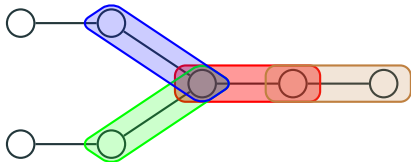
Treewidth

Treewidth by example:



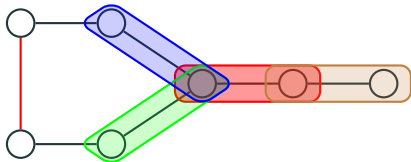
Treewidth

Treewidth by example:



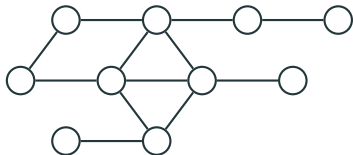
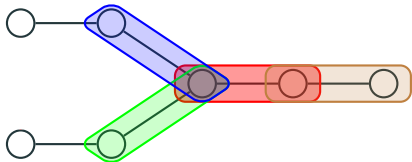
Treewidth

Treewidth by example:



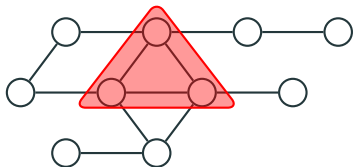
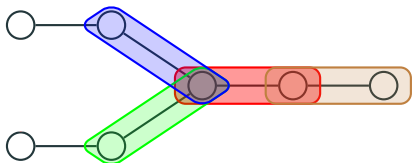
Treewidth

Treewidth by example:



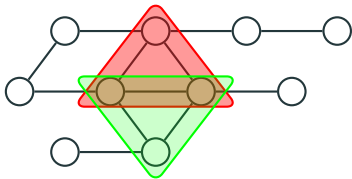
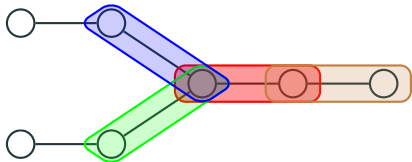
Treewidth

Treewidth by example:



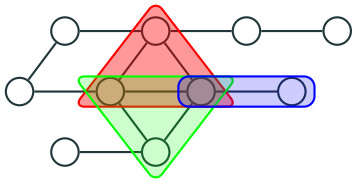
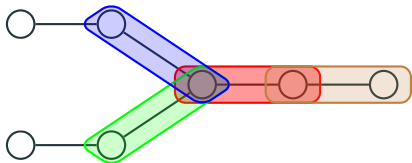
Treewidth

Treewidth by example:



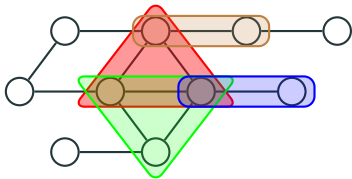
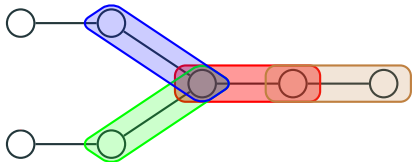
Treewidth

Treewidth by example:



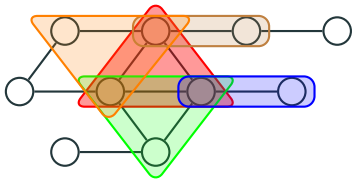
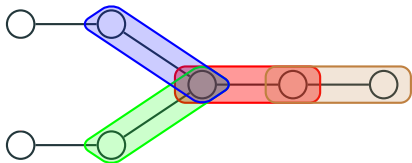
Treewidth

Treewidth by example:



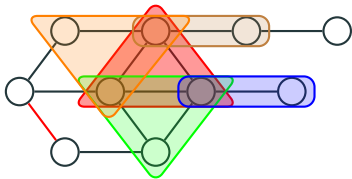
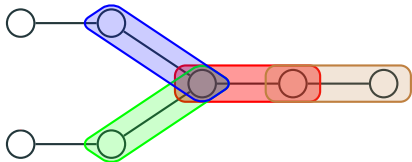
Treewidth

Treewidth by example:



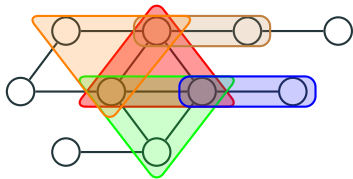
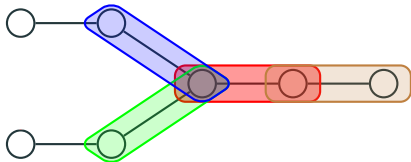
Treewidth

Treewidth by example:



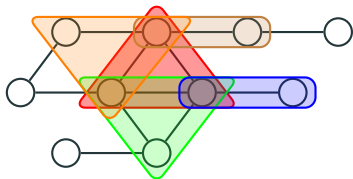
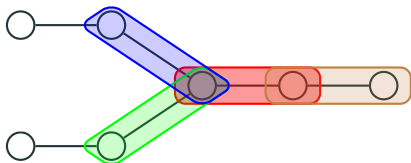
Treewidth

Treewidth by example:



Treewidth

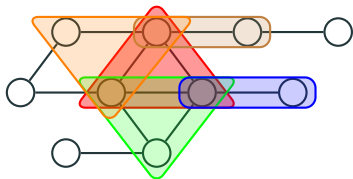
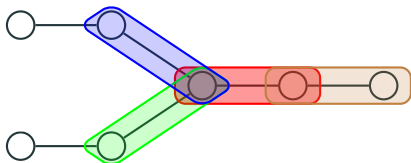
Treewidth by example:



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

Treewidth

Treewidth by example:

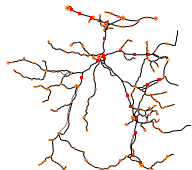


- **Trees** have treewidth **1**
- **Cycles** have treewidth **2**
- **k -cliques** and **$(k - 1)$ -grids** have treewidth **$k - 1$**

→ **Treelike**: the **treewidth** is bounded by a **constant**

Courcelle's theorem

Tree-like data

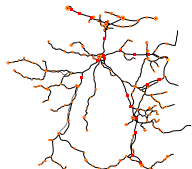


MSO query

$(\text{RER}|\text{metro})^*$
 $|(\text{bus}|\text{tram})^*$

Courcelle's theorem

Tree-like data

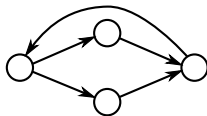


MSO query

$(\text{RER}|\text{metro})^*$
 $|(\text{bus}|\text{tram})^*$

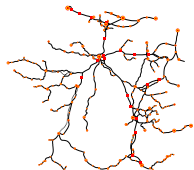


Tree automaton



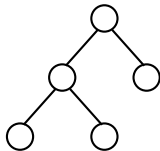
Courcelle's theorem

Tree**like** data



Tree **encoding**

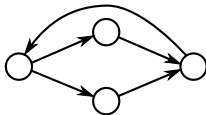
linear



MSO query

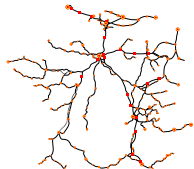
$(\text{RER}|\text{metro})^*$
 $|(\text{bus}|\text{tram})^*$

Tree **automaton**

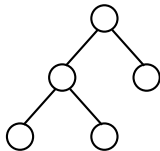


Courcelle's theorem

Tree**like** **data**



Tree **encoding**



linear

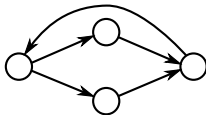
linear

Query
answer
TRUE

MSO query

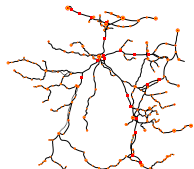
$(\text{RER}|\text{metro})^*$
 $|(\text{bus}|\text{tram})^*$

Tree **automaton**

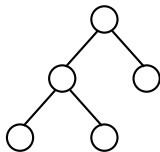


Courcelle's theorem

Tree**like data**



Tree **encoding**



linear

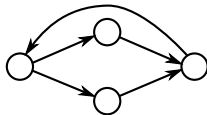
linear

Query
answer
TRUE

MSO query

$(\text{RER}|\text{metro})^*$
 $|(\text{bus}|\text{tram})^*$

Tree **automaton**



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

Table of contents

Introduction

Existing tools for non-probabilistic data

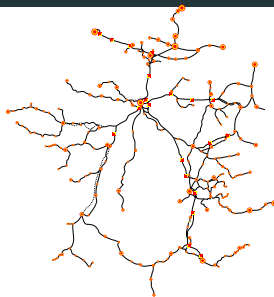
Provenance circuits and probabilistic query evaluation

Application to enumeration

Probabilistic query evaluation on treelike data



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



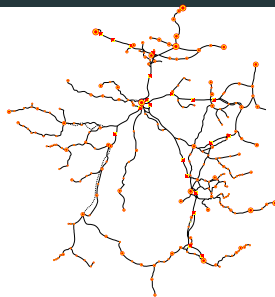
Probabilistic query evaluation on treelike data



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



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

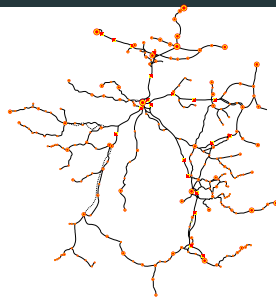


$(\text{Metro}|\text{RER})^*$
 $| (\text{Bus}|\text{Tram})^*$

Probabilistic query evaluation on treelike data



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



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

$(\text{Metro}|\text{RER})^*$
 $| (\text{Bus}|\text{Tram})^*$

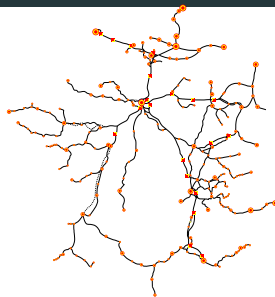


Result: **Probability** that the database D satisfies query Q

Probabilistic query evaluation on treelike data



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



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

$(\text{Metro}|\text{RER})^*$
 $| (\text{Bus}|\text{Tram})^*$

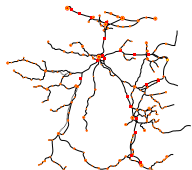


Result: **Probability** that the database D satisfies query Q

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

Roadmap

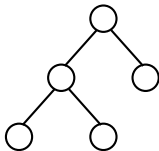
Tree**like data**



linear



Tree **encoding**

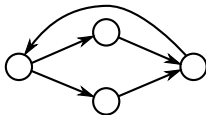


MSO query

$(\text{RER}|\text{metro})^*$
 $|(\text{bus}|\text{tram})^*$

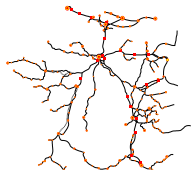


Tree **automaton**



Roadmap

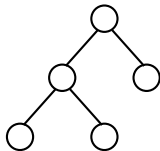
Tree-like **data**



linear



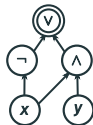
Tree **encoding**



linear



Provenance
circuit

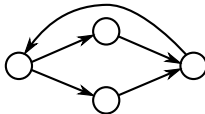


MSO query

$(\text{RER}|\text{metro})^*$
 $|(\text{bus}|\text{tram})^*$

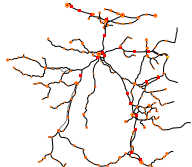


Tree **automaton**



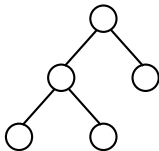
Roadmap

Probabilistic
treelike **data**



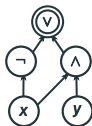
linear

Uncertain
tree **encoding**



linear

Provenance
circuit
+probabilities

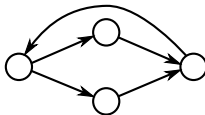


linear

MSO query

$(\text{RER}|\text{metro})^*$
 $|(\text{bus}|\text{tram})^*$

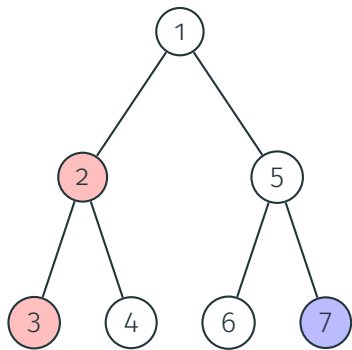
Tree **automaton**



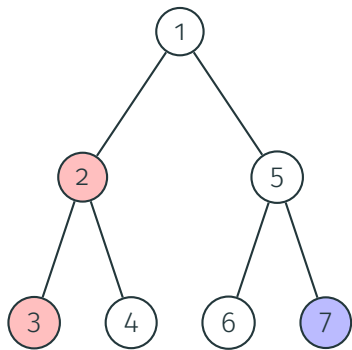
95%

Probability

Uncertain trees

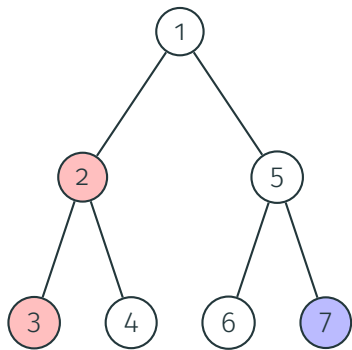


Uncertain trees



A **valuation** of a tree decides whether to **keep** (1) or **discard** (0) node labels

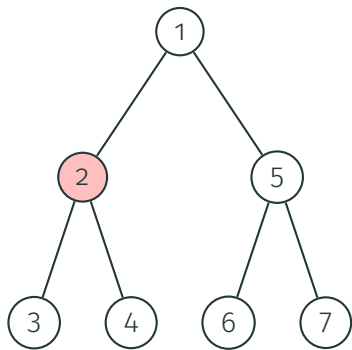
Uncertain trees



A **valuation** of a tree decides whether to **keep** (1) or **discard** (0) node labels

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

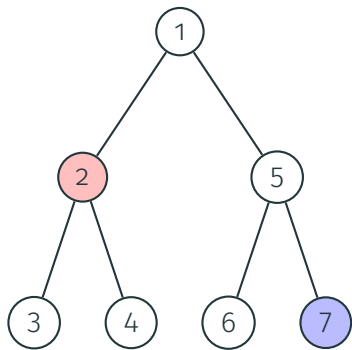
Uncertain trees



A **valuation** of a tree decides whether to **keep** (1) or **discard** (0) node labels

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

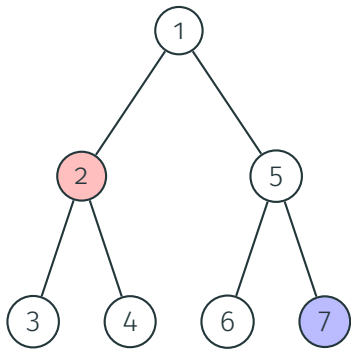
Uncertain trees



A **valuation** of a tree decides whether to **keep** (1) or **discard** (0) node labels

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

Uncertain trees

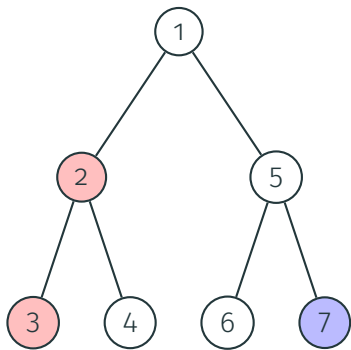


A **valuation** of a tree decides whether to **keep** (1) or **discard** (0) node labels

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

A: “Is there both a pink and a blue node?”

Uncertain trees



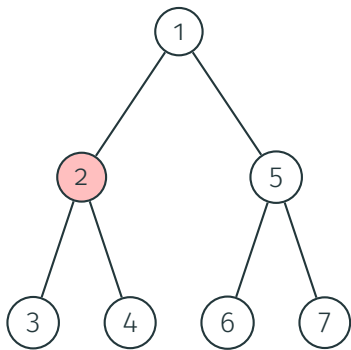
A **valuation** of a tree decides whether to **keep** (1) or **discard** (0) node labels

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

A: “Is there both a pink and a blue node?”

The tree automaton **A** **accepts**

Uncertain trees



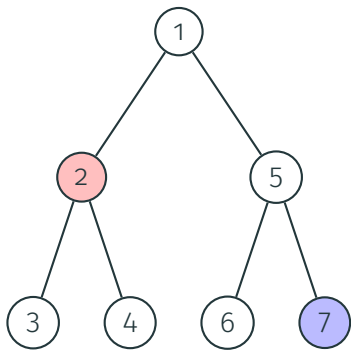
A **valuation** of a tree decides whether to **keep** (1) or **discard** (0) node labels

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

A: “Is there both a pink and a blue node?”

The tree automaton **A** **rejects**

Uncertain trees



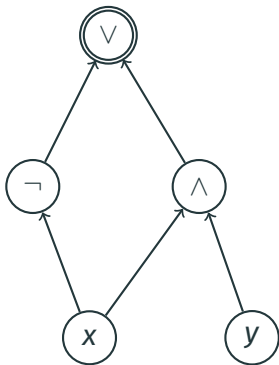
A **valuation** of a tree decides whether to **keep** (1) or **discard** (0) node labels

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

A: “Is there both a pink and a blue node?”

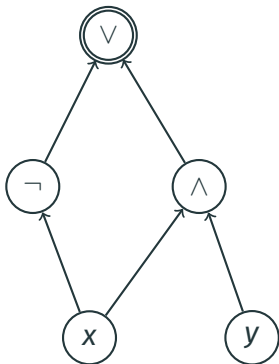
The tree automaton **A** **accepts**

Boolean circuit



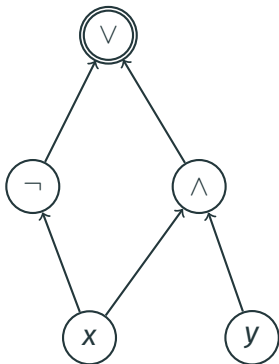
- Directed acyclic graph of **gates**

Boolean circuit



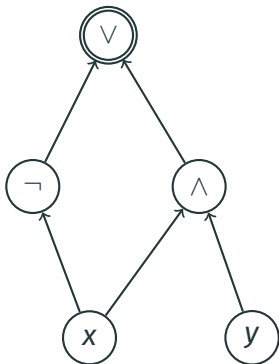
- Directed acyclic graph of **gates**
- **Output** gate: 

Boolean circuit



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

Boolean circuit



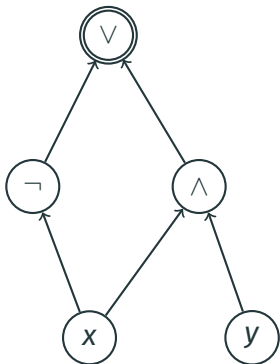
- Directed acyclic graph of **gates**

- **Output** gate:

- **Variable** gates:

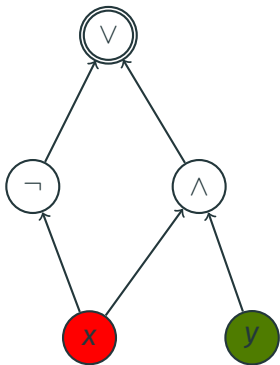
- **Internal** gates:






Boolean circuit



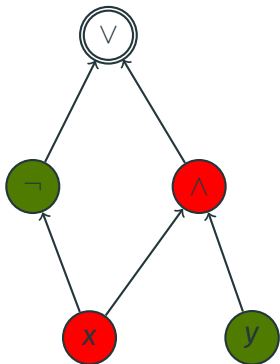
- Directed acyclic graph of **gates**
- **Output** gate:
- **Variable** gates:
- **Internal** gates:
- **Valuation**: function from variables to $\{\text{0}, \text{1}\}$
Example: $\nu = \{x \mapsto \text{0}, y \mapsto \text{1}\} \dots$

Boolean circuit



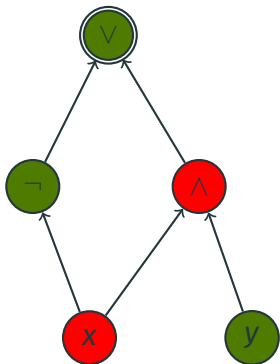
- Directed acyclic graph of **gates**
- **Output** gate: 
- **Variable** gates: 
- **Internal** gates:   
- **Valuation**: function from variables to $\{0, 1\}$
Example: $\nu = \{x \mapsto 0, y \mapsto 1\} \dots$




Boolean circuit



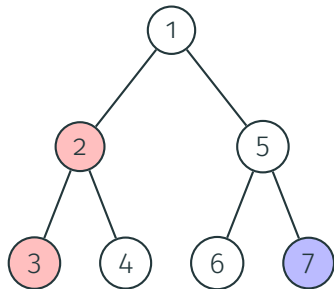
- Directed acyclic graph of **gates**
- **Output** gate:
- **Variable** gates:
- **Internal** gates:
- **Valuation**: function from variables to $\{0, 1\}$
Example: $\nu = \{x \mapsto 0, y \mapsto 1\} \dots$

Boolean circuit



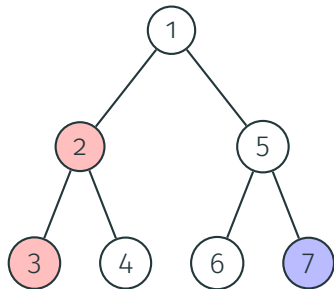
- Directed acyclic graph of **gates**
- **Output** gate: 
- **Variable** gates: 
- **Internal** gates: 
- **Valuation**: function from variables to $\{0, 1\}$
Example: $\nu = \{x \mapsto 0, y \mapsto 1\} \dots$ mapped to **1**

Provenance circuit



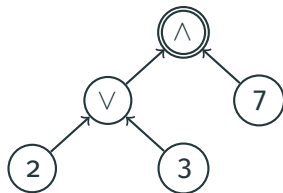
Query: *Is there both a pink and a blue node?*

Provenance circuit

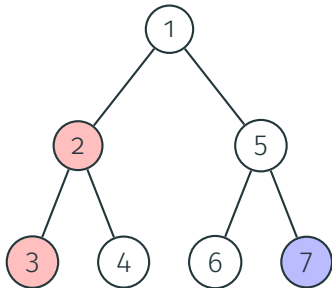


Query: *Is there both a pink and a blue node?*

Provenance circuit:

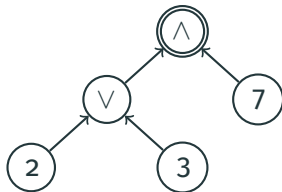


Provenance circuit



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

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 $\nu(T)$

Building provenance circuits on trees

Theorem

For any bottom-up *tree automaton* A and input *tree* T ,
we can build a *provenance circuit* of A on T in $O(|A| \times |T|)$

Building provenance circuits on trees

Theorem

For any bottom-up **tree automaton** A and input **tree** T , we can build a **provenance circuit** of A on T in $O(|A| \times |T|)$

- **Alphabet:** 

- **Automaton:** “Is there both a pink and a blue node?”

- **States:** $\{\perp, B, P, T\}$

- **Final:** $\{T\}$

- **Transitions:**



Building provenance circuits on trees

Theorem

For any bottom-up **tree automaton** A and input **tree** T , we can build a **provenance circuit** of A on T in $O(|A| \times |T|)$

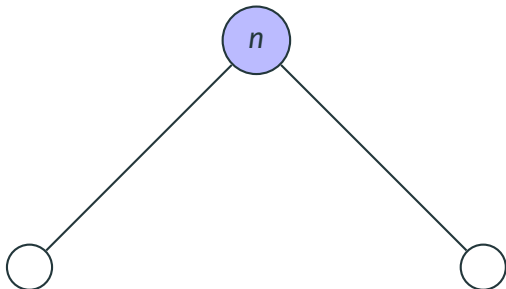
- **Alphabet:**   

- **Automaton:** “Is there both a pink and a blue node?”

- **States:** $\{\perp, B, P, \top\}$

- **Final:** $\{\top\}$

- **Transitions:**



Building provenance circuits on trees

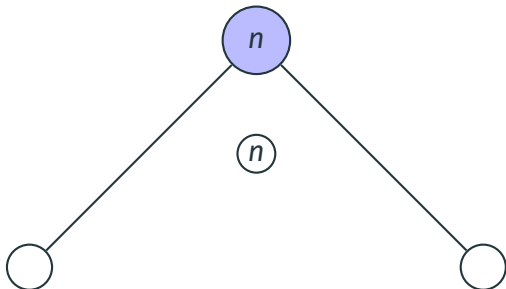
Theorem

For any bottom-up **tree automaton** A and input **tree** T , we can build a **provenance circuit** of A on T in $O(|A| \times |T|)$

- **Alphabet:** ○ ● ●
- **Automaton:** “Is there both a pink and a blue node?”

- **States:** $\{\perp, B, P, \top\}$
- **Final:** $\{\top\}$

- **Transitions:**



Building provenance circuits on trees

Theorem

For any bottom-up **tree automaton** A and input **tree** T , we can build a **provenance circuit** of A on T in $O(|A| \times |T|)$

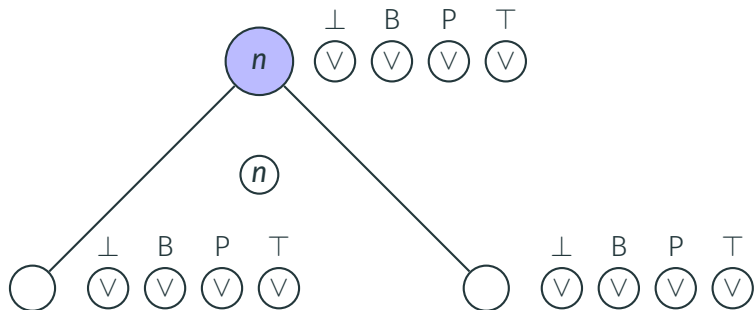
- Alphabet:** ○ ● ●

- Automaton:** “Is there both a pink and a blue node?”

- States:**
 $\{\perp, B, P, T\}$

- Final:** $\{T\}$

- Transitions:**



Building provenance circuits on trees

Theorem

For any bottom-up **tree automaton** A and input **tree** T , we can build a **provenance circuit** of A on T in $O(|A| \times |T|)$

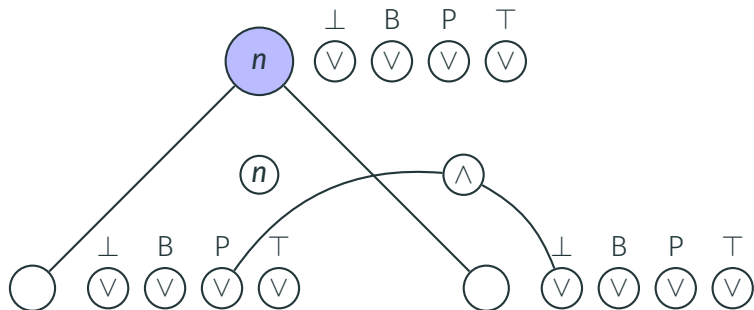
- Alphabet:** ○ ● ●

- Automaton:** “Is there both a pink and a blue node?”

- States:** $\{\perp, B, P, T\}$

- Final:** $\{T\}$

- Transitions:**



Building provenance circuits on trees

Theorem

For any bottom-up **tree automaton** A and input **tree** T , we can build a **provenance circuit** of A on T in $O(|A| \times |T|)$

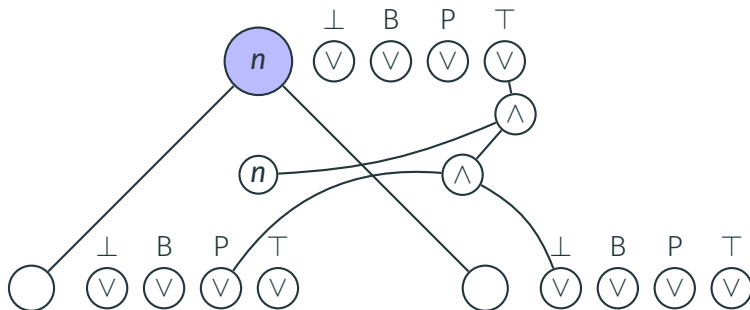
- Alphabet:** ○ ● ●

- Automaton:** “Is there both a pink and a blue node?”

- States:** $\{\perp, B, P, T\}$

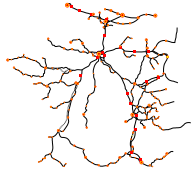
- Final:** $\{T\}$

- Transitions:**

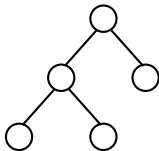


Probabilistic query evaluation

Probabilistic
treelike **data**



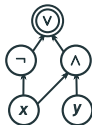
Uncertain
tree **encoding**



linear

linear

Provenance
circuit
+probabilities

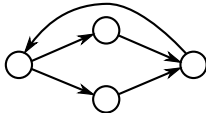


linear

MSO query

$(\text{RER}|\text{metro})^*$
 $|(\text{bus}|\text{tram})^*$

Tree **automaton**

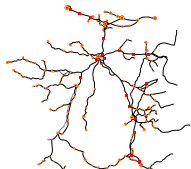


95%

Probability

Details of the approach

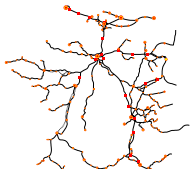
Probabilistic
treelike **data**



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

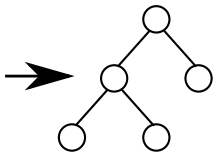
Details of the approach

Probabilistic
treelike **data**



Each **fact** can
disappear
with some
probability

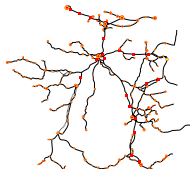
Uncertain
tree **encoding**



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

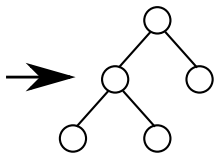
Details of the approach

Probabilistic
treelike **data**



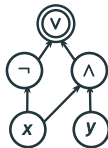
Each **fact** can
disappear
with some
probability

Uncertain
tree **encoding**



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

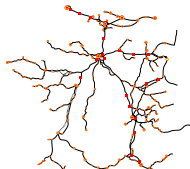
Provenance
circuit
+probabilities



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

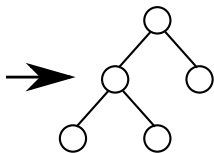
Details of the approach

Probabilistic
treelike **data**



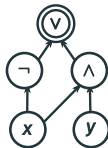
Each **fact** can **disappear** with some probability

Uncertain
tree **encoding**



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

Provenance
circuit
+probabilities



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

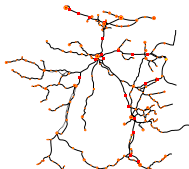
Probability

95%

Probability that the **circuit** evaluates to **true**

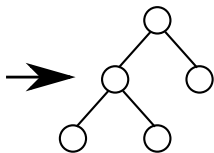
Details of the approach

Probabilistic
treelike **data**



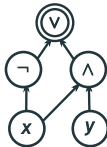
Each **fact** can **disappear** with some probability

Uncertain
tree **encoding**



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

Provenance
circuit
+probabilities



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

Probability

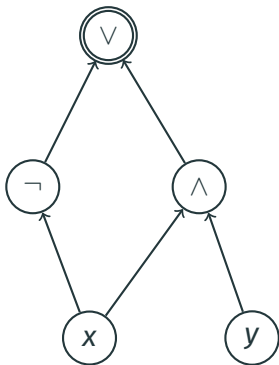
→ **95%**

Probability that the **circuit** evaluates to **true**

→ How to compute **efficiently** the probability of the circuit?

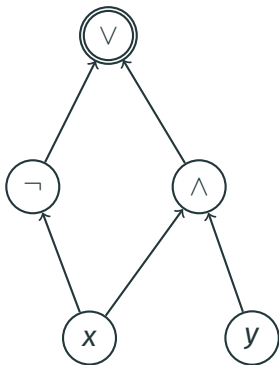
Computing the probability of a circuit

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



Computing the probability of a circuit

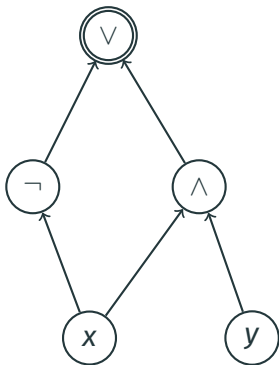
- We are given a **circuit** and a **probability** P for each variable



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

Computing the probability of a circuit

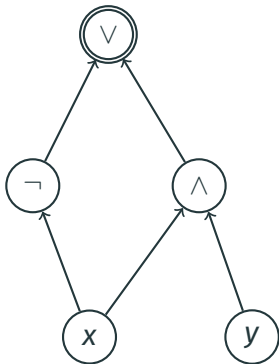
- We are given a **circuit** and a **probability** P for each variable
- Each variable x is true **independently** with probability $P(x)$



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

Computing the probability of a circuit

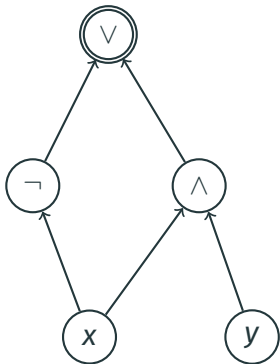
- We are given a **circuit** and a **probability** P for each variable
- Each variable x is true **independently** with probability $P(x)$
- What is the probability that the circuit **evaluates to true**?



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

Computing the probability of a circuit

- 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**?

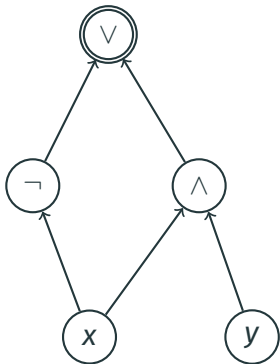


- In general, **#P-hard** (harder than SAT)

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

Computing the probability of a circuit

- 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**?

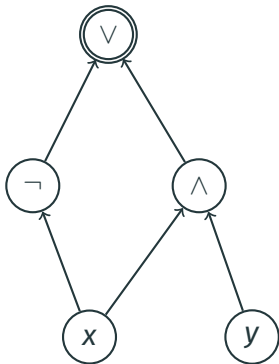


- In general, **#P-hard** (harder than SAT)
- Here it's **easy**:

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

Computing the probability of a circuit

- 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**?

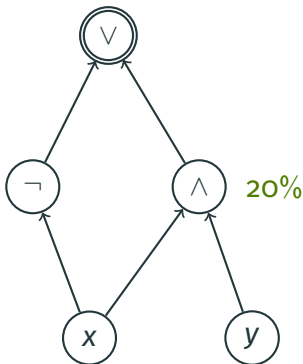


- In general, **#P-hard** (harder than SAT)
- Here it's **easy**:
 - The inputs to the **AND-gate** are **independent**

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

Computing the probability of a circuit

- 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**?

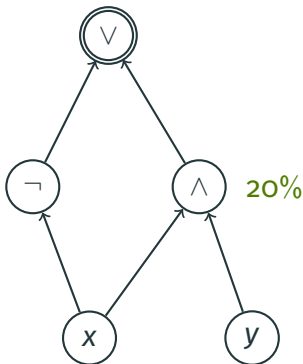


- In general, **#P-hard** (harder than SAT)
- Here it's **easy**:
 - The inputs to the **AND-gate** are **independent**

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

Computing the probability of a circuit

- 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**?

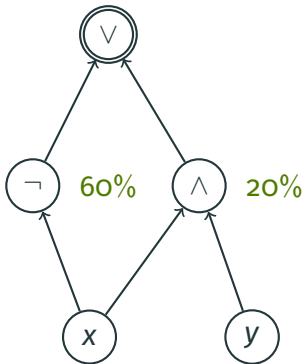


- In general, **#P-hard** (harder than SAT)
- Here it's **easy**:
 - The inputs to the **\wedge -gate** are **independent**
 - The **\neg -gate** has probability $1 - P(\text{input})$

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

Computing the probability of a circuit

- 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**?

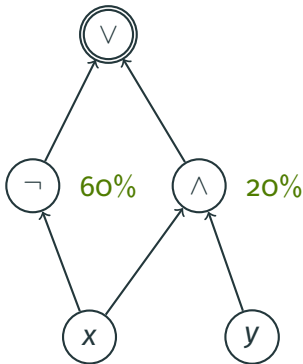


- In general, **#P-hard** (harder than SAT)
- Here it's **easy**:
 - The inputs to the **∧-gate** are **independent**
 - The **¬-gate** has probability $1 - P(\text{input})$

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

Computing the probability of a circuit

- 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**?

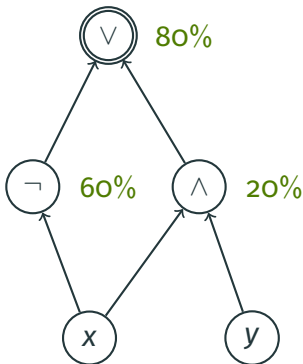


- In general, **#P-hard** (harder than SAT)
- Here it's **easy**:
 - The inputs to the \wedge -gate are **independent**
 - The \neg -gate has probability $1 - P(\text{input})$
 - The \vee -gate has **mutually exclusive** inputs

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

Computing the probability of a circuit

- 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**?

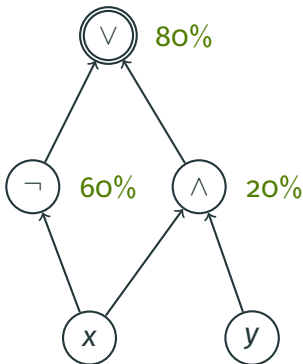


- In general, **#P-hard** (harder than SAT)
- Here it's **easy**:
 - The inputs to the **∧-gate** are **independent**
 - The **¬-gate** has probability $1 - P(\text{input})$
 - The **∨-gate** has **mutually exclusive** inputs

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

Computing the probability of a circuit

- 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**?



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

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

d-DNNFs

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



d-DNNFs

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

-  gates only have **variables** as inputs

d-DNNFs

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

-  gates only have **variables** as inputs
-  gates always have **mutually exclusive** inputs

d-DNNFs

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

- \neg gates only have **variables** as inputs
- \vee gates always have **mutually exclusive** inputs
- \wedge gates are all on **independent** inputs

d-DNNFs

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

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

- \neg gates only have **variables** as inputs
- \vee gates always have **mutually exclusive** inputs
- \wedge gates are all on **independent** inputs

d-DNNFs

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

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

- \neg gates only have **variables** as inputs
- \vee gates always have **mutually exclusive** inputs
- \wedge gates are all on **independent** inputs



d-DNNFs

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

- \neg gates only have **variables** as inputs
- \vee gates always have **mutually exclusive** inputs
- \wedge gates are all on **independent** inputs

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



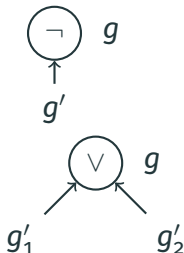
$$P(g) := 1 - P(g')$$

d-DNNFs

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

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

- \neg gates only have **variables** as inputs
- \vee gates always have **mutually exclusive** inputs
- \wedge gates are all on **independent** inputs



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

d-DNNFs

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

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

- \neg gates only have **variables** as inputs
- \vee gates always have **mutually exclusive** inputs
- \wedge gates are all on **independent** inputs



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



$$P(g) := P(g'_1) + P(g'_2)$$

d-DNNFs

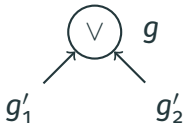
The circuit is a **d-DNNF**...

- \neg gates only have **variables** as inputs
- \vee gates always have **mutually exclusive** inputs
- \wedge gates are all on **independent** inputs

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



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



$$P(g) := P(g'_1) + P(g'_2)$$



d-DNNFs

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

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

- \neg gates only have **variables** as inputs
- \vee gates always have **mutually exclusive** inputs
- \wedge gates are all on **independent** inputs



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



$$P(g) := P(g'_1) + P(g'_2)$$



$$P(g) := P(g'_1) \times P(g'_2)$$

d-DNNFs

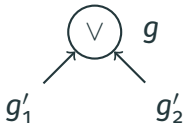
The circuit is a **d-DNNF**...

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

- \neg gates only have **variables** as inputs
- \vee gates always have **mutually exclusive** inputs
- \wedge gates are all on **independent** inputs



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



$$P(g) := P(g'_1) + P(g'_2)$$



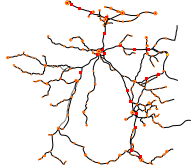
$$P(g) := P(g'_1) \times P(g'_2)$$

Lemma

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

Final result

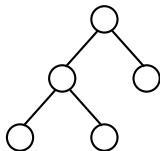
Probabilistic
treelike **data**



linear



Uncertain
tree **encoding**



linear



Provenance
circuit
+probabilities

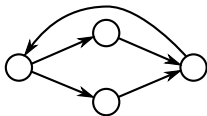


MSO query

$(\text{RER}|\text{metro})^*$
 $|(\text{bus}|\text{tram})^*$

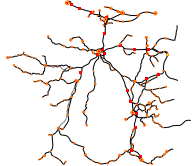


Tree **automaton**



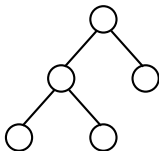
Final result

Probabilistic
treelike **data**



linear

Uncertain
tree **encoding**



linear

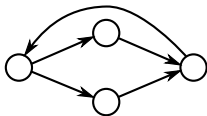
Provenance
d-DNNF
+probabilities



MSO query

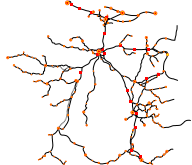
$(\text{RER}|\text{metro})^*$
 $|(\text{bus}|\text{tram})^*$

Tree **automaton**

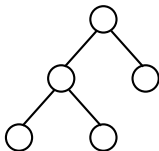


Final result

Probabilistic
treelike **data**



Uncertain
tree **encoding**



linear

linear

Provenance
d-DNNF
+probabilities



linear

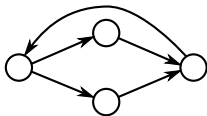
95%

Probability

MSO query

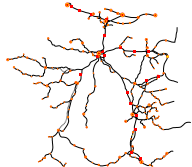
$(\text{RER}|\text{metro})^*$
 $|(\text{bus}|\text{tram})^*$

Tree **automaton**

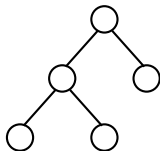


Final result

Probabilistic
treelike **data**



Uncertain
tree **encoding**



Provenance
d-DNNF
+probabilities



linear

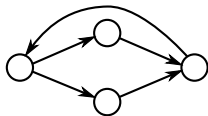
linear

linear

MSO query

$(\text{RER}|\text{metro})^*$
 $|(\text{bus}|\text{tram})^*$

Tree **automaton**



95%

Probability

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

Table of contents

Introduction

Existing tools for non-probabilistic data

Provenance circuits and probabilistic query evaluation

Application to enumeration

Non-Boolean queries

- We have studied **Boolean queries**:

“Is there both a pink and a blue node?”

$$Q() : \exists x y P_{\text{pink}}(x) \wedge P_{\text{blue}}(y)$$

Non-Boolean queries

- We have studied **Boolean queries**:

“Is there both a pink and a blue node?”

$$Q() : \exists x y P_{\text{pink}}(x) \wedge P_{\text{blue}}(y)$$

- In practice, queries often return some **results**:

“Find all pairs of a pink and a blue node?”

$$Q(x, y) : P_{\text{pink}}(x) \wedge P_{\text{blue}}(y)$$

Non-Boolean queries

- We have studied **Boolean queries**:

“Is there both a pink and a blue node?”

$$Q() : \exists x y P_{\text{pink}}(x) \wedge P_{\text{blue}}(y)$$

- In practice, queries often return some **results**:

“Find all pairs of a pink and a blue node?”

$$Q(x, y) : P_{\text{pink}}(x) \wedge P_{\text{blue}}(y)$$

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

Non-Boolean queries

- We have studied **Boolean queries**:

“Is there both a pink and a blue node?”

$$Q() : \exists x y P_{\text{pink}}(x) \wedge P_{\text{blue}}(y)$$

- In practice, queries often return some **results**:

“Find all pairs of a pink and a blue node?”

$$Q(x, y) : P_{\text{pink}}(x) \wedge P_{\text{blue}}(y)$$

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

Circuits as factorized representations of query results

- **Query:** $Q(X_1, \dots, X_n)$ with **free variables** X_1, \dots, X_n

Circuits as factorized representations of query results

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

Circuits as factorized representations of query results

- **Query:** $Q(X_1, \dots, X_n)$ with **free variables** X_1, \dots, X_n
 - **Goal:** find all tuples $\mathbf{a}_1, \dots, \mathbf{a}_n$ such that $Q(\mathbf{a}_1, \dots, \mathbf{a}_n)$ holds
- Add **special facts** to materialize all possible assignments
- e.g., $X_i(\mathbf{a}_j)$ means element \mathbf{a}_j is mapped to variable X_i

Circuits as factorized representations of query results

- **Query:** $Q(X_1, \dots, X_n)$ with **free variables** X_1, \dots, X_n
 - **Goal:** find all tuples $\mathbf{a}_1, \dots, \mathbf{a}_n$ such that $Q(\mathbf{a}_1, \dots, \mathbf{a}_n)$ holds
- Add **special facts** to materialize all possible assignments
- e.g., $X_i(\mathbf{a}_j)$ means element \mathbf{a}_j is mapped to variable X_i
- The **provenance circuit** of Q is now a **factorized representation** which describes all the tuples that make Q true

Circuits as factorized representations of query results

- **Query:** $Q(X_1, \dots, X_n)$ with **free variables** X_1, \dots, X_n
 - **Goal:** find all tuples $\mathbf{a}_1, \dots, \mathbf{a}_n$ such that $Q(\mathbf{a}_1, \dots, \mathbf{a}_n)$ holds
- Add **special facts** to materialize all possible assignments
- e.g., $X_i(\mathbf{a}_j)$ means element \mathbf{a}_j is mapped to variable X_i
- The **provenance circuit** of Q is now a **factorized representation** which describes all the tuples that make Q true

Example query:

$$Q(X_1, X_2) : P_{\text{red}}(x) \wedge P_{\text{blue}}(y)$$

Circuits as factorized representations of query results

- **Query:** $Q(X_1, \dots, X_n)$ with **free variables** X_1, \dots, X_n
 - **Goal:** find all tuples $\mathbf{a}_1, \dots, \mathbf{a}_n$ such that $Q(\mathbf{a}_1, \dots, \mathbf{a}_n)$ holds
- Add **special facts** to materialize all possible assignments
- e.g., $X_i(\mathbf{a}_j)$ means element \mathbf{a}_j is mapped to variable X_i
- The **provenance circuit** of Q is now a **factorized representation** which describes all the tuples that make Q true

Example query:

$$Q(X_1, X_2) : P_{\text{red}}(x) \wedge P_{\text{blue}}(y)$$

Database:



Circuits as factorized representations of query results

- **Query:** $Q(X_1, \dots, X_n)$ with **free variables** X_1, \dots, X_n
 - **Goal:** find all tuples $\mathbf{a}_1, \dots, \mathbf{a}_n$ such that $Q(\mathbf{a}_1, \dots, \mathbf{a}_n)$ holds
- Add **special facts** to materialize all possible assignments
- e.g., $X_i(\mathbf{a}_j)$ means element \mathbf{a}_j is mapped to variable X_i
- The **provenance circuit** of Q is now a **factorized representation** which describes all the tuples that make Q true

Example query:

$$Q(X_1, X_2) : P_{\text{red}}(x) \wedge P_{\text{blue}}(y)$$

Database:



Results:

X_1	X_2
1	3
1	5

Circuits as factorized representations of query results

- **Query:** $Q(X_1, \dots, X_n)$ with **free variables** X_1, \dots, X_n
 - **Goal:** find all tuples $\mathbf{a}_1, \dots, \mathbf{a}_n$ such that $Q(\mathbf{a}_1, \dots, \mathbf{a}_n)$ holds
- Add **special facts** to materialize all possible assignments
- e.g., $X_i(\mathbf{a}_j)$ means element \mathbf{a}_j is mapped to variable X_i
- The **provenance circuit** of Q is now a **factorized representation** which describes all the tuples that make Q true

Example query:

$Q(X_1, X_2) : P_{\text{red}}(x) \wedge P_{\text{blue}}(y)$

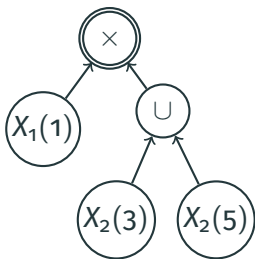
Database:



Results:

X_1	X_2
1	3
1	5

Provenance circuit:



Circuits as factorized representations of query results

- **Query:** $Q(X_1, \dots, X_n)$ with **free variables** X_1, \dots, X_n
 - **Goal:** find all tuples $\mathbf{a}_1, \dots, \mathbf{a}_n$ such that $Q(\mathbf{a}_1, \dots, \mathbf{a}_n)$ holds
- Add **special facts** to materialize all possible assignments
- e.g., $X_i(\mathbf{a}_j)$ means element \mathbf{a}_j is mapped to variable X_i
- The **provenance circuit** of Q is now a **factorized representation** which describes all the tuples that make Q true

Example query:

$Q(X_1, X_2) : P_{\text{red}}(x) \wedge P_{\text{blue}}(y)$

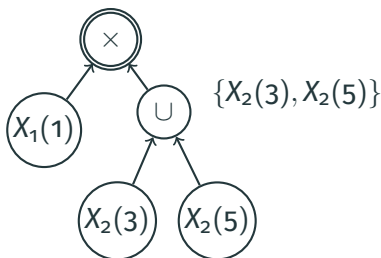
Database:



Results:

X_1	X_2
1	3
1	5

Provenance circuit:



Circuits as factorized representations of query results

- **Query:** $Q(X_1, \dots, X_n)$ with **free variables** X_1, \dots, X_n
 - **Goal:** find all tuples $\mathbf{a}_1, \dots, \mathbf{a}_n$ such that $Q(\mathbf{a}_1, \dots, \mathbf{a}_n)$ holds
- Add **special facts** to materialize all possible assignments
- e.g., $X_i(\mathbf{a}_j)$ means element \mathbf{a}_j is mapped to variable X_i
- The **provenance circuit** of Q is now a **factorized representation** which describes all the tuples that make Q true

Example query:

$Q(X_1, X_2) : P_{\text{red}}(x) \wedge P_{\text{blue}}(y)$

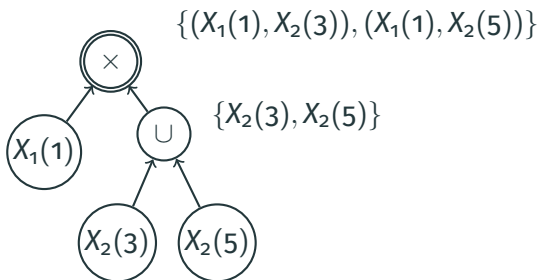
Database:



Results:

X_1	X_2
1	3
1	5

Provenance circuit:



Application: enumerating query results

We can compute a factorized representation of the query results in **linear time** in the data, even if there are **polynomially many** results

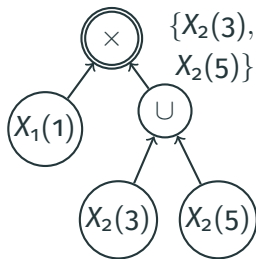
→ Application: **Constant-delay enumeration** of query results

Application: enumerating query results

We can compute a factorized representation of the query results in **linear time** in the data, even if there are **polynomially many** results

→ Application: **Constant-delay enumeration** of query results

$\{(X_1(1), X_2(3)),$
 $(X_1(1), X_2(5))\}$



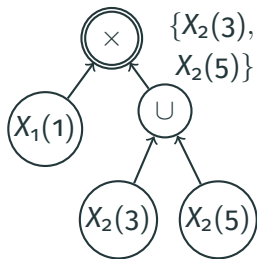
- First, preprocess the circuit in **linear time**
- Then, produce each result in **constant time**

Application: enumerating query results

We can compute a factorized representation of the query results in **linear time** in the data, even if there are **polynomially many** results

→ Application: **Constant-delay enumeration** of query results

$\{(X_1(1), X_2(3)),$
 $(X_1(1), X_2(5))\}$



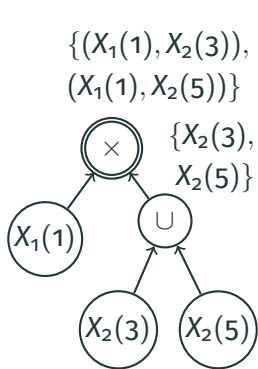
$\{X_2(3),$
 $X_2(5)\}$

- First, preprocess the circuit in **linear time**
- Then, produce each result in **constant time**
- Extends **existing results** on MSO enumeration [Bagan, 2006, Kazana and Segoufin, 2013]

Application: enumerating query results

We can compute a factorized representation of the query results in **linear time** in the data, even if there are **polynomially many** results

→ Application: **Constant-delay enumeration** of query results

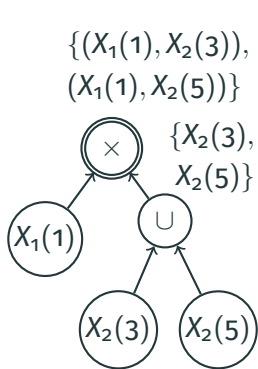


- First, preprocess the circuit in **linear time**
- Then, produce each result in **constant time**
- Extends **existing results** on MSO enumeration [Bagan, 2006, Kazana and Segoufin, 2013]
- Can even be done **tractably in the automaton**

Application: enumerating query results

We can compute a factorized representation of the query results in **linear time** in the data, even if there are **polynomially many** results

→ Application: **Constant-delay enumeration** of query results

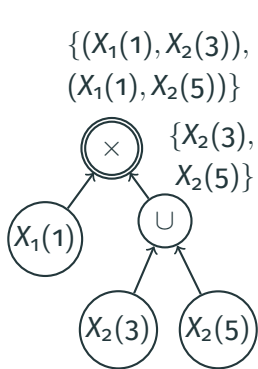


- First, preprocess the circuit in **linear time**
- Then, produce each result in **constant time**
- Extends **existing results** on MSO enumeration [Bagan, 2006, Kazana and Segoufin, 2013]
- Can even be done **tractably in the automaton**
- Applications to **information extraction** (“spanners”) [Amarilli et al., 2019a]

Application: enumerating query results

We can compute a factorized representation of the query results in **linear time** in the data, even if there are **polynomially many** results

→ Application: **Constant-delay enumeration** of query results



- First, preprocess the circuit in **linear time**
- Then, produce each result in **constant time**
- Extends **existing results** on MSO enumeration [Bagan, 2006, Kazana and Segoufin, 2013]
- Can even be done **tractably in the automaton**
- Applications to **information extraction** (“spanners”) [Amarilli et al., 2019a]
- Extensions to support **updates** on the database [Amarilli et al., 2019a, Amarilli et al., 2019b]

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., 2017a, Amarilli et al., 2017b]
 - **Semiring provenance** and explanations [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., 2017a, Amarilli et al., 2017b]
 - **Semiring provenance** and explanations [Green et al., 2007]
- 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., 2017a, Amarilli et al., 2017b]
 - **Semiring provenance** and explanations [Green et al., 2007]
- 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

Thanks for your attention!

References i



Amarilli, A., Bourhis, P., Mengel, S., and Niewerth, M. (2019a).
**Constant-Delay Enumeration for Nondeterministic Document
Spanners.**

In *ICDT*.



Amarilli, A., Bourhis, P., Mengel, S., and Niewerth, M. (2019b).
**Enumeration on Trees with Tractable Combined Complexity and
Efficient Updates.**

Under review.



Amarilli, A., Bourhis, P., Monet, M., and Senellart, P. (2017a).
**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. (2017b).

Conjunctive Queries on Probabilistic Graphs: Combined Complexity.

In *PODS*.



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: newspaper articles (fair use):
 - <http://www.leparisien.fr/transports/circulation-alternee-a-paris-et-en-banlieue-une-panne-de-rer-et-des-bouchons-06-12-2016-6419610.php>
 - <http://www.rtl.fr/actu/societe-faits-divers/paris-le-traffic-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-traffic-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 Dederling, 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://mikaël-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