

Enumeration for MSO Queries on Trees via Circuits

Antoine Amarilli¹ Pierre Bourhis², Louis Jachiet², Stefan Mengel³, Matthias Niewerth⁴

July 14, 2019

¹Télécom Paris, Institut Polytechnique de Paris

²CNRS CRIStAL

³CNRS CRIL

⁴Universität Bayreuth



Enumeration for MSO Queries on Trees via Circuits

Antoine Amarilli¹ Pierre Bourhis², Louis Jachiet², Stefan Mengel³, Matthias Niewerth⁴

July 14, 2019

¹Télécom Paris, Institut Polytechnique de Paris

²CNRS CRIStAL

³CNRS CRIL

⁴Universität Bayreuth





Circuits pour l'énumération de MSO sur des arbres

Antoine Amarilli¹ Pierre Bourhis², Louis Jachiet², Stefan Mengel³, Matthias Niewerth⁴

July 14, 2019

¹Télécom Paris, Institut Polytechnique de Paris

²CNRS CRIStAL

³CNRS CRIL

⁴Universität Bayreuth







Data D: what we know

? Query **Q**: question asked about the data



Data D: what we know

? Query **Q**: question asked about the data

(1) Result: all results of the query **Q** on the data **D**



Data D: what we know

? Query **Q**: question asked about the data

Result: all results of the query **Q** on the data **D**

Measure of efficiency: computational complexity:

- Combined complexity: D and Q are inputs
- Data complexity: Q is fixed, D is the input

Structured data: trees

- Several kinds of data are structured as a **tree** (HTML pages, XML, folder hierarchies...)
 - Important special case: text



Structured data: trees



- Several kinds of data are structured as a **tree** (HTML pages, XML, folder hierarchies...)
 - Important special case: **text**
- On this kind of data, we can use **more efficient** query evaluation techniques

Structured data: trees



- Several kinds of data are structured as a **tree** (HTML pages, XML, folder hierarchies...)
 - Important special case: text
- On this kind of data, we can use **more efficient** query evaluation techniques
- Natural query language: monadic second-order logic (MSO)
 - Very expressive
 - Corresponds to tree automata
 - Data complexity is in linear time







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



Query Q: a Boolean formula 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)



Data: 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)

I Result: YES/NO indicating if the tree au satisfies the query au



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



- Query Q: a Boolean formula 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)

 \mathbf{i} **Result**: YES/NO indicating if the tree *T* satisfies the query *Q*

- Data complexity: linear in T
- Combined complexity: cannot be elementary

Beyond Boolean queries

• Boolean queries only gives YES/NO answers

 $\rightarrow\,$ YES, there is both a pink and a blue node

Beyond Boolean queries

- Boolean queries only gives YES/NO answers
 → YES, there is both a pink and a blue node
- Better idea: non-Boolean queries
 → e.g., Q(x, y): "x is a pink node and y is a blue node"

- Boolean queries only gives YES/NO answers
 → YES, there is both a pink and a blue node
- Better idea: non-Boolean queries
 → e.g., Q(x, y): "x is a pink node and y is a blue node"
- Challenge: how to define complexity for non-Boolean queries?

- Boolean queries only gives YES/NO answers
 → YES, there is both a pink and a blue node
- Better idea: non-Boolean queries
 → e.g., Q(x, y): "x is a pink node and y is a blue node"
- Challenge: how to define complexity for non-Boolean queries?
 - Just writing the output is slow in the worst case

- Boolean queries only gives YES/NO answers
 → YES, there is both a pink and a blue node
- Better idea: non-Boolean queries

 → e.g., Q(x, y): "x is a pink node and y is a blue node"
- Challenge: how to define complexity for non-Boolean queries?
 - Just writing the output is slow in the worst case
 - Easy solution: the **naive algorithm** that tests all pairs

- Boolean queries only gives YES/NO answers
 → YES, there is both a pink and a blue node
- Better idea: non-Boolean queries

 → e.g., Q(x, y): "x is a pink node and y is a blue node"
- Challenge: how to define complexity for non-Boolean queries?
 - Just writing the output is slow in the worst case
 - Easy solution: the naive algorithm that tests all pairs
 - $\rightarrow\,$ We need a **new definition** of complexity

Q how to find patterns

Search

Q how to find patterns

Search

Results 1 - 20 of 10,514

Q how to find patterns

Search

Results 1 - 20 of 10,514

. . .

Q how to find patterns

Search

Results 1 - 20 of 10,514

View (previous 20 | next 20) (20 | 50 | 100 | 250 | 500)

. . .

Q how to find patterns

Search

Results 1 - 20 of 10,514

View (previous 20 | next 20) (20 | 50 | 100 | 250 | 500)

→ Formalization: **enumeration algorithms**



Input


















• We are given: the tree *T*, a non-Boolean query *Q* in MSO \rightarrow Example: *Q*(*x*, *y*) asks for pairs of a blue node *x* and a pink node *y*

- We are given: the tree T, a non-Boolean query Q in MSO \rightarrow Example: Q(x, y) asks for pairs of a blue node x and a pink node y
- We want: enumerate the results of the query

 $\rightarrow\,$ Here, the pairs of a pink node and blue node

- We are given: the tree T, a non-Boolean query Q in MSO \rightarrow Example: Q(x, y) asks for pairs of a blue node x and a pink node y
- We want: enumerate the results of the query
 - $\rightarrow\,$ Here, the pairs of a pink node and blue node

Theorem

For any **fixed** MSO query **Q**, given a tree **T**, we can preprocess **T** in linear time in **T** and then enumerate each result in linear time in the result

- We are given: the tree T, a non-Boolean query Q in MSO \rightarrow Example: Q(x, y) asks for pairs of a blue node x and a pink node y
- We want: enumerate the results of the query
 → Here, the pairs of a pink node and blue node

Theorem [Bagan, 2006, Kazana and Segoufin, 2013]

For any **fixed** MSO query **Q**, given a tree **T**, we can preprocess **T** in linear time in **T** and then enumerate each result in linear time in the result

This was already known, so what's new?

- We are given: the tree T, a non-Boolean query Q in MSO \rightarrow Example: Q(x, y) asks for pairs of a blue node x and a pink node y
- We want: enumerate the results of the query
 → Here, the pairs of a pink node and blue node

Theorem [Bagan, 2006, Kazana and Segoufin, 2013]

For any **fixed** MSO query **Q**, given a tree **T**, we can preprocess **T** in linear time in **T** and then enumerate each result in linear time in the result

This was already known, so what's new?

- Modular proof based on a notion of set circuits
- Tractable in combined complexity for Q given as an automaton
- Efficiently **update** the preprocessing when the tree **changes**

• Building a set circuit: given a tree *T* and automaton *A*, we can build a set circuit *C* that represents the results

- Building a set circuit: given a tree *T* and automaton *A*, we can build a set circuit *C* that represents the results
- Enumeration on set circuits: given a set circuit *C*, we can enumerate efficiently the results that it captures (under some assumptions on *C*)

- Building a set circuit: given a tree *T* and automaton *A*, we can build a set circuit *C* that represents the results
- Enumeration on set circuits: given a set circuit *C*, we can enumerate efficiently the results that it captures (under some assumptions on *C*)
- New stuff:
 - Tractability in the **automaton** and application to text
 - Efficient **updates** of the index

Building a set circuit



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



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



- $P_{\odot}(x)$ means "x is blue"; also $P_{\odot}(x)$, $P_{\odot}(x)$
- $x \rightarrow y$ means "x is the parent 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 ∀
 - $\cdot \exists x \ y \ P_{\bigcirc}(x) \land P_{\bigcirc}(y)$ means "There is both a pink and a blue node"



- $P_{\odot}(x)$ means "x is blue"; also $P_{\odot}(x)$, $P_{\odot}(x)$
- $x \rightarrow y$ means "x is the parent 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 ∀
 - $\cdot \exists x \ y \ P_{\bigcirc}(x) \land P_{\bigcirc}(y)$ 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 an ancestor of y"
 - $\forall x P_{\bigcirc}(x) \Rightarrow \exists y P_{\bigcirc}(y) \land x \rightarrow^{*} y$ means "There is a blue node below every pink node"





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



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



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



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

Theorem [Thatcher and Wright, 1968]

MSO and tree automata have the same expressive power on trees

 $\rightarrow\,$ Given a Boolean MSO query, we can compute a tree automaton that accepts precisely the trees on which the query holds

Theorem [Thatcher and Wright, 1968]

MSO and tree automata have the same expressive power on trees

- $\rightarrow\,$ Given a Boolean MSO query, we can compute a tree automaton that accepts precisely the trees on which the query holds
- ightarrow Complexity (in the query) is generally nonelementary

Theorem [Thatcher and Wright, 1968]

MSO and tree automata have the same expressive power on trees

- $\rightarrow\,$ Given a Boolean MSO query, we can compute a tree automaton that accepts precisely the trees on which the query holds
- ightarrow Complexity (in the query) is generally nonelementary

Corollary

Evaluating a Boolean MSO query on a tree is in linear time in the tree

A small hack for non-Boolean queries

Remember our problem on non-Boolean queries:



A small hack for non-Boolean queries

Remember our problem on non-Boolean queries:

Data: A tree T

"What are the pairs of **Query:** A **non-Boolean** MSO **Q(x)** formula a pink and blue node?"

A small hack for non-Boolean queries

Remember our problem on non-Boolean queries:

Data: A tree T

"What are the pairs of Query: A non-Boolean MSO Q(x) formula a pink and blue node?"

Result: All the **a** such that **Q**(**a**) holds



For technical reasons it will be simpler to:

• Write the choice for variables as colors on the tree



For **technical reasons** it will be simpler to:

• Write the choice for variables as colors on the tree



- Write the choice for variables as colors on the tree
- Replace the non-Boolean query ${\it Q}$ by a Boolean query ${\it Q}'$

Data: A tree T A tree T with more colors to select nodes

Query: A non Boolean MSO Q(x) formula A Boolean MSO formula Q'

"What are the pairs of a pink and blue node?" "Are the two selected nodes pink and blue?"

(**Result**: All the **a** such that **Q**(**a**) holds

- Write the choice for variables as colors on the tree
- Replace the non-Boolean query ${\it Q}$ by a Boolean query ${\it Q}'$

Data: A tree T A tree T with more colors to select nodes

Query: A non Boolean MSO Q(x) formula A Boolean MSO formula Q'

"What are the pairs of a pink and blue node?" "Are the two selected nodes pink and blue?"

Result: All the **a** such that **Q**(**a**) holds

- Write the choice for variables as colors on the tree
- Replace the non-Boolean query ${\it Q}$ by a Boolean query ${\it Q}'$
- Enumerate the ways to color the tree

Data: A tree T A tree T with more colors to select nodes

"What are the pairs of a pink and blue node?" "Are the two selected nodes pink and blue?"

(i) Result: All the**a** $such that <math>Q(\mathbf{a})$ holds All the ways ν to color T such that Q' holds on $\nu(T)$

Query: A non Boolean MSO Q(x) formula A Boolean MSO formula Q'

- Write the choice for variables as colors on the tree
- Replace the non-Boolean query ${\it Q}$ by a Boolean query ${\it Q}'$
- Enumerate the ways to color the tree

Now: Boolean query on a tree where the color of nodes is uncertain



Now: Boolean query on a tree where the color of nodes is uncertain



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


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

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



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

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



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

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



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?"



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



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



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

→ The results that we want to enumerate are all valuations of T that make A accept



• 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
- Output gate:
- Variable gates:
- Constant gates:





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



Every gate *g* captures a set *S*(*g*)



Every gate g captures a set S(g)

• **Variable** gate with label *x*: *S*(*g*) := {{*x*}}



Every gate **g** captures a set **S**(**g**)

- **Variable** gate with label *x*: *S*(*g*) := {{*x*}}
- \top -gates: $S(g) = \{\{\}\}$



Every gate **g** captures a set **S**(**g**)

- **Variable** gate with label *x*: *S*(*g*) := {{*x*}}
- \top -gates: $S(g) = \{\{\}\}$
- \perp -gates: $S(g) = \emptyset$



Every gate **g** captures a set **S**(**g**)

- **Variable** gate with label *x*: *S*(*g*) := {{*x*}}
- \top -gates: $S(g) = \{\{\}\}$
- \perp -gates: $S(g) = \emptyset$
- ×-gate with children g_1, g_2 :

 $S(g) := \{s_1 \cup s_2 \mid s_1 \in S(g_1), s_2 \in S(g_2)\}$



Every gate g captures a set S(g)

- **Variable** gate with label *x*: *S*(*g*) := {{*x*}}
- \top -gates: $S(g) = \{\{\}\}$
- \perp -gates: $S(g) = \emptyset$
- \times -gate with children g_1, g_2 : $S(g) := \{s_1 \cup s_2 \mid s_1 \in S(g_1), s_2 \in S(g_2)\}$
- \cup -gate with children g_1, g_2 : $S(g) := S(g_1) \cup S(g_2)$



Set circuit for automaton A on uncertain tree T:

• Variable gates: nodes of T



Set circuit for automaton A on uncertain tree T:

- Variable gates: nodes of T
- Condition: Let ν be a valuation of T, then A accepts $\nu(T)$ iff the set $S(g_0)$ of the output gate g_0 contains $\{n \in T \mid \nu(n) = 1\}$.



Set circuit for automaton A on uncertain tree T:

- Variable gates: nodes of T
- Condition: Let ν be a valuation of T, then A accepts $\nu(T)$ iff the set $S(g_0)$ of the output gate g_0 contains $\{n \in T \mid \nu(n) = 1\}$.

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



Set circuit for automaton A on uncertain tree T:

- Variable gates: nodes of T
- Condition: Let ν be a valuation of T, then A accepts $\nu(T)$ iff the set $S(g_o)$ of the output gate g_o contains $\{n \in T \mid \nu(n) = 1\}$.

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





Set circuit for automaton A on uncertain tree T:

- Variable gates: nodes of T
- Condition: Let ν be a valuation of T, then A accepts $\nu(T)$ iff the set $S(g_0)$ of the output gate g_0 contains $\{n \in T \mid \nu(n) = 1\}$.

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



Theorem

Theorem

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

- Transitions:
- $\begin{array}{c} \mathsf{P}^\top & \mathsf{P}^\mathsf{P} \\ \mathsf{P}^\top & \mathsf{P}^\mathsf{P}^\mathsf{P} \end{array}$

Theorem

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

- Transitions:
- $\begin{array}{c} \mathsf{P}^\top & \mathsf{P}^\mathsf{P} \\ \mathsf{P}^\top & \mathsf{P}^\mathsf{P}^\mathsf{P} \end{array}$



Theorem

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

- Transitions:
- $\begin{array}{c} \mathsf{P}^\top & \mathsf{P}^\mathsf{P} \\ \mathsf{P}^\top & \mathsf{P}^\mathsf{P}^\mathsf{T} \end{array}$



Theorem

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

- Transitions:



Theorem

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

- Transitions:



Theorem

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

- Transitions:
- $\begin{array}{c} \mathsf{P}^\top & \mathsf{P}^\mathsf{P} \\ \bot & \mathsf{P}^\top \end{array}$



Theorem

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

- Transitions:



To have efficient enumeration, we need restrictions on the set circuit:

d-DNNF set circuit:

- U are all **deterministic**:
- The inputs are **disjoint** (= no set is captured by two inputs)



To have efficient enumeration, we need restrictions on the set circuit:

d-DNNF set circuit:

- ① are all **deterministic**: The inputs are **disjoint** (= no set is captured by two inputs)
 - are all **decomposable**:

The inputs are **independent** (= no variable **x** has a path to two different inputs)



Building d-DNNF set circuits

Theorem

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

- Transitions:
- $\mathsf{P}_{\perp}^{\top} \quad \mathsf{P}_{\perp}^{\mathsf{P}}$



\rightarrow The set circuit of Q is now a factorized representation which describes all the tuples that make Q true

 \rightarrow The set circuit of Q is now a factorized representation which describes all the tuples that make Q true

Example query: $Q(X_1, X_2) : P_{\odot}(x) \land P_{\odot}(y)$
Example query: $Q(X_1, X_2) : P_{\odot}(x) \land P_{\odot}(y)$





Example query: $Q(X_1, X_2) : P_{\odot}(x) \land P_{\odot}(y)$

| Data: | Results: | |
|-------|-----------------|----------------|
| | X ₁ | X ₂ |
| | 1 | 2 |
| | 1 | 3 |

Example query: $Q(X_1, X_2) : P_{\odot}(x) \land P_{\odot}(y)$ Data: 2 3 $\frac{\text{Results:}}{X_1 X_2}$ 1 3 $\frac{X_1 X_2}{1 3}$







Enumeration for set circuits

The situation so far



- We started from our input tree T and query Q
- We have built a **circuit C** describing the results
- We know it is a **d-DNNF**
- We want to **enumerate** these results efficiently

The situation so far



- We started from our input tree ${\it T}$ and query ${\it Q}$
- We have built a **circuit C** describing the results
- We know it is a **d-DNNF**
- We want to enumerate these results efficiently

Theorem

Given a **d-DNNF set circuit C**, we can enumerate its **captured sets** with preprocessing **linear in** |**C**| and delay **linear in each set**

The situation so far



- We started from our input tree ${\it T}$ and query ${\it Q}$
- We have built a **circuit C** describing the results
- We know it is a **d-DNNF**
- We want to **enumerate** these results efficiently

Theorem

Given a **d-DNNF set circuit C**, we can enumerate its **captured sets** with preprocessing **linear in** |**C**| and delay **linear in each set**

- \rightarrow This is a **generic result** (does not talk about MSO or trees)
- $\rightarrow\,$ Any problems whose solutions can be coded as a d-DNNF can be efficiently enumerated via this method

Preprocessing phase:

♥ d-DNNF set circuit

Preprocessing phase:



Preprocessing phase:





Preprocessing phase:





ightarrow E.g., for $S(g) = \{\{x\}, \{x, y\}\}$, enumerate $\{x\}$ and then $\{x, y\}$

ightarrow E.g., for $S(g) = \{\{x\}, \{x, y\}\}$, enumerate $\{x\}$ and then $\{x, y\}$

Base case: variable (X) :

ightarrow E.g., for $S(g) = \{\{x\}, \{x, y\}\}$, enumerate $\{x\}$ and then $\{x, y\}$

Base case: variable (x) : enumerate $\{x\}$ and stop

ightarrow E.g., for $S(g) = \{\{x\}, \{x, y\}\}$, enumerate $\{x\}$ and then $\{x, y\}$

Base case: variable (x) : enumerate $\{x\}$ and stop



Concatenation: enumerate S(g)and then enumerate S(g')

ightarrow E.g., for $S(g) = \{\{x\}, \{x, y\}\}$, enumerate $\{x\}$ and then $\{x, y\}$

Base case: variable (x) : enumerate $\{x\}$ and stop



- Concatenation: enumerate S(g)and then enumerate S(g')
- Determinism: no duplicates

ightarrow E.g., for $S(g) = \{\{x\}, \{x, y\}\}$, enumerate $\{x\}$ and then $\{x, y\}$

Base case: variable (x) : enumerate $\{x\}$ and stop





Concatenation: enumerate S(g)and then enumerate S(g')

Determinism: no duplicates

Lexicographic product: enumerate S(g)and for each result t enumerate S(g')and concatenate t with each result

ightarrow E.g., for $S(g) = \{\{x\}, \{x, y\}\}$, enumerate $\{x\}$ and then $\{x, y\}$

Base case: variable (x) : enumerate $\{x\}$ and stop





Concatenation: enumerate S(g)and then enumerate S(g')

Determinism: no duplicates

Lexicographic product: enumerate S(g)and for each result t enumerate S(g')and concatenate t with each result

Decomposability: no duplicates











• **Problem:** if $S(g) = \emptyset$ we waste time



- **Problem:** if $S(g) = \emptyset$ we waste time
- Solution: in preprocessing
 - compute **bottom-up** if $S(g) = \emptyset$



- **Problem:** if $S(g) = \emptyset$ we waste time
- Solution: in preprocessing
 - compute **bottom-up** if $S(g) = \emptyset$
 - \cdot then get rid of the gate













• **Problem:** if *S*(*g*) contains {} we waste time in chains of ×-gates



- **Problem:** if *S*(*g*) contains {} we waste time in chains of ×-gates
- Solution:



- **Problem:** if *S*(*g*) contains {} we waste time in chains of ×-gates
- Solution:
 - split g between $S(g) \cap \{\{\}\}$ and $S(g) \setminus \{\{\}\}$ (homogenization)


- **Problem:** if *S*(*g*) contains {} we waste time in chains of ×-gates
- Solution:
 - split g between $S(g) \cap \{\{\}\}$ and $S(g) \setminus \{\{\}\}$ (homogenization)
 - remove inputs with $S(g) = \{\{\}\}$ for x-gates



- **Problem:** if *S*(*g*) contains {} we waste time in chains of ×-gates
- Solution:
 - split g between $S(g) \cap \{\{\}\}$ and $S(g) \setminus \{\{\}\}$ (homogenization)
 - remove inputs with $S(g) = \{\{\}\}$ for x-gates



• **Problem:** if *S*(*g*) contains {} we waste time in chains of ×-gates

• Solution:

- split g between $S(g) \cap \{\{\}\}$ and $S(g) \setminus \{\{\}\}$ (homogenization)
- remove inputs with $S(g) = \{\{\}\}$ for x-gates
- collapse ×-chains with fan-in 1



• **Problem:** if *S*(*g*) contains {} we waste time in chains of ×-gates

• Solution:

- split g between $S(g) \cap \{\{\}\}$ and $S(g) \setminus \{\{\}\}$ (homogenization)
- remove inputs with $S(g) = \{\{\}\}$ for x-gates
- collapse ×-chains with fan-in 1



• **Problem:** if *S*(*g*) contains {} we waste time in chains of ×-gates

• Solution:

- split g between $S(g) \cap \{\{\}\}$ and $S(g) \setminus \{\{\}\}$ (homogenization)
- remove inputs with $S(g) = \{\{\}\}$ for x-gates
- collapse ×-chains with fan-in 1
- → Now, traversing a ×-gate ensures that we make progress: it splits the sets non-trivially



• **Problem:** we waste time in ∪-hierarchies to find a **reachable exit** (non-∪ gate)



- **Problem:** we waste time in ∪-hierarchies to find a **reachable exit** (non-∪ gate)
- Solution: compute reachability index



- **Problem:** we waste time in ∪-hierarchies to find a **reachable exit** (non-∪ gate)
- Solution: compute reachability index



- **Problem:** we waste time in ∪-hierarchies to find a **reachable exit** (non-∪ gate)
- Solution: compute reachability index
- Problem: must be done in linear time



- **Problem:** we waste time in ∪-hierarchies to find a **reachable exit** (non-∪ gate)
- Solution: compute reachability index
- Problem: must be done in linear time

• Solution: Determinism ensures we have a multitree (we cannot have the pattern at the right)





- **Problem:** we waste time in ∪-hierarchies to find a **reachable exit** (non-∪ gate)
- Solution: compute reachability index
- Problem: must be done in linear time

- Solution: Determinism ensures we have a multitree (we cannot have the pattern at the right)
- Custom constant-delay reachability index for multitrees





- **Problem:** we waste time in ∪-hierarchies to find a **reachable exit** (non-∪ gate)
- Solution: compute reachability index
- Problem: must be done in linear time

- Solution: Determinism ensures we have a multitree (we cannot have the pattern at the right)
- Custom constant-delay reachability index for multitrees
- For MSO query evaluation: upwards-deterministic circuit so we have a tree: simpler constant-memory index



Summary of results

We have shown:

Theorem

Given a *d-DNNF* set circuit *C*, we can enumerate its captured sets with preprocessing linear in |*C*| and delay linear in each set

Summary of results

We have shown:

Theorem

Given a *d-DNNF* set circuit *C*, we can enumerate its captured sets with preprocessing linear in |*C*| and delay linear in each set

And for any **fixed** MSO query **Q**, given a tree **T**, we can...

- Construct a d-DNNF C representing the results in O(T)
- Apply to C the scheme above

Summary of results

We have shown:

Theorem

Given a **d-DNNF set circuit C**, we can enumerate its **captured sets** with preprocessing **linear in** |**C**| and delay **linear in each set**

And for any **fixed** MSO query **Q**, given a tree **T**, we can...

- **Construct** a d-DNNF **C** representing the results in **O**(**T**)
- Apply to C the scheme above

So we have **re-proved**:

Theorem [Bagan, 2006, Kazana and Segoufin, 2013]

For any **fixed** MSO query **Q**, given a tree **T**, we can preprocess **T** in linear time in **T** and then enumerate each result in linear time in the result

Application to text and combined complexity



Data: a text T

Antoine Amarilli Description Name Antoine Amarilli. Handle: a3nm. Identity Born 1990-02-07. French national. Appearance as of 2017. Auth OpenPGP. OpenId. Bitcoin. Contact Email and XMPP a3nm@a3nm.net Affiliation Associate professor of computer science (office C201-4) in the DIG team of Télécom ParisTech, 46 rue Barrault, F-75634 Paris Cedex 13, France. Studies PhD in computer science awarded by Télécom ParisTech on March 14, 2016. Former student of the École normale supérieure. test@example.com More Résumé Location Other sites Blogging: a3nm.net/blog Git: a3nm.net/git...



Data: a text T

Antoine Amarilli Description Name Antoine Amarilli. Handle: a3nm. Identity Born 1990-02-07. French national. Appearance as of 2017. Auth OpenPGP. OpenId. Bitcoin. Contact Email and XMPP a3nm@a3nm.net Affiliation Associate professor of computer science (office C201-4) in the DIG team of Télécom ParisTech, 46 rue Barrault, F-75634 Paris Cedex 13, France. Studies PhD in computer science avarded by Télécom ParisTech on March 14, 2016. Former student of the École normale supérieure. test@example.com More Résumé Location Other sites Blogging: a3nm.net/blog Git: a3nm.net/git ...



Query: a pattern P given as a regular expression

 $P := \sqcup [a-z0-9.]^* @ [a-z0-9.]^* \sqcup$



Data: a text T

Antoine Amarilli Description Name Antoine Amarilli. Handle: a3nm. Identity Born 1990-02-07. French national. Appearance as of 2017. Auth OpenPGP. OpenId. Bitcoin. Contact Email and XMPP a3nm@a3nm.net Affiliation Associate professor of computer science (office C201-4) in the DIG team of Télécom ParisTech, 46 rue Barrault, F-75634 Paris Cedex 13, France. Studies PhD in computer science avarded by Télécom ParisTech on March 14, 2016. Former student of the École normale supérieure. test@example.com More Résumé Location Other sites Blogging: a3nm.net/blog Git: a3nm.net/git ...

?

Query: a pattern P given as a regular expression

 $P := \Box [a-z0-9.]^* @ [a-z0-9.]^* \Box$

Output: the list of **substrings** of **T** that match **P**:

 $[186,200\rangle,\ [483,500\rangle,\ \ldots$



Data: a text T

Antoine Amarilli Description Name Antoine Amarilli. Handle: a3nm. Identity Born 1990-02-07. French national. Appearance as of 2017. Auth OpenPGP. OpenId. Bitcoin. Contact Email and XMPP a3nm@a3nm.net Affiliation Associate professor of computer science (office C201-4) in the DIG team of Télécom ParisTech, 46 rue Barrault, F-75634 Paris Cedex 13, France. Studies PhD in computer science avarded by Télécom ParisTech on March 14, 2016. Former student of the École normale supérieure. test@example.com More Résumé Location Other sites Blogging: a3nm.net/blog Git: a3nm.net/git ...

Query: a pattern P given as a regular expression

 $P := \sqcup [a-z0-9.]^* @ [a-z0-9.]^* \sqcup$

(1) Output: the list of substrings of *T* that match *P*: [186, 200), [483, 500), ...

Goal:

- be very efficient in T (constant-delay)
- be reasonably efficient in P (polynomial-time)

• A text is just a tree with a simpler shape

- A text is just a tree with a simpler shape
- A regular expression pattern can be expressed in MSO

- A text is just a tree with a simpler shape
- A regular expression pattern can be expressed in MSO
 - $\rightarrow\,$ More generally: regular expressions with <code>variables</code>

 \rightarrow Example: $P := \bullet^* \alpha a^* \beta b^* \gamma \bullet^*$

• Translate to a word automaton (with capture variables)

- A text is just a tree with a simpler shape
- A regular expression pattern can be expressed in MSO
 - $\rightarrow\,$ More generally: regular expressions with variables
 - \rightarrow Example: $P := \bullet^* \alpha a^* \beta b^* \gamma \bullet^*$
- Translate to a word automaton (with capture variables)



- A text is just a tree with a simpler shape
- A regular expression pattern can be expressed in MSO
 - $\rightarrow\,$ More generally: regular expressions with variables
 - \rightarrow Example: $P := \bullet^* \alpha a^* \beta b^* \gamma \bullet^*$
- Translate to a word automaton (with capture variables)



ightarrow The MSO result implies:

Theorem [Florenzano et al., 2018]

We can enumerate all matches of a regular expression pattern on a tree with linear preprocessing and constant delay

- A text is just a tree with a simpler shape
- A regular expression pattern can be expressed in MSO
 - $\rightarrow\,$ More generally: regular expressions with variables
 - \rightarrow Example: $P := \bullet^* \alpha a^* \beta b^* \gamma \bullet^*$
- Translate to a word automaton (with capture variables)



ightarrow The MSO result implies:

Theorem [Florenzano et al., 2018]

We can enumerate all matches of a regular expression pattern on a tree with linear preprocessing and constant delay

→ The resulting set circuit is a binary decision diagram, i.e., each ×-gate has only one input which is not a variable

We have shown linear preprocessing and constant delay in the **data**; but what about the **query**?

• For general MSO queries: nonelementary complexity

We have shown linear preprocessing and constant delay in the **data**; but what about the **query**?

- For general MSO queries: nonelementary complexity
- For regular expressions: exponential (determinization)

We have shown linear preprocessing and constant delay in the **data**; but what about the **query**?

- For general MSO queries: nonelementary complexity
- For regular expressions: exponential (determinization)
- In fact: our methods adapt to nondeterministic automata

We have shown linear preprocessing and constant delay in the **data**; but what about the **query**?

- For general MSO queries: nonelementary complexity
- For regular expressions: exponential (determinization)
- In fact: our methods adapt to nondeterministic automata

Theorem [Amarilli et al., 2019a, Amarilli et al., 2019b]

We can enumerate all matches of a nondeterministic tree automaton on a tree with

- Preprocessing linear in the tree and polynomial in the automaton
- Delay **constant** in the tree and **polynomial** in the automaton

We have shown linear preprocessing and constant delay in the **data**; but what about the **query**?

- For general MSO queries: nonelementary complexity
- For regular expressions: exponential (determinization)
- In fact: our methods adapt to nondeterministic automata

Theorem [Amarilli et al., 2019a, Amarilli et al., 2019b]

We can enumerate all matches of a nondeterministic tree automaton on a tree with

- Preprocessing linear in the tree and polynomial in the automaton
- Delay **constant** in the tree and **polynomial** in the automaton

Corollary

Given a **regular expression pattern P** and text **T**, we can enumerate all matches of **P** on **T** with the complexity above

Implementation (ongoing internship by Rémi Dupré)

- Prototype to find matches of a **regular expression** in a **text**
- https://github.com/remi-dupre/enum-spanner-rs
- Work-in-progress

Implementation (ongoing internship by Rémi Dupré)

- Prototype to find matches of a **regular expression** in a **text**
- https://github.com/remi-dupre/enum-spanner-rs
- Work-in-progress
- Open questions / projects:
 - What about **memory usage**? (we cannot keep the whole index)
 - Output matches in streaming? (problem: duplicates)
 - · Can we enumerate other notions of matches?
 - $\rightarrow~$ factors of <code>maximal/minimal size</code>
 - \rightarrow distinct matching strings
 - \rightarrow etc.

Implementation (ongoing internship by Rémi Dupré)

- Prototype to find matches of a **regular expression** in a **text**
- https://github.com/remi-dupre/enum-spanner-rs
- Work-in-progress
- Open questions / projects:
 - What about **memory usage**? (we cannot keep the whole index)
 - Output matches in streaming? (problem: duplicates)
 - · Can we enumerate other notions of matches?
 - $\rightarrow~$ factors of <code>maximal/minimal size</code>
 - $\rightarrow~$ distinct matching strings
 - \rightarrow etc.
 - Which application domains need this?
 - Are there good **benchmarks**?

Handling updates

Updates



• The input data can be **modified** after the preprocessing


• The input data can be **modified** after the preprocessing



• The input data can be **modified** after the preprocessing



- The input data can be **modified** after the preprocessing
- If this happen, we must rerun the preprocessing from scratch



- The input data can be **modified** after the preprocessing
- If this happen, we must rerun the preprocessing from scratch
- \rightarrow Can we **do better**?

| Work | Data | Preproc. | Delay | Updates |
|-----------------------------|-------|-----------------------|--------------|-----------------------|
| [Bagan, 2006], | trees | <i>O</i> (<i>T</i>) | <i>O</i> (1) | <i>O</i> (<i>T</i>) |
| [Kazana and Segoufin, 2013] | | | | |

| Work | Data | Preproc. | Delay | Updates |
|------------------------------|-------|-----------------------|---------------|-----------------------|
| [Bagan, 2006], | trees | <i>O</i> (<i>T</i>) | <i>O</i> (1) | <i>O</i> (<i>T</i>) |
| [Kazana and Segoufin, 2013] | | | | |
| [Losemann and Martens, 2014] | trees | O(T) | $O(\log^2 T)$ | $O(\log^2 T)$ |

| Work | Data | Preproc. | Delay | Updates |
|------------------------------|-------|-----------------------|---------------|-----------------------|
| [Bagan, 2006], | trees | <i>O</i> (<i>T</i>) | <i>O</i> (1) | <i>O</i> (<i>T</i>) |
| [Kazana and Segoufin, 2013] | | | | |
| [Losemann and Martens, 2014] | trees | O(T) | $O(\log^2 T)$ | $O(\log^2 T)$ |
| [Losemann and Martens, 2014] | text | O(T) | $O(\log T)$ | $O(\log T)$ |

....

-

| Work | Data | Preproc. | Delay | Updates |
|-------------------------------|-------|----------|---------------|---------------|
| [Bagan, 2006], | trees | O(T) | O(1) | O (T) |
| [Kazana and Segoufin, 2013] | | | | |
| [Losemann and Martens, 2014] | trees | O(T) | $O(\log^2 T)$ | $O(\log^2 T)$ |
| [Losemann and Martens, 2014] | text | O(T) | $O(\log T)$ | $O(\log T)$ |
| [Niewerth and Segoufin, 2018] | text | O(T) | O(1) | $O(\log T)$ |

| Work | Data | Preproc. | Delay | Updates |
|-------------------------------|-------|-----------------------|---------------|-----------------------|
| [Bagan, 2006], | trees | <i>O</i> (<i>T</i>) | O(1) | <i>O</i> (<i>T</i>) |
| [Kazana and Segoufin, 2013] | | | | |
| [Losemann and Martens, 2014] | trees | O(T) | $O(\log^2 T)$ | $O(\log^2 T)$ |
| [Losemann and Martens, 2014] | text | O(T) | $O(\log T)$ | $O(\log T)$ |
| [Niewerth and Segoufin, 2018] | text | O(T) | O(1) | $O(\log T)$ |
| [Amarilli et al., 2019b] | trees | O(T) | O(1) | $O(\log T)$ |

Summary and open problems

Summary and conclusion

Theorem

Given a deterministic tree automaton A and a tree T, we can build in $O(|A| \times |T|)$ a d-DNNF set circuit capturing the results of A on T.

Theorem

Given a d-DNNF set circuit **C**, we can enumerate its results with linear preprocessing and delay linear in each result

Summary and conclusion

Theorem

Given a deterministic tree automaton A and a tree T, we can build in $O(|A| \times |T|)$ a d-DNNF set circuit capturing the results of A on T.

Theorem

Given a d-DNNF set circuit **C**, we can enumerate its results with linear preprocessing and delay linear in each result

- If A is nondeterministic, this still works with O(Poly(A))
- If T is **updated**, we can handle the change in $O(\log |T|)$

Theorem

Given a deterministic tree automaton A and a tree T, we can build in $O(|A| \times |T|)$ a d-DNNF set circuit capturing the results of A on T.

Theorem

Given a d-DNNF set circuit **C**, we can enumerate its results with linear preprocessing and delay linear in each result

- If A is nondeterministic, this still works with O(Poly(A))
- If T is **updated**, we can handle the change in $O(\log |T|)$

Open problems:

- Implementation use cases?
- Lower bounds?
- Enumeration with order?

- Memory usage?
- Connection to tuple testing?
- Generic indexes?

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.

In PODS.



MSO queries on Tree Decomposable Structures Are Computable with Linear Delay.

In CSL.

- Florenzano, F., Riveros, C., Ugarte, M., Vansummeren, S., and Vrgoc, D. (2018).

Constant Delay Algorithms for Regular Document Spanners. In *PODS*.

- Kazana, W. and Segoufin, L. (2013).
 Enumeration of Monadic Second-Order Queries on Trees. TOCL, 14(4).
 - Losemann, K. and Martens, W. (2014).
 MSO queries on trees: Enumerating answers under updates. In CSL-LICS.



Niewerth, M. and Segoufin, L. (2018).

Enumeration of MSO queries on strings with constant delay and logarithmic updates.

In *PODS*. To appear.

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.

• 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 nodes**: for each node *n* and variable *x_i*, add a node *n_i* which is colored **red** iff *x_i* is the node *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
- \rightarrow Add **special nodes**: for each node *n* and variable *x_i*, add a node *n_i* which is colored **red** iff *x_i* is the node *n*





- Rewrite the query to a Boolean query which uses the new nodes n_i to read the valuation of x_i
- This can be done in linear time in the input tree

- 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 nodes**: for each node *n* and variable *x_i*, add a node *n_i* which is colored **red** iff *x_i* is the node *n*





- Rewrite the query to a Boolean query which uses the new nodes n_i to read the valuation of x_i
- This can be done in linear time in the input tree
- $\rightarrow\,$ Now, the results are all ways to color the special nodes red and make the Boolean query true

Existential Marked Ancestor Queries

- Given: Tree **t** with some marked nodes
- Query: Does node **v** have a marked ancestor?
- Updates: Mark or unmark a node

Theorem

$$t_{query} \in \Omega\left(rac{\log(n)}{\log(t_{update}\log(n))}
ight)$$

Reduction to Query Enumeration

Fixed Query **Q**: Return all special nodes with a marked ancestor For every marked ancestor query **v**:

- 1. Mark node **v** special
- 2. Enumerate **Q** and return "yes", iff **Q** produces some result
- 3. Mark \mathbf{v} as non-special again

Theorem

$$\max(t_{delay}, t_{update}) \in \Omega\left(\frac{\log(n)}{\log\log(n)}\right)$$

Title slide: Eiffel tower image by Yann Caradec https: //www.flickr.com/photos/la_bretagne_a_paris/35118647963, license CC-BY-SA 2.0.