Worst-case Analysis for Interactive Evaluation of Boolean Provenance

Antoine Amarilli, Yael Amsterdamer





Department of Computer Science Faculty of Exact Sciences Bar-Ilan University

Slides courtesy of Yael Amsterdamer



- Incremental maintenance: maintain languages on words/trees under updates
 - For the language a^*b^* , can you maintain membership on a word under updates?
 - ICALP'21 with Louis Jachiet and Charles Paperman, internship of Jakub Różycki

- Incremental maintenance: maintain languages on words/trees under updates
 - For the language a^*b^* , can you maintain membership on a word under updates?
 - ICALP'21 with Louis Jachiet and Charles Paperman, internship of Jakub Różycki
- Enumeration: enumerate words in a language with bounded edit distance
 - For the language $(a + b)^*$, can you enumerate all words while bounding the edit distance? how about $a^* + b^*$?
 - STACS'23 with Mikaël Monet, project of Mohamed Taha Bayoumi

- Incremental maintenance: maintain languages on words/trees under updates
 - For the language a^*b^* , can you maintain membership on a word under updates?
 - ICALP'21 with Louis Jachiet and Charles Paperman, internship of Jakub Różycki
- Enumeration: enumerate words in a language with bounded edit distance
 - For the language $(a + b)^*$, can you enumerate all words while bounding the edit distance? how about $a^* + b^*$?
 - STACS'23 with Mikaël Monet, project of Mohamed Taha Bayoumi
- Probabilistic query evaluation: for a fixed query *Q*, given a database *I*, count how many subsets of *I* satisfy *Q*
 - ICDT'23 paper, presented on Wednesday

- Incremental maintenance: maintain languages on words/trees under updates
 - For the language a^*b^* , can you maintain membership on a word under updates?
 - ICALP'21 with Louis Jachiet and Charles Paperman, internship of Jakub Różycki
- Enumeration: enumerate words in a language with bounded edit distance
 - For the language $(a + b)^*$, can you enumerate all words while bounding the edit distance? how about $a^* + b^*$?
 - STACS'23 with Mikaël Monet, project of Mohamed Taha Bayoumi
- Probabilistic query evaluation: for a fixed query *Q*, given a database *I*, count how many subsets of *I* satisfy *Q*
 - ICDT'23 paper, presented on Wednesday
- Document spanners for declarative information extraction: ongoing work
 - With Benny Kimelfeld, Sébastien Labbé, Stefan Mengel, on skylines (maximal matches)
 - With Louis Jachiet, Martín Muñoz, on the complexity of core spanners

- Incremental maintenance: maintain languages on words/trees under updates
 - For the language a^*b^* , can you maintain membership on a word under updates?
 - ICALP'21 with Louis Jachiet and Charles Paperman, internship of Jakub Różycki
- Enumeration: enumerate words in a language with bounded edit distance
 - For the language $(a + b)^*$, can you enumerate all words while bounding the edit distance? how about $a^* + b^*$?
 - STACS'23 with Mikaël Monet, project of Mohamed Taha Bayoumi
- Probabilistic query evaluation: for a fixed query *Q*, given a database *I*, count how many subsets of *I* satisfy *Q*
 - ICDT'23 paper, presented on Wednesday
- Document spanners for declarative information extraction: ongoing work
 - With Benny Kimelfeld, Sébastien Labbé, Stefan Mengel, on skylines (maximal matches)
 - With Louis Jachiet, Martín Muñoz, on the complexity of core spanners
- Bounds on query answering with guarded TGDs with Michael Benedikt

- Incremental maintenance: maintain languages on words/trees under updates
 - For the language a^*b^* , can you maintain membership on a word under updates?
 - ICALP'21 with Louis Jachiet and Charles Paperman, internship of Jakub Różycki
- Enumeration: enumerate words in a language with bounded edit distance
 - For the language $(a + b)^*$, can you enumerate all words while bounding the edit distance? how about $a^* + b^*$?
 - STACS'23 with Mikaël Monet, project of Mohamed Taha Bayoumi
- Probabilistic query evaluation: for a fixed query *Q*, given a database *I*, count how many subsets of *I* satisfy *Q*
 - ICDT'23 paper, presented on Wednesday
- Document spanners for declarative information extraction: ongoing work
 - With Benny Kimelfeld, Sébastien Labbé, Stefan Mengel, on skylines (maximal matches)
 - With Louis Jachiet, Martín Muñoz, on the complexity of core spanners
- Bounds on query answering with guarded TGDs with Michael Benedikt
- HDR defense next Tuesday

Worst-case Analysis for Interactive Evaluation of Boolean Provenance

Antoine Amarilli, Yael Amsterdamer





Department of Computer Science Faculty of Exact Sciences Bar-Ilan University

Slides courtesy of Yael Amsterdamer



Boolean Provenance

Acquisitions				Roles				Education			
Acquired	Acquiring	Date		Organization	Role	Member		Alumni	Institute	Year	
A2Bdone microBarg fPharm Optobest	Zazzer Fiffer Fiffer microBarg	7/11/2020 1/5/2017 1/2/2016 8/8/2015	$egin{array}{c} a_0 \ a_1 \ a_2 \ a_3 \end{array}$	A2Bdone A2Bdone A2Bdone microBarg microBarg microBarg	Founder Founding member Founding member Co-founder Co-founder CTO	Usha Koirala Pavel Lebede Nana Alvi Nana Alvi Gao Yawen Amaal Kader	$\begin{array}{c c} r_0 \\ r_1 \\ r_2 \\ r_3 \\ r_4 \\ r_5 \end{array}$	Usha Koirala Pavel Lebedev Nana Alvi Nana Alvi Gao Yawen Amaal Kader	U. Melbourne U. Melbourne U. Sau Paolo U. Melbourne U. Sau Paolo U. Cape Town	2017 2017 2010 2017 2010 2010 2005	$ \begin{array}{c} e_0\\ e_1\\ e_2\\ e_3\\ e_4\\ e_5 \end{array} $
1 SELEC 2 FROM 3 WHERE 4 5	CT DISTI Acqui E a.Acq r.Mem r.Rol	NCT a.A sitions uired = ber = e e LIKE	ACQU AS r. Al	ired, e.In a, Roles Organizat: umni AND a ound%' ANN	nstitute AS r, Educat ion AND a.Date >= 201 D e.YEAR <= y	tion AS e 17.01.01 /ear(a.Da	AND te)		Input data Output re	abase	1
						Acquired	Institute	2			
						A2Bdone A2Bdone microBarg microBarg	U. Melb U. Sau F U. Melb U. Sau F	ourne $(a_0 \wedge r_0 / r_0)$ Paolo $(a_0 \wedge r_2 / r_0)$ ourne $(a_1 \wedge r_3 / r_0)$ Paolo $(a_1 \wedge r_3 / r_0)$	$(e_0) \lor (a_0 \land r_1 \land e_2)$ (e_2) (e_3) (e_2) \lor (a_1 \land r_4 \land e_2)	$(e_1) \lor (e_4)$	$(a_0 \wedge r_2 \wedge e_3)$

Boolean Provenance

Acquisitions				Roles				Education			
Acquired	Acquiring	Date		Organization	Role	Member		Alumni	Institute	Year	
<mark>A2Bdone</mark> microBarg fPharm Optobest	Zazzer Fiffer Fiffer microBarg	7/11/2020 1/5/2017 1/2/2016 8/8/2015	$\begin{array}{c} a_0\\ a_1\\ a_2\\ a_3 \end{array}$	A2Bdone A2Bdone A2Bdone microBarg microBarg microBarg microBarg	Founder Founding member Founding member Co-founder Co-founder CTO	Usha Koirala Pavel Lebedev Nana Alvi Nana Alvi Gao Yawen Amaal Kader	$r_0 \\ r_1 \\ r_2 \\ r_3 \\ r_4 \\ r_5$	Usha Koirala Pavel Lebedev Nana Alvi Nana Alvi Gao Yawen Amaal Kader	U. Melbourne U. Melbourne U. Sau Paolo U. Melbourne U. Sau Paolo U. Cape Town	2017 2017 2010 2017 2010 2010 2005	$e_0 \\ e_1 \\ e_2 \\ e_3 \\ e_4 \\ e_5$
1 SELECT DISTINCT a.Acquired, e.Institute2 FROMAcquisitions AS a, Roles AS r, Education AS e3 WHEREa.Acquired = r.Organization AND4r.Member = e.Alumni AND a.Date >= 2017.01.01 AND5r.Role LIKE '%found%' AND e.YEAR <= year(a.Date)											
						Acquired Ir A2Bdone U	nstitute <mark>J. Melbo</mark> J. Sau P	e ourne $(a_0 \wedge r_0)$ Paolo $(a_0 \wedge r_2)$	$(e_0) \lor (a_0 \land r_1 \land a_0)$	$e_1) \lor ($	$a_0 \wedge r_2 \wedge e_3$

Boolean Provenance

Acquisitions				Roles			Education	Education			
Acquired	Acquiring	Date		Organization	Role	Member		Alumni	Institute	Year	
A2Bdone microBarg fPharm Optobest	Zazzer Fiffer Fiffer microBarg	7/11/2020 1/5/2017 1/2/2016 8/8/2015	$\begin{array}{c} a_0\\ a_1\\ a_2\\ a_3 \end{array}$	A2Bdone A2Bdone A2Bdone microBarg microBarg microBarg	Founder Founding member Founding member Co-founder Co-founder CTO	Usha Koirala Pavel Lebeda Nana Alvi Nana Alvi Gao Yawen Amaal Kaden	$r_0 = r_0$ $r_1 = r_2$ $r_3 = r_4$ $r_5 = r_5$	Usha Koirala Pavel Lebedev Nana Alvi Nana Alvi Gao Yawen Amaal Kader	U. Melbourne U. Melbourne U. Sau Paolo U. Melbourne U. Sau Paolo U. Cape Town	$\begin{array}{ccc} 2017 & e_0 \\ 2017 & e_1 \\ 2010 & e_2 \\ 2017 & e_3 \\ 2010 & e_4 \\ 2005 & e_5 \end{array}$	
1 SELECT DISTINCT a.Acquired, e.Institute 2 FROM Acquisitions AS a, Roles AS r, Education AS e 3 WHERE a.Acquired = r.Organization AND 4 r.Member = e.Alumni AND a.Date >= 2017.01.01 AND 5 r Role IKE '%found%' AND e YEAR <= vear(a Date)											
AcquiredInstituteA2BdoneU. MelbourneA2BdoneU. MelbourneA2BdoneU. Sau PaolomicroBargU. Melbourne $(a_0 \wedge r_2 \wedge e_2)$ $(a_1 \wedge r_3 \wedge e_3)$											

Boolean Provenance: Possible Worlds

Worst-case Analysis for Interactive Evaluation of Boolean Provenance

Acquisitio	ns			$a_0 = Fa$	lse, others=T	rue Education				_		
Acquired	Acquiring	Date							lumni	Institute	Year	
A2Bdone microBarg fPharm Optobest	Zazzer Fiffer Fiffer microBarg	7/11/2020 1/5/2017 1/2/2016 8/8/2015	$ \begin{array}{c} a_0 \\ a_1 \\ a_2 \\ a_3 \end{array} $	A2Bdone A2Bdone A2Bdone microBarg microBarg microBarg	Founder Founding member Founding member Co-founder Co-founder CTO	Usha Koin Pavel Leb Nana Alv Nana Alv Gao Yawe Amaal Ka	rala r_0 pedev r_1 ri r_2 ri r_3 en r_4 ader r_5		Jsha Koirala Pavel Lebedev Jana Alvi Jana Alvi Gao Yawen Amaal Kader	U. Melbourne U. Melbourne U. Sau Paolo U. Melbourne U. Sau Paolo U. Cape Town	$\begin{array}{c c} 2017 & e_0 \\ 2017 & e_1 \\ 2010 & e_2 \\ 2017 & e_3 \\ 2010 & e_4 \\ 2005 & e_5 \end{array}$	
For any truth valuation <i>val</i> : an output tuple <i>t</i> evaluates to true iff it appears in the possible world of <i>val</i>												
						$(a_0 \wedge \cdot)$	$r_0 \wedge e_0$) V	$(a_0 \wedge r_1 \wedge$	$e_1)$ V (a_0 A	$r_2 \wedge e_3$)=False
										Output re	lation	
						Acquired	d Instit	tute				
						A2Bdone	e U. Ma	elbou	$rme (a_0 \wedge r_0) $	$(e_0) \lor (a_0 \land r_1 \land$	$e_1) \vee (a_0$	$\wedge r_2 \wedge e_3)$
						A2Bdone microBa	e U.Sa ro II Mo	u Pac albou	$\frac{10}{10}$ $(a_0 \wedge r_2)$	(e_2)		
						microBa	rg U. Sa	u Pac	blo $(a_1 \wedge r_3)$	(e_3) $(a_1 \wedge r_4 \wedge r_4 \wedge r_4)$	e4)	

Boolean Provenance: Uses

Worst-case Analysis for Interactive Evaluation of Boolean Provenance



Deletion propagation







Consent Management*

Worst-case Analysis for Interactive Evaluation of Boolean Provenance

ante 🦲



Data owners are probed **on a need basis** for fine-grained consent – per tuple

*Managing Consent for Data Access in Shared Databases [ICDE 2021, Drien, A., Amsterdamer]

Consent Management

- We can use the output iff we can derive it from input tuples with consent
- We can choose **which** variables truth values to probe
- Effectiveness depends on the answer and Boolean expressions structure

Acquired	Institute	
A2Bdone	U. Melbourne	$(a_0 \wedge r_0 \wedge e_0) \vee (a_0 \wedge r_1 \wedge e_1) \vee (a_0 \wedge r_2 \wedge e_3)$
A2Bdone	U. Sau Paolo	$(a_0 \wedge r_2 \wedge e_2)$
microBarg	U. Melbourne	$(a_1 \wedge r_3 \wedge e_3)$
microBarg	U. Sau Paolo	$(a_1 \wedge r_3 \wedge e_2) \vee (a_1 \wedge r_4 \wedge e_4)$

Example Evaluation

Worst-case Analysis for Interactive Evaluation of Boolean Provenance



We can use an output tuple iff we can derive it from input tuples with consent

Optimizing the Worst-case Evaluation

Worst-case Analysis for Interactive Evaluation of Boolean Provenance

• We are interested in a "cautious" probing strategy that minimizes the number of probed variables for any valuation



Three Problem Definitions (Intuitive)

Input: a set of Boolean provenance expressions

- **OPT-BDD-DEPTH:** minimize the worst-case number of probes
 - (there is always a trivial strategy that queries all variables in order)
- **DEC-BDD-DEPTH:** decide whether there exists a strategy making at most k probes
- DEC-BDD-EVASIVE: decide whether the expressions are evasive = no strategy is better than the trivial one (making less than n probes over n variables)

Used in Boolean Function Learning



Previous Work

- <u>Expected</u> depth optimization by testing variables of Boolean formulas
 - Interactive Boolean Evaluation, Sequential System Testing, Active Learning, Consent management
- Worst-case BDD Analysis
 - Graph/ String properties
 - Construction based on input-output pairs
- Other metrics





BDDs for Expression Sets

Worst-case Analysis for Interactive Evaluation of Boolean Provenance

 $x \land \neg x$



 $\varphi_0: x \land \neg x \land y$ $\varphi_1:$ False $\varphi_2: y \lor \neg y$



 $\Phi: (a_0 \wedge r_0 \wedge e_0) \vee (a_0 \wedge r_1 \wedge e_1) \vee (a_0 \wedge r_2 \wedge e_3)$ $(a_0 \wedge r_2 \wedge e_2)$ $(a_1 \wedge r_3 \wedge e_3)$ $(a_1 \wedge r_3 \wedge e_2) \vee (a_1 \wedge r_4 \wedge e_4)$ r_2 false true **BDD** for **BDD** for $\Phi_{r_2=True}$ $\Phi_{r_2=\text{False}}$ $\varphi_0 \mapsto \text{True}$ $\varphi_1 \mapsto \text{True}$. . .



- **BDD Depth:** maximal path length from the root to a leaf
- Expression Set Depth: minimal BDD depth
 - Constant expression set ⇔ depth = 0







- Proposition: DEC-BDD-DEPTH is coNP-hard, even if the input Boolean expression is in DNF/CNF and the depth upper bound is k = 0.
- Proof: by reduction from CNF satisfiability / DNF falsifiability.
 A non satisfiable CNF ⇒ constant False ⇒ depth 0

 $x \land \neg x$ false



$$\Phi: (a_0 \wedge r_0 \wedge e_0) \vee (a_0 \wedge r_1 \wedge e_1) \vee (a_0 \wedge r_2 \wedge e_3)$$

$$(a_0 \wedge r_2 \wedge e_2)$$

$$(a_1 \wedge r_3 \wedge e_3)$$

$$(a_1 \wedge r_3 \wedge e_2) \vee (a_1 \wedge r_4 \wedge e_4)$$

Not read-once: variables repeat within/across expressions

Previous work: query classes yielding read-once provenance or compiling provenance to read-once form. E.g., SP queries Worst-case Analysis for Interactive Evaluation of Boolean Provenance

$$\Phi: (a_0 \wedge r_0 \wedge e_0) \vee (a_0 \wedge r_1 \wedge e_1) \vee (a_0 \wedge r_2 \wedge e_3)$$

 $(a_1 \wedge r_3 \wedge e_2) \vee (a_1 \wedge r_4 \wedge e_4)$

Read once: no variable repetitions (in equivalent)

$$\Phi: \underline{a_0 \land ((r_0 \land e_0) \lor (r_1 \land e_1) \lor (r_2 \land e_3))}$$

$$a_1 \wedge ((r_3 \wedge e_2) \vee (r_4 \wedge e_4))$$



- **Proposition:** Sets of read-once of Boolean expressions (without constants), and their equivalents, are **evasive**.
- Proof: by induction
- This result does not hold if variables repeat **across** expressions





- Monotone k-DNF expressions: no negation, every term (conjunction) contains up to k unique variables
- In the paper: we show a 2-way correspondence between k-DNF expressions and SPJU queries
- **Question:** monotone expressions are satisfiable and falsifiable. What is the minimal depth for monotone Boolean expressions?



- Lower bound on depth: maximal term in DNF/clause in CNF
 - Each can be a minimal 0/1 certificate
- **Theorem:** for arbitrarily large n there exists a monotone Boolean expression with a BDD of depth **linear** in this bound
 - Term/clause size is $O(\log n)$ exponentially smaller than "trivial" solution.
 - The BDD is optimal in this case

$$(\psi_{i-1} \wedge u_i) \vee (u_i \wedge v_i) \vee (v_i \wedge \psi'_{i-1})$$



- Recursively define: $\psi_i = (\psi_{i-1} \wedge u_i) \vee (u_i \wedge v_i) \vee (v_i \wedge \psi'_{i-1})$ where u_i, v_i are fresh variables and ψ'_{i-1} is a copy of ψ_{i-1} using fresh variables. Let $\psi_0 = (w_0 \wedge x_0) \vee (x_0 \wedge y_0) \vee (x_0 \wedge y_0)$
- **Observation:** ψ_i cannot be evaluated without probing at least one of u_i , v_i
 - If $u_i = v_i$ we're done by probing both
 - Otherwise, we need to evaluate either ψ_{i-1} or ψ_{i-1}' but not both
- **Observation:** ψ_i includes 2^i copies of ψ_0 and $n = \Theta(2^i)$ variables
- "Bad" algorithm: evaluate all copies of ψ_0 first. Each copy requires 2-4 probes.
- "Good" algorithm: evaluate u_i , v_i first, then if needed proceed to one of the ψ_{i-1} and continue recursively. We query at most $2i + 3 = O(\log n)$



Monotone Acyclic Graph DNF

- When each term is of size 2, terms can be viewed as edges
- When the resulting graph is acyclic, we have the following
- **Theorem:** Given a monotone acyclic graph DNF, DEC-BDD-EVASIVE is in PTIME.
- Proof: We define an **non-evasiveness pattern**, which exists iff the provenance is not evasive



Proof Sketch

Worst-case Analysis for Interactive Evaluation of Boolean Provenance

X

Isolated vertex = non-evasive



Evasive (e.g., if all are true)



= non-evasive

Probe every y_i . If all are false – no need to probe x. Assume w.l.o.g y_0 is true.

$$z_0 \wedge \text{True} = z_0 \text{ absorbs } z_0 \wedge w_0$$

 w_0 is the new root. By recursive argument – it is nonevasive!

The other direction is by induction on the tree structure, showing having no pattern entails that any probe and any answer yields remaining sub-graphs without our pattern

Conclusion and Future Work

- Overview
 - Optimizing the BDD depth for deciding the truth value of Boolean provenance expressions
 - Results for different classes of queries and provenance shapes
 - Many open questions
- Further application domains, further query classes





Thank you!



