

# Distinct Shortest Walk Enumeration for RPQs

Claire DAVID, Nadime FRANCIS and Victor MARSAULT

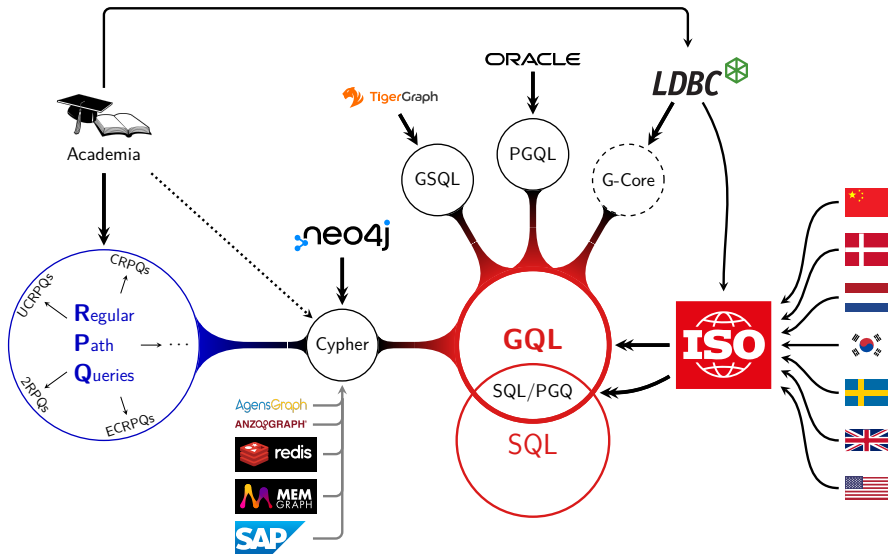
LIGM, Université Gustave Eiffel, CNRS - France

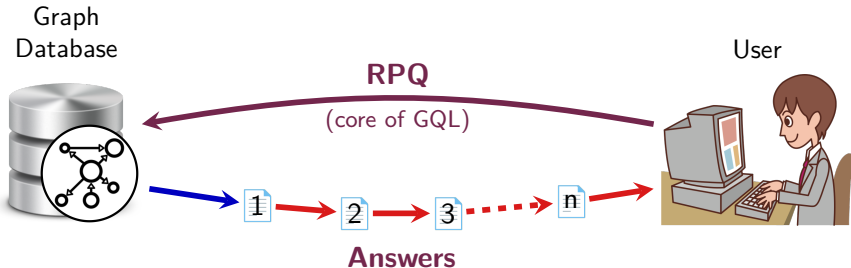


ACM SIGMOD/PODS'24

9-14 June 2024

Santiago, Chile



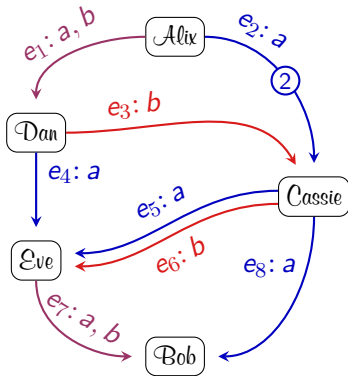


## How to enumerate answers

- Without duplicate
- With good complexity
  - **Preprocessing**: time before first answer
  - **Delay**: time between consecutive answers
  - Memory Usage

- Finite label alphabet:  $\Sigma = \{a, b\}$
- Vertices:  $\{Alix, Bob, \dots, Eve\}$
- Labelled edges:  $\{e_1, \dots, e_7\}$ 
  - $e_1$  carries two labels:  $a$  and  $b$
  - $e_2$  carries one label:  $a$
  - ...

! Graphs are multi-edge and multi-labeled like GQL data model



- **Walk** = consistent sequence of edges  
 $e_1 e_4 e_7$  is a walk     $e_1 e_8$  is not a walk
- **Labels of a walk** = concatenation of its edge labels  
 $e_1 e_4 e_7$  carries four labels:  $aaa$   $baa$   $aab$   $bab$

$Q ::= \mathbf{A}$       Atoms  
 $Q \cdot Q$       Concatenation  
 $Q + Q$       Disjunction  
 $Q^*$       Kleene star

where  $\mathbf{A}$  is a label in the graph.

Ex:  $Q = b^* \cdot a \cdot (a + b)^*$

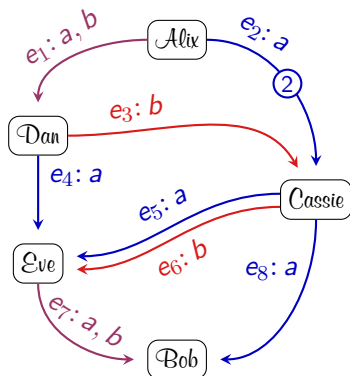
- Match to  $Q =$  walk in  $D$  with a label that conforms to  $Q$

$e_1 e_4 e_7$  matches  $Q$  : it carries  $bab$

$e_1 e_3 e_5 e_7$  matches  $Q$  : it carries  $bbab$

⚠  $e_1 e_4 e_7$  matches  $Q$  three times and  $e_1 e_3 e_5 e_7$  four times

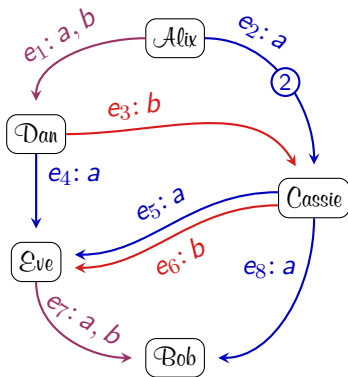
⚠ Matches may have different length



Enumerate all **shortest** walks in  $D$  from  $s$  to  $t$  matching  $Q$  **without duplicates**.

Ex:  $s = \text{Alix}$ ,  $t = \text{Bob}$ ,  $Q = b^* \cdot a \cdot (a + b)^*$

- No matches of length 1 or 2
  - $e_1 e_4 e_7$
  - $e_1 e_3 e_8$
  - $e_2 e_5 e_7$
  - $e_2 e_6 e_7$
  - $e_1 e_3 e_5 e_7$  has length  $> 3$
- } Shortest walks : length 3



⚠ Output each answer only once

## Martens, Trautner 2018

Distinct enumeration of all shortest walks can be done with

- Polynomial time preprocessing
- Polynomial time delay

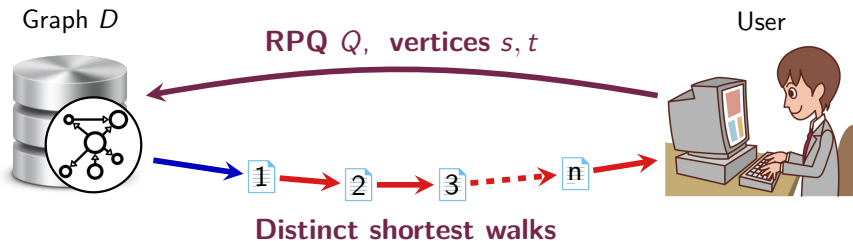
Based on Ackerman, Shallit, 2009

## Folklore

When duplicates cannot occur, enumeration can be done with

- $O(|Q| \times |D|)$  preprocessing
- $O(\lambda)$  delay

where  $\lambda$  is the length of one output.



An algorithm to enumerate distinct shortest walks with

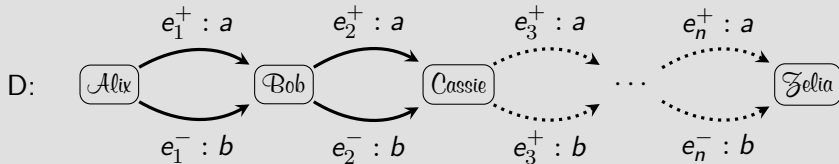
- $O(|Q| \times |D|)$  preprocessing (time before first answer)
- $O(|Q| \times |\lambda|)$  delay (time between two answers)
- $O(|Q| \times |D|)$  memory usage

where  $\lambda$  is the length of one output.



There are  $2^n$  walks from *Alix* to *Zelia* matching  $Q$

Q:  $(a + b)^*$



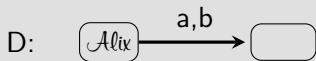
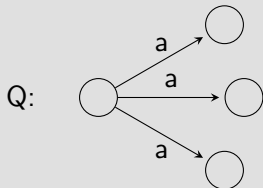
Memory usage must not be linear in the size of output



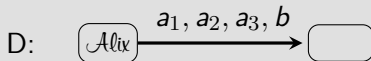
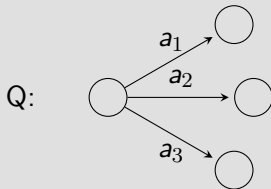


# Where do duplicates comes from?

From the query



From the data



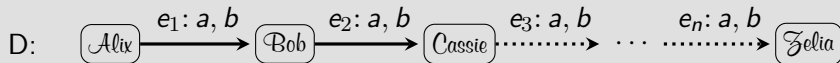
Determinizing the query/automaton does not help





The walk from *Alix* to *Zelia* matching  $Q$  has  $2^n$  duplicates

Q:  $(a + b)^*$

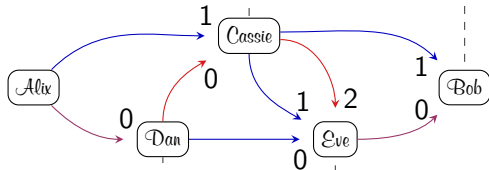


⚠ The algorithm cannot check whether a walk was already output ⚠

$q$	$L_C[q]$	$i$	$B_C[q][i]$	$C_C[q]$
0	1	0	$\square$	$\bullet \leftarrow$
		1	$[0]$	
1	2	0	$[0,1]$	$\bullet \leftarrow$
		1	$\square$	

$q$	$L_B[q]$	$i$	$B_B[q][i]$	$C_B[q]$
0	2	0	$\square$	$\bullet \leftarrow$
		1	$[0]$	
1	3	0	$[1,0,1]$	$\bullet \leftarrow$ $\bullet \leftarrow$
		1	$[1]$	

$q$	$L_A[q]$
0	0
1	$\perp$



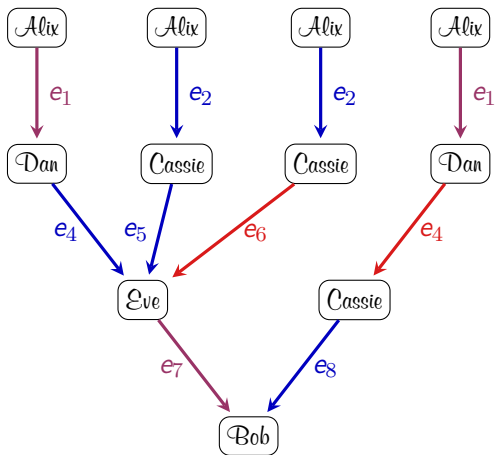
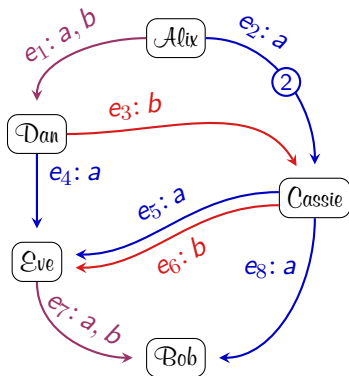
$q$	$L_D[q]$	$i$	$B_D[q][i]$	$C_D[q]$
0	1	0	$[0]$	$\bullet \leftarrow$
		1	$\square$	
1	1	0	$[0]$	$\bullet \leftarrow$
		1	$\square$	

$q$	$L_E[q]$	$i$	$B_E[q][i]$	$C_E[q]$
0	2	0	$[0]$	$\bullet \leftarrow$ $\bullet \leftarrow$
		1	$[0]$	
		2	$\square$	
1	2	0	$[1]$	$\bullet \leftarrow$ $\bullet \leftarrow$
		1	$\square$	
		2	$[0]$	

Inputs:

$$Q = b^* \cdot a \cdot (a + b)^*$$

$$s = \text{Alix} \quad t = \text{Bob}$$

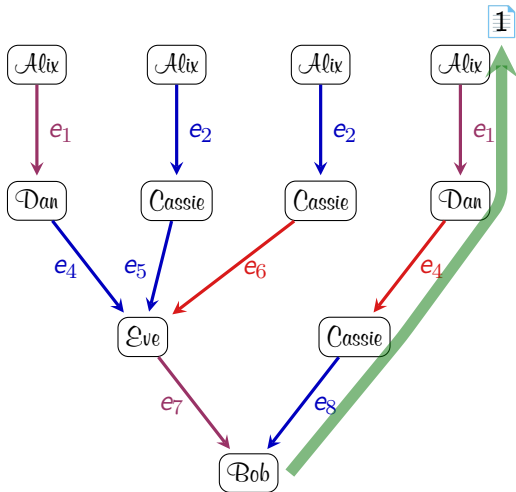
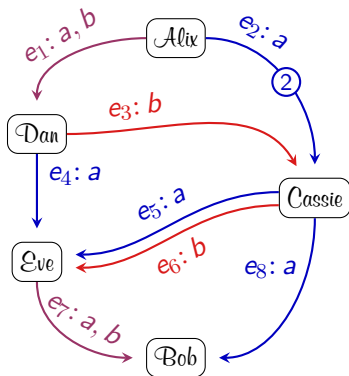


Backward tree of answers  
(up to exponential in  $|D|$  and  $|Q|$ )

Inputs:

$$Q = b^* \cdot a \cdot (a + b)^*$$

$$s = \text{Alix} \quad t = \text{Bob}$$

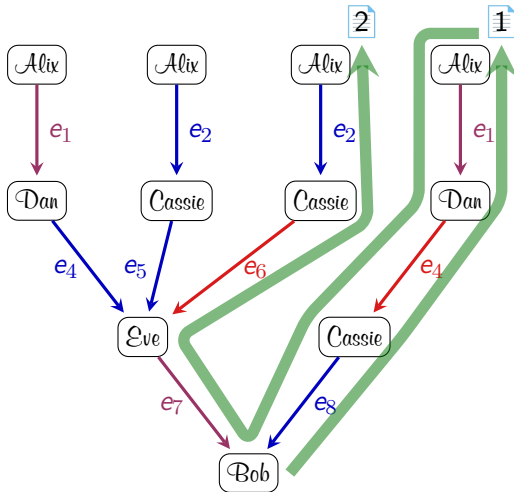
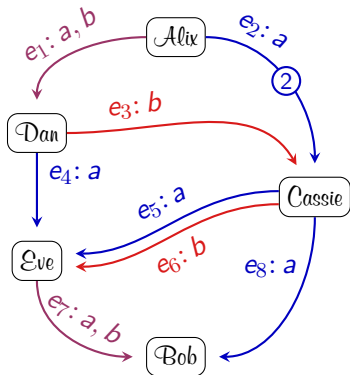


Backward tree of answers  
(up to exponential in  $|D|$  and  $|Q|$ )

Inputs:

$$Q = b^* \cdot a \cdot (a + b)^*$$

$$s = \text{Alix} \quad t = \text{Bob}$$

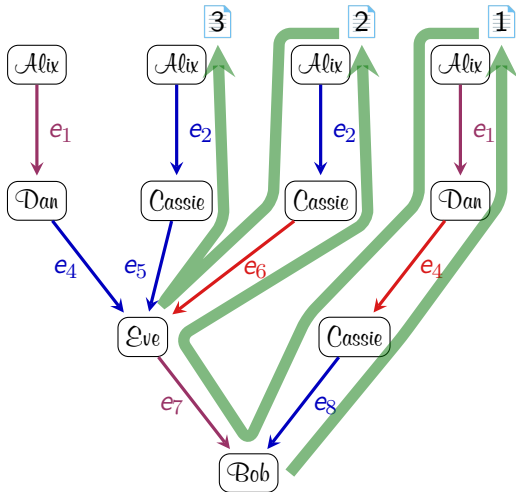
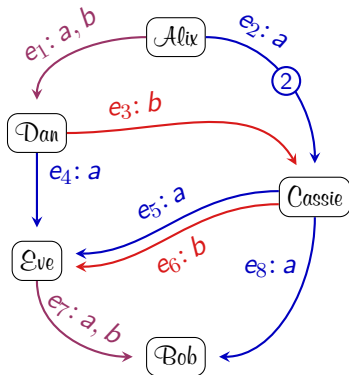


Backward tree of answers  
(up to exponential in  $|D|$  and  $|Q|$ )

Inputs:

$$Q = b^* \cdot a \cdot (a + b)^*$$

$$s = \text{Alix} \quad t = \text{Bob}$$



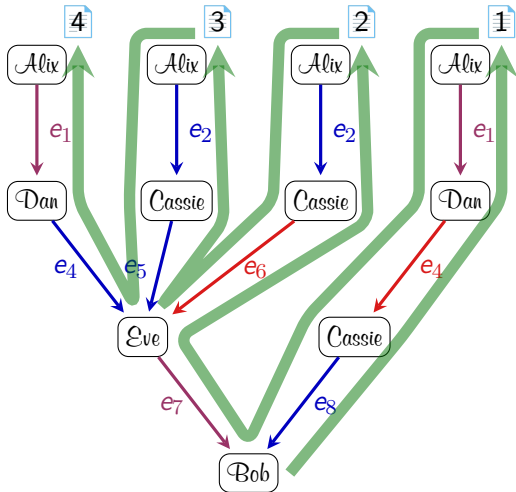
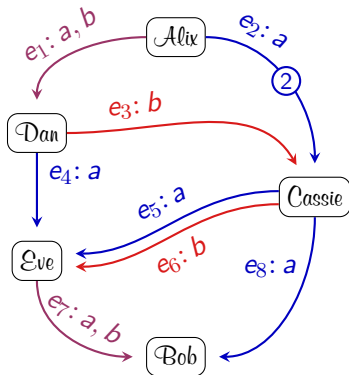
Backward tree of answers  
(up to exponential in  $|D|$  and  $|Q|$ )



Inputs:

$$Q = b^* \cdot a \cdot (a + b)^*$$

$s = \text{Alix}$     $t = \text{Bob}$

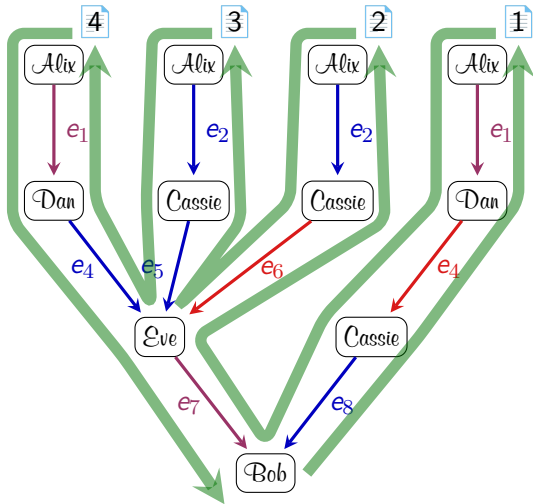
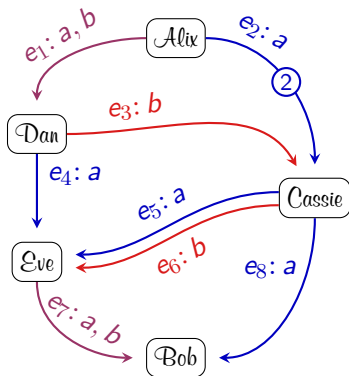


Backward tree of answers  
(up to exponential in  $|D|$  and  $|Q|$ )

Inputs:

$$Q = b^* \cdot a \cdot (a + b)^*$$

$s = \text{Alix}$     $t = \text{Bob}$



Backward tree of answers  
(up to exponential in  $|D|$  and  $|Q|$ )

## Preprocessing phase

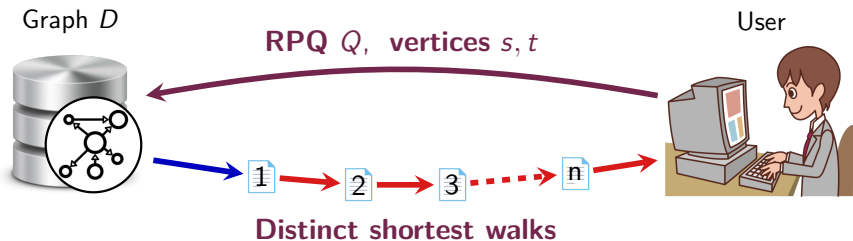
- Annotate  $D$  using a BFS of product graph  $D \times Q$
- Reindex the annotation to encode the backward tree

## Enumeration phase

- DFS of the backward tree, computed on-the-fly
- No dead-ends
- Find next children independently of local topology
- Annotation is almost read-only

## Carefull use of classical data structures

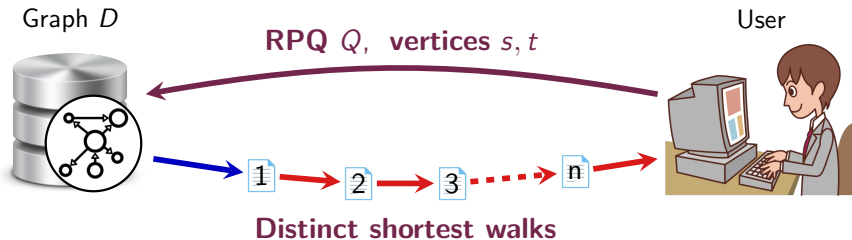
- "Sorting" in linear time via bucket-sort
- Shortcuts in the data structure via LinkedLists
- N-ary "zipper" to merge sorted lists



An algorithm to enumerate distinct shortest walks with

- $O(|Q| \times |D|)$  preprocessing (time before first answer)
- $O(|Q| \times |\lambda|)$  delay (time between two answers)
- $O(|Q| \times |D|)$  memory usage

where  $\lambda$  is the length of one output.



- Reducing overhead when nondeterminism does not occur
- Add other GQL features
- Are data structures well-behaved **in practice**?
- Gather information about nondeterminism in real-life settings