

Distinct Shortest Walk Enumeration for RPQs

Claire DAVID, Nadime FRANCIS and Victor MARSAULT

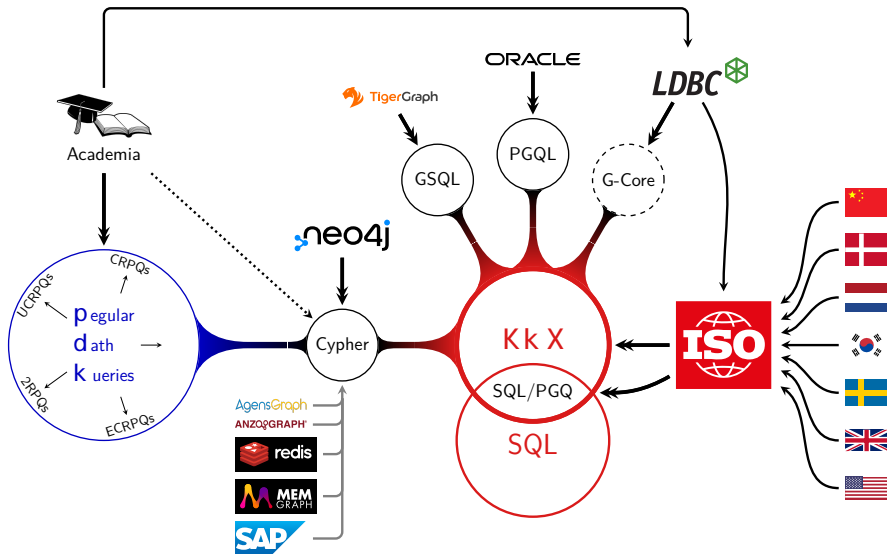
LIGM, Université Gustave Eiffel, CNRS - France



ACM SIGMOD/PODS'24

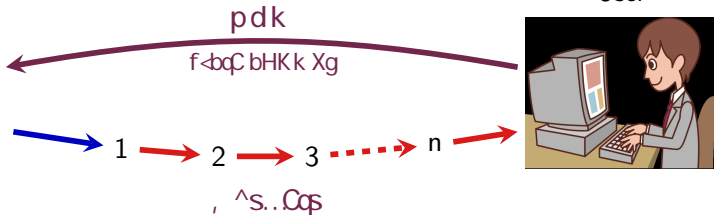
9-14 June 2024

Santiago, Chile



Graph
Database

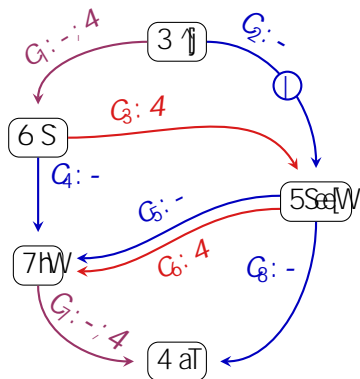
User



How to enumerate answers

- Without duplicate
- With good complexity
 - $d q c e p < C s s S ^ L$: time before first answer
 - $? C Y \% o$: time between consecutive answers
 - Memory Usage

- Finite label alphabet: $\Sigma = \{f, -, 4, g\}$
- Vertices: $\{3 \uparrow, 4 a\bar{T}, \dots, 7 \text{HW}\}$
- Labelled edges: $\{G_1, \dots, G_8\}$
 - G_1 carries two labels: $-$ and 4
 - G_2 carries one label: $-$
 -



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

- $\text{walk} = \text{consistent sequence of edges}$
 $G_1 G_4 G_7$ is a walk $G_1 G_8$ is not a walk
- $X = G_1 G_4 G_7$ carries four labels: $\{-, -, 4, 4\}$

$k ::= ,$ Atoms
 $k k$ Concatenation
 $k + k$ Disjunction
 k^* Kleene star
 where $,$ is a label in the graph.

Ex: $k = 4^* - (- + 4)^*$

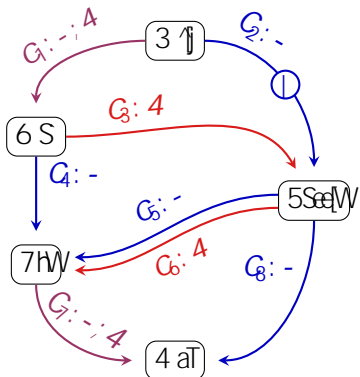
- $[-z < P z b k = \dots W \text{ in } ?$ with a label that conforms to k

$G_1 G_4 G_7$ matches k : it carries $4-4$

$G_1 G_3 G_5 G_7$ matches k : it carries $44-4$

$G_1 G_4 G_7$ matches k three times and $G_1 G_3 G_5 G_7$ four times

Matches may have different length



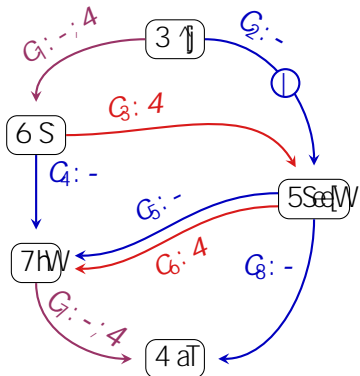
Problem: Distinct Shortest Walk Enumeration



Enumerate all shortest walks in G from s to z matching k

Ex: $s = 3 \uparrow; z = 4 a \bar{T}; k = 4 - (- + 4)$

- No matches of length 1 or 2
 - $G_1 G_4 G$
 - $G_1 G_3 G$
 - $G_2 G_5 G$
 - $G_2 G_6 G$
 - $G_1 G_3 G_5 G$ 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

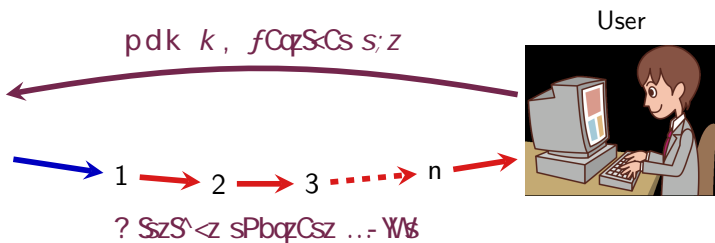
Folklore

When duplicates cannot occur, enumeration can be done with

- a $jkj \ j?j$ preprocessing
- a $()$ delay

where k is the length of one output.

Graph ?



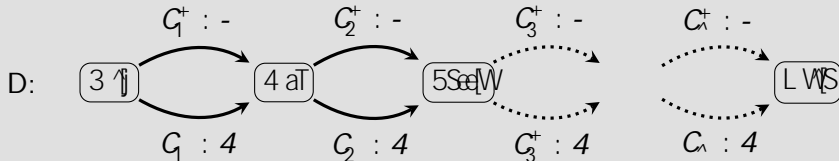
An algorithm to enumerate distinct shortest walks with

- a jkj $j?j$ preprocessing (time before first answer)
- a jkj $j j$ delay (time between two answers)
- a jkj $j?j$ memory usage

where l is the length of one output.

There are 2^k walks from $3 \uparrow$ to $L \downarrow$ matching k

Q: $(- + 4)$

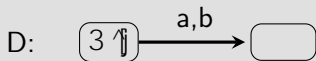
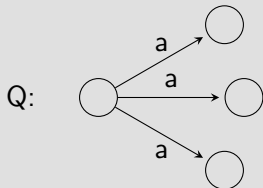


Memory usage must not be linear in the size of output

Where do duplicates comes from?

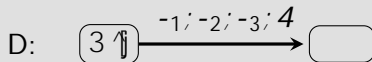
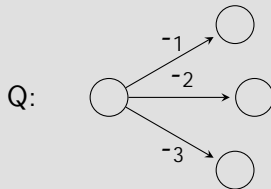


From the query



()

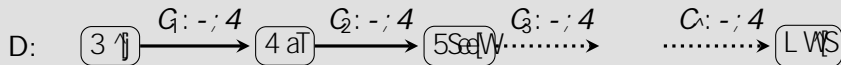
From the data



Determinizing the query/automaton does not help

The walk from $3 \uparrow$ to LWS matching k has 2^k duplicates

Q: $(- + 4)$

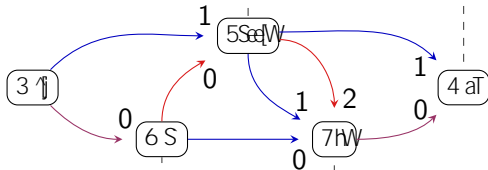


The algorithm cannot check whether a walk was already output

I	$X_5[I]$	S	$3_5[I][\$; 5[I]$
⊆	c	⊆	9
		c	⊆
c		⊆	⊆
		c	9

I	$X_4[I]$	S	$3_4[I][\$; 4[I]$
⊆		⊆	9
		c	⊆
c	{	⊆	⊆
		c	⊆

I	$X_3[I]$
⊆	⊆
c	?



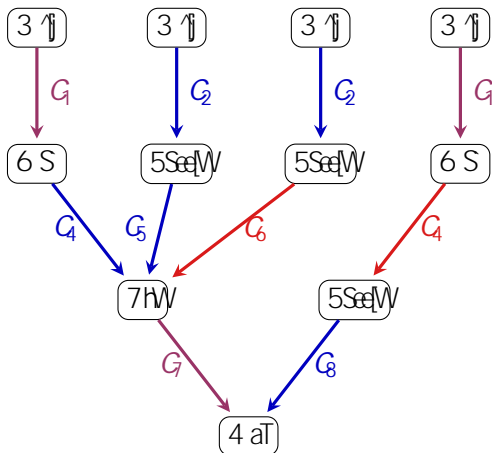
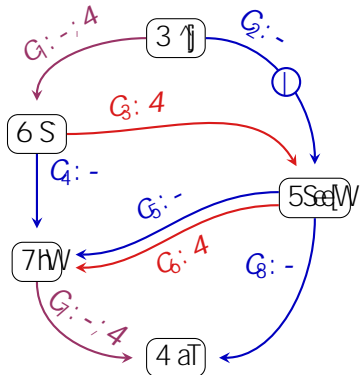
I	$X_6[I]$	S	$3_6[I][\$; 6[I]$
⊆	c	⊆	⊆
c	c	⊆	⊆

I	$X_7[I]$	S	$3_7[I][\$; 7[I]$
⊆		⊆	⊆
		c	⊆
c		⊆	⊆
		c	9
			⊆

Inputs:

$k = 4$ - (- + 4)

$s = 3 \uparrow$ $z = 4 a\bar{T}$



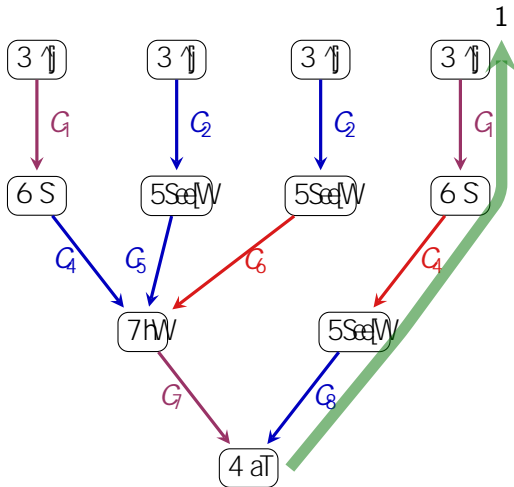
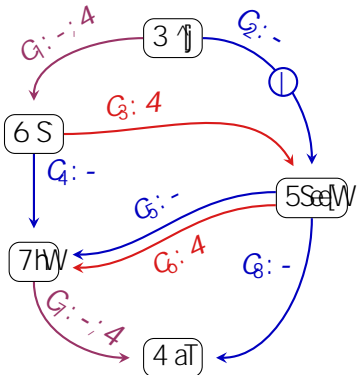
Backward tree of answers

f~e zb Cteb^C^zSYS^ j? j - ^@jk jg

Inputs:

$k = 4$ - (- + 4)

$s = 3 \uparrow$ $z = 4 a\bar{T}$



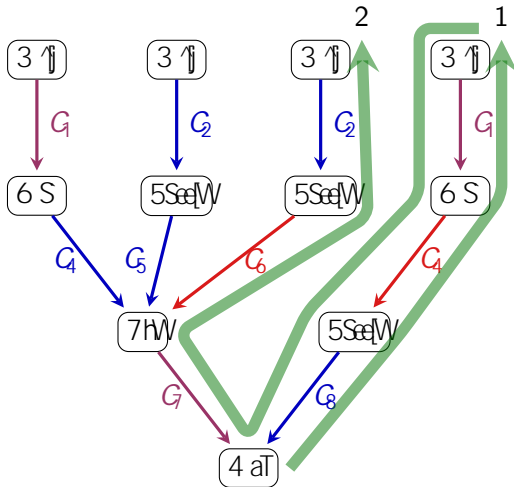
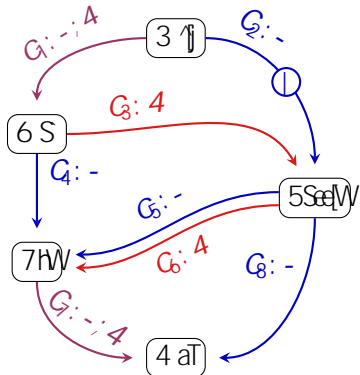
Backward tree of answers

f~e zb Cteb^C^zSYS^ j? j - ^@/jk jg

Inputs:

$k = 4$ - (- + 4)

$s = 3 \uparrow$ $z = 4 a\bar{T}$



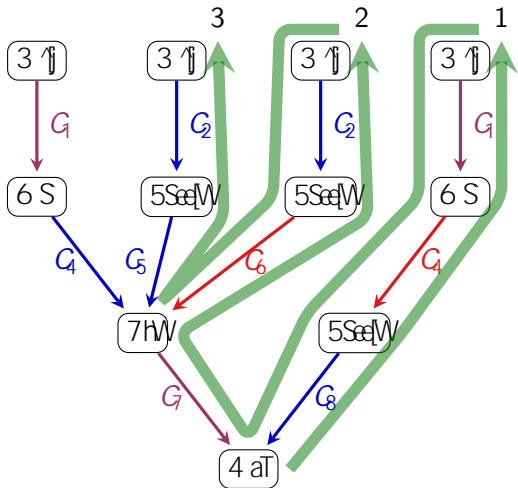
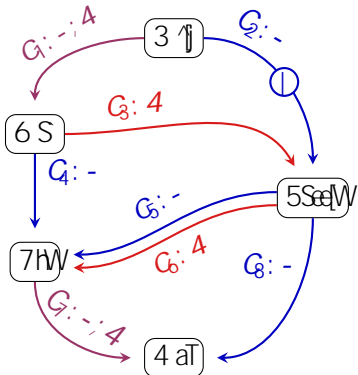
Backward tree of answers

f~e zb Cteb^C^zSYS^ j? j - ^@/jk jg

Inputs:

$k = 4$ - (- + 4)

$s = 3 \uparrow$ $z = 4 a\bar{T}$



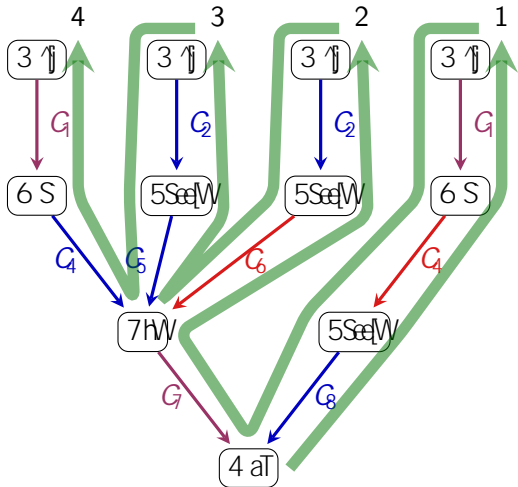
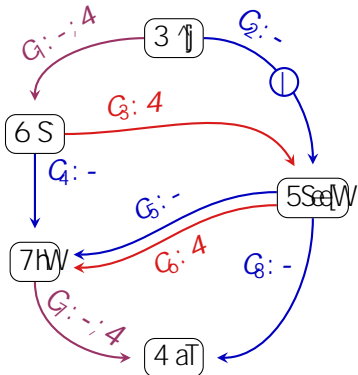
Backward tree of answers

f~e zb Cteb^C^zSYS^ j? j - ^@jk jg

Inputs:

$k = 4$ - (- + 4)

$s = 3 \uparrow$ $z = 4 a\bar{T}$



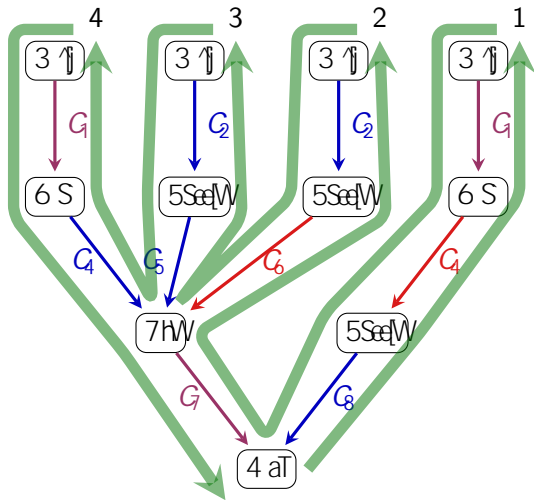
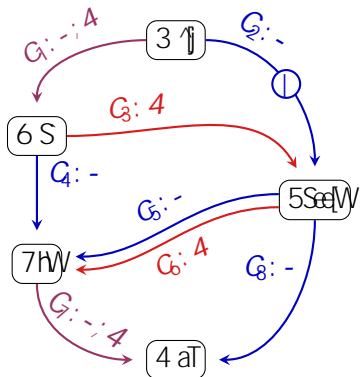
Backward tree of answers

f-e zb Cteb^C^zSYS^ j? j - ^@/jk jg

Inputs:

$k = 4$ - (- + 4)

$s = 3 \uparrow$ $z = 4 a\bar{T}$



Backward tree of answers

f-e zb Cteb^C^zSYS^ j? j - ^@/jk jg

Preprocessing phase

- Annotate τ using a BFS of product graph $\tau \times k$
- Reindex the annotation to encode the backward tree

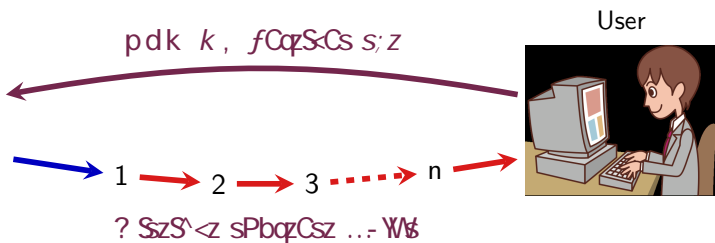
Enumeration phase

- 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

Graph ?

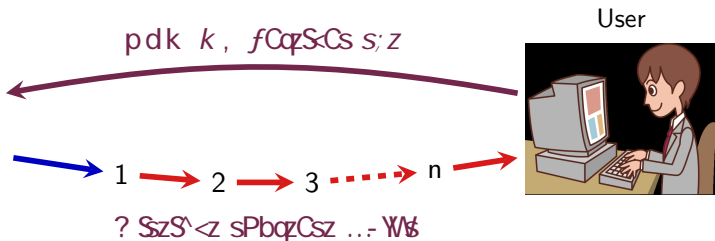


An algorithm to enumerate distinct shortest walks with

- a jkj $j?j$ preprocessing (time before first answer)
- a jkj $j j$ delay (time between two answers)
- a jkj $j?j$ memory usage

where k is the length of one output.

Graph ?



- Reducing overhead when nondeterminism does not occur
- Add other GQL features
- Are data structures well-behaved $S^{\wedge} eq \langle zS:C?$
- Gather information about nondeterminism in real-life settings