

Run-based semantics for RPQs

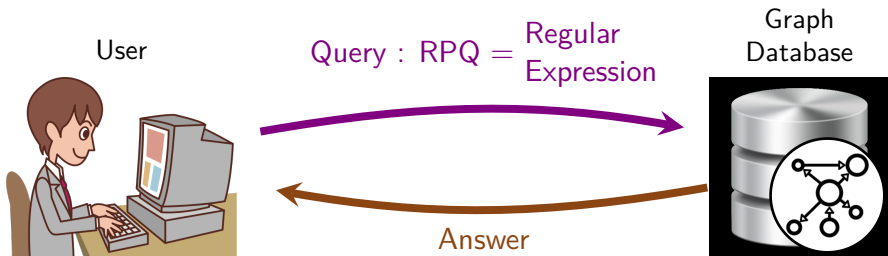
Victor MARSAULT*

joint work with Claire DAVID* and Nadime FRANCIS*

* Université Gustave-Eiffel, CNRS, LIGM

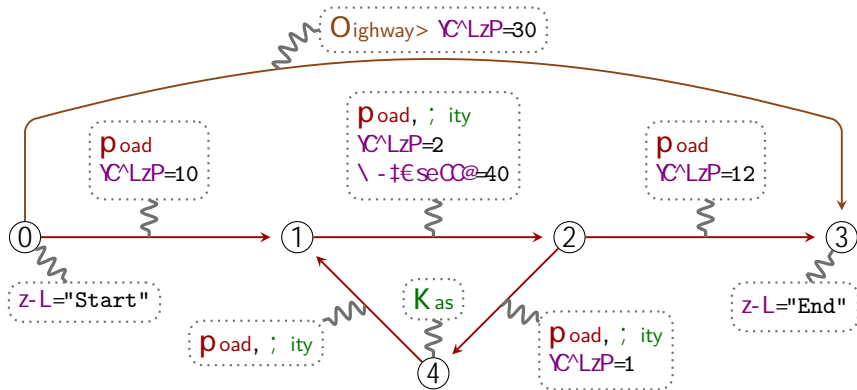
Verigraph Meeting, 2024-04-30, Grenoble, France

(mostly Claire slides used for KR'23)



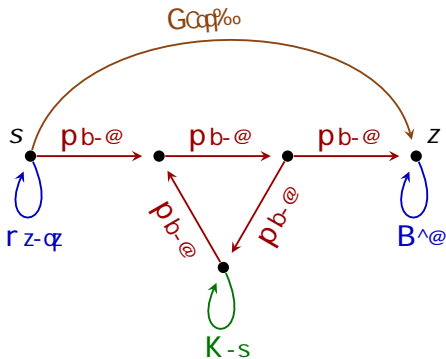
What is a good semantics for RPQ's over graph DB's ?

- Meaningful answers
- Good complexity

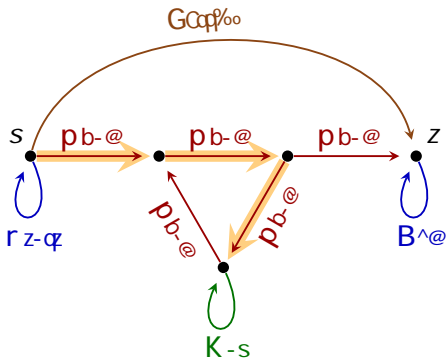


RDF is another data model used in practice...

- Finite label alphabet:
 $= \{r; p; G; K; B\}$
- Vertices
- Edges labelled over



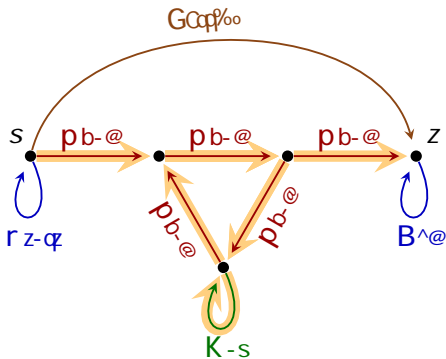
- Finite label alphabet:
 $= \{r; p; G; K; B\}$
- Vertices
- Edges labelled over



Terminology: Walk

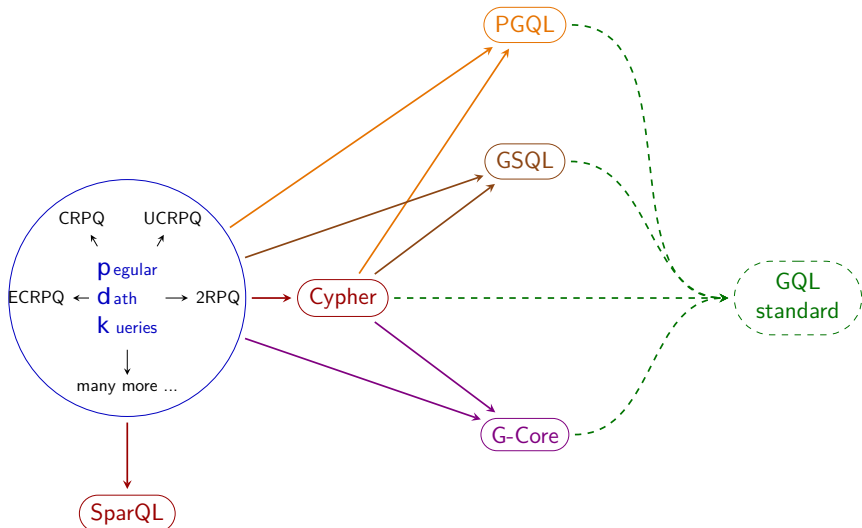
- Consistent sequence of edges
- Repetitions are allowed

- Finite label alphabet:
 $= \{r; p; G; K; B\}$
- Vertices
- Edges labelled over



Terminology: Walk

- Consistent sequence of edges
- Repetitions are allowed



$k ::= \mathbf{A}$ Atoms
 kk Concatenation
 $k + k$ Disjunction
 k^* Kleene star
 where \mathbf{A} is a label in the graph.

Label of a walk

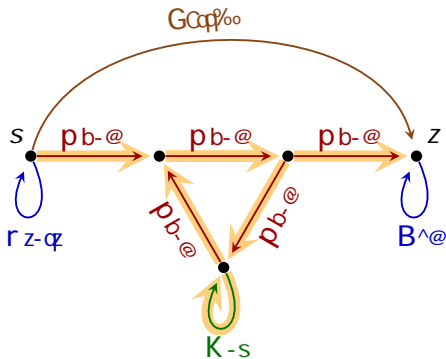
Concatenation of labels of edges

Ex : $pppKppp$

Definition: match for k

A walk ... such that the label of ... matches k .

Ex: match for p^*Kp^*

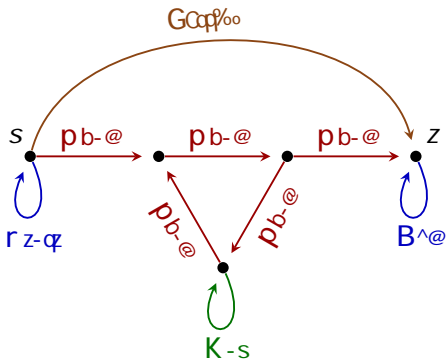


Our two running RPQs

$GS^{\setminus} @ - \dots \% bdp \setminus s z b z$

$$k_1 = r (p + G)^* B$$

Which walks match k_1 ?

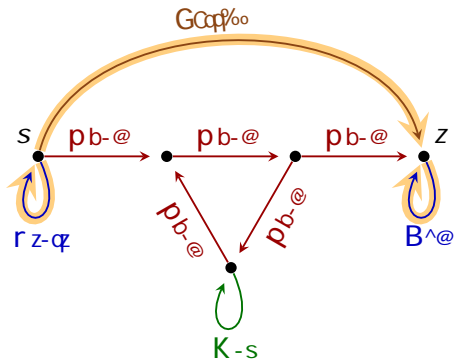


$G \cup \{r, p, z\}$

$$k_1 = r (p + G)^* B$$

Which walks match k_1 ?

- The ferry

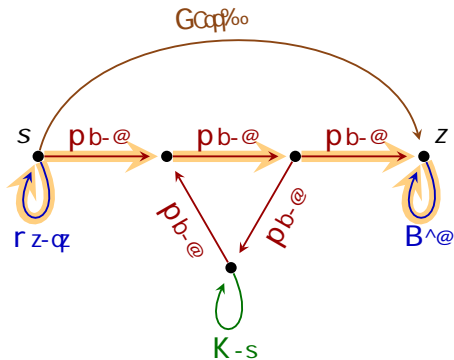


$GS^{\setminus} @ - \dots \% b p \setminus s z b z$

$$k_1 = r (p + G)^* B$$

Which walks match k_1 ?

- The ferry
- The straight road

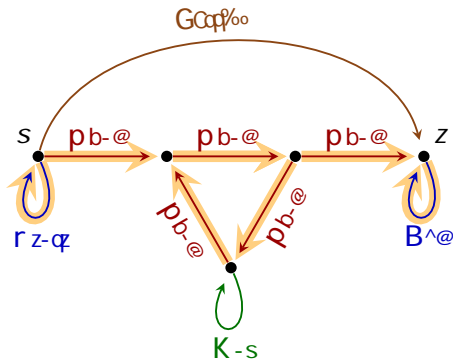


$G S^{\setminus} @ - \dots \% b p \setminus s z b z$

$$k_1 = r (p + G)^* B$$

Which walks match k_1 ?

- The ferry
- The straight road
- Walks with some circuit laps



$G S^* \dots$

$$k_1 = r (p + G)^* B$$

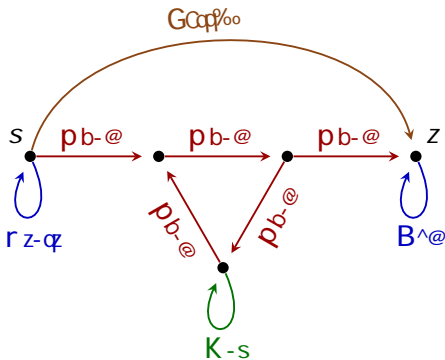
Which walks match k_1 ?

- The ferry
- The straight road
- Walks with some circuit laps

$i i i \dots$

$$k_2 = r (p + G)^* K (p + G)^* B$$

Which walks match k_2 ?



$GS^{\setminus} @ - \dots \% b p \setminus s z b z$

$$k_1 = r (p + G)^* B$$

Which walks match k_1 ?

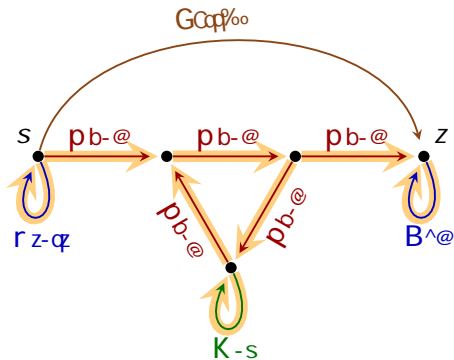
- The ferry
- The straight road
- Walks with some circuit laps

$iii.. \mathcal{S} P \setminus - \wedge @ z b q \% b - s s z b e$

$$k_2 = r (p + G)^* K (p + G)^* B$$

Which walks match k_2 ?

- Walks with some circuit laps



$GS^* \dots$

$$k_1 = r (p + G)^* B$$

Which walks match k_1 ?

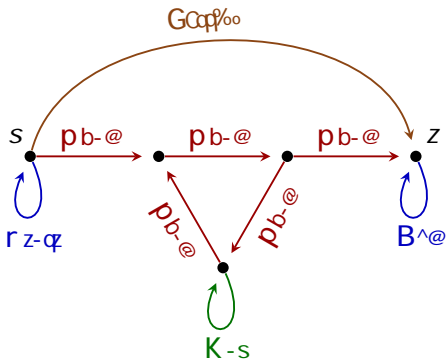
- The ferry
- The straight road
- Walks with some circuit laps

$iii..SP \dots$

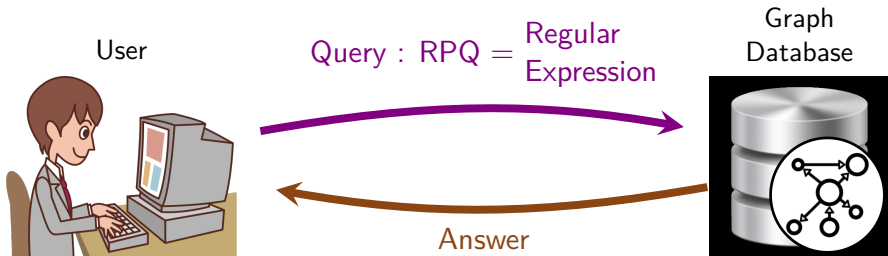
$$k_2 = r (p + G)^* K (p + G)^* B$$

Which walks match k_2 ?

- Walks with some circuit laps



) Infinitely many matches



What is a good semantics for RPQ's over graph DB's ?

- Meaningful answer
- Good complexity

⚠ Infinitely many matches but users expect a finite answer ⚠

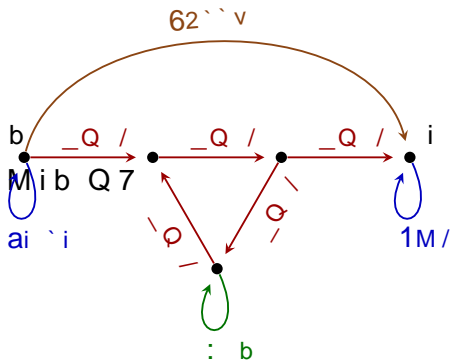
> QKQKQ`T?BbK b2K MiB+b



J BM i?2Q`2iB+ H b2K MiB+b- aS _ZG U_.6V

.2}MBiBQM

- $_2im`Mbi?2 2M/TQB$
K i+?BM; r HFb



$$Z_1 = a \underline{U} Y \theta \quad 1$$

$$Z_2 = a \underline{U} Y \theta \quad : \underline{U} Y \theta \quad 1$$

- "Qi? `2im`M QM T B` ,

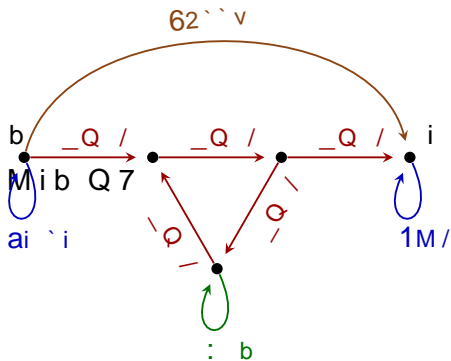
> QKQKQ`T?BbK b2K MiB+b



J BM i?2Q`2iB+ H b2K MiB+b- aS _ZG U_.6V

.2}MBiBQM

- $_2im`Mb i?2 2M/TQB$
K i+?BM; r HFb



$$Z_1 = a \underline{U} Y \theta \quad 1$$

$$Z_2 = a \underline{U} Y \theta \quad : \underline{U} Y \theta \quad 1$$

- "Qi? `2im`M QM T B`

- $q2HH ; `QmM/2/ i?2Q$
- $1\{+B2Mi H;Q`Bi?Kb$
- $_2im`Mb HBiiH2 BM7$

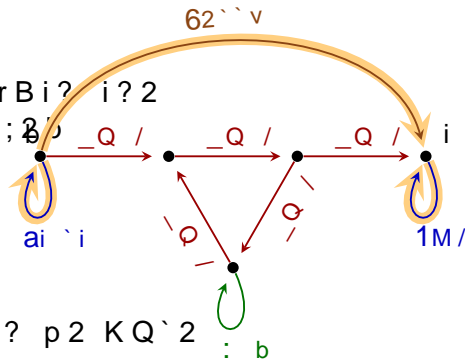
a?Q`i2bi@r HF b2K MiB+b



S:ZG UP` +H2V- :aZG UhB;2`:` T?V- :@*Q`2 (M;H2b 2i

.2}MbiBQM

- $_2im`Mb i?2 r HF rBi? i?2$
 $H2 bi MmK\#2` Q7 2/;2/$



Z₁ = a U_Y 6 1

- $_2im`Mb i?2 72` `v$
- $q HFb mbBM; `Q /b ? p2 KQ`2$
 $2/;2b$

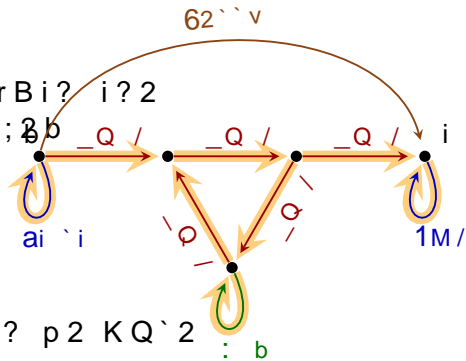
a?Q`i2bi@r HF b2K MiB+b



S:ZG UP` +H2V- :aZG UhB;2`:` T?V- :@*Q`2 (M;H2b 2i

.2}MBiBQM

- $_2im`Mb i?2 r HF rBi? i?2$
 $H2 bi MmK\#2` Q7 2/;2b$



$Z_1 = a \underline{U} Y \theta 1$

- $_2im`Mb i?2 72` `v$
- $q HFb mbBM; `Q /b ? p2 KQ`2$
 $2/;2b$

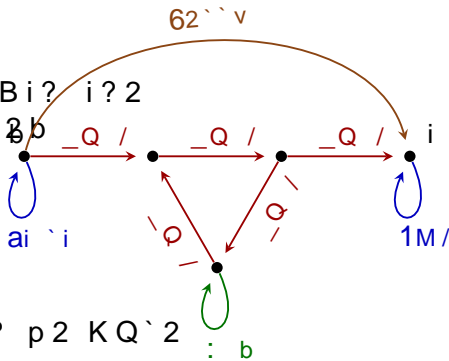
$Z_2 = a \underline{U} Y \theta : \underline{U} Y \theta 1$

- $_2im`Mb i?2 r HF rBi? QM2$
 $+B`+mBi H T$



.2}M BiBQM

- $_2im`Mb i?2 r HF rBi? i?2$
 $H2 bi MmK\#2` Q7 2/;2b$



$Z_1 = a \underline{U} \underline{Y} \theta 1$

- $_2im`Mb i?2 72` `v$
- $q HFb mbBM; `Q /b ? p2 KQ`2$
 $2/;2b$

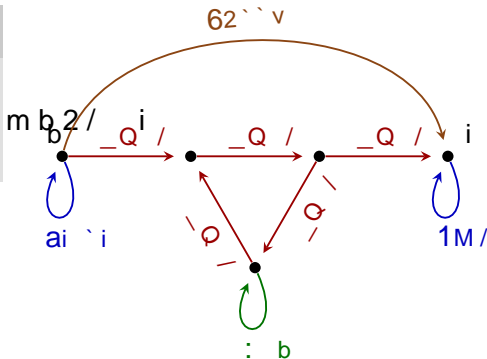
$Z_2 = a \underline{U} \underline{Y} \theta : \underline{U} \underline{Y} \theta 1$

- $_2im`Mb i?2 r HF rBi$
 $+B`+mBi H T$

- $1\{+B2Mi H;Q`Bi?Kb$
- $`\#Bi` `v +?QB+2 Q7$
- $LQ p2`iB+ H TQbi@T$

.2}MBiBQM

- $_2im`Mbr HFb$
- $1 + ? 2 / ; 2 + M \# 2$
KQbi QM+2



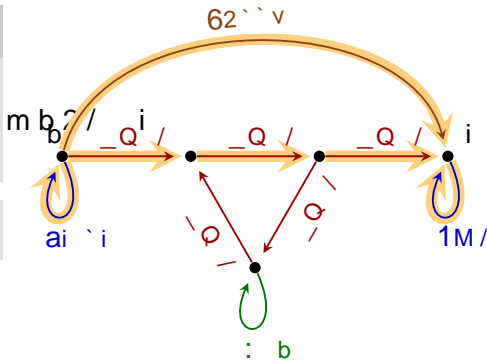
*vT?2` UL2Q9DV- :ZG

.2}MBiBQM

- $_2im`Mb r HFb$
- $1 + ? 2 / ; 2 + M \# 2$
KQbi QM+2

$$Z_1 = a \underline{U} Y 6) 1$$

- $Z_1`2im`Mb$
 - $i?2 72``v$
 - $i?2 bi` B;?i`Q /$



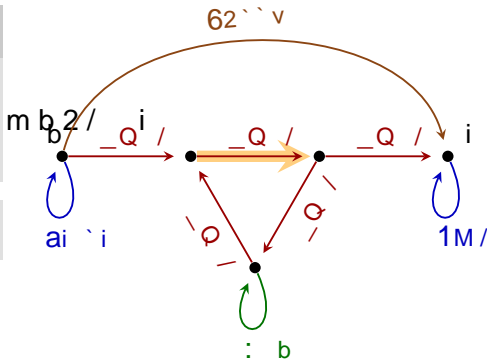
*vT?2` UL2Q9DV- :ZG

.2}MBiBQM

- $_2im`Mb r HFb$
- $1 + ? 2 / ; 2 + M \# 2$
KQbi QM+2

$$Z_1 = a \underline{U} Y 6 1$$

- $Z_1`2im`Mb$
 - $i?2 72``v$
 - $i?2 bi` B ; ?i` Q /$
- $q HFb rBi? +B`+mBi H Tb$
) `2T2 ii?2 KB//H2 2/;2



*vT?2` UL2Q9DV- :ZG

.2}MBiBQM

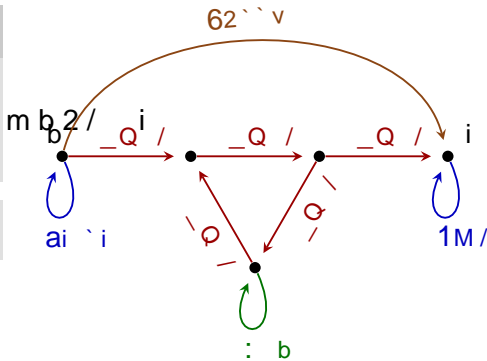
- $_2im`Mbr HFb$
- $1 +? 2/;2 + M \#2$
KQbi QM+2

$$Z_1 = a \underline{U} Y \theta 1$$

- $Z_1`2im`Mb$
 - $i?2 72``v$
 - $i?2 bi`B;?i`Q /$
- $q HFb rBi? +B`+mBi H Tb$
) `2T2 ii?2 KB//H2 2/;2

$$Z_2 = a \underline{U} Y \theta : \underline{U} Y \theta 1$$

- $Z_2`2im`Mb MQ`2bmHib 55$





*vT?2` UL2Q9DV- :ZG

.2}MBiBQM

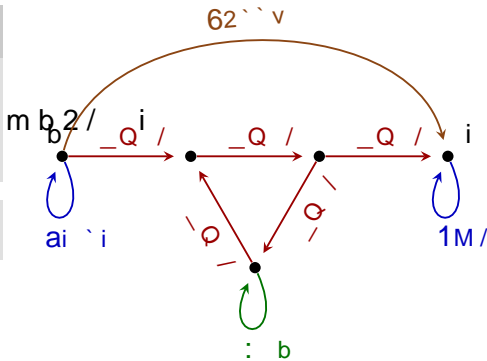
- $_2im`Mb r HFb$
- $1 + ? 2 / ; 2 + M \# 2$
KQbi QM+2

$$Z_1 = a \underline{U} Y \theta 1$$

- $Z_1`2im`Mb$
 - $i?2 72``v$
 - $i?2 bi`B;?i`Q /$
- $q HFb rBi? +B`+mBi$

$$Z_2 = a \underline{U} Y \theta : \underline{U} Y \theta 1$$

- $Z_2`2im`Mb MQ`2bmf$



- $1M \# H2b TQbi@T`Q+$
- $S`Q\#H2Kb`2 mMi`+$
- $.Bb+`/b K2 MBM;7m$
K i+?BM; r HFb

q? i B b ; Q Q / b 2 K M i B + b 7 Q ` _ S Z R

l b 2 `

Z m 2 ` v , - S z 2 ; m H ` : ` T ?
T t T ` 2 b b B Q M i # b 2



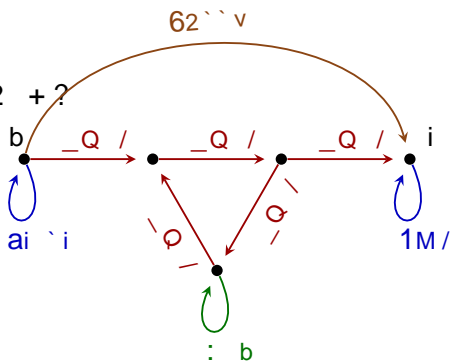
A M } M B i 2 H v K M v K i + ? 2 b # m i i ? 2 m b 2 ` 2 t T

- > Q K Q K Q ` T ? B B K i 2 ` b Q m i K Q b i B M 7 Q ` K i B Q M
- a ? Q ` i 2 b i ! @ ! H 5 + Q p 2 ` ; 2 ô Q 7 K i + ? B M ; r H F b
- h ` B H S ` Q # H 2 K b ` 2 + Q K T m i i B Q M H H v ? ` /
J v / B b + ` / K 2 M B M ; 7 m H K i + ? B M ; r H F

L Q b Q H m i B Q M B b + H 2 ` H v b m T 2 ` B

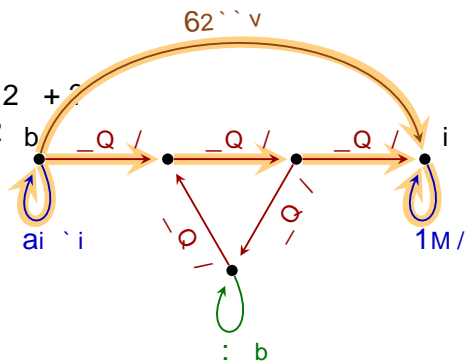
.2}MBiBQM

- $_2im`Mbr HFb$
- $1 + ? 2 / ; 2 + MK i + ? 2 + ?$
 $iQK Q7i KQbi QM+2$



.2}MBiBQM

- $_2im`Mb r HFb$
- $1 + ? 2 / ; 2 + MK i + ? 2 + ?$
 $i Q K Q 7 i K Q b i Q M + 2$



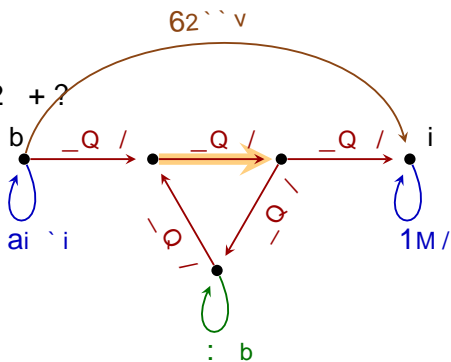
$$Z_1 = a U_Y 6) 1$$

- $_2im`Mb$
 - $i ? 2 7 2 ``v$
 - $i ? 2 b i ` B ; ? i ` Q /$



.2}M BiBQM

- $_2im`Mb r HFb$
- $1 + ? 2 / ; 2 + MK i + ? 2 + ?$
 $i Q K Q 7 i K Q b i Q M + 2$



$$Z_1 = a U_Y 6) 1$$

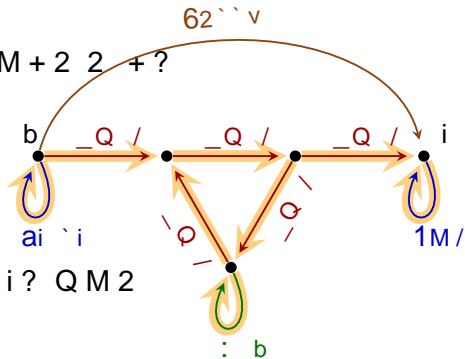
- $_2im`Mb$
 - $i ? 2 7 2 `` v$
 - $i ? 2 b i ` B ; ? i ` Q /$
- $AM r HFb r Bi ? + B ` + m Bi H T b$
 $4) i ? 2 KB // H 2 2 / ; 2 ` 2 m b 2 b$

.2}MBiBQM

- $_2im`Mbr HFb$
- $1 + ? 2 / ; 2 K v mb 2 QM + 2 2 + ?$
iQK BM Z

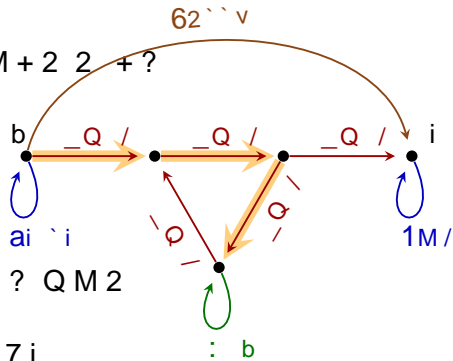
$$Z_2 = a U_Y \theta : U_Y \theta 1$$

- $_2im`Mbi?2 r HF rBi? QM2$
 $+B`+mBi H T$



.2}MBiBQM

- $_2im`Mbr HFb$
- $1 + ? 2 / ; 2 K v mb 2 QM + 2 2 + ?$
iQK BM Z

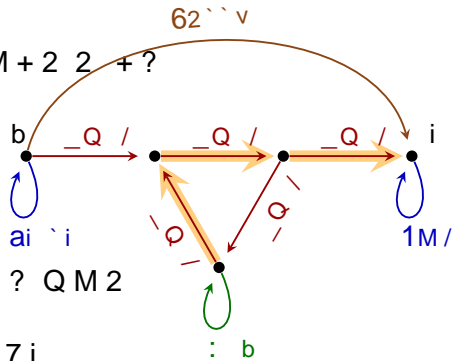


$$Z_2 = a U_Y \theta : U_Y \theta 1$$

- $_2im`Mbi?2r HF rBi? QM2$
 $+B`+mBi H T$
■ "27Q:`2 mb2 i?2_H27i

.2}MBiBQM

- $_2im`Mbr HFb$
- $1 + ? 2 / ; 2 K v mb 2 QM + 2 2 + ?$
iQK BM Z



$$Z_2 = a U_Y \theta : U_Y \theta 1$$

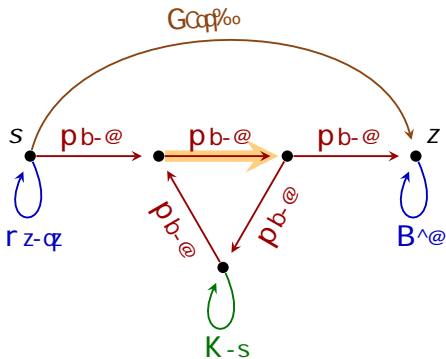
- $_2im`Mbi?2 r HF rBi? QM2$
 $+B`+mBi H T$
 - "27Q:`2 mb2 i?2_H27i
 - 7i2:`! mb2 i?2 `B;?i

Definition

- Returns walks
- Each edge may use once each atom in Q

$$k_2 = r (p + G)^* K (p + G)^* B$$

- Returns the walk with one circuit lap
 - Before K / use the left p
 - After K / use the right p
- In walks with 2+ circuit laps
 - E) the middle edge reuses the left p or the right p

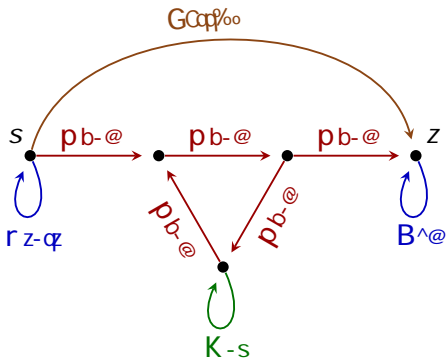


$$k_1 = r (p + G)^* B$$

- The ferry
- The straight road

$$k_2 = r (p + G)^* K (p + G)^* B$$

- The walk with one circuit lap



$$k_1 = r (p + G)^* B$$

- The ferry
- The straight road

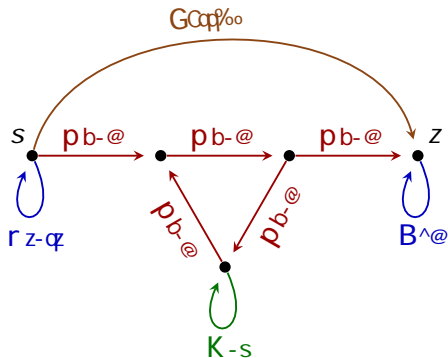
$$k_2 = r (p + G)^* K (p + G)^* B$$

- The walk with one circuit lap

Lemma

\mathcal{B} match ... of k

E) some subwalk $s..$ returned



$$k_1 = r (p + G)^* B$$

- The ferry
- The straight road

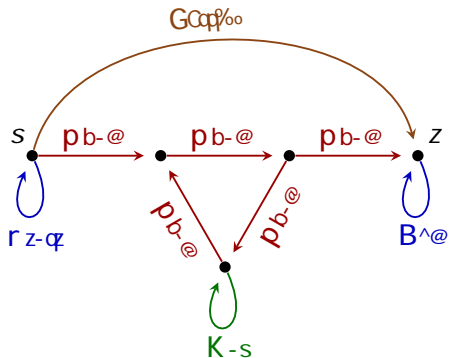
$$k_2 = r (p + G)^* K (p + G)^* B$$

- The walk with one circuit lap

Lemma

\mathcal{B} match ... of k

E) some subwalk $s..$ returned



Open questions

- Define “coverage”
- Define “good” coverage

Binding-trail is compatible with Homomorphism

Homomorphism semantics returns $(s; z)$

() Binding-trail semantics returns some walk ... from s to z .

E) If users need endpoints only, use algorithm for homomorphism

Binding-trail is compatible with Shortest-walk

Shortest matching walks Binding-trail matching walks.

E) If users need only one walk, use algorithm for shortest-walk

Binding-trail is compatible with Trail

Matching trails Binding-trail matching walks.

Tractable problems

- 1 $B \setminus \text{ez}^{\text{S}} \text{C}^{\text{SS}}$ is NL-complete.
- 2 $B^{\wedge} \sim \setminus \text{Cq}^{\text{z}} \text{S}^{\text{L}}$ the 4-L of answers is Poly-delay.

Untractable problems

- 3 ; $b \sim \wedge \text{z}^{\text{S}} \text{L}$ the number of matched walk is $\sim d$ complete.
- 4 [$\text{C} \setminus \text{4C}^{\text{PS}}$ of a given walk is NP-complete.

About 3 : Counting is $\sim d$ complete for any reasonable semantics.

About 4 : Mostly a theoretical concern.

Open problem

- 5 $B^{\wedge} \sim \setminus \text{Cq}^{\text{z}} \text{S}^{\text{L}}$ the SCz of answers.

The 🐘 in the room

yPC b-ze-z @CeC^@s b^ zPC s%z- ‡ bHzPC | ~Cq%o

p^*

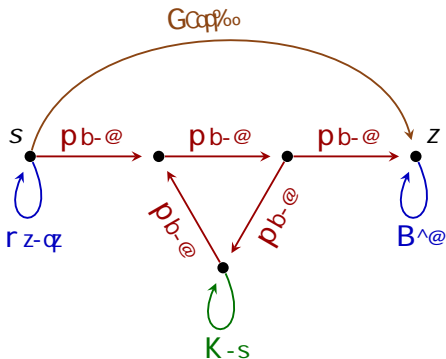
- allows no lap in the circuit

p^*p^*

- allows 1 lap in the circuit

$(p + p)^*$

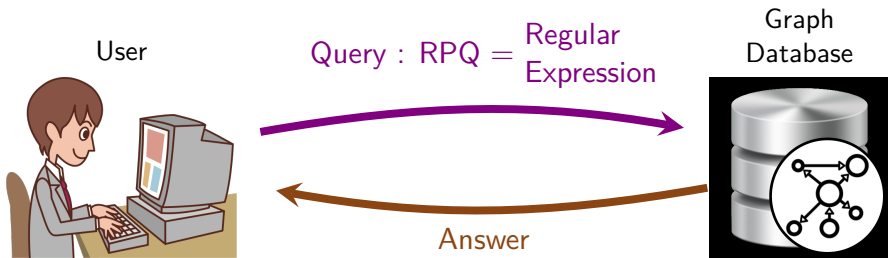
- allows 1 lap in the circuit
- In general, $\notin p^*p^*$



Unusual from theoretical point of view 🤖

The user has finer control on the output 😊

This kind of syntax quirks exists in practice 🤔



- Coverage of matches
- Tractable emptiness and enumeration
- Syntax-dependent

Perspectives

- Deduplicated enumeration
- Containment
- How to deal with data values?
- Get binding-trail into GQL 😄

- Theoreticians mostly worry about complexity 🎓
- But complexity is not the only relevant criterion
Why? Trail semantics is the most popular in industry 🧑💻

Ideas for a comparison framework 🏢

- Coverage (horizontal aggregation)
- Monotonicity (distributed databases)
- Compatibility with operators (result predictability)
- Kinds of definition (selector/restrictor in GQL)

Promising semantics 🙌

- Undominated semantics (minimal for the subwalk order)
- Shortest coverage semantics

Appendix

Semantics	Shortest-walk	Trail	Run-based
Existence	■ Tractable	■ Untractable	■ Tractable
Enumeration	■ Tractable	■ Untractable	■ Tractable
Distinct Enum	■ Tractable	■ Untractable	Open
Counting	■ Meaningless	■ Untractable	■ Untractable
Walk Memb.	■ Tractable	■ Tractable	■ Untractable
Coverage	■ None	■ No guarantee	■ "Subwalk" guarantee

Simple run

- Query given as an automaton
- Outputs simple walks in the product ? A_Q
- Good formal setting for theory

Binding trail

- Query given as a RegExp
- Outputs matching walks that do not reuse edge on a same atom
- Closer to practice

Theorem

The two semantics are computationally equivalent.

Key idea for \Rightarrow

From an automaton A , we build a regular expression B such that runs of A are encoded into runs of the Glushkov automaton of B .