# Path Logics for Querying Graphs
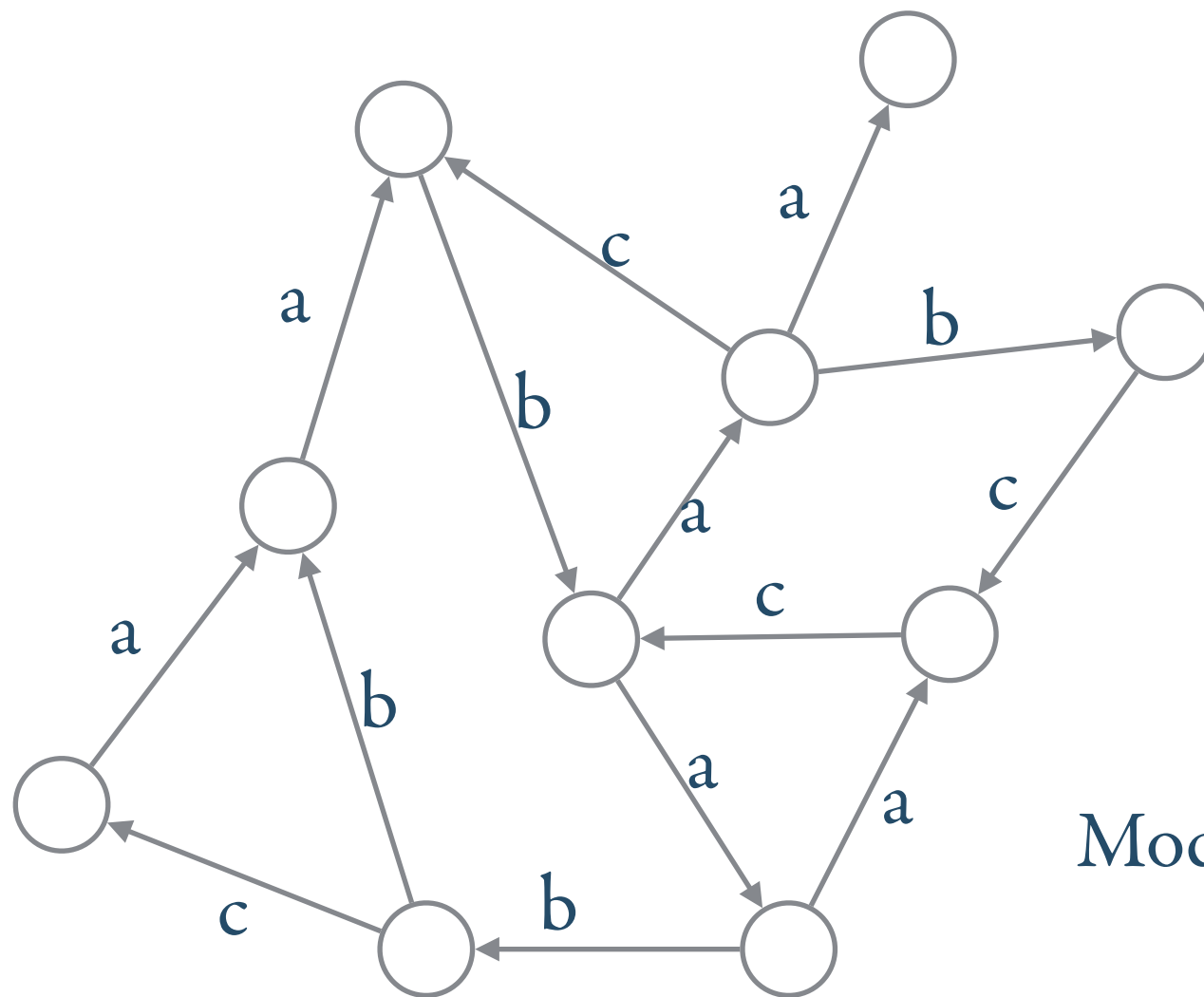
## combining expressiveness and efficiency

**Diego Figueira**

CNRS, LaBRI

France

# Graph*databases*

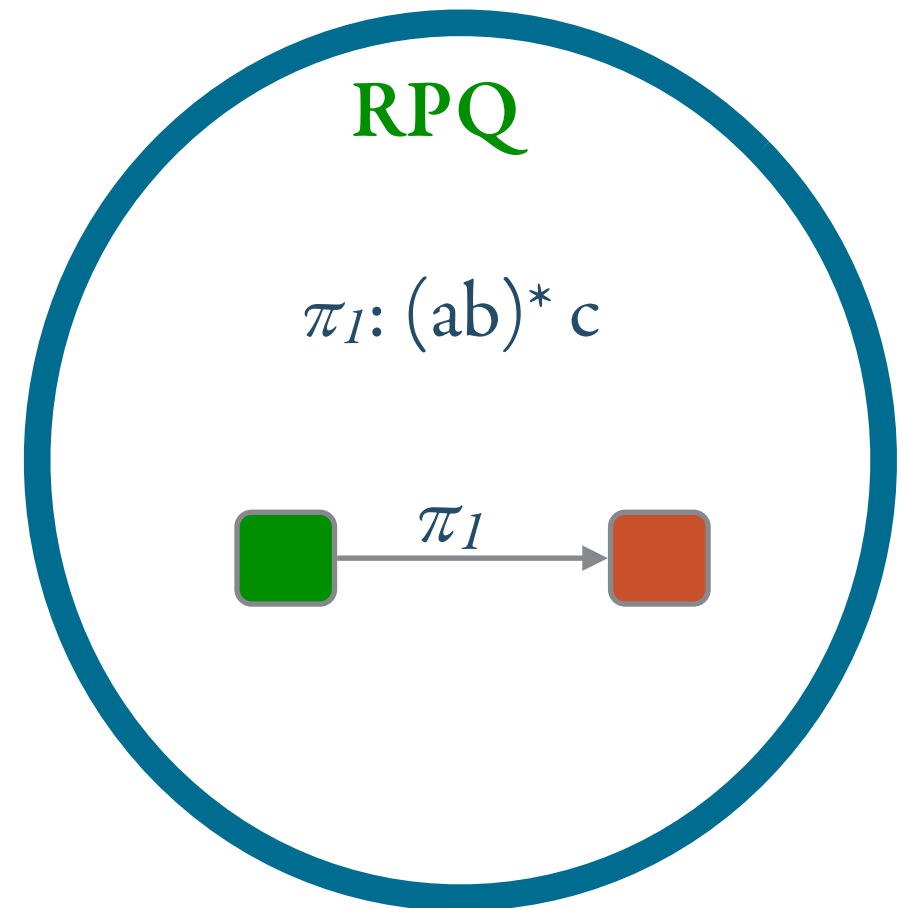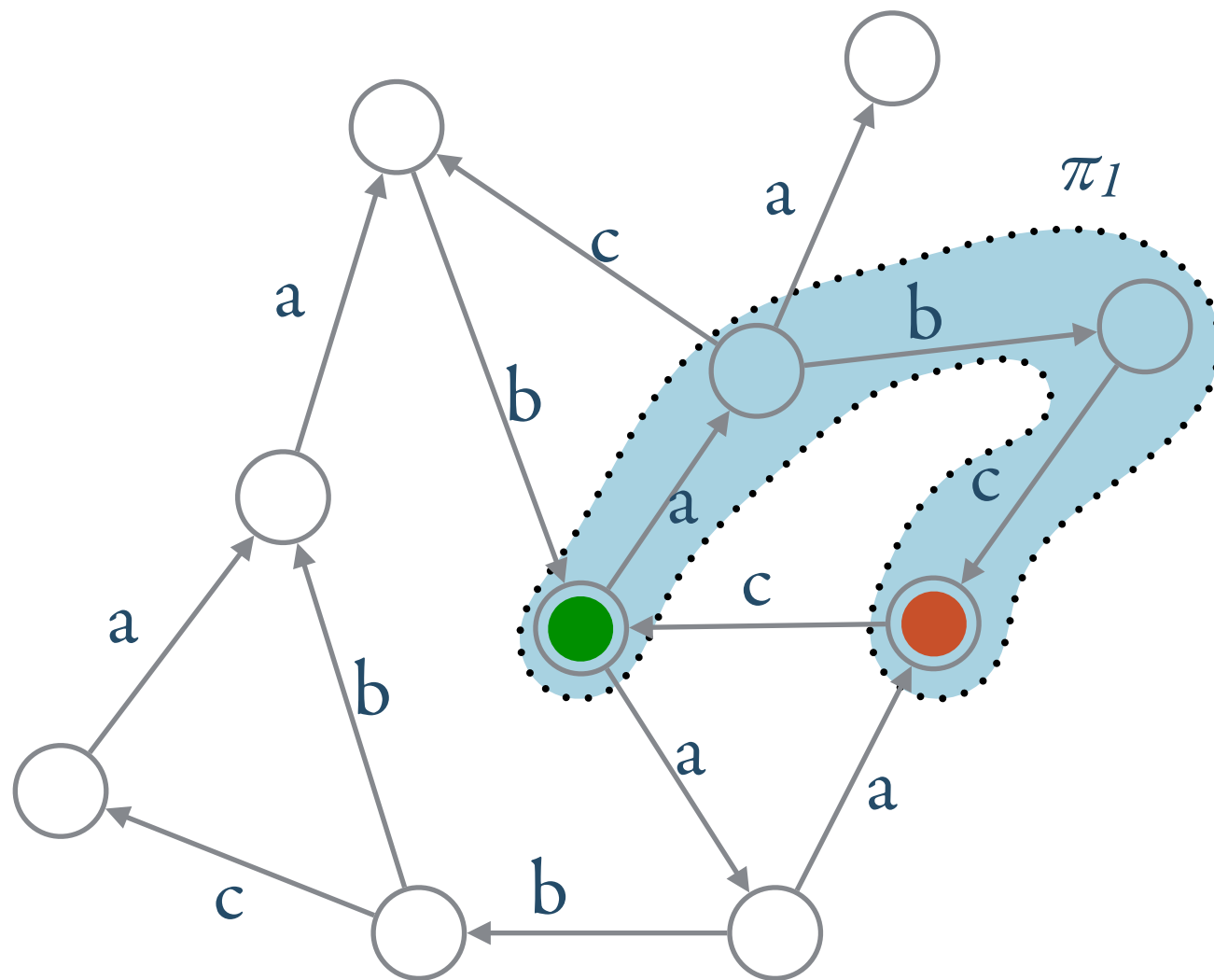Semantic web / RDF / social networks / …

"Entities + Relations"

Modelled as: edge-labelled directed graphs

Notion of *path* of central importance

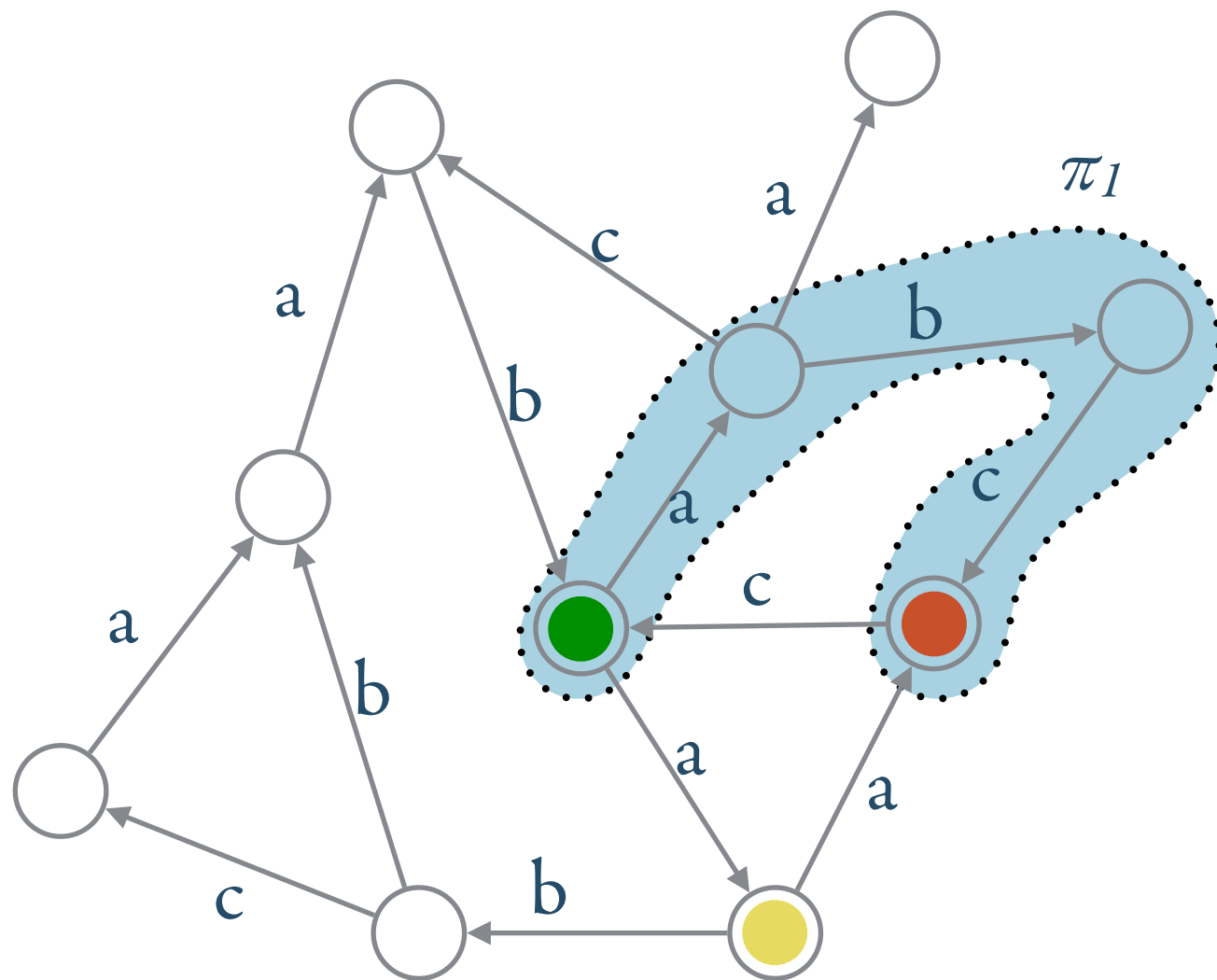# Graph*databases*



RPQ

$\pi_1$: (ab)* c

$\pi_1$

Evaluation: P (combined)
NL (data)

# Graph*databases*



**CRPQ**

$\pi_1$: (ab)* c

$\pi_2$: (ac)*

$\pi_3$: a c*

**Acyclic**

**P**

Evaluation: ~~NP~~ (combined)

NL (data)

Unions, inverse

# Graph*databases*



**CRPQ**

$\pi_1$: (ab)* c

$\pi_2$: (ac)*

$\pi_3$: a c*

What about...

*"All the pairs (u,v) that can reach some node z in the same number of steps"*

# Graph*databases*

What about testing for **relations** on the paths?

- $|\pi_i| = |\pi_j|$

- $\pi_i$ is a **prefix** of $\pi_j$

- $\pi_i$ is a **subsequence** of $\pi_j$

- $\pi_i$ is a **factor** of $\pi_j$

- $\pi_i = \pi_j$ projected onto A

**CRPQ(S)**

**CRPQ**

$\pi_1$: (ab)* c

$\pi_2$: (ac)*

$\pi_3$: a c*



$\pi_1$

$\pi_2$   $\pi_3$

R($\pi_1, \pi_2$), R∈S

Motivations from: entity resolution, semantic associations, crime detection,...

# Graph*databases*

What about testing for **relations** on the paths?

$$\text{CRPQ}(S) = \begin{array}{l} \text{CRPQ} + \\ \text{tests } R(\pi_{i_1,\dots,i_n}), R \in S \end{array}$$

**S**: Class of well-behaved word relations...

**CRPQ(S)**

**CRPQ**

$\pi_1$: (ab)* c

$\pi_2$: (ac)*

$\pi_3$: a c*

$\pi_1$

$\pi_2$  $\pi_3$

$R(\pi_1, \pi_2), R \in S$

*binary relations*

$R \subseteq \mathbb{A}^* \times \mathbb{A}^*$

REC₂ *recognizable*

REG₂ *regular*

RAT₂ *rational*

prefix, equal, equal length, ...

suffix, infix, projection, subsequence, ...

# Graph*databases*

$CRPQ(S) =$ **CRPQ +**
tests $R(\pi_{i_1,\ldots,}\pi_{i_n})$, $R \in S$

**CRPQ(REC) NP/NL complexity**

**Can this be extended?**

**CRPQ(REG) PSPACE/NL complexity**

**CRPQ(RAT)** *undecidable*

**Related to the Intersection Problem:**
Given relations $R_1,\ldots,R_n$, whether $R_1 \cap \cdots \cap R_n \neq \varnothing$

# intersection *problem*

$R \cap S = \varnothing$ ?

$R, S$ : classes of binary relations

it has been studied...

$(u_{i_1} \ldots u_{i_n}, v_{i_1} \ldots v_{i_n})$

PCP

REG $\cap$ RAT $= \varnothing$ ?
already undecidable

*input:* $R \in R$, $S \in S$
*output:* $R \cap S = \varnothing$ ?

...but what about
real world relations?

like

$u \sqsubseteq v$

| $u$ | a | b | a | c | a | b | a | c |
|---|---|---|---|---|---|---|---|---|

| $v$ | a | a | b | a | b | a | c | a | c | b | c | b | a | c |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

subsequence

**subsequence**...?

**subword**...?

**suffix**...?

# Can we extend CRPQ beyond REG relations?

| Language | Data complexity | Combined complexity |
|---|---|---|
| CRPQ(REG$_k$) | NL | PSPACE |
| CRPQ(RAT$_k$) | Undecidable | Undecidable |
| CRPQ(REG$_k$ + suffix) | Undecidable | Undecidable |
| CRPQ(REG$_k$ + factor) | Undecidable | Undecidable |
| CRPQ(REG$_k$ + subsequence) | non-elementary | non-PR |
| CRPQ(suffix) | NL | PSPACE |
| CRPQ(factor) | PSPACE | PSPACE |
| CRPQ(subsequene) | PSPACE | NEXPTIME |

$\forall \, k > 1$

# Can we extend CRPQ beyond REG relations?

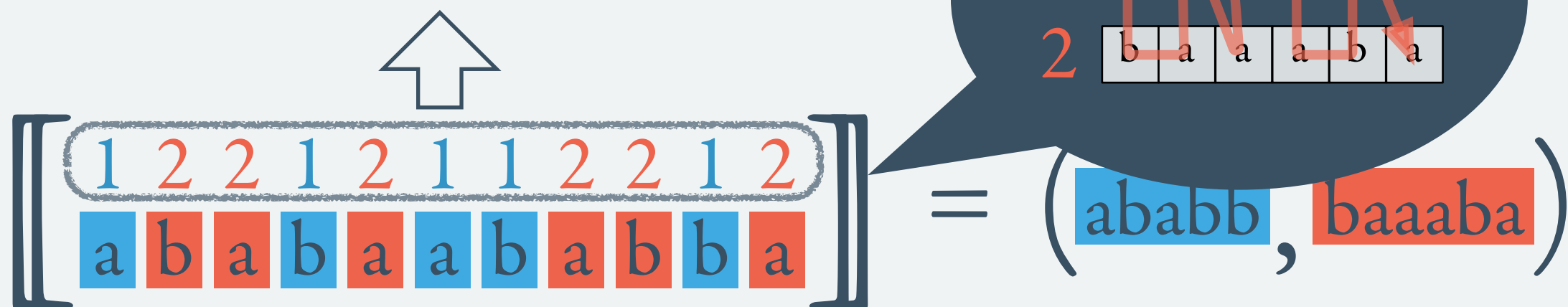Proposed alternative: **approximate** RAT through REG + counters

How?
1) take a an NFA
2) add counters
3) use it to read $k$-tuples of words

# 2 tapes over $\mathbb{A}$ ≈ 1 tape over $\mathbb{A} \times \{1,2\}$

control word

$$\left[\!\!\left[\begin{array}{c} 1\ 2\ 2\ 1\ 2\ 1\ 1\ 2\ 2\ 1\ 2 \\ a\ b\ a\ b\ a\ a\ b\ a\ b\ b\ a \end{array}\right]\!\!\right] = \left(\text{ababb}, \text{baaaba}\right)$$



$$\left[\!\!\left[ (\mathbb{A}\times\{1,2\})^* \right]\!\!\right] = \mathbb{A}^* \times \mathbb{A}^*$$

**(1|2)\*-controlled**

$$\left[\!\!\left[ ((a,1)(a,2)|(b,1)(b,2))^* \right]\!\!\right] = \text{equality}$$

**(12)\*-controlled**

$L \subseteq \{1,2\}^*$

$$\mathbf{Rel}(L) = \{\, \left[\!\left[ S \right]\!\right] \mid S \in \mathrm{REG}(\mathbb{A}\times\{1,2\}) \text{ is } L\text{-controlled} \,\}$$

# Eg:

$$\mathbf{Rel}((1|2)^*) = \mathrm{RAT}_2$$

$$\mathbf{Rel}((12)^*(1^*|2^*)) = \mathrm{REG}_2$$

$$\mathbf{Rel}(1^*2^*) = \mathrm{REC}_2$$

$$\mathbf{Rel}((12)^*) = \text{length-preserving } \mathrm{REG}_2$$

$$\mathbf{Rel}((1^*|2^*)(12)^*) = \mathrm{REG}_2^{rev}$$

Instead of **regular languages**...

$$\mathbf{Rel}(L) = \{ [\![ S ]\!] \mid S \in \mathbf{\color{red}REG}(\mathbb{A} \times \{1,2\}) \text{ is } L\text{-controlled} \}$$

...use **automata with counting**

Idea

**Evaluation of CRPQ with counting is feasible**

**PSPACE** in combined complexity

**NL** in data complexity

# Parikh Automata[*]

[Klaedtke & Rueß]

dimension

NFA with **n** counters $c_1,...,c_n$ and a semilinear set $S \subseteq \mathbb{N}^n$

$(\mathbb{A}, Q, q_0, \delta, F, \mathbf{n}, \mathbf{S})$

Transitions of $\delta$: $(q, a, (x_1,...,x_n), q') \in Q \times \mathbb{A} \times \mathbb{N}^n \times Q$

Run:

- Initial configuration: $(q_0, (0,...,0)) \in Q \times \mathbb{N}^n$

counters can only be **incremented**

- $(q, \mathbf{x}) \xrightarrow{(q,a,\mathbf{y},p)} (p, (\mathbf{x+y}))$ $\in \delta$

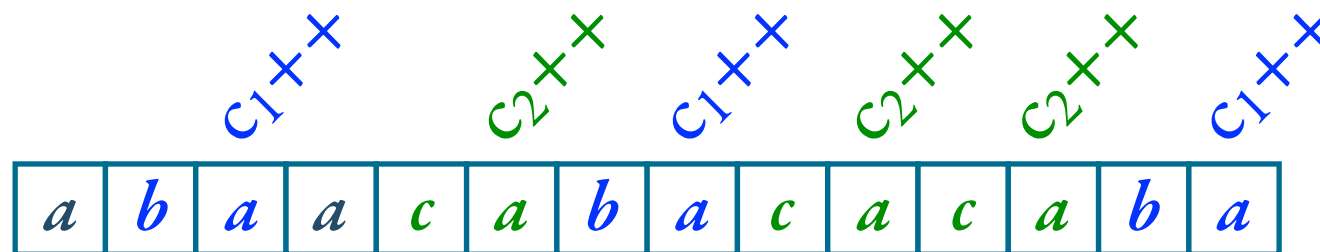- Acceptance: last configuration in $F \times \mathbf{S}$

✱ Many equivalent definitions (eg. reversal-bounded counter systems)

# Parikh Automata

Eg: $L_{ba=ca} = \left\{ w \;\middle|\; \begin{array}{c} \text{number of } a\text{'s after a } b \\ = \\ \text{number of } a\text{'s after a } c \end{array} \right\}$



Parikh Automaton $A = (\mathbb{A}, Q, q_0, \delta, F, 2, \{(k,k) \mid k \in \mathbb{N}\})$

- dimension 2 (2 counters)

- increment $c_1$ whenever we see "**ba**"

- increment $c_2$ whenever we see "**ca**"

- F=Q

- Semilinear set assures that counters must be equal to accept a word

# Parikh Automata

## Decidable

non-emptiness,

membership

## Closed under

intersection,

union,

(inverse) homomorphisms,

concatenation

(not complementation/iteration)

$$\mathbf{Rel}^{\mathbf{PA}}(L) = \{\; [\![S]\!] \mid S \in \mathbf{PA}(\mathbb{A} \times \{1,2\}) \text{ is } L\text{-controlled}\;\}$$

PA relations

Eg:

$$\text{REG}_2^{\mathbf{PA}} = \mathbf{Rel}^{\mathbf{PA}}((12)^*(1^*|2^*))$$

$$\text{REG}_{2\,rev}^{\mathbf{PA}} = \mathbf{Rel}^{\mathbf{PA}}((1^*|2^*)(12)^*)$$

$$\text{RAT}_2^{\mathbf{PA}} = \mathbf{Rel}^{\mathbf{PA}}((1|2)^*)$$

$$\vdots$$

Word relations

RAT$_k$
*rational*

REG$_k$
*regular*

REC$_k$
*recognizable*

REG$_k^{PA}$
*Parikh-regular*

**Theorem:** Evaluation of CRPQ($\mathbf{REG^{PA}}$) is
$\qquad$ **PSPACE** in combined complexity
$\qquad$ **NL** in data complexity

Proof ingredients:

- **Intersection problem** for Parikh Automata

  Given PA's $A_1,...,A_n$, is $L(A_1) \cap \cdots \cap L(A_n) \neq \varnothing$ ?

  is PSPACE-complete

- **Intersection closure** for $REG^{PA}$

  For all $R,S \in REG^{PA}$, $R \cap S \in REG^{PA}$
  it suffices to intersect the automata representing them

- **Closure under product** of $REG^{PA}$

**Theorem:** Evaluation of CRPQ$^{PA}$ (no relations) is
$\qquad$ **NP** in combined complexity
$\qquad$ **NL** in data complexity

# Approximating rational relations

$u \sim_k v$ are **k-similar** iff ~~for all $w$ with $|w| \leq k$, they have the same number of appearances of $w$ (as factor)~~ (as subsequence)

Given $R \in \mathrm{RAT}$,

$$R_k = \{(u,v) \mid u \sim_k u', v \sim_k v', (u', v') \in R\} \in \mathrm{REG}^{\mathrm{PA}}$$

# Alternative: Syntactic restrictions

E.g.

$\pi_1$: (ab)* c    $R(\pi_1, \pi_3)$

Maximum cardinality of connected component

$\pi_2$: (ac)*    $S(\pi_3, \pi_2)$

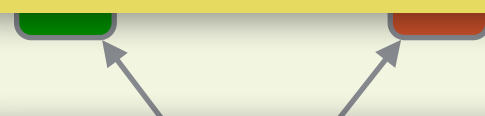$\pi_3$: a c*    $R(\pi_3, \pi_2)$

**cyclic**

Theorem: Evaluation of **acyclic**-CRPQ(**RAT$^{PA}$**) is

**PSPACE** in combined complexity

**NL** in data complexity

If also fixed **join size**: **NP** combined complexity

E.g.

$\pi_2$: (ac)*

$R(\pi_1, \pi_3)$

**acyclic**

If also fixed **PA dimension** and **unary representation**:

**PTIME** combined complexity

# Conclusion

**Counting does not increase complexity**

**Avoid the curse of of rational relations**

Approximating by regular relations with counting

Or staying away from cycles in path relations

Thank you