# Designing and Comparing RPQ Semantics

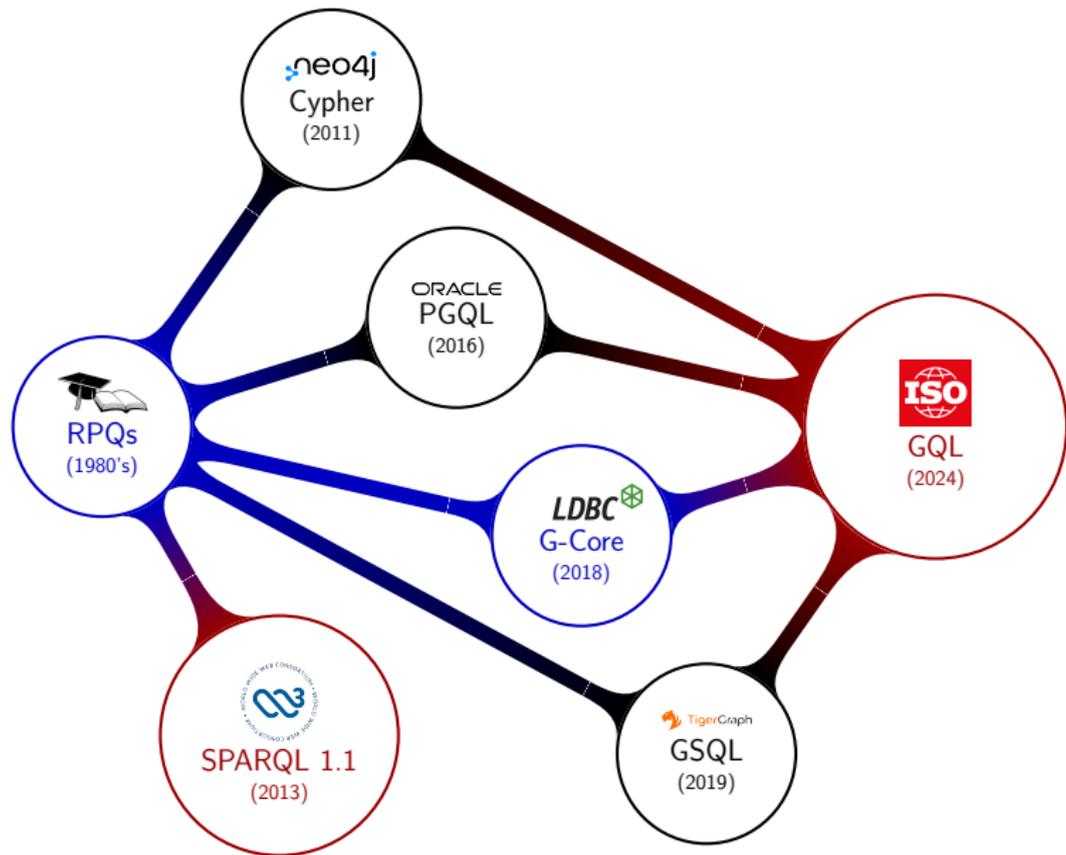Victor Marsault                    Antoine Meyer
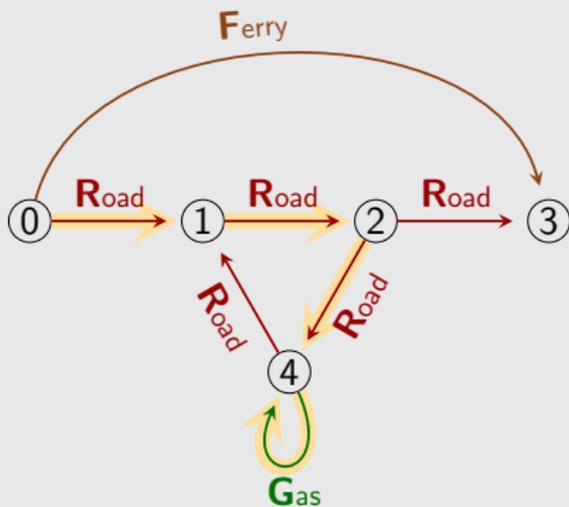
Univ. Gustave Eiffel, CNRS, LIGM – France

Labeled graph D

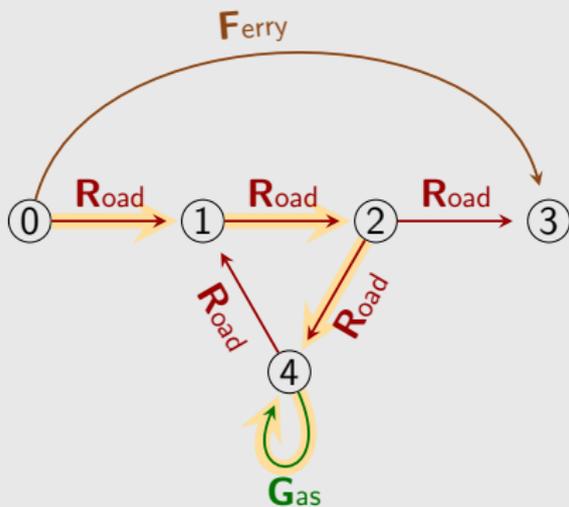Nodes, Edges, Labels, Paths

## Labeled graph $D$

## Regular Path Query $Q$

- Input: a regexp over labels
  Ex: $(\mathbf{R} + \mathbf{F})^* \mathbf{G}$

- Match: path with a label that conforms to $Q$
  Ex: $0 \xrightarrow{\mathbf{R}} 1 \xrightarrow{\mathbf{R}} 2 \xrightarrow{\mathbf{R}} 4 \xrightarrow{\mathbf{G}} 4$

$\textsc{Matches}(D, Q)$: match set

- Issue: $\textsc{Matches}(D, Q)$ is sometimes infinite
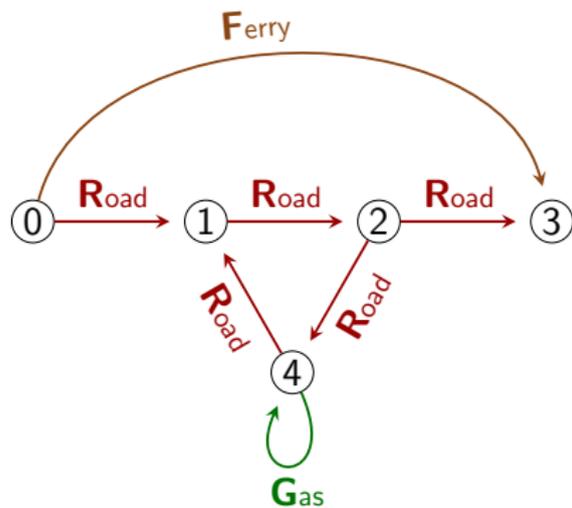
Nodes, Edges, Labels, Paths

**Endpoint** semantics
▶ Outputs the **endpoints** of matches

**Trail** semantics
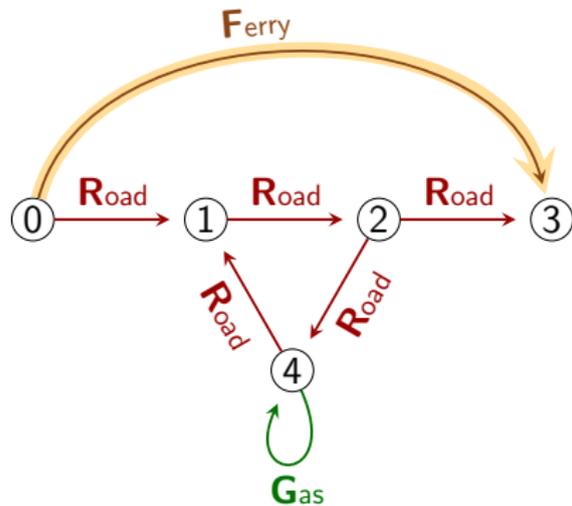▶ Outputs **matches** with **no repeated edges**

**Shortest** semantics
▶ Outputs **matches** with the **minimum number of edges**

Three kinds of $0 \rightsquigarrow 3$ paths:

▶ The ferry
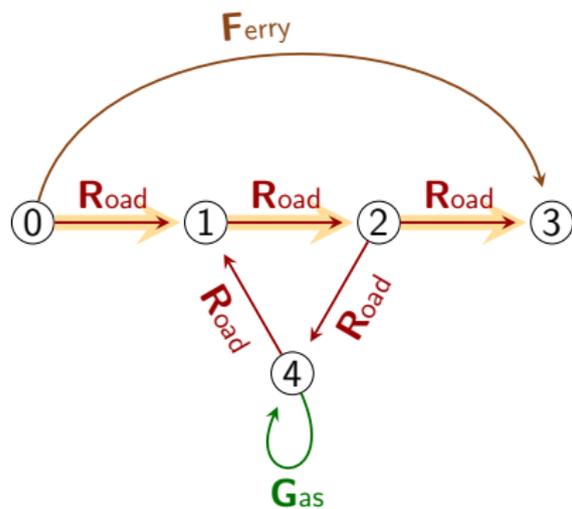▶ The straight road
▶ Some road path with laps

Three kinds of $0 \rightsquigarrow 3$ paths:
- ▶ The ferry
- ▶ The straight road
- ▶ Some road path with laps

Three kinds of $0 \rightsquigarrow 3$ paths:
- ▶ The ferry
- ▶ The straight road
- ▶ Some road path with laps

Three kinds of $0 \rightsquigarrow 3$ paths:
▶ The ferry
▶ The straight road
▶ Some road path with laps

$$Q_1 = (\mathbf{R} + \mathbf{F})^*$$

- Endpoint outputs $(0, 3)$
- Shortest outputs the ferry
- Trail outputs the ferry and the straight road

Three kinds of $0 \rightsquigarrow 3$ paths:

- The ferry
- The straight road
- Some road path with laps

$$Q_1 = (\mathbf{R} + \mathbf{F})^*$$
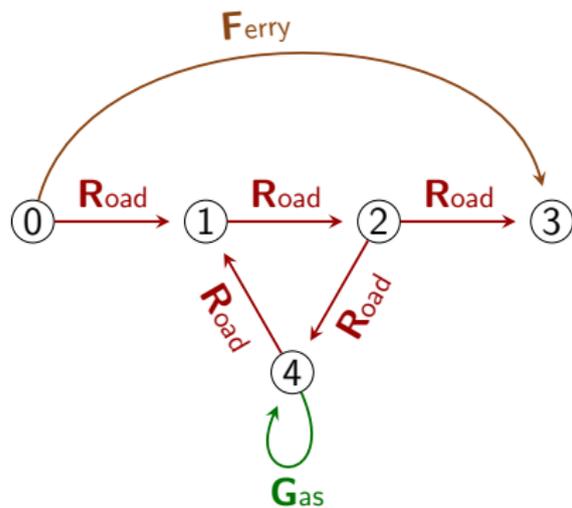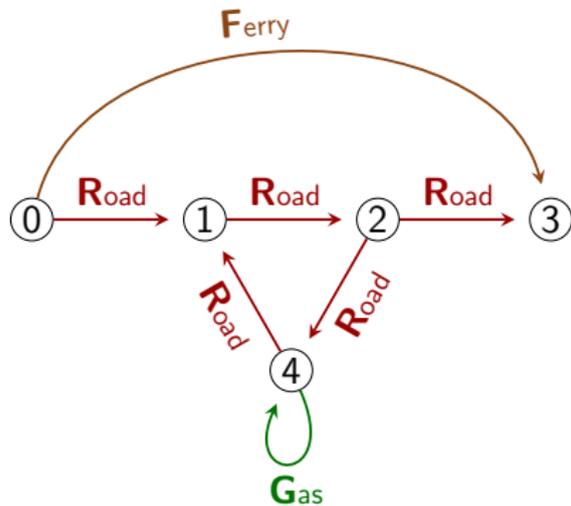
- Endpoint outputs $(0, 3)$
- Shortest outputs the ferry
- Trail outputs the ferry and the straight road

$$Q_2 = (\mathbf{R} + \mathbf{F})^* \mathbf{G} (\mathbf{R} + \mathbf{F})^*$$

- Endpoint outputs $(0, 3)$
- Shortest outputs the road with one lap
- Trail has no output

Three kinds of $0 \rightsquigarrow 3$ paths:
- The ferry
- The straight road
- Some road path with laps

Complexity perspective

|  | Shortest | Trail |
|---|---|---|
| Existence<br>*"Is there a path in the output?"* | ■ Tractable | ■ Untractable |
| Enumeration<br>*"Enumerate all paths in the output"* | ■ Tractable | ■ Untractable |
| Membership<br>*"Is a given path in the output?"* | ■ Tractable | ■ Tractable |

Complexity perspective

|  | Shortest | Trail |
|---|---|---|
| Existence<br>*"Is there a path in the output?"* | ■ Tractable | ■ Untractable |
| Enumeration<br>*"Enumerate all paths in the output"* | ■ Tractable | ■ Untractable |
| Membership<br>*"Is a given path in the output?"* | ■ Tractable | ■ Tractable |

### Questions
▶ Why was trail (Neo4j/Cypher) the most successful?
▶ What makes a semantics good?
▶ Can we design better semantics?

RPQ semantics = function $S$

$$S : (D, Q) \longmapsto O \quad \subseteq_{\text{fin}} \text{MATCHES}(D, Q)$$

- Database $D$
- RPQ $Q$
- **Finite** set $O$ of outputs

Ex: Trail and Shortest

RPQ semantics $=$ function $S$

$$S : (D, Q) \longmapsto O \quad \subseteq_{\text{fin}} \text{MATCHES}(D, Q)$$

▶ Database $D$
▶ RPQ $Q$
▶ **Finite** set $O$ of outputs

Ex:   Trail   and   Shortest

Filter-based semantics

*"Keeps the **good** matches w.r.t. some boolean function $f$"*

$f(p) = \text{good} \iff p$ is a trail
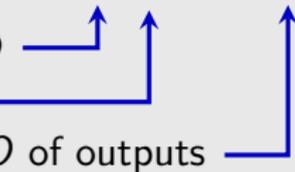
RPQ semantics = function $S$

$$S \colon (D, Q) \longmapsto O \quad \subseteq_{\mathsf{fin}} \text{MATCHES}(D, Q)$$

▶ Database $D$
▶ RPQ $Q$
▶ **Finite** set $O$ of outputs

Ex:  Trail  and  Shortest

### Filter-based semantics

*"Keeps the **good** matches w.r.t. some boolean function f"*

$f(p) = \text{good} \iff p$ is a trail

### Order-based semantics

*"Keeps the **best** matches w.r.t. some well-quasi-order $\prec$"*

$p_1 \prec p_2 \iff p_1$ is shorter than $p_2$

Given $S: (D, Q) \longmapsto O \subseteq_{\mathsf{fin}} \textrm{MATCHES}(D, Q)$

Given $\quad S \colon (D, Q) \longmapsto O \quad \subseteq_{\mathsf{fin}} \mathrm{MATCHES}(D, Q)$

$S$ is compatible with $+$

$$S(D, Q_1 + Q_2) = S(D, Q_1) \cup S(D, Q_2)$$

$S$ is compatible with $\cdot$

$\cdots$

$S$ is compatible with $*$

$\cdots$

$\Rightarrow$ Bottom-up computation

$\Rightarrow$ User Explainability

Given $S\colon (D, Q) \longmapsto O \subseteq_{\mathsf{fin}} \mathrm{MATCHES}(D, Q)$

$S$ is compatible with $+$
$S(D, Q_1 + Q_2) = S(D, Q_1) \cup S(D, Q_2)$

$S$ is monotonous
$D_1 \subseteq D_2 \Rightarrow S(D_1, Q) \subseteq S(D_2, Q)$

$\Rightarrow$ Distributed databases

$S$ is compatible with $\cdot$
$\cdots$

$S$ is compatible with $*$
$\cdots$

$\Rightarrow$ Bottom-up computation
$\Rightarrow$ User Explainability

Under various mild hypotheses,

- A semantics is never ·-compatible, nor ∗-compatible
- Compatible with $+$ $\iff$ Filter-based

Under various mild hypotheses,

▶ A semantics is never $\cdot$-compatible, nor $*$-compatible

▶ Compatible with $+$ $\iff$ Filter-based

▶ Filter-based semantics are monotonous

▶ Order-based semantics are **not** monotonous
(Adding new elements might create a shortcut)

Complexity & properties

| | Shortest | Trail |
|---|---|---|
| Existence<br>*"Is there a path in the output?"* | ■ Tractable | ■ Untractable |
| Enumeration<br>*"Enumerate all paths in the output"* | ■ Tractable | ■ Untractable |
| Membership<br>*"Is a given path in the output?"* | ■ Tractable | ■ Tractable |
| Monotonous | ■ No | ■ Yes |
| Compatible with $+$ | ■ No | ■ Yes |
| Compatible with $\cdot/*$ | ■ No | ■ No |

### Subpath Minimal semantics

- Order-based semantics
- $p_1 \preceq p_2 \iff p_1$ is a subsequence of $p_2$

## Subpath Minimal semantics

- Order-based semantics
- $p_1 \preceq p_2 \iff p_1$ is a subsequence of $p_2$

## Properties

- It is monotonous
  (unlike most order-based semantics)
- It guarantees some "coverage" of $\mathrm{MATCHES}(D, Q)$
  (only removes cycles, and only if they are unnecessary)

- ▶ Practical RPQ semantics keep a finite number of matches
- ▶ Few ad-hoc semantics were used in practice
- ▶ ...and then query evaluation efficiency was studied

## Systematic study of RPQ semantics as functions

- ▶ Classes and properties
  ⇒ Compare existing semantics
- ▶ Relations between them (implications, incompatibilies, etc.)
  ⇒ Design of future RPQ semantics

## Future work

- ▶ Many ways to extend the framework
- ▶ New interesting semantics to study