

# Scalable Multi-Query Optimization for SPARQL

Wangchao Le<sup>1</sup> Anastasios Kementsietsidis<sup>2</sup> Songyun Duan<sup>2</sup> Feifei Li<sup>1</sup>



<sup>1</sup>University of Utah



<sup>2</sup>IBM Research

April 2, 2012

- 1 Introduction
- 2 Preliminary
- 3 Our approach
- 4 Experiments
- 5 Conclusions

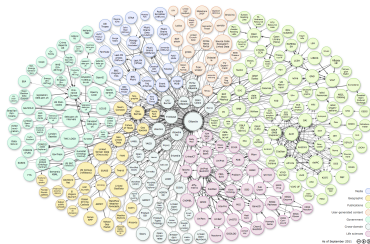
- 1 Introduction
- 2 Preliminary
- 3 Our approach
- 4 Experiments
- 5 Conclusions

# Introduction

- We are inundated with a large collection of RDF (Resource Description Framework) data.

# Introduction

- We are inundated with a large collection of RDF (Resource Description Framework) data.
  - DBpedia, Uniprot, Freebase etc

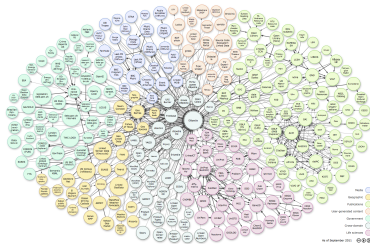


## Internally ...

```
<rdf:RDF
  xmlns:rdf=http://www.w3.org/1999/02/22-rdf-syntax-ns#
  xmlns:dcterms="http://purl.org/dc/terms/" >
  <rdf:Description rdf:about="urn:x-states:New York">
    <dcterms:alternative>NY</dcterms:alternative>
  </rdf:Description>
</rdf:RDF>
```

# Introduction

- We are inundated with a large collection of RDF (Resource Description Framework) data.
  - DBpedia, Uniprot, Freebase etc



## Internally ...

```
<rdf:RDF
  xmlns:rdf=http://www.w3.org/1999/02/22-rdf-syntax-ns#
  xmlns:dcterms="http://purl.org/dc/terms/" >
  <rdf:Description rdf:about="urn:x-states:New York">
    <dcterms:alternative>NY</dcterms:alternative>
  </rdf:Description>
</rdf:RDF>
```

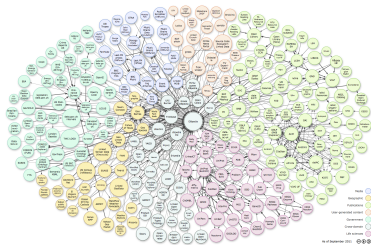
Triple format:

```
<http://.../New York> <http://purl.org/dc/terms/alternative> "NY" .
```

subject predicate object

# Introduction

- We are inundated with a large collection of RDF (Resource Description Framework) data.
  - DBpedia, Uniprot, Freebase etc
  - A large graph and encode rich semantics



## Internally ...

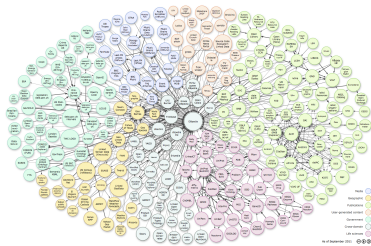
```
<rdf:RDF
  xmlns:rdf=http://www.w3.org/1999/02/22-rdf-syntax-ns#
  xmlns:dcterms="http://purl.org/dc/terms/" >
  <rdf:Description rdf:about="urn:x-states:New York">
    <dcterms:alternative>NY</dcterms:alternative>
  </rdf:Description>
</rdf:RDF>
```

Triple format:

```
<http://.../New York> <http://purl.org/dc/terms/alternative> "NY" .
```

subject predicate object

- We are inundated with a large collection of RDF (Resource Description Framework) data.
  - DBpedia, Uniprot, Freebase etc
  - A large graph and encode rich semantics



Internally ...

```
<rdf:RDF
  xmlns:rdf=http://www.w3.org/1999/02/22-rdf-syntax-ns#
  xmlns:dcterms="http://purl.org/dc/terms/" >
  <rdf:Description rdf:about="urn:x-states:New York">
    <dcterms:alternative>NY</dcterms:alternative>
  </rdf:Description>
</rdf:RDF>
```

Triple format:

```
<http://.../New York> <http://purl.org/dc/terms/alternative> "NY" .
```

subject predicate object

Query language: SPARQL



# Introduction

- We are inundated with a large collection of RDF (Resource Description Framework) data.
  - DBpedia, Uniprot, Freebase etc
  - A large graph and encode rich semantics
- Available engines to manage RDF data?

- We are inundated with a large collection of RDF (Resource Description Framework) data.
  - DBpedia, Uniprot, Freebase etc
  - A large graph and encode rich semantics
- Available engines to manage RDF data?
  - **RDBMS**: Migrate RDF, e.g., Sesame, JenaSDB etc.
  - **Generic RDF stores**: e.g., RDF3X, JenaTDB etc.

[VP07] D.J. Abadi, et al. Scalable semantic web data management using vertical partitioning. In *VLDB*, 2007.

[HEX08] C. Weiss, et al. Hexastore: sextuple indexing for semantic web data management. In *VLDB*, 2008.

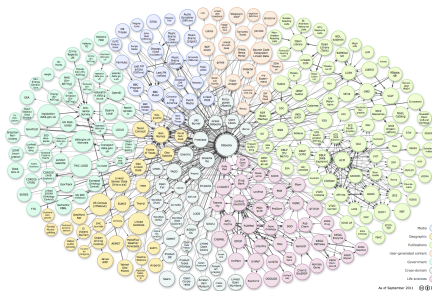
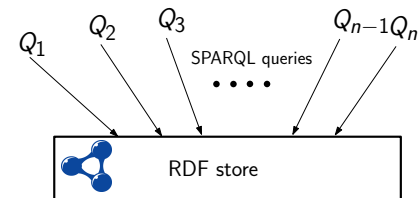
[RDF3X] T. Neumann, G. Weikum. RDF-3X: a RISC-style engine for RDF. In *VLDB*, 2008.

[SJP09] T. Neumann, G. Weikum. Scalable Join Processing on Very Large RDF Graphs. In *SIGMOD*, 2009.

[BM10] M. Atre, et al. Matrix "Bit" loaded: A Scalable Lightweight Join Query Processor for RDF Data. In *WWW*, 2010.

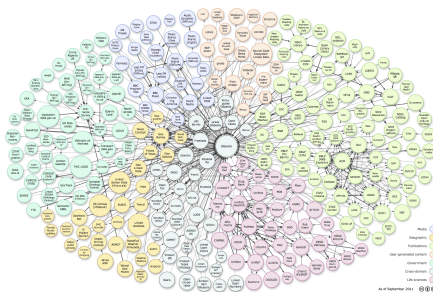
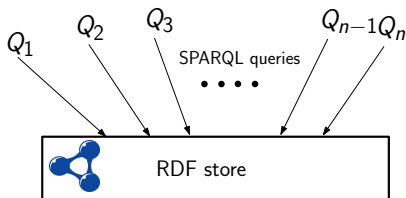
[SSQ11] J. Huang, et al. Scalable SPARQL Querying of Large RDF Graphs. In *VLDB*, 2011.

# Introduction



# Introduction

- Observation: queries share common parts
- Multi-query optimization



- A tempting choice: turn to MQO in relational databases  
[\[MQO88\]](#)[\[MQO90\]](#)[\[MQO00\]](#)
  - SPARQL $\leftrightarrow$ relational algebra [\[EPS08\]](#)[\[FSR07\]](#).
  - Exist quite a few relational solutions for RDF store.

[\[MQO90\]](#) T. Sellis, et al. On the Multiple-Query Optimization Problem. In *TKDE*, 1990.

[\[MQO88\]](#) T. Sellis, et al. Multiple-query optimization. In *TODS*, 1988.

[\[MQO00\]](#) P. Roy, et al. Efficient and extensible algorithms for multi query optimization. In *SIGMOD*, 2000.

[\[EPS08\]](#) R. Angles, et al. The Expressive Power of SPARQL. In *ISWC*, 2008.

[\[FSR07\]](#) A. Polleres, et al. From SPARQL to rules (and back). In *WWW*, 2007.

- A tempting choice: turn to MQO in relational databases  
[\[MQO88\]](#)[\[MQO90\]](#)[\[MQO00\]](#)
  - SPARQL $\leftrightarrow$ relational algebra [\[EPS08\]](#)[\[FSR07\]](#).
  - Exist quite a few relational solutions for RDF store.
- For SPARQL and RDF, new issues arise in practice.

- A tempting choice: turn to MQO in relational databases  
[\[MQO88\]](#)[\[MQO90\]](#)[\[MQO00\]](#)
  - SPARQL $\leftrightarrow$ relational algebra [\[EPS08\]](#)[\[FSR07\]](#).
  - Exist quite a few relational solutions for RDF store.
- For SPARQL and RDF, new issues arise in practice.
  - Convert SPARQL to SQL: not all engines use RDBMS
  - Conversion to SQL  $\rightarrow$  a large number of joins

- A tempting choice: turn to MQO in relational databases  
[\[MQO88\]](#)[\[MQO90\]](#)[\[MQO00\]](#)
  - SPARQL $\leftrightarrow$ relational algebra [\[EPS08\]](#)[\[FSR07\]](#).
  - Exist quite a few relational solutions for RDF store.
- For SPARQL and RDF, new issues arise in practice.
  - Convert SPARQL to SQL: not all engines use RDBMS
  - Conversion to SQL  $\rightarrow$  a large number of joins
  - Store dependent solution



- 1 Introduction
- 2 Preliminary**
- 3 Our approach
- 4 Experiments
- 5 Conclusions

- We focus on two types of queries

- We focus on two types of queries

**Type 1:**  $Q := \text{SELECT RD WHERE GP}$

**Type 2:**  $Q_{\text{OPT}} := \text{SELECT RD WHERE GP (OPTIONAL GP}_{\text{OPT}})^+$

- We focus on two types of queries

**Type 1:**  $Q := \text{SELECT RD WHERE GP}$

**Type 2:**  $Q_{\text{OPT}} := \text{SELECT RD WHERE GP (OPTIONAL GP}_{\text{OPT}})^+$

subj	pred	obj
p1	name	"Alice"
p1	zip	10001
p1	mbox	alice@home
p1	mbox	alice@work
p1	www	http://home/alice
p2	name	"Bob"
p2	zip	10001
p3	name	"Ella"
p3	zip	10001
p3	www	http://work/ella
p4	name	"Tim"
p4	zip	"11234"

(a) triple table D

```
SELECT ?name
WHERE { ?x name ?name, ?x zip 10001,
}
```

(b) Example query  $Q_{\text{OPT}}$

name
"Alice"
"Bob"
"Ella"

- We focus on two types of queries

**Type 1:**  $Q := \text{SELECT RD WHERE GP}$

**Type 2:**  $Q_{\text{OPT}} := \text{SELECT RD WHERE GP (OPTIONAL GP}_{\text{OPT}})^+$

subj	pred	obj
p1	name	"Alice"
p1	zip	10001
p1	mbox	alice@home
p1	mbox	alice@work
p1	www	http://home/alice
p2	name	"Bob"
p2	zip	10001
p3	name	"Ella"
p3	zip	10001
p3	www	http://work/ella
p4	name	"Tim"
p4	zip	"11234"

(a) triple table D

```
SELECT ?name , ?mail, ?hpage
WHERE { ?x name ?name, ?x zip 10001,
        OPTIONAL { ?x mbox ?mail }
        OPTIONAL { ?x www ?hpage } }
```

(b) Example query  $Q_{\text{OPT}}$

name
"Alice"
"Bob"
"Ella"

- We focus on two types of queries

**Type 1:**  $Q := \text{SELECT RD WHERE GP}$

**Type 2:**  $Q_{\text{OPT}} := \text{SELECT RD WHERE GP (OPTIONAL GP}_{\text{OPT}})^+$

subj	pred	obj
p1	name	"Alice"
p1	zip	10001
p1	mbox	alice@home
p1	mbox	alice@work
p1	www	http://home/alice
p2	name	"Bob"
p2	zip	10001
p3	name	"Ella"
p3	zip	10001
p3	www	http://work/ella
p4	name	"Tim"
p4	zip	"11234"

(a) triple table D

```
SELECT ?name , ?mail, ?hpage
WHERE { ?x name ?name, ?x zip 10001,
        OPTIONAL { ?x mbox ?mail }
        OPTIONAL { ?x www ?hpage } }
```

(b) Example query  $Q_{\text{OPT}}$

name	mail	hpage
"Alice"	alice@home	http://home/alice
"Alice"	alice@work	http://home/alice
"Bob"		
"Ella"		http://work/ella

(c) Output  $Q_{\text{OPT}}(D)$

- We focus on two types of queries

**Type 1:**  $Q := \text{SELECT RD WHERE GP}$

**Type 2:**  $Q_{\text{OPT}} := \text{SELECT RD WHERE GP (OPTIONAL GP}_{\text{OPT}})^+$

- Problem statement.

- We focus on two types of queries

**Type 1:**  $Q := \text{SELECT RD WHERE GP}$

**Type 2:**  $Q_{\text{OPT}} := \text{SELECT RD WHERE GP (OPTIONAL GP}_{\text{OPT}})^+$

- Problem statement.

- Input: a set  $Q$  of **Type 1** queries and a data graph  $G$



- We focus on two types of queries

**Type 1:**  $Q := \text{SELECT RD WHERE GP}$

**Type 2:**  $Q_{\text{OPT}} := \text{SELECT RD WHERE GP (OPTIONAL GP}_{\text{OPT}})^+$

- Problem statement.
  - Input: a set  $Q$  of **Type 1** queries and a data graph  $G$
  - Output: a set of **rewritten** queries,  $Q_{\text{OPT}}$  of **Type 1** and **Type 2** queries

- We focus on two types of queries

**Type 1:**  $Q := \text{SELECT RD WHERE GP}$

**Type 2:**  $Q_{\text{OPT}} := \text{SELECT RD WHERE GP (OPTIONAL GP}_{\text{OPT}})^+$

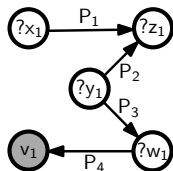
- Problem statement.

- Input: a set  $Q$  of **Type 1** queries and a data graph  $G$
- Output: a set of **rewritten** queries,  $Q_{\text{OPT}}$  of **Type 1** and **Type 2** queries
- Requirements:
  - soundness and completeness:  $Q_{\text{OPT}}(G) \equiv Q(G)$ .
  - cost:  $\frac{T_r(Q) + T_e(Q_{\text{opt}})}{T_e(Q)} \leq 1$

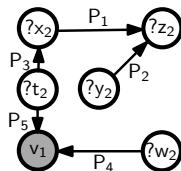
# Our approach

- 1 Introduction
- 2 Preliminary
- 3 Our approach**
- 4 Experiments
- 5 Conclusions

# Motivating example

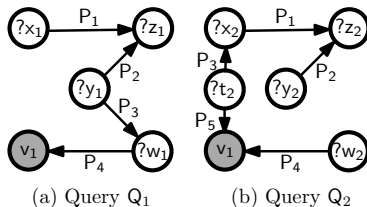


(a) Query  $Q_1$



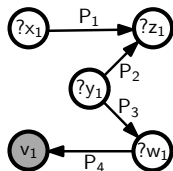
(b) Query  $Q_2$

# Motivating example

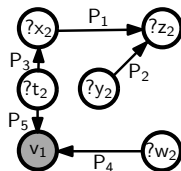


●: constant      ○: variable

# Motivating example

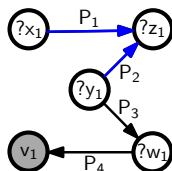


(a) Query  $Q_1$

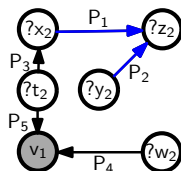


(b) Query  $Q_2$

# Motivating example



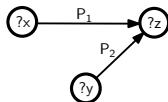
(a) Query  $Q_1$



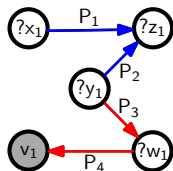
(b) Query  $Q_2$

```
SELECT *  
WHERE { ?x P1 ?z, ?y P2 ?z,
```

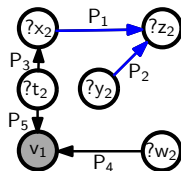
```
}
```



# Motivating example



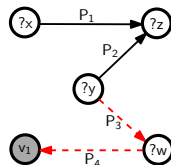
(a) Query  $Q_1$



(b) Query  $Q_2$

```

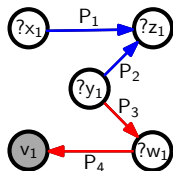
SELECT *
WHERE { ?x P1 ?z, ?y P2 ?z,
        OPTIONAL { ?y P3 ?w, ?w P4 v1 }
}
    
```



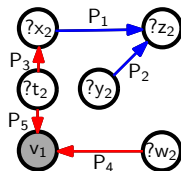
(l) Structure only  $Q_{OPT}$



# Motivating example



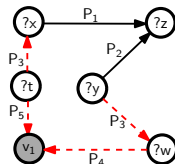
(a) Query  $Q_1$



(b) Query  $Q_2$

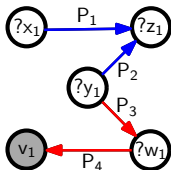
```

SELECT *
WHERE { ?x P1 ?z, ?y P2 ?z,
        OPTIONAL { ?y P3 ?w, ?w P4 v1 }
        OPTIONAL { ?t P3 ?x, ?t P5 v1, ?w P4 v1 }
}
    
```

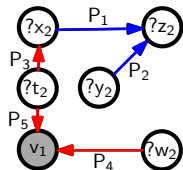


(l) Structure only  $Q_{OPT}$

# Motivating example



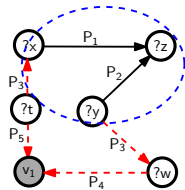
(a) Query  $Q_1$



(b) Query  $Q_2$

```
SELECT *  
WHERE { ?x P1 ?z, ?y P2 ?z,  
        OPTIONAL { ?y P3 ?w, ?w P4 v1 }  
        OPTIONAL { ?t P3 ?x, ?t P5 v1, ?w P4 v1 }  
}
```

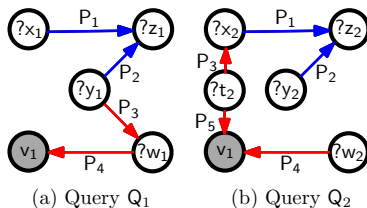
Evaluated once → potential saving



(I) Structure only  $Q_{OPT}$

**OPTIONALS** are evaluated on top of the common substructures  
(intermediate results cached by engine).

# Motivating example

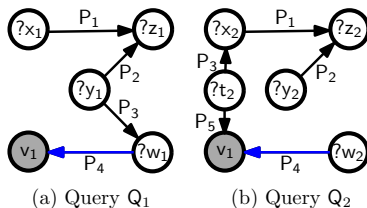


pattern $p$	$\alpha(p)$
$?x P_1 ?z$	30%
$?y P_2 ?z$	20%
$?y P_3 ?w$	18%
$?w P_4 v_1$	1%
$?t P_5 v_1$	2%

\*Max common subquery is not selective

(II) Using cost in optimization

# Motivating example

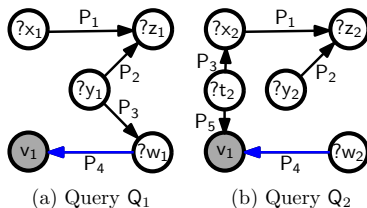


pattern $p$	$\alpha(p)$
$?x P_1 ?z$	30%
$?y P_2 ?z$	20%
$?y P_3 ?w$	18%
$?w P_4 v_1$	1%
$?t P_5 v_1$	2%

\*Max common subquery is not selective

(II) Using cost in optimization

# Motivating example



```
SELECT *  
WHERE { ?w P4 v1,  
        OPTIONAL { ?x1 P1 ?z1, ?y1 P2 ?z1, ?y1 P3 ?w }  
        OPTIONAL { ?x2 P1 ?z2, ?y2 P2 ?z2, ?t2 P3 ?x2, ?t2 P5 v1 }  
}
```

(II) Using cost in optimization

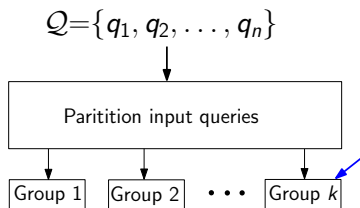
# Our approach

$$\mathcal{Q} = \{q_1, q_2, \dots, q_n\}$$

# Our approach

$\mathcal{Q} = \{q_1, q_2, \dots, q_n\}$  ← They often do not share one common subquery

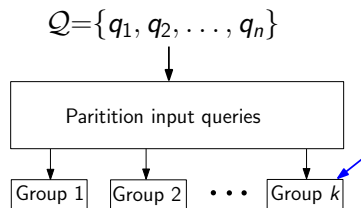
# Our approach



- Similar queries can be optimized together

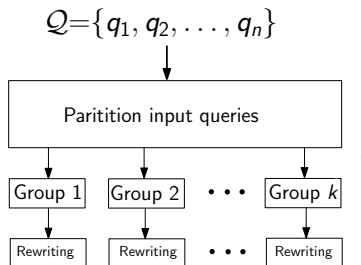


# Our approach

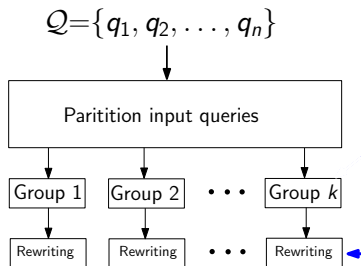


- Similar queries can be optimized together
  - finding structure similarity is expensive
  - group by predicates
  - distance: Jaccard similarity of predicate sets

# Our approach

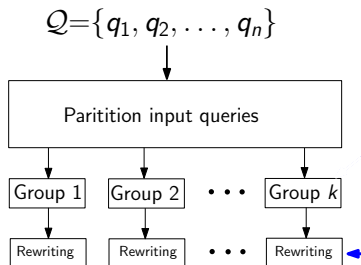


# Our approach



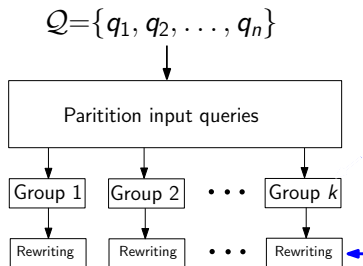
- Recursively rewrite a subset of type 1 queries (hierarchically)  $\rightarrow$  a set of type 2 queries

# Our approach



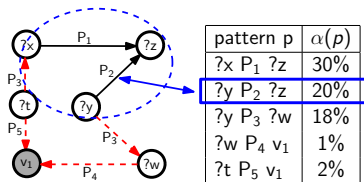
- Recursively rewrite a subset of type 1 queries (hierarchically)  $\rightarrow$  a set of type 2 queries
  - finding common edge subgraphs
  - optimizations to avoid bad efficiency
  - **cost**: guard against bad rewritings
  - approx. by the min selectivity in common subquery

# Our approach

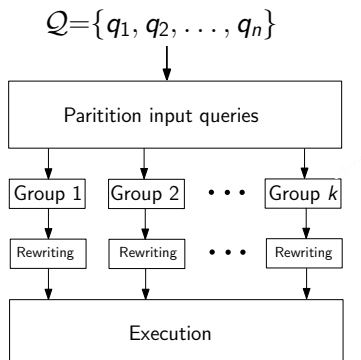


• Recursively rewrite a subset of type 1 queries (hierarchically) → a set of type 2 queries

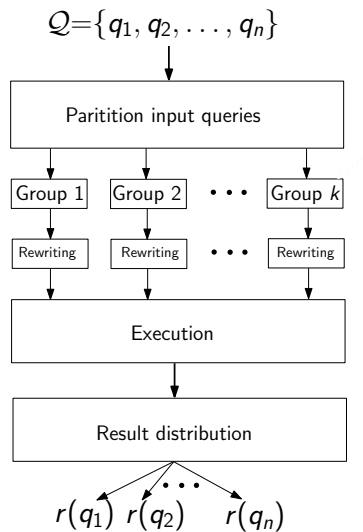
- finding common edge subgraphs
- optimizations to avoid bad efficiency
- **cost**: guard against bad rewritings
- approx. by the min selectivity in common subquery



# Our approach



# Our approach



# Our approach

- Related issues



# Our approach

- Related issues
  - Distributing results, *i.e.*, **Type 2** query  $\rightarrow$  **Type 1** queries

X	Y	Z
a		
b		c
d	e	
f		g

RD of a **Type 1** query: e.g., X and Z

$\updownarrow$

columns from results of the **Type 2** rewriting

# Our approach

- Related issues

- Distributing results, *i.e.*, **Type 2** query  $\rightarrow$  **Type 1** queries

X	Y	Z
a		
b		c
d	e	
f		g

RD of a **Type 1** query: e.g., X and Z

$\updownarrow$

columns from results of the **Type 2** rewriting

- Soundness and completeness

# Our approach

- Related issues
  - Distributing results, *i.e.*, **Type 2** query  $\rightarrow$  **Type 1** queries

X	Y	Z
a		
b		c
d	e	
f		g

RD of a **Type 1** query: e.g., X and Z

$\updownarrow$

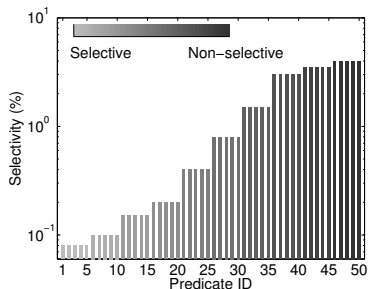
columns from results of the **Type 2** rewriting

- Soundness and completeness
- Extensibility of the solution: more general queries
  - handle variable predicates
  - OPTIONAL queries

- Implementation highlights
  - C++
  - 64-bit Linux, 2GHz Xeon(R) CPU, 4GB memory

# Experiments

- Implementation highlights
  - C++
  - 64-bit Linux, 2GHz Xeon(R) CPU, 4GB memory
- Dataset
  - Extend LUBM benchmark generator:  
randomness in structure, variances of sel.



- Implementation highlights
  - C++
  - 64-bit Linux, 2GHz Xeon(R) CPU, 4GB memory
- Dataset
  - Extend LUBM benchmark generator:  
randomness in structure, variances of sel.
- RDF stores: Jena TDB 0.85 etc

- Implementation highlights
  - C++
  - 64-bit Linux, 2GHz Xeon(R) CPU, 4GB memory
- Dataset
  - Extend LUBM benchmark generator:  
randomness in structure, variances of sel.
- RDF stores: Jena TDB 0.85 etc
- Queries  
Ensure randomness in structure, e.g., star, chain and circle

- Implementation highlights
  - C++
  - 64-bit Linux, 2GHz Xeon(R) CPU, 4GB memory
- Dataset
  - Extend LUBM benchmark generator:  
randomness in structure, variances of sel.
- RDF stores: Jena TDB 0.85 etc
- Queries

Parameter	Symbol	Default	Range
Dataset size	D	4M	3M to 9M
Number of queries	$ Q $	100	60 to 160
Size of query (num of triple patterns)	$ Q $	6	5 to 9
Number of seed queries	$\kappa$	6	5 to 10
Size of seed queries	$ q_{cmn} $	$\sim  Q /2$	1 to 5
Max selectivity of patterns in Q	$\alpha_{max}(Q)$	random	0.1% to 4%
Min selectivity of patterns in Q	$\alpha_{min}(Q)$	1%	0.1% to 4%



- Implementation highlights
  - C++
  - 64-bit Linux, 2GHz Xeon(R) CPU, 4GB memory
- Dataset
  - Extend LUBM benchmark generator:  
randomness in structure, variances of sel.
- RDF stores: Jena TDB 0.85 etc
- Queries

Parameter	Symbol	Default	Range
Dataset size	D	4M	3M to 9M
Number of queries	$ Q $	100	60 to 160
Size of query (num of triple patterns)	$ Q $	6	5 to 9
Number of seed queries	$\kappa$	6	5 to 10
Size of seed queries	$ q_{cmn} $	$\sim  Q /2$	1 to 5
Max selectivity of patterns in Q	$\alpha_{max}(Q)$	random	0.1% to 4%
Min selectivity of patterns in Q	$\alpha_{min}(Q)$	1%	0.1% to 4%

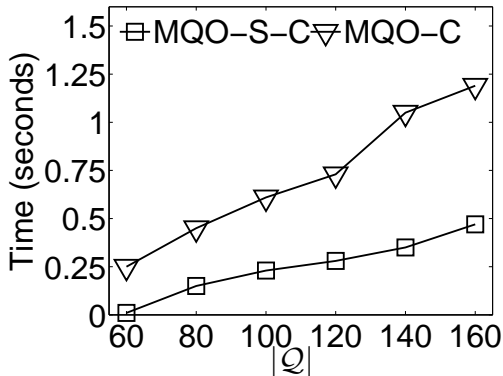
- Rewriting w/ structure: MQO-S; rewriting w/ structure and cost: MQO

# Experiments

- Time on rewriting
  - MQO-S-C: structure based rewriting
  - MQO-C: rewriting integrating with cost

# Experiments

- Time on rewriting
  - MQO-S-C: structure based rewriting
  - MQO-C: rewriting integrating with cost

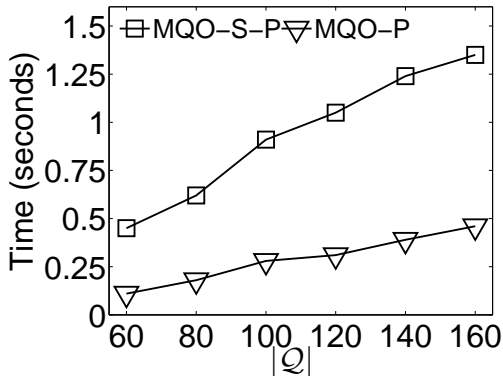


\*Costly/bad rewritings are rejected  
→ more rounds of comparisons.

- Time on distributing results
  - MQO-S-P: parsing results from MQO-S
  - MQO-P: parsing results with MQO

# Experiments

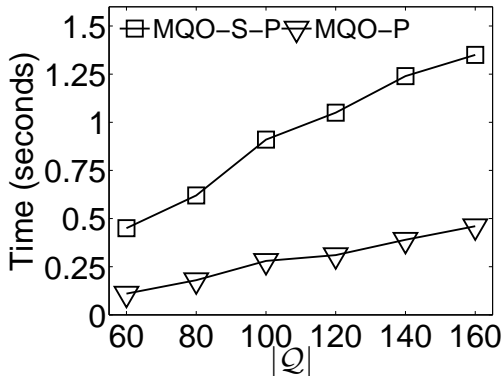
- Time on distributing results  
MQO-S-P: parsing results from MQO-S  
MQO-P: parsing results with MQO



\*Non-selective common subqueries increase the set of results.

# Experiments

- Time on distributing results  
MQO-S-P: parsing results from MQO-S  
MQO-P: parsing results with MQO



\*Non-selective common subqueries increase the set of results.

\*Both **rewriting** and **parsing** are efficiently doable

# Experiments

- Varying num of queries in a batch
  - No-MQO: no optimization
  - MQO-S: optimization based on structural rewriting
  - MQO: integrating cost

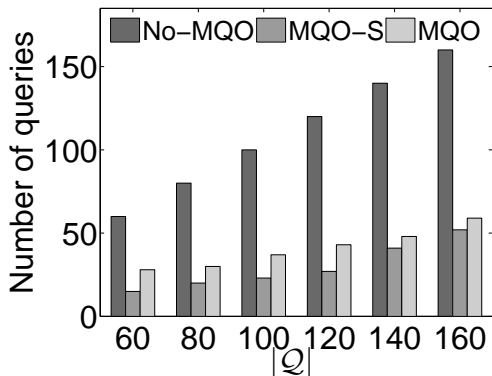
# Experiments

- Varying num of queries in a batch

No-MQO: no optimization

MQO-S: optimization based on structural rewriting

MQO: integrating cost

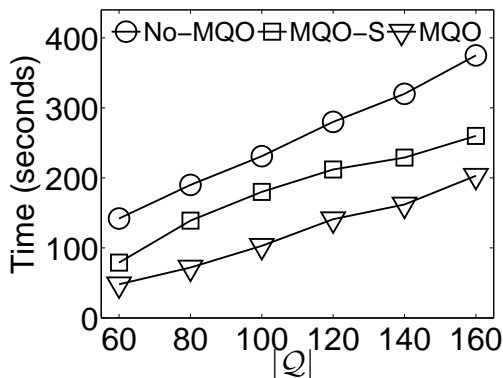


\*Both reduce the num of queries to be executed



# Experiments

- Varying num of queries in a batch
  - No-MQO: no optimization
  - MQO-S: optimization based on structural rewriting
  - MQO: integrating cost



- Varying min. selectivity in seed queries
  - No-MQO: no optimization
  - MQO-S: optimization based on structural rewriting
  - MQO: integrating cost

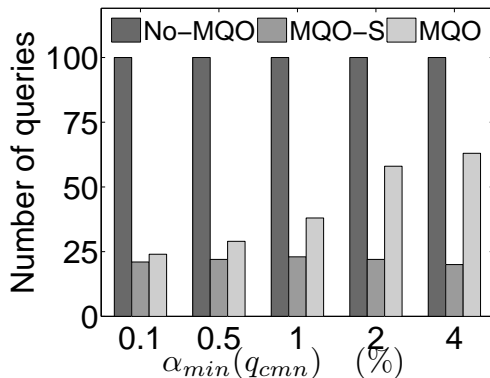
# Experiments

- Varying min. selectivity in seed queries

No-MQO: no optimization

MQO-S: optimization based on structural rewriting

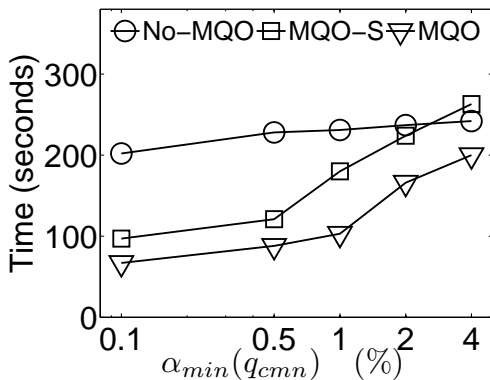
MQO: integrating cost



\*MQO: reject more bad rewritings;  
MQO-S: not sensitive

# Experiments

- Varying min. selectivity in seed queries
  - No-MQO: no optimization
  - MQO-S: optimization based on structural rewriting
  - MQO: integrating cost

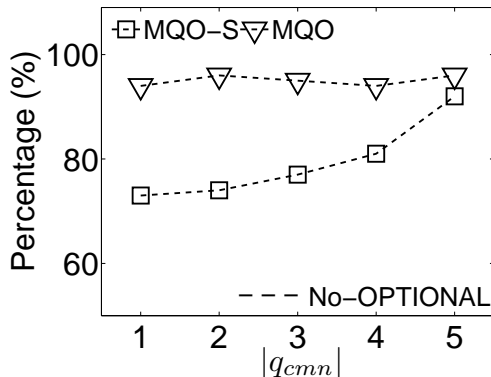


- Varying seed size
    - MQO-S: optimization based on structural rewriting
    - MQO: integrating cost
- percentage =  $\frac{\mathcal{T}_e(\text{common subquery})}{\mathcal{T}_e(Q_{opt})} \times 100\%$

# Experiments

- Varying seed size  
MQO-S: optimization based on structural rewriting  
MQO: integrating cost

$$\text{percentage} = \frac{\mathcal{T}_e(\text{common subquery})}{\mathcal{T}_e(Q_{opt})} \times 100\%$$



\*MQO-S: up to 25% time on optional

- In dealing RDF data on the Web, store independency is important
- Combining SPARQL language and graph algorithms can achieve MQO, *i.e.*, by rewriting queries
- Cost must be taken in consideration during rewriting

Thank You

Q and A



# Our approach

- Partition input queries

# Our approach

- Partition input queries
  - Object: similar queries can be optimized together in rewriting

# Our approach

- Partition input queries
  - Object: similar queries can be optimized together in rewriting

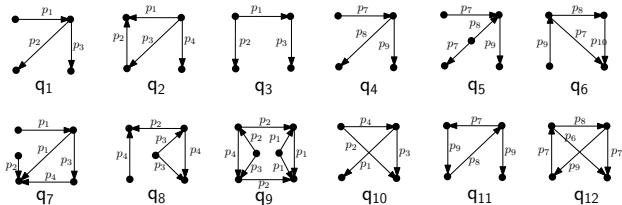
# Our approach

- Partition input queries
  - Object: similar queries can be optimized together in rewriting



# Our approach

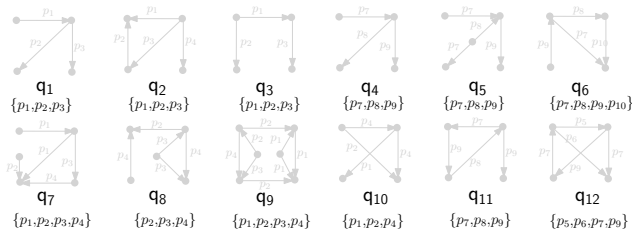
- Partition input queries
  - Object: similar queries can be optimized together in rewriting



Distance: Jaccard similarity on predicates

# Our approach

- Partition input queries
  - Object: similar queries can be optimized together in rewriting

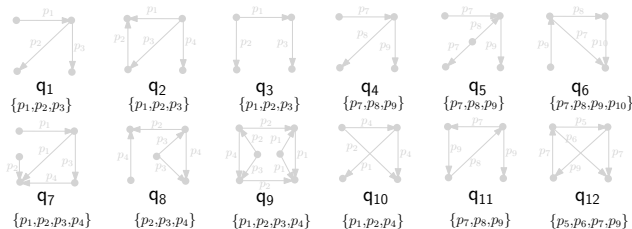


Distance: Jaccard similarity on predicates

- Represent each query as a set of predicates.
- Measure the similarity of a pair of queries by set similarity

# Our approach

- Partition input queries
  - Object: similar queries can be optimized together in rewriting

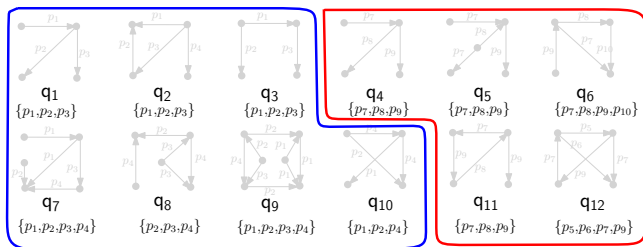


Distance: Jaccard similarity on predicates

- Represent each query as a set of predicates.
- Measure the similarity of a pair of queries by set similarity
- Grouping: k-means

# Our approach

- Partition input queries
  - Object: similar queries can be optimized together in rewriting



Distance: Jaccard similarity on predicates

- Represent each query as a set of predicates.
- Measure the similarity of a pair of queries by set similarity
- Grouping: k-means



# Our approach

- Hierarchical rewriting and clustering (inside a group)

# Our approach

- Hierarchical rewriting and clustering (inside a group)

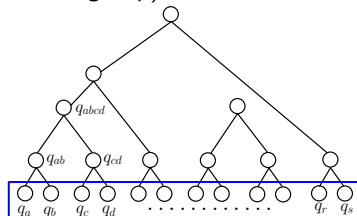
Rewrite pairs of queries bottom up



# Our approach

- Hierarchical rewriting and clustering (inside a group)

Rewrite pairs of queries bottom up

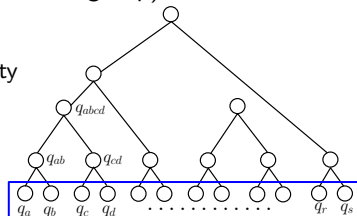


# Our approach

- Hierarchical rewriting and clustering (inside a group)

Rewrite pairs of queries bottom up

Pair up queries with max Jaccard similarity

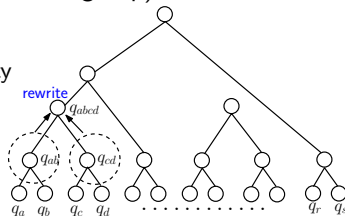


# Our approach

- Hierarchical rewriting and clustering (inside a group)

Rewrite pairs of queries bottom up

Pair up queries with max Jaccard similarity

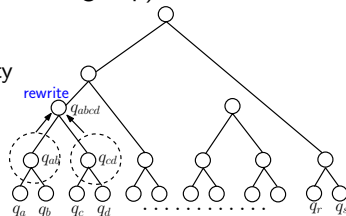


# Our approach

- Hierarchical rewriting and clustering (inside a group)

Rewrite pairs of queries bottom up

Pair up queries with max Jaccard similarity



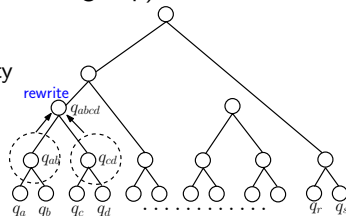
- Rewriting  $\rightarrow$  finding maximal common triple patterns
- In the language of graph ...

# Our approach

- Hierarchical rewriting and clustering (inside a group)

Rewrite pairs of queries bottom up

Pair up queries with max Jaccard similarity



- Rewriting  $\rightarrow$  finding maximal common triple patterns
- In the language of graph . . .

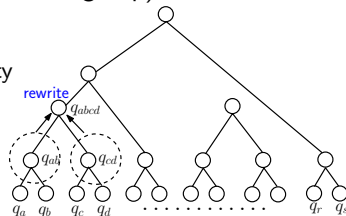
[CLQ01] I. Koch. Enumerating all connected maximal common subgraphs in two graphs. In *Theoretical Computer Science*, 2001.

# Our approach

- Hierarchical rewriting and clustering (inside a group)

Rewrite pairs of queries bottom up

Pair up queries with max Jaccard similarity



- Rewriting  $\rightarrow$  finding maximal common triple patterns
- In the language of graph . . .

[CLQ01] I. Koch. Enumerating all connected maximal common subgraphs in two graphs. In *Theoretical Computer Science*, 2001.

- maximal common connected edge subgraphs

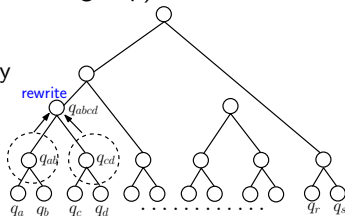


# Our approach

- Hierarchical rewriting and clustering (inside a group)

Rewrite pairs of queries bottom up

Pair up queries with max Jaccard similarity



- Rewriting  $\rightarrow$  finding maximal common triple patterns
- In the language of graph . . .

[CLQ01] I. Koch. Enumerating all connected maximal common subgraphs in two graphs. In *Theoretical Computer Science*, 2001.

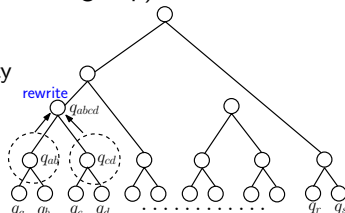
- maximal common connected edge subgraphs  
 $\rightarrow$  maximal common connected *induced* subgraphs in *linegraphs*

# Our approach

- Hierarchical rewriting and clustering (inside a group)

Rewrite pairs of queries bottom up

Pair up queries with max Jaccard similarity

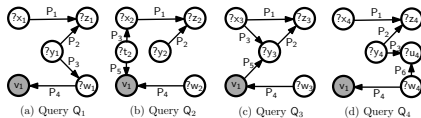


- Rewriting  $\rightarrow$  finding maximal common triple patterns
- In the language of graph . . .

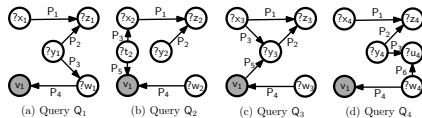
[CLQ01] I. Koch. Enumerating all connected maximal common subgraphs in two graphs. In *Theoretical Computer Science*, 2001.

- maximal common connected edge subgraphs
  - $\rightarrow$  maximal common connected *induced* subgraphs in *linegraphs*
  - $\rightarrow$  maximal cliques in the product graph

# Our approach

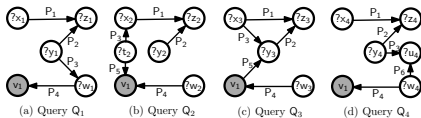


# Our approach

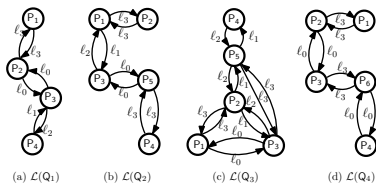


- Linegraph: invert vertices and edges

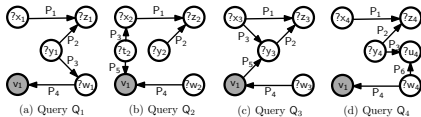
# Our approach



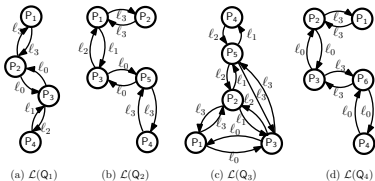
- Linegraph: invert vertices and edges
- sub-sub:  $\ell_0$ , sub-obj:  $\ell_1$ , obj-sub:  $\ell_2$ , obj-obj:  $\ell_3$



# Our approach

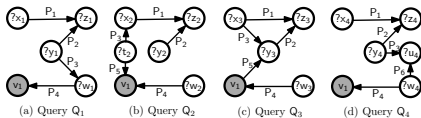


- Linegraph: invert vertices and edges
- sub-sub:  $\ell_0$ , sub-obj:  $\ell_1$ , obj-sub:  $\ell_2$ , obj-obj:  $\ell_3$

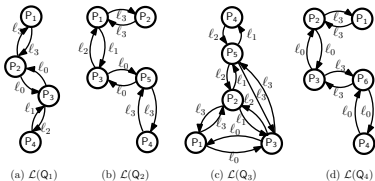


- product graph: simultaneous walk

# Our approach

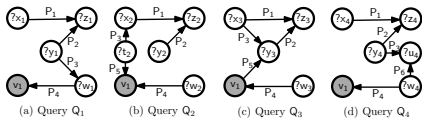


- Linegraph: invert vertices and edges
- sub-sub:  $\ell_0$ , sub-obj:  $\ell_1$ , obj-sub:  $\ell_2$ , obj-obj:  $\ell_3$

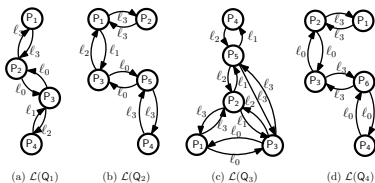


- product graph: simultaneous walk

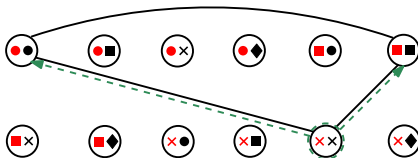
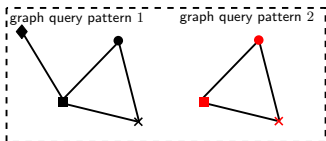
# Our approach



- Linegraph: invert vertices and edges
- sub-sub:  $\ell_0$ , sub-obj:  $\ell_1$ , obj-sub:  $\ell_2$ , obj-obj:  $\ell_3$



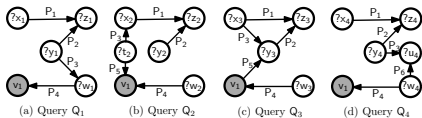
- product graph: simultaneous walk



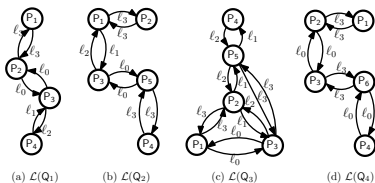
The triangle (clique) highlights the common subgraph composed by ■ × ●



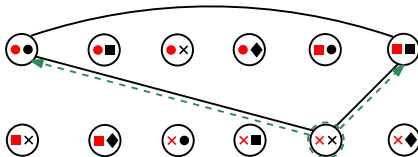
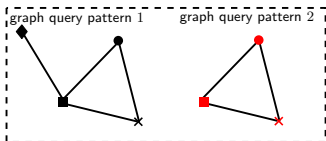
# Our approach



- Linegraph: invert vertices and edges
- sub-sub:  $\ell_0$ , sub-obj:  $\ell_1$ , obj-sub:  $\ell_2$ , obj-obj:  $\ell_3$

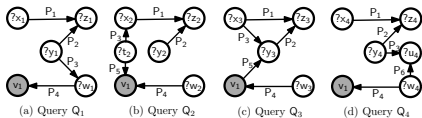


- product graph: simultaneous walk
- blowup in size, esp.  $> 2$  queries affect clique detection

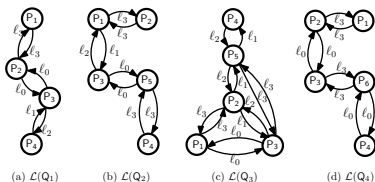


The triangle (clique) highlights the common subgraph composed by ■ × ●

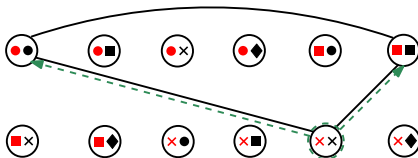
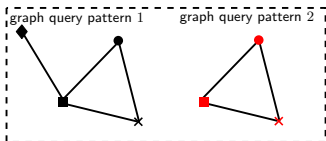
# Our approach



- Linegraph: invert vertices and edges
- sub-sub:  $\ell_0$ , sub-obj:  $\ell_1$ , obj-sub:  $\ell_2$ , obj-obj:  $\ell_3$

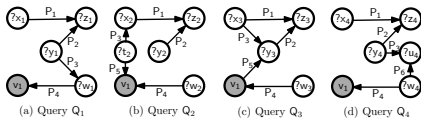


- product graph: simultaneous walk
- blowup in size, esp.  $> 2$  queries affect clique detection
- optimize the product graph

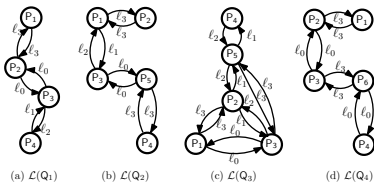


The triangle (clique) highlights the common subgraph composed by ■ × ●

# Our approach



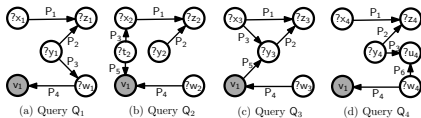
- Linegraph: invert vertices and edges
- sub-sub:  $\ell_0$ , sub-obj:  $\ell_1$ , obj-sub:  $\ell_2$ , obj-obj:  $\ell_3$



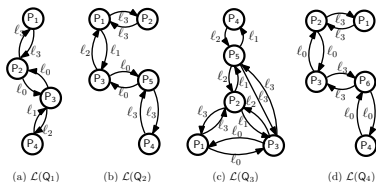
- product graph: simultaneous walk
- blowup in size, esp.  $> 2$  queries affect clique detection
- optimize the product graph

- prune non-common predicates
- check the constants

# Our approach



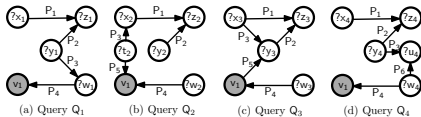
- Linegraph: invert vertices and edges
- sub-sub:  $\ell_0$ , sub-obj:  $\ell_1$ , obj-sub:  $\ell_2$ , obj-obj:  $\ell_3$



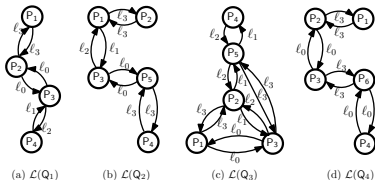
- product graph: simultaneous walk
- blowup in size, esp.  $> 2$  queries affect clique detection
- optimize the product graph

- prune non-common predicates
- check the constants
- prune vertices with non-common neighborhoods

# Our approach

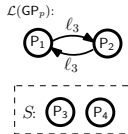


- Linegraph: invert vertices and edges
- sub-sub:  $\ell_0$ , sub-obj:  $\ell_1$ , obj-sub:  $\ell_2$ , obj-obj:  $\ell_3$



- product graph: simultaneous walk
- blowup in size, esp.  $> 2$  queries affect clique detection
- optimize the product graph

- prune non-common predicates
- check the constants
- prune vertices with non-common neighborhoods



- Find maximal cliques in the product graph [CLQ02][CLQ03]

[CLQ02] Patric R.J. Östergård. A fast algorithm for the maximum clique problem. In *Discrete Applied Mathematics*, 2002.

[CLQ03] E. Tomita et al. An efficient branch-and-bound algorithm for finding a maximum clique. In *Discrete Mathematics and Theoretical Computer Science, LNCS*, 2003.

- Find maximal cliques in the product graph [CLQ02][CLQ03]

[CLQ02] Patric R.J. Östergård. A fast algorithm for the maximum clique problem. In *Discrete Applied Mathematics*, 2002.

[CLQ03] E. Tomita et al. An efficient branch-and-bound algorithm for finding a maximum clique. In *Discrete Mathematics and Theoretical Computer Science, LNCS*, 2003.

# Our approach

- Find maximal cliques in the product graph [CLQ02][CLQ03]

[CLQ02] Patric R.J. Östergård. A fast algorithm for the maximum clique problem. In *Discrete Applied Mathematics*, 2002.

[CLQ03] E. Tomita et al. An efficient branch-and-bound algorithm for finding a maximum clique. In *Discrete Mathematics and Theoretical Computer Science, LNCS*, 2003.

- Integrate cost into rewriting



- Find maximal cliques in the product graph [CLQ02][CLQ03]

[CLQ02] Patric R.J. Östergård. A fast algorithm for the maximum clique problem. In *Discrete Applied Mathematics*, 2002.

[CLQ03] E. Tomita et al. An efficient branch-and-bound algorithm for finding a maximum clique. In *Discrete Mathematics and Theoretical Computer Science, LNCS*, 2003.

- Integrate cost into rewriting
  - **Structure**: maximize size of the common subquery in a rewriting
  - Evaluation on cost: guard against bad rewritings

- Find maximal cliques in the product graph [CLQ02][CLQ03]

[CLQ02] Patric R.J. Östergård. A fast algorithm for the maximum clique problem. In *Discrete Applied Mathematics*, 2002.

[CLQ03] E. Tomita et al. An efficient branch-and-bound algorithm for finding a maximum clique. In *Discrete Mathematics and Theoretical Computer Science, LNCS*, 2003.

- Integrate cost into rewriting
  - **Structure**: maximize size of the common subquery in a rewriting
  - Evaluation on cost: guard against bad rewritings
  - Measure: min selectivity in the common subquery for approximation

- Find maximal cliques in the product graph [CLQ02][CLQ03]

[CLQ02] Patric R.J. Östergård. A fast algorithm for the maximum clique problem. In *Discrete Applied Mathematics*, 2002.

[CLQ03] E. Tomita et al. An efficient branch-and-bound algorithm for finding a maximum clique. In *Discrete Mathematics and Theoretical Computer Science, LNCS*, 2003.

- Integrate cost into rewriting
  - **Structure**: maximize size of the common subquery in a rewriting
  - Evaluation on cost: guard against bad rewritings
  - Measure: min selectivity in the common subquery for approximation
  - **Cost**: discard bad rewritings, keep good ones in hierarchical rewriting

- Find maximal cliques in the product graph [CLQ02][CLQ03]

[CLQ02] Patric R.J. Östergård. A fast algorithm for the maximum clique problem. In *Discrete Applied Mathematics*, 2002.

[CLQ03] E. Tomita et al. An efficient branch-and-bound algorithm for finding a maximum clique. In *Discrete Mathematics and Theoretical Computer Science, LNCS*, 2003.

- Integrate cost into rewriting

- **Structure**: maximize size of the common subquery in a rewriting
- Evaluation on cost: guard against bad rewritings
- Measure: min selectivity in the common subquery for approximation
- **Cost**: discard bad rewritings, keep good ones in hierarchical rewriting

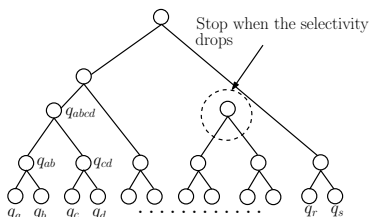
# Our approach

- Find maximal cliques in the product graph [CLQ02][CLQ03]

[CLQ02] Patric R.J. Östergård. A fast algorithm for the maximum clique problem. In *Discrete Applied Mathematics*, 2002.

[CLQ03] E. Tomita et al. An efficient branch-and-bound algorithm for finding a maximum clique. In *Discrete Mathematics and Theoretical Computer Science, LNCS*, 2003.

- Integrate cost into rewriting
  - Structure**: maximize size of the common subquery in a rewriting
  - Evaluation on cost**: guard against bad rewritings
  - Measure**: min selectivity in the common subquery for approximation
  - Cost**: discard bad rewritings, keep good ones in hierarchical rewriting



# Our approach

- Related issues

# Our approach

- Related issues
  - Distributing results, *i.e.*, **Type 2** query  $\rightarrow$  **Type 1** queries

name	mail	hpage
"Alice"	alice@home	http://home/alice
"Alice"	alice@work	http://home/alice
"Bob"		
"Ella"		http://work/ella

RD of a **Type 1** query



columns from results of the **Type 2** rewriting

# Our approach

- Related issues

- Distributing results, *i.e.*, **Type 2** query  $\rightarrow$  **Type 1** queries

name	mail	hpage
"Alice"	alice@home	http://home/alice
"Alice"	alice@work	http://home/alice
"Bob"		
"Ella"		http://work/ella

RD of a **Type 1** query



columns from results of the **Type 2** rewriting

- Soundness and completeness



- Related issues

- Distributing results, *i.e.*, **Type 2** query  $\rightarrow$  **Type 1** queries

name	mail	hpage
"Alice"	alice@home	http://home/alice
"Alice"	alice@work	http://home/alice
"Bob"		
"Ella"		http://work/ella

RD of a **Type 1** query



columns from results of the **Type 2** rewriting

- Soundness and completeness
- Extensibility of the solution: more general queries
  - handle variable predicates
  - nested OPTIONALs