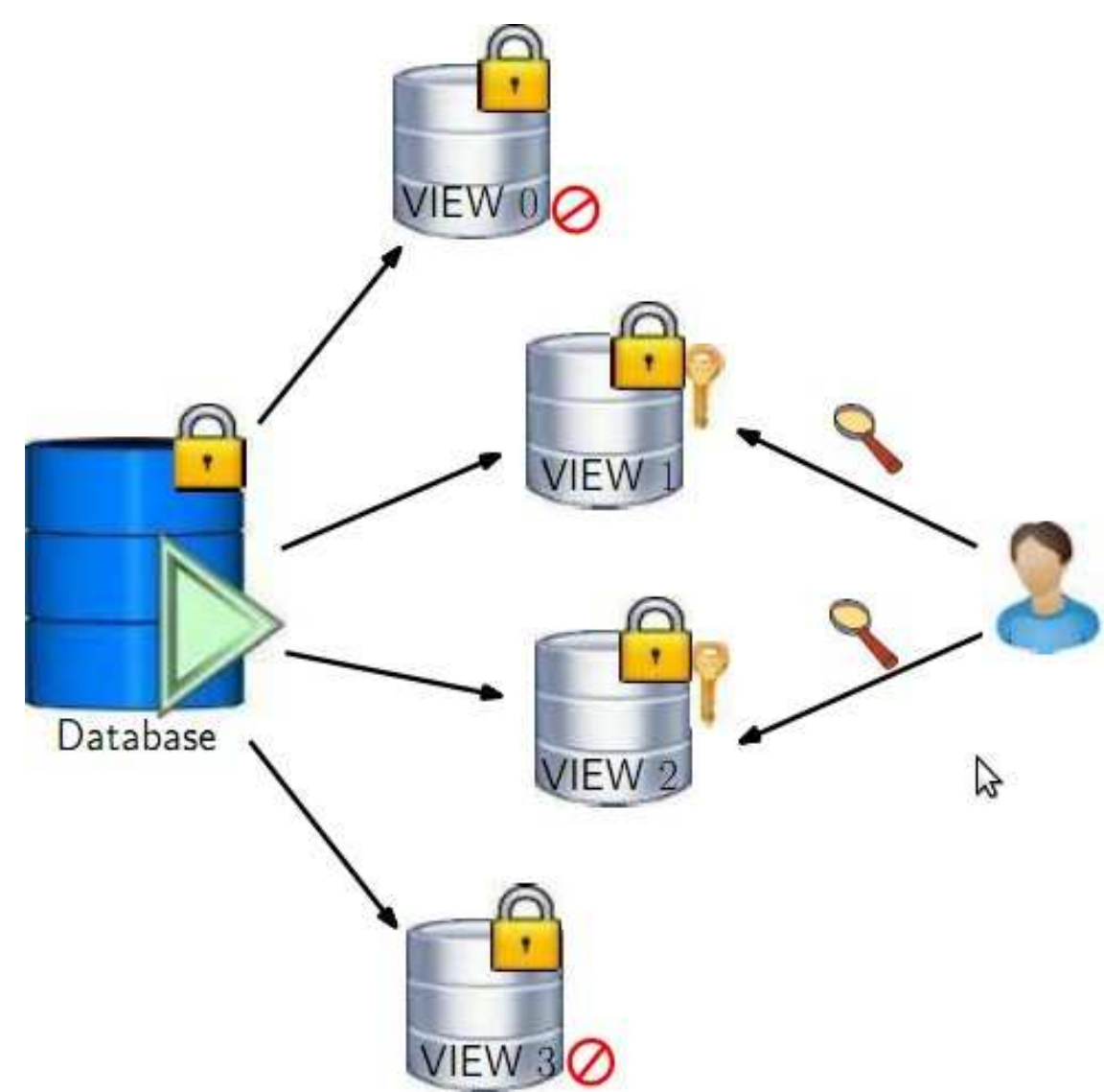


Scalable Multi-Query Optimization for SPARQL

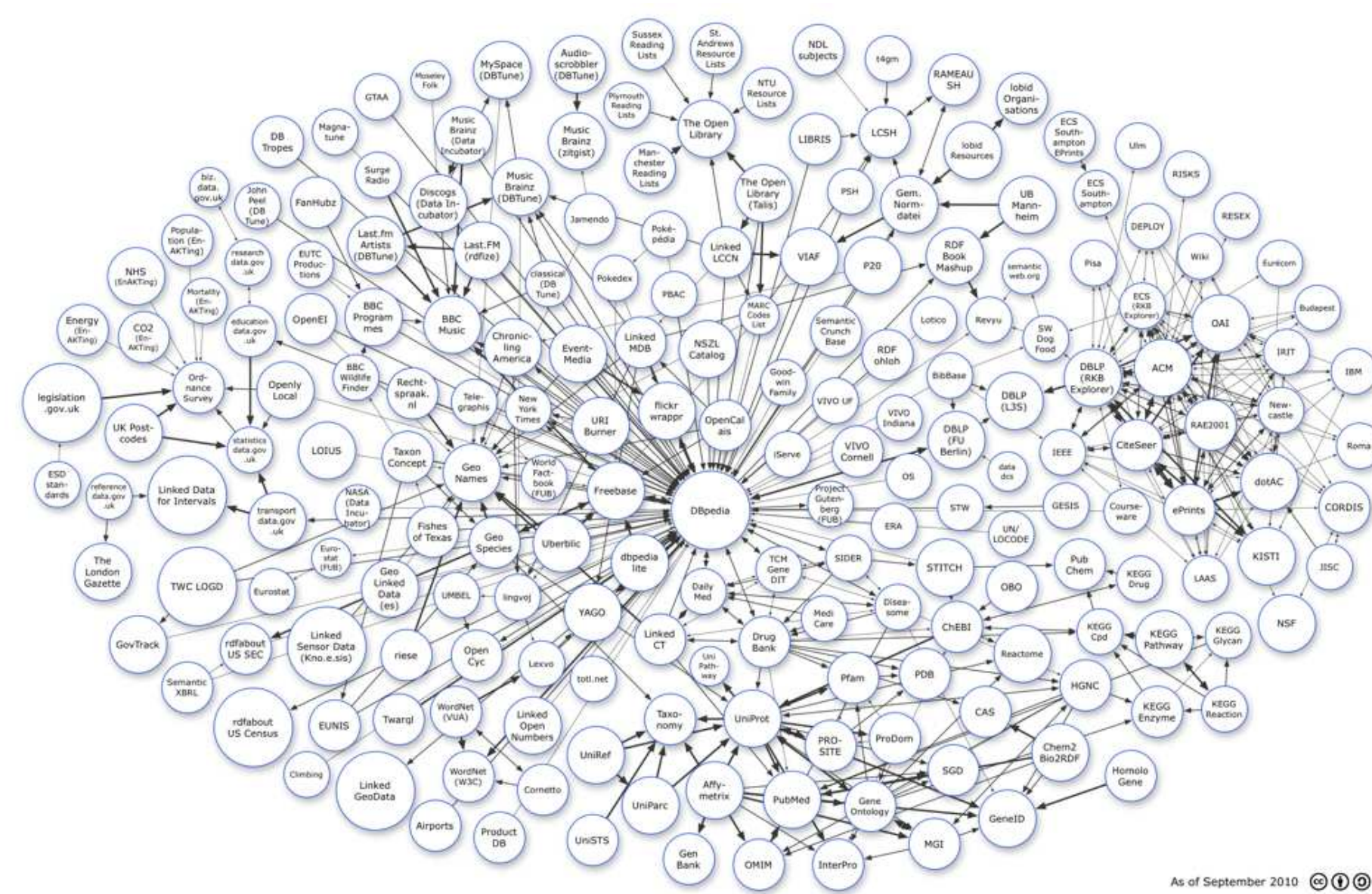
Wangchao Le, Anastasios Kementsietsidis, Songyun Duan, Feifei Li



Motivation 1: Access Control for RDF data



Motivation 2: Web Data Integration on Query Endpoints



Finding Maximal Common Connected Substructures for SPARQL MQO

Theorem 1 Given two graphs, finding the maximal common connected subgraphs amounts to finding the maximal common connected induced subgraphs in their linegraphs.

Theorem 2 Given two graphs, finding the maximal common connected induced subgraphs amounts to finding the maximal cliques with strong covering trees in their product graph.

***Challenges:** (I) deal with hundreds of graphs in one shot; (II) blend selectivity into the structure-based MQO.

Two Types of SPARQL Queries for RDF Data

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

*GP: encoded by conjunctive triples: (sub pred obj).

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

* GP_{OPT}^+ : left-join in relational database.

An Example to Use Left-Join in SPARQL

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"

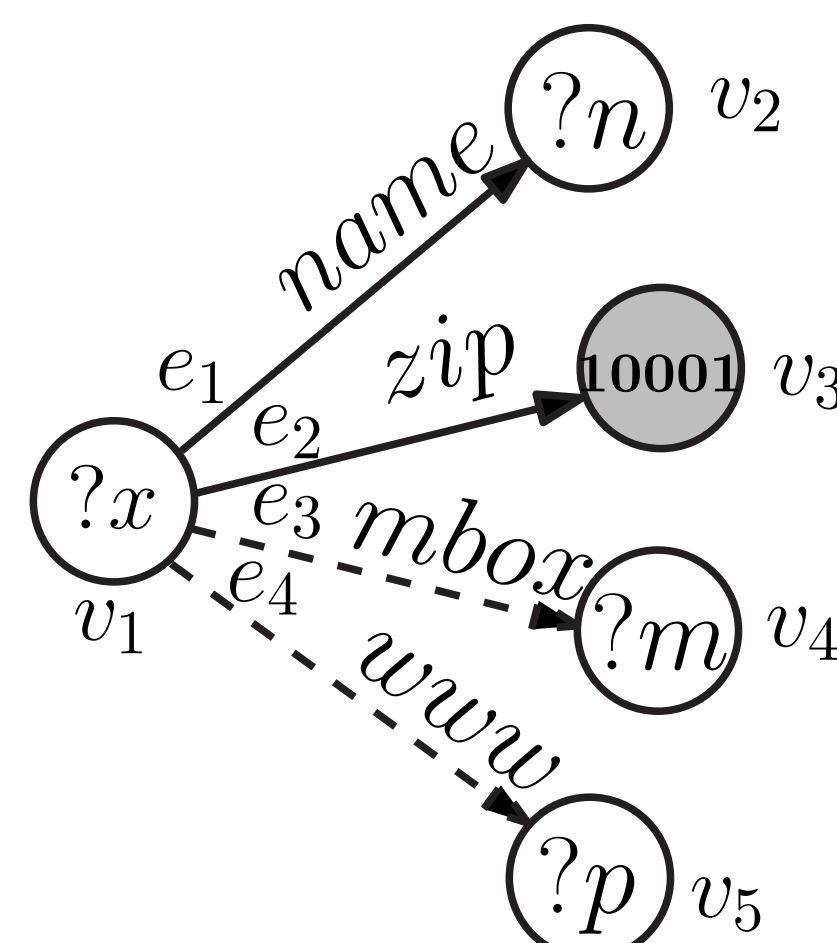
(a) Input data D

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

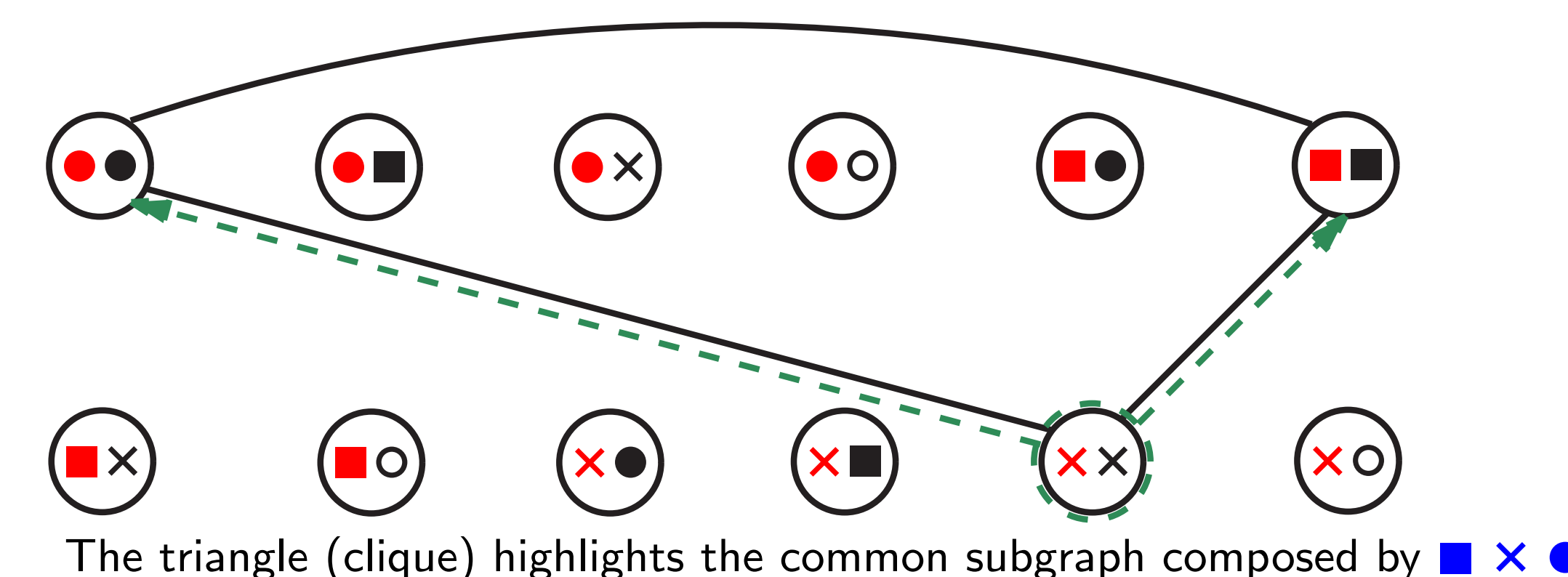
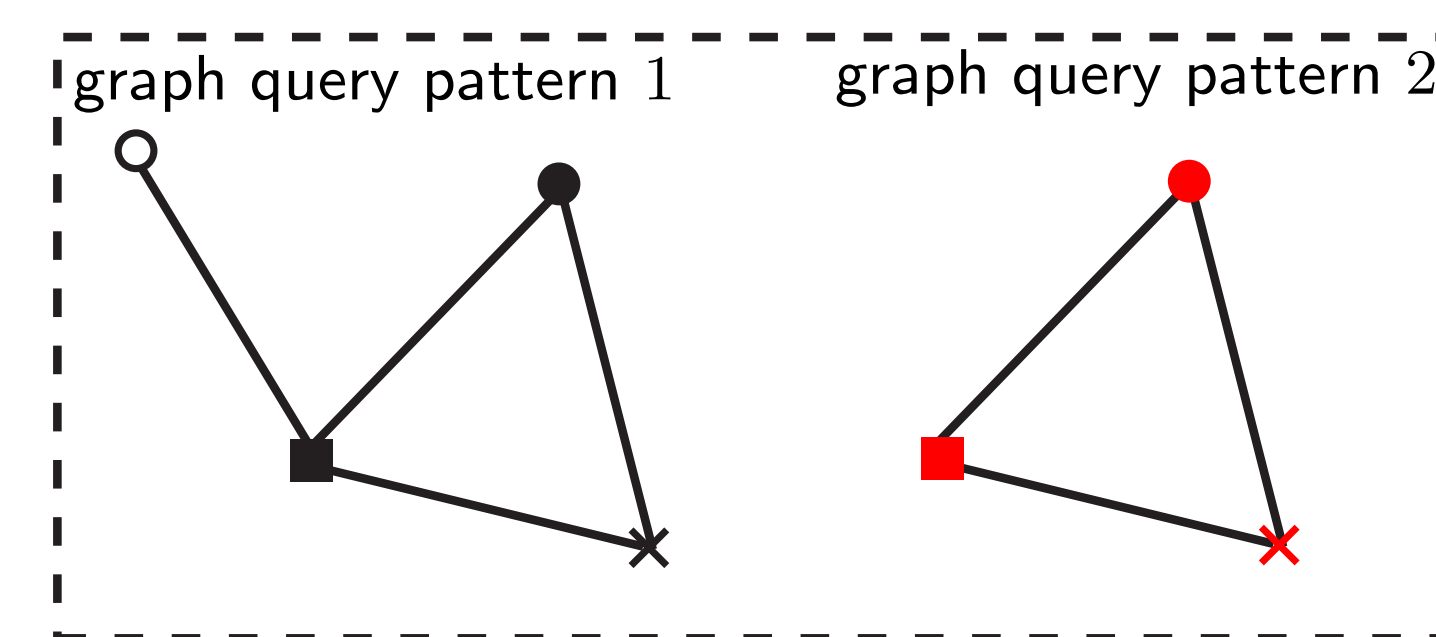
(b) Example query Q_{OPT}

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

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



Clique, Maximal Common Induced Subgraph and Strong Covering Tree

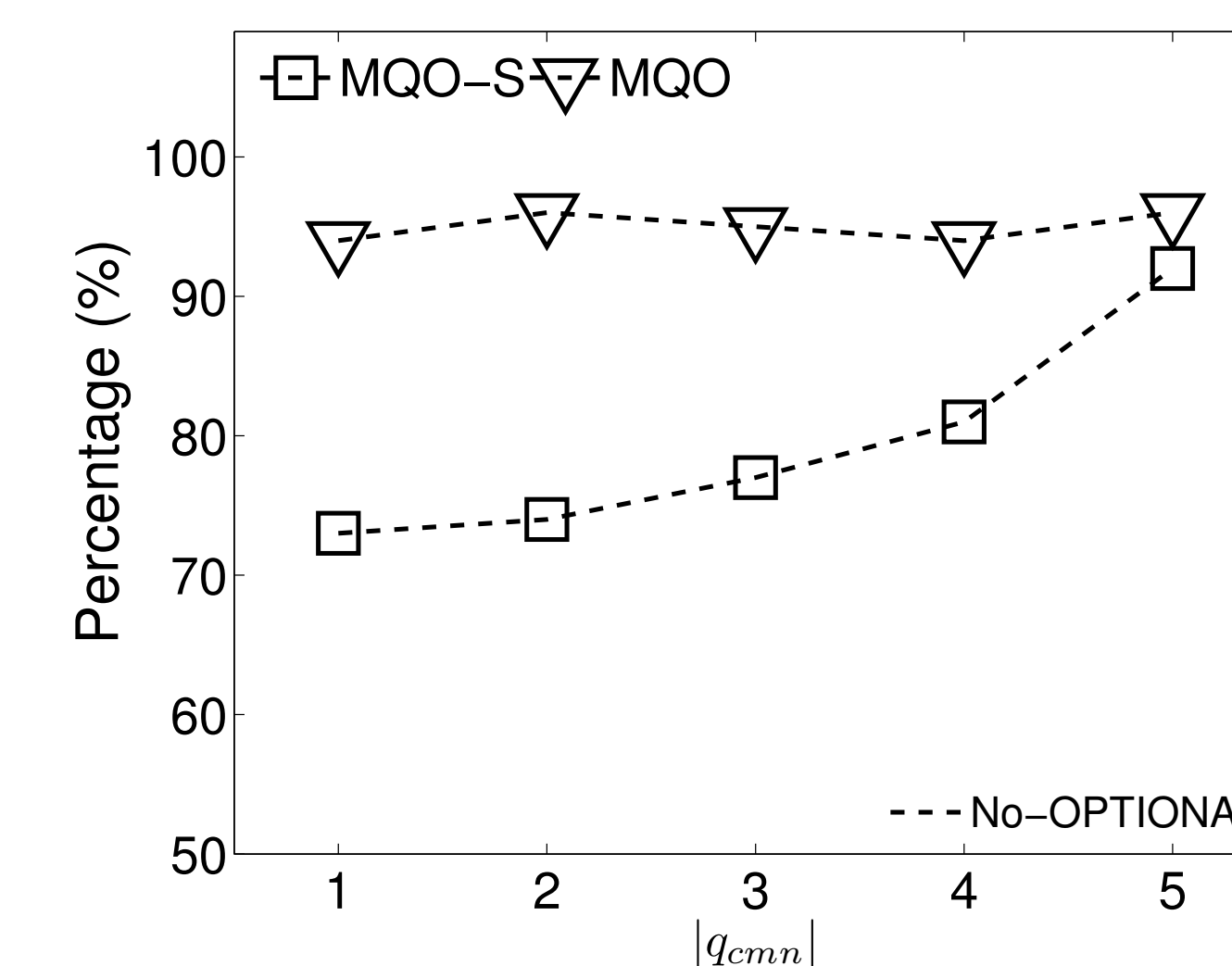


The triangle (clique) highlights the common subgraph composed by ■ × ●

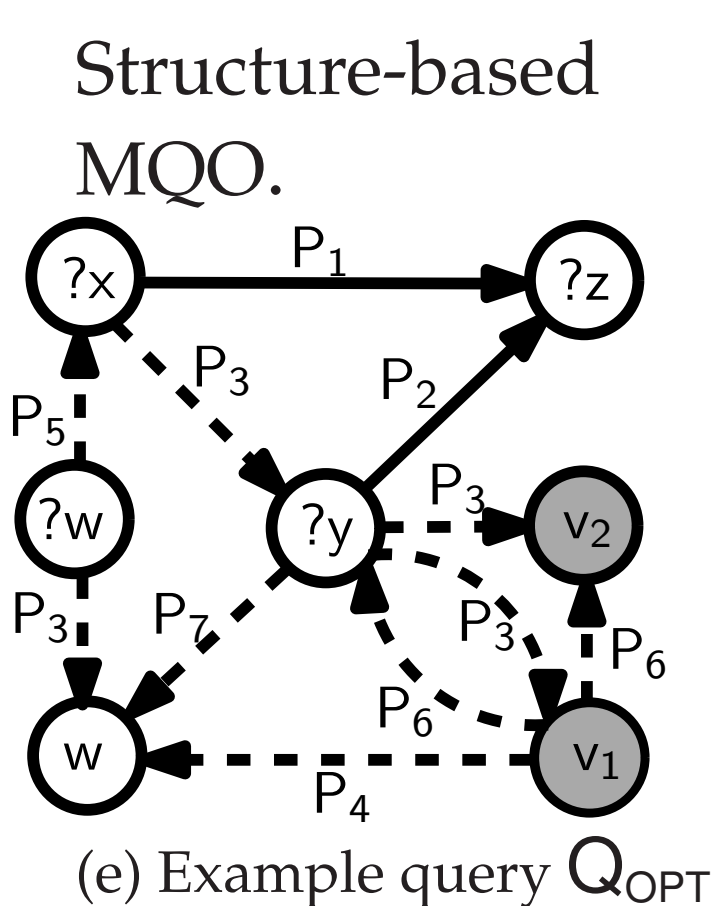
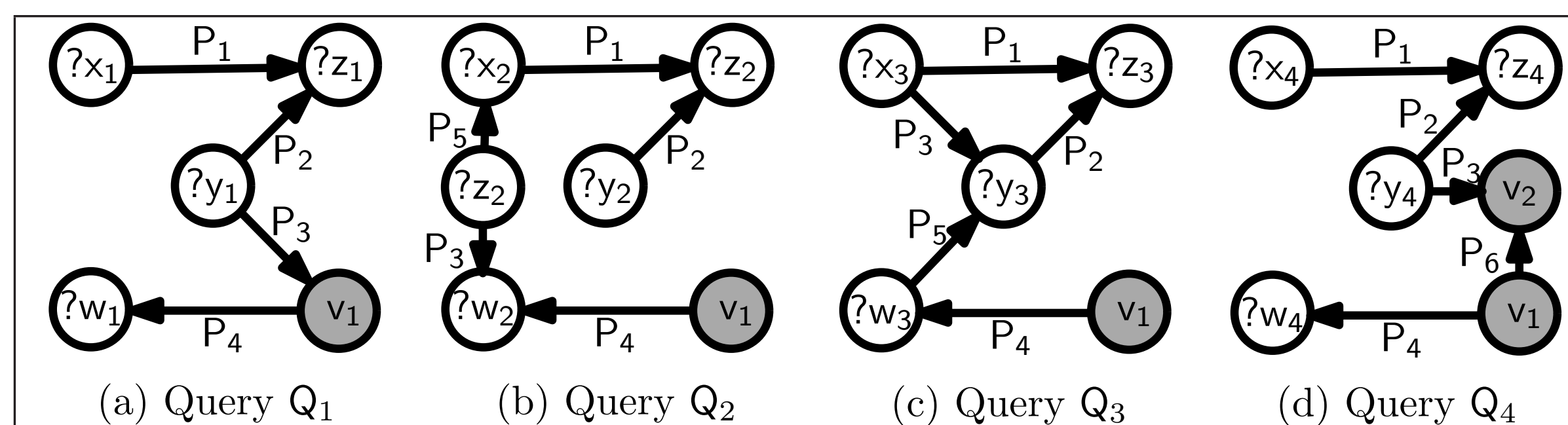
Blend Cost into MQO

$$\text{Cost}(Q) = \begin{cases} \text{Min}(\text{sel}(t)) & Q \text{ is a Type 1 query, } t \in GP \\ \text{Min}(\text{sel}(t)) + \Delta & Q \text{ is a Type 2 query, } t \in GP \end{cases}$$

Observation: >90% of query evaluation time for our MQO is on evaluating the common structure (do it once for all queries), resulting in less time in evaluating the non-common substructures; while the pure structure-based MQO (MQO-S) is sensitive to the variances of common substructures rewritten, leading to more overhead in evaluating the non-common substructures.



An Example to Evison The Optimization



(e) Example query Q_{OPT}

pattern p	$\alpha(p)$
?x P1 ?z	15%
?y P2 ?z	9%
?y P3 ?w	18%
?w P4 v1	4%
?t P5 v1	2%
v1 P5 ?t	7%
?w P6 ?u	13%

(f) Selectivity

Rewrite 4 queries in 1 SPARQL query.

```
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 }
        OPTIONAL { ?x P3 ?y, v1 P5 ?y, ?w P4 v1 }
        OPTIONAL { ?y P3 ?u, ?w P6 ?u, ?w P4 v1 } }
```

Blending selectivity in rewriting.

```
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 }
        OPTIONAL { ?x3 P1 ?z3, ?y3 P2 ?z3, ?x3 P3 ?y3, v1 P5 ?y3 }
        OPTIONAL { ?x4 P1 ?z4, ?y4 P2 ?z4, ?y4 P3 ?u4, ?w P6 ?u4 } }
```

Experiments on varying selectivity

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

MQO-S: MQO based on structure; MQO: Cost-based MQO; No-MQO: No MQO.

