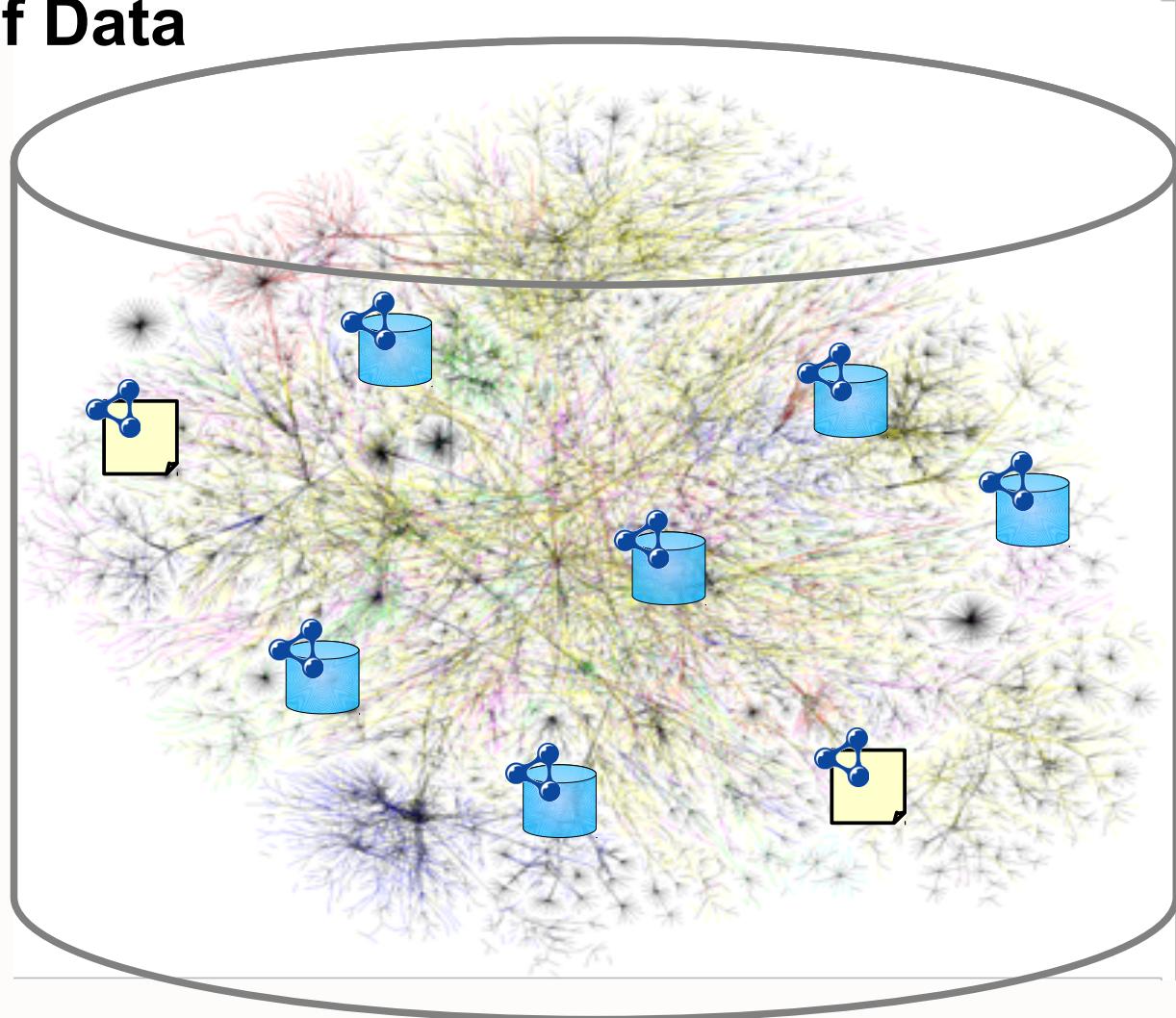
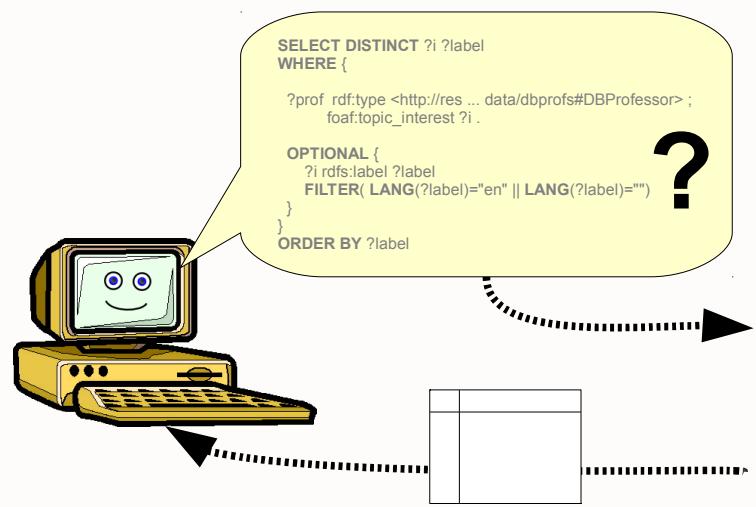


# How Caching Improves Efficiency and Result Completeness for Querying Linked Data

**Olaf Hartig**  
<http://olafhartig.de/foaf.rdf#olaf>  
[@olafhartig](mailto:@olafhartig)

**Database and Information Systems Research Group**  
**Humboldt-Universität zu Berlin**

# Can we query the Web of Data as of it were a single, giant database?

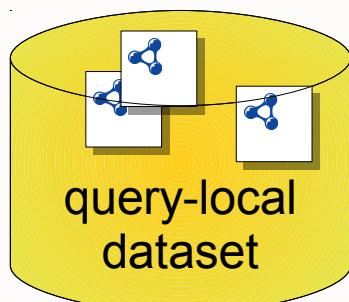


## Our approach: Link Traversal Based Query Execution

[ISWC'09]

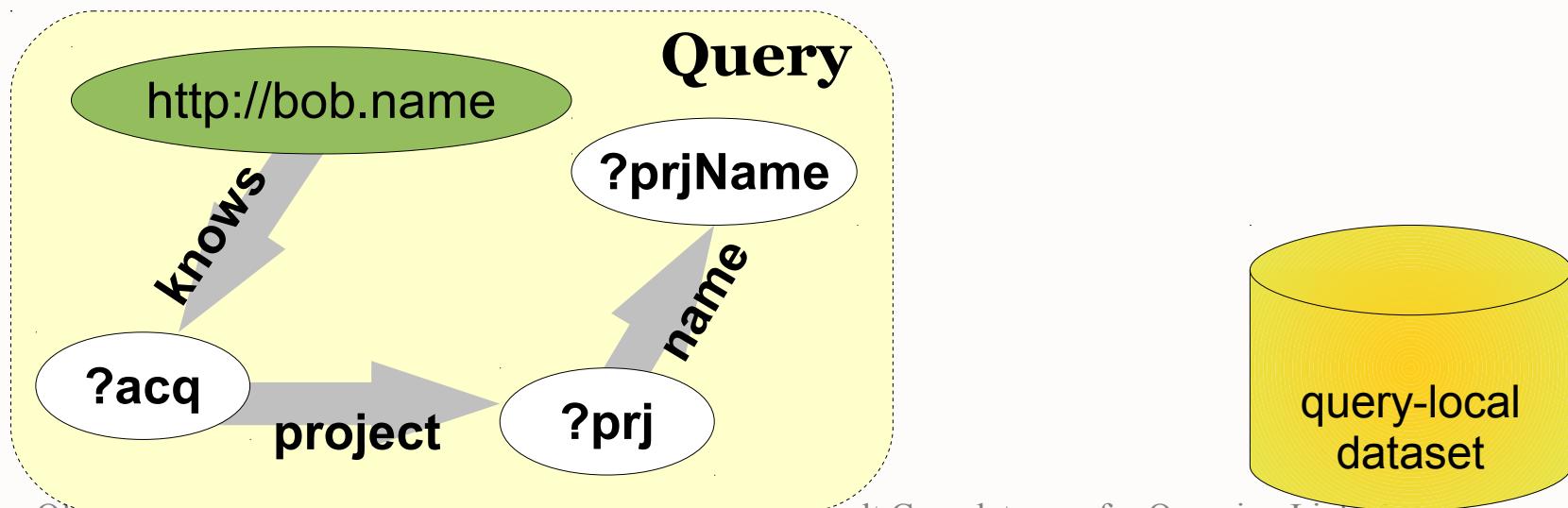
# Main Idea

- **Intertwine query evaluation with traversal of data links**
- **We alternate between:**
  - Evaluate parts of the query (triple patterns) on a continuously augmented set of data
  - Look up URIs in intermediate solutions and add retrieved data to the query-local dataset



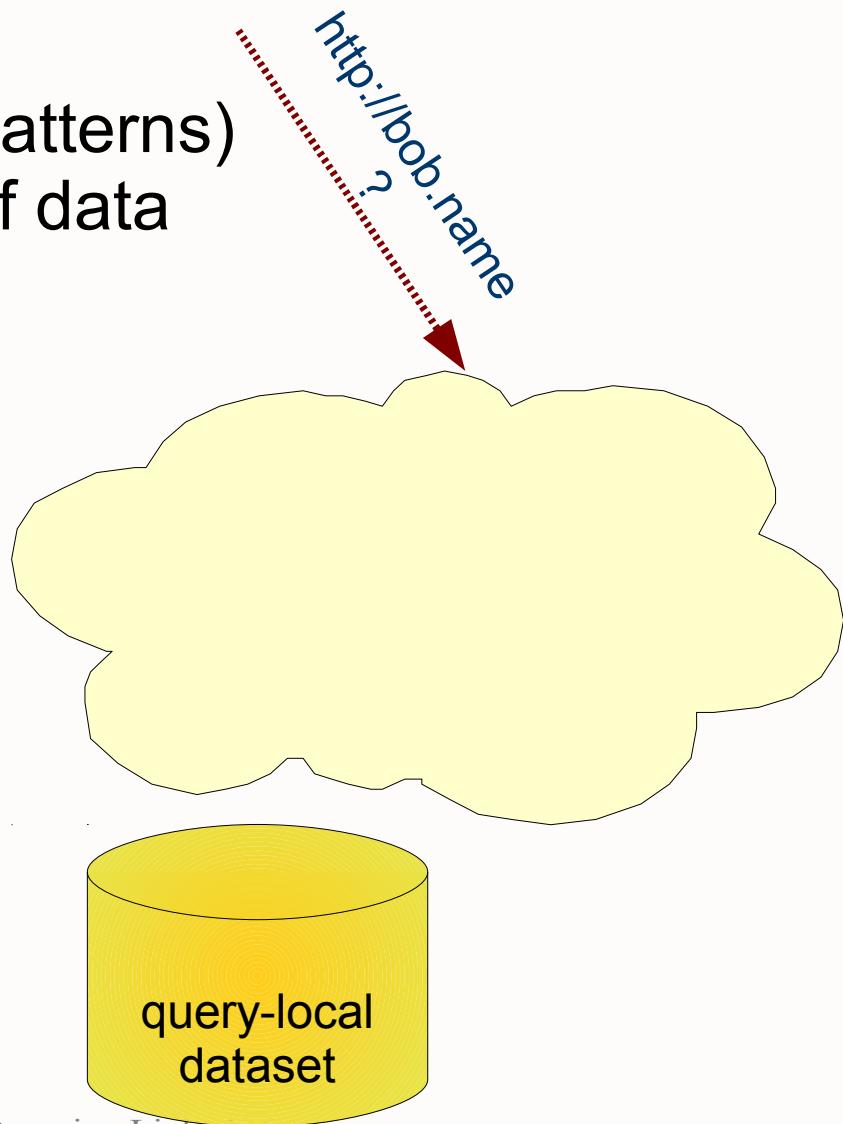
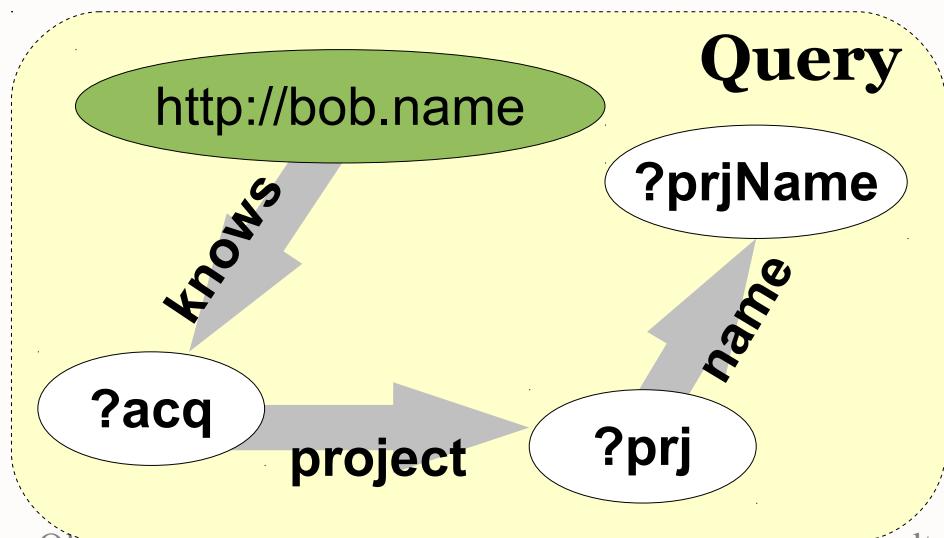
# Main Idea

- Intertwine query evaluation with traversal of data links
- We alternate between:
  - Evaluate parts of the query (triple patterns) on a continuously augmented set of data
  - Look up URIs in intermediate solutions and add retrieved data to the query-local dataset



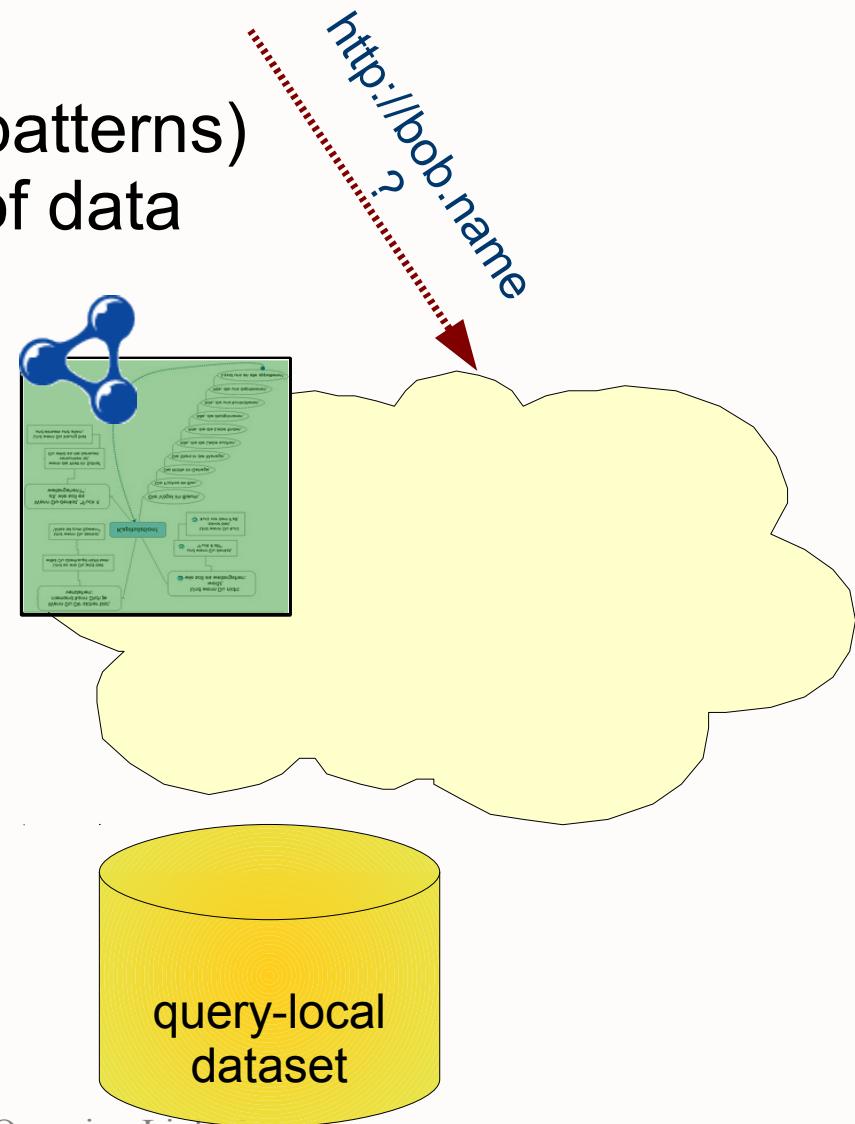
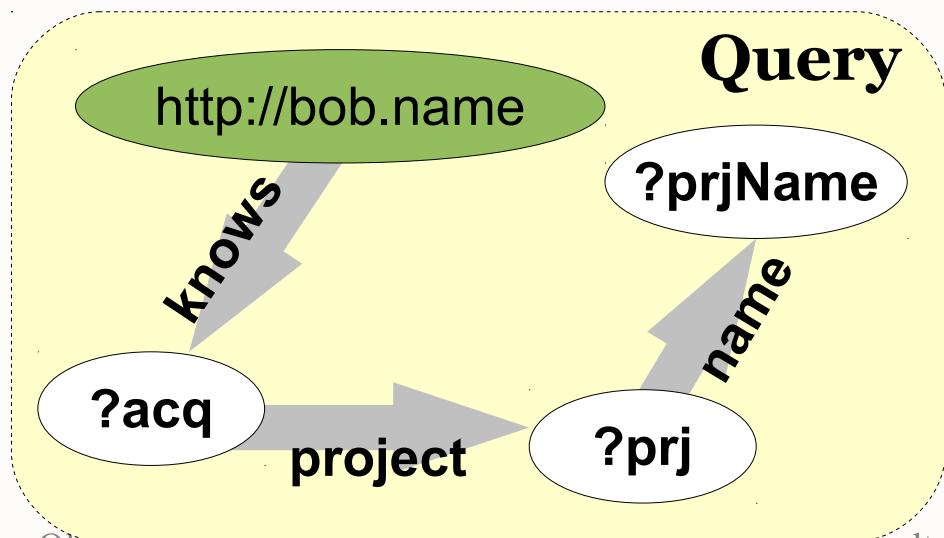
# Main Idea

- Intertwine query evaluation with traversal of data links
- We alternate between:
  - Evaluate parts of the query (triple patterns) on a continuously augmented set of data
  - Look up URIs in intermediate solutions and add retrieved data to the query-local dataset



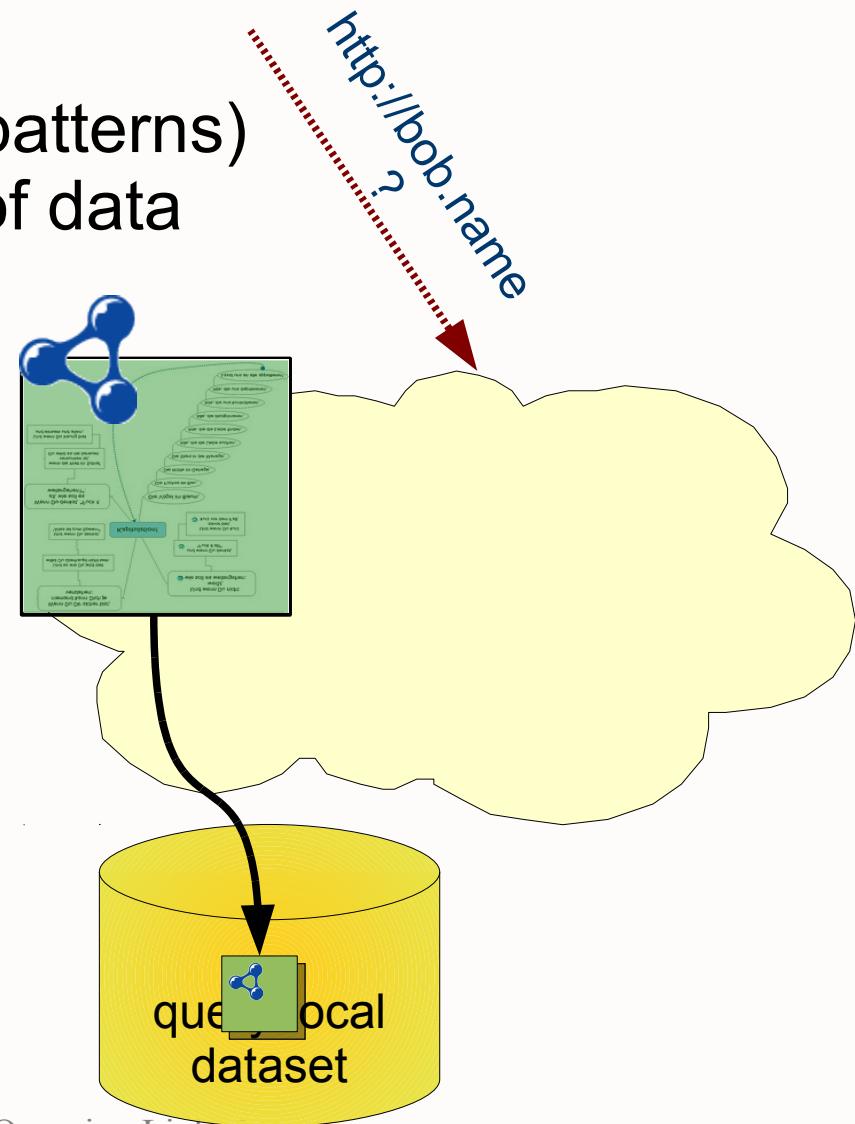
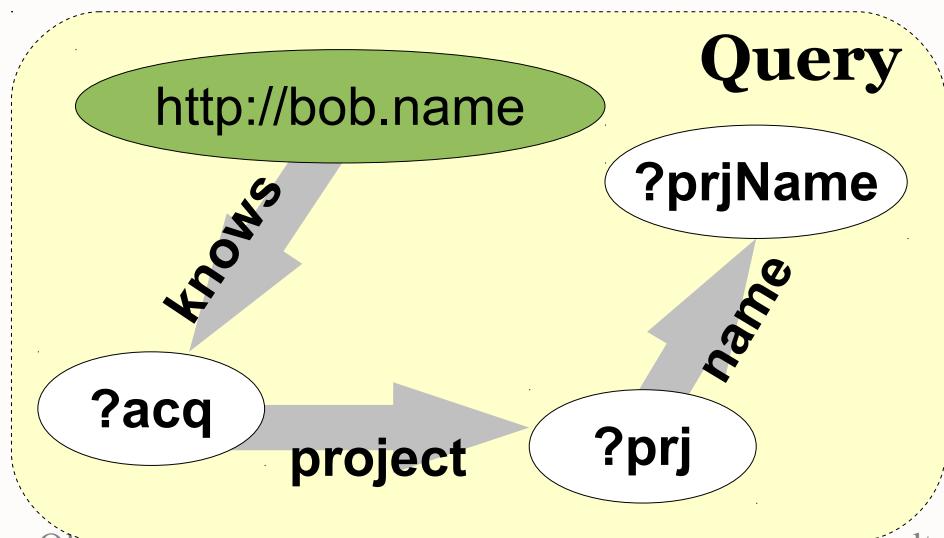
# Main Idea

- Intertwine query evaluation with traversal of data links
- We alternate between:
  - Evaluate parts of the query (triple patterns) on a continuously augmented set of data
  - Look up URIs in intermediate solutions and add retrieved data to the query-local dataset



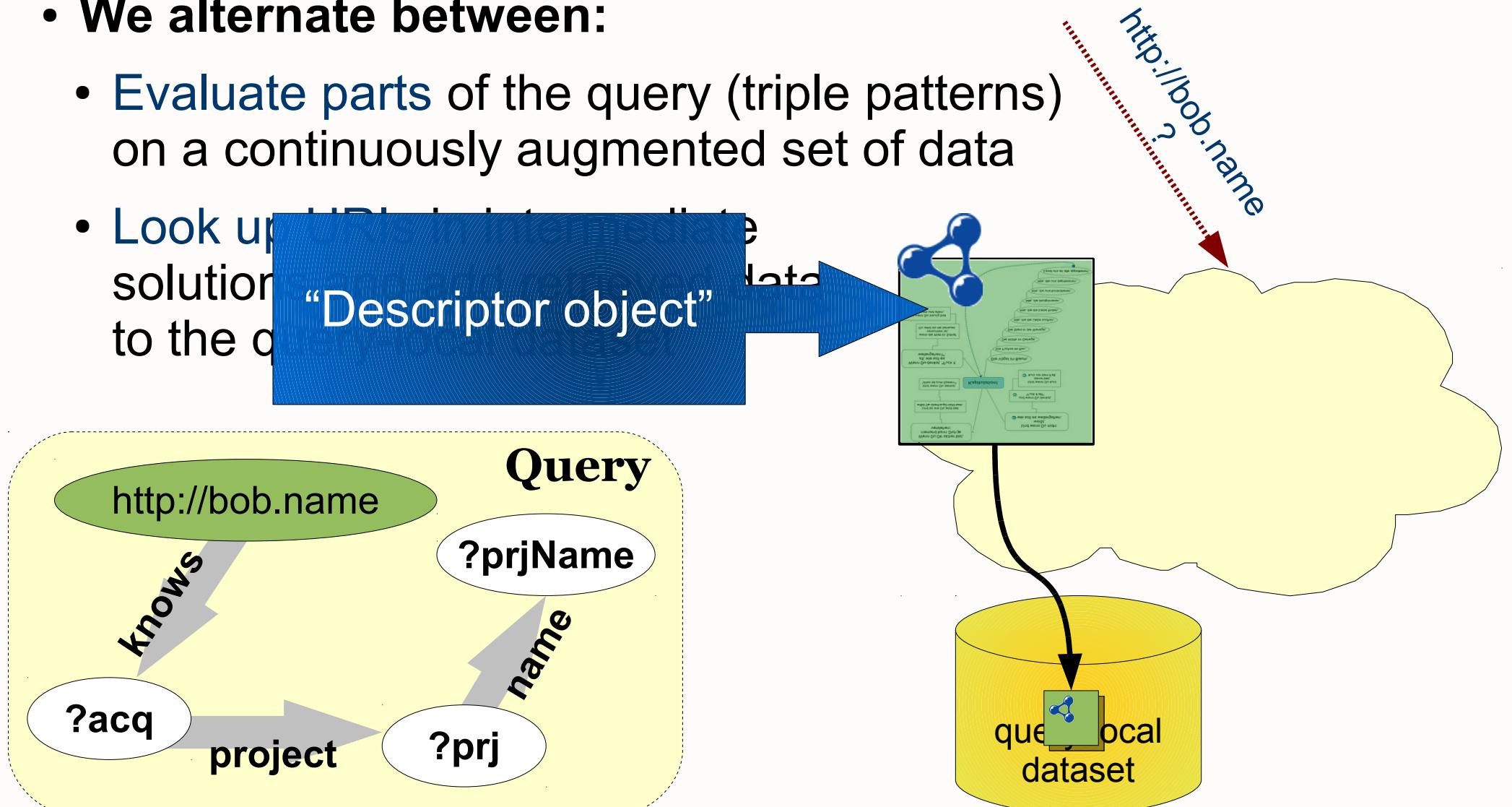
# Main Idea

- Intertwine query evaluation with traversal of data links
- We alternate between:
  - Evaluate parts of the query (triple patterns) on a continuously augmented set of data
  - Look up URIs in intermediate solutions and add retrieved data to the query-local dataset



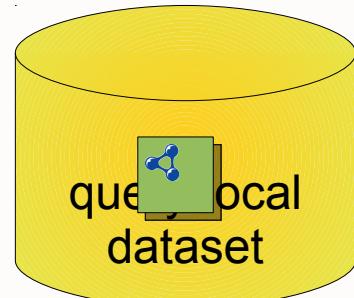
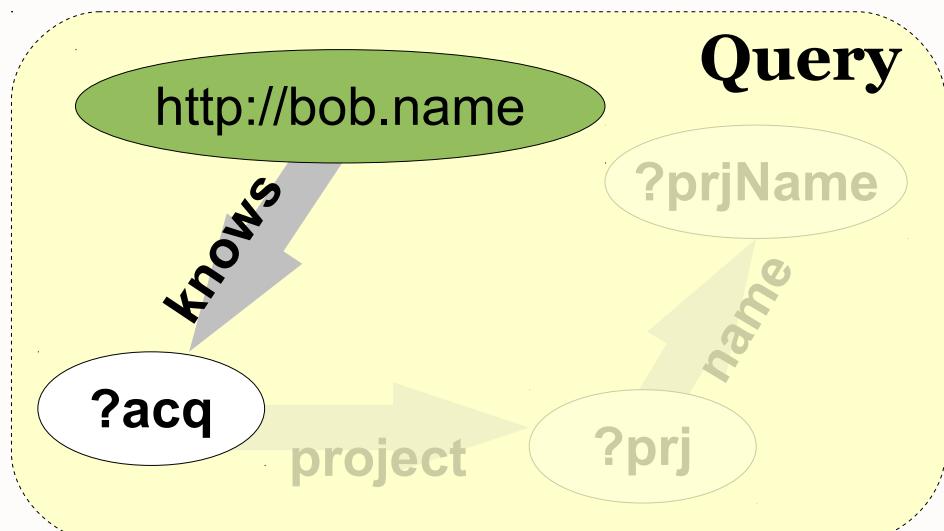
# Main Idea

- Intertwine query evaluation with traversal of data links
- We alternate between:
  - Evaluate parts of the query (triple patterns) on a continuously augmented set of data
  - Look up ~~HTTP~~ triple patterns in the ~~data~~ solution space to the query



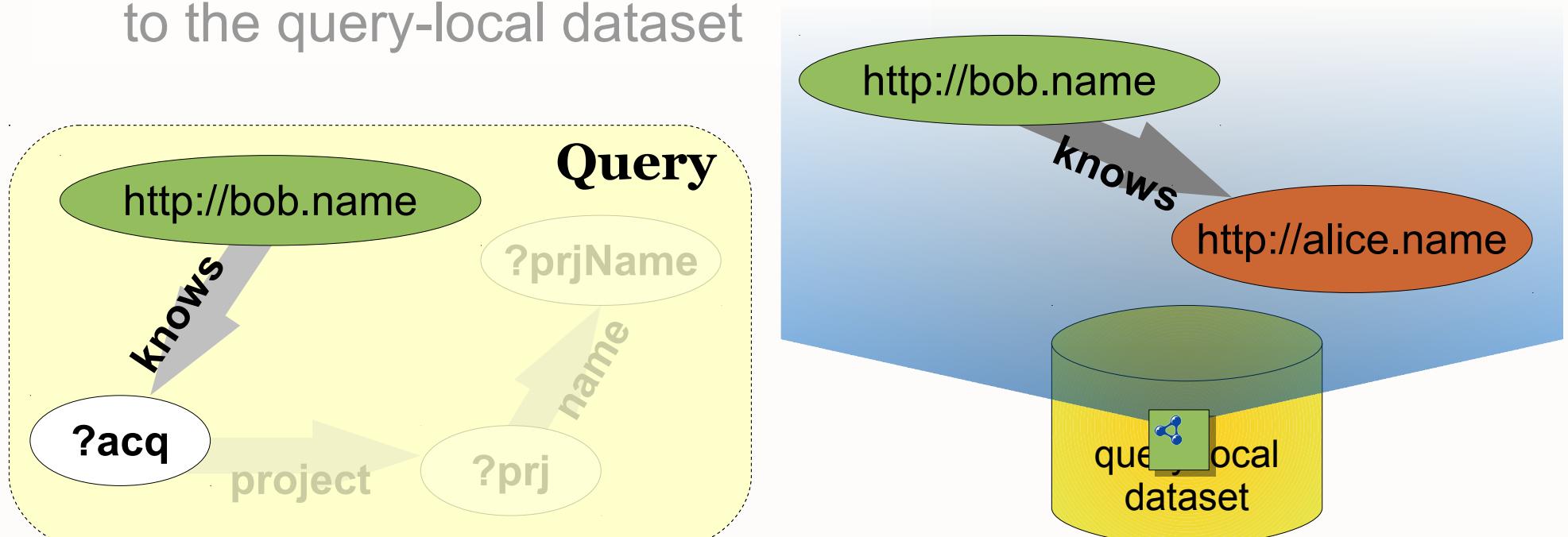
# Main Idea

- Intertwine query evaluation with traversal of data links
- We alternate between:
  - Evaluate parts of the query (triple patterns) on a continuously augmented set of data
  - Look up URIs in intermediate solutions and add retrieved data to the query-local dataset



# Main Idea

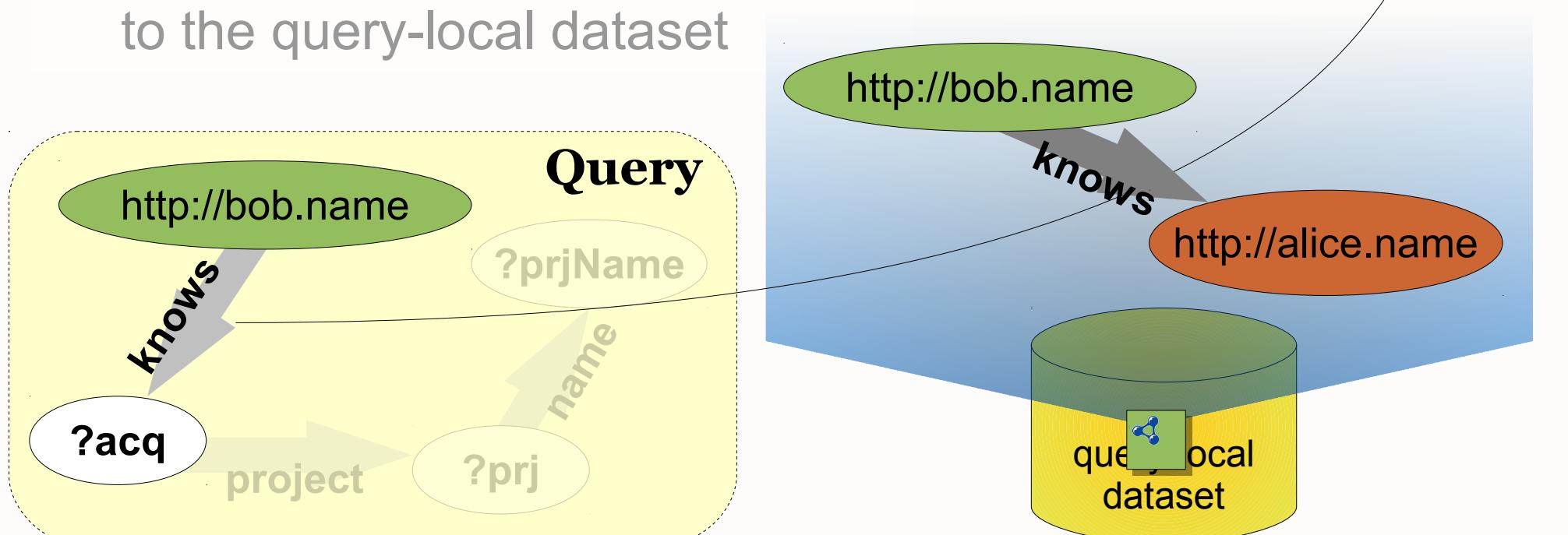
- **Intertwine query evaluation with traversal of data links**
- **We alternate between:**
  - Evaluate parts of the query (triple patterns) on a continuously augmented set of data
  - Look up URIs in intermediate solutions and add retrieved data to the query-local dataset



# Main Idea



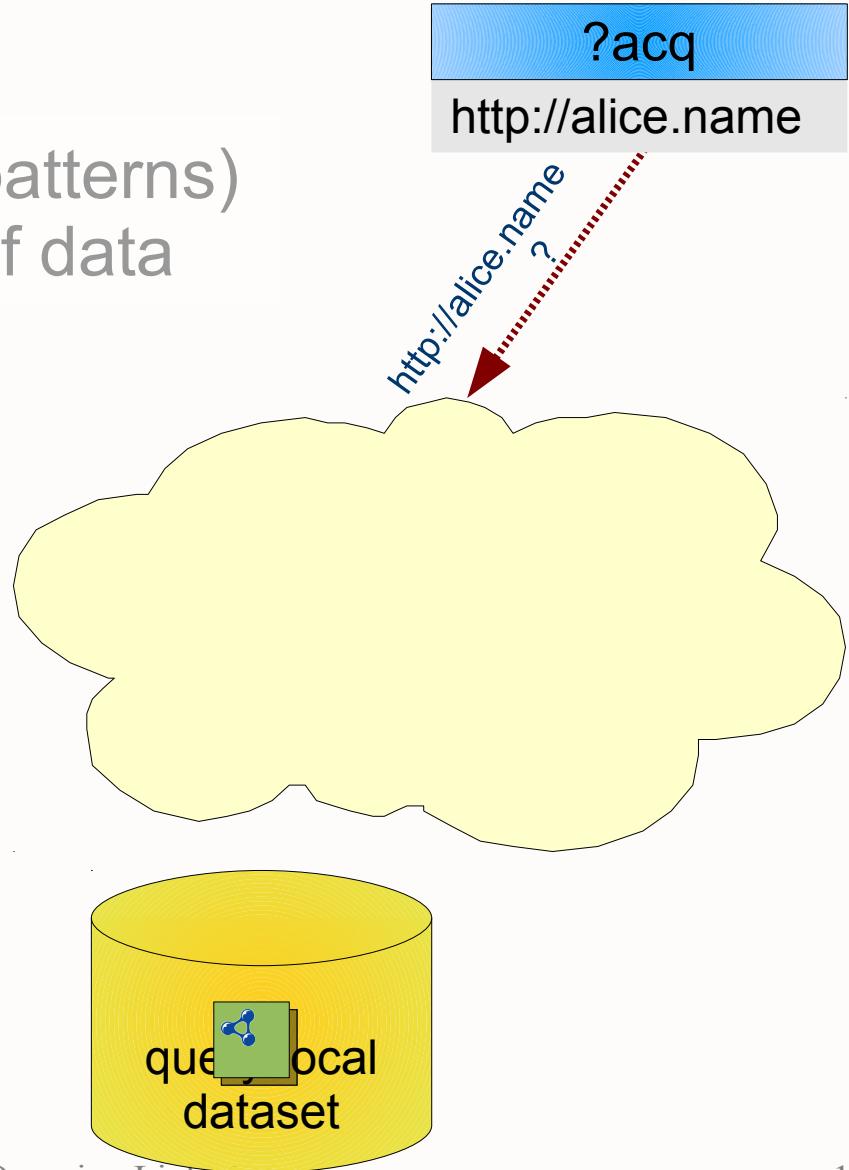
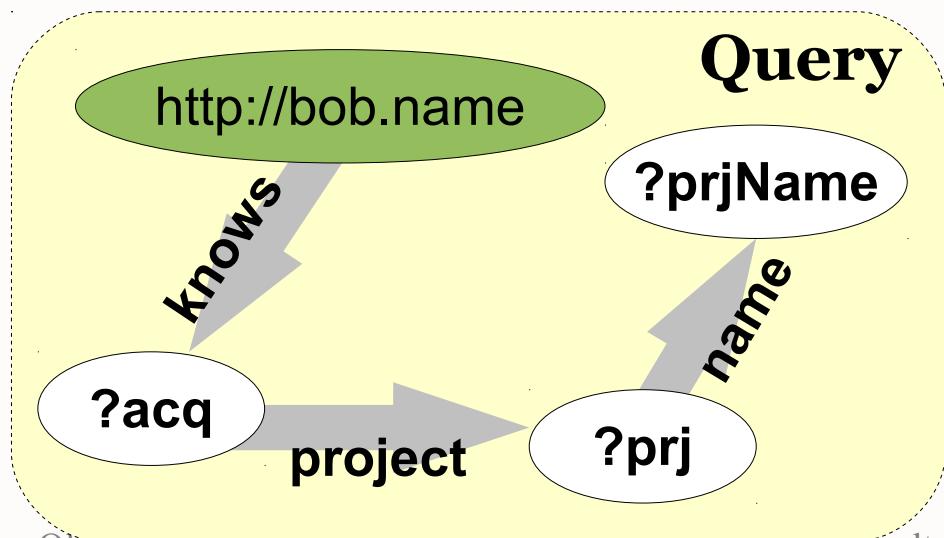
- **Intertwine query evaluation with traversal of data links**
- **We alternate between:**
  - Evaluate parts of the query (triple patterns) on a continuously augmented set of data
  - Look up URIs in intermediate solutions and add retrieved data to the query-local dataset



# Main Idea

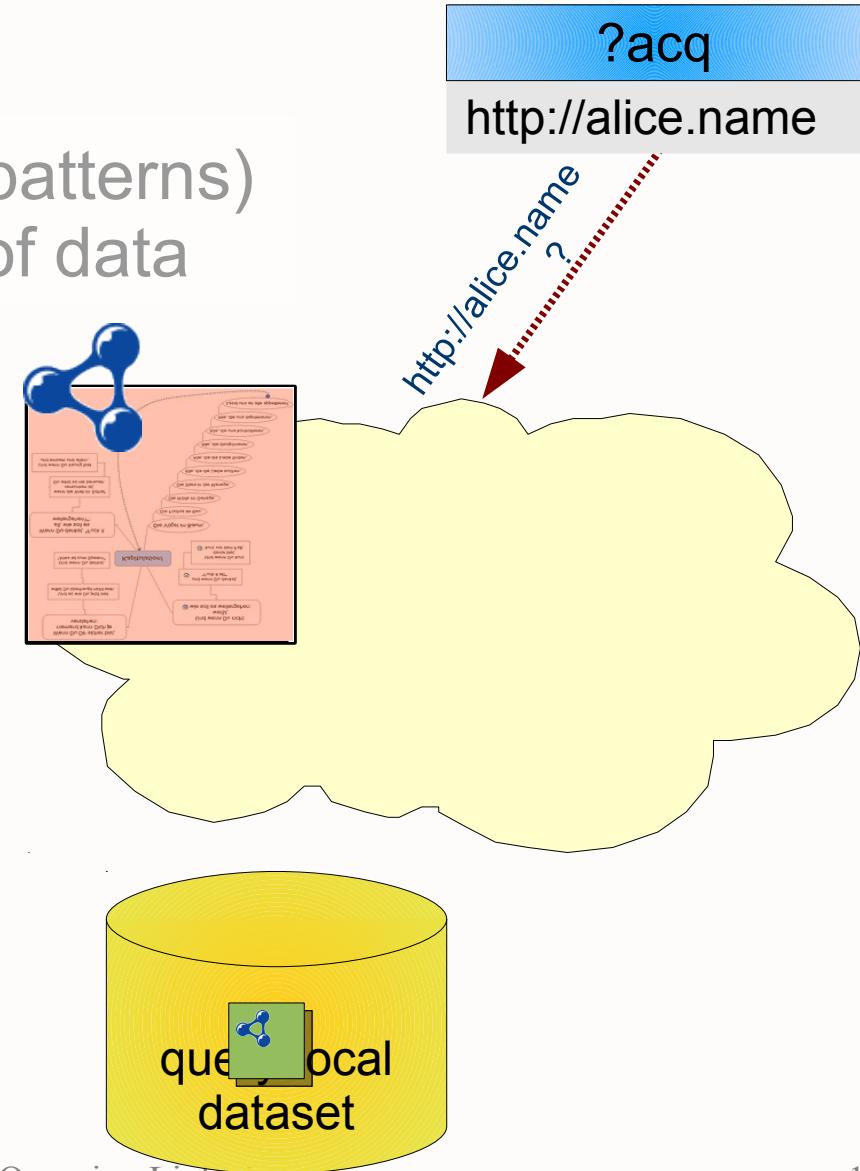
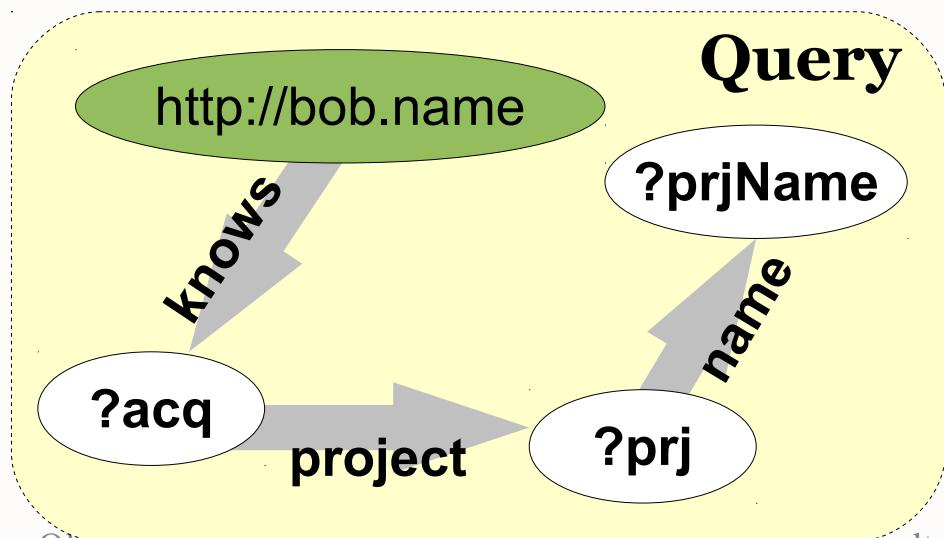


- **Intertwine query evaluation with traversal of data links**
- **We alternate between:**
  - Evaluate parts of the query (triple patterns) on a continuously augmented set of data
  - Look up URIs in intermediate solutions and add retrieved data to the query-local dataset



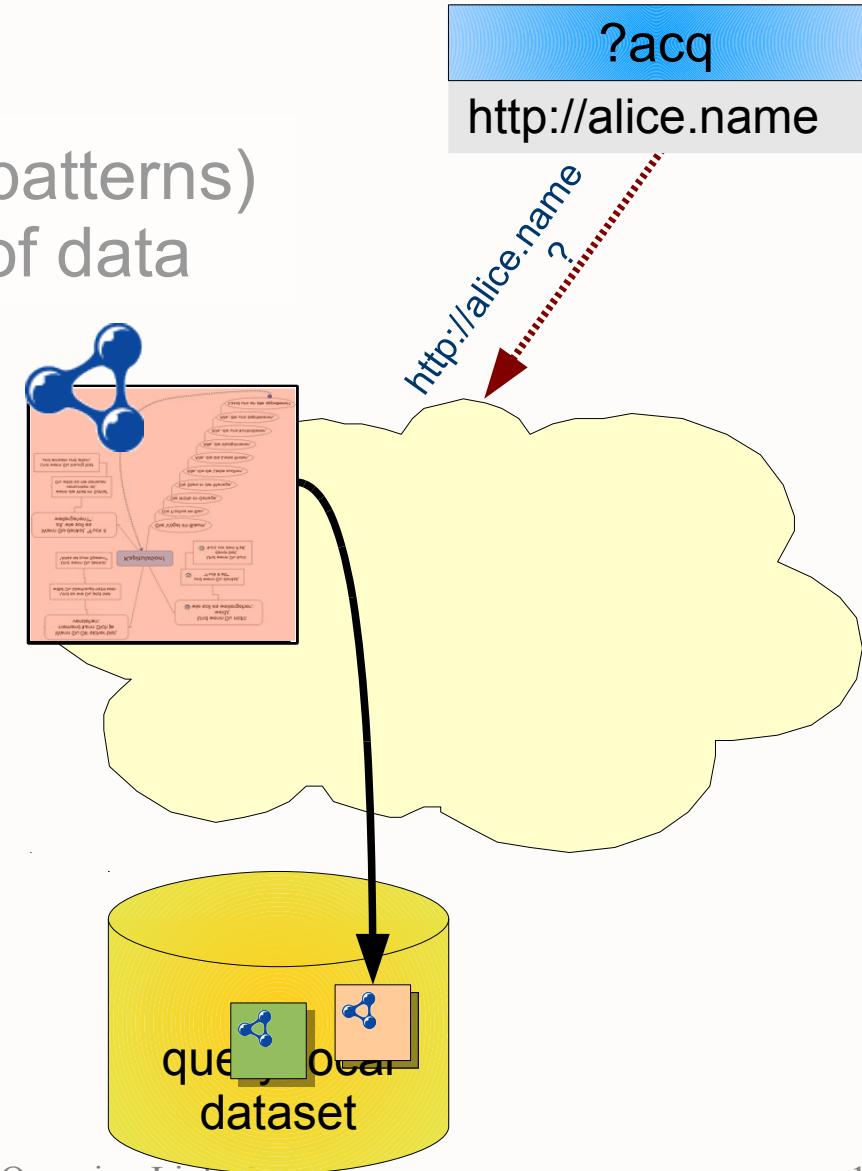
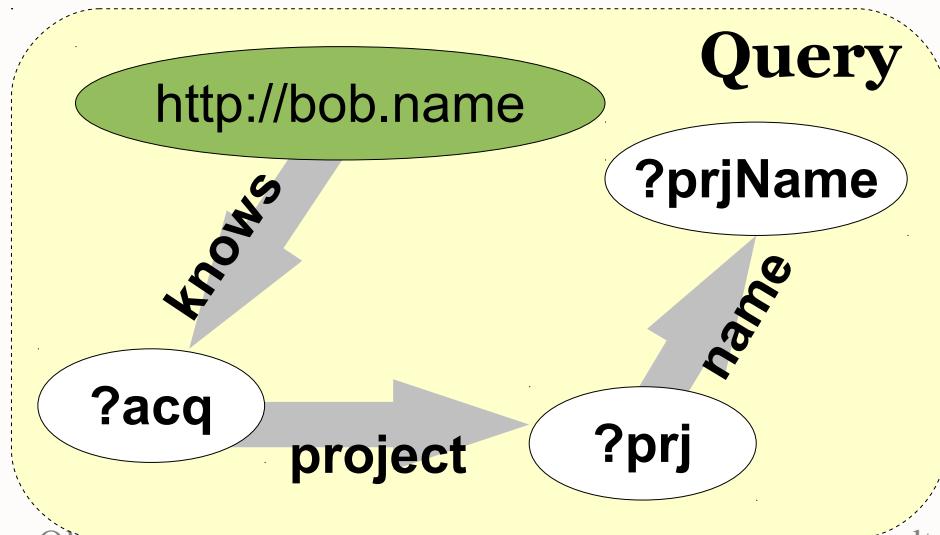
# Main Idea

- Intertwine query evaluation with traversal of data links
- We alternate between:
  - Evaluate parts of the query (triple patterns) on a continuously augmented set of data
  - Look up URIs in intermediate solutions and add retrieved data to the query-local dataset



# Main Idea

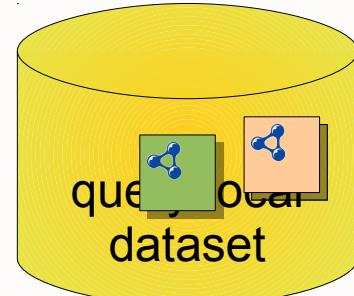
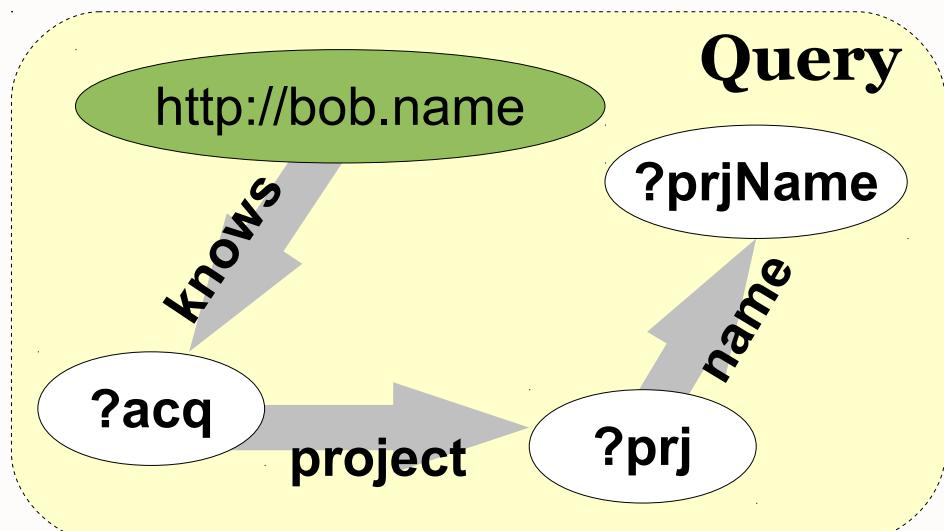
- Intertwine query evaluation with traversal of data links
- We alternate between:
  - Evaluate parts of the query (triple patterns) on a continuously augmented set of data
  - Look up URIs in intermediate solutions and add retrieved data to the query-local dataset



# Main Idea

- **Intertwine query evaluation with traversal of data links**
- **We alternate between:**
  - Evaluate parts of the query (triple patterns) on a continuously augmented set of data
  - Look up URIs in intermediate solutions and add retrieved data to the query-local dataset

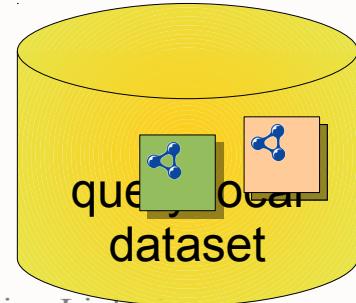
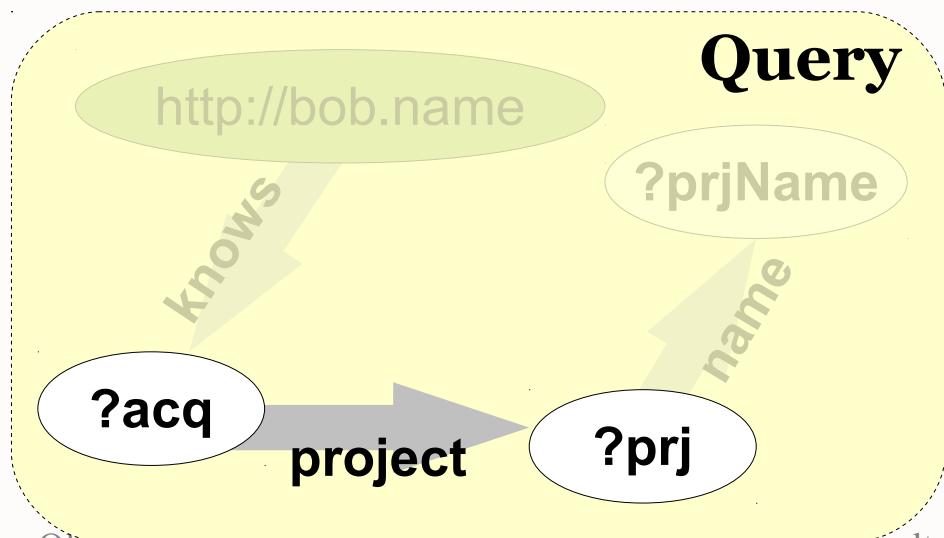
?acq  
http://alice.name



# Main Idea

- **Intertwine query evaluation with traversal of data links**
- **We alternate between:**
  - Evaluate parts of the query (triple patterns) on a continuously augmented set of data
  - Look up URIs in intermediate solutions and add retrieved data to the query-local dataset

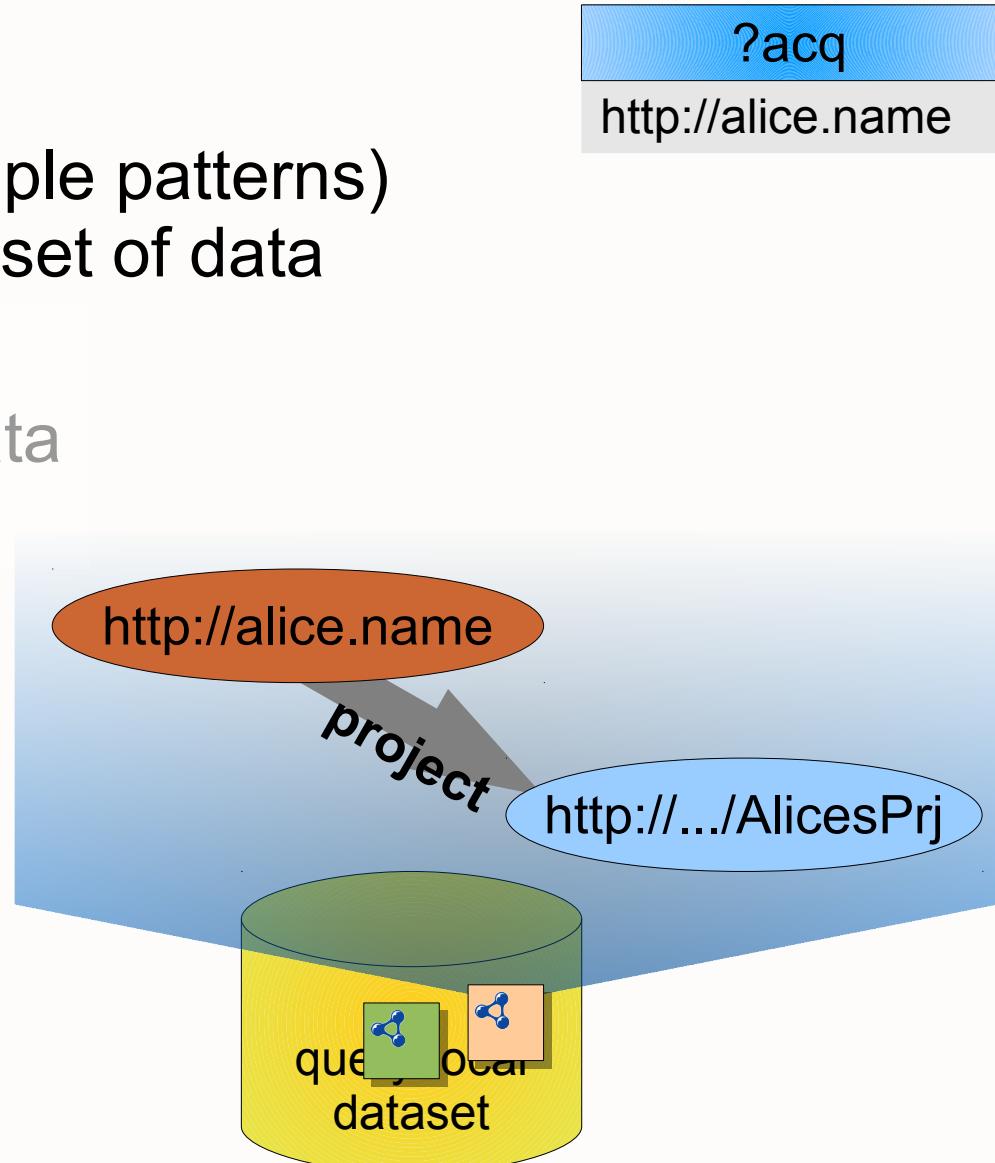
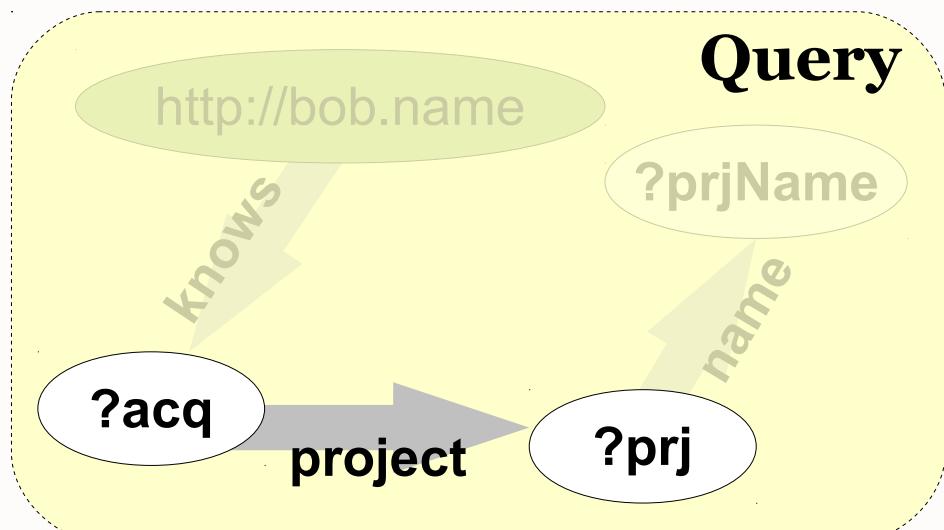
?acq  
http://alice.name



# Main Idea



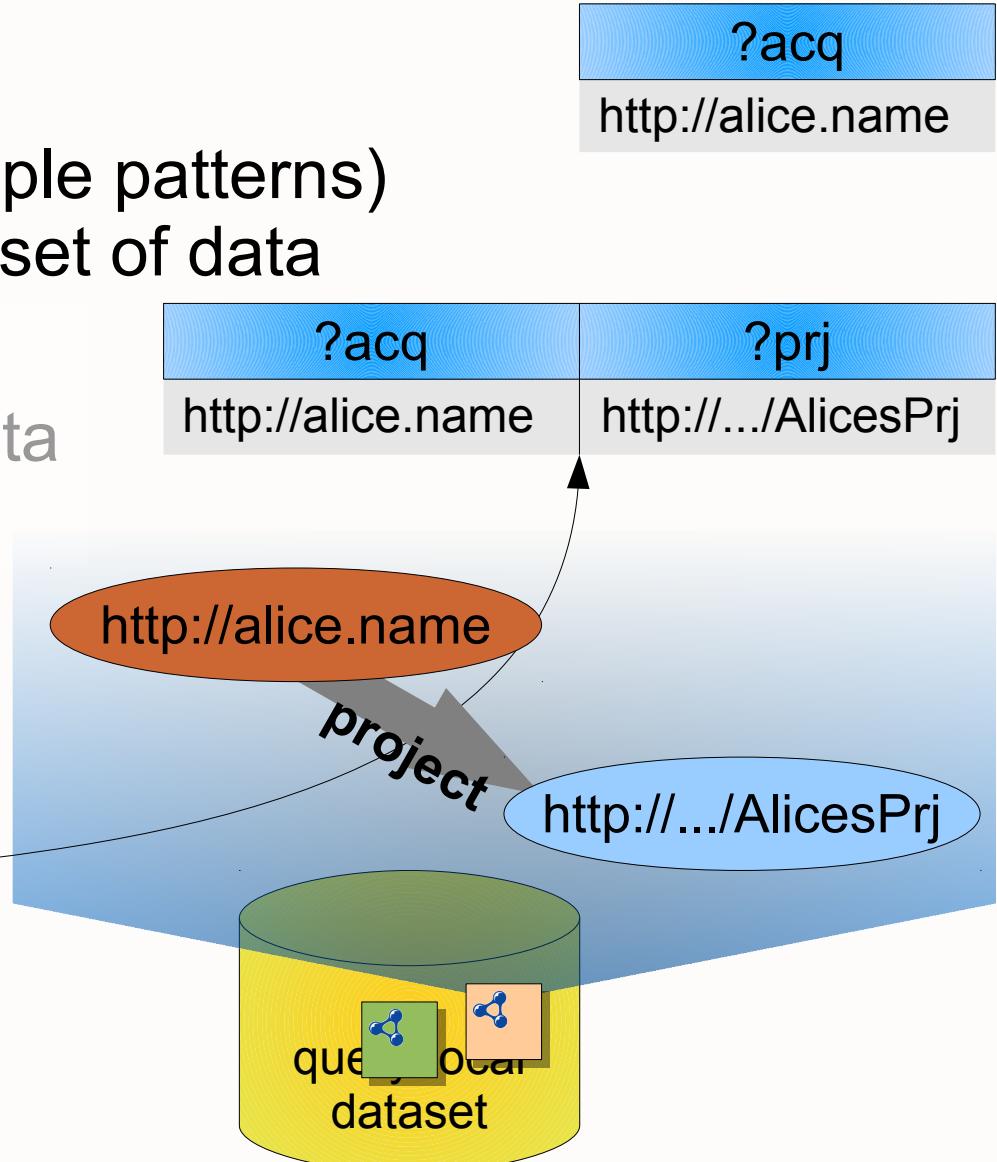
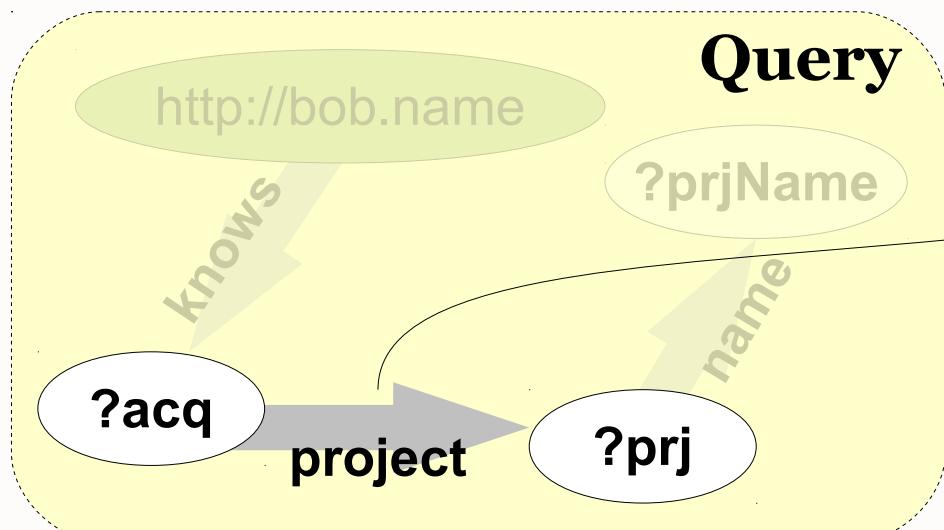
- **Intertwine query evaluation with traversal of data links**
- **We alternate between:**
  - Evaluate parts of the query (triple patterns) on a continuously augmented set of data
  - Look up URIs in intermediate solutions and add retrieved data to the query-local dataset



# Main Idea

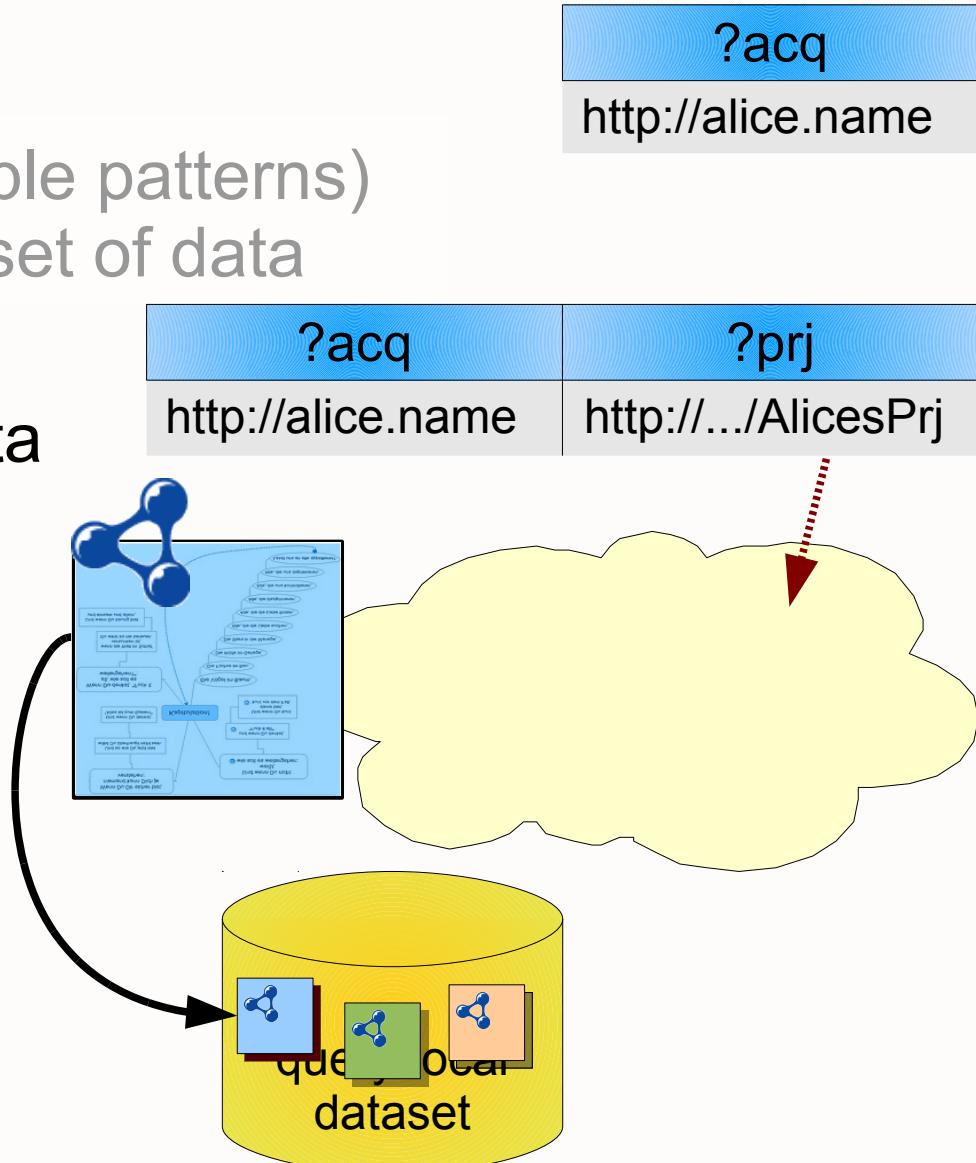
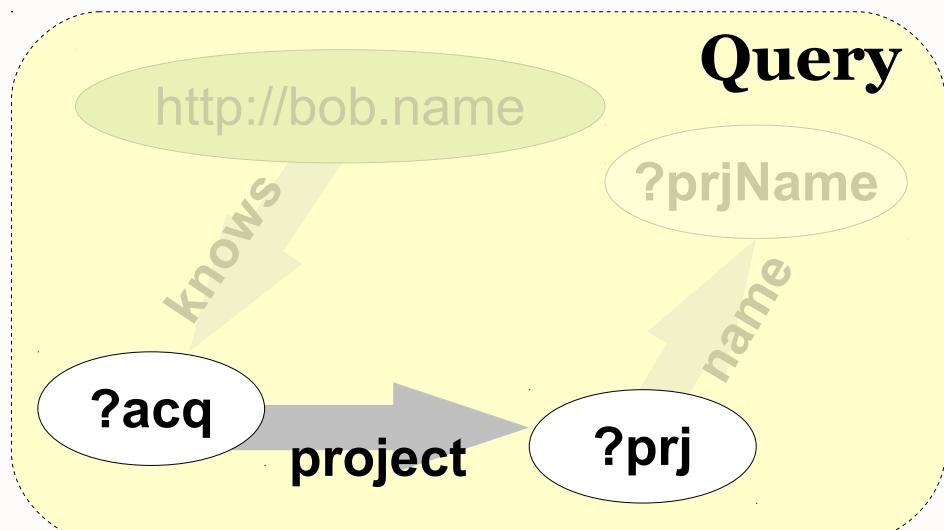


- Intertwine query evaluation with traversal of data links
- We alternate between:
  - Evaluate parts of the query (triple patterns) on a continuously augmented set of data
  - Look up URIs in intermediate solutions and add retrieved data to the query-local dataset



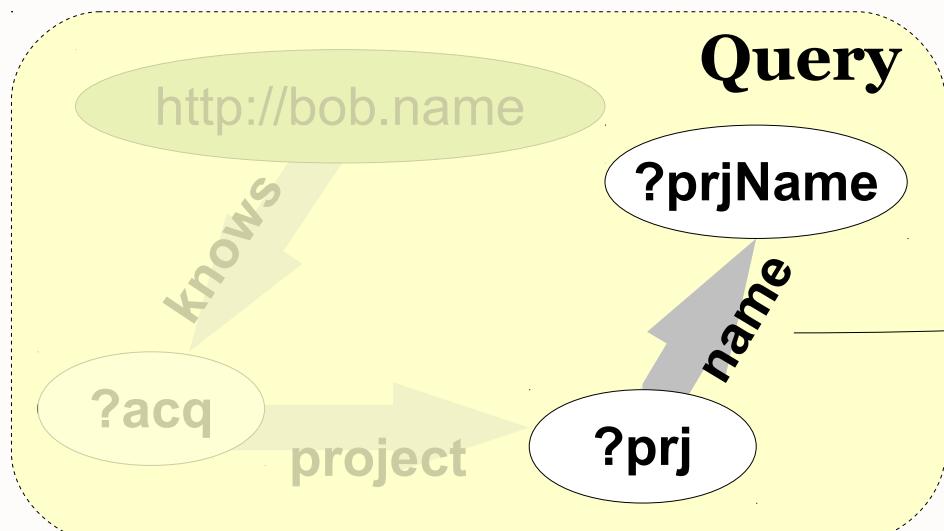
# Main Idea

- **Intertwine query evaluation with traversal of data links**
- **We alternate between:**
  - Evaluate parts of the query (triple patterns) on a continuously augmented set of data
  - Look up URIs in intermediate solutions and add retrieved data to the query-local dataset

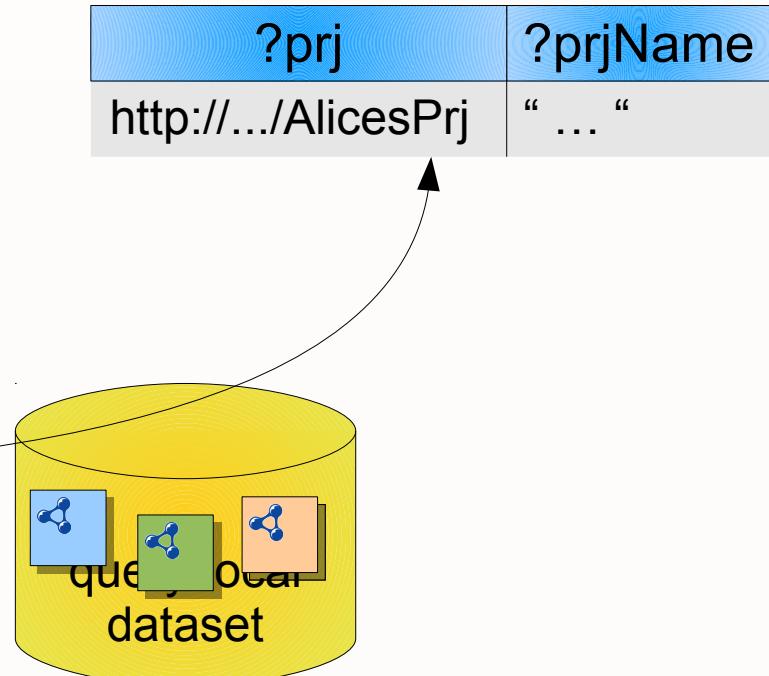


# Main Idea

- Intertwine query evaluation with traversal of data links
- We alternate between:
  - Evaluate parts of the query (triple patterns) on a continuously augmented set of data
  - Look up URIs in intermediate solutions and add retrieved data to the query-local dataset

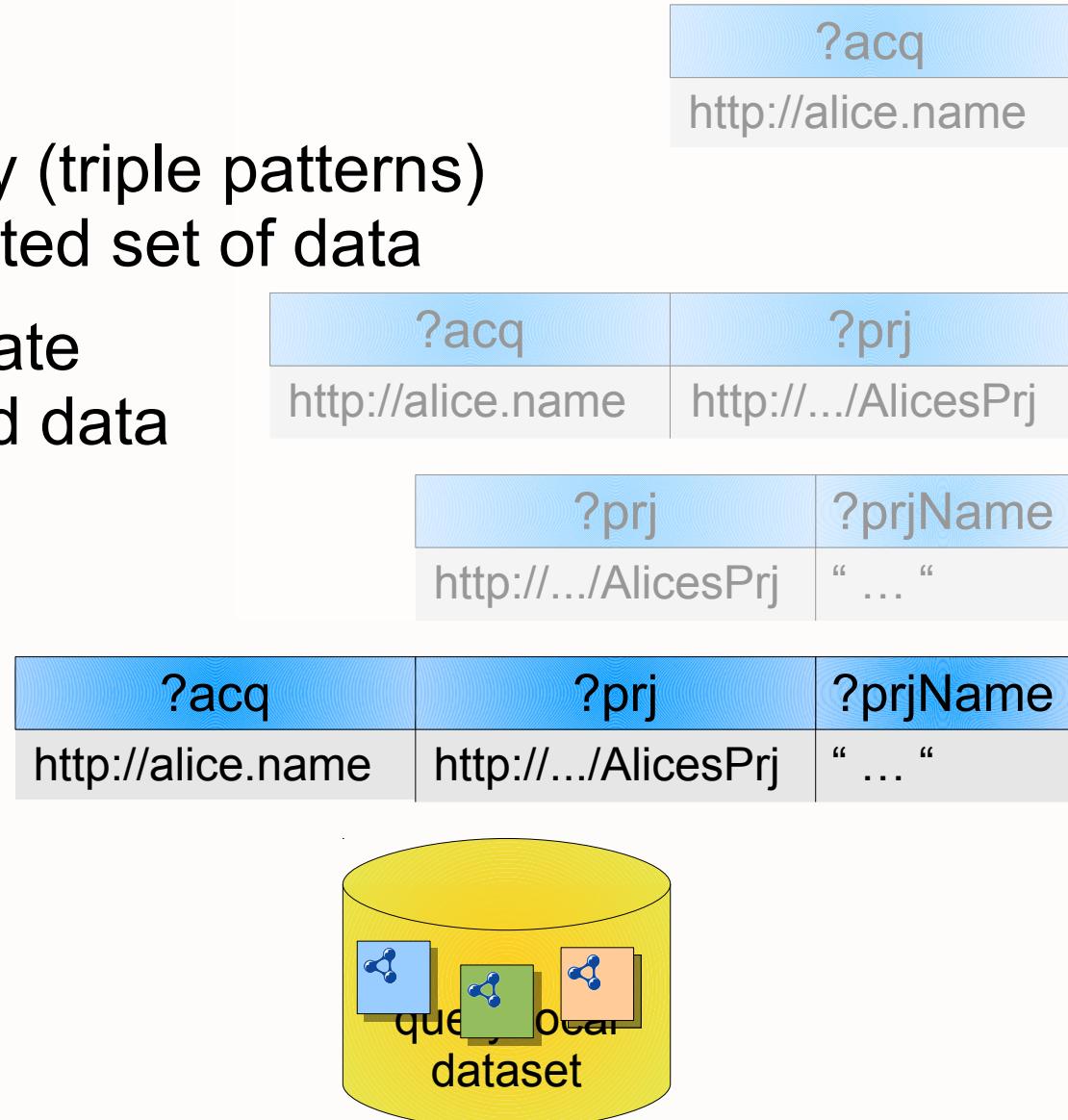
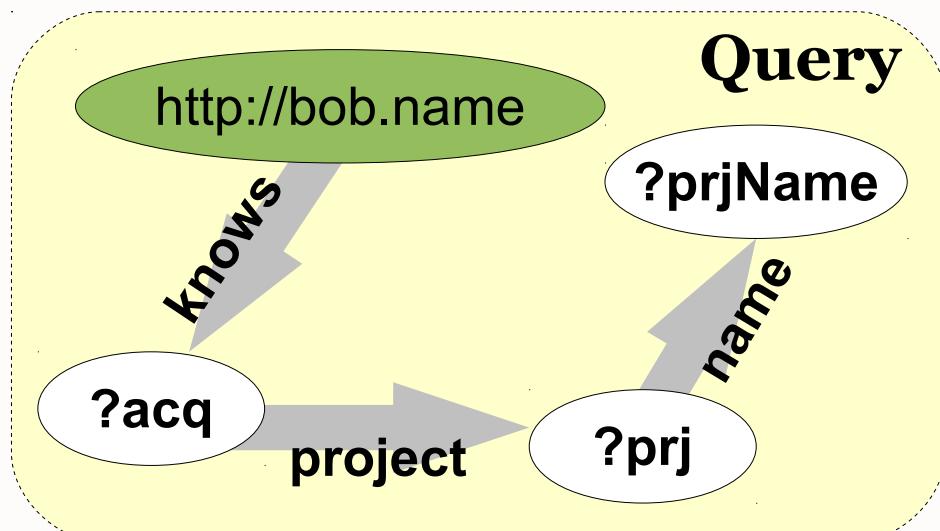


?	acq	?	prj
http://alice.name		http://.../AlicesPrj	
?	prj	?	prjName
http://.../AlicesPrj		" ... "	



# Main Idea

- Intertwine query evaluation with traversal of data links
- We alternate between:
  - Evaluate parts of the query (triple patterns) on a continuously augmented set of data
  - Look up URIs in intermediate solutions and add retrieved data to the query-local dataset



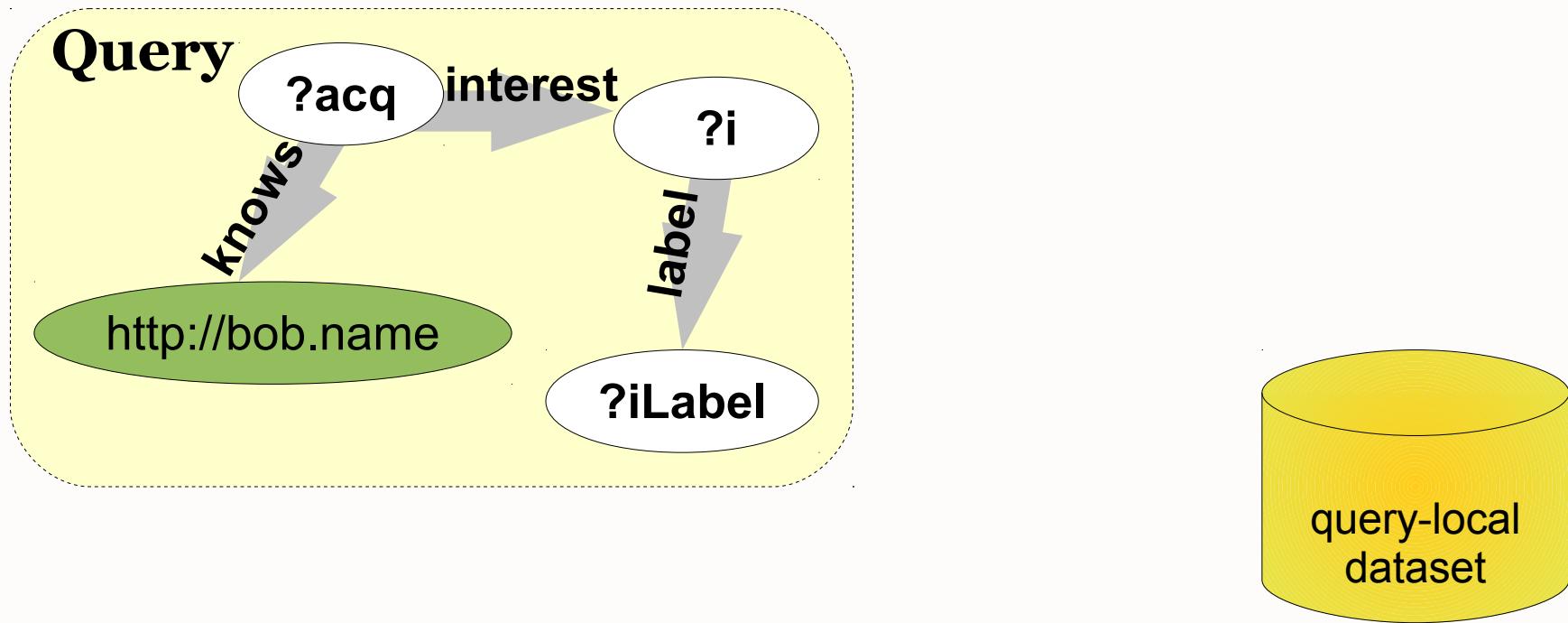


# Characteristics

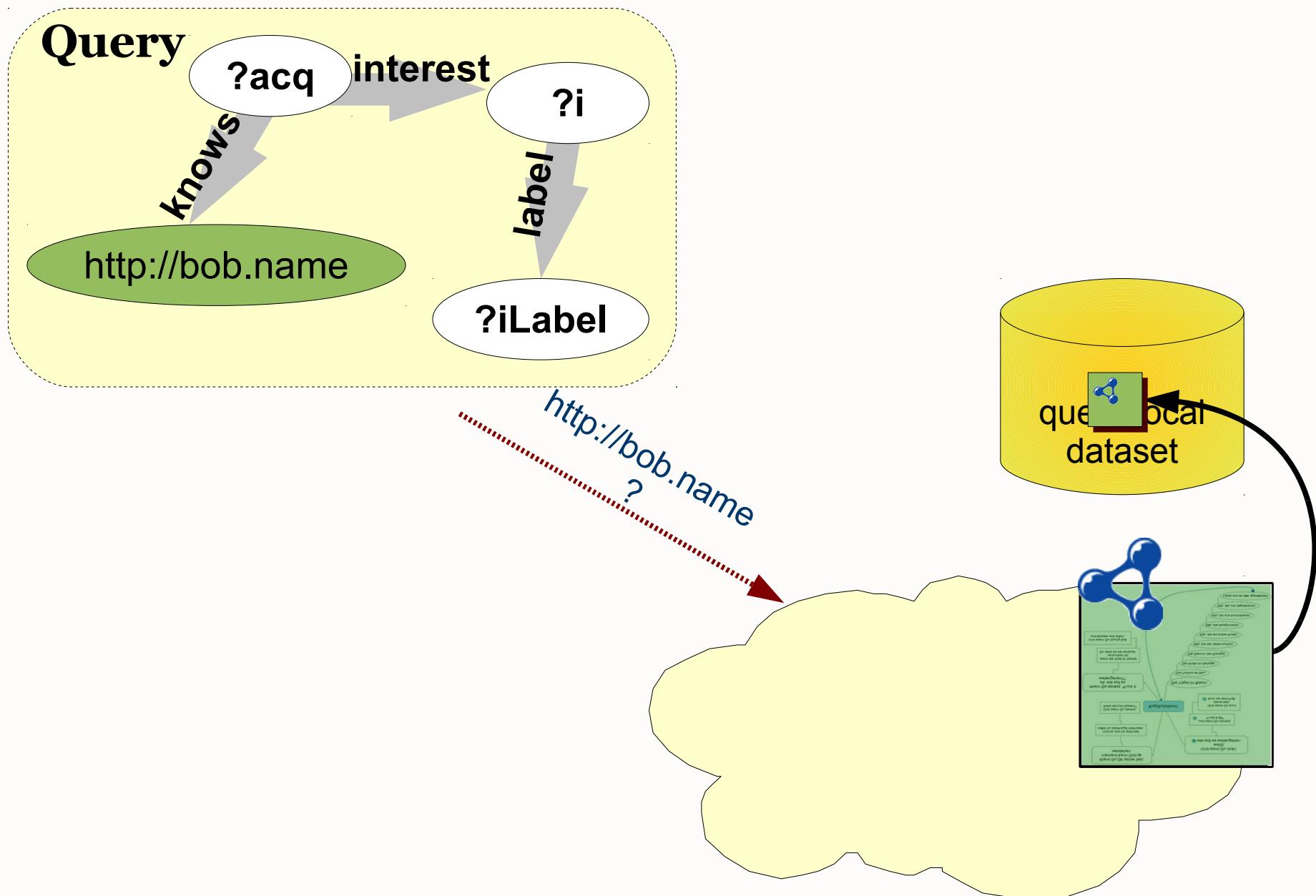
- **Link traversal based query execution:**
  - Evaluation on a continuously augmented dataset
  - Discovery of potentially relevant data during execution
  - Discovery driven by intermediate solutions
- **Main advantage:**
  - No need to know all data sources in advance
- **Limitations:**
  - Query has to contain a URI as a starting point
  - Ignores data that is not *reachable*<sup>\*</sup> by the query execution

<sup>\*</sup> formal definition in the paper

# The Issue

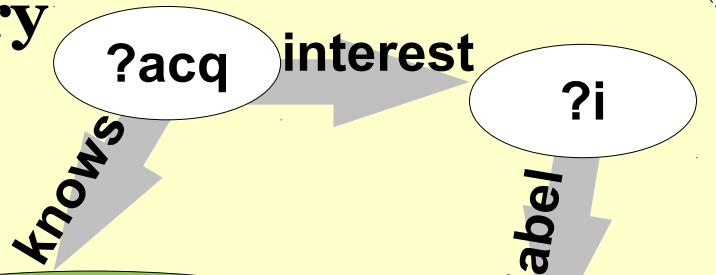


# The Issue



# The Issue

**Query**



http://bob.name

?iLabel

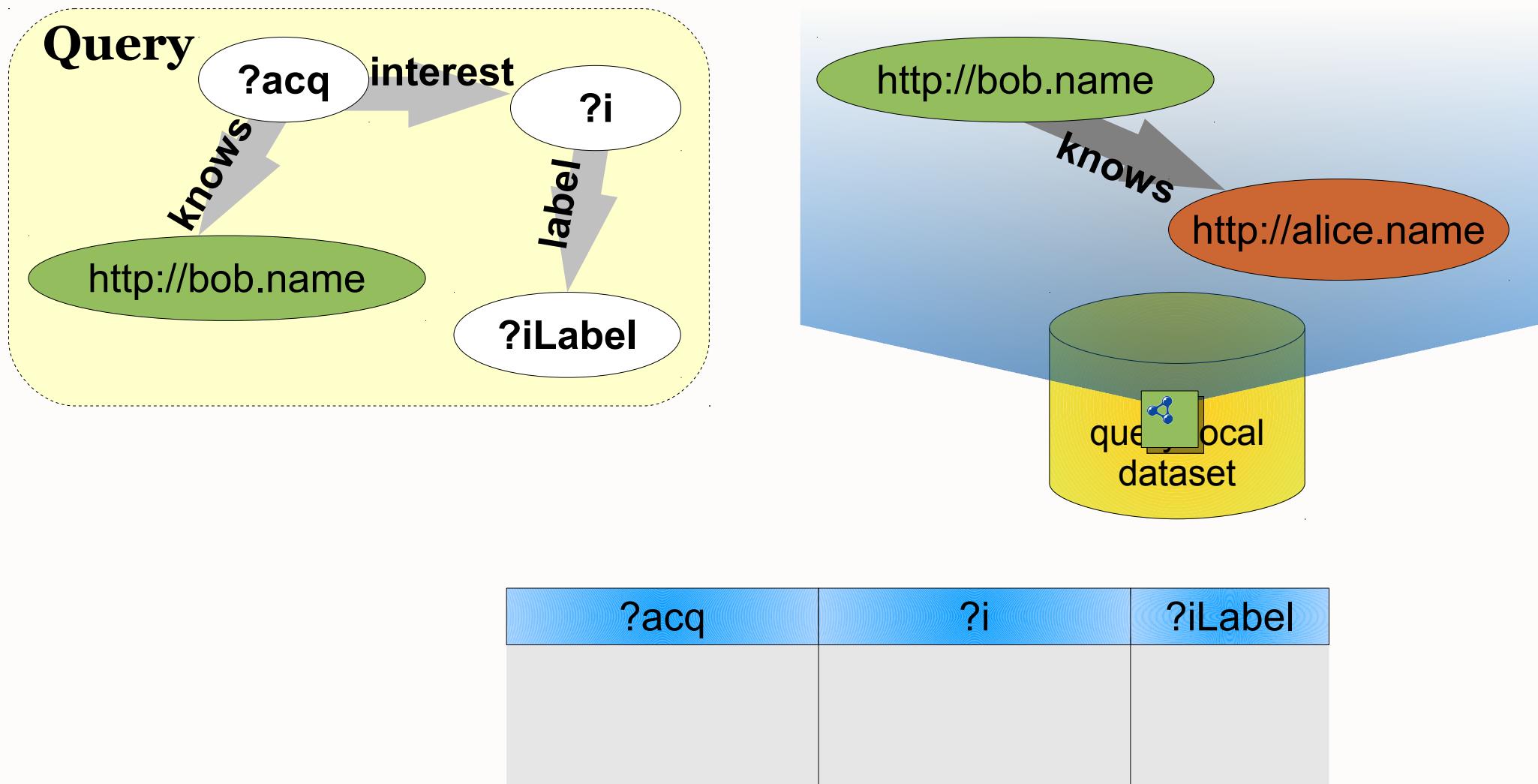
http://bob.name

knows

http://alice.name

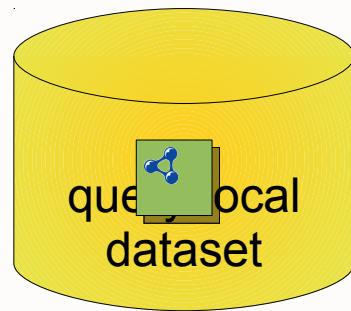
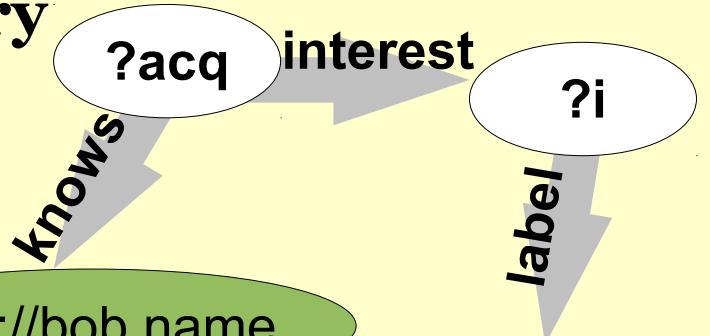
query local  
dataset

# The Issue

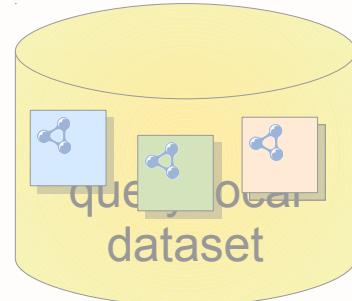
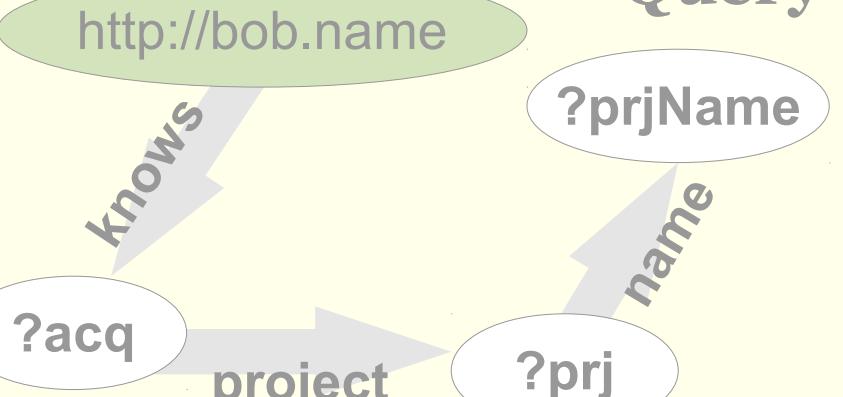


# The Issue

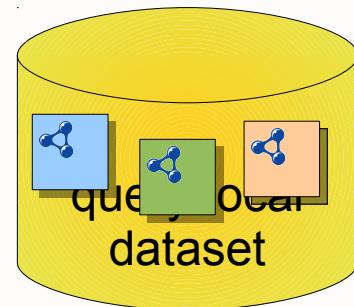
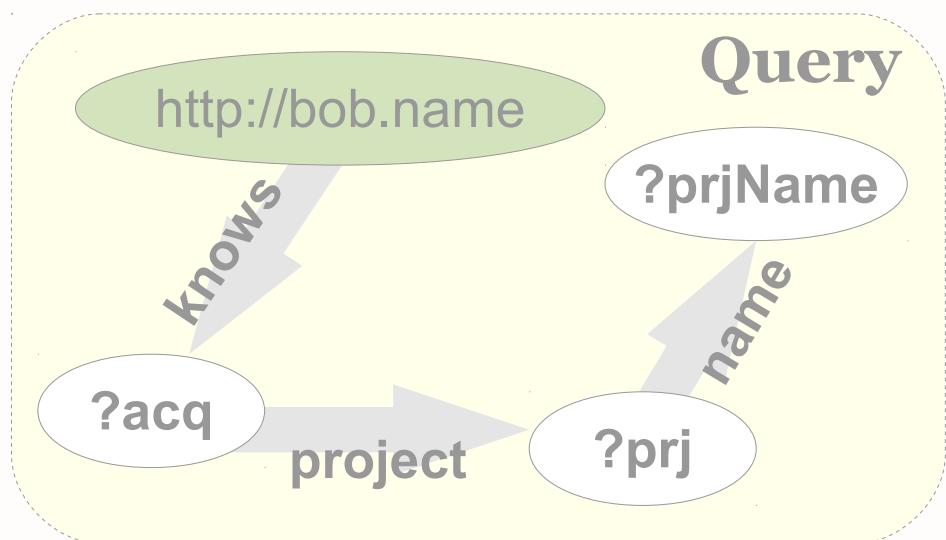
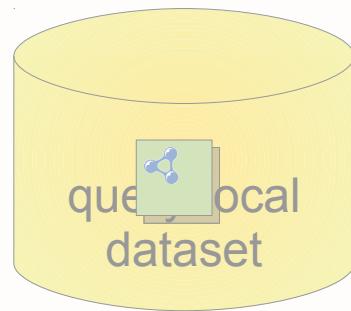
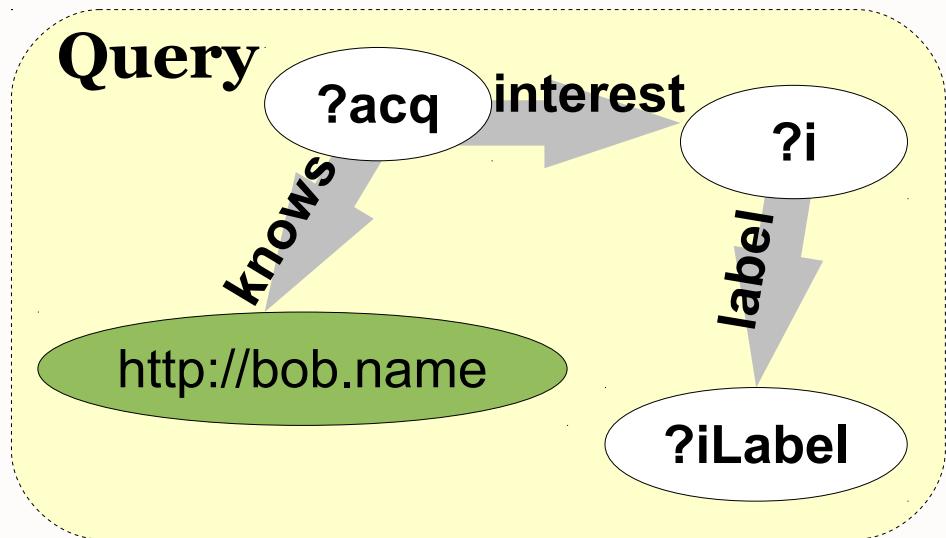
**Query**



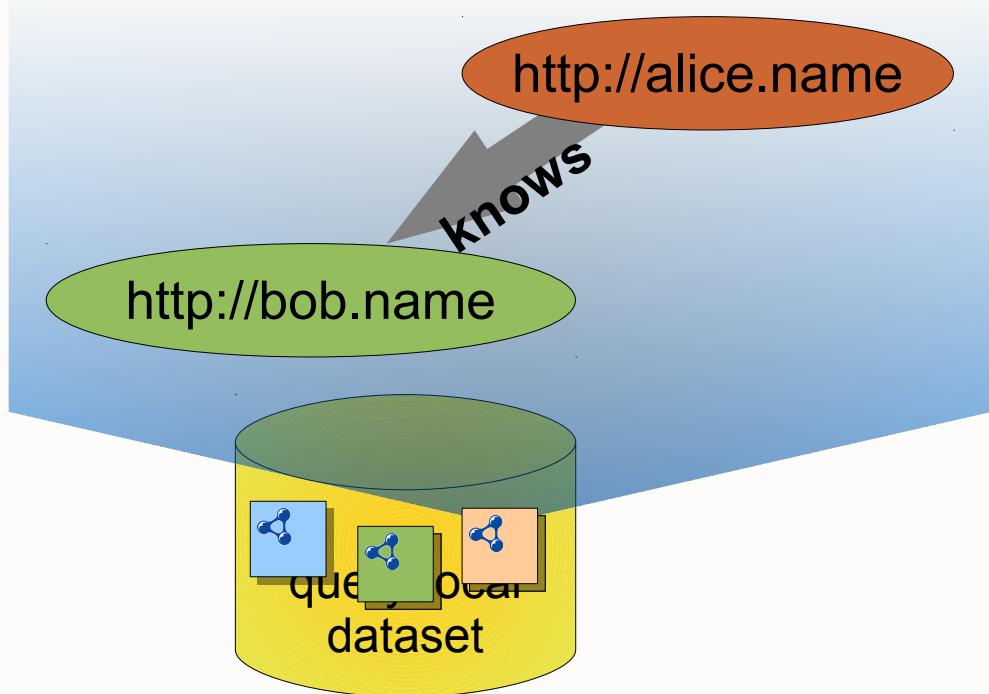
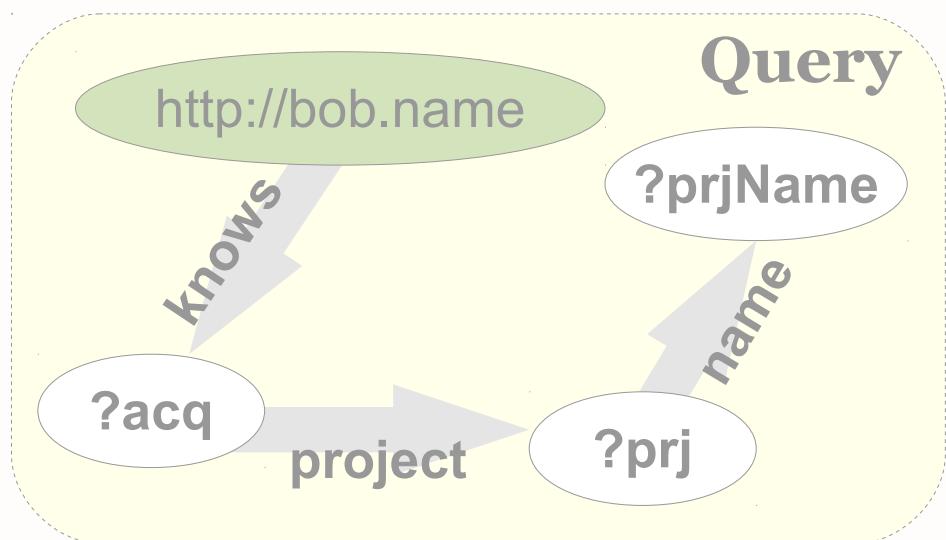
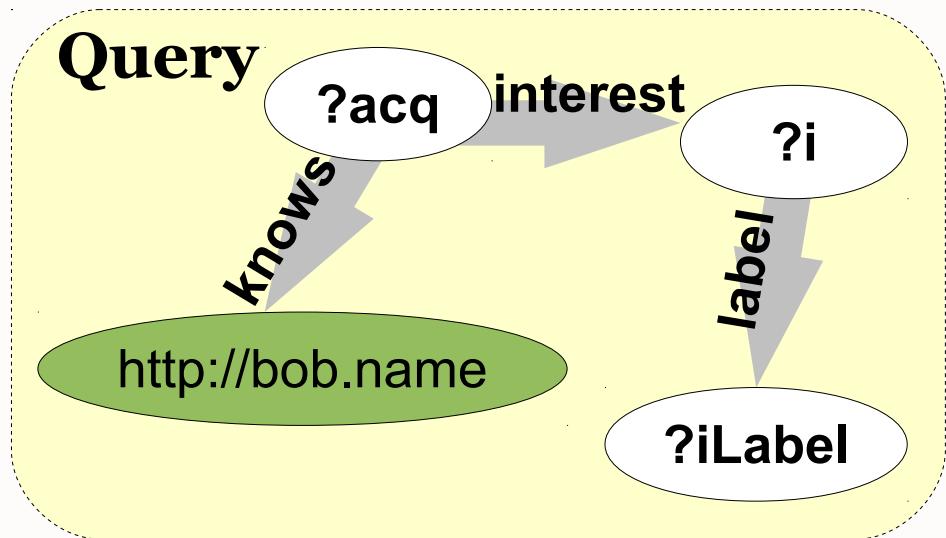
**Query**



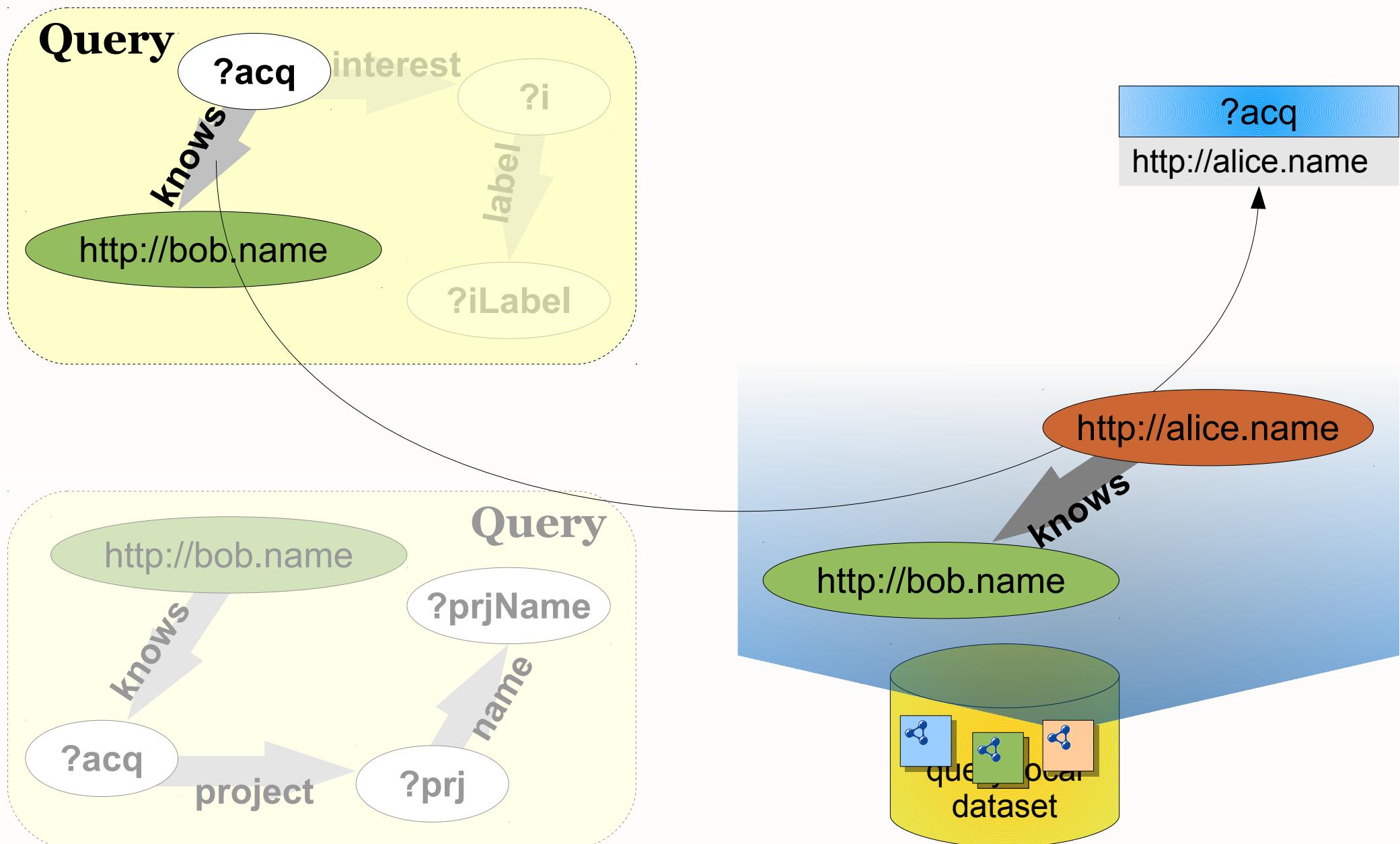
# Reusing the Query-Local Dataset



# Reusing the Query-Local Dataset



# Reusing the Query-Local Dataset



# Hypothesis



**Re-using the query-local dataset (a.k.a. data caching)  
may benefit  
query performance + result completeness**

# Contributions

- **Systematic analysis of the impact of data caching**
  - Theoretical foundation\*
  - Conceptual analysis\*
  - Empirical evaluation of the potential impact

\* see paper

- **Out of scope:** Caching strategies (replacement, invalidation)

# Experiment – Scenario



## FOAF Letter

Home | Options | Mobile |  Search

Welcome, Olaf Hartig!

You have no new incoming contacts.

You claim to know these 3 people, but according to their FOAF profile, they do not seem to know you:

- Juan Sequeda
- Christian Bizer
- Michael Hausenblas

[View potential new contacts](#) (ALPHA warning: might take ages)

There are no upcoming birthdays.

Improve your profile by adding the following information:

- [schoolHomepage](#): A homepage of a school attended by the person
- [myersBriggs](#): A Myers Briggs (MBTI) personality classification.

Privacy policy. We will not use any data you provide here or in your FOAF profile for purposes other than improving this service. We will not store any data in your social network and notifying other users of this service who you claim to know. To protect your privacy, we will remove all our records of your data once you unsubscribe from this newsletter service.

- **Information about the distributed social network of FOAF profiles**

- 5 types of queries

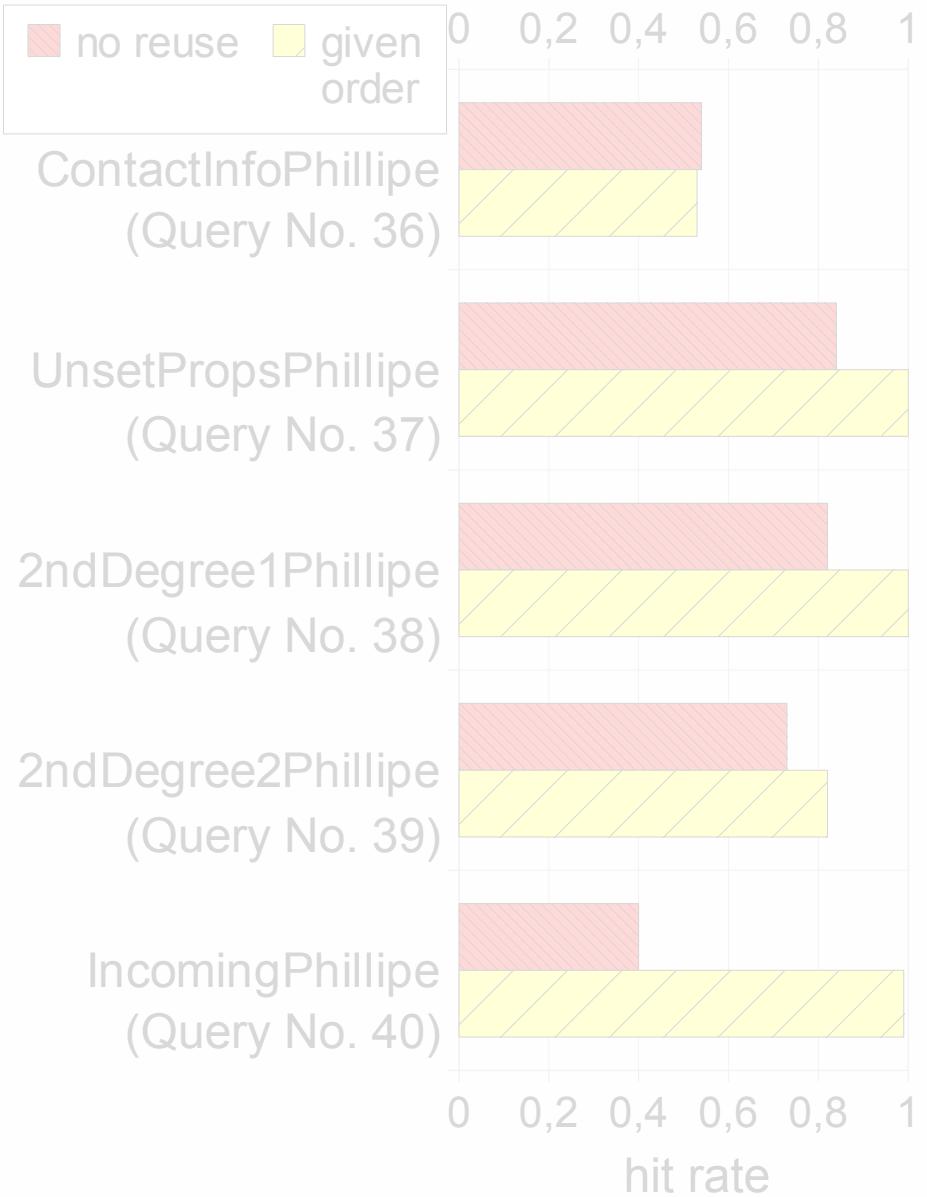
- **Experiment Setup:**

- 23 persons

- Sequential use

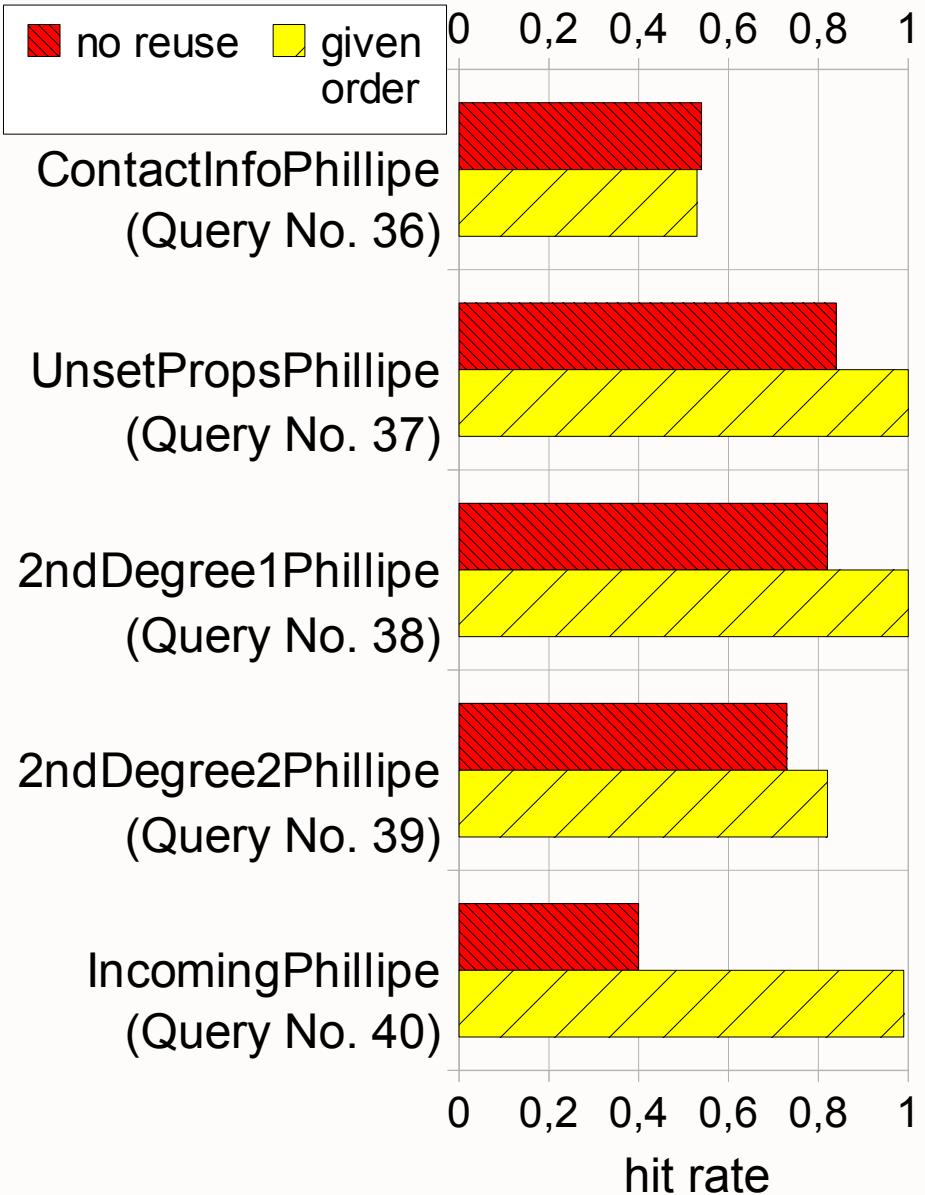
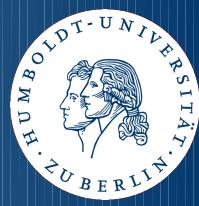
- 115 queries

# Experiment – Complete Sequence



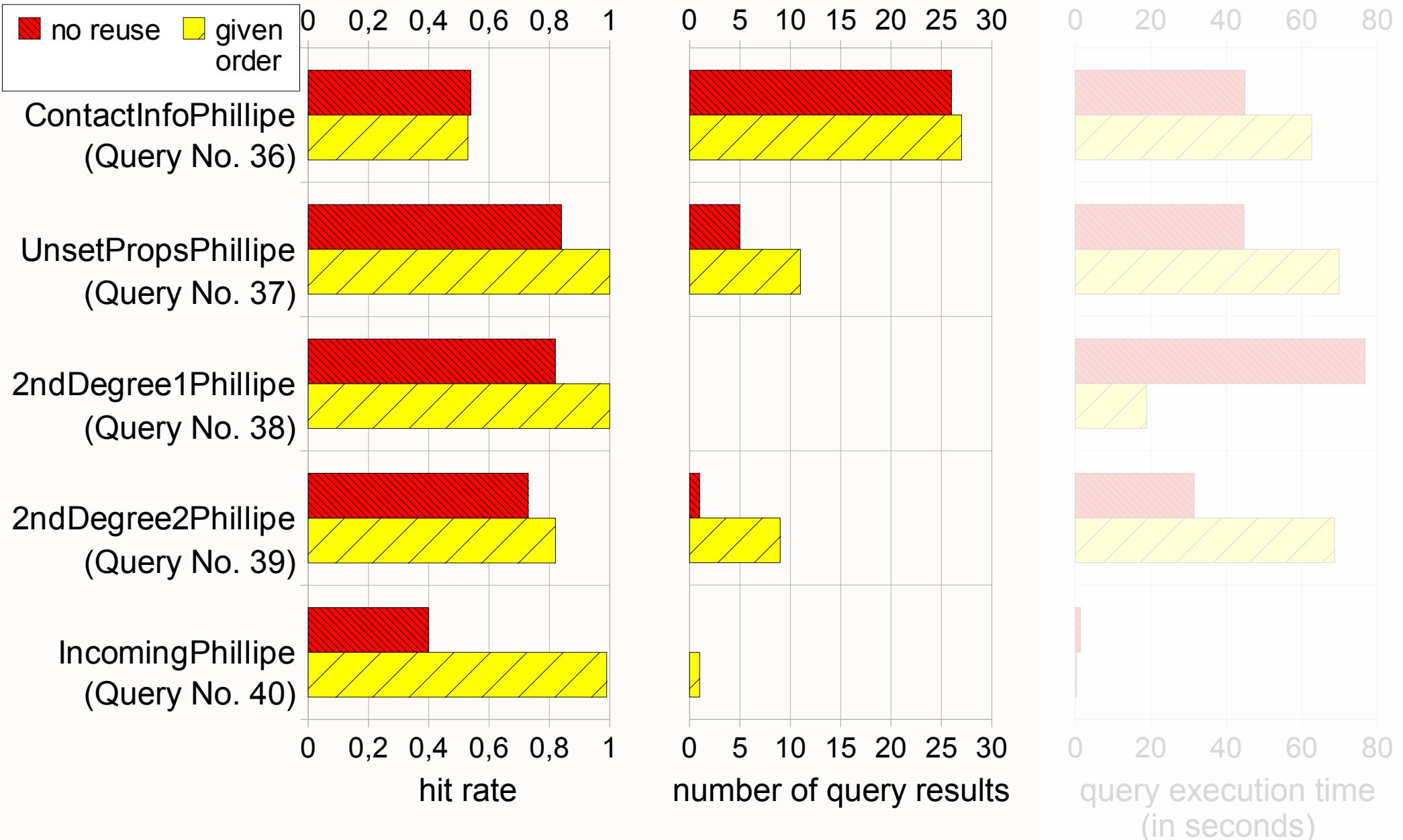
- ***no reuse experiment:***
  - No data caching
- ***given order experiment***
  - Reuse of the query-local dataset for the complete sequence of all 115 queries
- **Hit rate:**  
$$\frac{\text{look-ups answered from cache}}{\text{all look-up requests}}$$

# Experiment – Complete Sequence

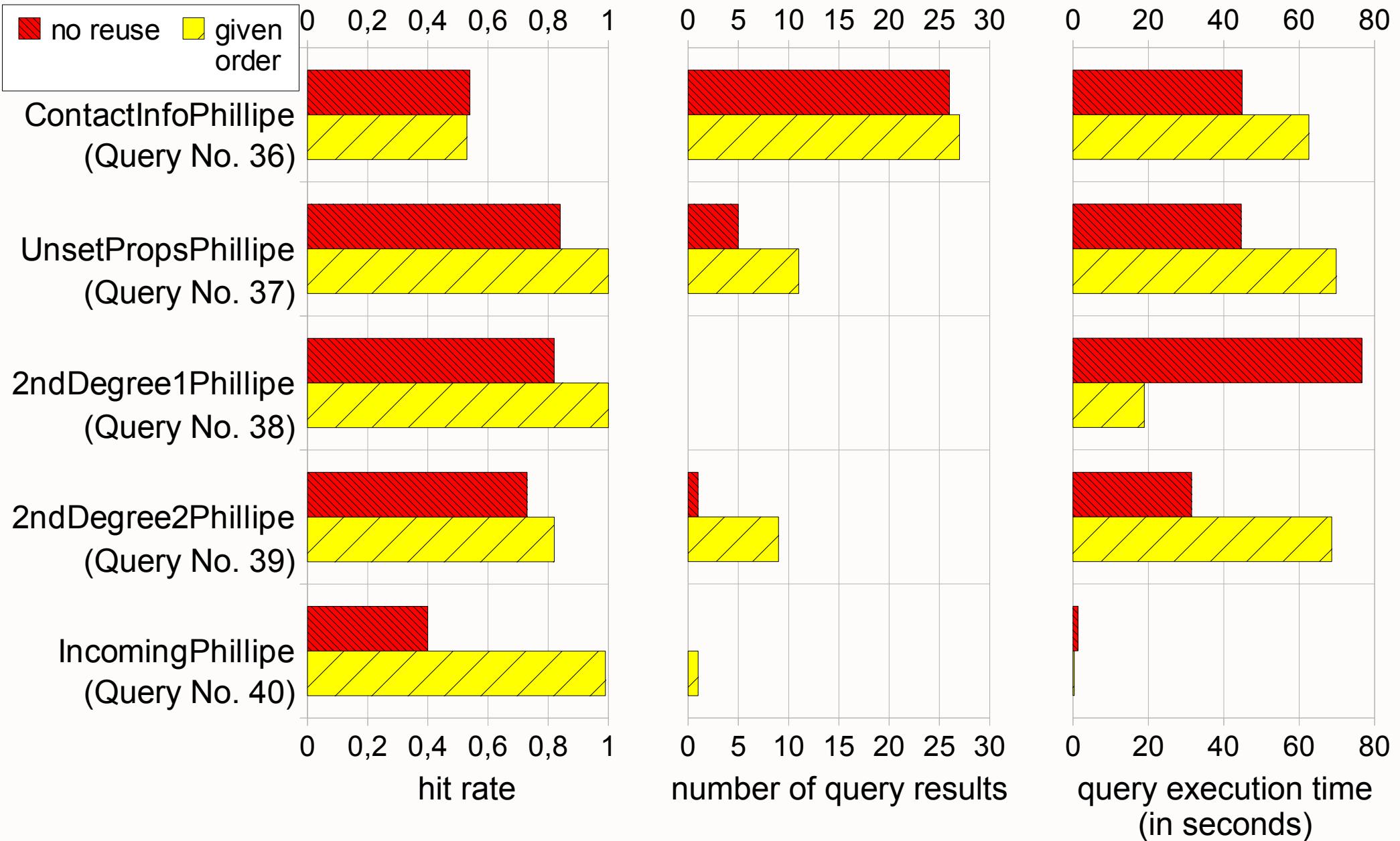


- ***no reuse experiment:***
  - No data caching
- ***given order experiment***
  - Reuse of the query-local dataset for the complete sequence of all 115 queries
- **Hit rate:**  
look-ups answered from cache  
all look-up requests

# Experiment – Complete Sequence



# Experiment – Complete Sequence



# Summary

- **Contributions:**
  - Theoretical foundation
  - Conceptual analysis
  - Empirical evaluation
- **Main findings:**
  - Additional results possible (for semantically similar queries)
  - Impact on performance may be positive but also negative
- **Future work:**
  - Analysis of caching strategies in our context
  - Main issue: invalidation



# Backup Slides

# Contributions

- **Theoretical foundation (extension of the original definition)**
  - Reachability by a  $D_{\text{seed}}$ -initialized execution of a BGP query  $b$
  - $D_{\text{seed}}$ -dependent solution for a BGP query  $b$
  - Reachability  $R(B)$  for a serial execution of  $B = b_1, \dots, b_n$ 
    - Each solution for  $b_{\text{cur}}$  is also  $R(B)$ -dependent solution for  $b_{\text{cur}}$
- **Conceptual analysis of the impact of data caching**
  - Performance factor:  $p( b_{\text{cur}}, B ) = c( b_{\text{cur}}, [] ) - c( b_{\text{cur}}, B )$
  - Serendipity factor:  $s( b_{\text{cur}}, B ) = b( b_{\text{cur}}, B ) - b( b_{\text{cur}}, [] )$
- **Empirical verification of the potential impact**
- **Out of scope:** Caching strategies (replacement, invalidation)

# Query Template Contact

```
SELECT * WHERE {   <PERSON> foaf:knows ?p .  
                    OPTIONAL { ?p foaf:name ?name }  
                    OPTIONAL { ?p foaf:firstName ?firstName }  
                    OPTIONAL { ?p foaf:givenName ?givenName }  
                    OPTIONAL { ?p foaf:givenname ?givenname }  
                    OPTIONAL { ?p foaf:familyName ?familyName }  
                    OPTIONAL { ?p foaf:family_name ?family_name }  
                    OPTIONAL { ?p foaf:lastName ?lastName }  
                    OPTIONAL { ?p foaf:surname ?surname }  
  
                    OPTIONAL { ?p foaf:birthday ?birthday }  
  
                    OPTIONAL { ?p foaf:img ?img }  
  
                    OPTIONAL { ?p foaf:phone ?phone }  
                    OPTIONAL { ?p foaf:aimChatID ?aimChatID }  
                    OPTIONAL { ?p foaf:icqChatID ?icqChatID }  
                    OPTIONAL { ?p foaf:jabberID ?jabberID }  
                    OPTIONAL { ?p foaf:msnChatID ?msnChatID }  
                    OPTIONAL { ?p foaf:skypeID ?skypeID }  
                    OPTIONAL { ?p foaf:yahooChatID ?yahooChatID }  
}
```

# Query Template UnsetProps

```
SELECT DISTINCT ?result ?resultLabel WHERE
{
    ?result rdfs:isDefinedBy <http://xmlns.com/foaf/0.1/> .
    ?result rdfs:domain foaf:Person .

OPTIONAL { <PERSON> ?result ?var0 }
FILTER ( !bound(?var0) )

<PERSON> foaf:knows ?var2 .
?var2 ?result ?var3 .
?result rdfs:label ?resultLabel .
?result vs:term_status ?var1 .
}
ORDER BY ?var1
```

# Query Template Incoming

```
SELECT DISTINCT ?result WHERE
{
    ?result foaf:knows <PERSON> .

OPTIONAL
{
    ?result foaf:knows ?var1 .
    FILTER ( <PERSON> = ?var1 )
    <PERSON> foaf:knows ?result .
}
FILTER ( !bound(?var1) )
}
```

# Query Template 2ndDegree1

```
SELECT DISTINCT ?result WHERE
{
  <PERSON> foaf:knows ?p1 .
  <PERSON> foaf:knows ?p2 .
  FILTER ( ?p1 != ?p2 )

  ?p1 foaf:knows ?result .
  FILTER ( <PERSON> != ?result )
  ?p2 foaf:knows ?result .

  OPTIONAL {
    <PERSON> ?knows ?result .
    FILTER ( ?knows = foaf:knows )
  }
  FILTER ( !bound(?knows) )
}
```

# Query Template 2ndDegree2

```
SELECT DISTINCT ?result WHERE
```

```
{
```

```
    <PERSON> foaf:knows ?p1 .
```

```
    <PERSON> foaf:knows ?p2 .
```

```
    FILTER ( ?p1 != ?p2 )
```

```
    ?result foaf:knows ?p1 .
```

```
    FILTER ( <PERSON> != ?result )
```

```
    ?result foaf:knows ?p2 .
```

```
OPTIONAL {
```

```
    <PERSON> ?knows ?result .
```

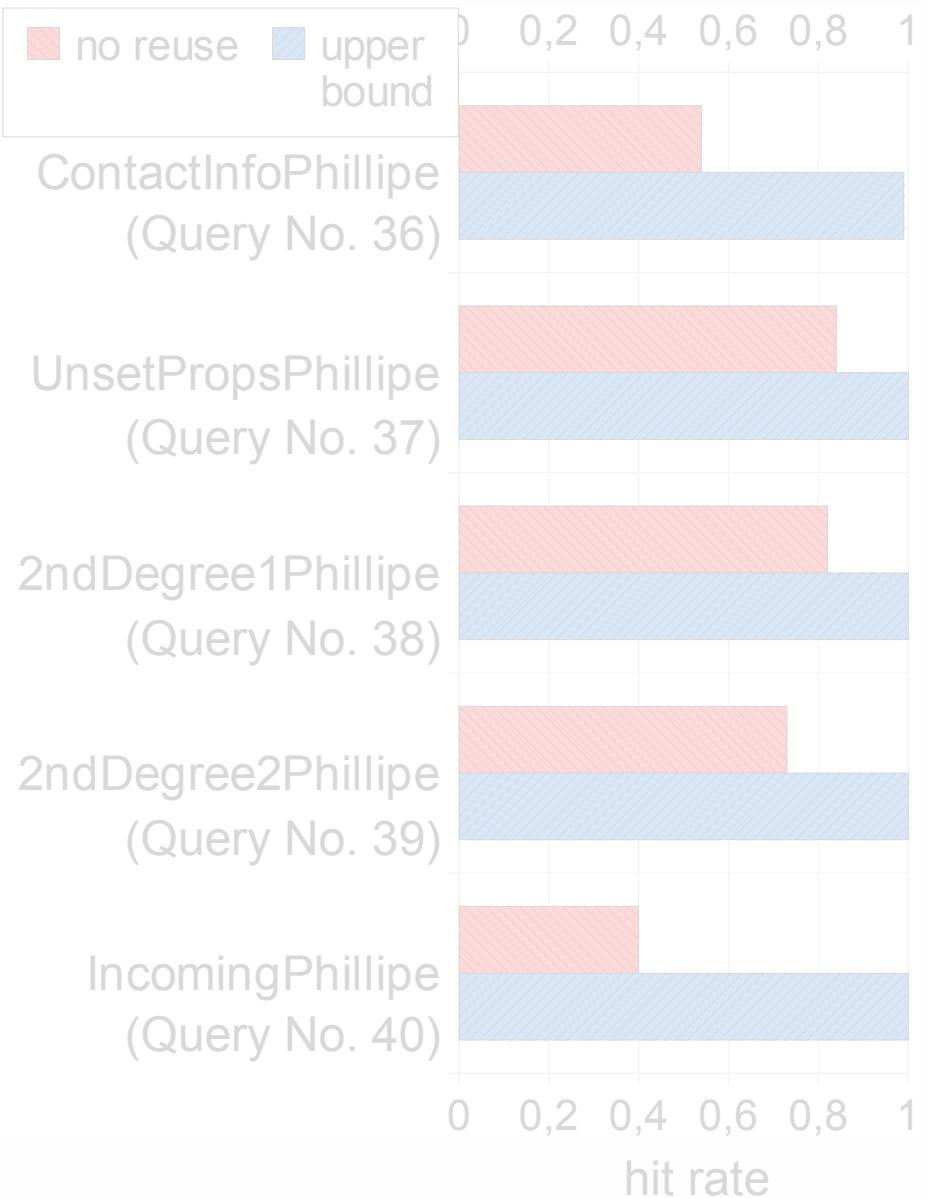
```
    FILTER ( ?knows = foaf:knows )
```

```
}
```

```
    FILTER ( !bound(?knows) )
```

```
}
```

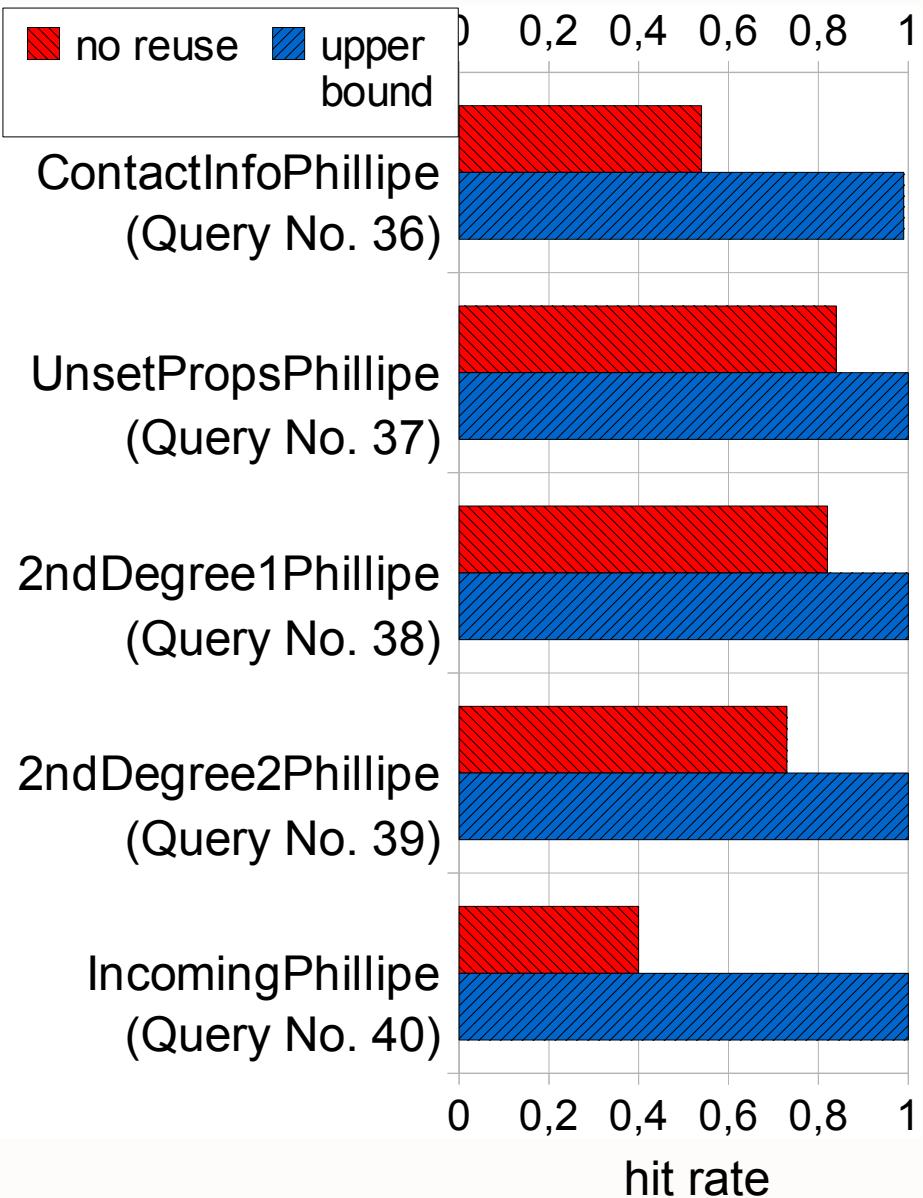
# Experiment – Single Query



- ***no reuse experiment:***
  - No data caching
- ***upper bound experiment***
  - Reuse of query-local dataset for 3 executions of each query
  - Third execution measured
- ***Hit rate:***

$$\frac{\text{look-ups answered from cache}}{\text{all look-up requests}}$$

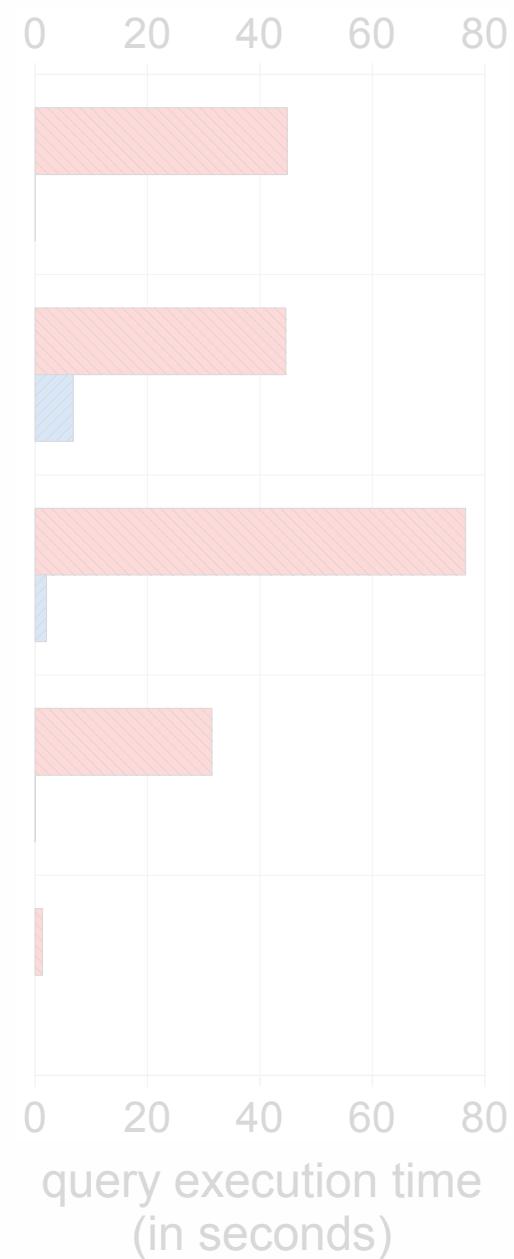
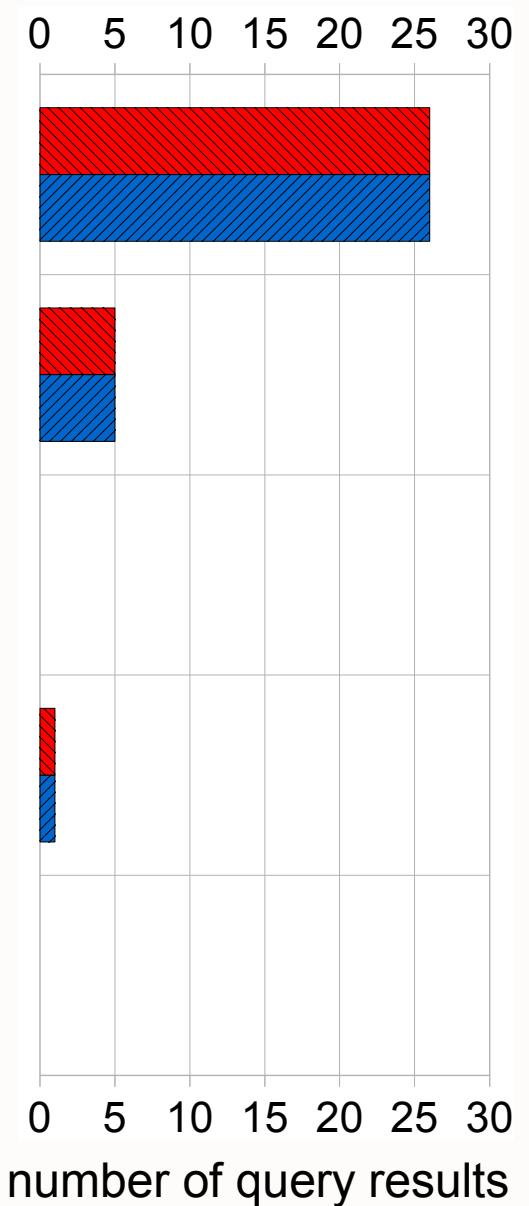
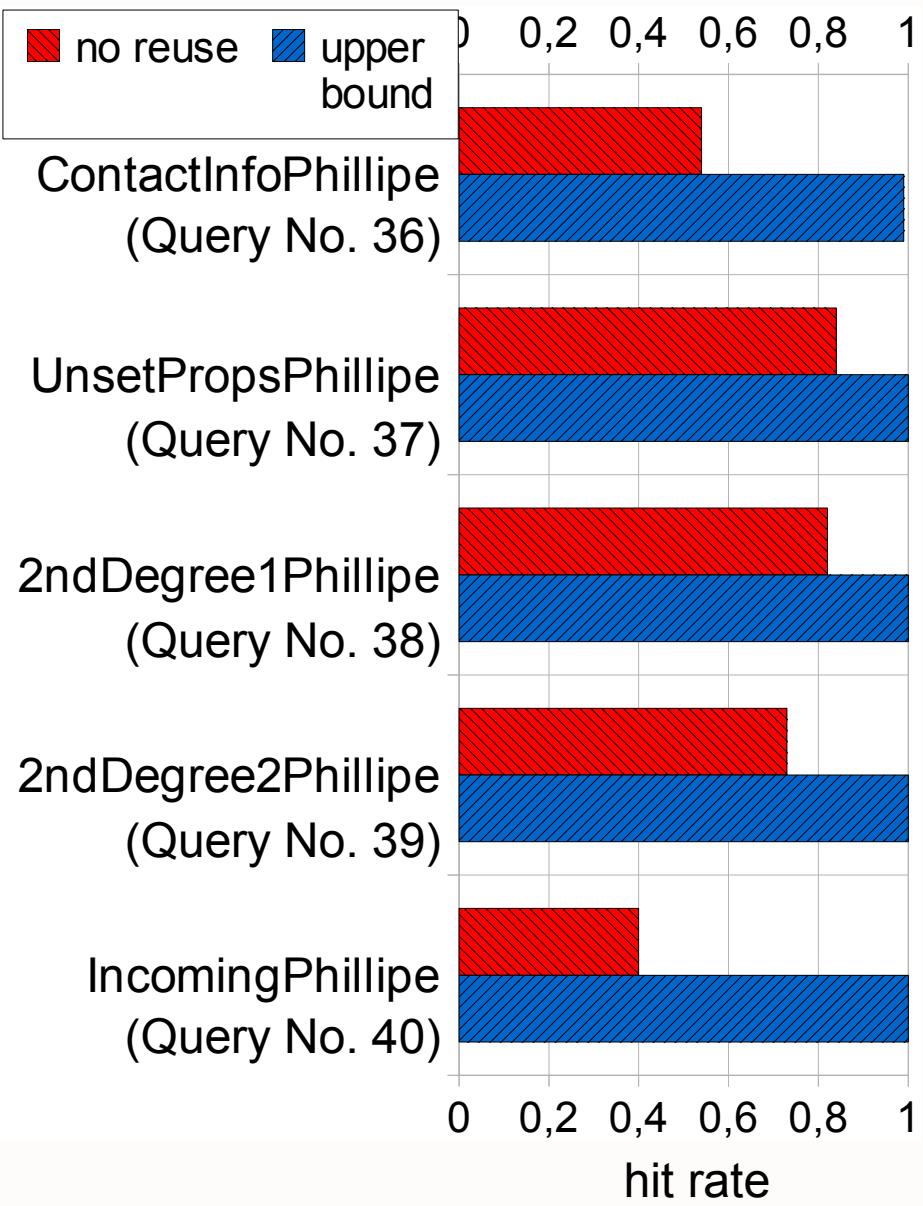
# Experiment – Single Query



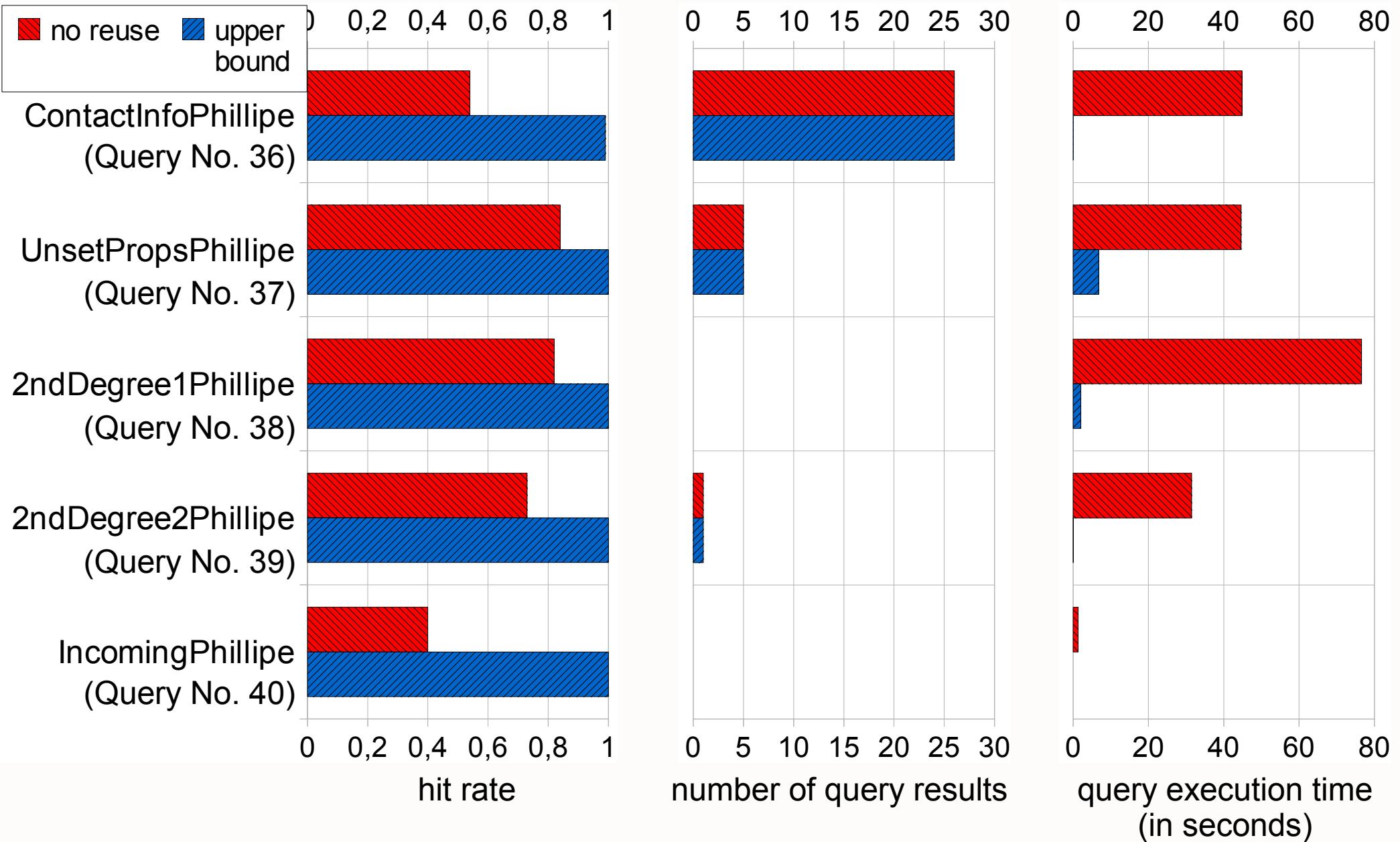
- ***no reuse experiment:***
  - No data caching
- ***upper bound experiment***
  - Reuse of query-local dataset for 3 executions of each query
  - Third execution measured
- ***Hit rate:***

$$\frac{\text{look-ups answered from cache}}{\text{all look-up requests}}$$

# Experiment – Single Query



# Experiment – Single Query



# Experiment – Single Query

Experiment	Avg. <sup>1</sup> number of Query Results (std.dev.)	Average <sup>1</sup> Hit Rate (std.dev.)	Avg. <sup>1</sup> query Execution Time (std.dev.)
no reuse	4.983 (11.658)	0.576 (0.182)	30.036 s (46.708)
upper bound	5.070 (11.813)	0.996 (0.017)	1.943 s (11.375)

<sup>1</sup> Averaged over all 115 queries

- In the ideal case for  $B_{\text{upper}} = [ b_{\text{cur}}, b_{\text{cur}} ]$ :

- $p_{\text{upper}}( b_{\text{cur}}, B_{\text{upper}} ) = c( b_{\text{cur}}, [] ) - c( b_{\text{cur}}, B_{\text{upper}} ) = c( b_{\text{cur}}, [] )$
- $s_{\text{upper}}( b_{\text{cur}}, B_{\text{upper}} ) = b( b_{\text{cur}}, B_{\text{upper}} ) - b( b_{\text{cur}}, [] ) = 0$

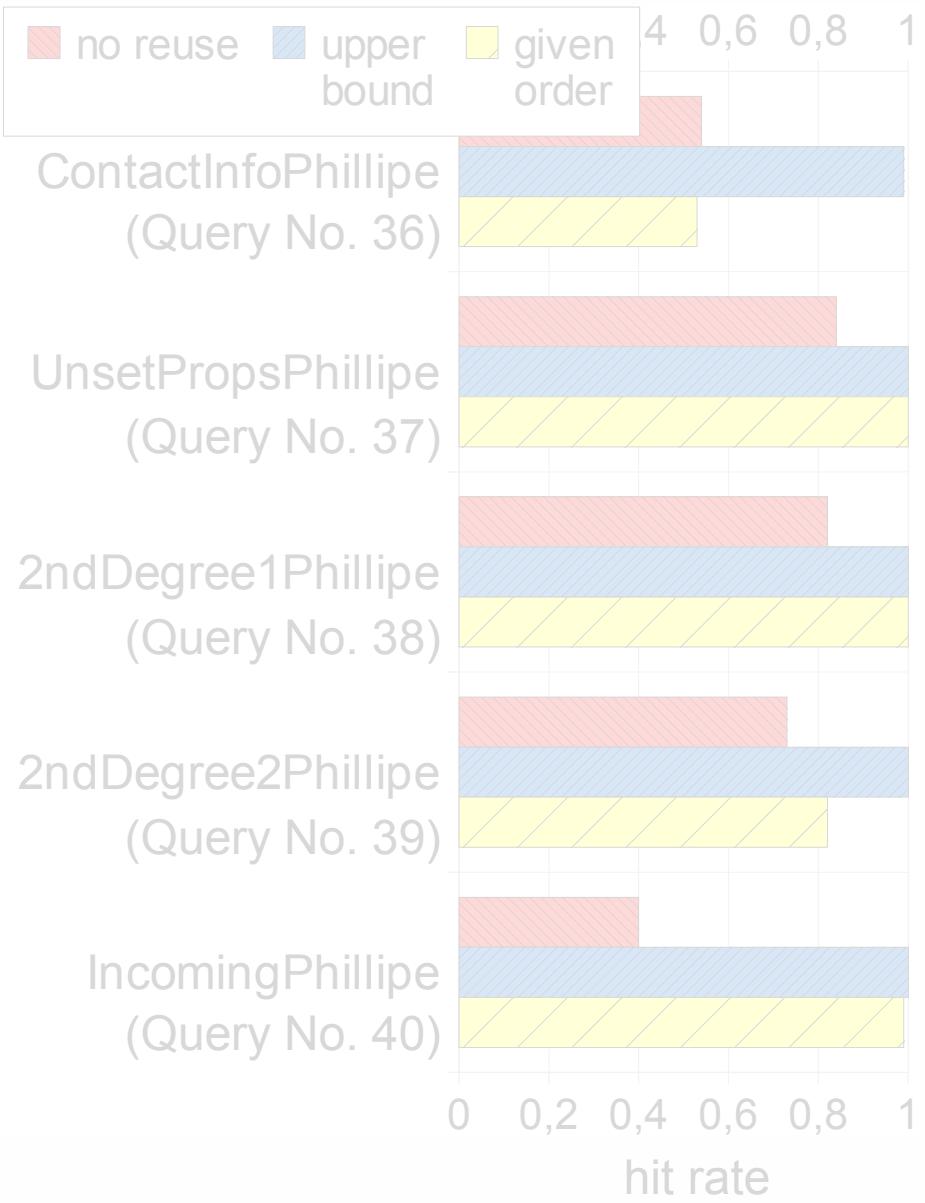
# Experiment – Single Query

Experiment	Avg. <sup>1</sup> number of Query Results (std.dev.)	Average <sup>1</sup> Hit Rate (std.dev.)	Avg. <sup>1</sup> query Execution Time (std.dev.)
no reuse	4.983 (11.658)	0.576 (0.182)	30.036 s (46.708)
upper bound	5.070 (11.813)	0.996 (0.017)	1.943 s (11.375)

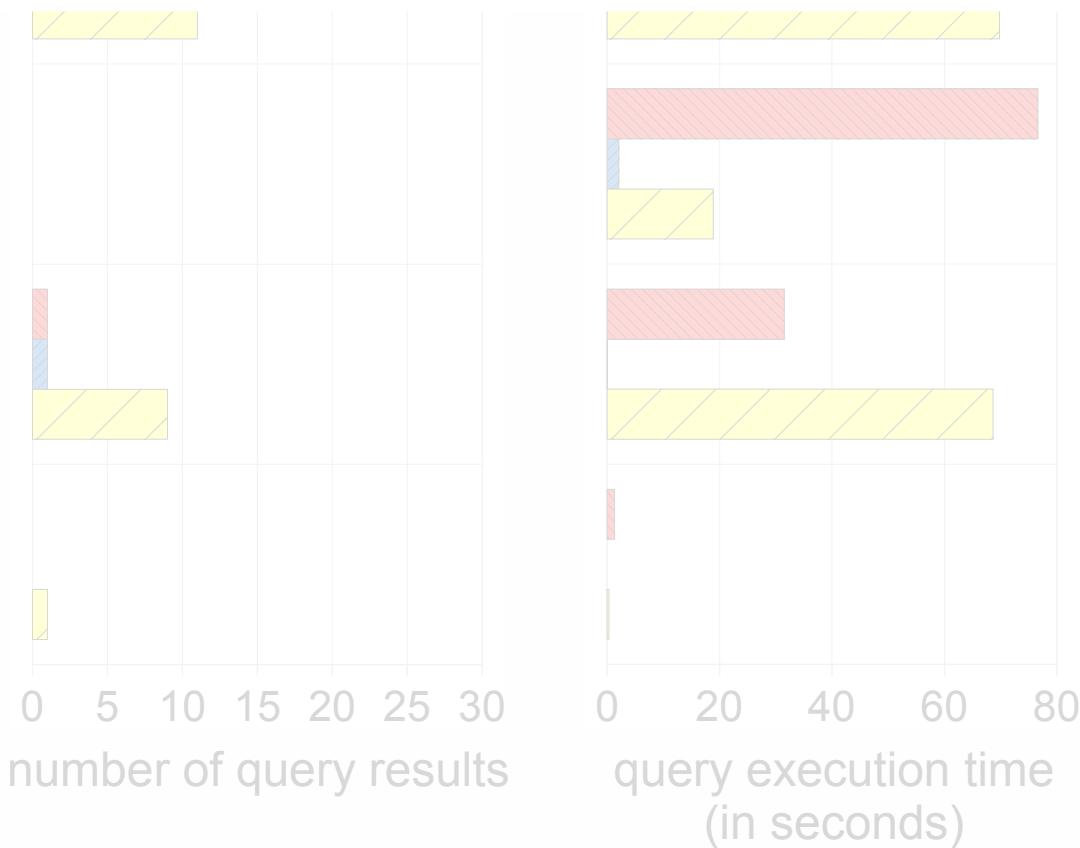
<sup>1</sup> Averaged over all 115 queries

- **Summary (measurement errors aside):**
  - Same number of query results
  - Significant improvements in query performance

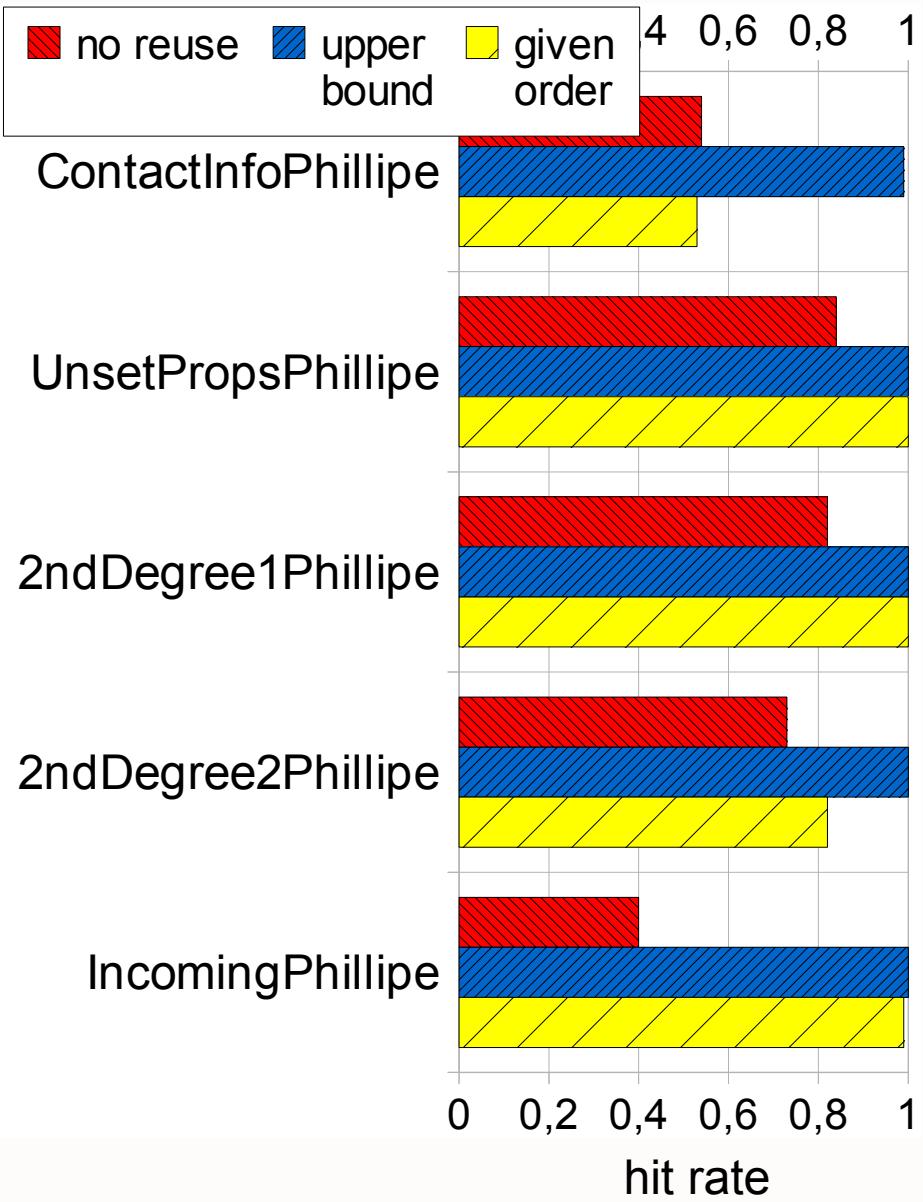
# Experiment – Complete Sequence



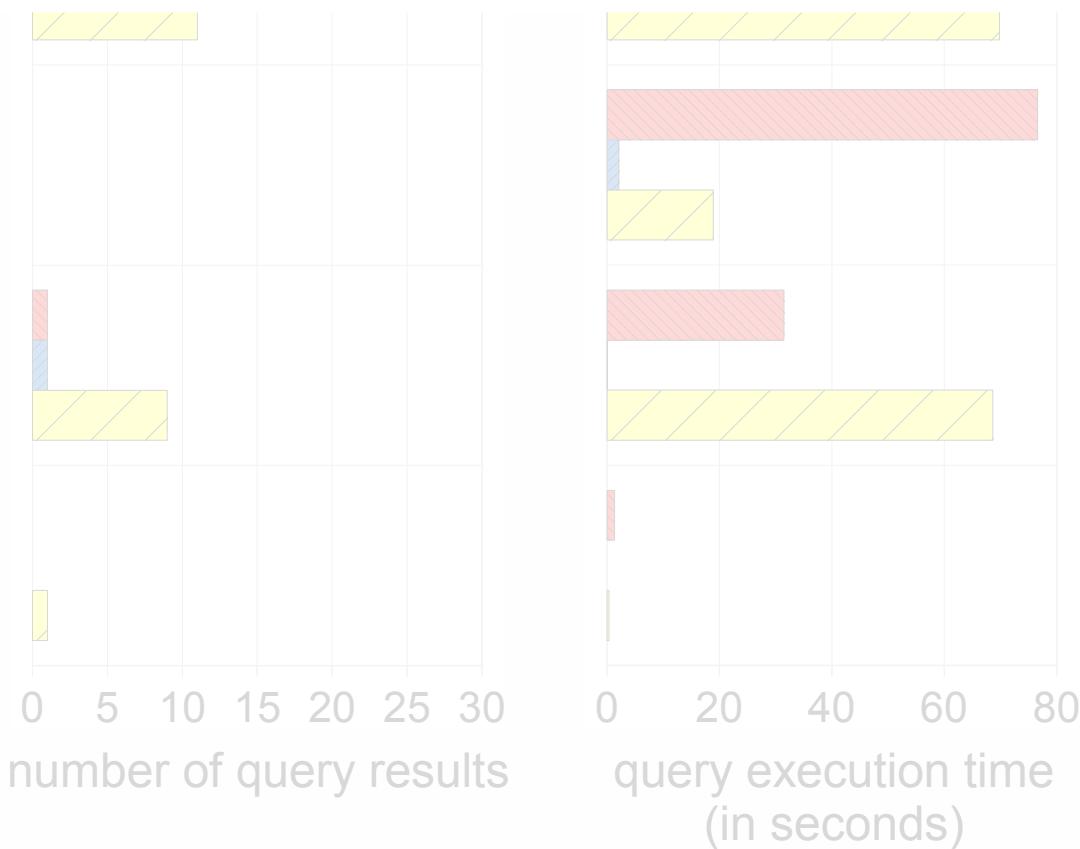
- ***given order experiment:***
- Reuse of the query-local dataset for the complete sequence of all 115 queries



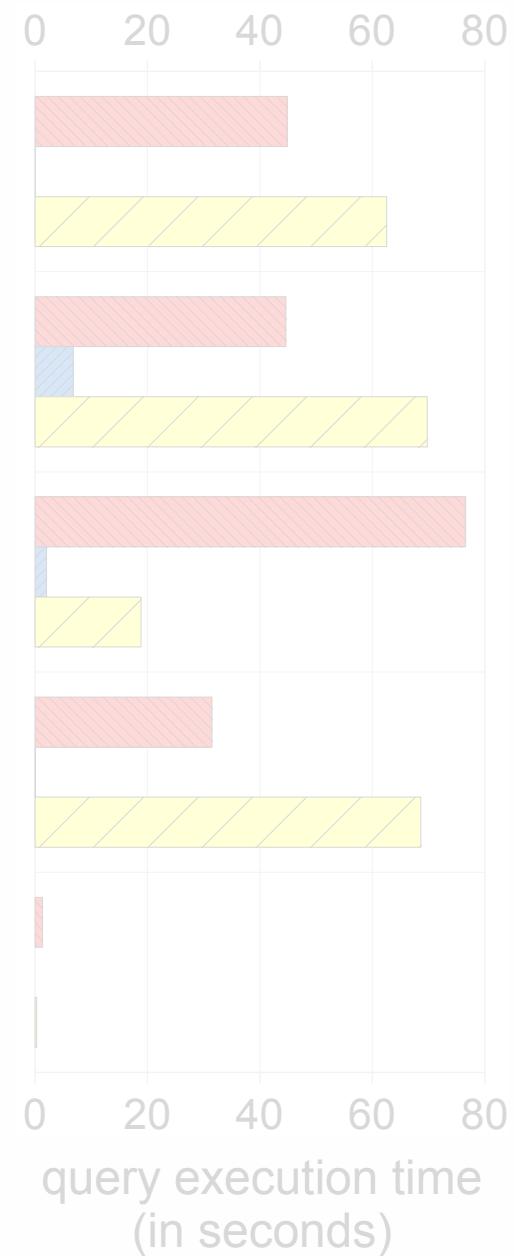
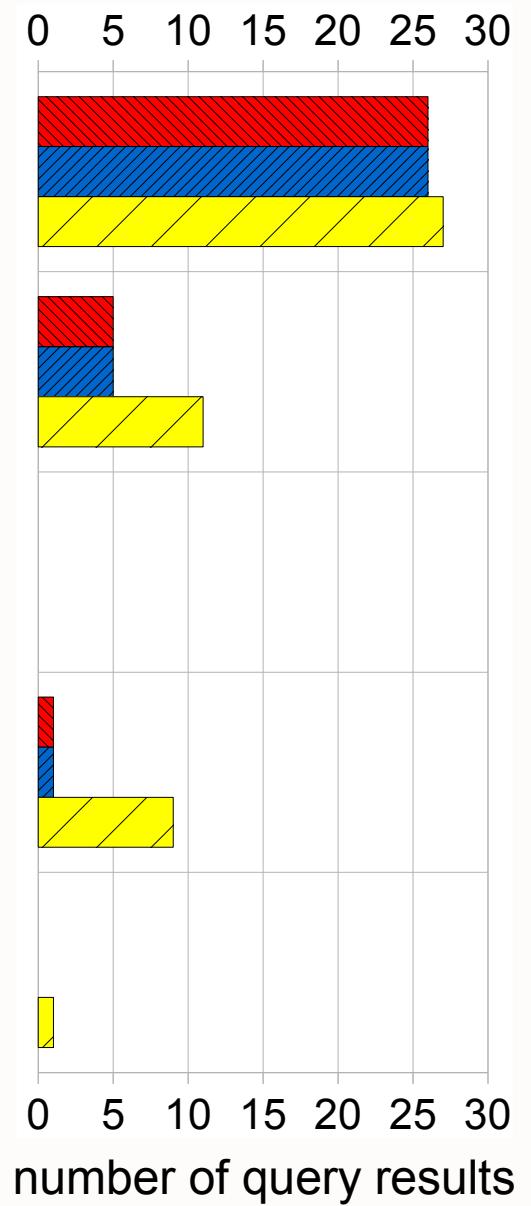
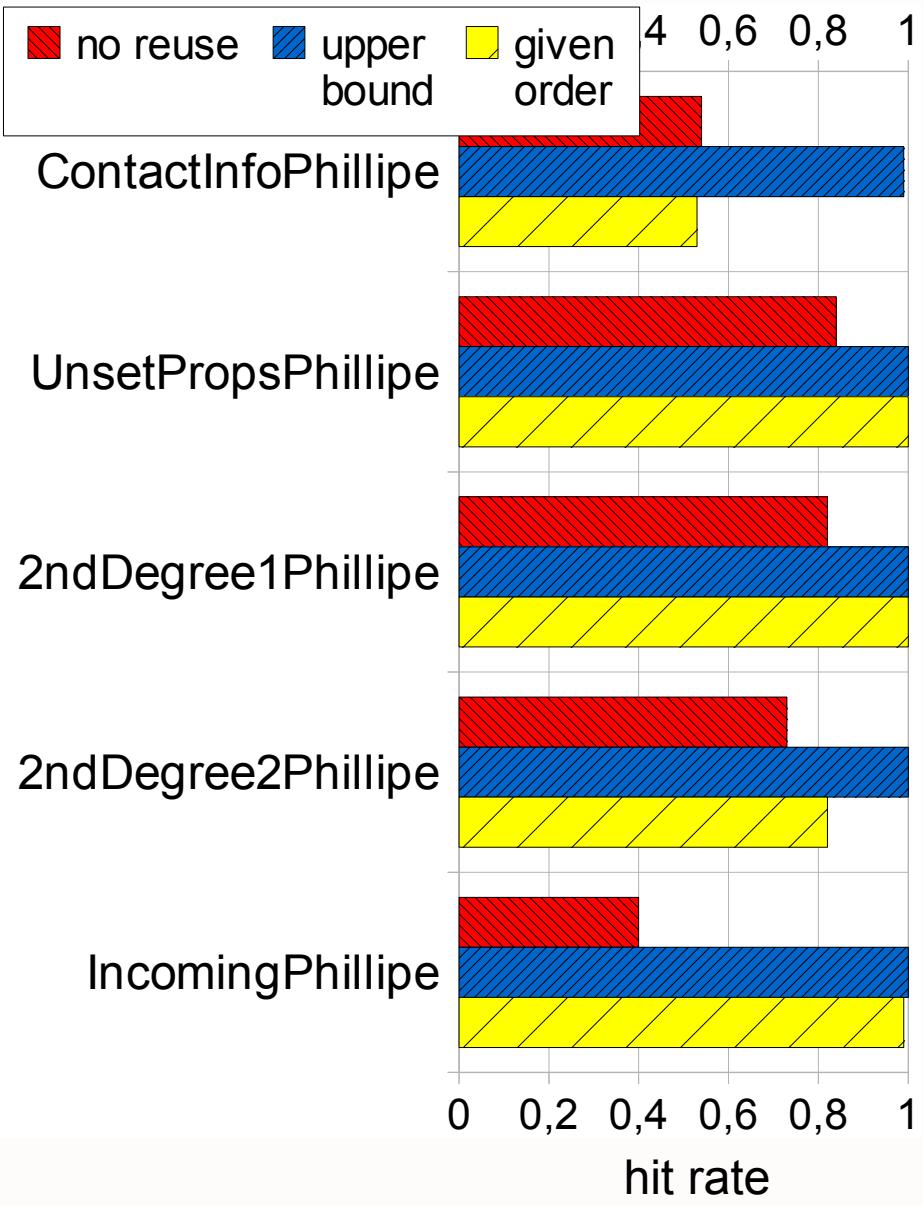
# Experiment – Complete Sequence



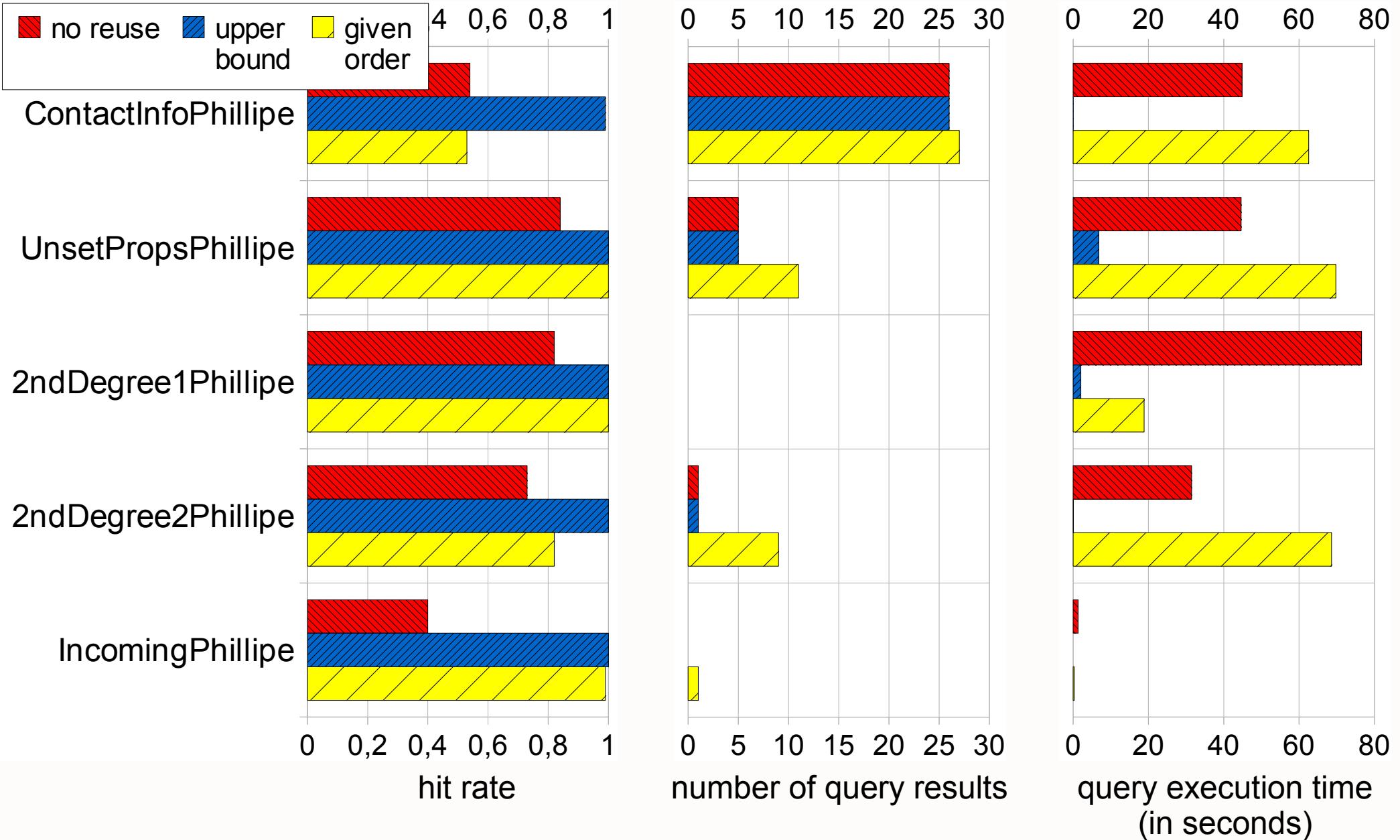
- ***given order experiment:***
- Reuse of the query-local dataset for the complete sequence of all 115 queries



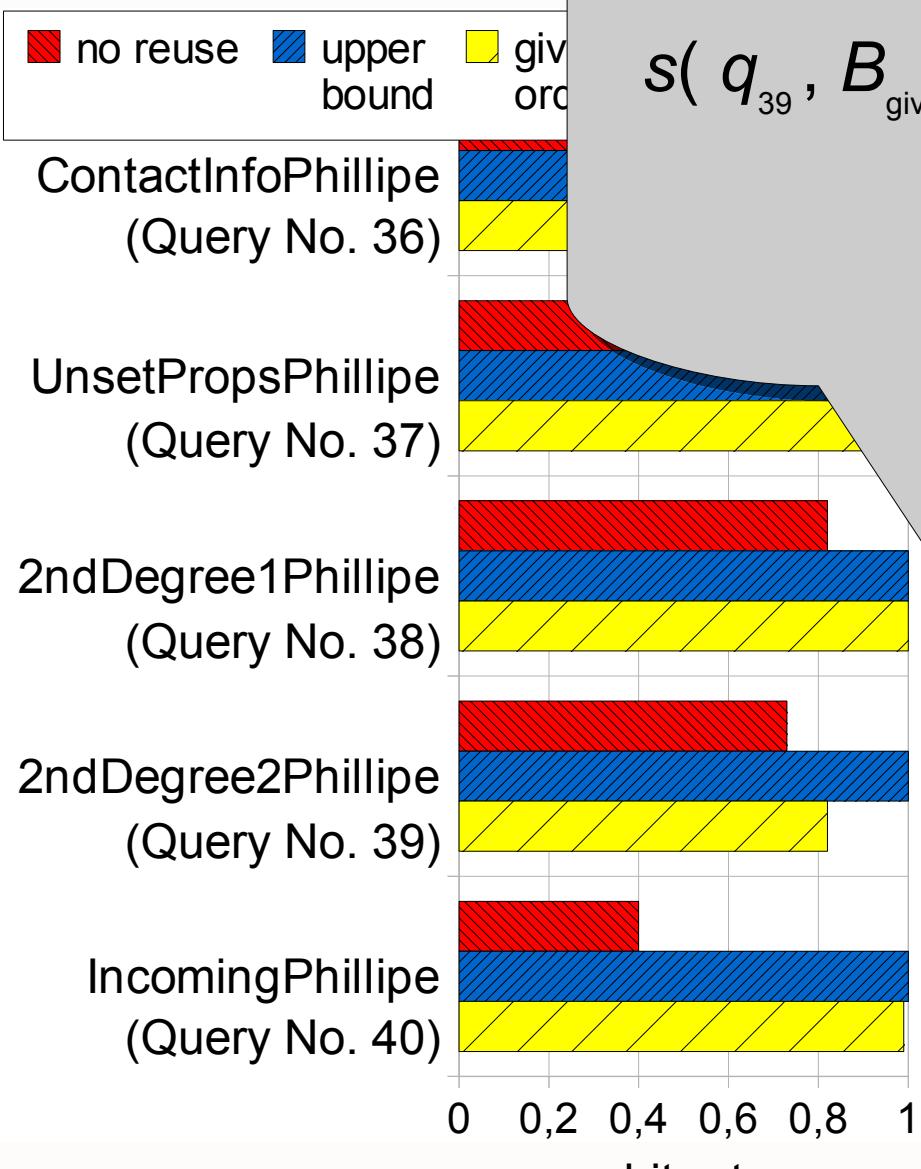
# Experiment – Complete Sequence



# Experiment – Complete Sequence

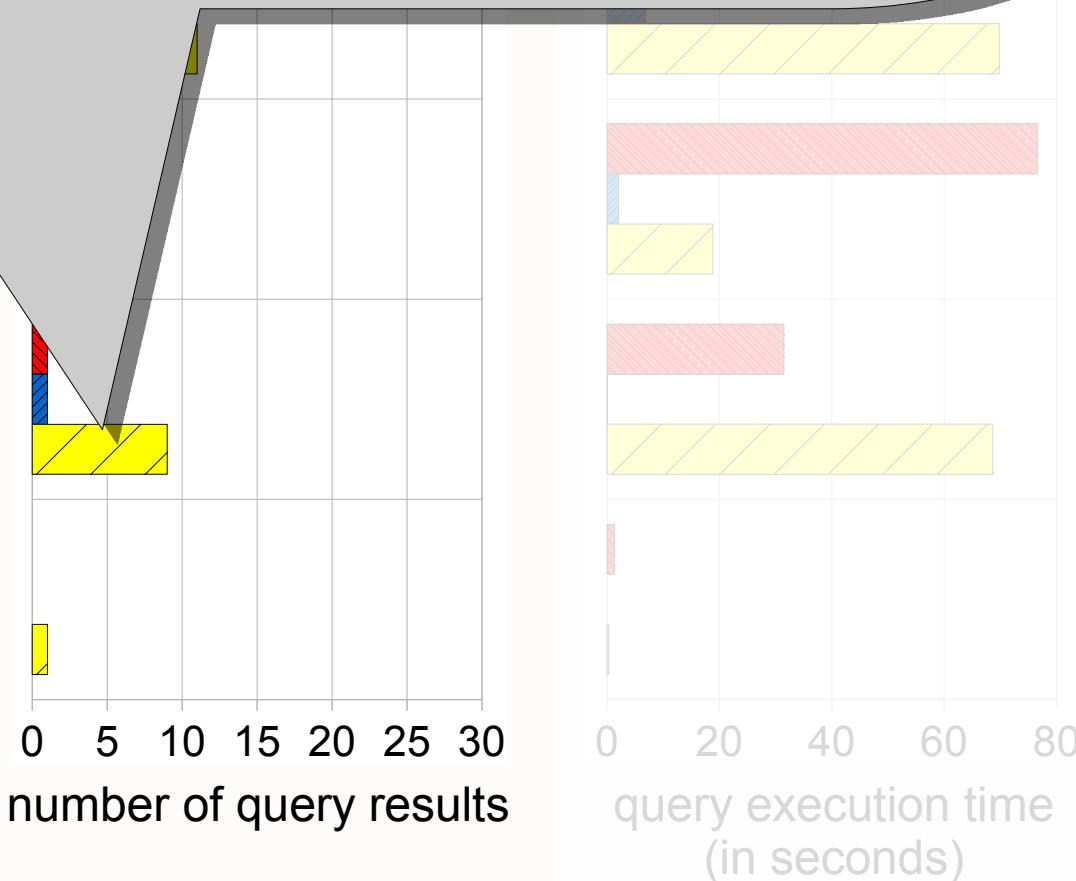


# Experiment - Complete Sequences



$$B_{\text{given order}} = [q_1, \dots, q_{38}]$$

$$\begin{aligned}
 s(q_{39}, B_{\text{given order}}) &= b(q_{39}, B_{\text{given order}}) - b(q_{39}, [ ]) \\
 &= 9 - 1 \\
 &= 8
 \end{aligned}$$



# Experiment Complete Sequence



$$B_{\text{given order}} = [ q_1, \dots, q_{38} ]$$

$$p'( q_{39}, B_{\text{given order}} ) = c'( q_{39}, [] ) - c'( q_{39}, B_{\text{given order}} )$$

$$= 31.48 \text{ s} - 68.64 \text{ s} \\ = -37.16 \text{ s}$$

■ no reuse ■ reuse

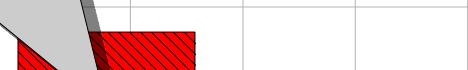
ContactInfoFor  
(Query N)

UnsetPropsPhillipe  
(Query No. 37)

2ndDegree1Phillipe  
(Query No. 38)

2ndDegree2Phillipe  
(Query No. 39)

IncomingPhillipe  
(Query No. 40)



# Experiment – Complete Sequence

Experiment	Avg. <sup>1</sup> number of Query Results (std.dev.)	Average <sup>1</sup> Hit Rate (std.dev.)	Avg. <sup>1</sup> query Execution Time (std.dev.)
no reuse	4.983 (11.658)	0.576 (0.182)	30.036 s (46.708)
upper bound	5.070 (11.813)	0.996 (0.017)	1.943 s (11.375)
given order	6.878 (12.158)	0.932 (0.139)	39.845 s (145.898)

<sup>1</sup> Averaged over all 115 queries

- **Summary:**

- Data cache may provide for additional query results
- Impact on performance may be positive but also negative

# Experiment – Complete Sequence



Experiment	Avg. <sup>1</sup> number of Query Results (std.dev.)	Average <sup>1</sup> Hit Rate (std.dev.)	Avg. <sup>1</sup> query Execution Time (std.dev.)
no reuse	4.983 (11.658)	0.576 (0.182)	30.036 s (46.708)
upper bound	5.070 (11.813)	0.996 (0.017)	1.943 s (11.375)
given order	6.878 (12.158)	0.932 (0.139)	39.845 s (145.898)
random orders	6.652 (11.966)	0.954 (0.036)	36.994 s (118.700)

- Executing the query sequence in a random order results in measurements similar to the given order.



**These slides have been created by  
Olaf Hartig**

<http://olafhartig.de>

**This work is licensed under a  
Creative Commons Attribution-Share Alike 3.0 License  
(<http://creativecommons.org/licenses/by-sa/3.0/>)**

