



## Systematic Generation of SPARQL Benchmark Queries for Linked Open Data

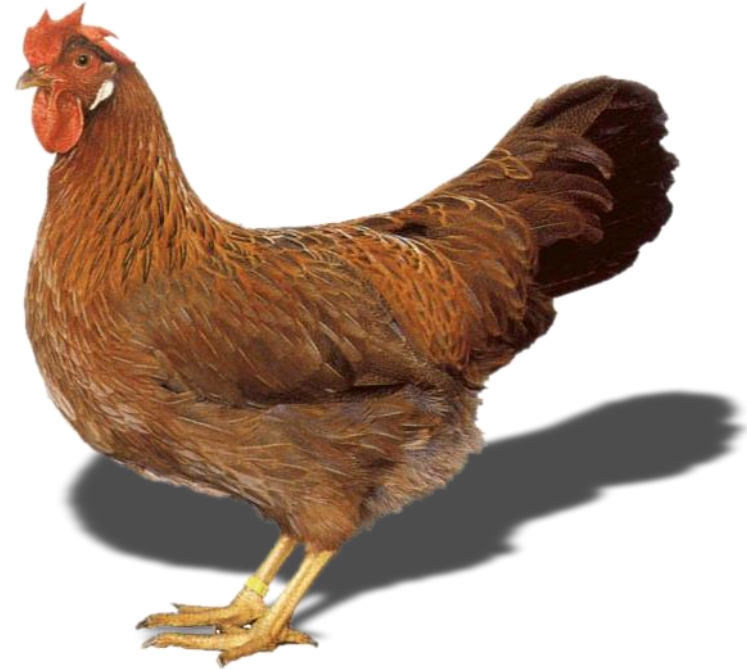
Olaf Görlitz, Matthias Thimm, Steffen Staab



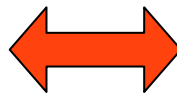
Why not use  
benchmark  
queries?



distributed  
queries



federation  
implementation



LUBM, BSBM, SP<sup>2</sup>B, ...

- Synthetic datasets
- Domain-specific
- Highly structured
- Sophisticated queries

FedBench (ISWC'11)

- 10 Linked Data sets (~170M triples)
- 25 handpicked distributed queries

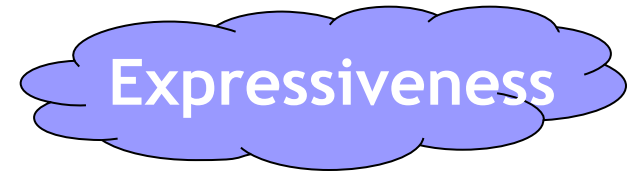
**Centralized**



**Fixed**

**Scalable, Flexible, Expressive  
Linked Data Benchmark**

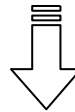
- ◆ Benchmark Idea
- ◆ Methodology
- ◆ Evaluation



Real Linked Data Sets

Customization

Typical+Complex Queries



Scalability

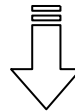
Flexibility

Expressiveness

Real Linked Data Sets

Customization

Typical+Complex Queries



Systematic SPARQL Benchmark Query Generator  
for Linked Open Data

Scalability

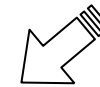
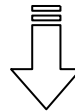
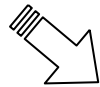
Flexibility

Expressiveness

Real Linked Data Sets

Customization

Typical+Complex Queries

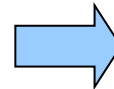


Systematic **SP**ARQL Benchmark Query **Ge**nerator  
for **L**inked **O**pen **D**ata



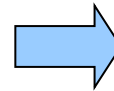
## What we want:

1. Define Query Characteristics



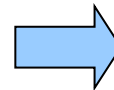
Customize Benchmark

2. Automatic Query Generation



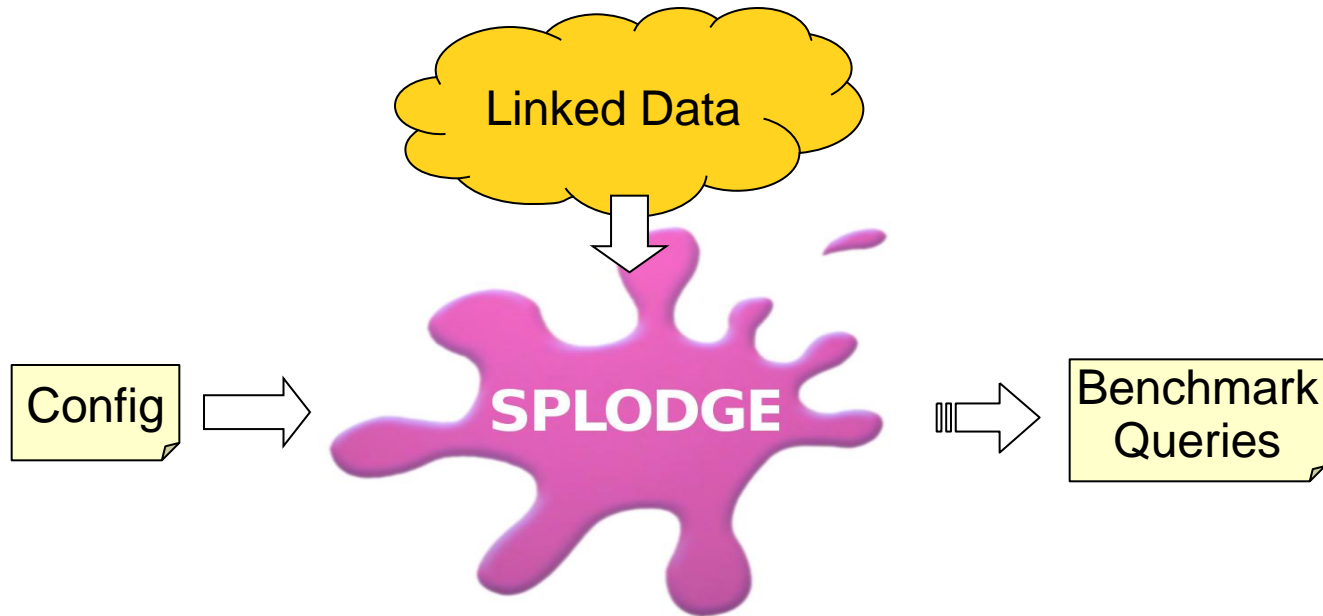
Random Queries

3. Query Validation



#results > 0

Methodology and toolset for systematic query generation



Parameterization

Query Generation

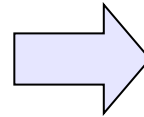
Query Validation

- ◆ Benchmark Idea
- ◆ Methodology
- ◆ Evaluation



Define typical + challenging distributed queries

No federation query logs available



Analyze queries of benchmarks

```
SELECT ?drug ?keggUrl ?chebiImage WHERE {  
  ?drug rdf:type drugbank:drugs .  
  ?drug drugbank:keggCompoundId ?keggDrug .  
  ?keggDrug bio2rdf:url ?keggUrl .  
  ?drug drugbank:genericName ?drugBankName .  
  ?chebiDrug purl:title ?drugBankName .  
  ?chebiDrug chebi:image ?chebiImage . }
```



DrugBank

FedBench/LifeScience#5



## Algebra

- Query Form  
(*Select, Construct, ...*)
- Join Type  
(*conj. / disj. / left-join*)
- Result Modifiers  
(*limit, offs, order by*)

## Structure

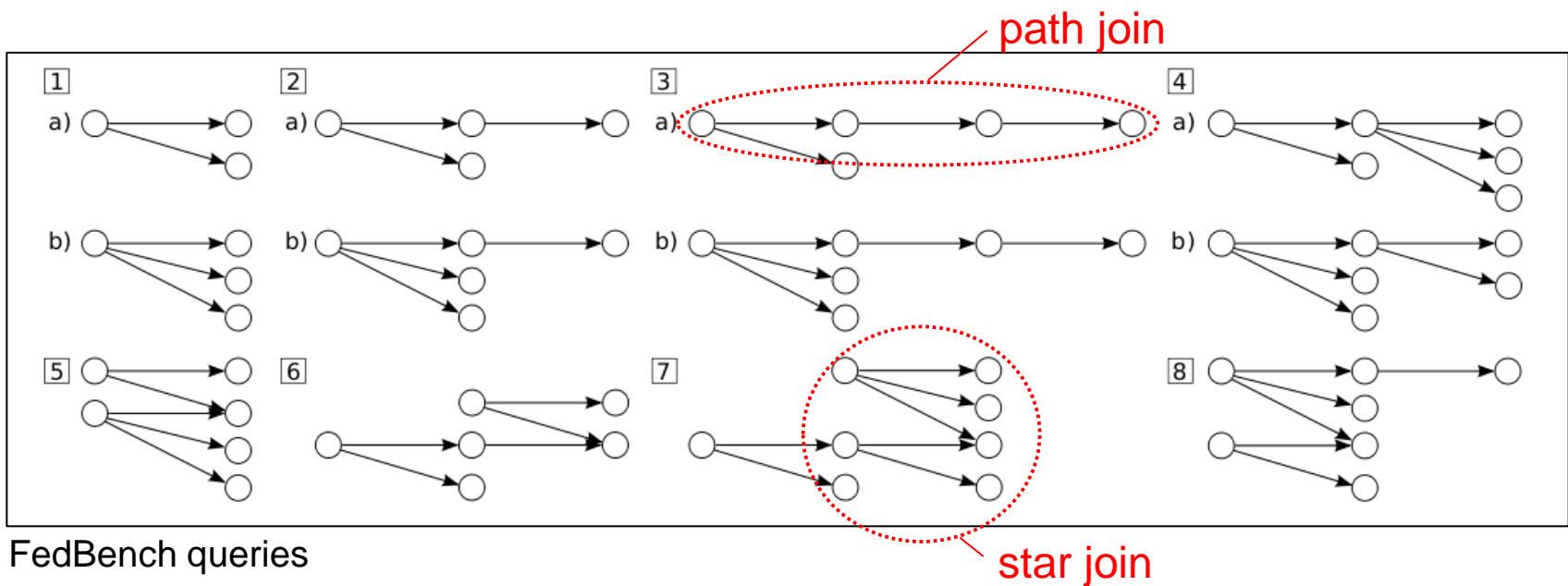
- Variable Patterns  
(*s, o, s+o, ...*)
- Join Patterns  
(*star, path*)
- Cross Product

## Cardinality

- # Data Sources
- # Joins/ Patterns
- # Results



Main query parameter: join structure

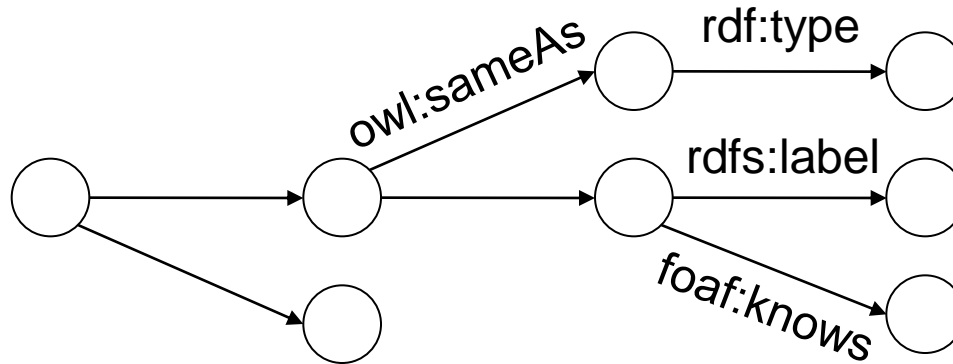
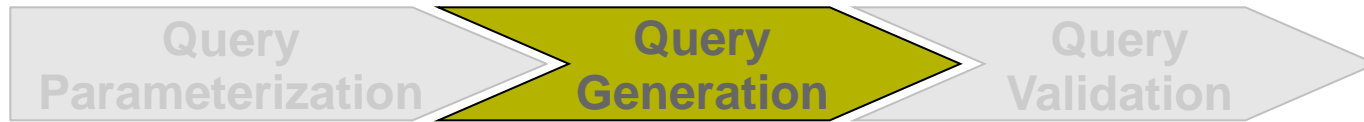




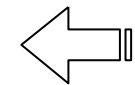
Additional query parameters: # triple patterns  
# data sources  
result size  
...

Path-join:  $n$  triple patterns,  
 $m$  sources ( $m \leq n$ )

Star-join:  $n$  triple pattern,  
anchor node (s/o)



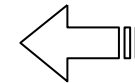
Iteratively add random triple pattern



#results > 0 ?



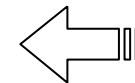
Need background knowledge



level of detail?

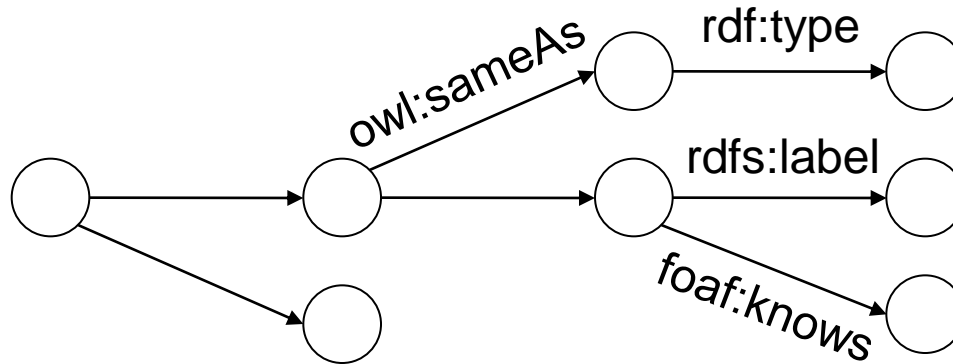


Predicate combinations



how provided?





## Linked Predicates

*(owl:sameAs → rdf:type)*

DBpedia → geonames (43, 58)

freebase → DBpedia (86, 72)

...

## Characteristics Sets\*

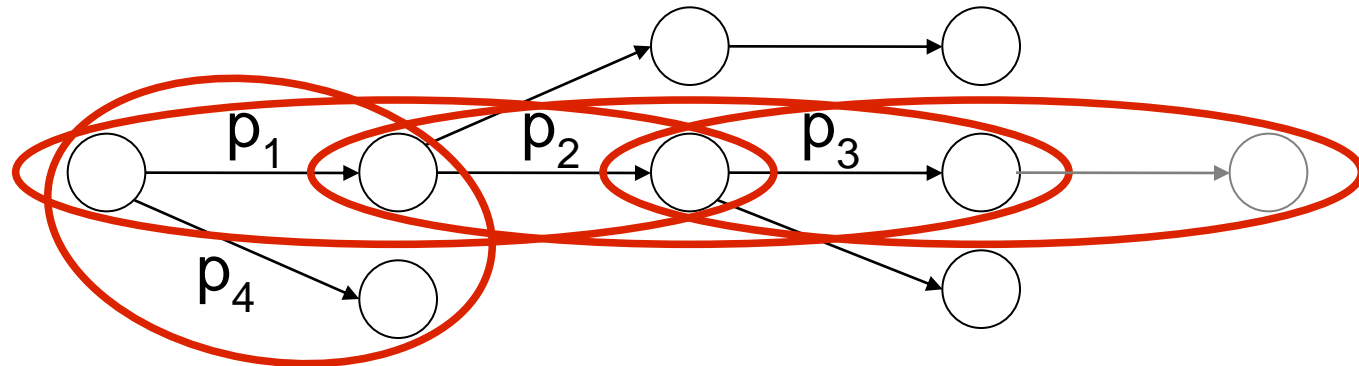
*{rdfs:label, foaf:knows, ...}*

DBpedia (322), rdfs:label (437)

foaf:knows (322)

...

\*[Neumann, Moerkotte, ICDE 2011]



## Linked Predicates

## Characteristics Sets

$$(p_1 \rightarrow p_2) \otimes (p_2 \rightarrow p_3)$$

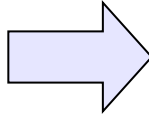
$$\{p_1, p_4\}$$

$$\otimes (p_3 \rightarrow p_i)$$

$$\{p_1, p_4, \dots\}$$



Verify generated queries (#results > 0)

How to evaluate?  Compute confidence value

**minimum join selectivity > e**

- ◆ Benchmark Idea
- ◆ Methodology
- ◆ Evaluation

- ◆ Verify generation of valid queries (#results >0)
- ◆ Compare variations of query generation algorithms

Baseline

“random”  
predicate

SPLODGE *lite*

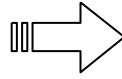
background  
knowledge

SPLODGE

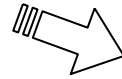
+ minimum  
join selectivity  
( $> 10^{-4}/10^{-3}/10^{-2}$ )

- ◆ Metrics:
  - ◆ #queries with non-empty results
  - ◆ #result per query

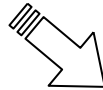
- ◆ Real Linked Data
- ◆ Random queries
- ◆ Triple Store



Billion Triple Challenge Dataset



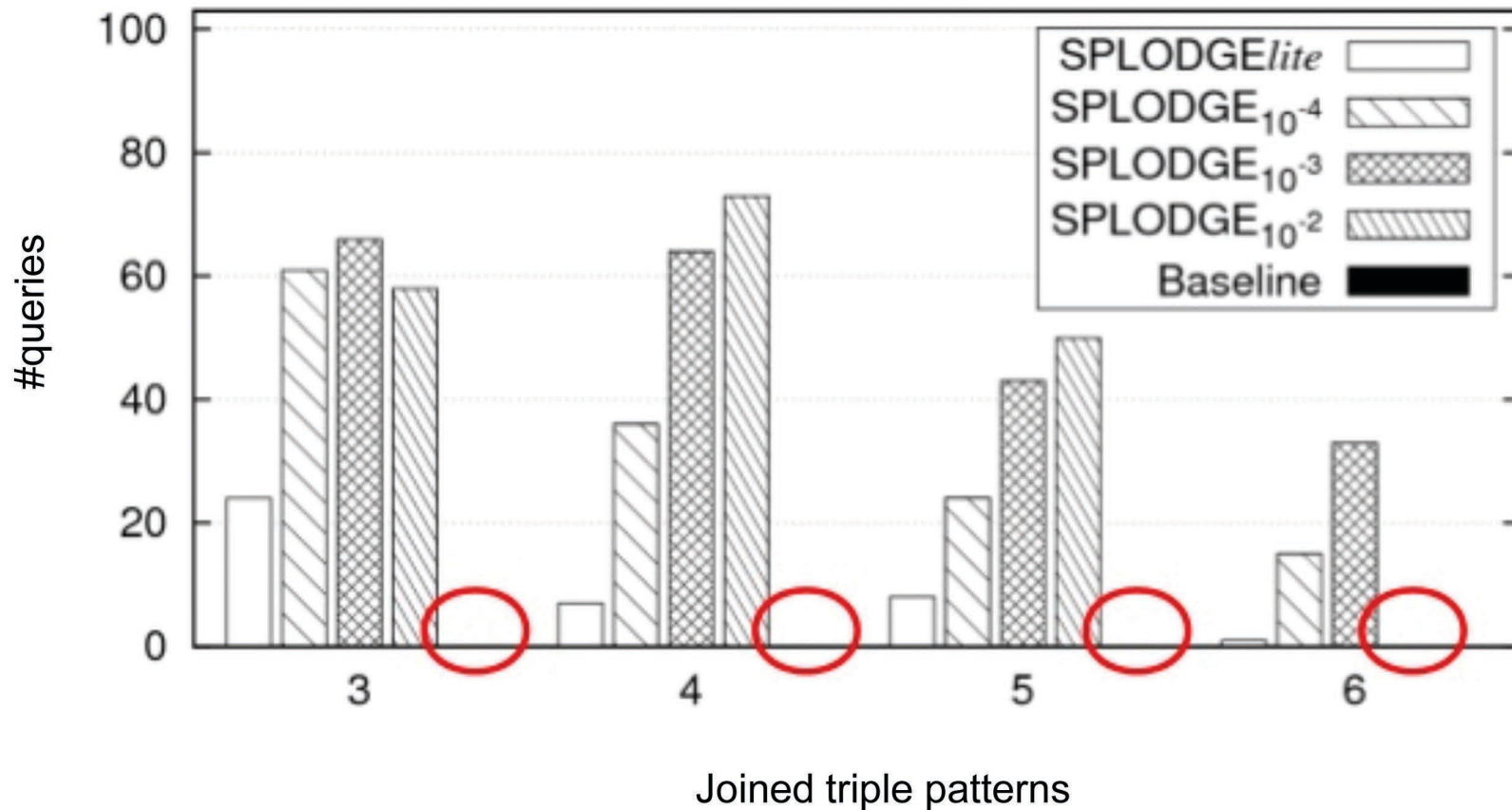
- Path-joins across data sources
- 3-6 patterns, bound predicates
- 100 queries per batch



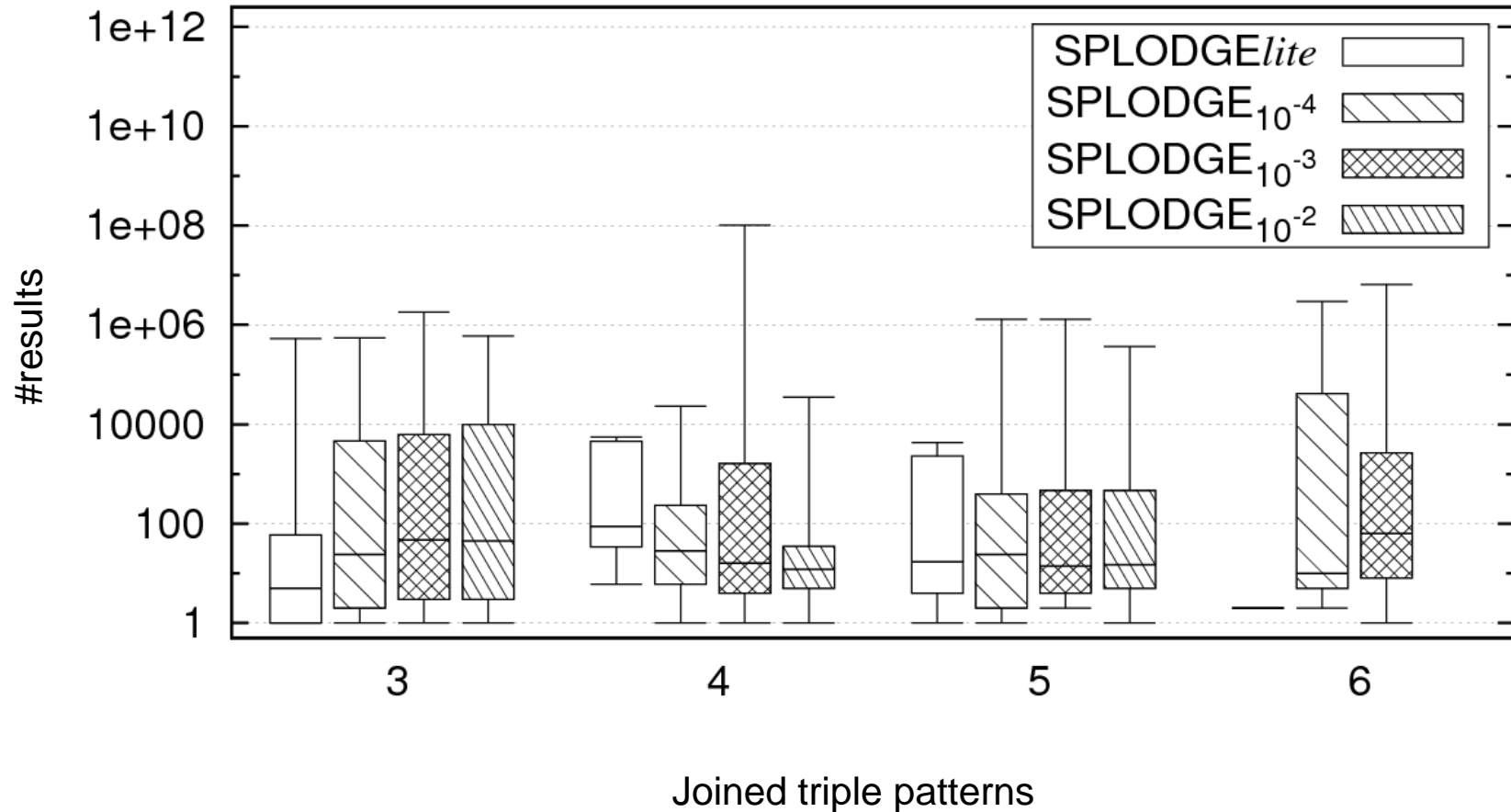
RDF3X

```
SELECT * WHERE {  
  ?var1 <http://dbpedia.org/property/description> ?var2 .  
  ?var2 <http://www.w3.org/1999/02/22-rdf-syntax-ns#type> ?var3 .  
  ?var3 <http://www.w3.org/2002/07/owl#disjointWith> ?var4 .  
  ?var4 <http://www.w3.org/2002/07/owl#disjointWith> ?var5 .  
  ?var5 <http://semantic-mediawiki.org/swivt/1.0#wikiPageModificationDate> ?var6  
}
```

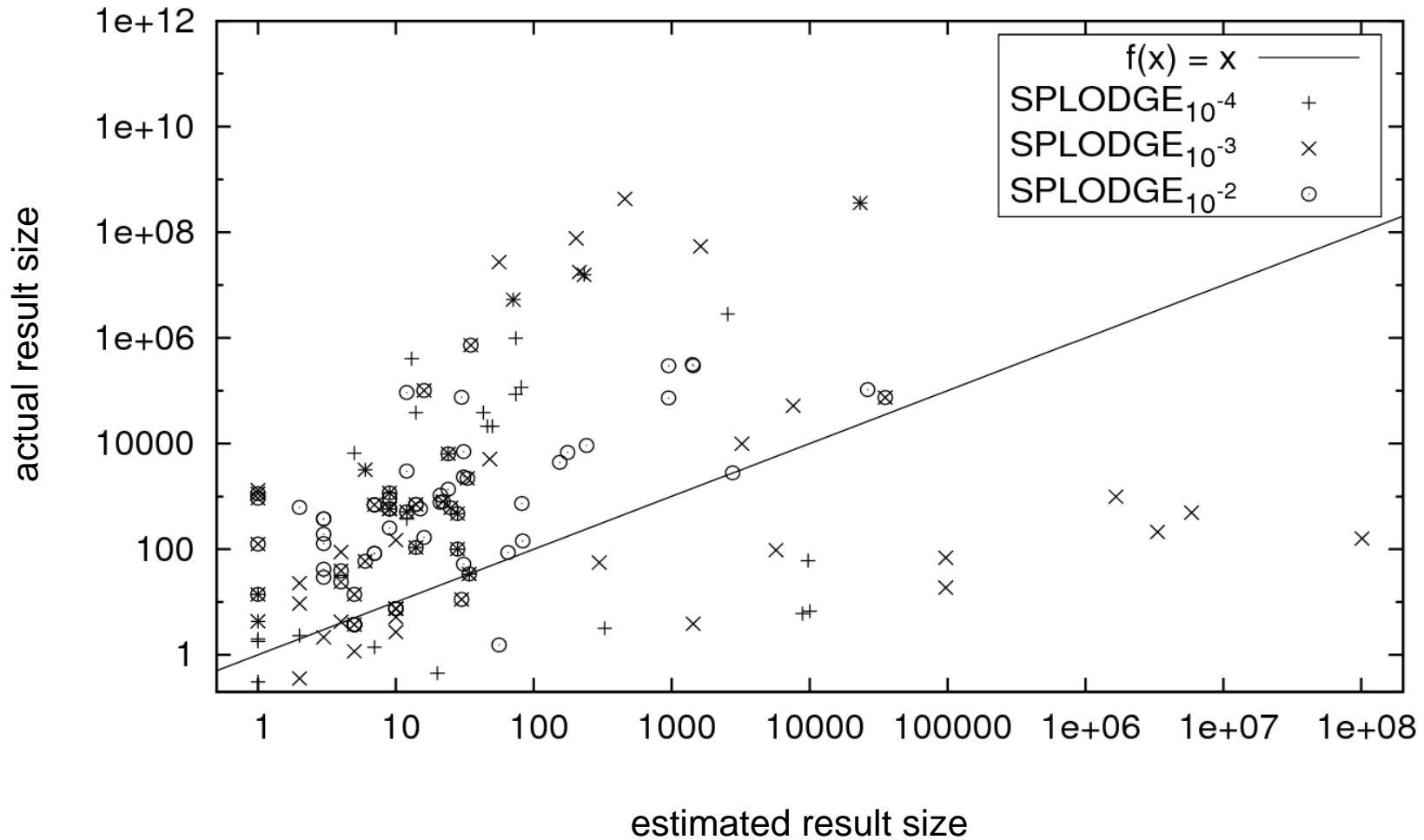
### Number of non-empty results for path-joins with 3-6 patterns



## Min/Max result sizes and quantiles for path-joins with 3-6 patterns







	all	SPLODGE <sub>lite</sub>	SPLODGE <sub>10<sup>-4</sup></sub>	SPLODGE <sub>10<sup>-3</sup></sub>	SPLODGE <sub>10<sup>-2</sup></sub>
rdf:type	19.1 %	15.7 % (2)	18.6 % (1)	21.8 % (1)	21.0 % (1)
rdfs:seeAlso	8.4 %	16.3 % (1)	6.2 % (3)	5.3 % (4)	4.6 % (3)
owl:disjointWith	5.4 %	2.6 % (6)	3.7 % (5)	7.0 % (2)	9.5 % (2)
rdfs:subClassOf	3.9 %	3.1 % (5)	2.4 % (8)	6.5 % (3)	3.5 % (5)
rdfs:isDefinedBy	3.7 %	1.7 % (9)	6.5 % (2)	4.0 % (5)	1.9 % (10)
owl:equivalentClass	2.7 %	2.3 % (7)	2.7 % (7)	3.2 % (6)	2.7 % (8)
foaf:primaryTopic	2.4 %	3.9 % (4)	0.8 % (16)	1.9 % (7)	3.1 % (7)
rdfs:label	1.9 %	1.5 % (10)	1.6 % (10)	1.7 % (9)	3.2 % (6)
skos:exactMatch	1.8 %	1.3 % (11)	4.3 % (4)	0.9 % (18)	–
owl:sameAs	1.7 %	1.9 % (8)	3.2 % (6)	0.7 % (21)	0.8 % (14)
...					
zemanta:targetType	1.2 %	0.1 % (91)	0.5 % (29)	1.3 % (14)	3.8 % (4)
distinct predicates	687	402	353	283	191
total predicates	6600	1800	1800	1800	1200

SPLODGE provides

- ◆ Flexible query characterization + parameterization
- ◆ Methodology for Systematic & Scalable Query Generation
- ◆ Toolset as Open Source (<http://code.google.com/p/splodge/>)

Future Work:

- ◆ Create a LOD Federation Benchmark
- ◆ Interactive SPARQL query construction

## Questions?

## BTC 2011 dataset in RDF3X

- ◆ pure triples, no context
- ◆ 160 GB repository file  
(14h loading, 200 GB tmp mem)

	BTC 2010	BTC 2011	BTC 2012
Total Size	624 GB	450 GB	303 GB
Total Quads	3.171.793.030	2.178.395.469	1.436.545.545
Unique Quads	3.154.896.097	2.145.122.248	1.311.765.894
Unique Triples	1.426.831.520	1.968.347.976	1.056.184.911
Contexts/Documents	8.132.721	7.423.477	9.283.829
Common Domains	22.299	789	837
Types	168.482	314.448	296.607
Predicates	95.589	47.738	57.257

