

Optimization of SPARQL Queries Processing Using PRoST

Guilherme Schievelbein

Databases and Information Systems
Master's Thesis, March 2019

Concepts

Partitioned RDF on Spark Tables (PRoST)

Joined Wide Property Table

Dynamic ExtVP Database

Evaluation

Discussion

Acknowledgments

Concepts

- Resource Description Framework

- SPARQL

- SPARQL

- Data Models

Partitioned RDF on Spark Tables (PProST)

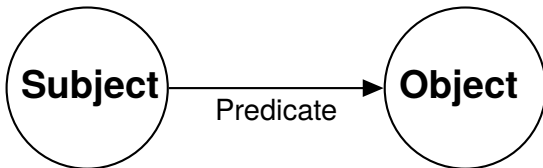
Joined Wide Property Table

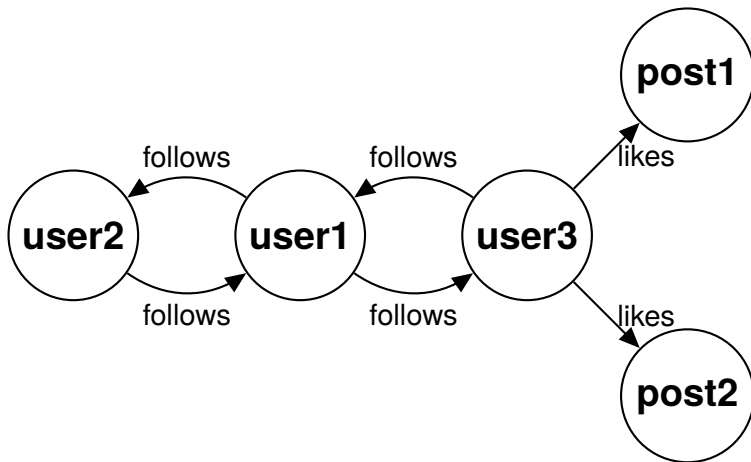
Dynamic ExtVP Database

Evaluation

Discussion

Acknowledgments





```
SELECT ?user
      ?post
WHERE {
    <user1> <follows> ?user .
    ?user   <likes>   ?post
}
```

```
SELECT ?user  
      ?post  
WHERE {  
  <user1> <follows> ?user .  
  ?user   <likes>   ?post  
}
```

?user	?post
user3	post1
user3	post2

Triplestore		
subject	predicate	object
user1	follows	user2
user1	follows	user3
user2	follows	user1
user3	follows	user1
user3	likes	post1
user3	likes	post2

follows	
subject	object
user1	user2
user1	user3
user2	user1
user3	user1

likes	
subject	object
user3	post1
user3	post2

Wide Property Table

Wide Property Table		
subject	follows	likes
user1	[user2, user3]	NULL
user2	user1	NULL
user3	user1	[post1, post2]

Concepts

Partitioned RDF on Spark Tables (PRoST)

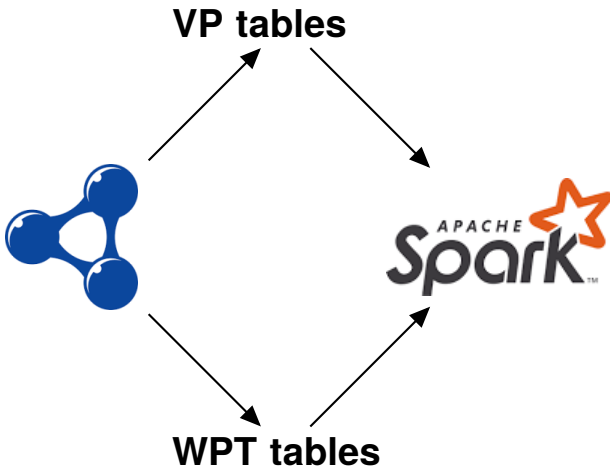
Joined Wide Property Table

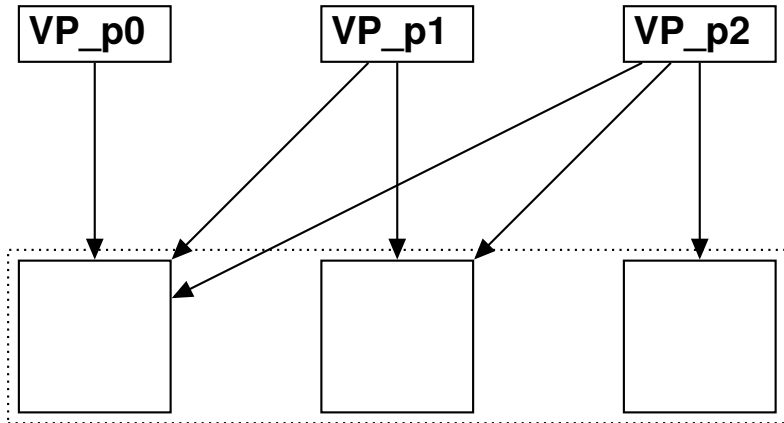
Dynamic ExtVP Database

Evaluation

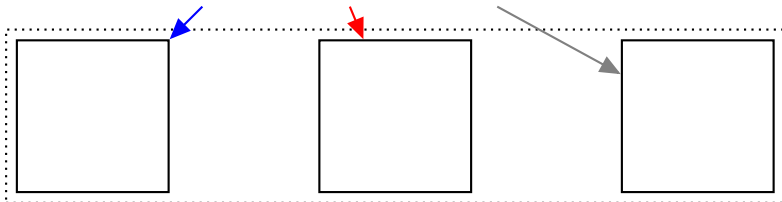
Discussion

Acknowledgments



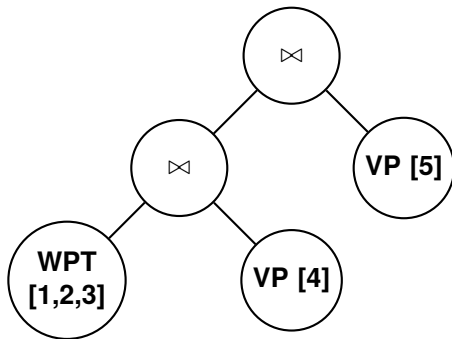


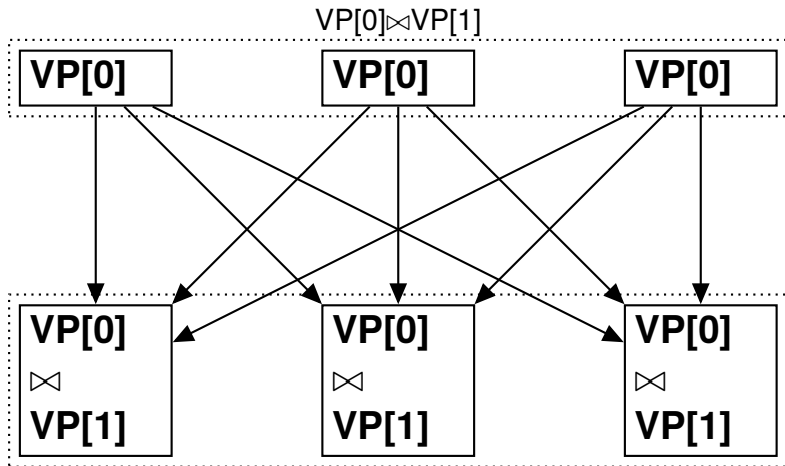
Wide Property Table		
subject	follows	likes
user1	[user2, user3]	NULL
user2	user1	NULL
user3	user1	[post1, post2]

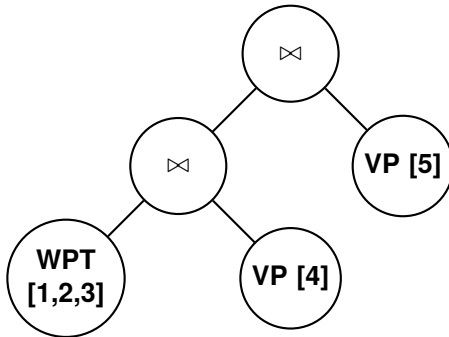


- 1 ?v0 <p0> <l0> .
- 2 ?v0 <p1> <l1> .
- 3 ?v0 <p2> ?v1 .
- 4 ?v1 <p3> ?v2 .
- 5 ?v2 <p4> ?v3

- 1 ?v0 <p0> <l0> .
- 2 ?v0 <p1> <l1> .
- 3 ?v0 <p2> ?v1 .
- 4 ?v1 <p3> ?v2 .
- 5 ?v2 <p4> ?v3







Concepts

Partitioned RDF on Spark Tables (PProST)

Joined Wide Property Table

- Supported Wide Property Tables

- Joined Wide Property Table

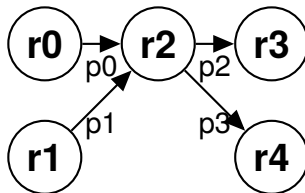
Dynamic ExtVP Database

Evaluation

Discussion

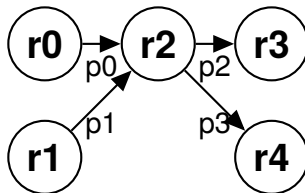
Acknowledgments

Wide Property Table and Inverse Wide Property Table



Wide Property Table				
s	p0	p1	p2	p3
r0	r2	NULL	NULL	NULL
r1	NULL	r2	NULL	NULL
r2	NULL	NULL	r3	r4

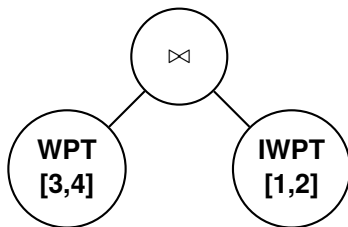
Wide Property Table and Inverse Wide Property Table



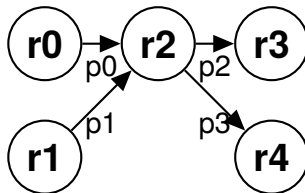
Inverse Wide Property Table				
o	p0	p1	p2	p3
r2	r0	r1	NULL	NULL
r3	NULL	NULL	r2	NULL
r4	NULL	NULL	NULL	r2

Wide Property Table and Inverse Wide Property Table

- 1 ?v0 <p0> **?v2** .
- 2 ?v1 <p1> **?v2** .
- 3 **?v2** <p3> ?v3 .
- 4 **?v2** <p4> ?v4 .



Joined Wide Property Table



Joined Wide Property Table								
p3-i	p2-i	p1-i	p0-i	r	p0	p1	p2	p3
NULL	NULL	NULL	NULL	r0	r2	NULL	NULL	NULL
NULL	NULL	NULL	NULL	r1	NULL	r2	NULL	NULL
NULL	NULL	r1	r0	r2	NULL	NULL	r3	r4
NULL	r2	NULL	NULL	r3	NULL	NULL	NULL	NULL
r2	NULL	NULL	NULL	r4	NULL	NULL	NULL	NULL

Joined Wide Property Table

- 1 ?v0 <p0> **?v2** .
- 2 ?v1 <p1> **?v2** .
- 3 **?v2** <p3> ?v3 .
- 4 **?v2** <p4> ?v4 .



Concepts

Partitioned RDF on Spark Tables (PProST)

Joined Wide Property Table

Dynamic ExtVP Database

- S2RDF

- Dynamic ExtVP Database

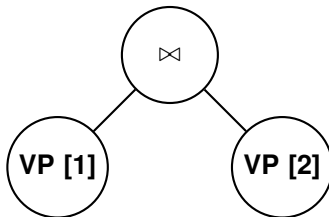
- Dynamic ExtVP Database

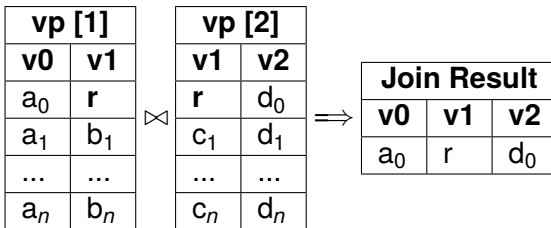
Evaluation

Discussion

Acknowledgments

- 1 ?v0 <p0> ?v1 .
- 2 ?v1 <p1> ?v2



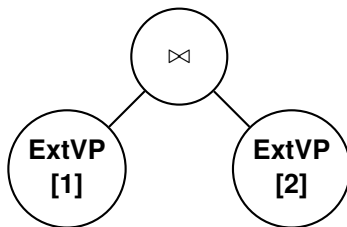


vp_p0		vp_p1	
s	o	s	o
a_0	r	r	d_0
a_1	b_1	c_1	d_1
...
a_n	b_n	c_n	d_n

\Rightarrow

extvp_OS_p0 p1		extvp_SO_p1 p0	
s	o	s	o
a_0	r	r	d_0

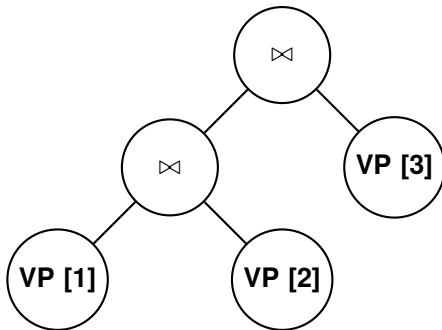
- 1 ?v0 <p0> ?v1 .
- 2 ?v1 <p1> ?v2



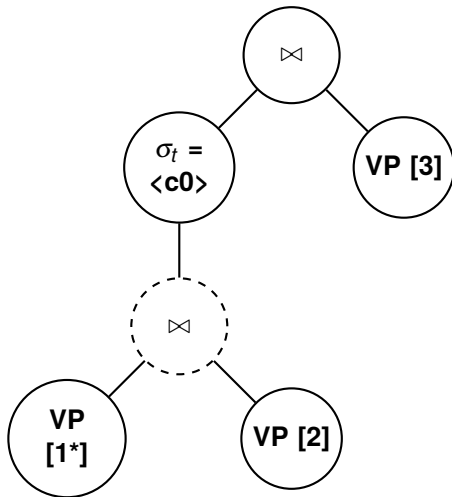
100 million triples:

- **PRoST**: less than 1h
- **S2RDF**: more than 16h

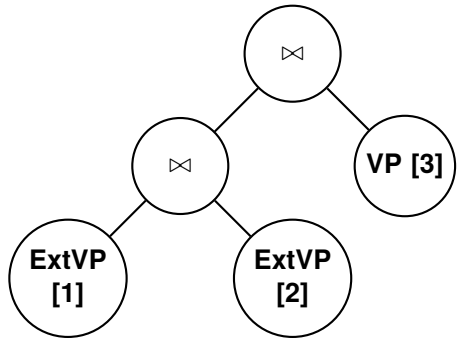
- 1 $\langle c0 \rangle \langle p0 \rangle ?v1 .$
- 2 $?v1 \langle p1 \rangle ?v2 .$
- 3 $?v2 \langle p2 \rangle ?v3 .$



- 1 $\langle c0 \rangle \langle p0 \rangle ?v1 .$
- 2 $?v1 \langle p1 \rangle ?v2 .$
- 3 $?v2 \langle p2 \rangle ?v3 .$



- 1 $\langle c1 \rangle \langle p0 \rangle ?v1 .$
- 2 $?v1 \langle p1 \rangle ?v2 .$
- 3 $?v2 \langle p2 \rangle ?v3 .$



Concepts

Partitioned RDF on Spark Tables (P_{Ro}ST)

Joined Wide Property Table

Dynamic ExtVP Database

Evaluation

- Benchmark Environment

- JWPT Evaluation

- Dynamic ExtVP Evaluation

Discussion

Acknowledgments

10 machines:

- 6 Core Intel Xeon E5-2420
- 32GB RAM
- 4TB HD
- Spark 2.2.0

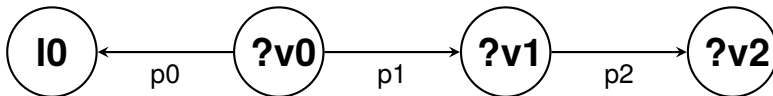
WatDiv synthetic dataset:

- 100 million triples
- 5 million subjects
- 80 distinct predicates
- 9 million objects

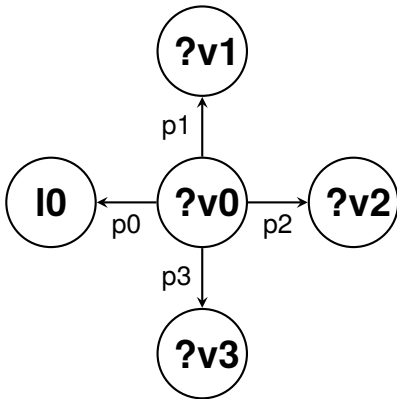
88 queries:

- 25 linear-shaped queries
- 35 star-shaped queries
- 25 snow-shaped queries
- 3 complex-shaped queries

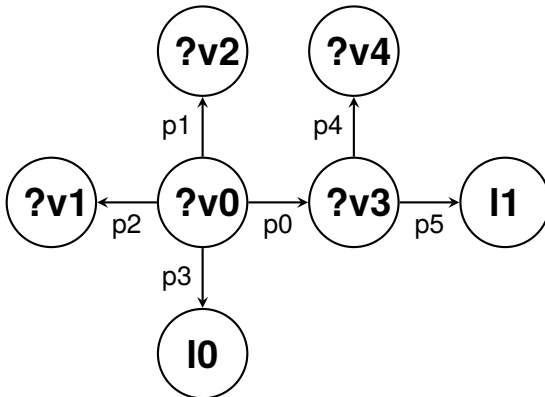
Linear-shaped query



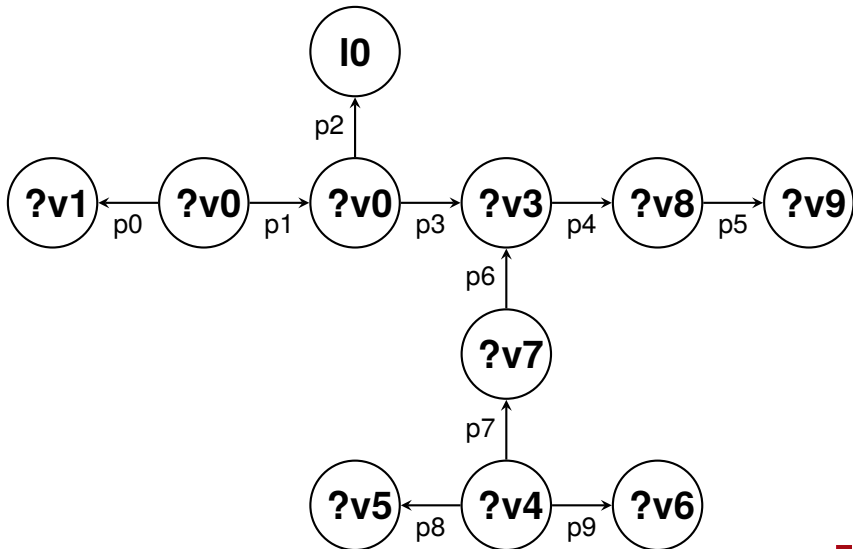
Star-shaped query



Snow-shaped query

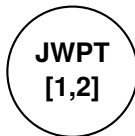


Complex-shaped query



- VP + (I)WPT (**baseline**)
- VP + JWPT
- VP + (I/J)WPT

- 1 ?v0 <p0> ?v1 .
- 2 ?v0 <p1> ?v2 .



Joined Wide Property Table

p3-i	p2-i	p1-i	p0-i	r	p0	p1	p2	p3
NULL	NULL	NULL	NULL	r0	r2	NULL	NULL	NULL
NULL	NULL	NULL	NULL	r1	NULL	r2	NULL	NULL
NULL	NULL	r1	r0	r2	NULL	NULL	r3	r4
NULL	r2	NULL	NULL	r3	NULL	NULL	NULL	NULL
r2	NULL	NULL	NULL	r4	NULL	NULL	NULL	NULL

- 1 ?v0 <p0> ?v1 .
- 2 ?v0 <p1> ?v2 .



Wide Property Table				
s	p0	p1	p2	p3
r0	r2	NULL	NULL	NULL
r1	NULL	r2	NULL	NULL
r2	NULL	NULL	r3	r4

Less join operations than the baseline:

- linear-shaped: 1.5s faster
- star-shaped: 200ms - 1.7s faster
- snow-shaped: 5s - 8s faster

Same number of join operations as the baseline:

- 200ms - 1.1s slower

Less join operations than the baseline:

- linear-shaped: 1.5s faster
- star-shaped: 200ms - 1.7s faster
- snow-shaped: 5s - 8s faster

Same number of join operations as the baseline:

- no statistically significant difference

VP + JWPT:

- Less join operations => faster processing
- More tuples => slower processing

VP + (I/J)WPT:

- At least as fast as the baseline

- VP (**baseline**)
- Dynamic ExtVP
- Dynamic ExtVP (table creation only)

Dynamic ExtVP:

■ 2s - 15s faster

Dynamic ExtVP (table creation only):

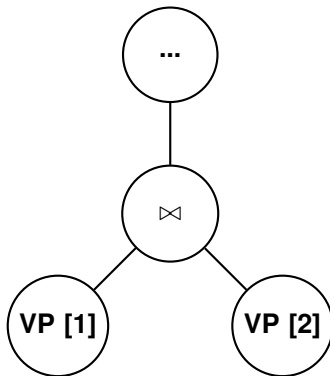
■ Two ExtVP tables created: no statistically significant difference

■ More tables created: 1.5s to 5.9s slower

- On average, faster processing of queries
- Slower processing, when many ExtVP tables are created in a single execution
- Tables only need to be created once
- Faster loading time than S2RDF
- Slower processing than P_{Ro}ST with VP +(I)WPT

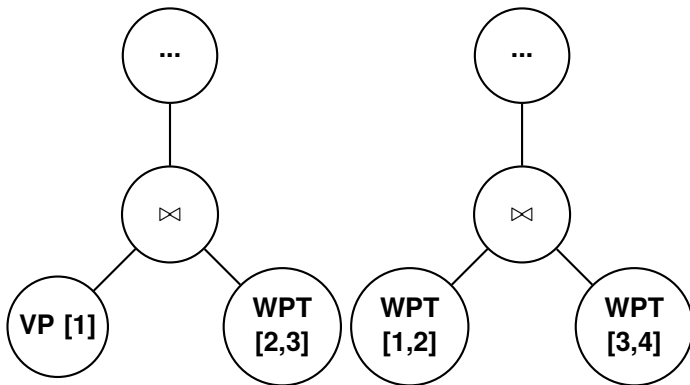
- VP + (I)WPT (**baseline**)
- Dynamic ExtVP + (I)WPT
- Dynamic ExtVP (table creation only) + (I)WPT

No statistically significant execution time difference.



=> Two ExtVP tables created

Why?



Concepts

Partitioned RDF on Spark Tables (PRoST)

Joined Wide Property Table

Dynamic ExtVP Database

Evaluation

Discussion

Dynamic ExtVP

Joined Wide Property Table

Acknowledgments

- Reduced tables are faster
- No extra loading cost
- Creation time of tables is not very high

Future: persist reduced tables from join operations between arbitrary node types

- Faster when it reduces the number of required join operations
- Not always efficient by itself

Future: reduce the number of tuples in the JWPT

Joined Wide Property Table

p3-i	p2-i	p1-i	p0-i	r	p0	p1	p2	p3
NULL	NULL	NULL	NULL	r0	r2	NULL	NULL	NULL
NULL	NULL	NULL	NULL	r1	NULL	r2	NULL	NULL
NULL	NULL	r1	r0	r2	NULL	NULL	r3	r4
NULL	r2	NULL	NULL	r3	NULL	NULL	NULL	NULL
r2	NULL	NULL	NULL	r4	NULL	NULL	NULL	NULL

Joined Wide Property Table

p3-i	p2-i	p1-i	p0-i	r	p0	p1	p2	p3
NULL	NULL	r1	r0	r2	NULL	NULL	r3	r4

Concepts

Partitioned RDF on Spark Tables (PRoST)

Joined Wide Property Table

Dynamic ExtVP Database

Evaluation

Discussion

Acknowledgments

Thank you.

Guilherme Schievelbein
schieveg@tf.uni-freiburg.de