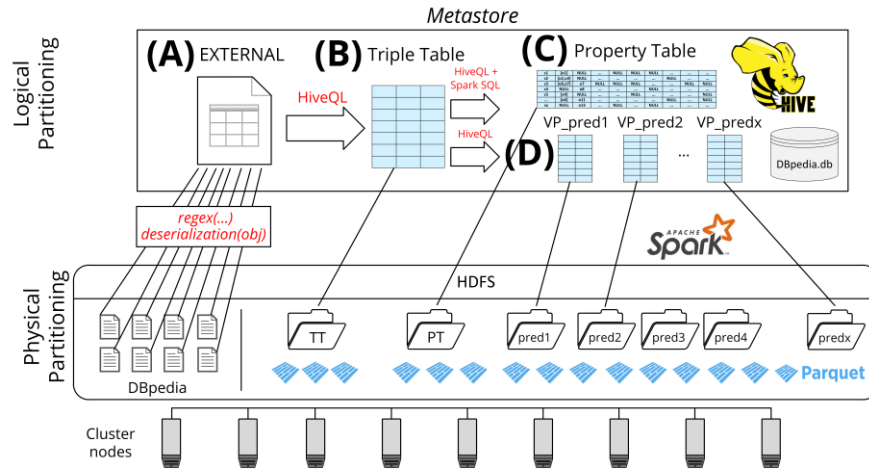


Topic 2: Extending P_{Ro}ST Features

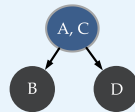


- Based on hadoop-based technologies such as Spark, and Hive
- To evaluate SPARQL queries against large RDF-graphs
- Considers multiple data partitioning strategies
- Query optimizer leverages these different partitioning strategies

1) SPARQL QUERY

```
SELECT ?v1 ?v3 ?v3
WHERE {
  A ?v1 p1 ?v2 .
  B ?v2 p2 "c1" .
  C ?v1 p3 ?v3 .
  D ?v3 p2 "c2"
}
```

2) JOIN TREE



3) SPARK SQL

```
SELECT s AS ?v1, p1 AS ?v2, p3 AS ?v3
FROM property_table

SELECT s AS ?v2
FROM vp_p2
WHERE object="c1"

SELECT s AS ?v3
FROM vp_p2
WHERE object="c2"
```

Publications (Recommended Literature)



- Matteo Cossu, Michael Färber, Georg Lausen: P_{Ro}ST: Distributed Execution of SPARQL Queries Using Mixed Partitioning Strategies. EDBT 2018: 469-472
- Matteo Cossu, Michael Färber, Georg Lausen: P_{Ro}ST: Distributed Execution of SPARQL Queries Using Mixed Partitioning Strategies. CoRR abs/1802.05898 (2018)
- Victor Anthony Arrascue Ayala, Georg Lausen: A Flexible N-Triples Loader for Hadoop. International Semantic Web Conference (P&D/Industry/BlueSky) 2018

Features (potential topics)



1. Implementing a SPARQL-endpoint
2. Extending query evaluation strategies
3. Comparing PProST with other systems

Thank you!



UNI
FREIBURG

Any questions?