

Exceptional service in the national interest



Using Vertex-Centric Programming Platforms to Implement SPARQL Queries on Large Graphs

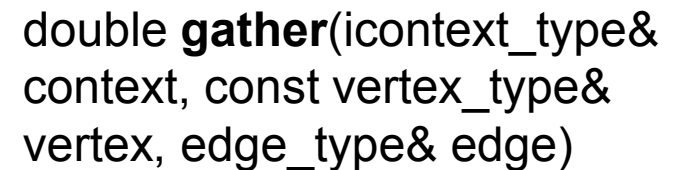
Eric Goodman^{1,2}, Dirk Grunwald²

1 – Sandia National Laboratories, 2 – University of Colorado at Boulder



Sandia National Laboratories is a multi-program laboratory managed and operated by Sandia Corporation, a wholly owned subsidiary of Lockheed Martin Corporation, for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-AC04-94AL85000. SAND NO. 2011-XXXXP

Vertex-centric Programming

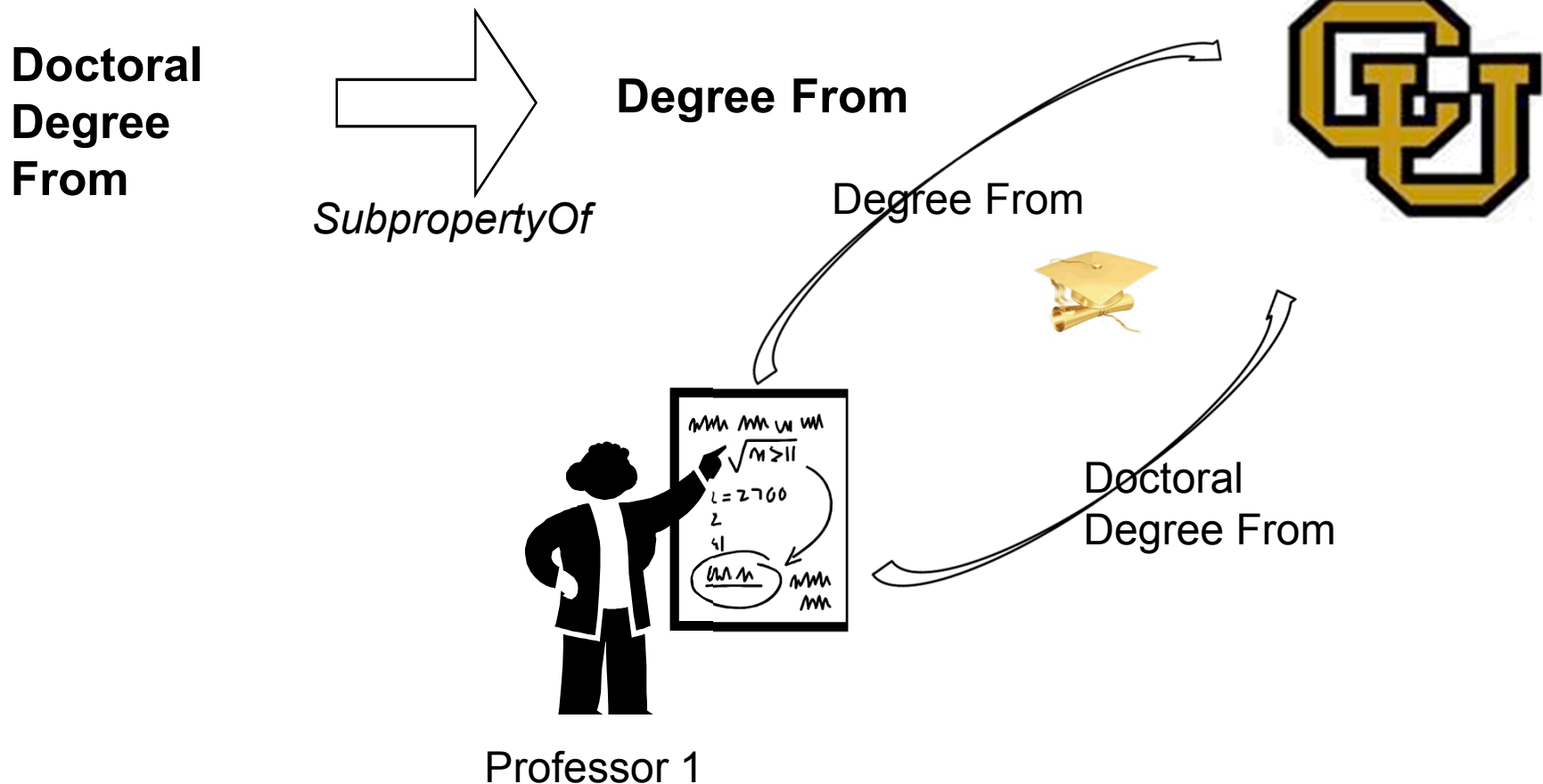


```
edge_dir_type
scatter_edges(icontext_type&
context, const vertex_type&
vertex)
```

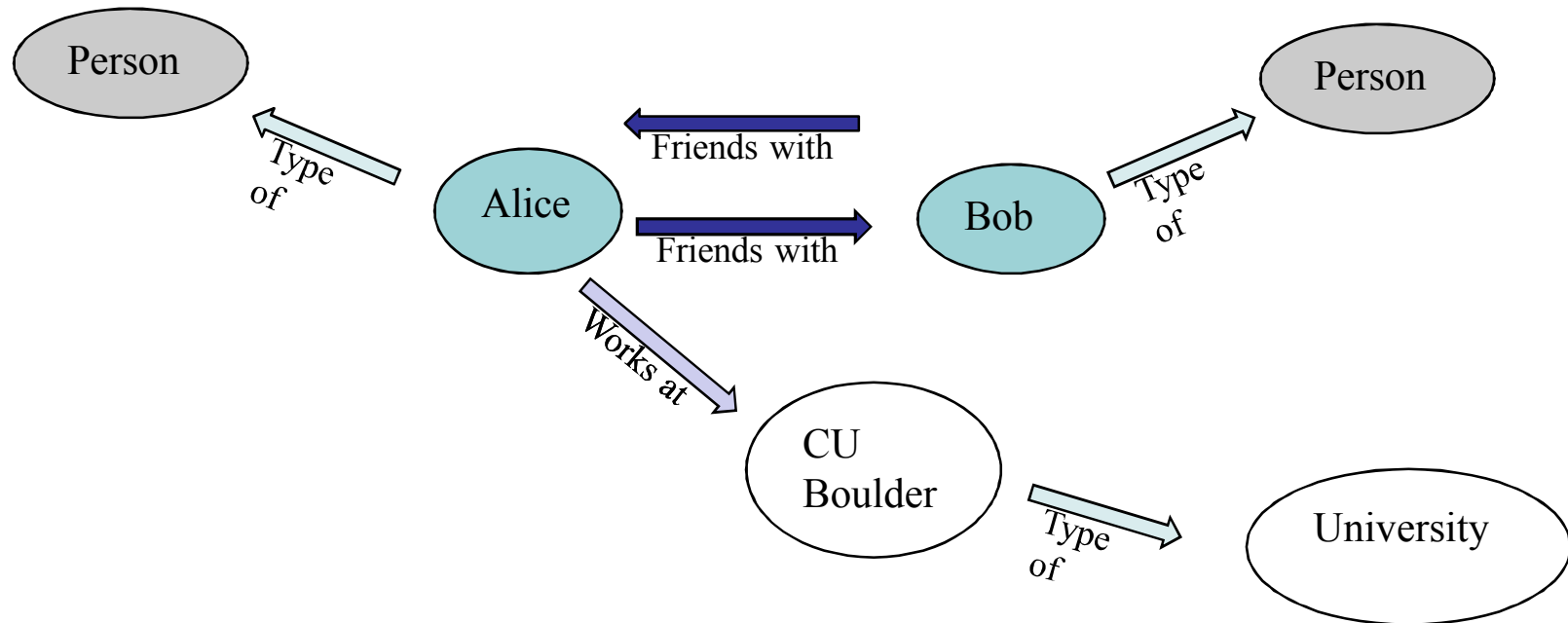
Linking Open Data cloud diagram 2014, by Max Schmachtenberg, Christian Bizer, Anja Jentzsch and Richard Cyganiak. <http://lod-cloud.net/>

What is the Semantic Web?

- Term coined by Berners-Lee, Hendler, and Lassila in 2001
- Data formatted in way such that computers can reason about data much the way we do.



Semantic Web as a Graph



RDF

- Resource Description Format
- Format: **Subject Predicate Object**

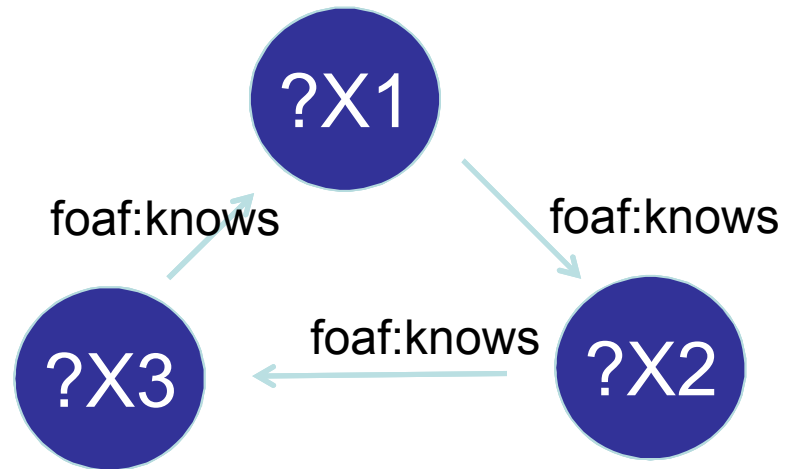
<http://www.Department9.University3272.edu/UndergraduateStudent354>

<http://www.lehigh.edu/~zhp2/2004/0401/univ-bench.owl#emailAddress>

"UndergraduateStudent354@Department9.University3272.edu"

SPARQL

```
SELECT ?X1 ?X2 ?X3
WHERE {
  ?X1 foaf:knows ?X2 .
  ?X2 foaf:knows ?X3 .
  ?X3 foaf:knows ?X1 }
```



Vertex-Centric Computing

- Pregel
- Giraph
- GraphX
- GraphLab
- GraphChi
- Graph Processing System (GPS)

Bulk Synchronous Parallel

- Has three phases:
 1. Concurrent computation: Each node performs computation independent of all other nodes. Memory access occurs only from the nodes' local memory.
 2. Communication: The nodes communicate with each other.
 3. Barrier Synchronization: All nodes are stalled until the communication phase is over.

Vertex-centric Computation

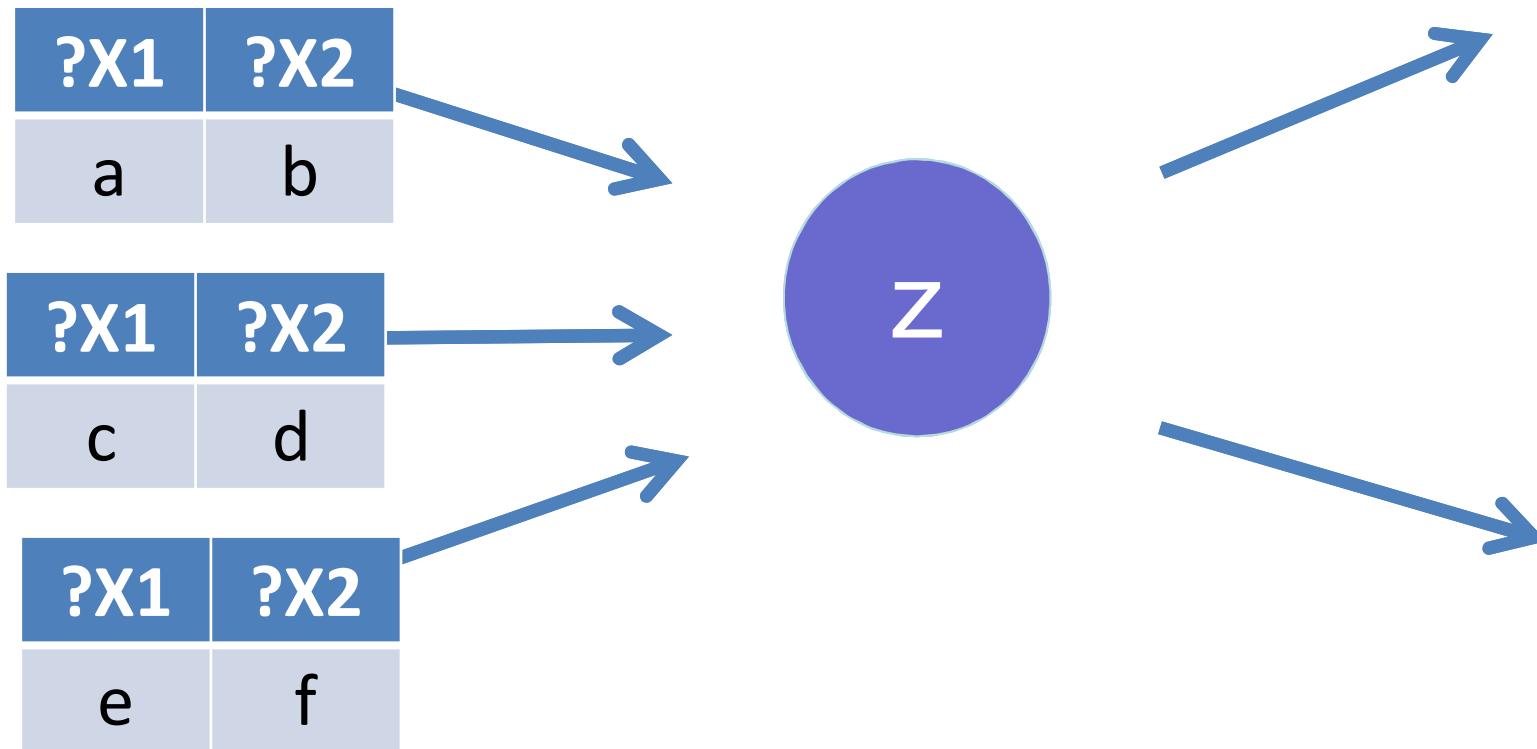
- For example PageRank:
 1. Each vertex takes its current pagerank score, divides by the number of outgoing edges, and then sends that value to all of its neighbors. This is the communication phase of BSP.
 2. Progress halts until all vertices have received messages with updated page rank values from their neighbors. This is the synchronization phase of BSP.
 3. Each vertex re-computes the pagerank score based upon the updated values from its neighbors. This is the concurrent computation phase of BSP.

General Approach

- Take a walk along the subgraph of interest, accumulating results as we go.

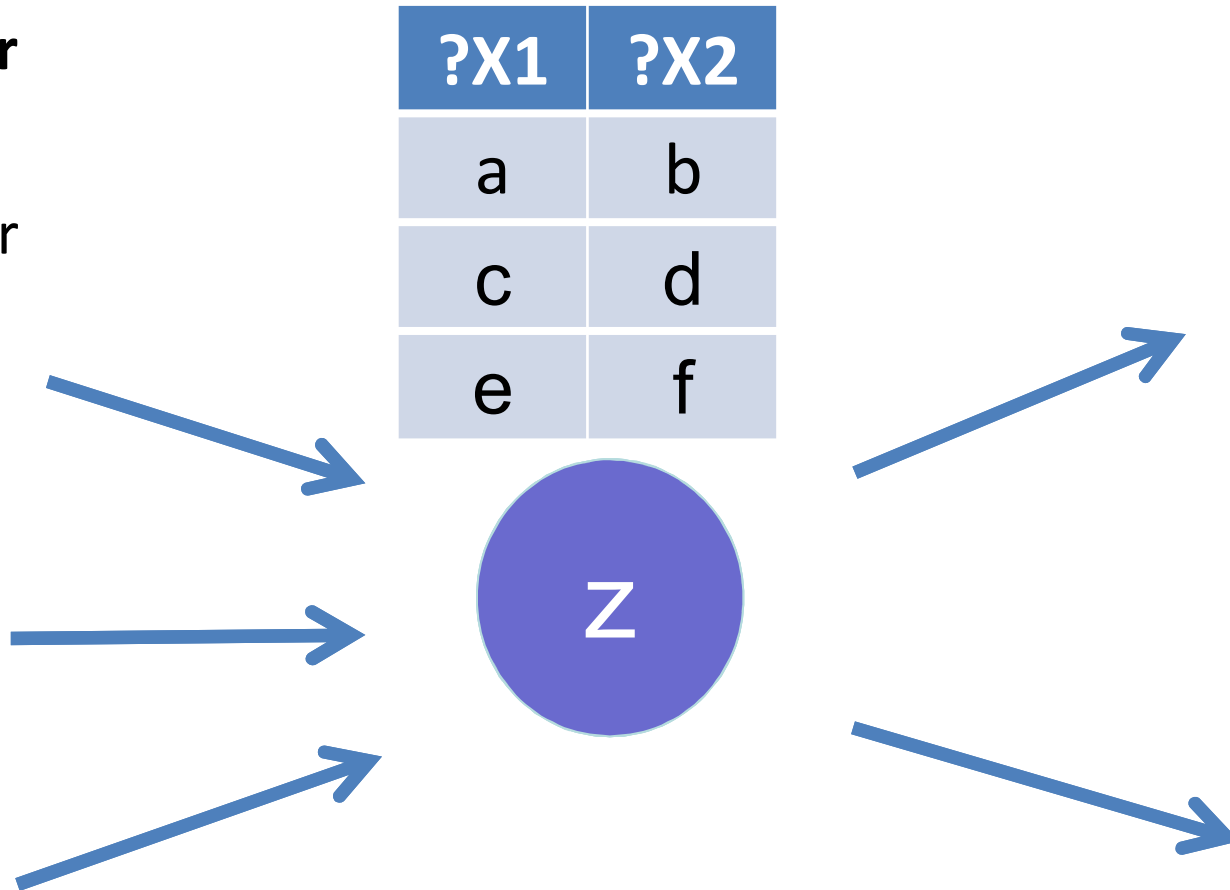
The Pattern

- Gather
- Apply
- Scatter



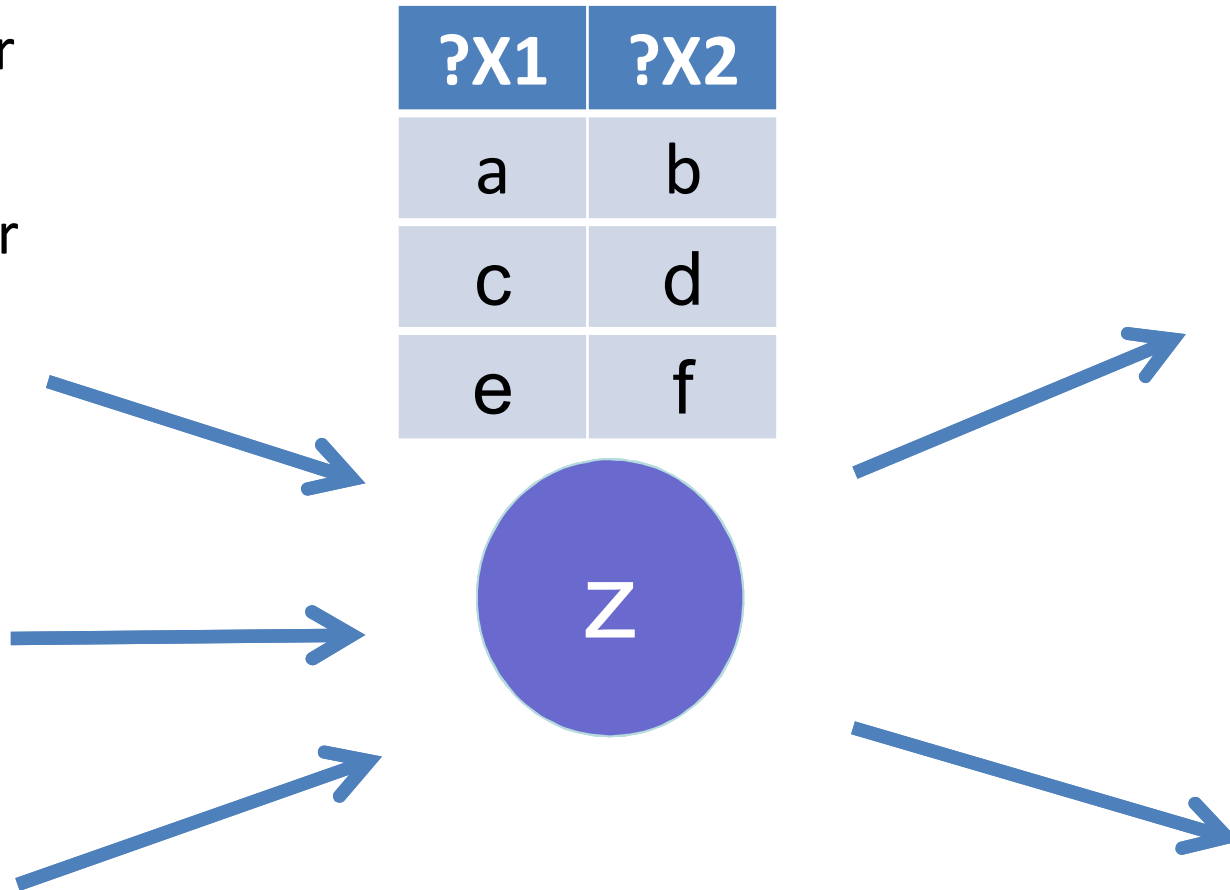
The Pattern

- **Gather**
- Apply
- Scatter



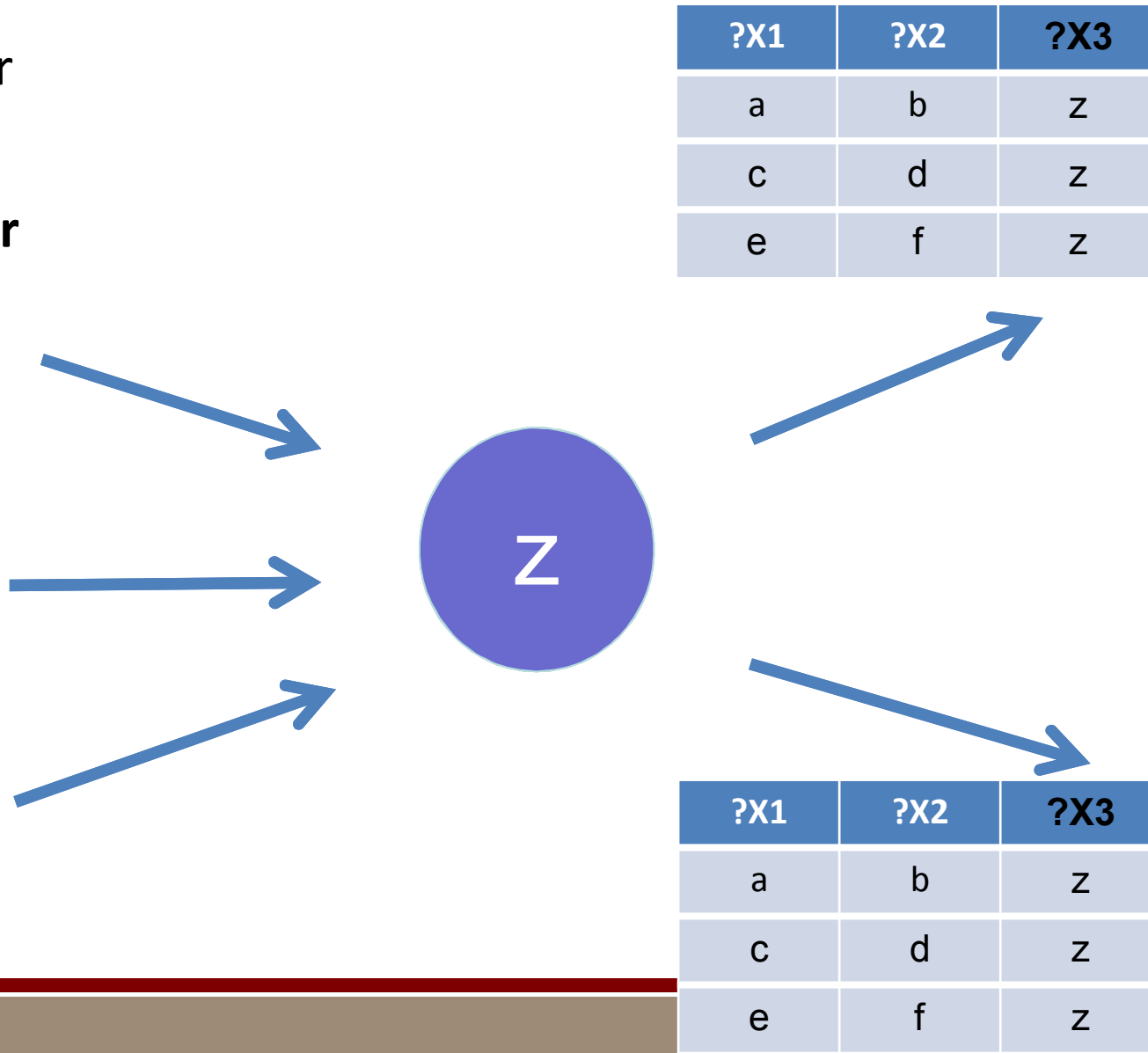
The Pattern

- Gather
- **Apply**
- Scatter



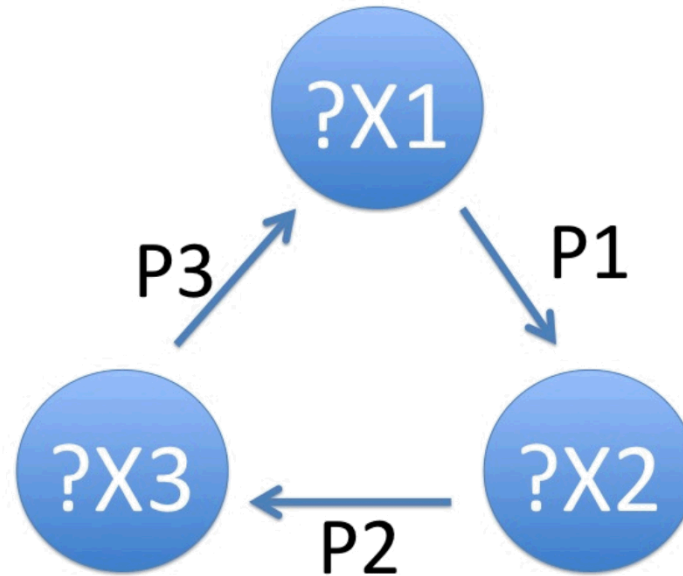
The Pattern

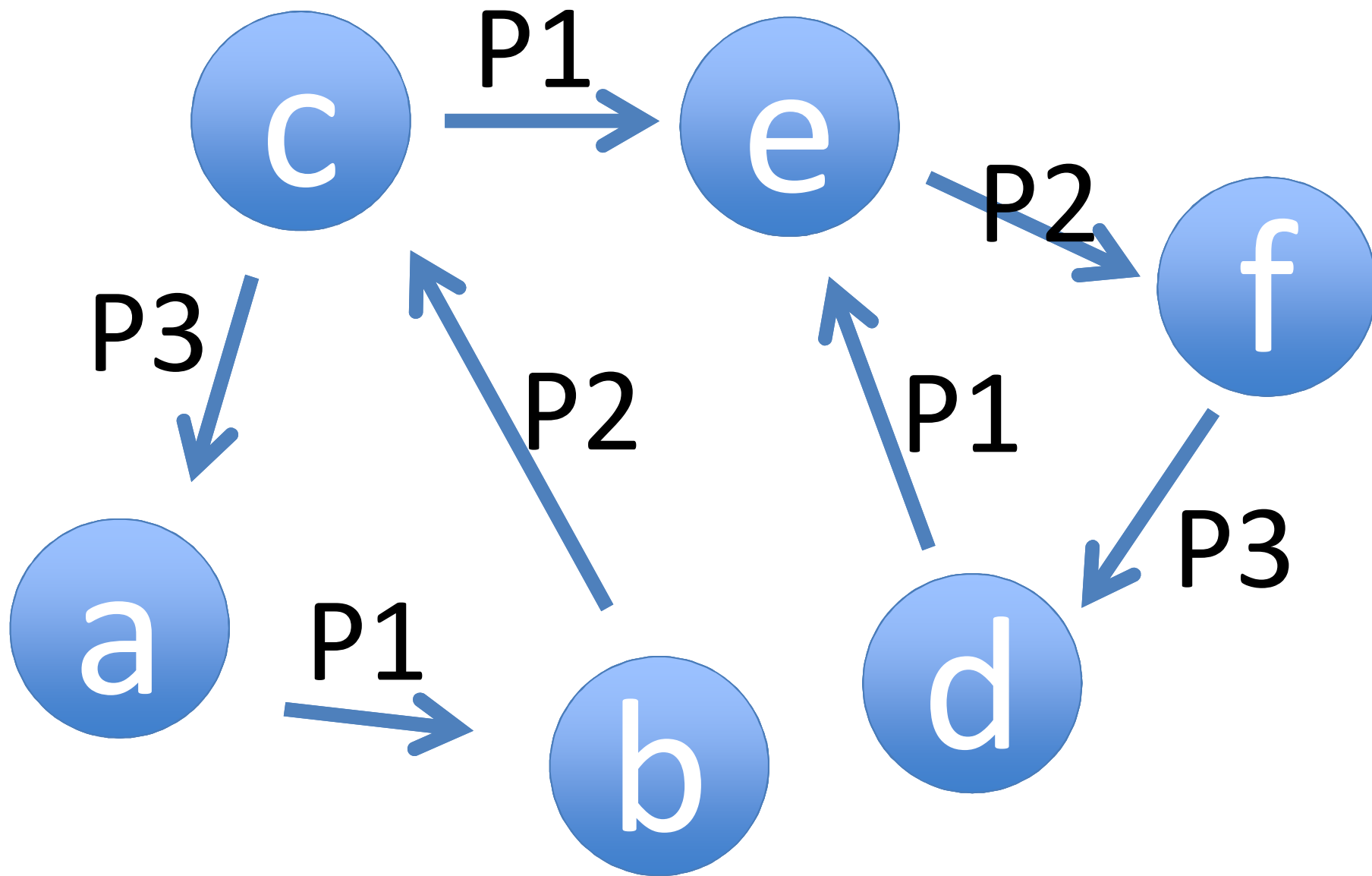
- Gather
- Apply
- **Scatter**

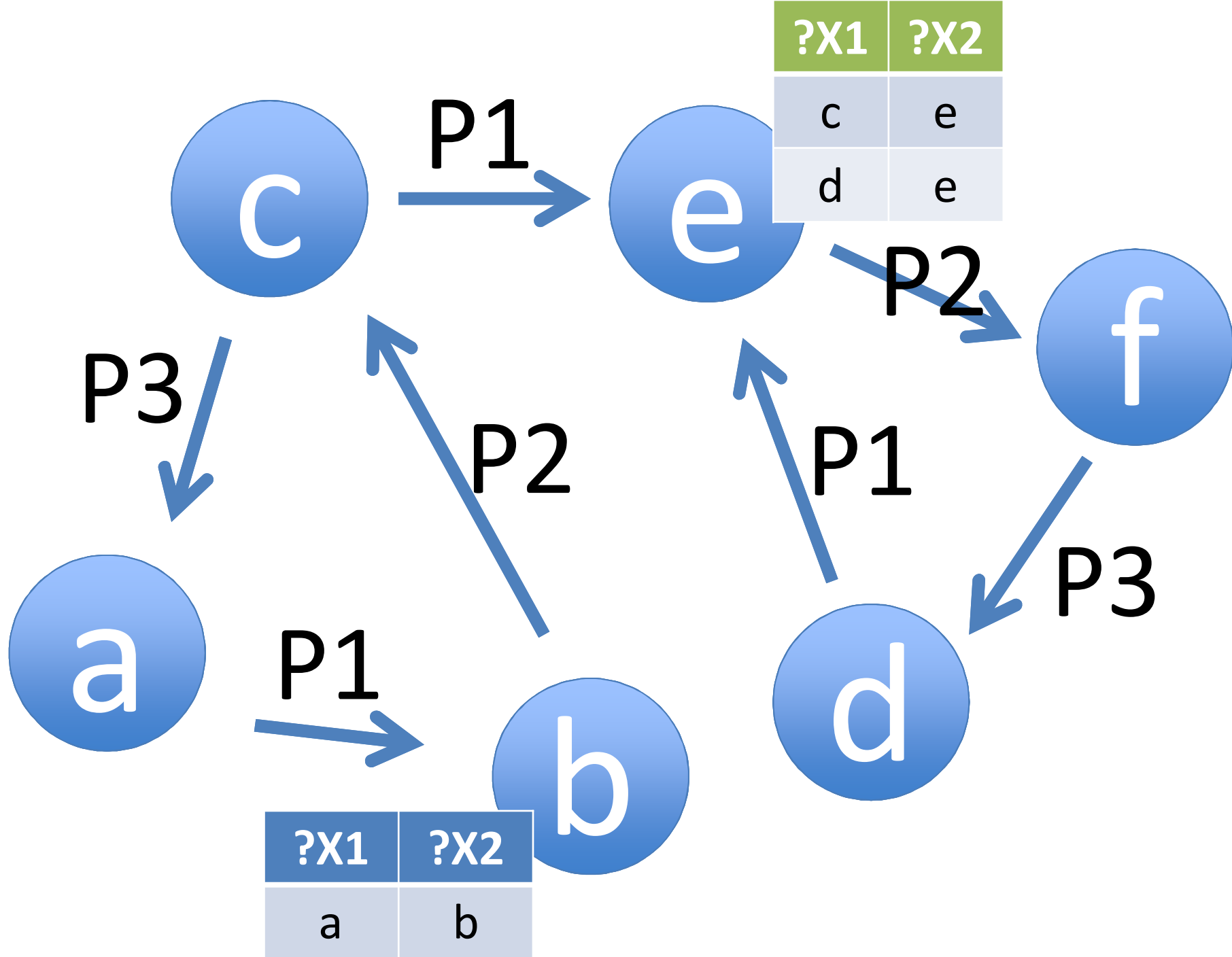


Example Query

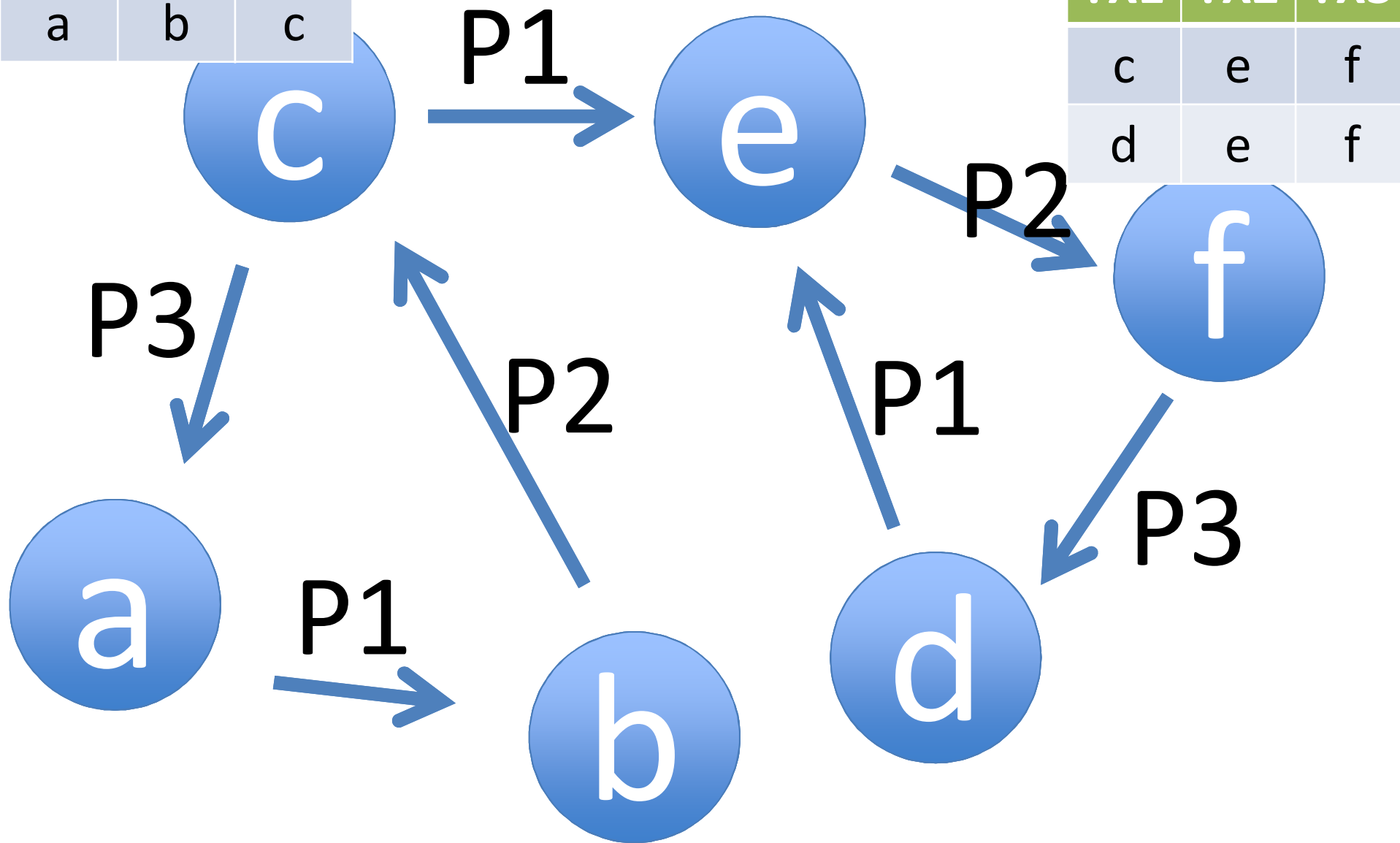
```
SELECT ?X1 ?X2 ?X3  
WHERE {  
  ?X1 P1 ?X2 .  
  ?X2 P2 ?X3 .  
  ?X3 P3 ?X1 }  
}
```

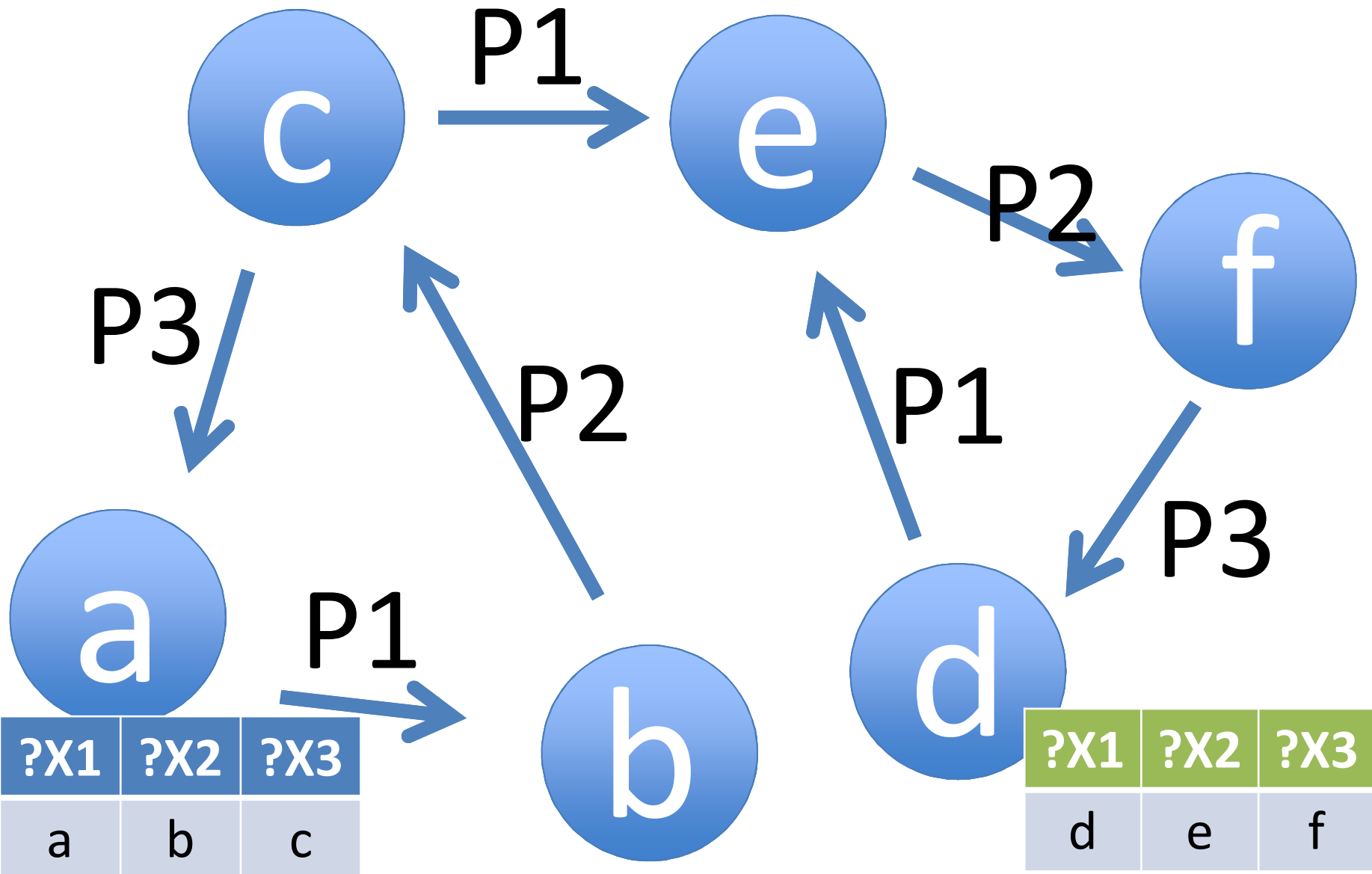






?X1	?X2	?X3
a	b	c





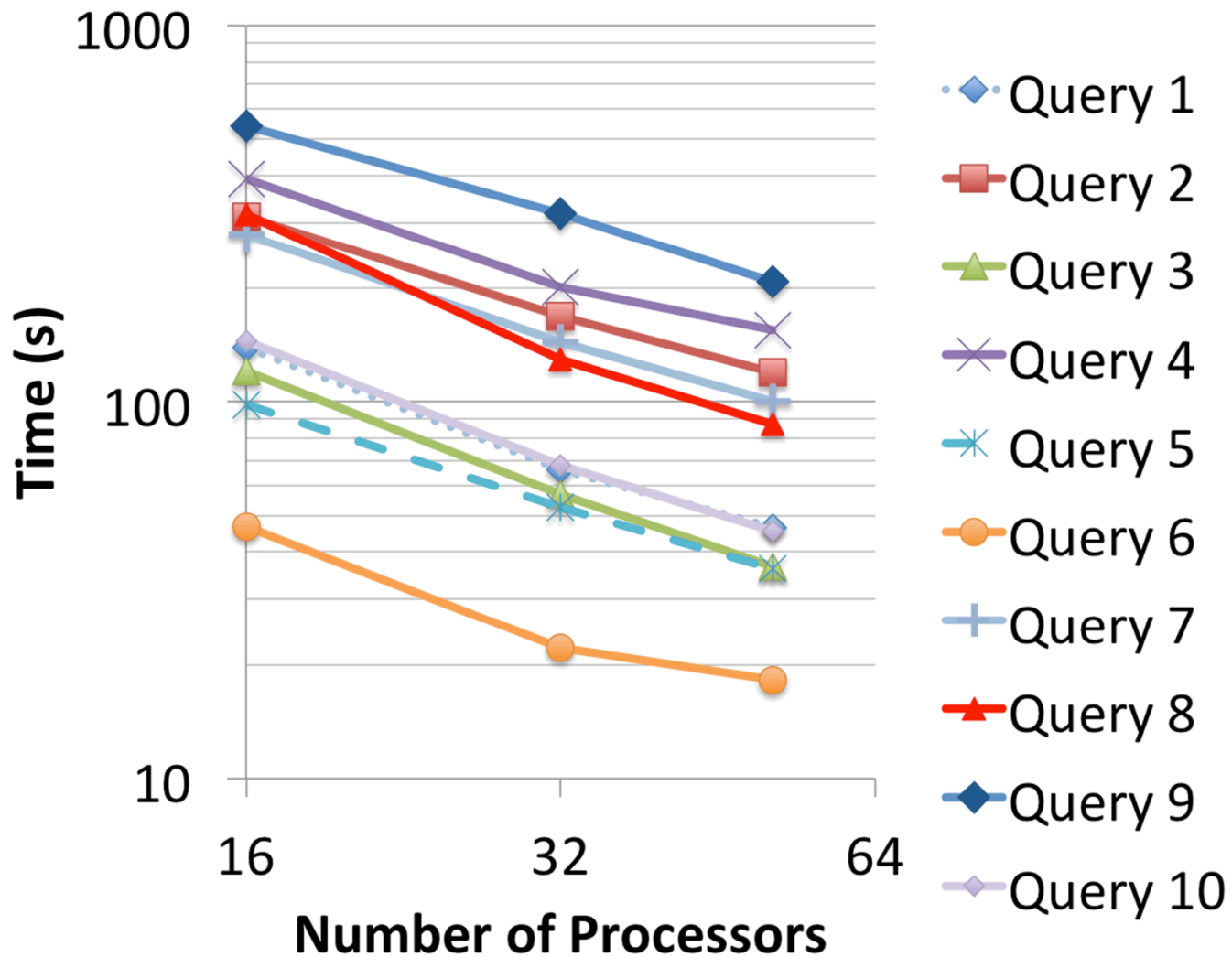
Test Dataset

- Lehigh University Benchmark
- Artificial data set
- Creates
 - Universities
 - Students
 - Teachers
 - Departments
 - Classes
 - Etc
- LUBM(8000) – 8000 universities and ~1.35 billion edges

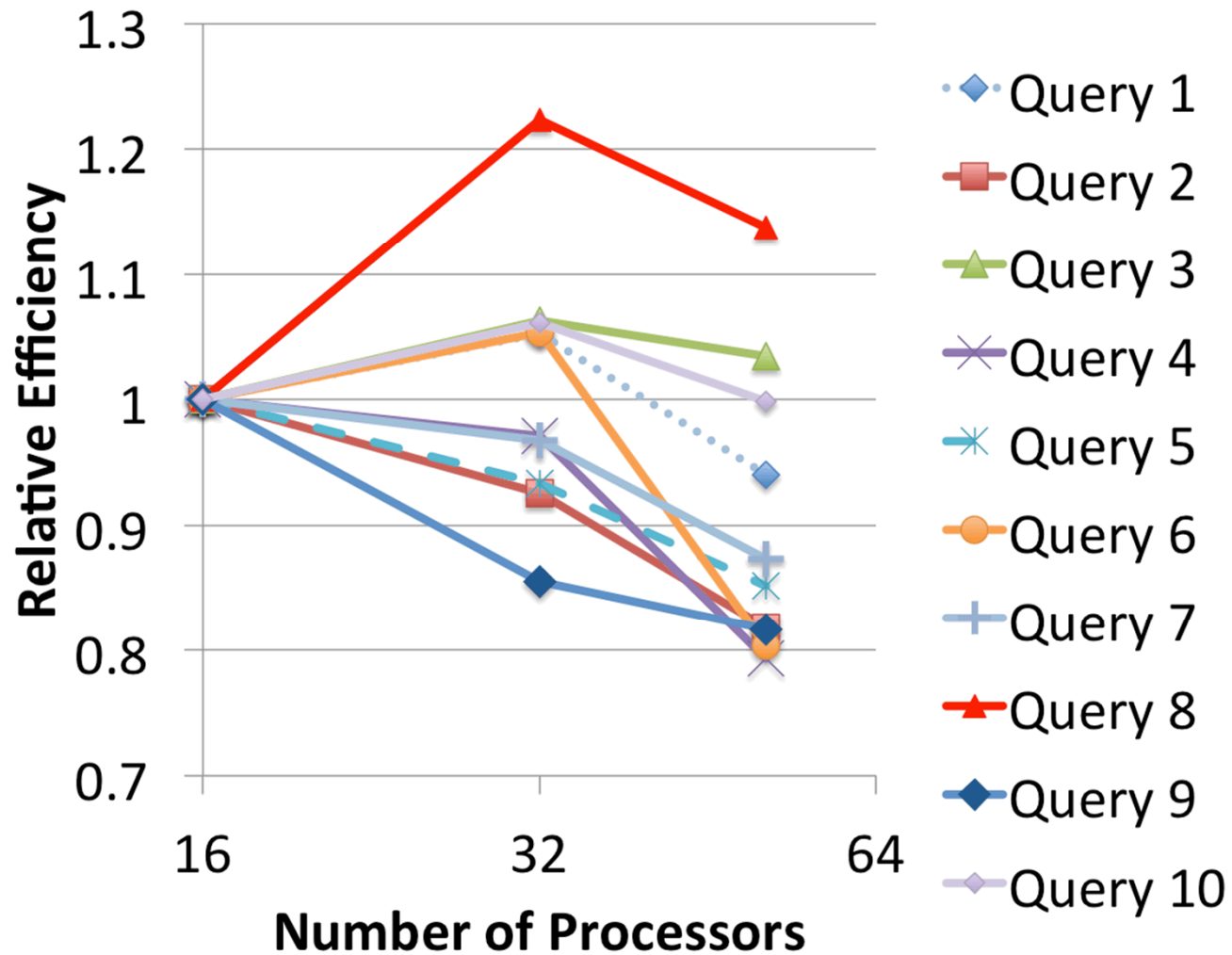
Platform: Plato

- 51 nodes with
 - 2x Intel E5-2470 2.30GHz 8-core
 - 96 GB memory
 - 15x 2TB 6G SAS HDD
 - 10Gb Ethernet

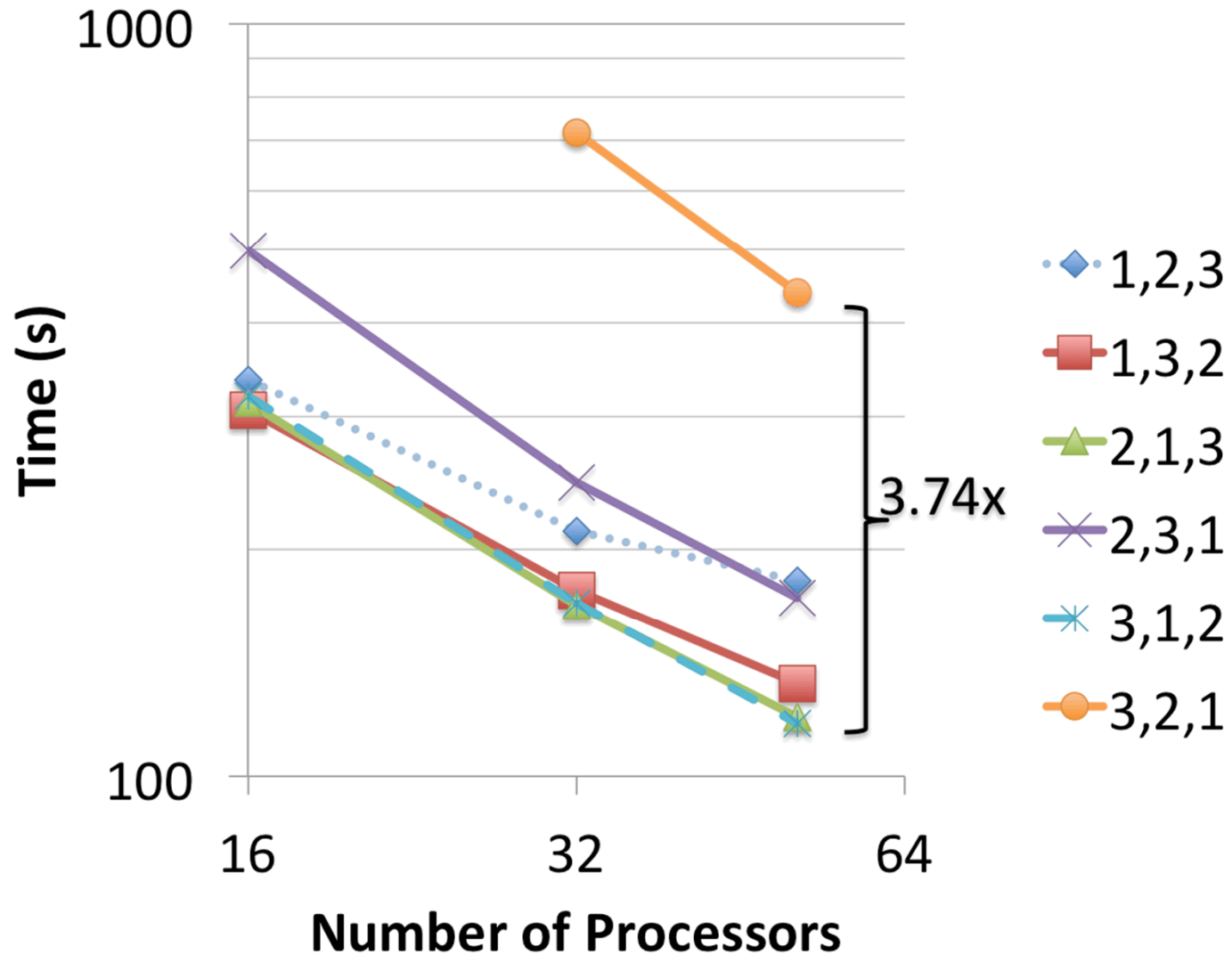
Time of LUBM Queries



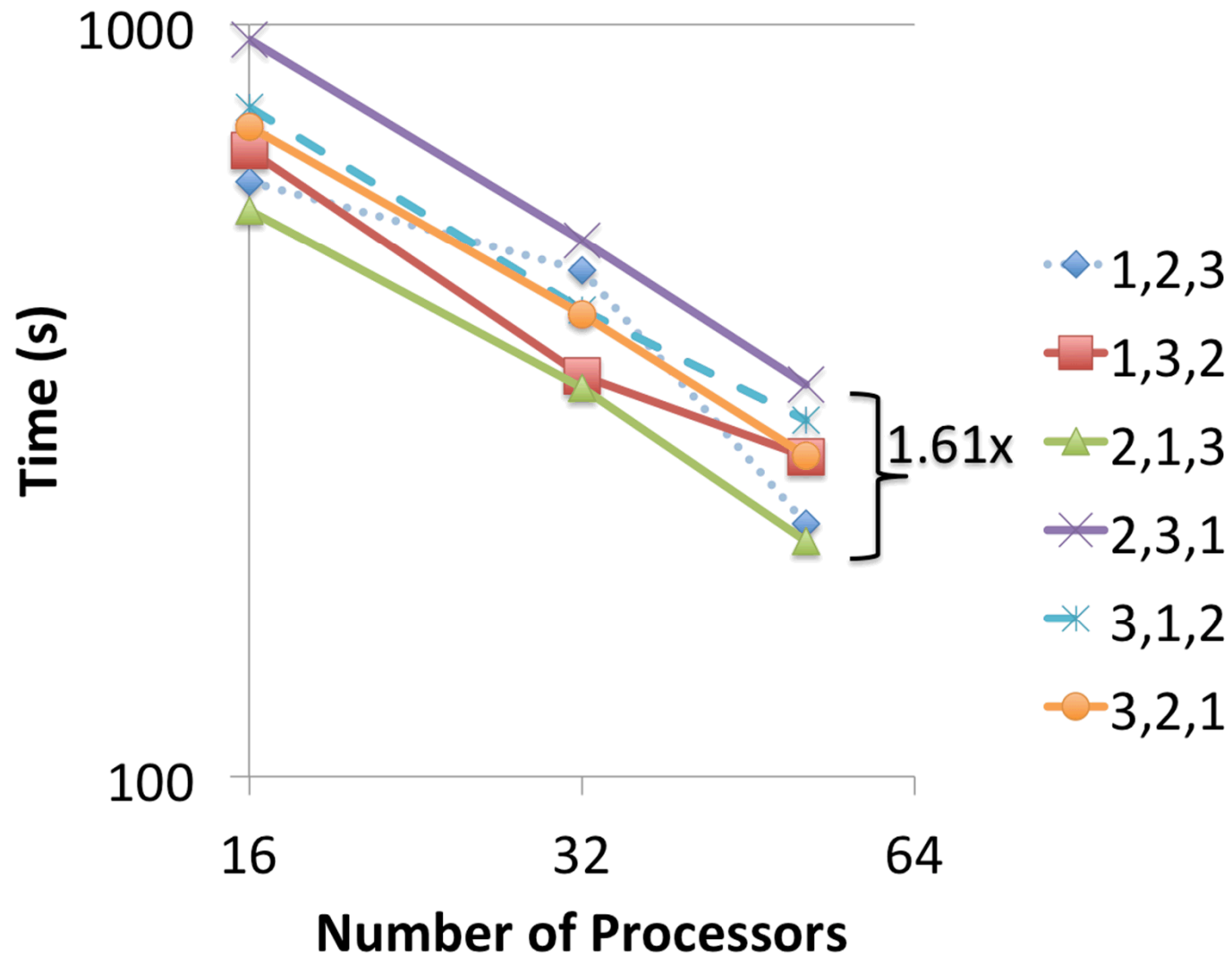
Efficiency of LUBM Queries



Variations with LUBM Query 2



Variations with LUBM Query 9



Limitations

- No execution plan
 - Need to add indices/statistics
- Looks at every edge the first iteration
- No dictionary encoding
- Limited set of SPARQL

Conclusions

- Vertex-centric computing is a viable option for parallelizing SPARQL queries
 - Straightforward implementation
 - Good scalability
- Semantic Web can benefit from additional functionality of vertex-centric paradigm