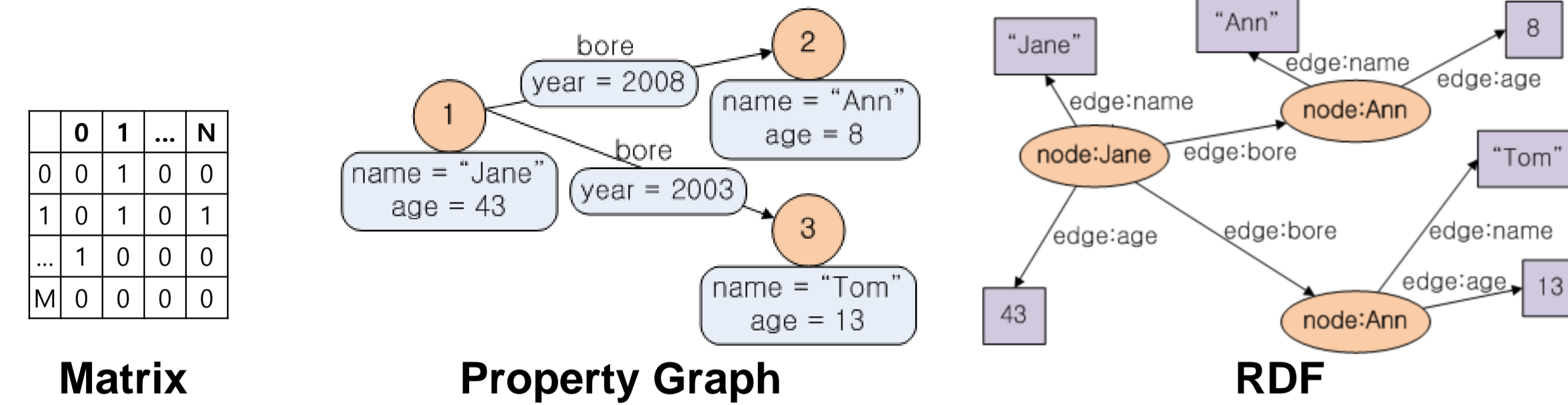


# Benchmarking High Performance Graph Analysis Systems with Graph Mining and Pattern Matching Workloads

Seokyoung Hong<sup>§</sup>, Seung-Hwan Lim, Sangkeun Lee, Sreenivas R. Sukumar,, and Ranga R. Vatsavai<sup>§</sup>  
 North Carolina State University<sup>§</sup>  
 Oak Ridge National Laboratory

## Problem

### ❖ Different Data Model



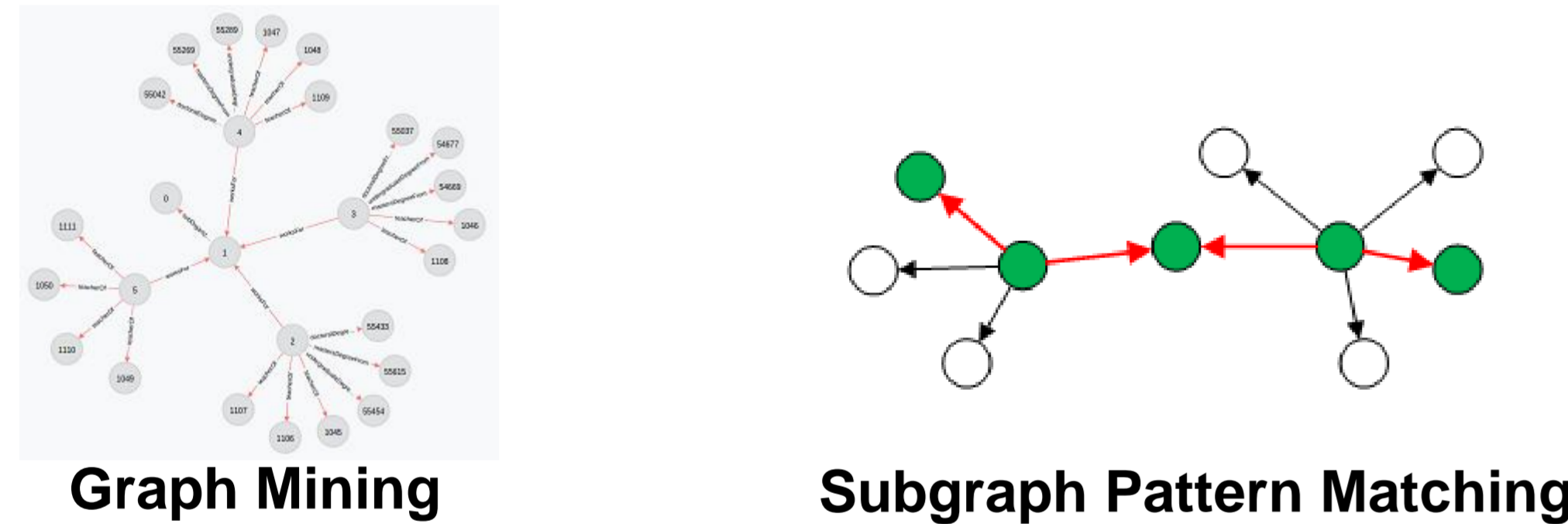
### ❖ Various Graph Processing Systems



### ❖ Different Querying Interfaces



### ❖ Different Graph Analysis Workloads



- Which systems deal with which workloads better or worse?
- What are the advantages/disadvantages of each system?
- What are the best query languages to describe graph analysis tasks?

## Graph Analysis Workloads

### ❖ Datasets

Graph Mining	Graph Pattern Matching
<b>SNAP Datasets</b> <ul style="list-style-type: none"> <li>• DBLP (0.3 million vertices, 1 million edges)</li> <li>• Patents (3.8 million vertices, 16.5 million edges)</li> <li>• Friendster (65.6 million vertices, 1.81 billion edges)</li> </ul>	<b>Extended LUBM Benchmark Datasets</b> <ul style="list-style-type: none"> <li>• 100 Universities (2.1 million vertices, 6.6 million edges)</li> <li>• 1,000 Universities (21.7 million vertices, 66.2 million edges)</li> <li>• 10,000 Universities (0.2 billion vertices, 0.6 billion edges)</li> </ul>

### ❖ Benchmark Queries

Graph Mining	Graph Pattern Matching
<b>Graph Mining Operations</b> <ul style="list-style-type: none"> <li>• Degree Distribution (non-iterative)</li> <li>• Connected Components (iterative)</li> <li>• PageRank (iterative)</li> </ul>	<b>LUBM Benchmark Queries<sup>[*]</sup></b> <ul style="list-style-type: none"> <li>• Query6</li> <li>• Query7</li> <li>• Query9</li> </ul>

## Benchmarked Systems

Graph Mining				Graph Pattern Matching			
System	Processing Platform	Graph Model	Query Interface	System	Processing Platform	Graph Model	Query Interface
Pegasus	MapReduce	Matrix	Java/GIM-V	Titan	Titan/HBase	Property Graph	Gremlin
GraphX	Spark	Property Graph	Scala/Pregel	GraphX	Spark	Property Graph	Scala/GraphX Library
uRIKA	Cray XMT	RDF	SPARQL	uRIKA	Cray XMT	RDF	SPARQL

## Evaluation Environment

Graph Mining			Graph Pattern Matching		
System	Version	Environment	System	Version	Environment
Pegasus	Hadoop 2.4.0	65 AWS m3.2xlarge (64 compute nodes)	Titan	Titan 0.5.1 HBase 1.0	9 32-core, 64G-RAM, 500GB-disk nodes (8 compute nodes)
GraphX	Spark 1.0.0	65 AWS m3.2xlarge (64 compute nodes)	GraphX	Spark 1.3.1	9 32-core, 64G-RAM, 500GB-disk nodes (8 compute nodes)
uRIKA	-	ORNL 2TB shared memory	uRIKA	-	ORNL 2TB shared memory

## Comparisons and Conclusion

### ❖ Performance

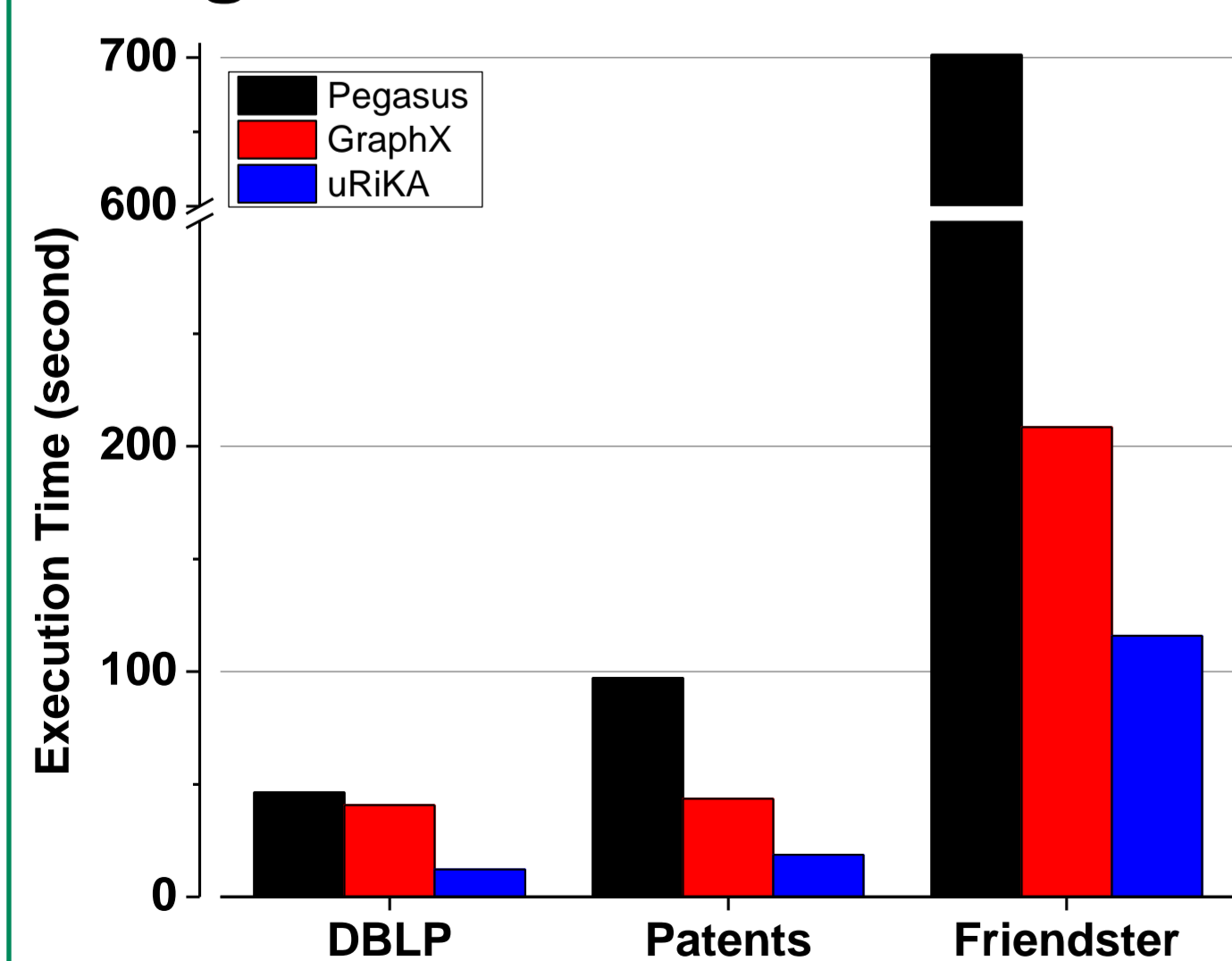
Graph Mining	Graph Pattern Matching
<ul style="list-style-type: none"> <li>• GraphX showed the best performance for iterative processing.</li> <li>• uRIKA showed the best performance for non-iterative processing (Degree Distribution) and good performance with small and moderate size graphs.</li> <li>• Pegasus showed the worst performance in most cases.</li> </ul>	<ul style="list-style-type: none"> <li>• uRIKA showed the best performance in most cases.</li> <li>• GraphX showed fairly good performance while it could not handle 10K universities on our nine-node cluster.</li> <li>• Titan/HBase showed poor performance compared to the others.</li> </ul>

### ❖ Querying Interfaces and APIs

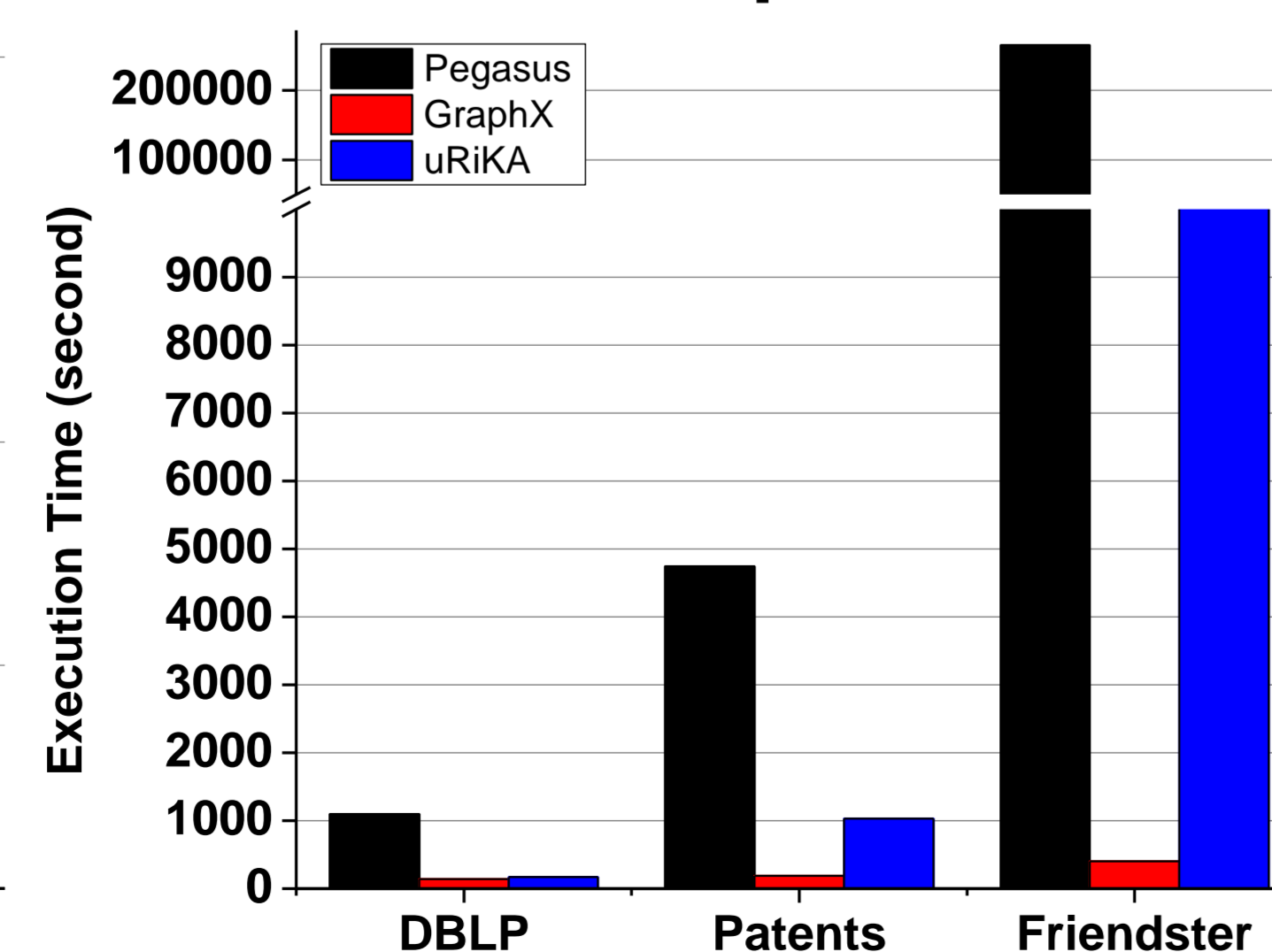
System	Description
Pegasus	<ul style="list-style-type: none"> <li>• It provides a set of matrix multiplication operations which enable graph mining tasks.</li> <li>• With those operations, it is not straightforward to describe pattern matching tasks.</li> </ul>
GraphX	<ul style="list-style-type: none"> <li>• It provides the BSP processing model with related primitives along with optimized graph mining operations which allow graph mining operations to be described effectively.</li> <li>• A set of operations (e.g., Triplet and Join) provides a good way to describe pattern matching tasks.</li> </ul>
Titan	<ul style="list-style-type: none"> <li>• Its default query language, Gremlin, provides an intuitive way to write graph traversal queries but it has limitations on describing pattern matching queries, which did not allow us to write LUBM query 9.</li> </ul>
uRIKA	<ul style="list-style-type: none"> <li>• Its query language, SPARQL, does not provide a way to describe iterative mining tasks in a single query but several wrapper languages (Python and Java) can be used for that.</li> <li>• SPARQL provides a fairly easy and intuitive way to describe pattern matching tasks.</li> </ul>

## Graph Mining Result

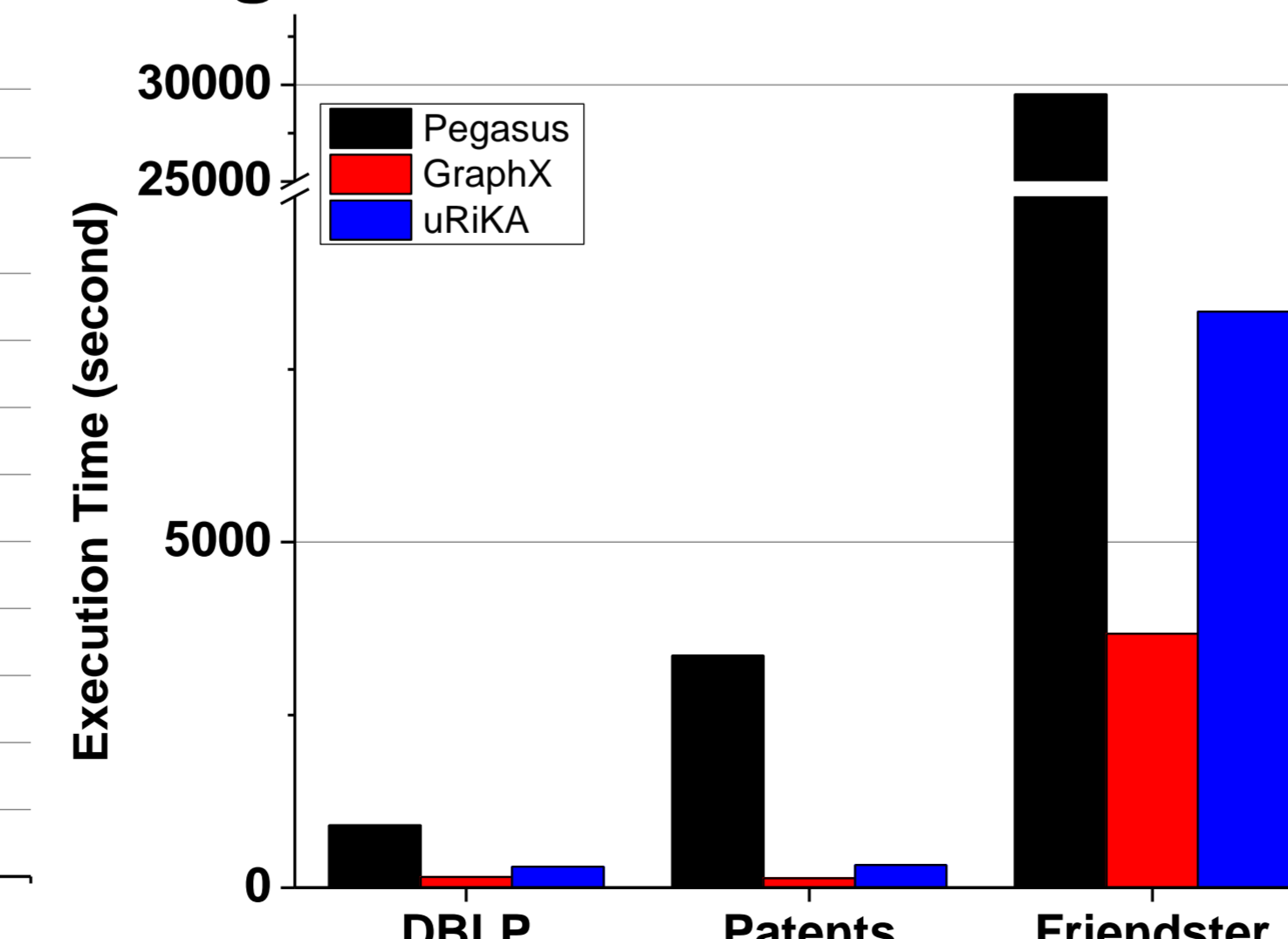
### ❖ Degree Distribution



### ❖ Connected Components

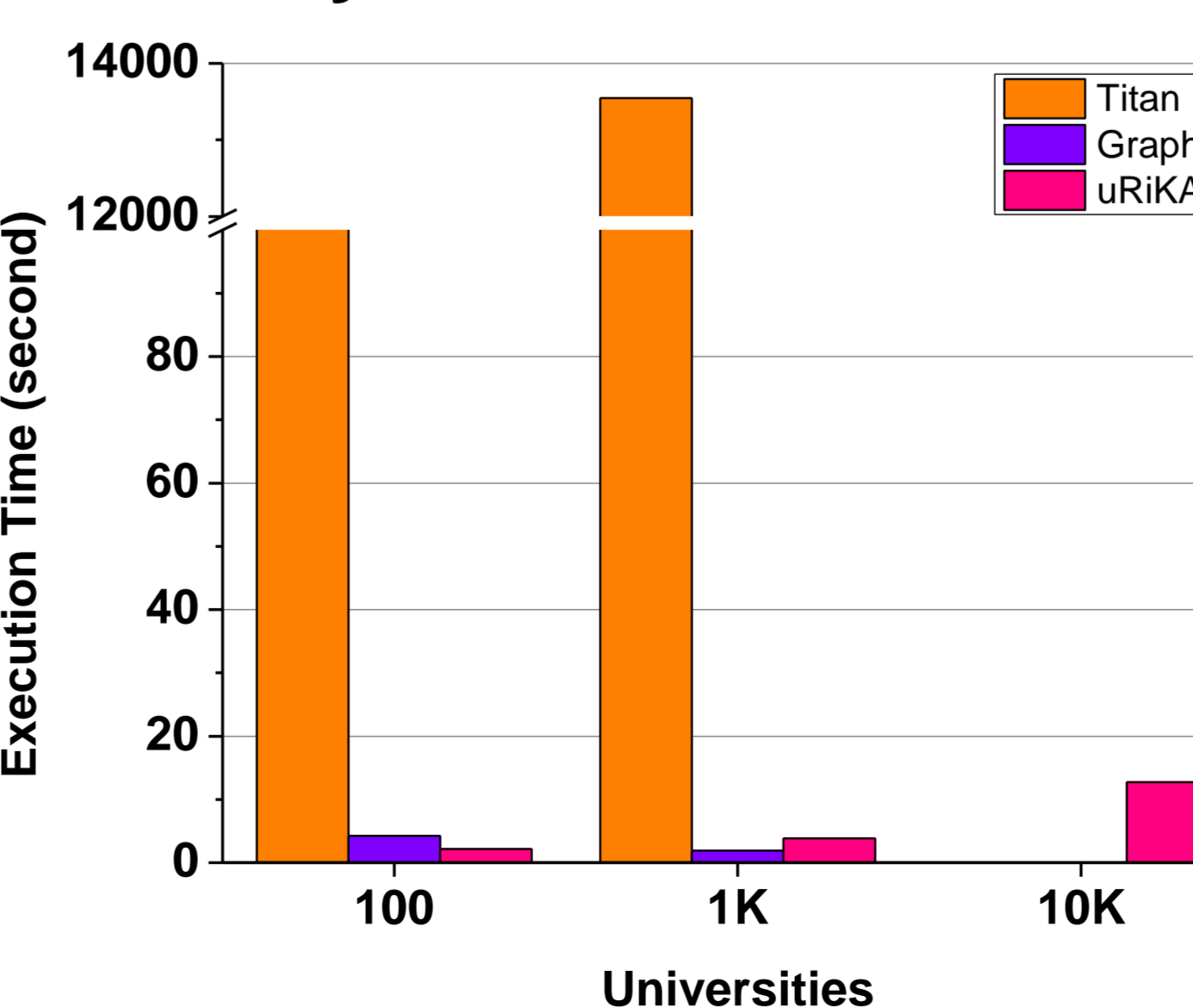


### ❖ PageRank

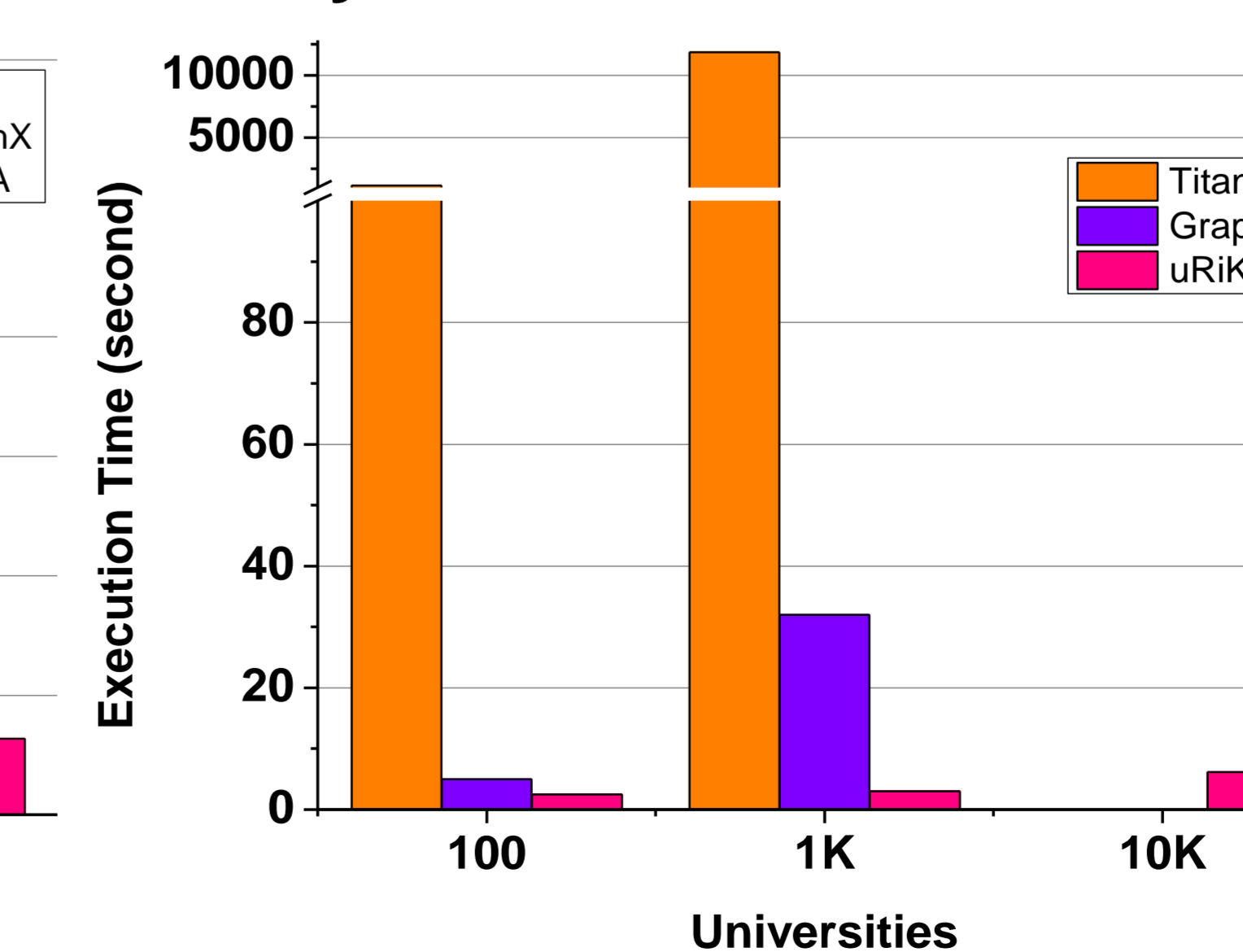


## Pattern Matching Result<sup>[\*]</sup>

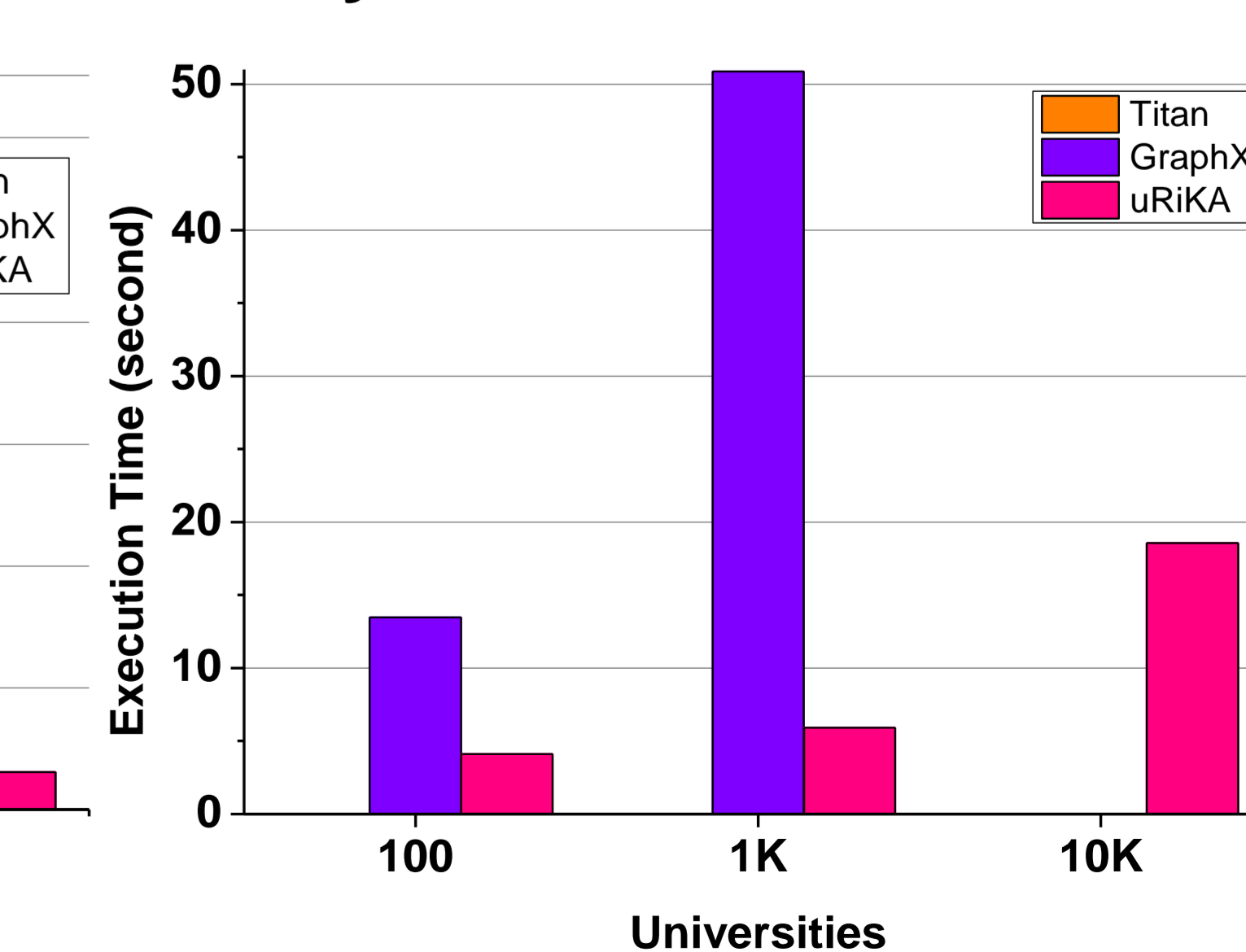
### ❖ Query6



### ❖ Query7



### ❖ Query9



[\*] We re-wrote 9 out of 14 original LUBM queries and tested all the nine queries. Due to the page limitation, we presented the execution times of three queries each of which has different selectivity and a selectivity pattern as graph size increases.