

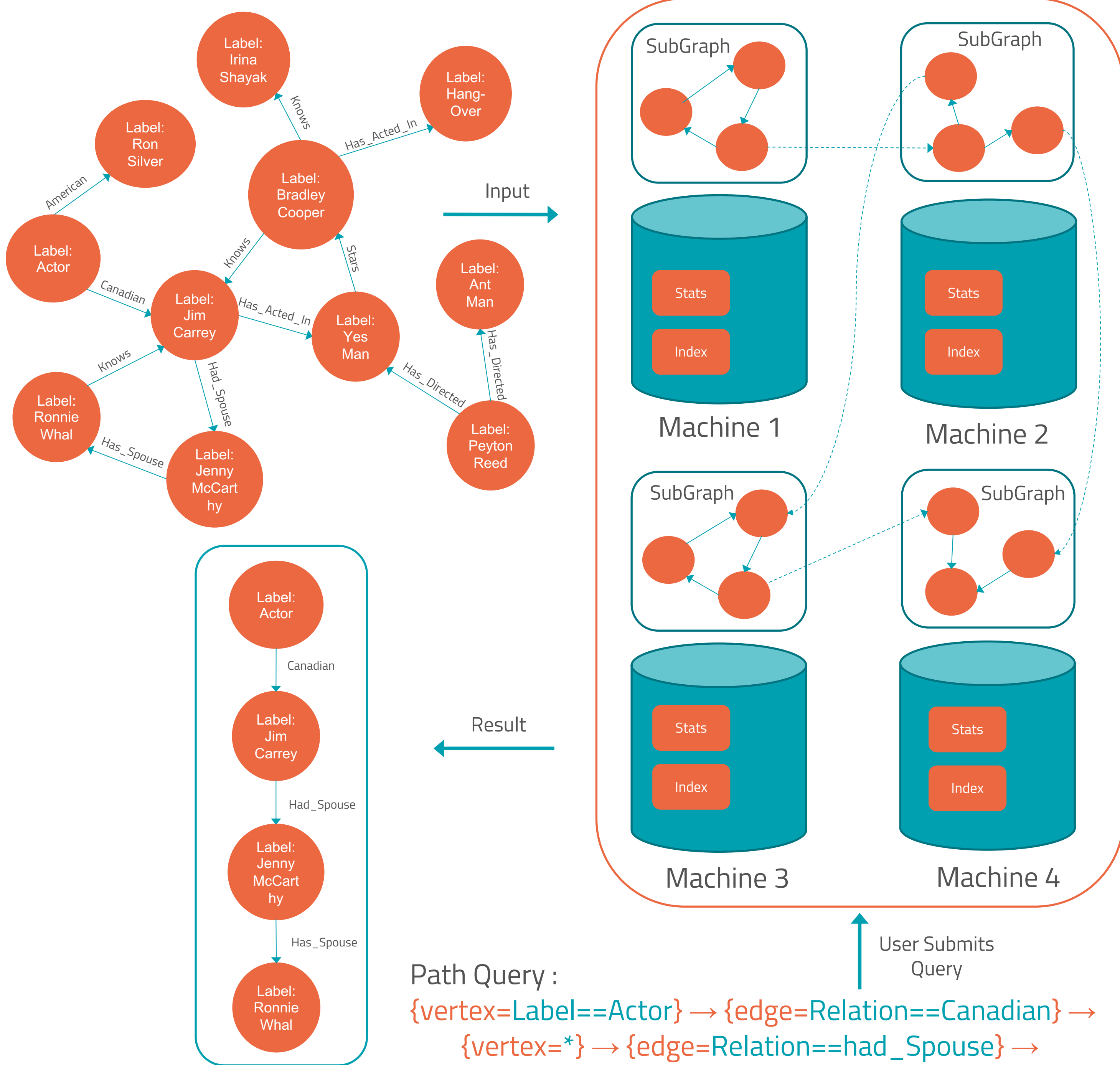


Distributed Querying over Compressed Property Graphs

Swapnil Gandhi, Sayandip Sarkar, Abhilash Sharma and Yogesh Simmhan

INTRODUCTION

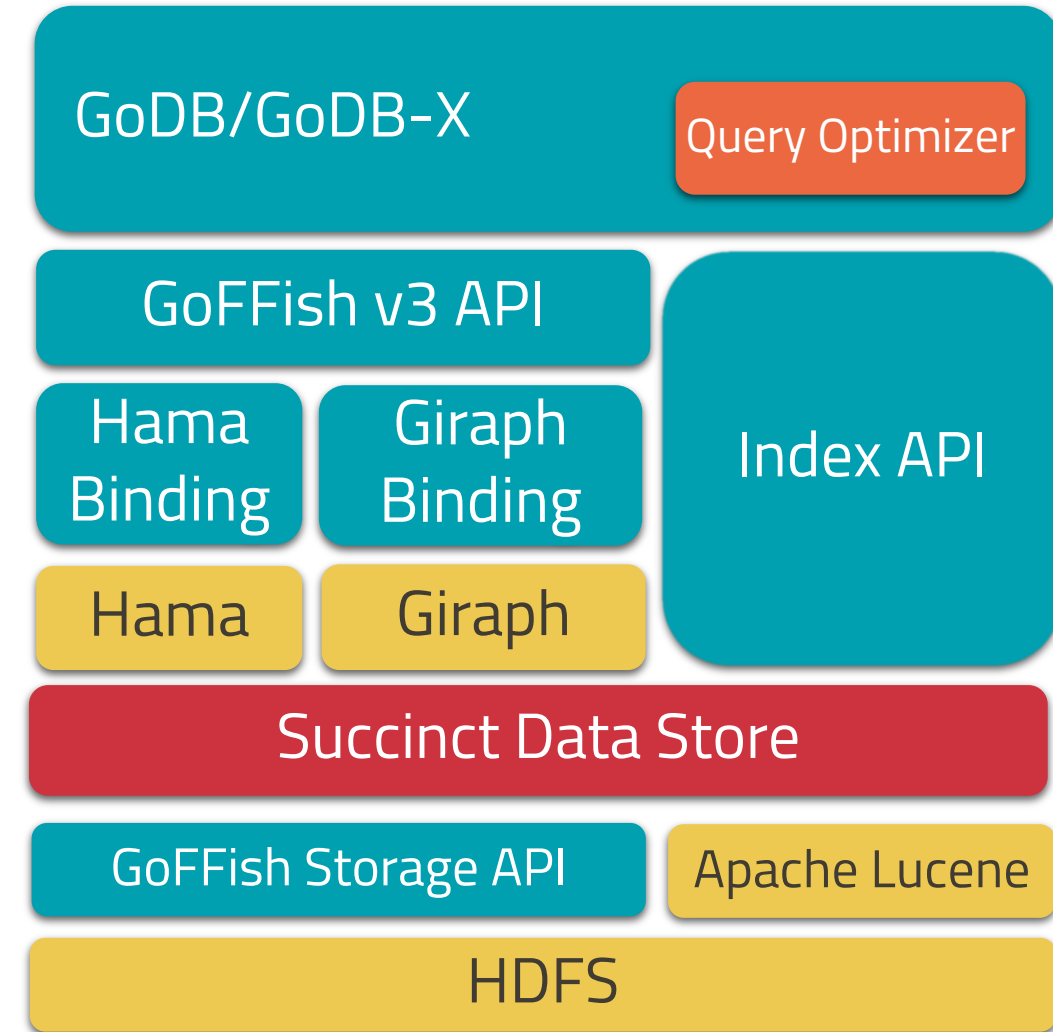
- Increasing trend toward representing semi-structured data as **Property Graphs** using compact in-memory store.



- Challenges:
 - Low latency interactive querying on compressed data
 - Scaling to Large Graphs (Billions of Vertices and Edges)
 - High Throughput
 - Graphs Keep Growing in Size

GoDB

- Builds upon **GoFFish** subgraph-centric batch processing framework
- Bulk Sync Parallel** Execution Model
- Common graph query types: VE, Path, BFS, Reachability
- Heuristics based **Cost Model** to select best distributed query execution plan



GoDB-X

- Harnesses query execution model of GoDB
- Uses **Succinct** as an underlying Data Store
- Reduces memory footprint of Java objects but incurs **higher traversal costs**
- Offers an interesting **trade-off** between distributed memory utilized and query latency time

DATA-MODEL

VertexID	Attribute List	AttributeID	(Order, Delimiter)
4061105	countryID: US, numberOfClaims: 40	countryID	(0,)
132768	countryID: FR, numberOfClaims: 7	numberOfClaims	(1, &
696696	countryID: SE, numberOfClaims: 92		

Figure 1: VertexFile layout storing vertexID and property name-value pairs

```

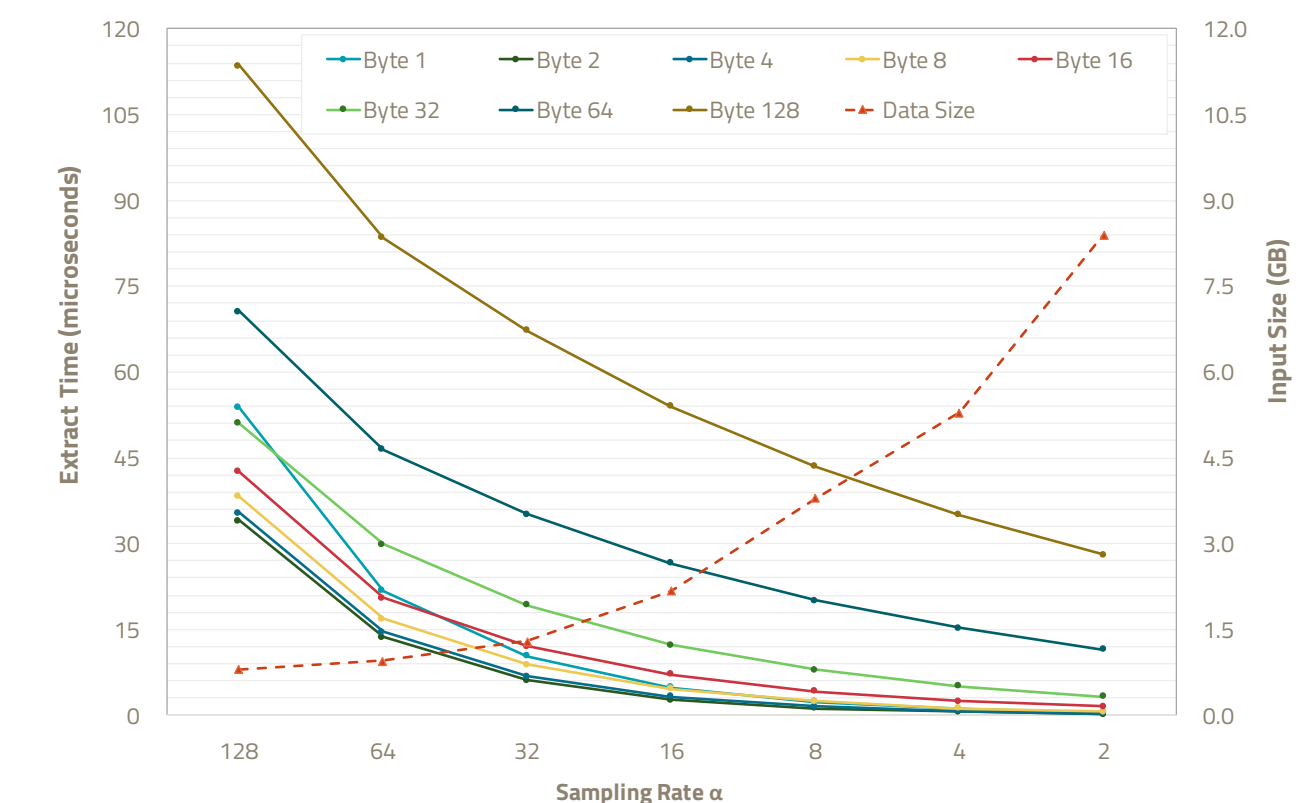
$S1•localEdgeCount•ID1:ID2:ID3...IDMθrD1:rD2:...rDN#
$S2•localEdgeCount•ID1:ID2:ID3...IDMθrD1:rD2:...rDN#
$S3•localEdgeCount•ID1:ID2:ID3...IDMθrD1:rD2:...rDN#
    
```

Figure 2: EdgeFile layout storing vertex adjacency list

DATASET

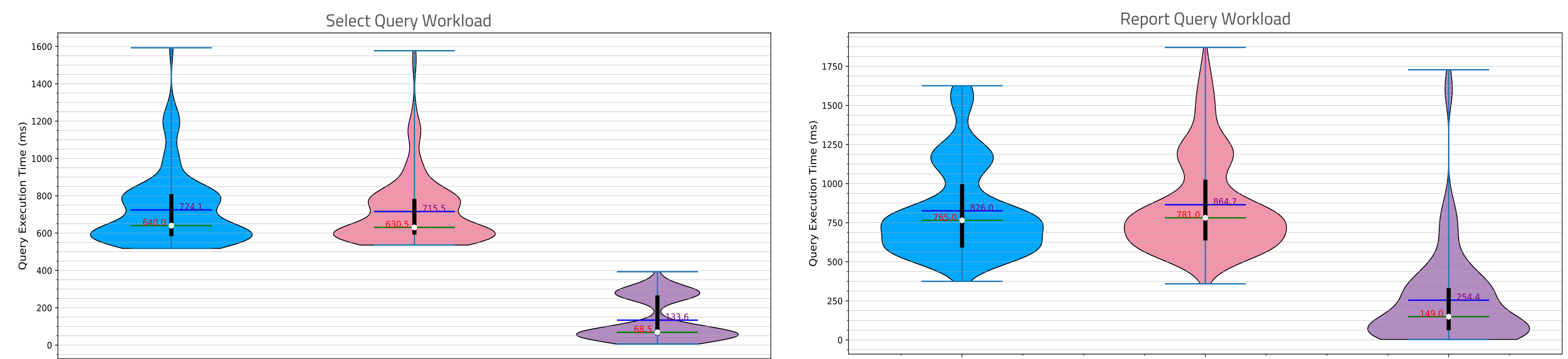
Dataset	Vertices	Edges	Properties	On-disk Size	Number of Select Queries	Number of Report Queries
CITP	3.7 M	16.5 M	2	1 GB	300	300
GITR	170 M	1.2 B	4	38 GB	120	480

EXPERIMENTAL EVALUATION



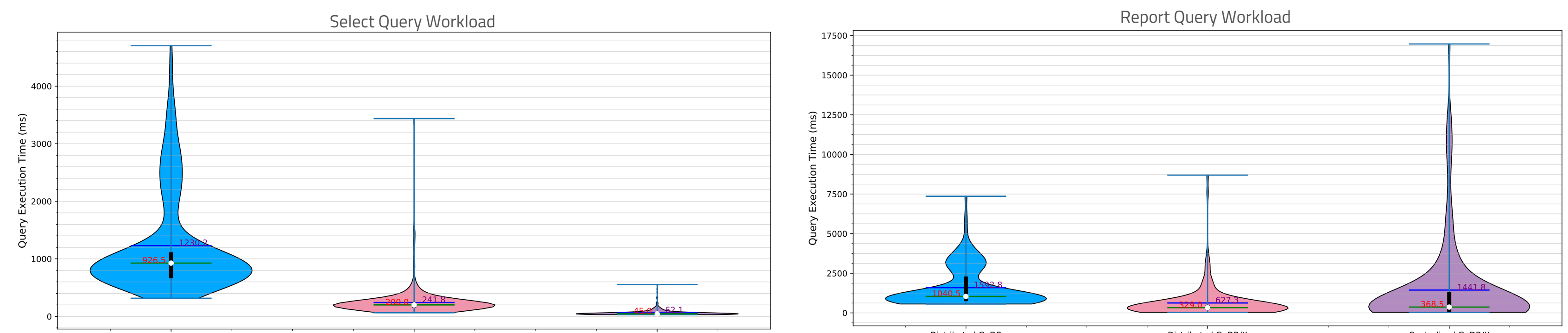
Space-Time trade-off of Succinct Data Store for different Sampling Rate α . We adopt $\alpha=32$ in our experiments

CITP DATASET



Though by using compressed data-model we were able to see **>75% Reduction** in Query Execution Time for Centralized GoDB-X, for Distributed (4 VMs) GoDB-X we see execution time to be marginally slower than GoDB. Static Overheads of querying compressed data + network communication cost **out-strips** any advantages offered by parallel execution.

GITR DATASET



>50% Reduction in Query Execution Time for Distributed GoDB-X. As more operations per Machine are required for GITR, we see advantages of parallel execution out-weighing communication costs and static overheads.

ON-GOING WORK

- Improvements in Data Model

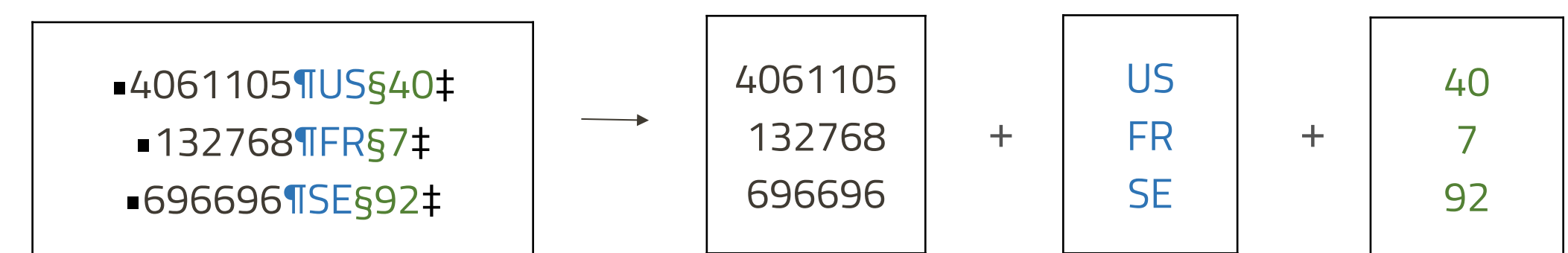


Figure 3: VertexFile further divided into separate columnar files for efficient Extract

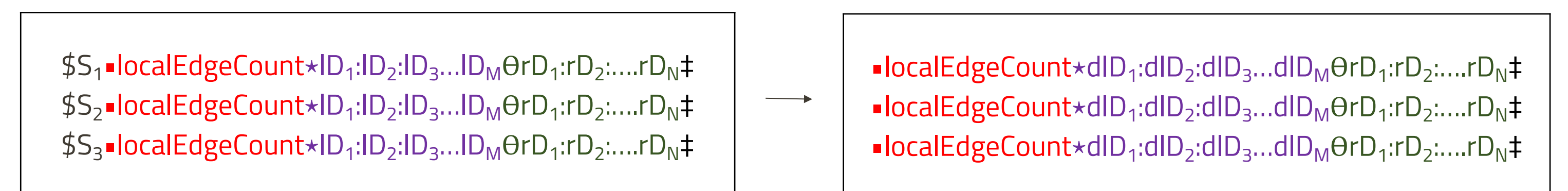
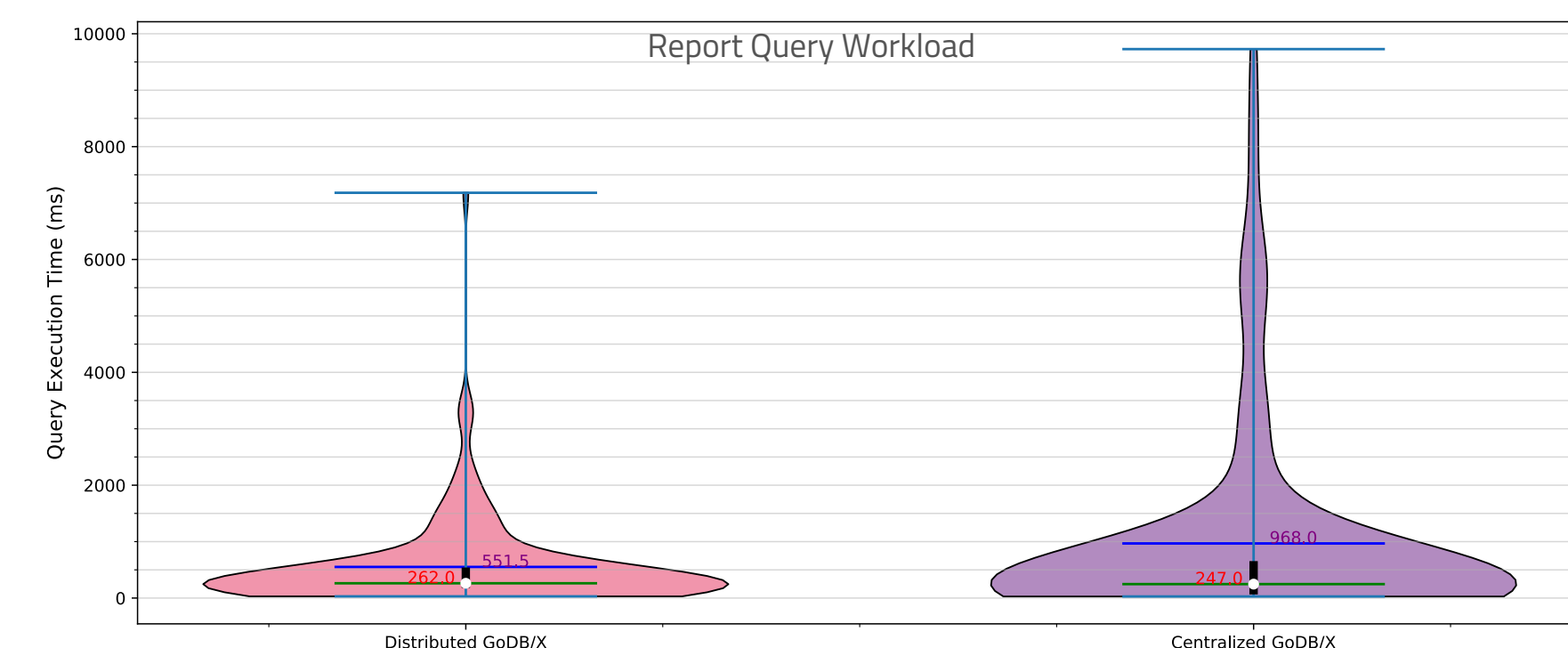


Figure 4: EdgeFile Layout with Implicit VertexIDs



20%-32% Reduction in Query Execution Time compared to previous Data-Model and **>65% Reduction** compared to GoDB. Improvements in Data Model cause **reduction in static overheads** of Succinct Data Store Operations, though communication costs remain the same.

- Cost Model Inclusion
- Support query over Time-Series Graph with **compound predicates**
- Concurrent Queries** Optimization

SUMMARY

- Fewer Machines \Rightarrow Smaller Memory Footprint \Rightarrow Reduced Communication Costs \Rightarrow Lower Query Latency

