

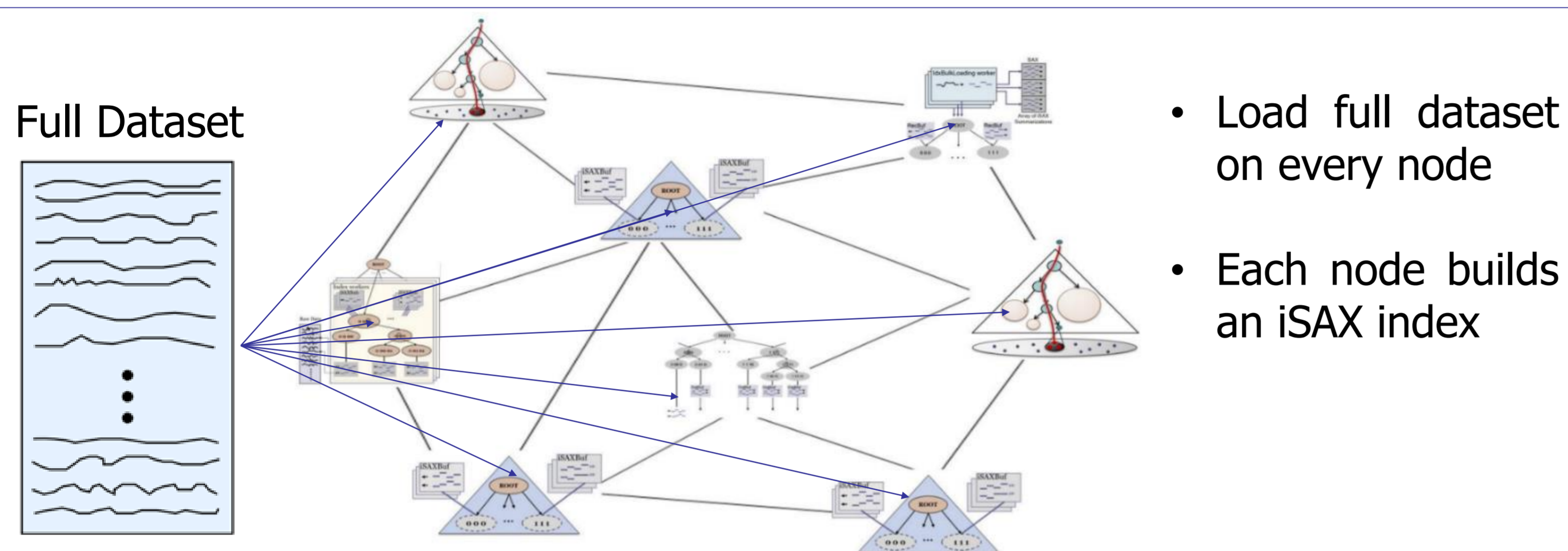
Distributed Query Answering for Large Data Series Collections

Manos Chatzakis, Panagiota Fatourou, Eleftherios Kosmas and Themis Palpanas

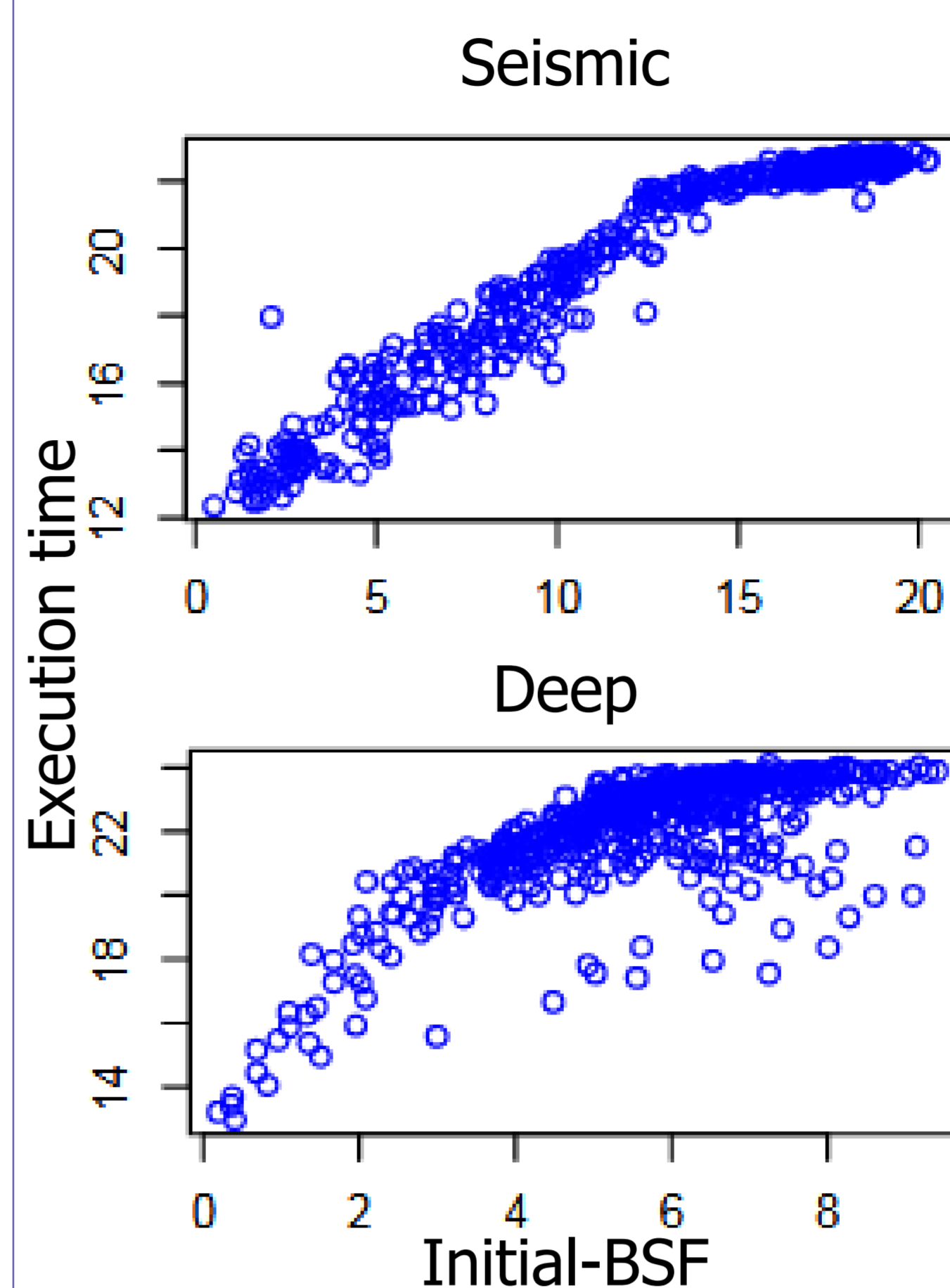
Motivation

- Achieve very fast query answering
- Assume that nodes have enough memory to maintain the whole dataset
- Design efficient scheduling algorithms to assign the queries to the available nodes

Full Data Replication



Predictions of Query Execution Time

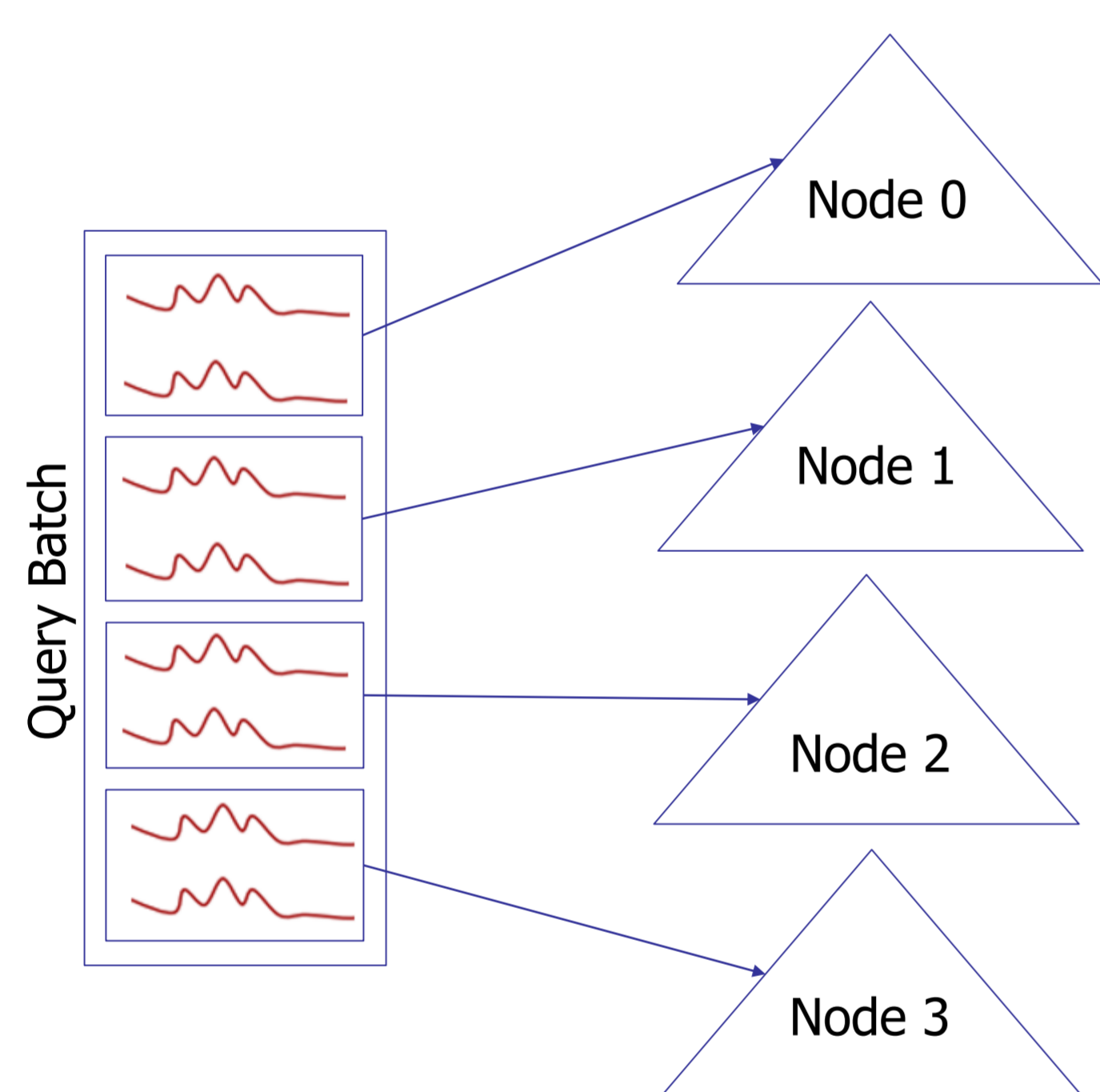


- Correlation between initial BSF and total execution time
- Predict total execution time per query
- Query scheduling based on these predictions



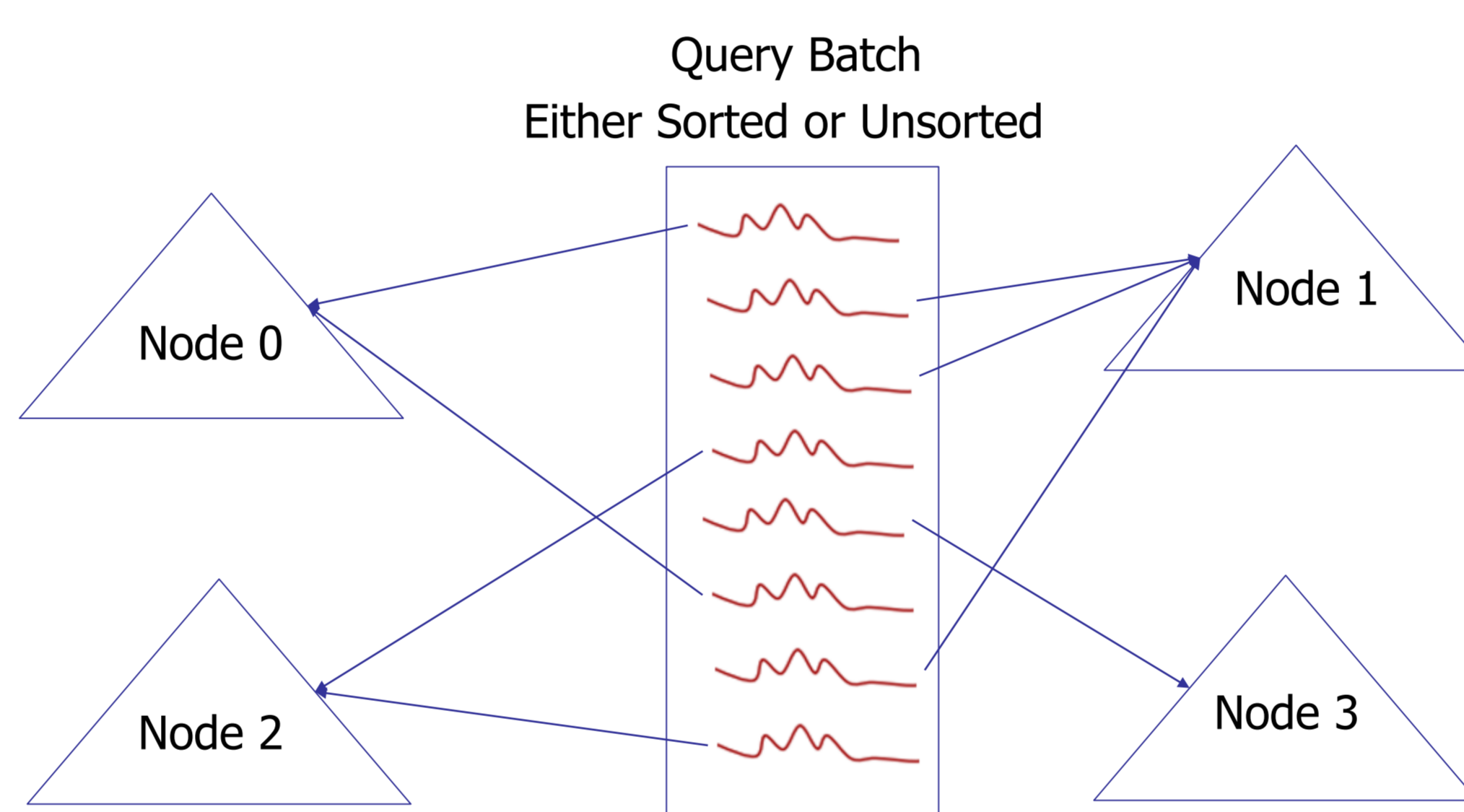
Distributed Query Answering

Static Scheduling



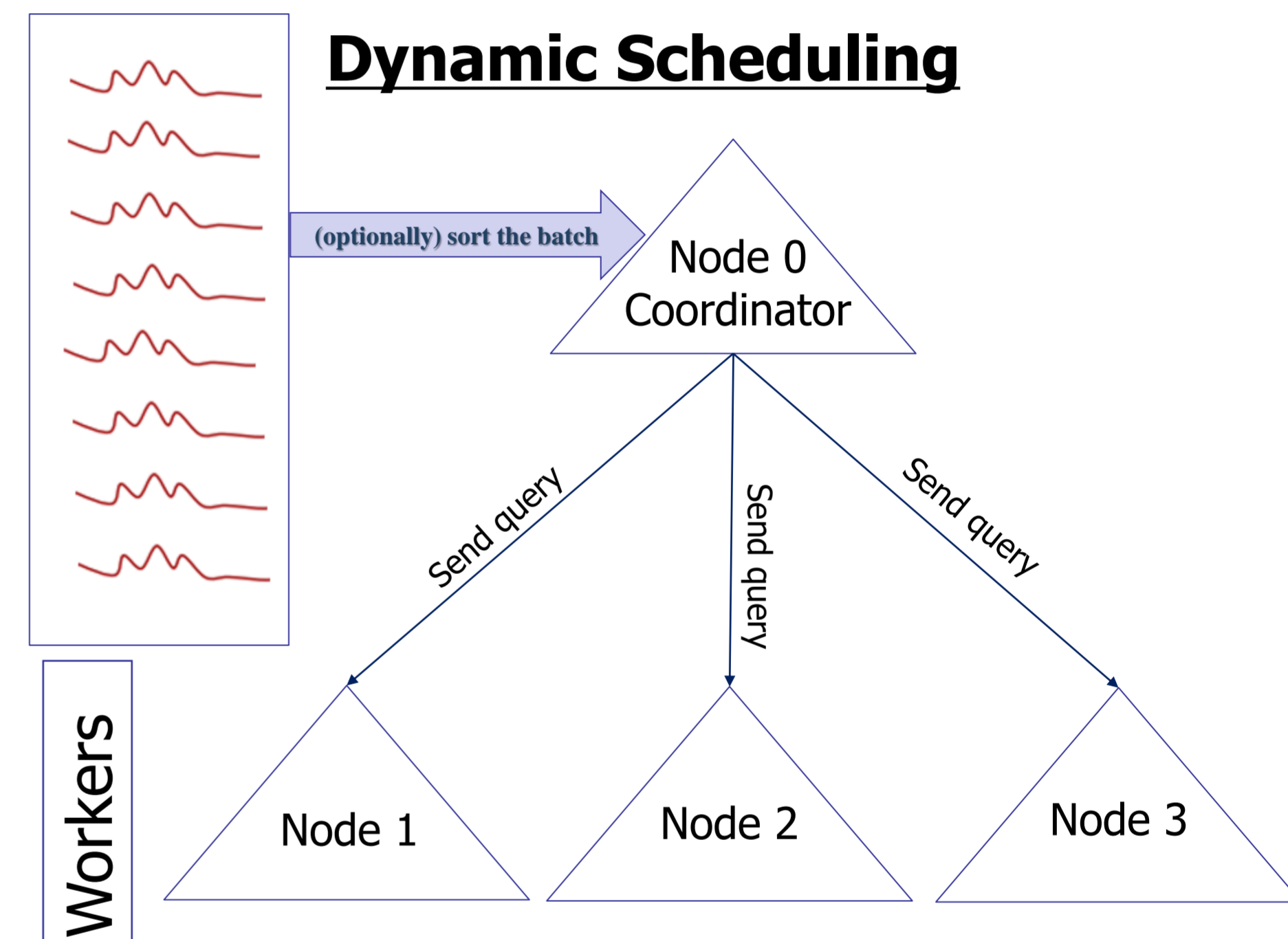
- Split the query batch into chunks

Prediction-based Scheduling



- Static scheduling based on predictions to balance workload among different nodes

Dynamic Scheduling



- Each worker node requests a new query from a coordinator, which
- optionally sorts the batch based on predictions
- assigns queries to the workers

Evaluation

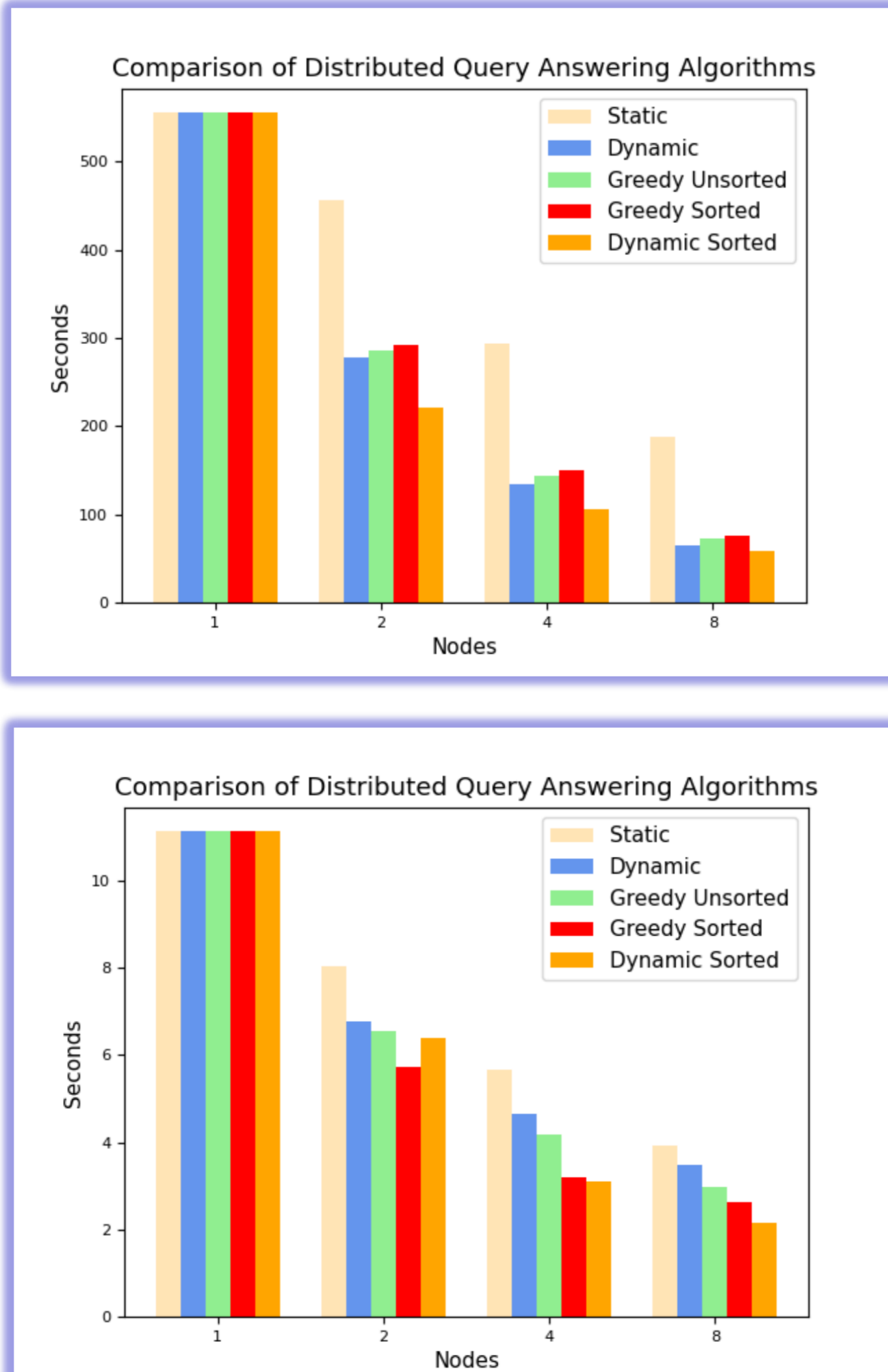
Results

Configurations and Datasets

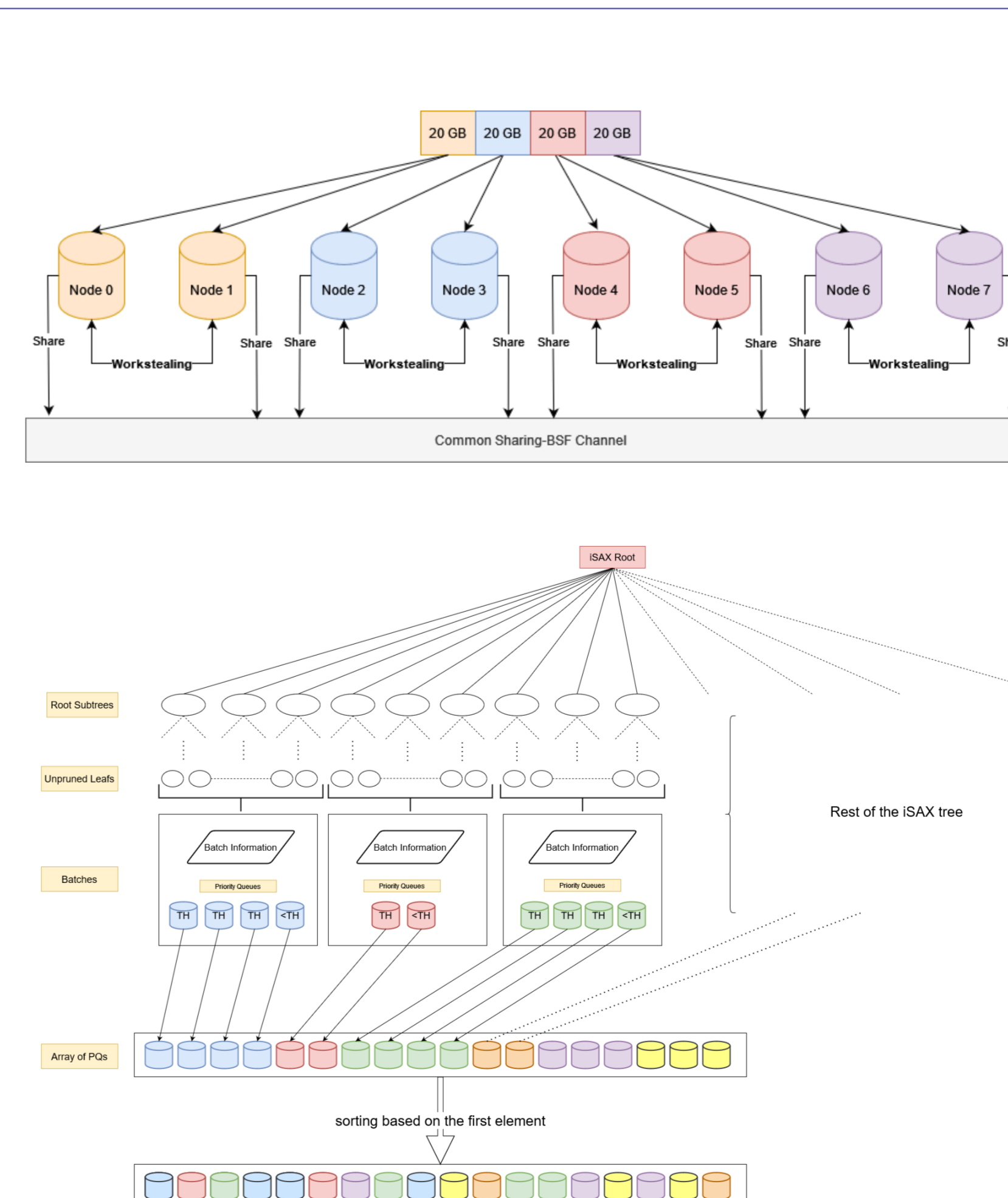
- Seismic dataset (seismic records)
- Query batches:
 - Big Batch: ~1000 queries
 - Small Batch: ~200 queries

Experimental Findings

- Dynamic scheduling is better than static
- Scheduling based on predictions performs better.
- Linear Speedup



Limitations and Current work



Current Work

- Memory limitations
- Extend the algorithms for Partial Data Replication
- There exists query batches for which all algorithms have poor performance
- Design work-stealing schemes to tackle load imbalances between nodes