

Workload-Driven Horizontal Partitioning and Pruning for Large HTAP Systems

Paris, 16 April 2018

Martin Boissier & Daniel Kurzynski

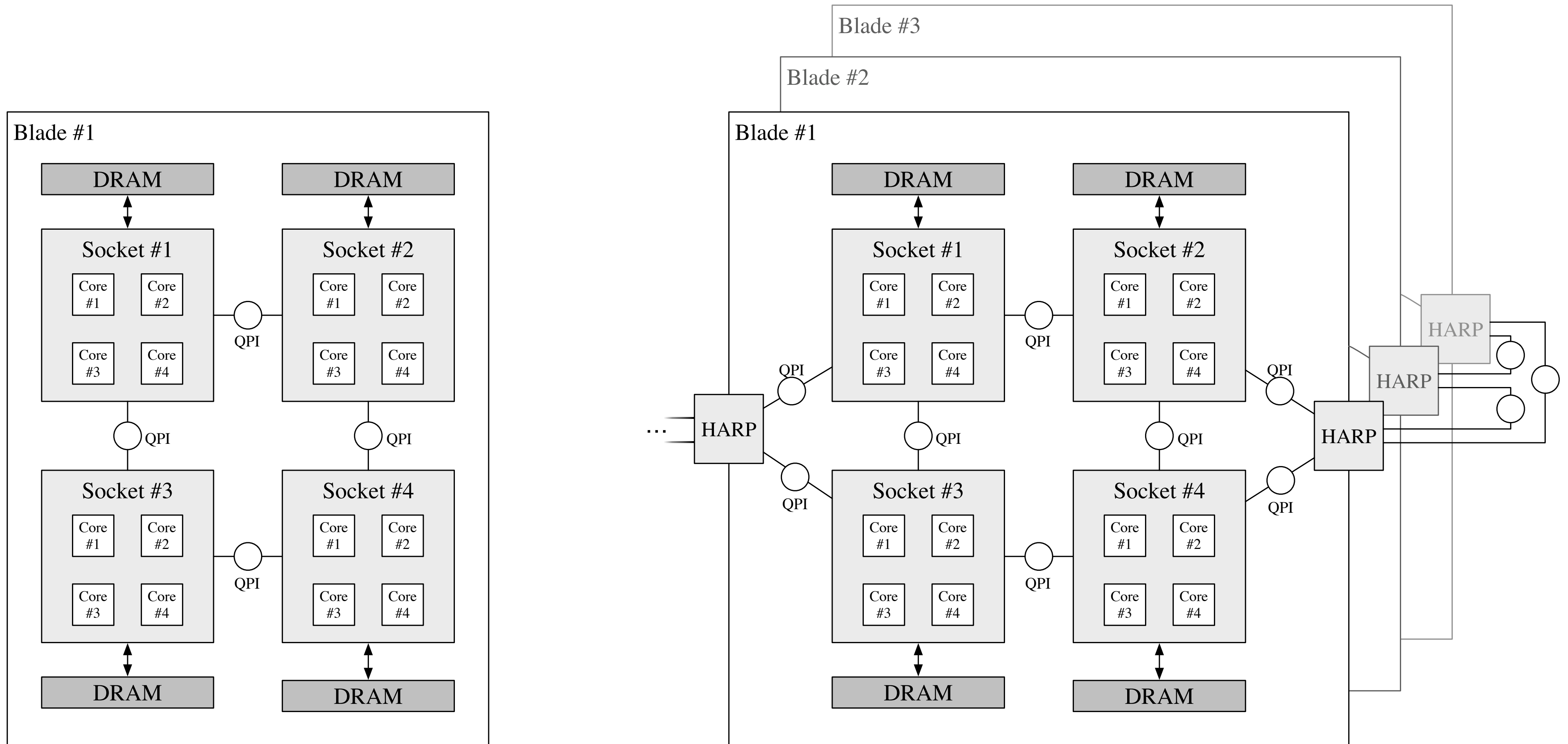
Project Goal

- **Background:**
 - In our research, we focus HTAP-optimized main memory-resident databases for modern enterprise systems
 - Own research database called **Hyrise**^{*}, an open source columnar HTAP database
- **Current Situation:**
 - NUMA architectures pose new challenges for data distribution, scheduling, and query execution
 - Optimizing for NUMA has large impact (cf. [1, 2])
- **Our Goal:**
 - Evaluate various partitioning approaches to distribute a table's data to n NUMA nodes with the goal to *maximize tuples skipped* (i.e., data skipping)

[1] Leis et al., Morsel-Driven Parallelism: A NUMA-Aware Query Evaluation Framework for the Many-Core Age. SIGMOD '14

[2] Pandis et al., Data-Oriented Transaction Execution. VLDB '10

Modern NUMA Systems



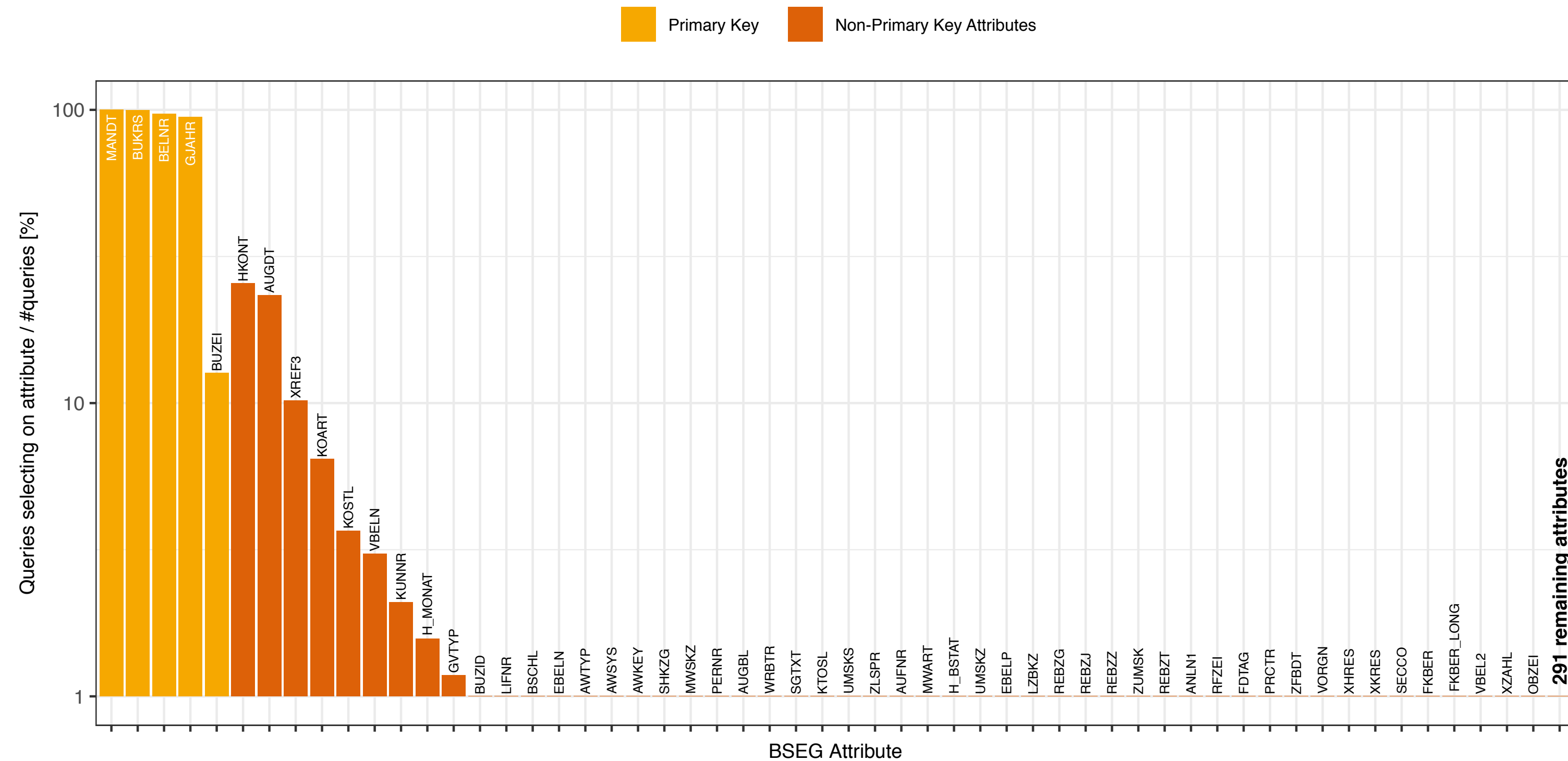
Modern NUMA Systems

- For optimal performance on NUMA systems:
 - data shall be equally distributed
 - processing shall be data-local *
- We see **two problems** with this statement:
 - Equal distribution works fine for both TPC-C and TPC-H
it does not for real-world systems
 - Partitioning elimination/pruning is considered an orthogonal topic
it should not be
- This projects evaluates means to “combine both worlds”

* Cf. “NUMA commandments”, Albutui et al., Massively Parallel Sort-Merge Joins in Main Memory Multi-Core Database Systems. PVLDB 2012

Workload-Driven Partitioning

- We evaluated several partitioning approaches
 - almost all of them to be *too simple*
 - Aggressive Data Skipping* by Sun et al. is one exception [3]



Aggressive Data Skipping

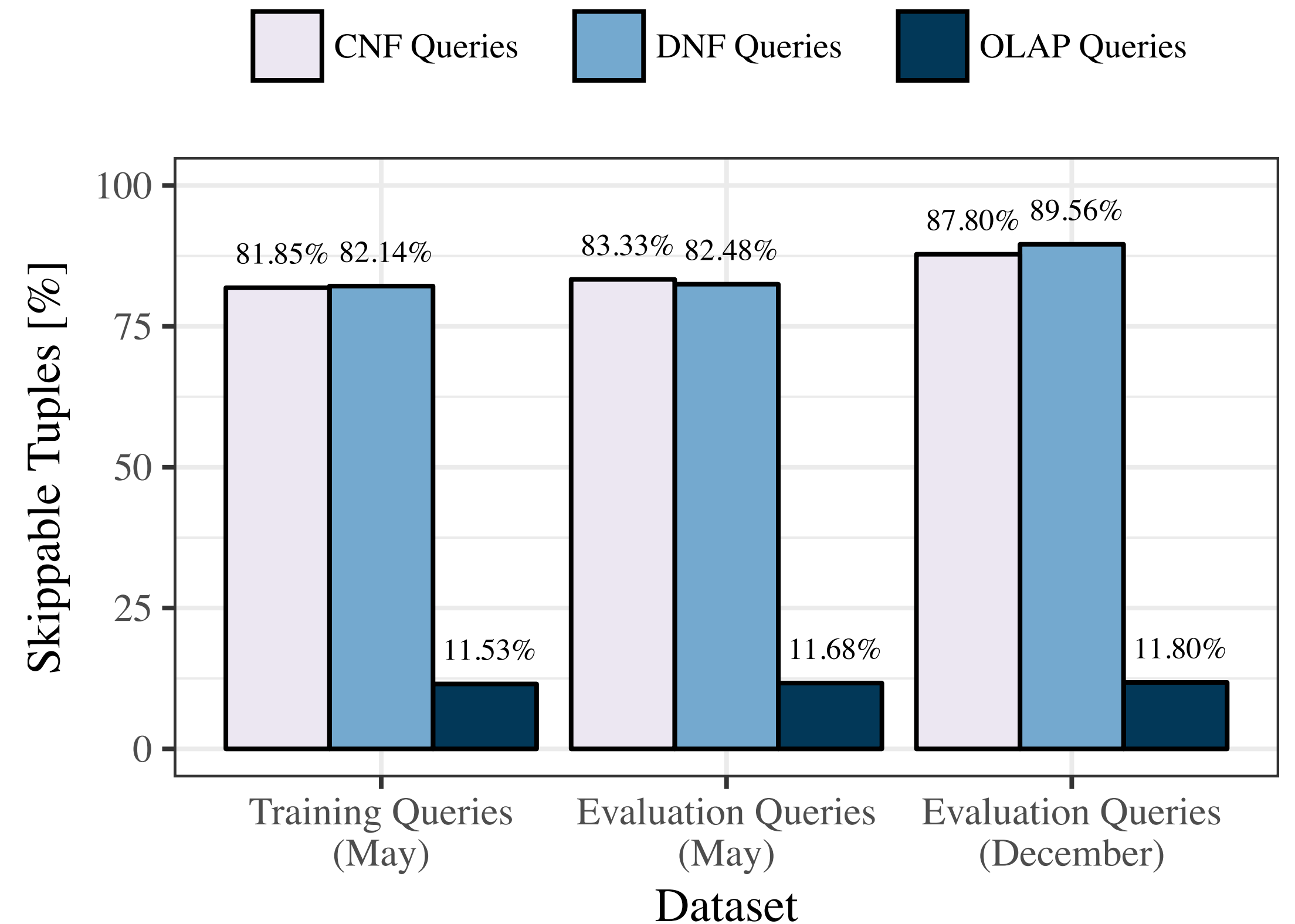
- Approach initially motivated for large-scale systems like Spark
- We misused the approach to create partitioning schemes
 - Configurable by the number of partitions to yield
 - We limit partition count to number of NUMA nodes
- The process
 - Parse workload and extract relevant selections + frequent item set mining
 - Scan data for distribution of features
 - Merge features to create partitions

Created Partitioning Scheme

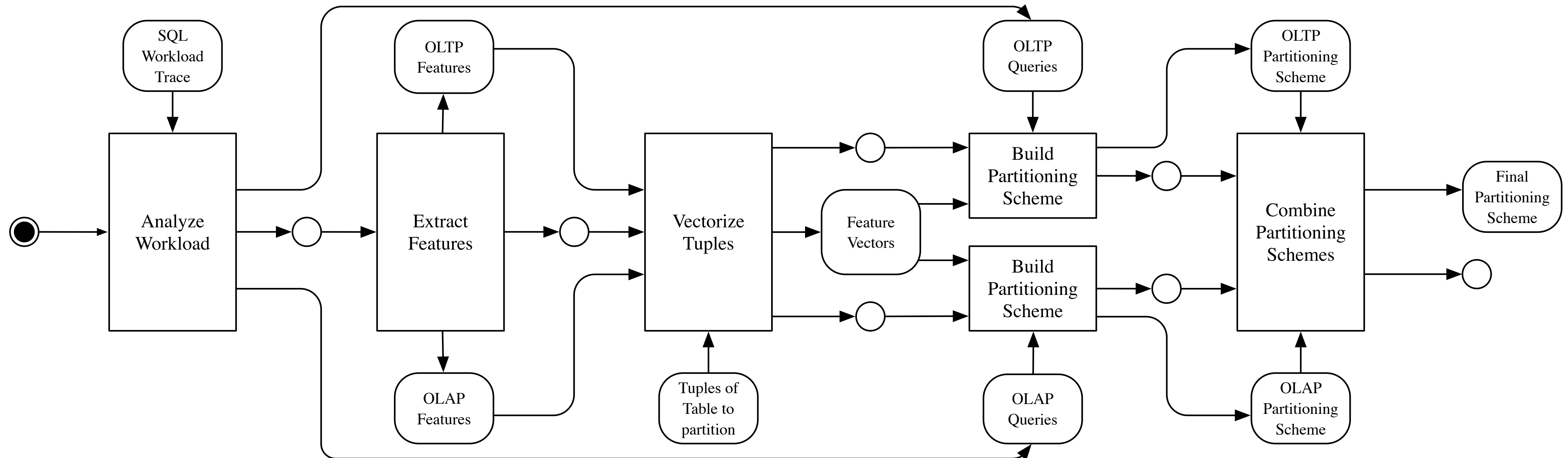
- clustered and non-trivial
 - Feature 1:
`mandt = 2 & koart <> 'k' & koart = 'd'`
 - Feature 2:
`mandt = 2 & koart = 'k' & bukrs = '9999'`
 - ...
- Created partitions are freely defined by 15 features

Aggressive Data Skipping

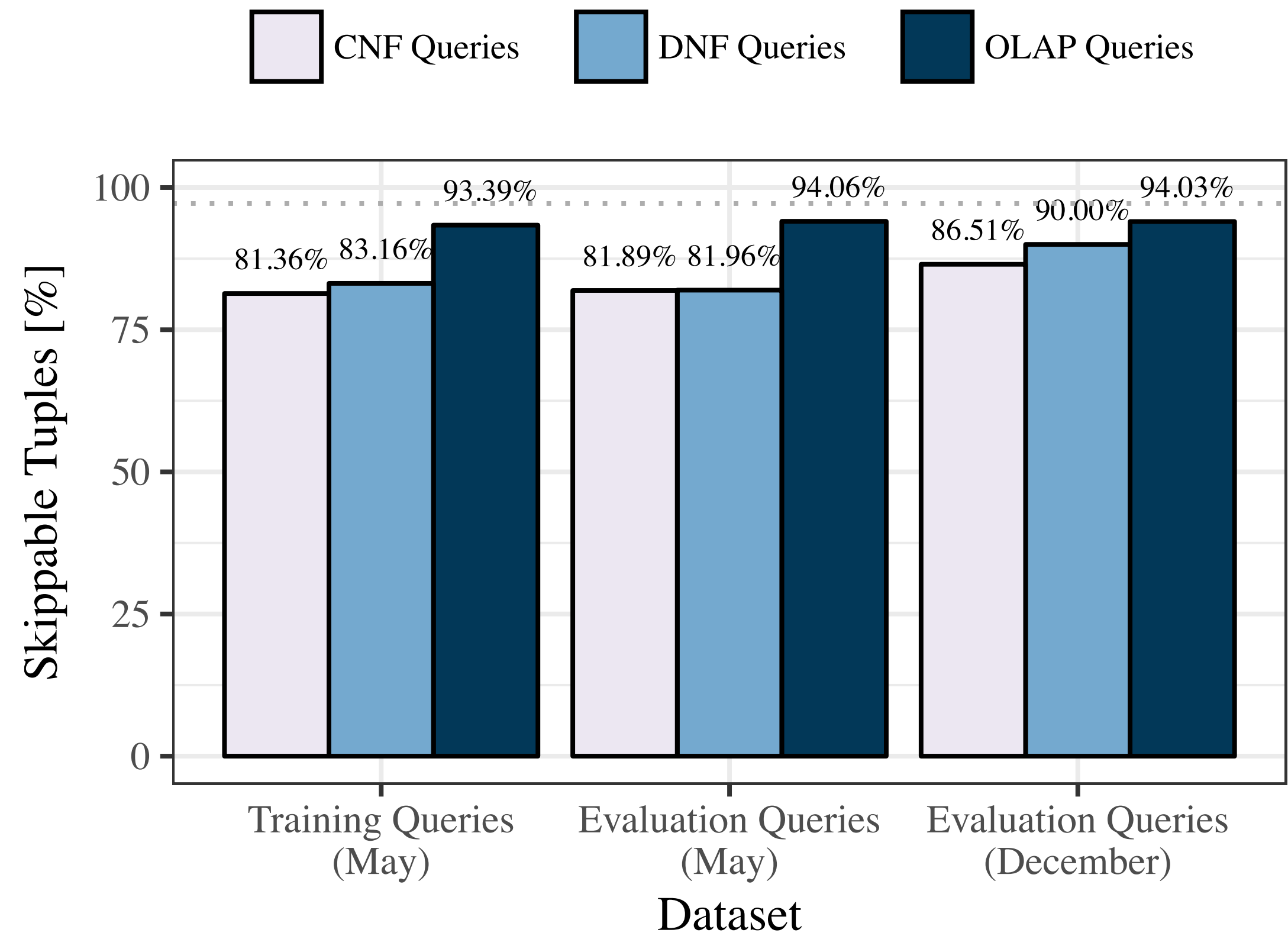
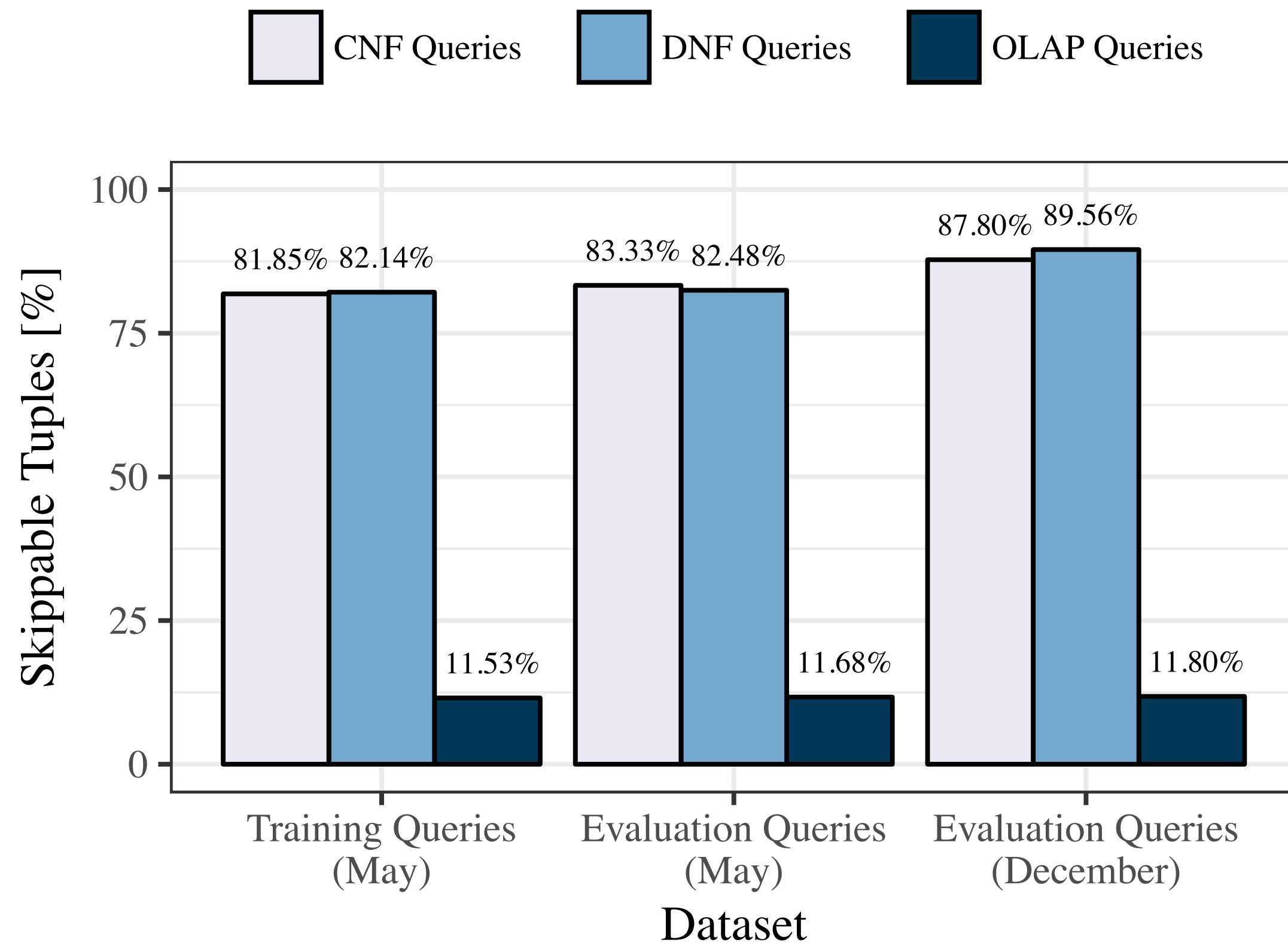
- **Problem:**
 - Objective is the number of pruned tuples
 - Heavily favors selective and frequent Queries
 - i.e., OLTP queries
- Many ways to adapt for multiple workload classes
 - Weighting by runtime, what-if based query costs, ...
- **Idea:** execute partitioning twice and merge



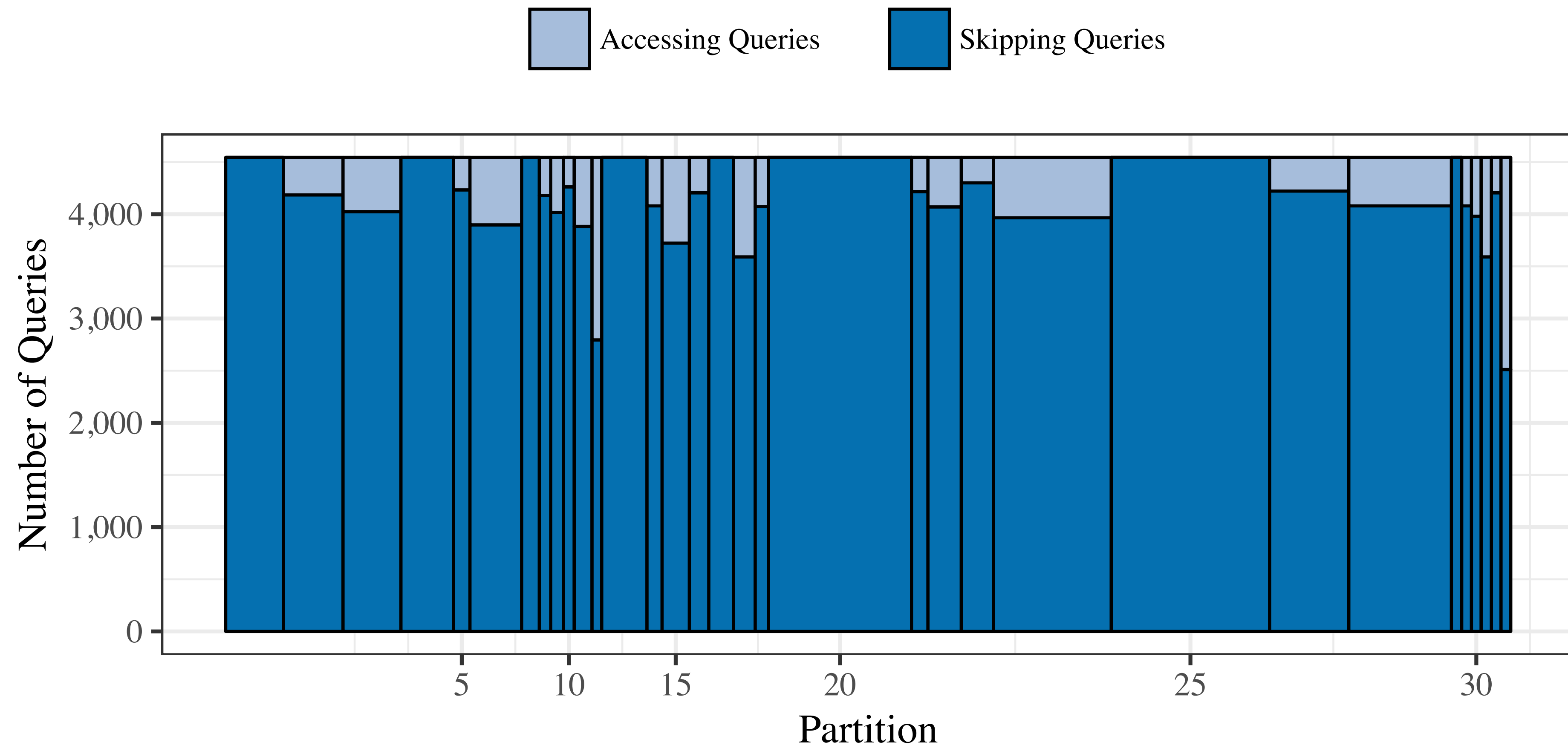
Merging Partition Schemes



Results (i)



Results (ii)



My Personal Outlook

- **Self-driving trend & distributed systems** force us to put more emphasis on (re-)partitioning
- **What is missing?**
 - More work on **skew-aware and pruning-optimized partitioning**
 - **Proper cost models** for adept partitioning schemes
 - Personally, **I doubt AI & DL will solve these problems** for us
 - More emphasis on repartitioning (not only for NUMA systems! [4])