

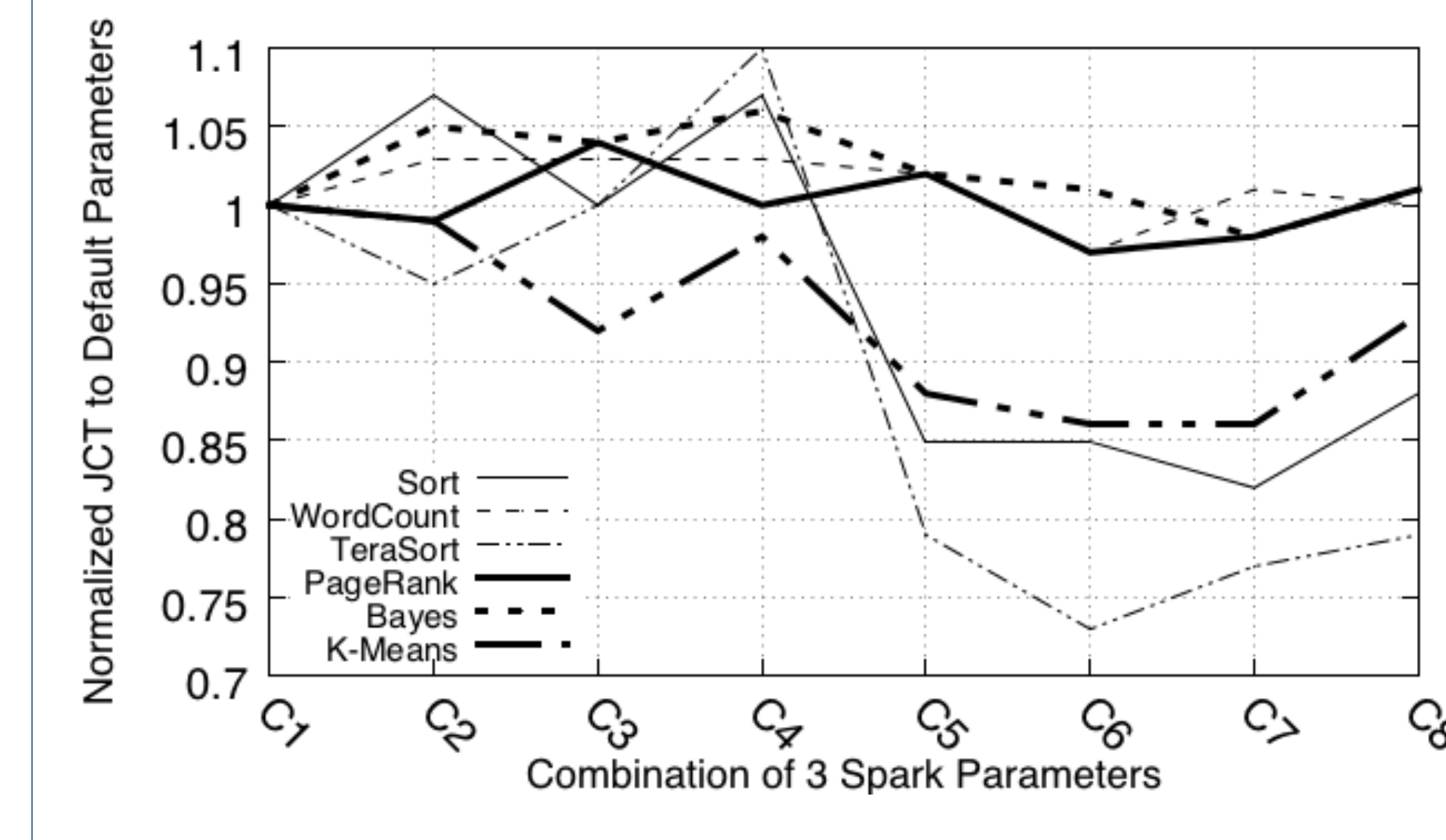
Leveraging Resource Bottleneck Awareness and Optimizations for Data Analytics Performance

Tiago Barreto Goes Perez, Xiaobo Zhou (Advisor)
 University of Colorado, Colorado Spring
 University of Colorado
 Colorado Springs

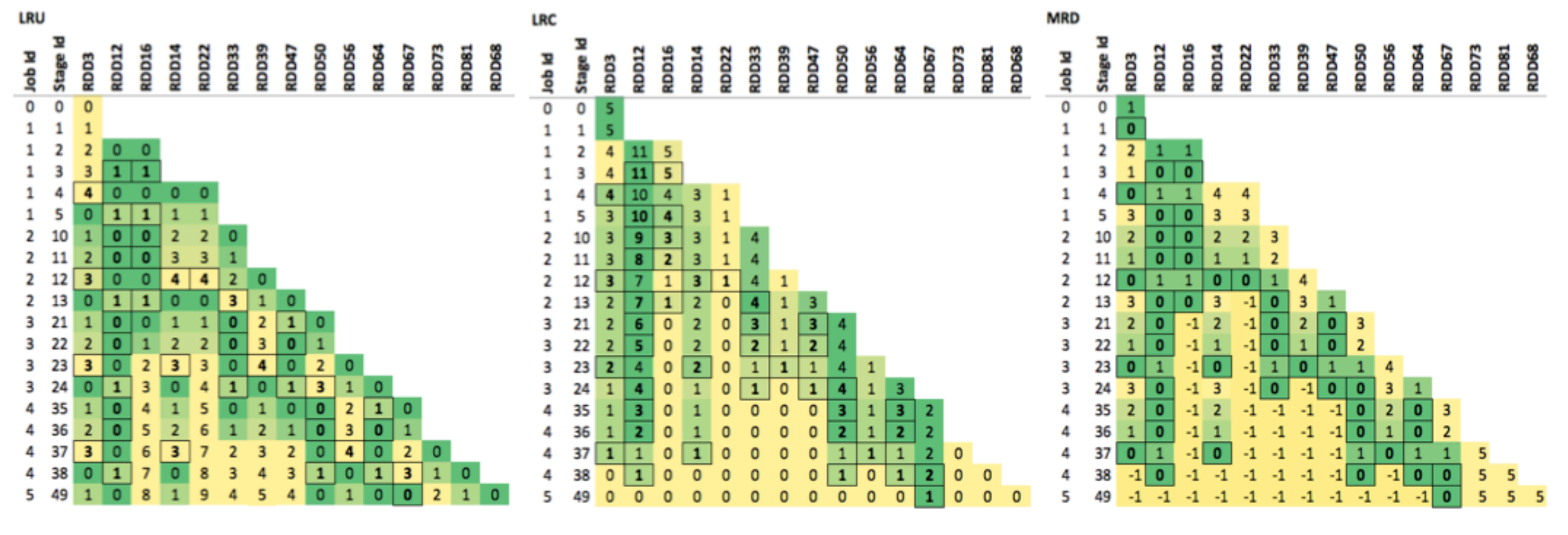
Introduction

- Spark Tuning:**
- Current tuning work is resource-oblivious
 - Spark has dozens of performance affecting parameters, previous tuning work modifies one at a time
- Spark Memory Management:**
- Spark LRU caching policy is DAG-oblivious
 - Current DAG-aware solutions do not take into account the reference-distance and gaps between data usage

Motivation

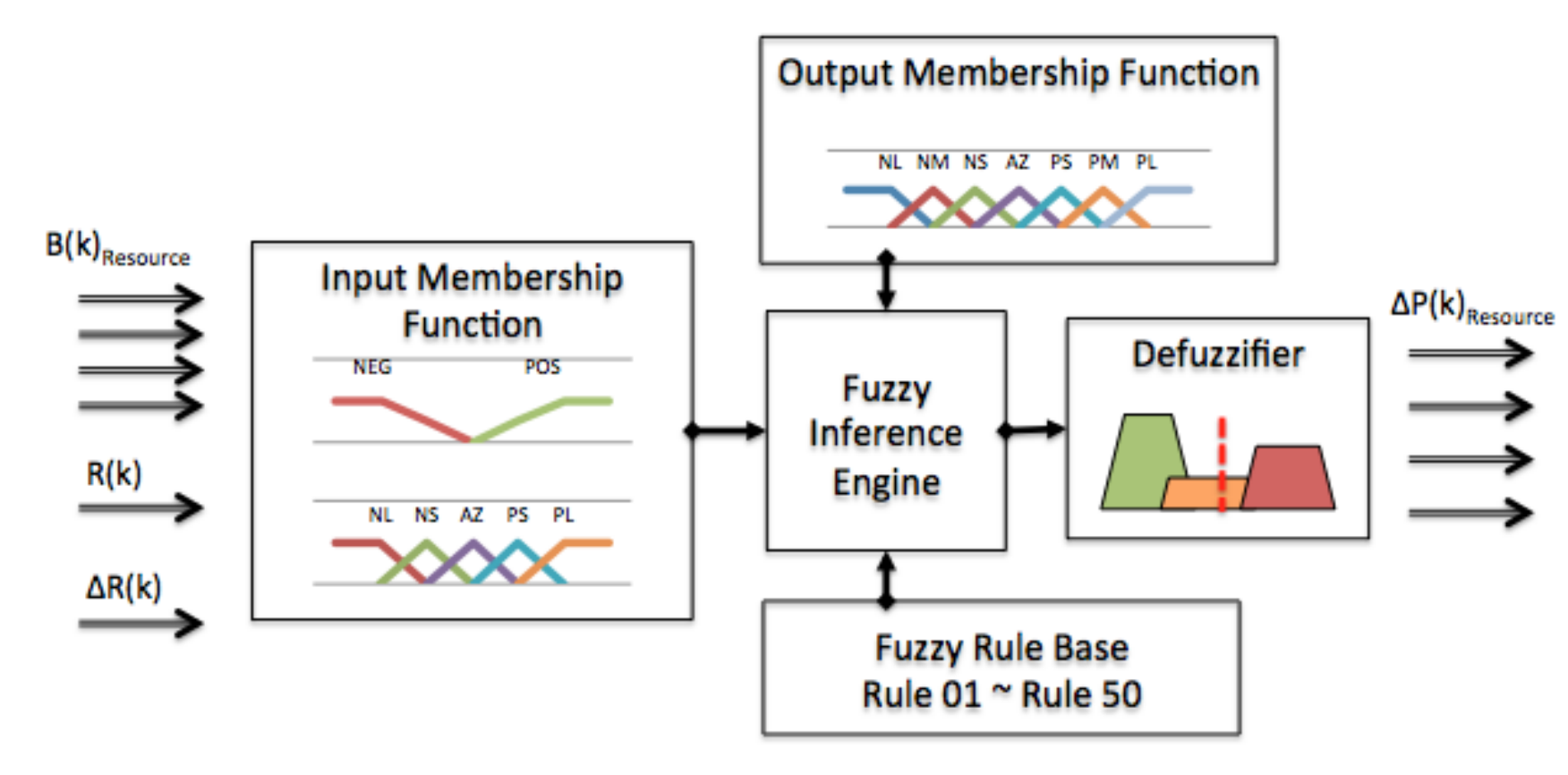


- PETS:**
- Comparison of effect of 3 Spark parameter combinations*
- MRD:**
- Comparison of different cache management policies for ConnectedComponents*

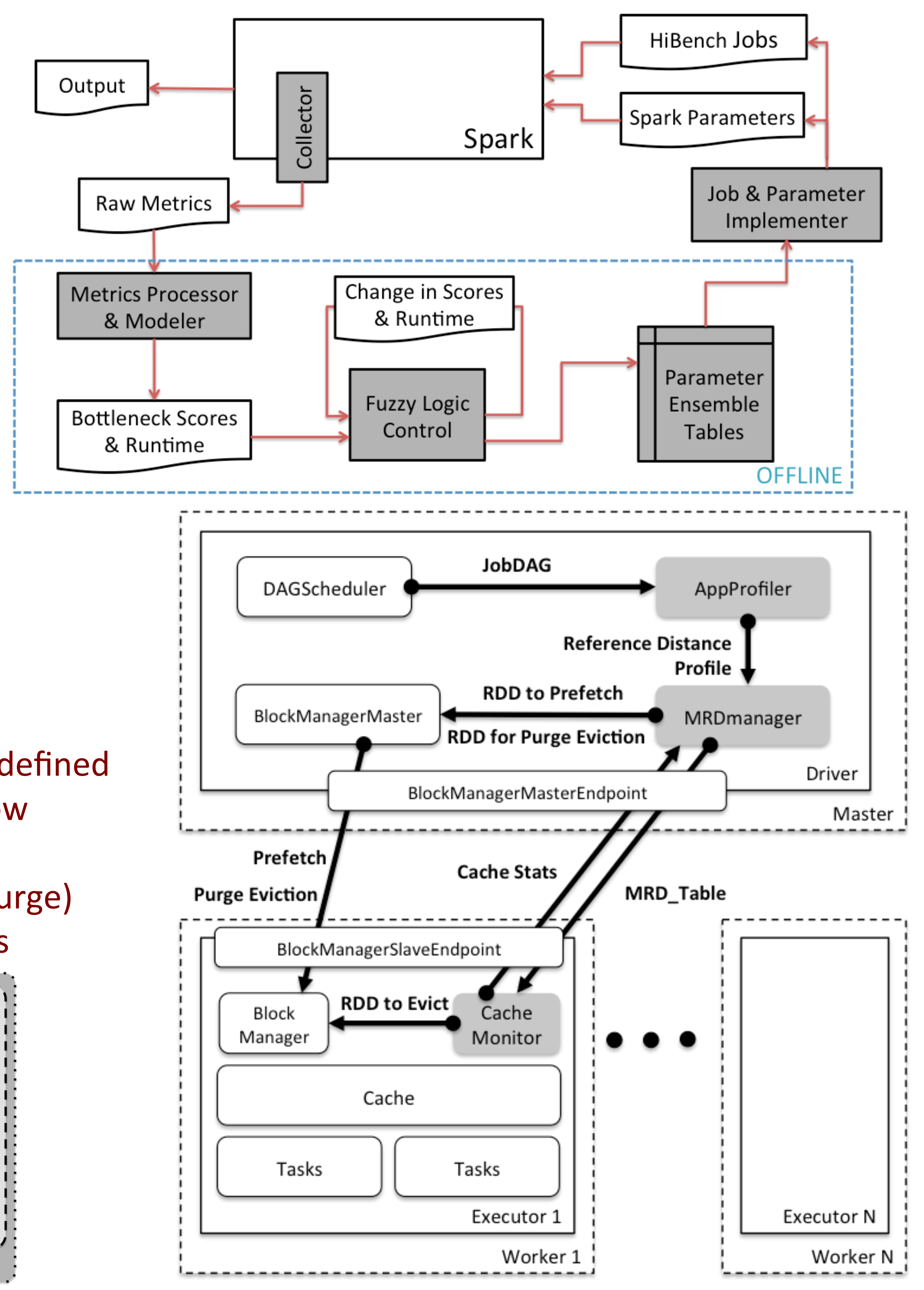
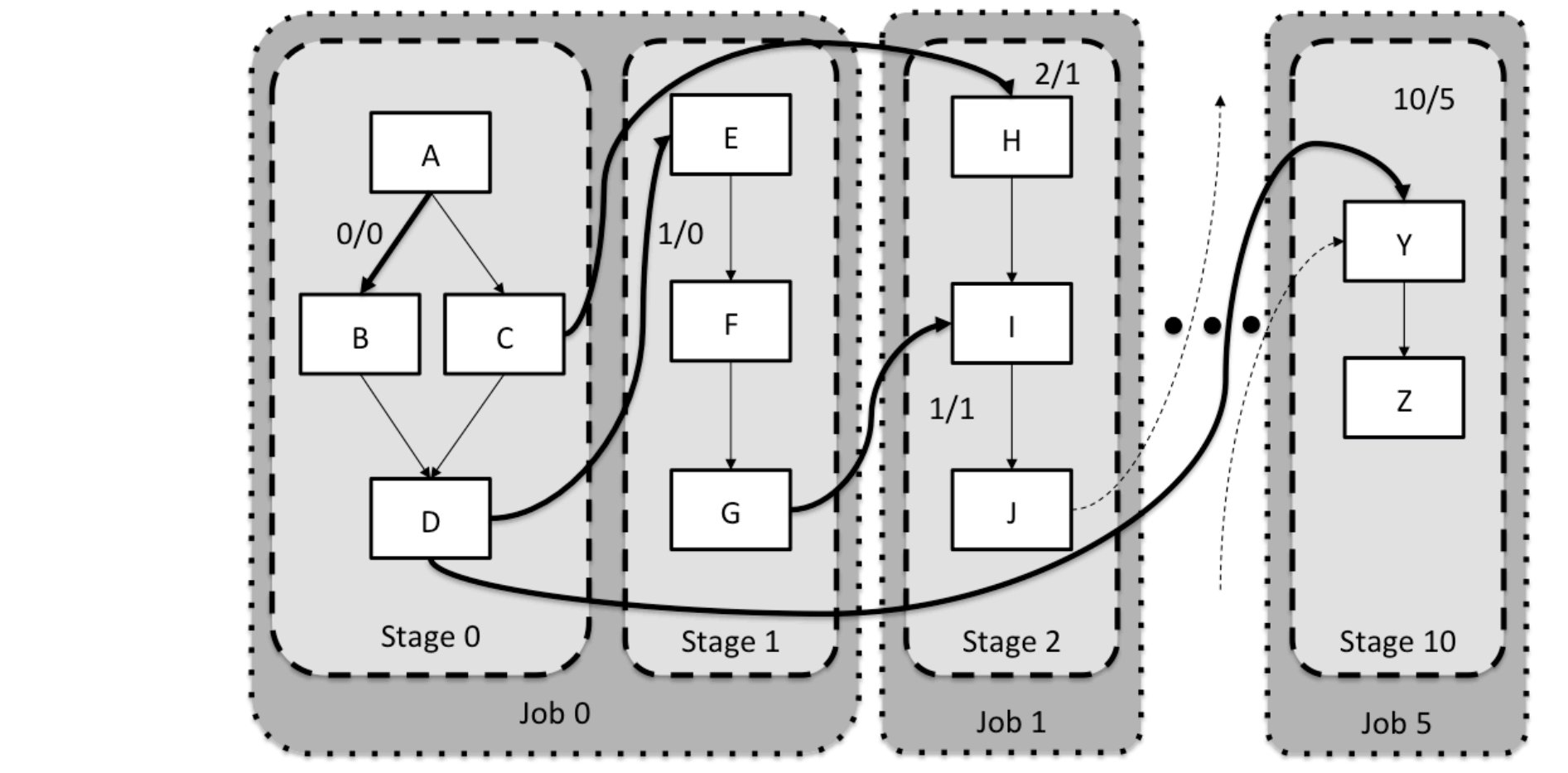


Design

- PETS:**
- Use of Fuzzy Logic with resource awareness feedback
 - Tuning is expedited by the use of Parameter Ensemble Tables, which allow multiple parameters to be tuned simultaneously

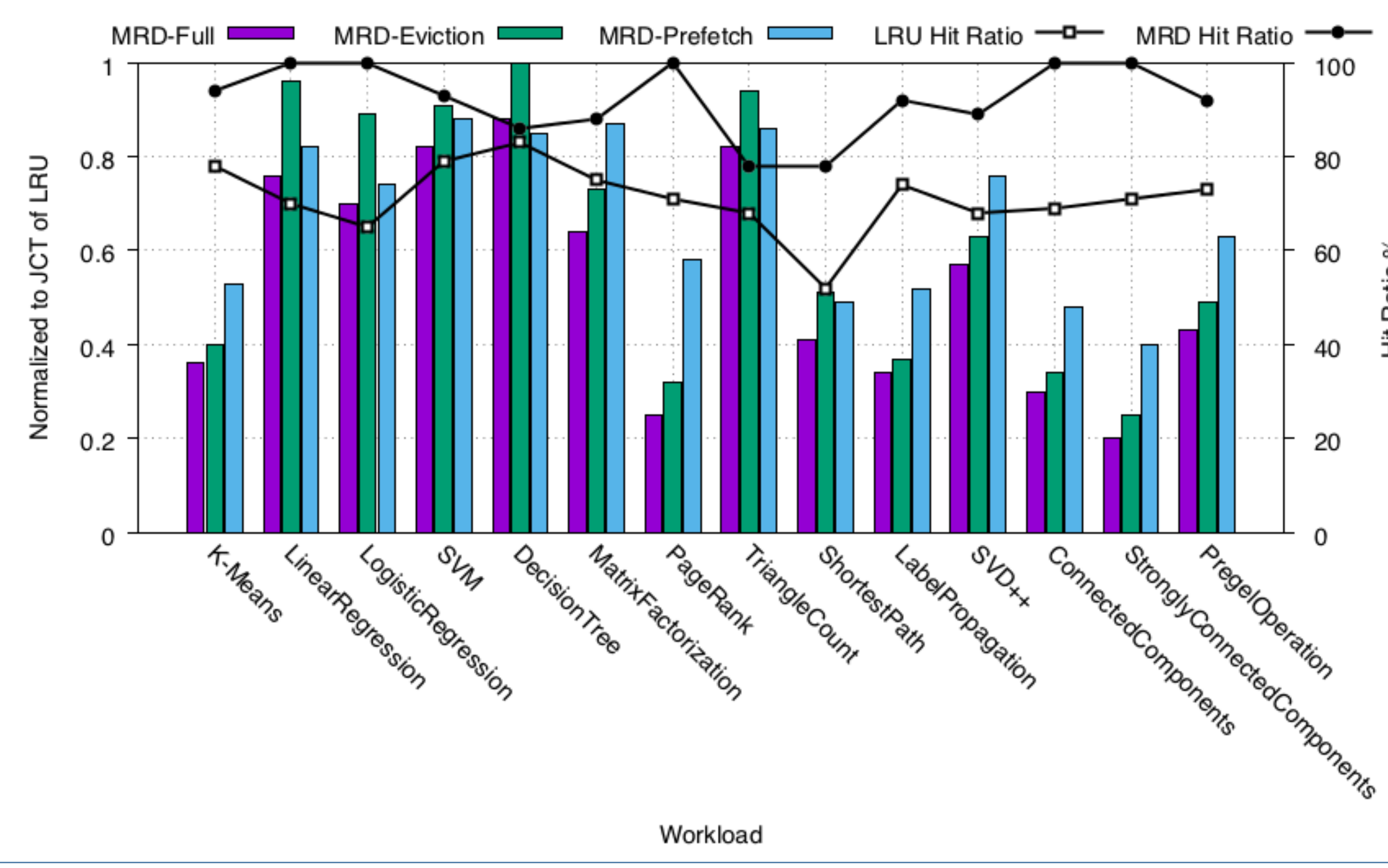
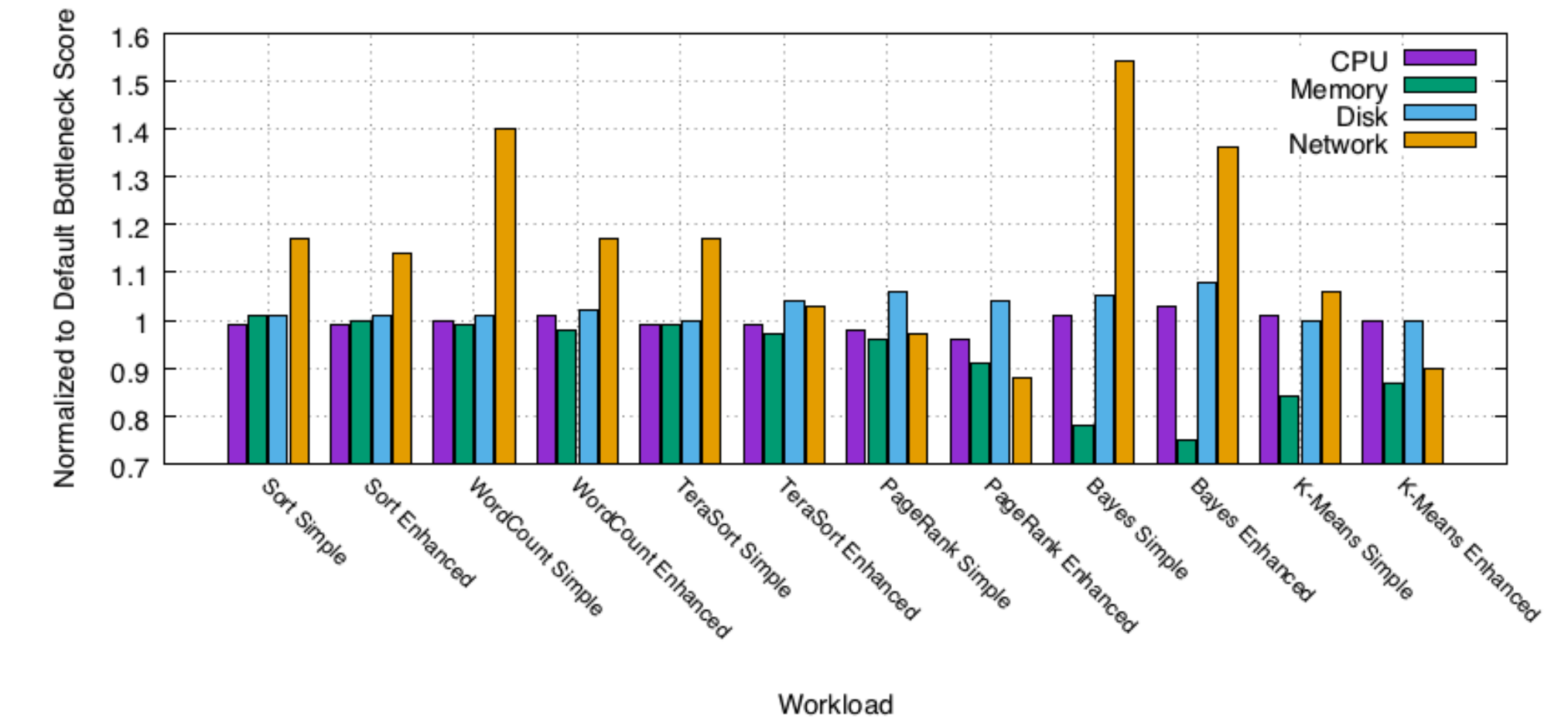
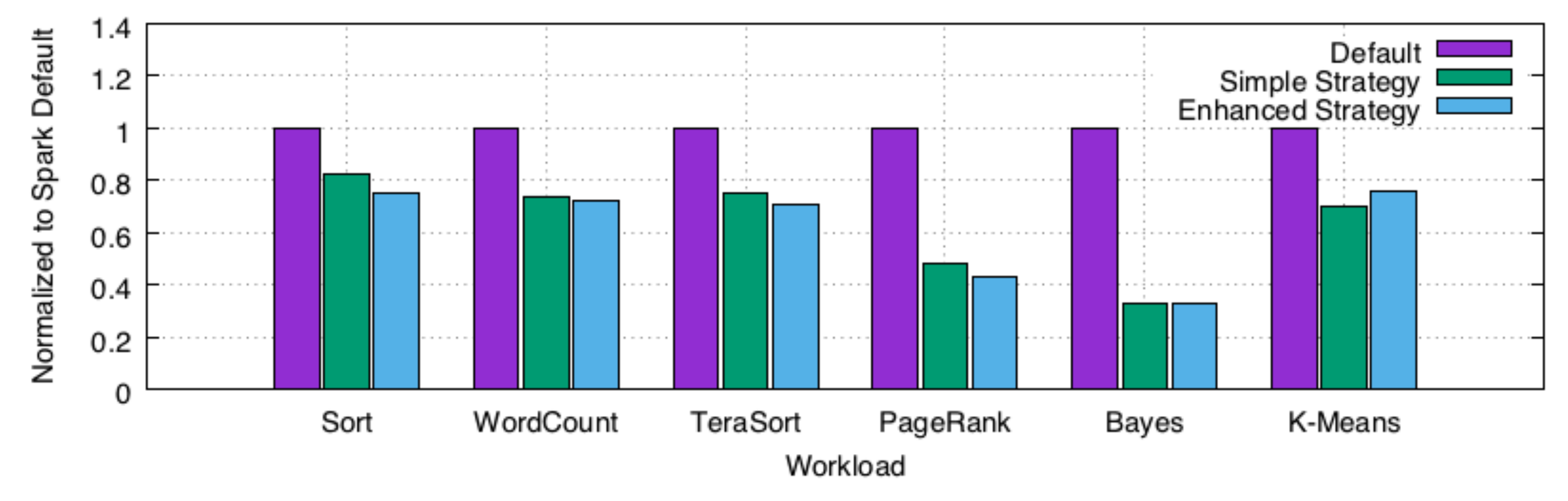


- MRD:**
- Reference-distance (job and stage) is defined the distance between current workflow processing and data block usage
 - MRD has centralized (pre-fetch and purge) and distributed (eviction) components



Evaluation

- PETS:**
- Speedups of up to x4.78;
 - Convergence as low as 2 iterations;
 - Performance stable with varying workload data sizes, homogenous and heterogeneous clusters, and varying initial parameters.



- MRD:**
- Average performance improvement over LRU by 53% and up to 4x faster;
 - Improvement over other DAG-aware caching policies up to 68%;
 - Best results with workloads that are I/O intensive, have high stage-reference distance and high reference per stage values.

