

Background and Motivation

Performance in multi-tenant data-analytics frameworks

- ❖ Disparity of performance according to configuration of data-analytics frameworks
- ❖ Minimize execution time of data-analytics jobs (user perspective)
- ❖ Maximize resource utilization of clusters (provider perspective)

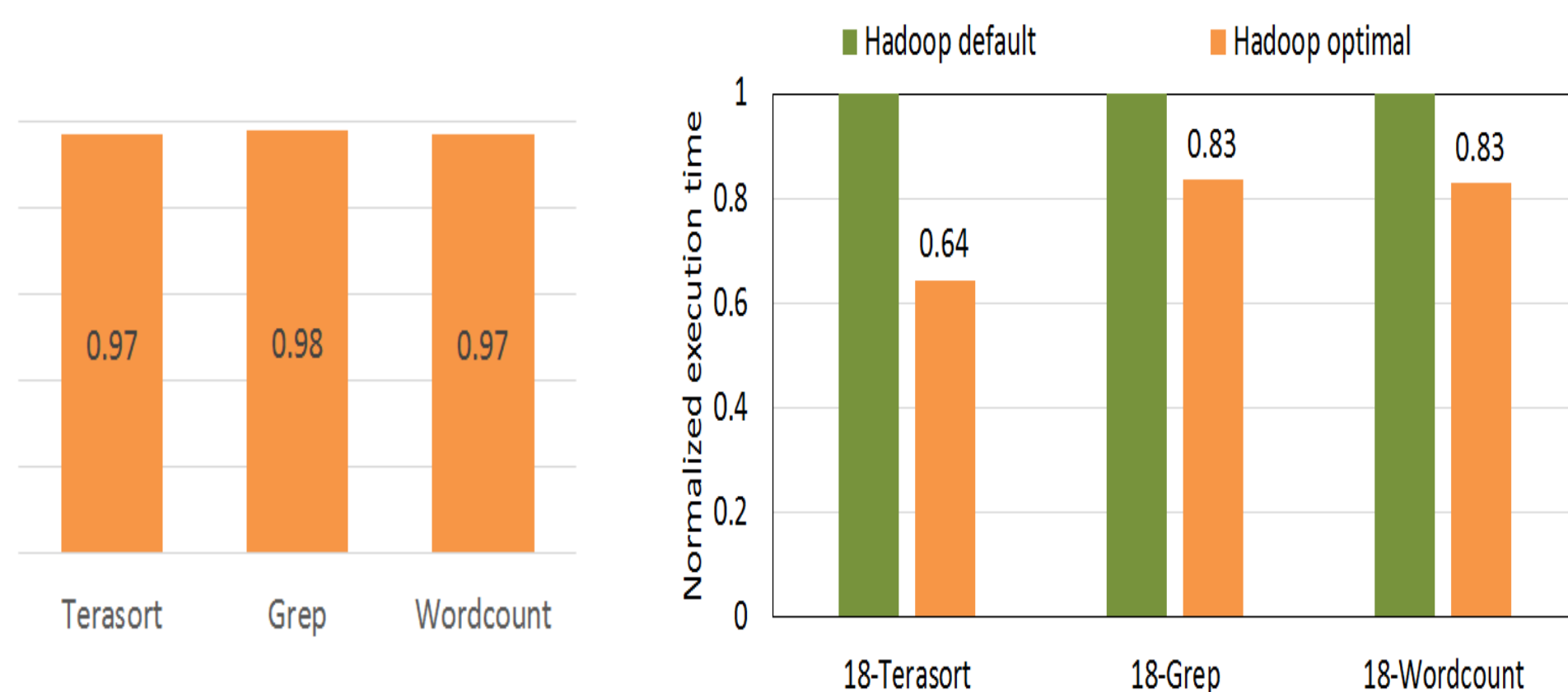


Tuning challenges for performance in data-analytics frameworks

- ❖ Too many tunable configuration parameters
- ❖ Parameter interdependencies
- ❖ Domain (application, system) specific knowledge requirement

Hadoop behavior

- ❖ CPU usage vs. MapReduce
- ❖ Makespan vs. job concurrency (MR) tasks



Self-Tuning Approaches and Challenges

Approach	Features	Challenges
Individual MR job	<ul style="list-style-type: none"> Performance prediction model Various techniques (e.g., Machine learning) 	<ul style="list-style-type: none"> Time-consuming training Renew or retraining models
MR jobs concurrency	<ul style="list-style-type: none"> Dynamic adjustment of concurrent MR jobs using log-based resource usage Dynamic adjustment of concurrent MR tasks Monitoring system-metric resource usage 	<ul style="list-style-type: none"> Insufficient relationship between workload and resource usage of MR jobs Modification of existing Hadoop systems
Slot utilization	<ul style="list-style-type: none"> Dynamic slot allocation to MR tasks in Hadoop 1 	<ul style="list-style-type: none"> Inapplicable to recent container-based Hadoop 2 system

Autonomic Control of MR Job Concurrency

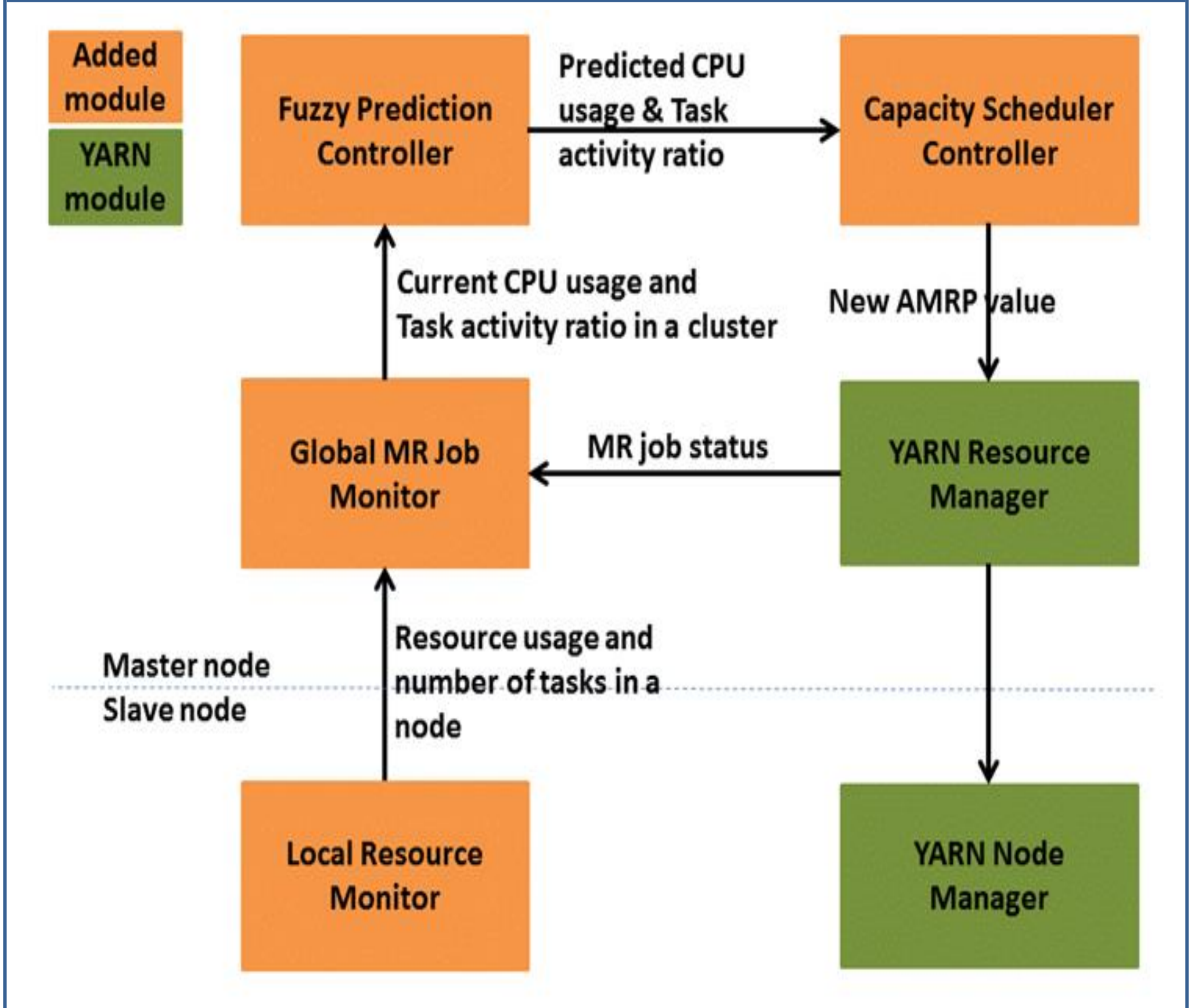
Goals

- ❖ Minimize overall execution time of multiple concurrent MR jobs
- ❖ Use unmodified Hadoop software

Scope

- ❖ Hadoop version 2.7.x
- ❖ Default Hadoop Capacity scheduler
- ❖ Use of available monitoring metrics

Architecture Overview



Results

Performance near static optimal

Scalability over size of datasets and types of MR jobs



Conclusion and Future Work

- ❑ Autonomic approach to adapt the level of concurrency in a multi-tenant cluster
- ❑ Demonstrated in context of Hadoop
- ❑ Consistent performance improvement over static solutions
- ❑ More improvement possible
 - ❖ Accounting for I/O and network behaviors
 - ❖ Task concurrency control
 - ❖ Distributed computing environments (e.g., Chameleon, CloudLab, PRAGMA-ENT)