

Reducing Tail Latencies For Datacenters with High Workload Arrival Rates

Smita Vijayakumar
Third Year PhD Student
Department of Computer Science

Evangelia Kalyvianaki
Supervisor

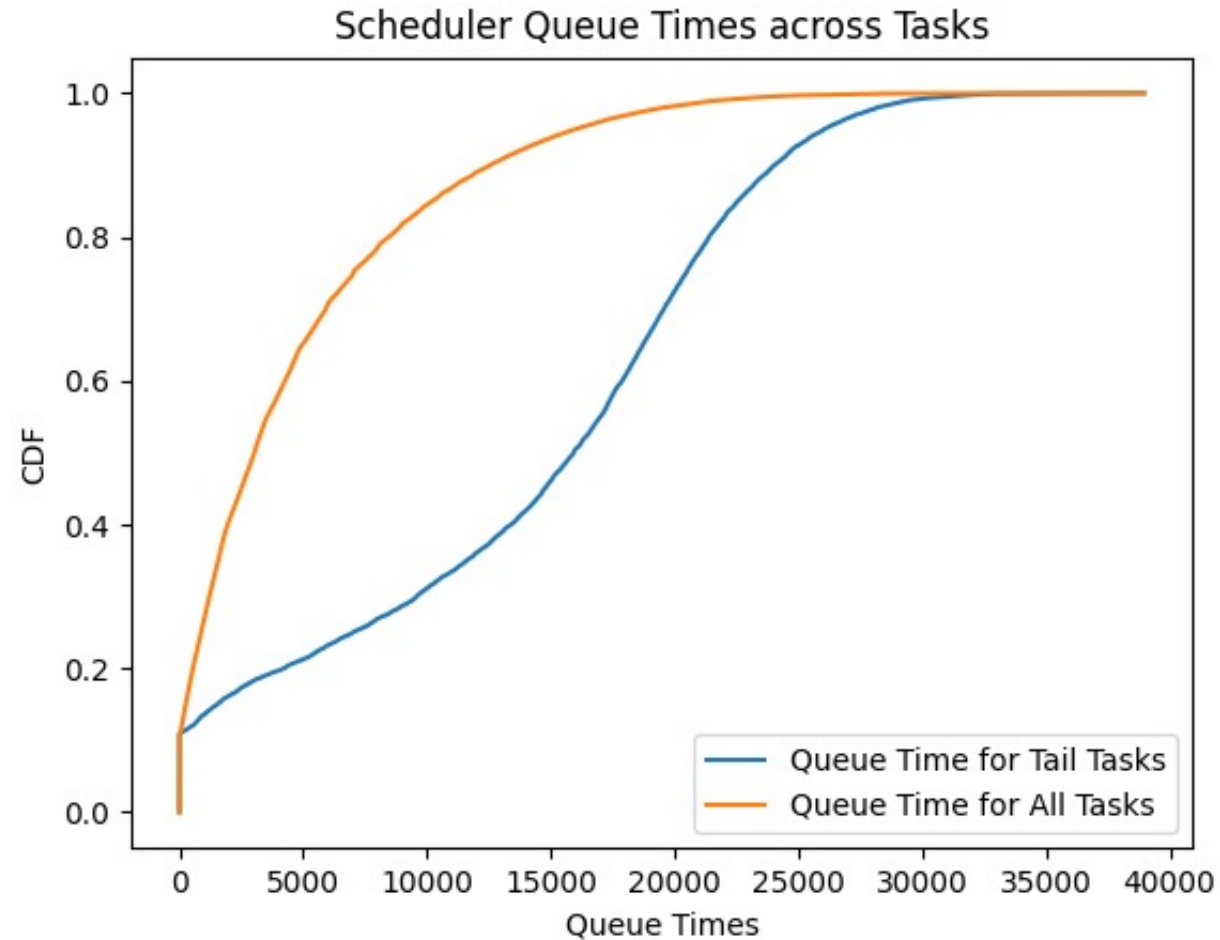
Anil Madhavapeddy
Second Supervisor

Presentation to Huawei • 21st Sept, 2022

Research Motivation

Minimize *Job Completion Times* at Cloud and Workload Scale

Kubernetes Centralized Scheduler



Large scheduler queue times for tail tasks



Large task completion times for tail tasks

Goals

- Small Tail Task Latency
- Better Overall Scheduler Placement Quality
- Scheduling Scale

Design Elements

- Small Scheduler Queue Time



Smaller Tail Task Latency

- Job Level Scheduling
- Estimate Worker Wait Times



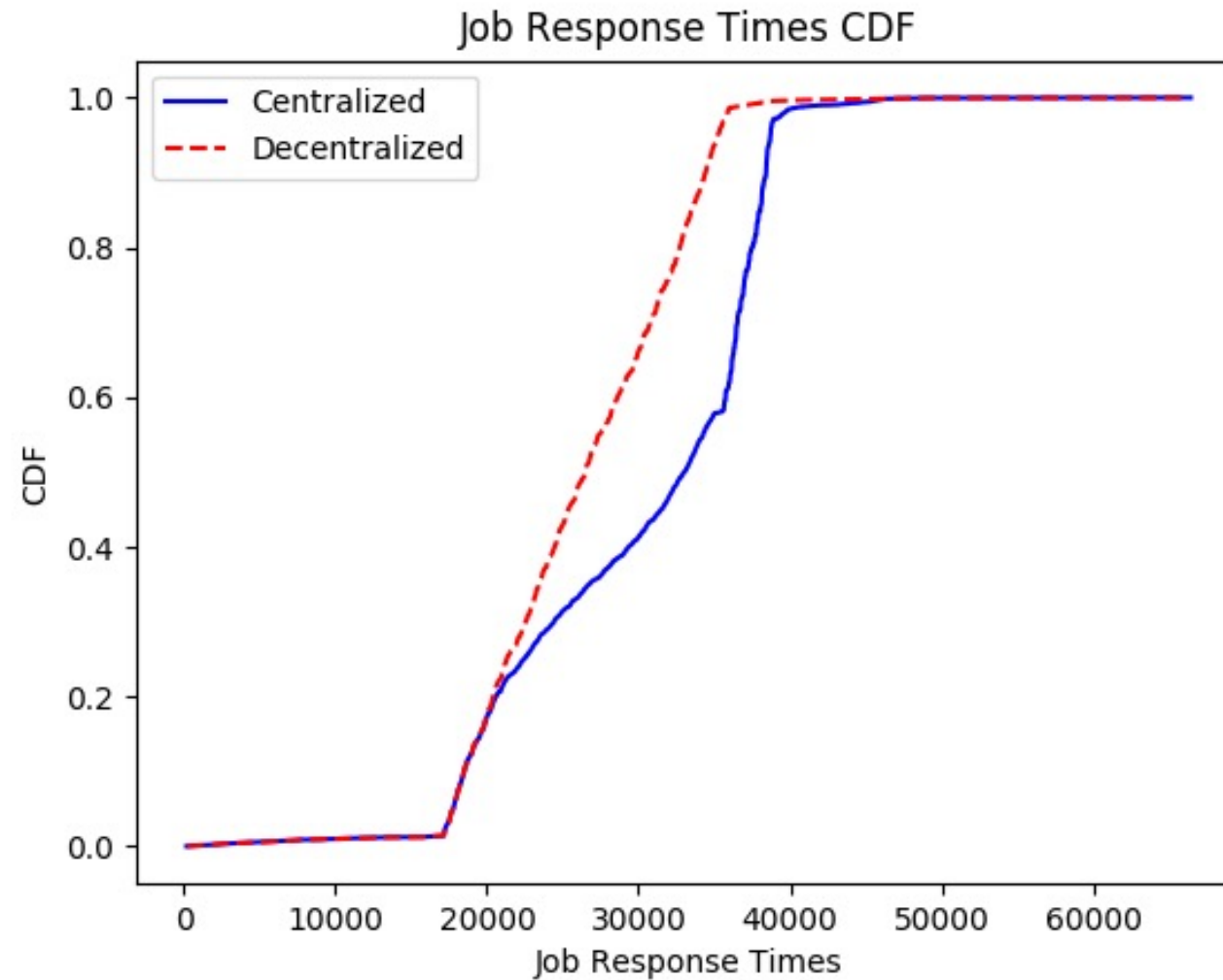
Better Placement Quality

- Loosely coordinating schedulers

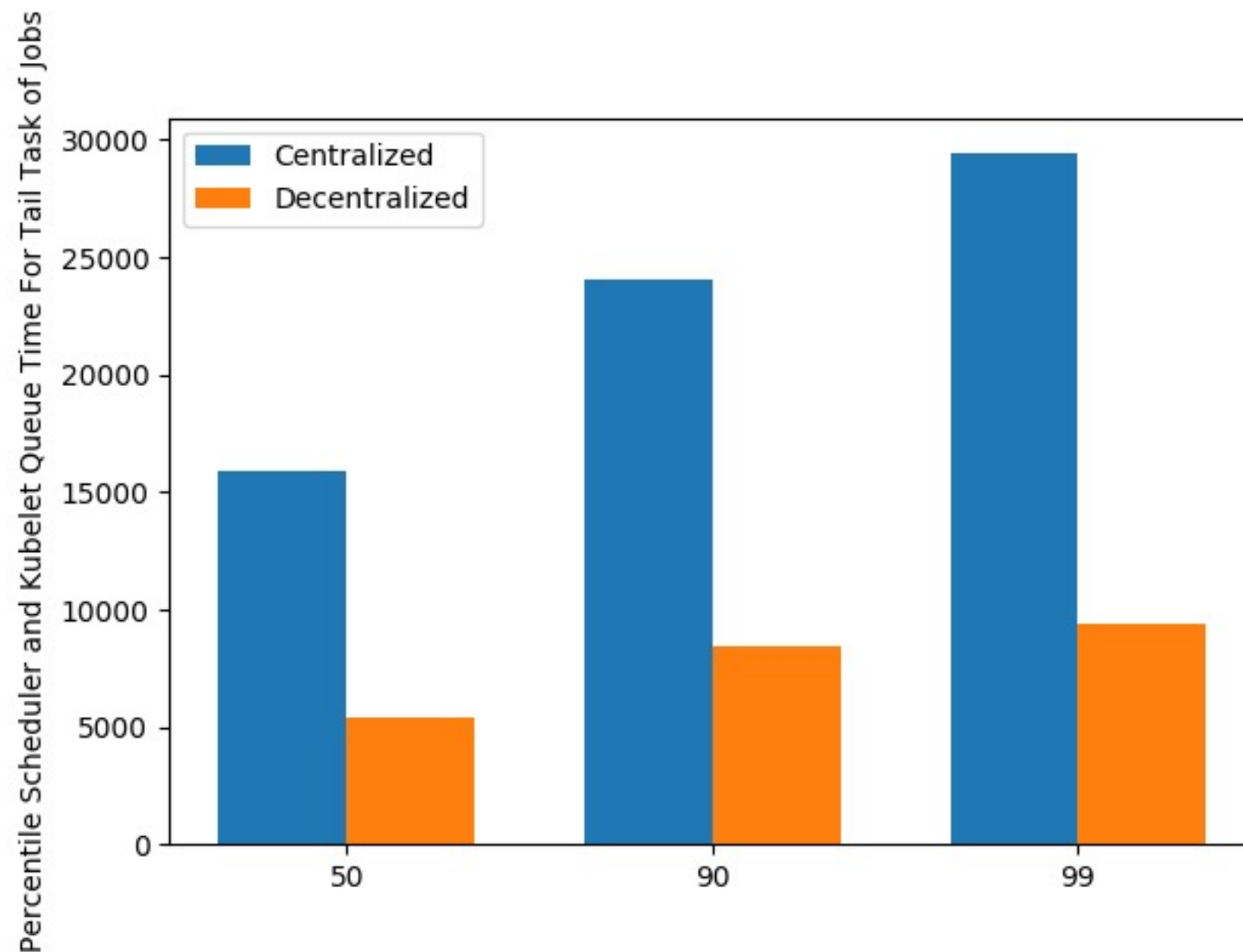


Scalability

Placement Quality



Tail Task Latencies



Conclusion

1. Presented drawbacks of centralized scheduling
 - Bottleneck single incoming job request queue
 - Large tail task latencies
2. Highlight the advantages of job level decentralized scheduling
 - Better overall job completion times
 - Shorter tail task queue times
 - Scalability