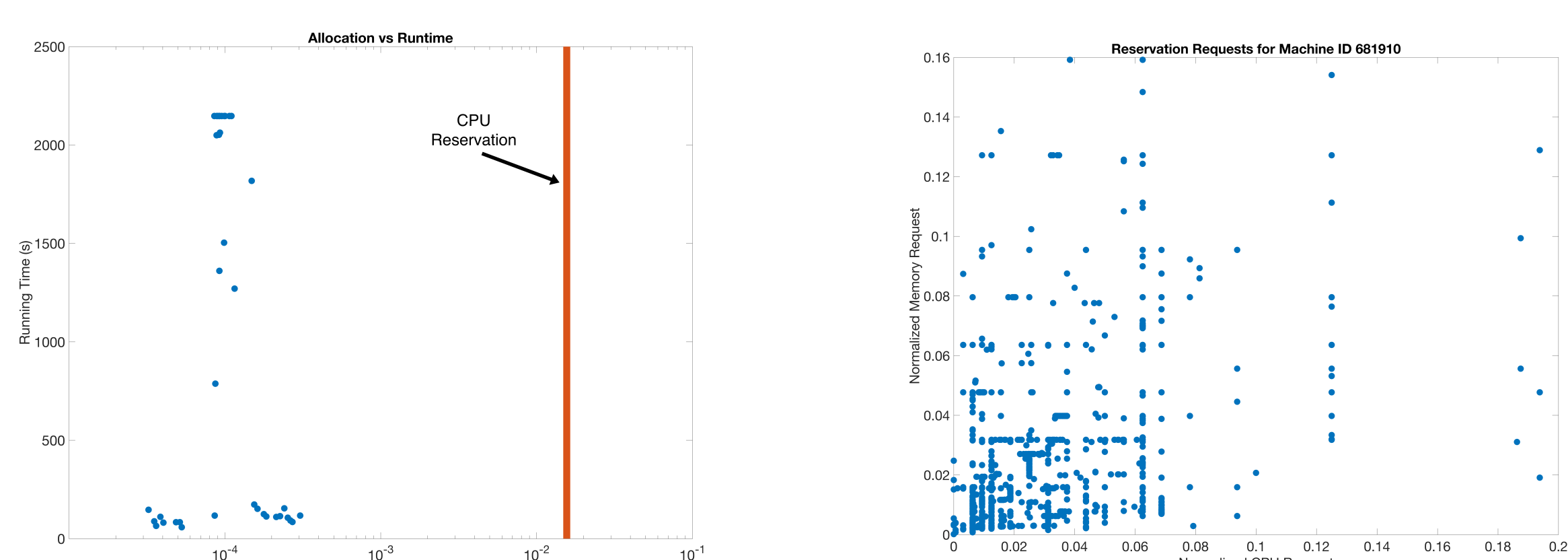


Heterogeneous Resource Allocation in Warehouse Scale Computers

Zach Ankenman, Ben Berg

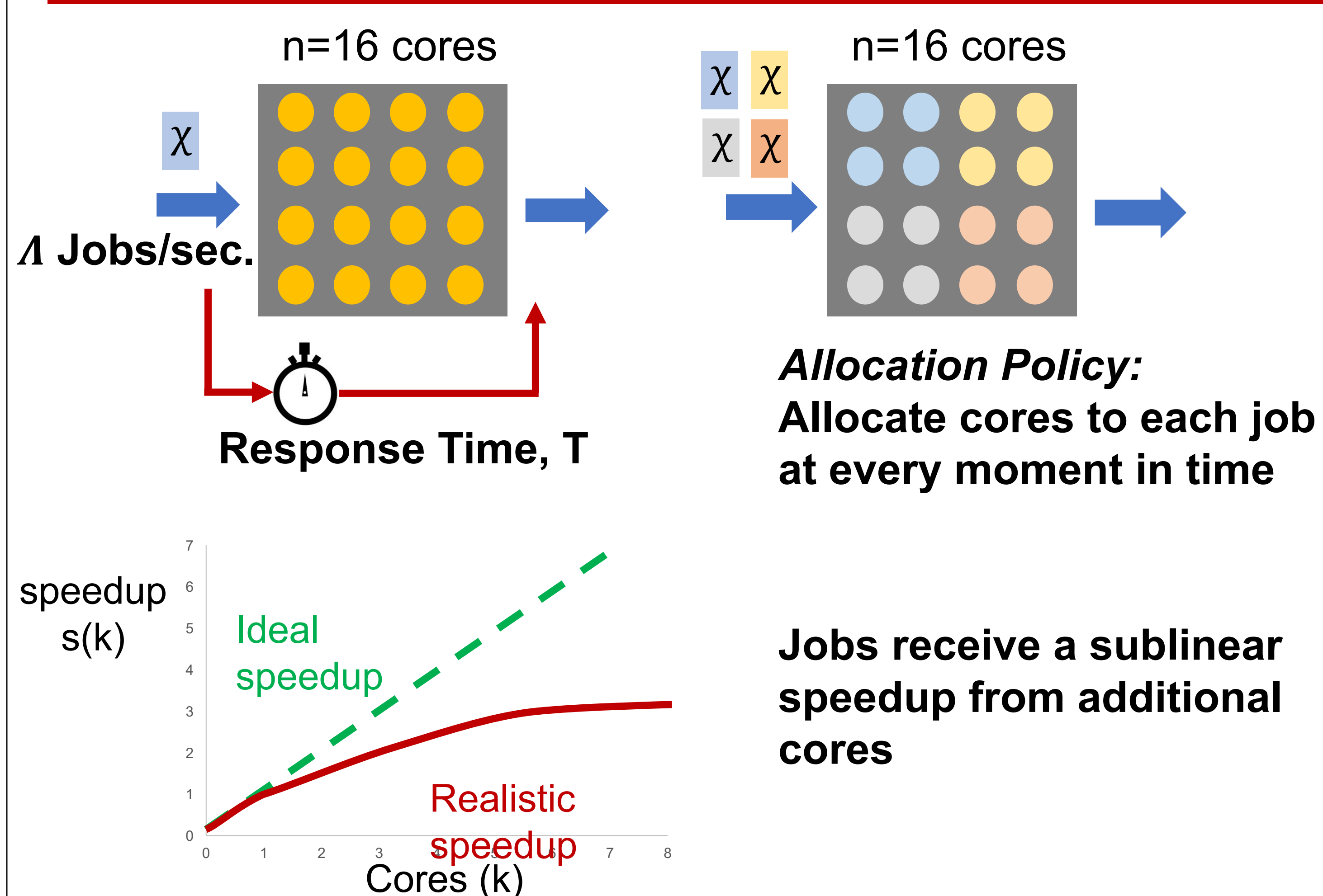
Motivation

- Warehouse scale computers (WSC's) used reservation-based management
- WSC workload is diverse
- Utilization is low, so *latency* is more important than *throughput*
- Current trend is towards *performance-based* management and *colocation*

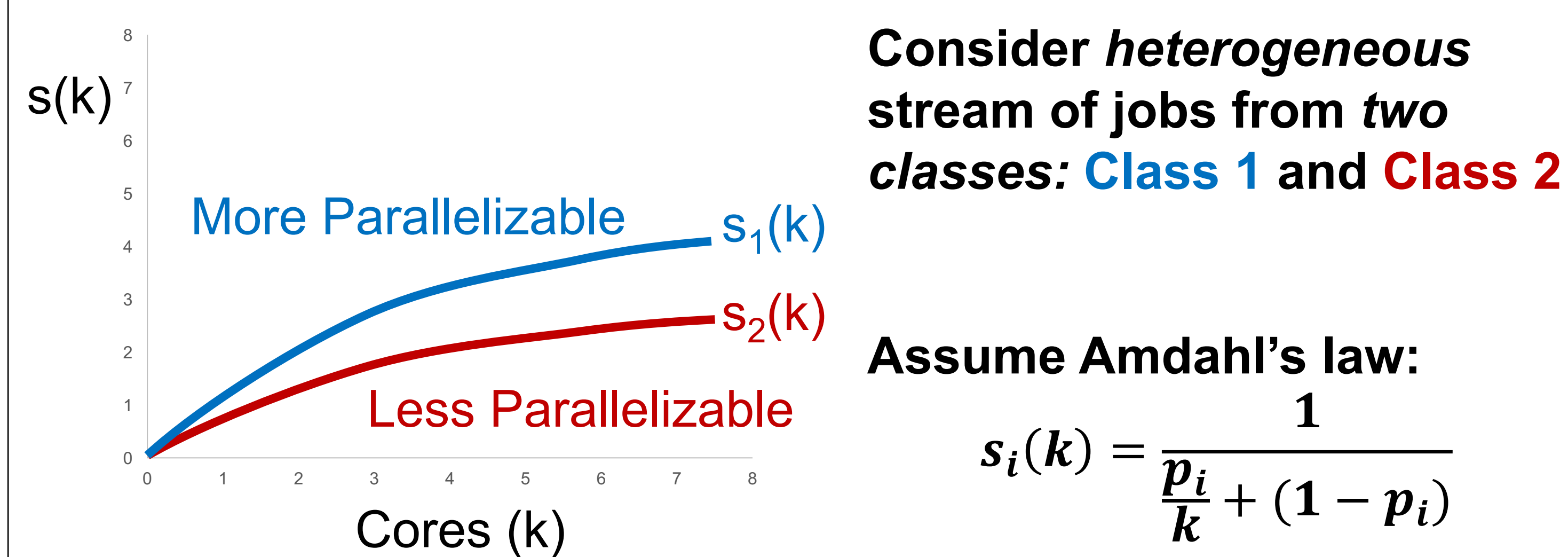


Question: How should we allocate resources to a diverse, latency-critical workload?

Model



Problem

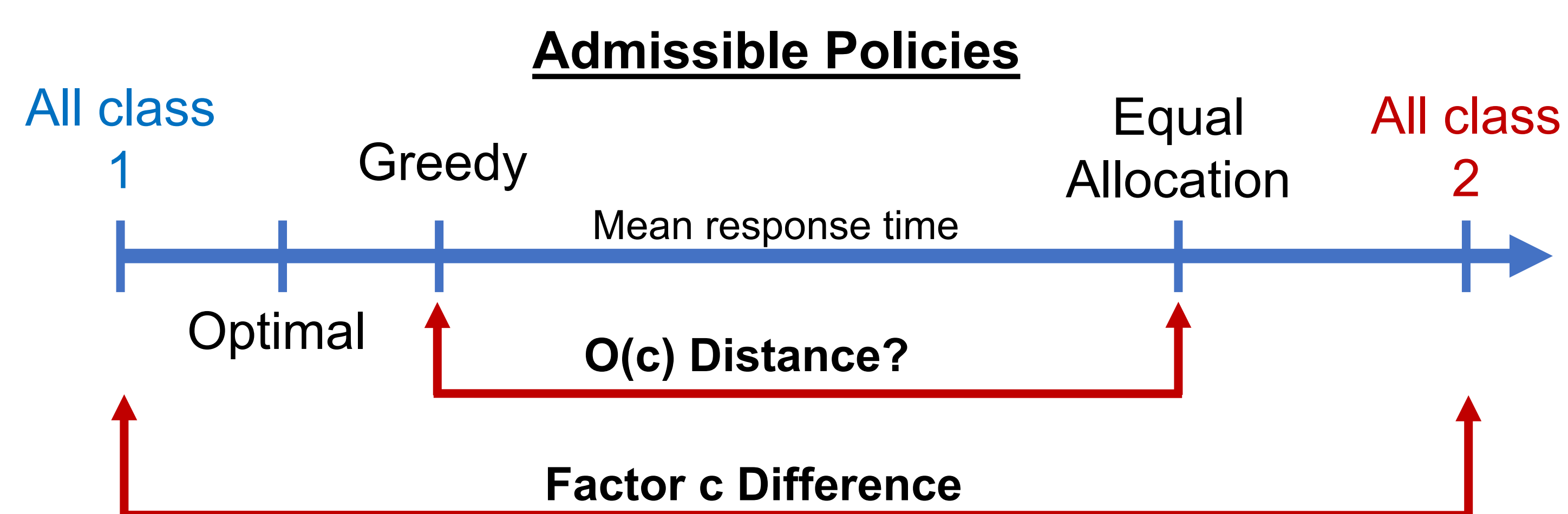


Goal: Describe allocation policy to minimize mean response time

Carnegie Mellon University

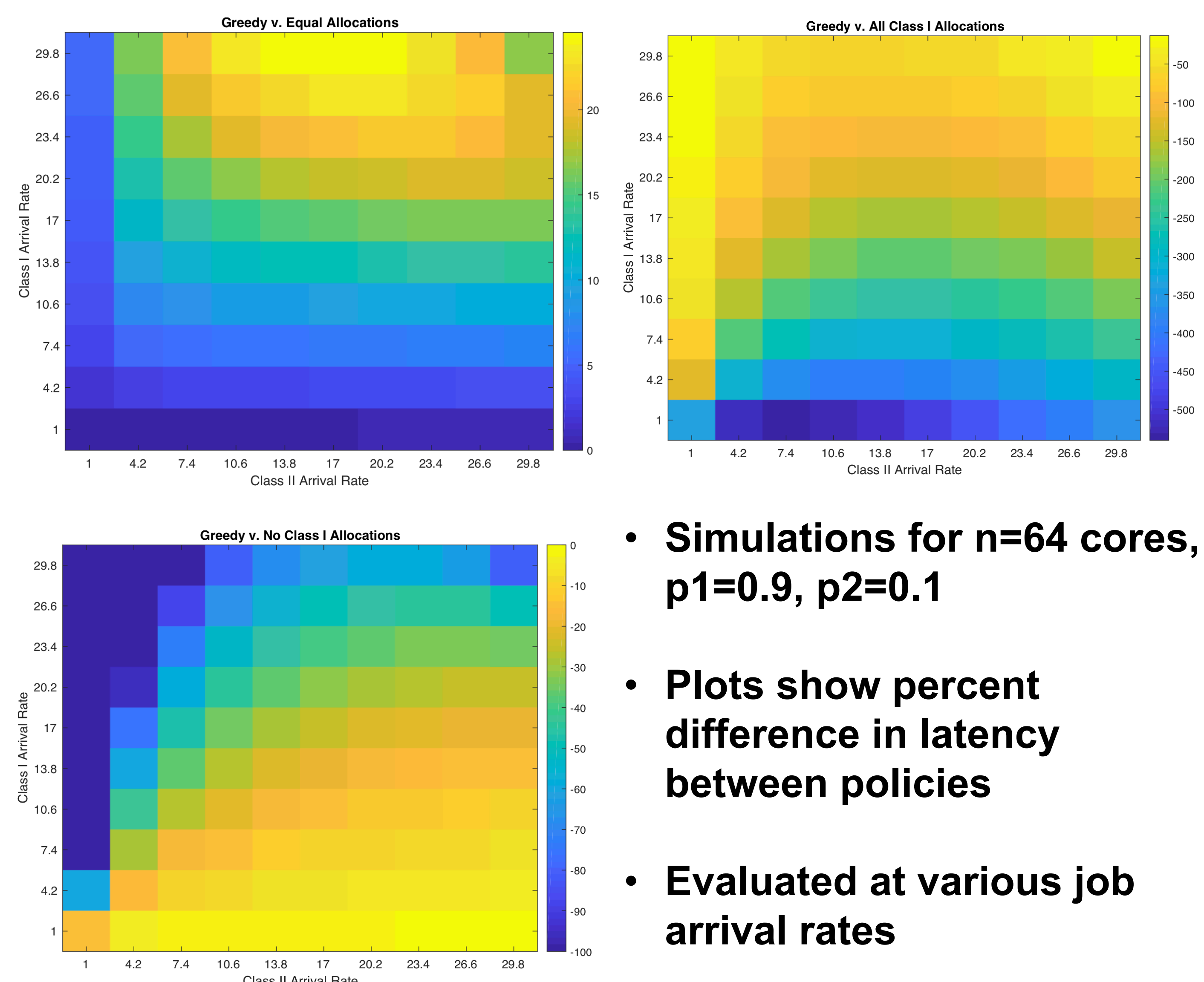
Theoretical Results

- In practice, cluster managers use *Greedy policy*: maximize throughput at all times
 - Provably suboptimal, but how bad is it?
- An *Admissible Policy* is:
 - Better than if *all* jobs were **class 2**
 - Worse than if *all* jobs were **class 1**
 - Greedy is admissible
- Theorem:** If $\frac{s_1(k)}{s_2(k)} \leq c$, any admissible policy becomes a *c*-approximation of Optimal as *n* becomes large



- Is Greedy much better than the Equal Allocation policy?

Simulation Results



Future Work

- Based on above comparisons, try and and prove bounds
- Adapt to handle additional job classes
- Use this information to make colocation decisions
 - What workloads can be collocated efficiently?