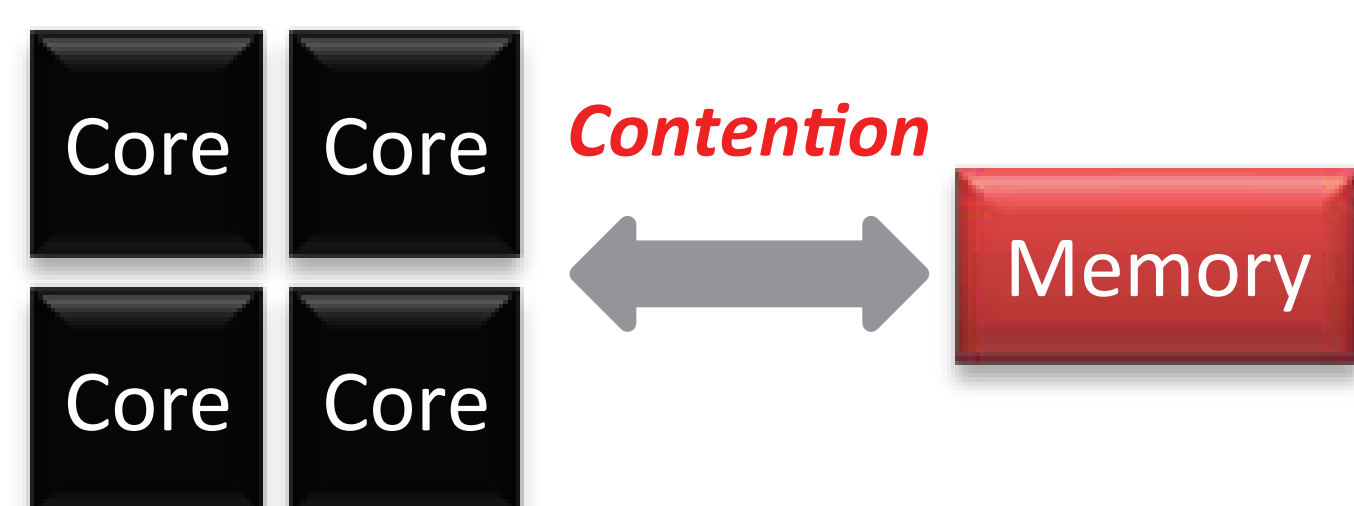


THREAD CLUSTER MEMORY SCHEDULING

Yoongu Kim, Michael Papamichael, Onur Mutlu, Mor Harchol-Balter (CMU)

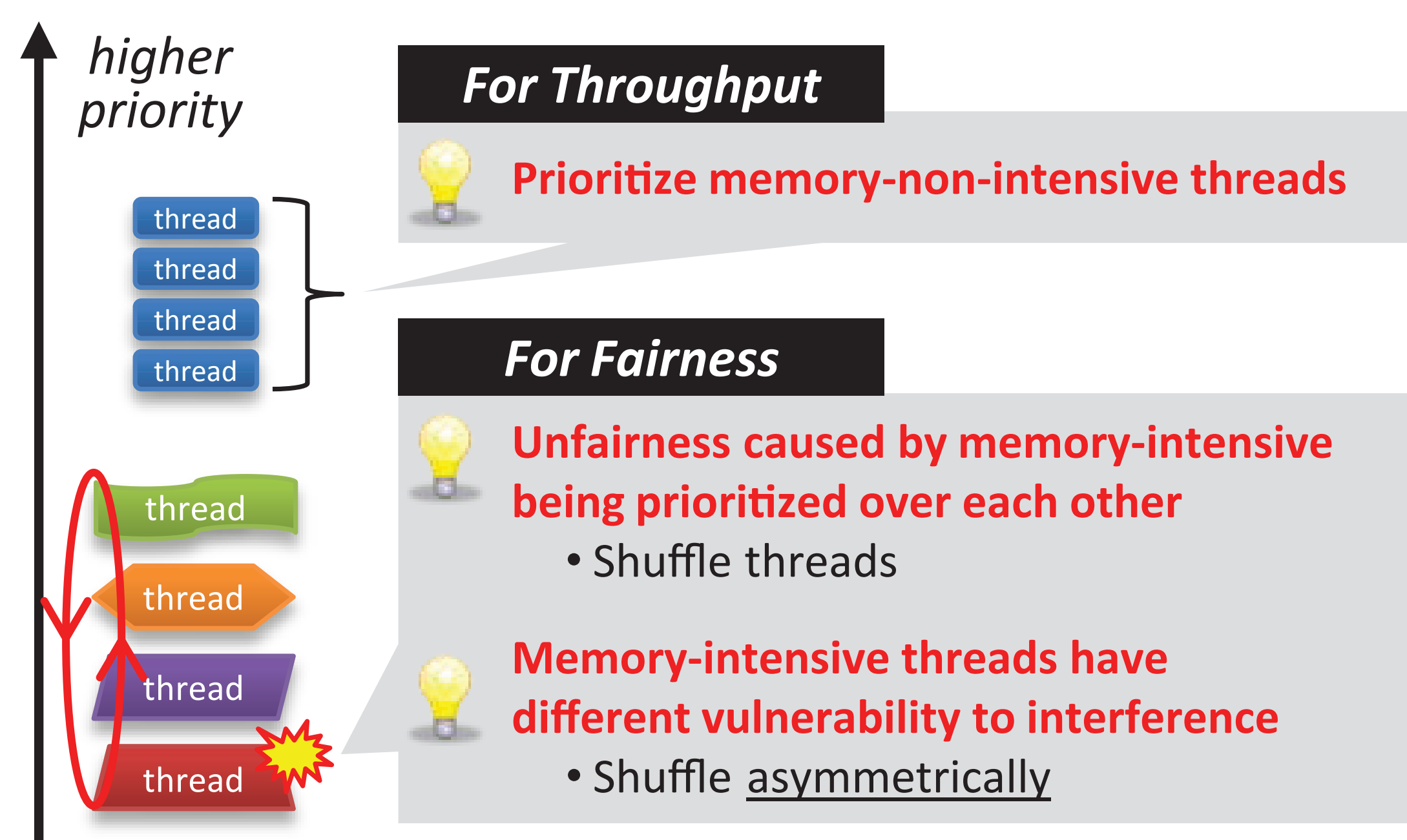
MOTIVATION

- Memory is a key shared resource in CMPs
- Contention for memory access leads to:
 - Degradation in single-thread performance
 - Starvation



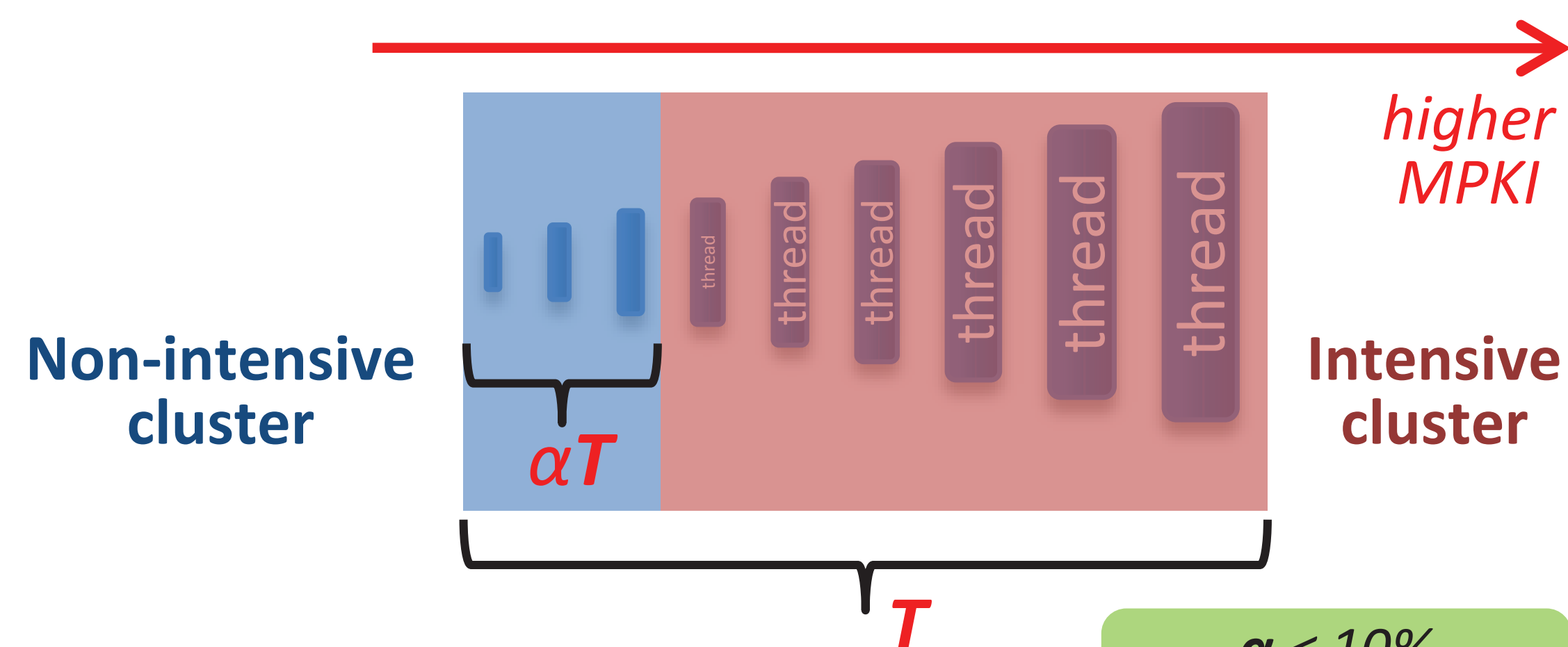
How to achieve both system throughput and fairness?

INSIGHT: BEST OF BOTH WORLDS



HOW TO CLUSTER THREADS?

Step1 Sort threads by MPKI (misses per kiloinstruction)



T = Total memory bandwidth usage

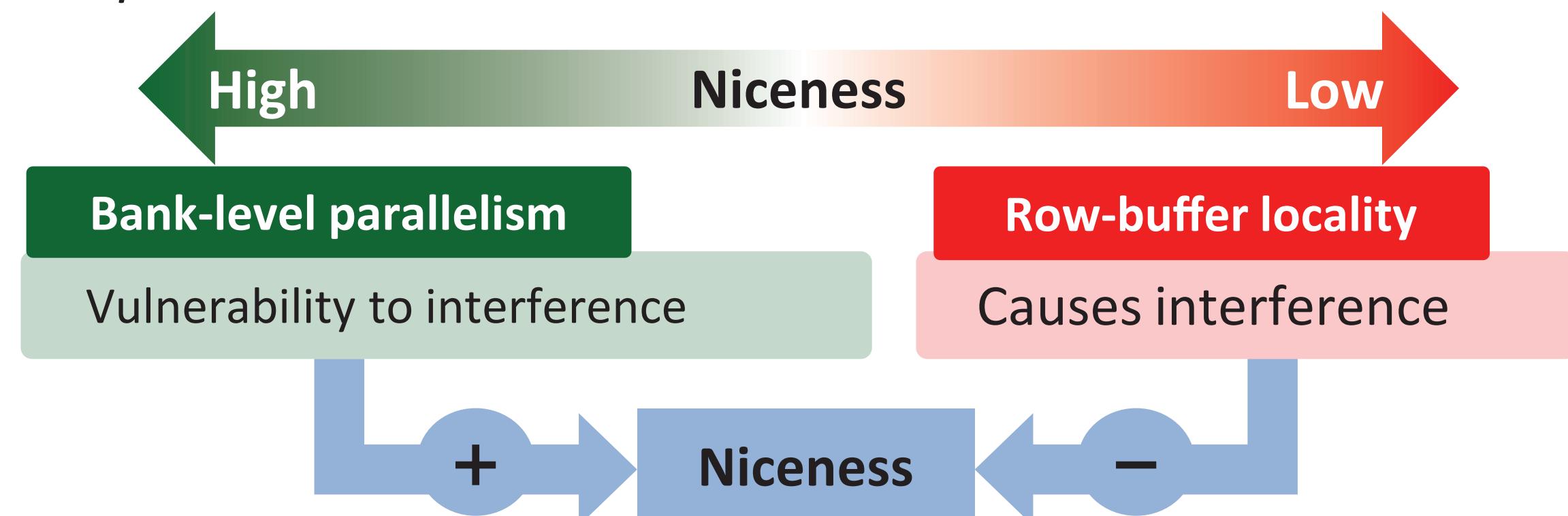
Step2 Memory bandwidth usage αT divides clusters

$\alpha < 10\%$
ClusterThreshold

PRIORITIZATION WITHIN INTENSIVE CLUSTER - II

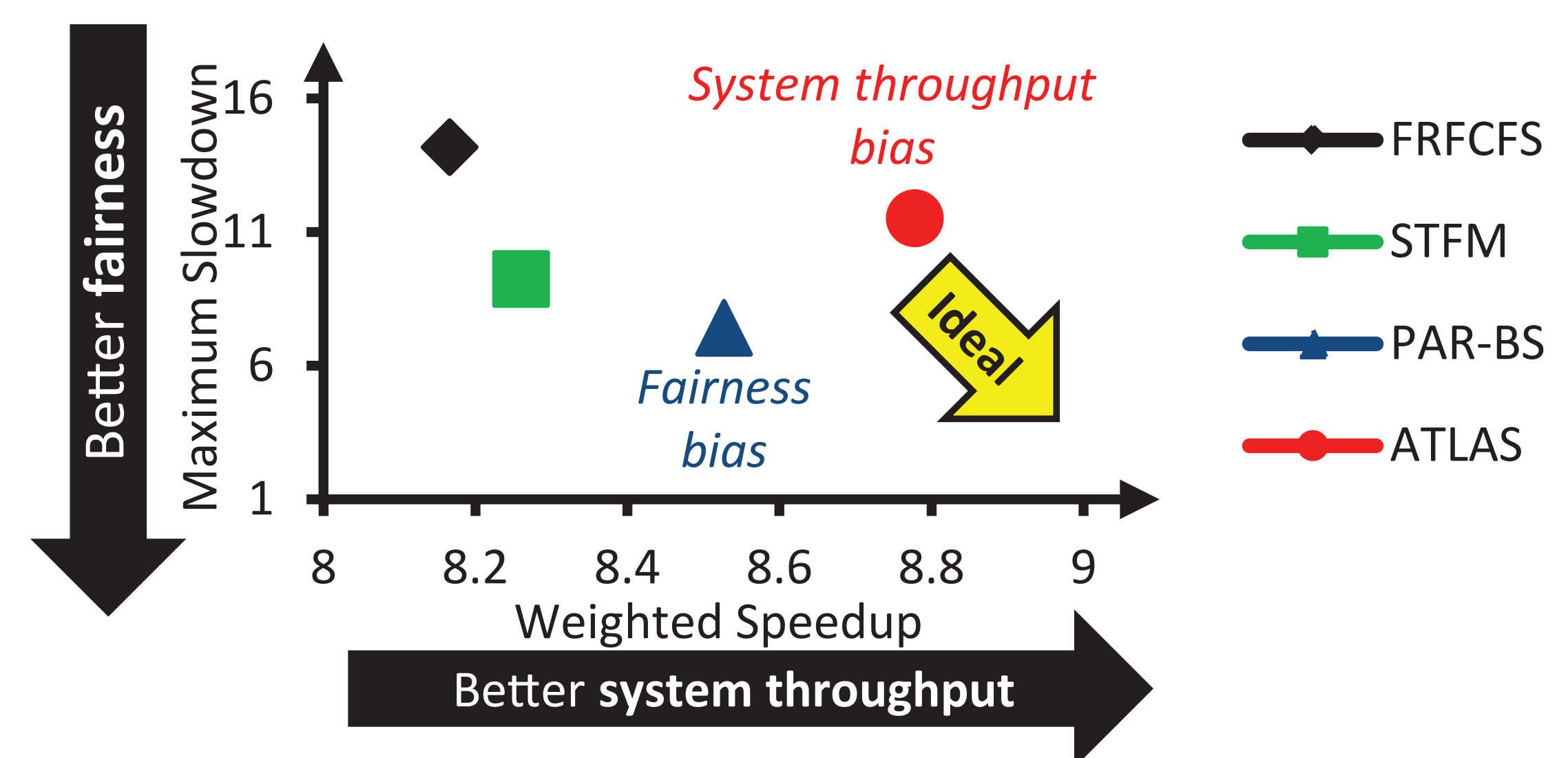
How to quantify difference between threads?

Proposed metric: Niceness



Shuffle priorities in a niceness-aware, asymmetric manner

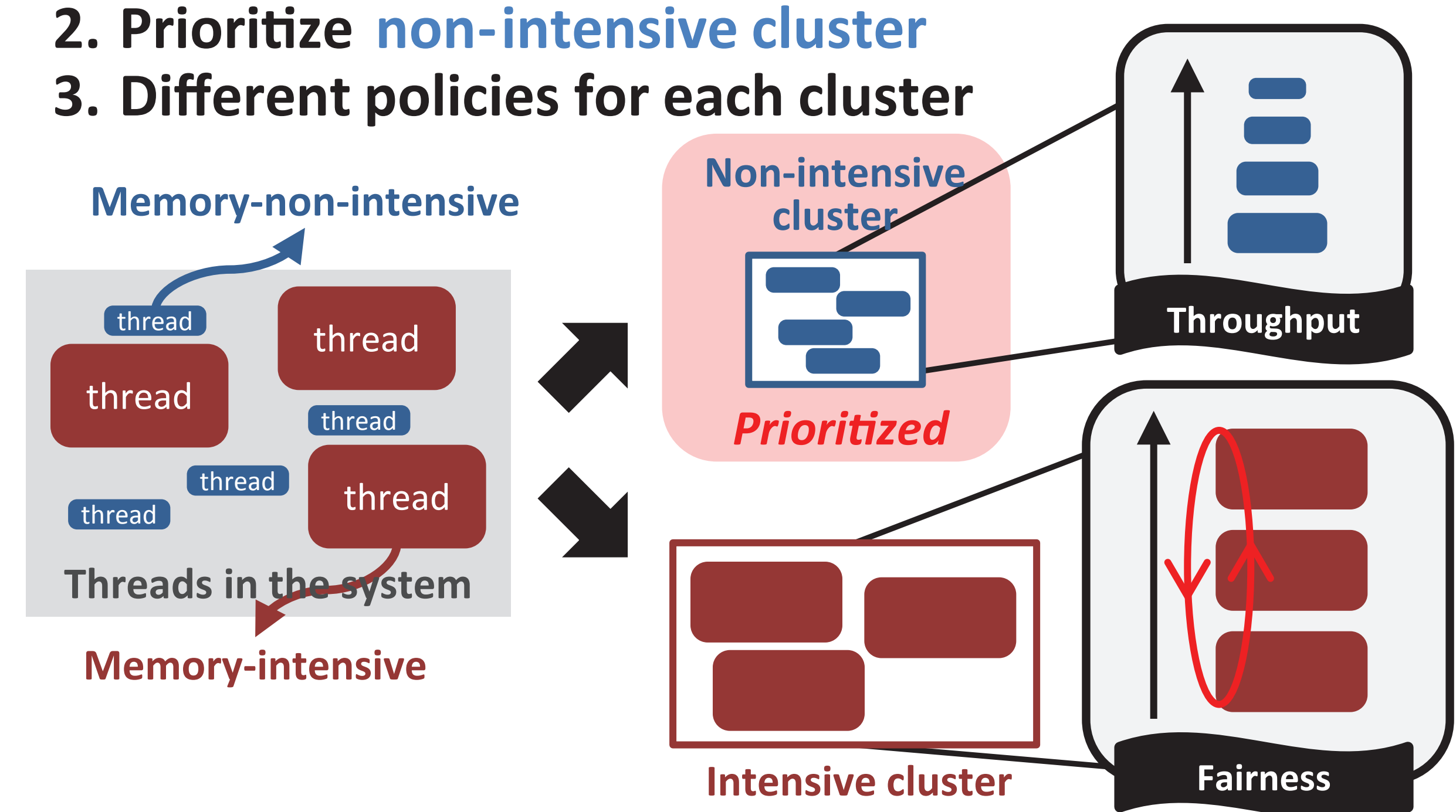
SHORTCOMINGS OF PREVIOUS WORK



Previous approaches are biased: can't achieve both

THREAD CLUSTER MEMORY SCHEDULING

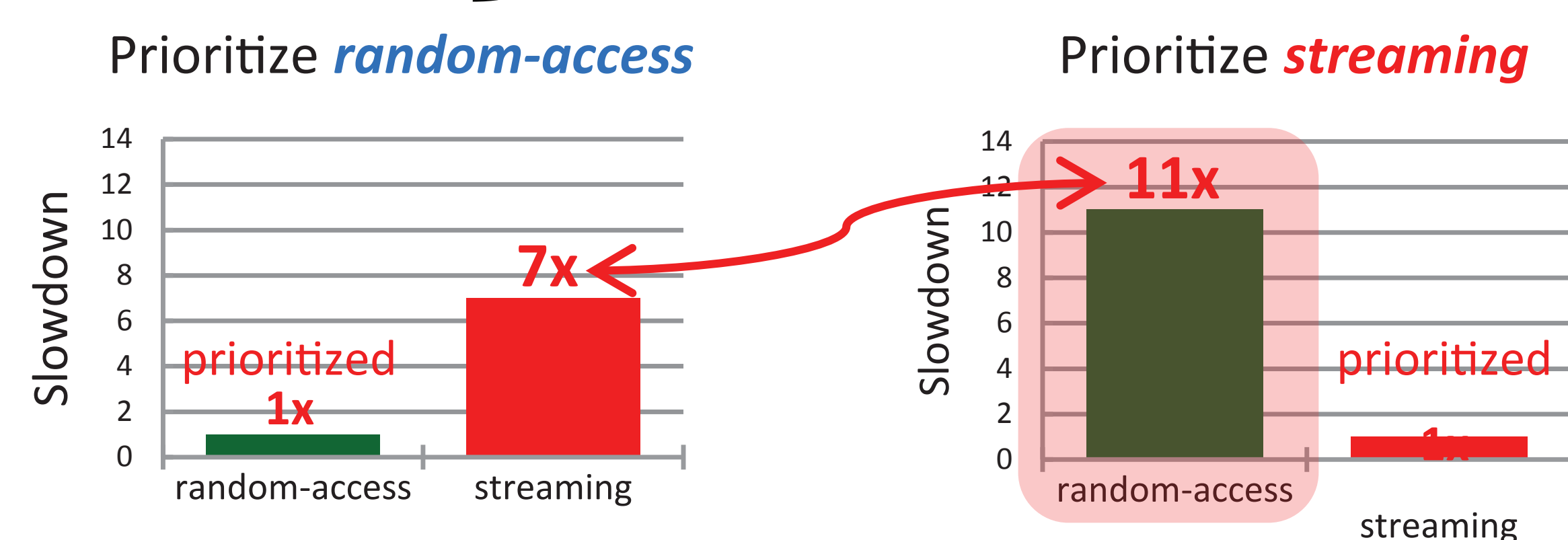
- Group threads into two clusters
- Prioritize non-intensive cluster
- Different policies for each cluster



PRIORITIZATION WITHIN INTENSIVE CLUSTER - I

Case Study: Two intensive threads contending

- random-access
 - streaming
- Which is slowed down more easily?



random-access thread is more easily slowed down

- Vulnerable to interference
- Causes less interference to other threads

RESULTS

