

DON'T CRY OVER SPILLED RECORDS

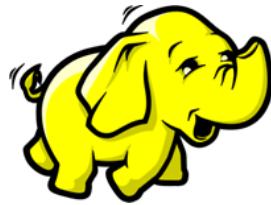
Memory elasticity of data-parallel applications and its application to cluster scheduling

Călin Iorgulescu (EPFL), Florin Dinu (EPFL), Aunn Raza (NUST Pakistan),
Wajih Ul Hassan (UIUC), Willy Zwaenepoel (EPFL)

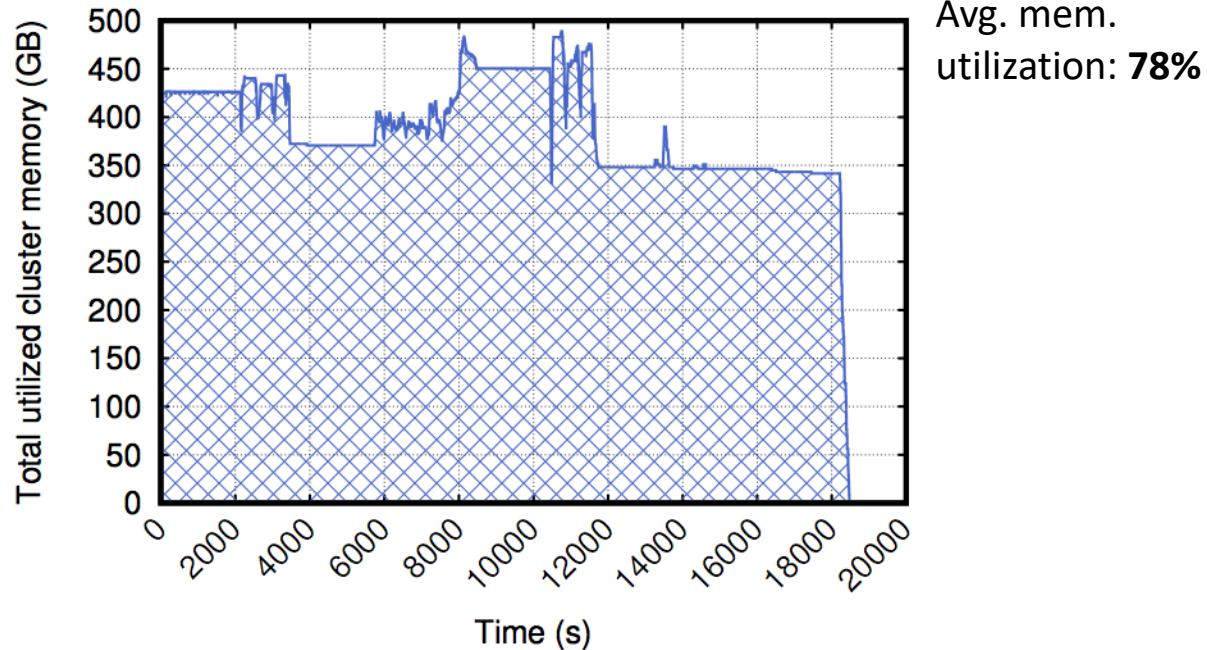


Cluster operators care about resource utilization

- ✓ Best bang for your buck!
- ✓ Maximize performance of data-parallel applications
- Idea: Efficient resource utilization through under-provisioning

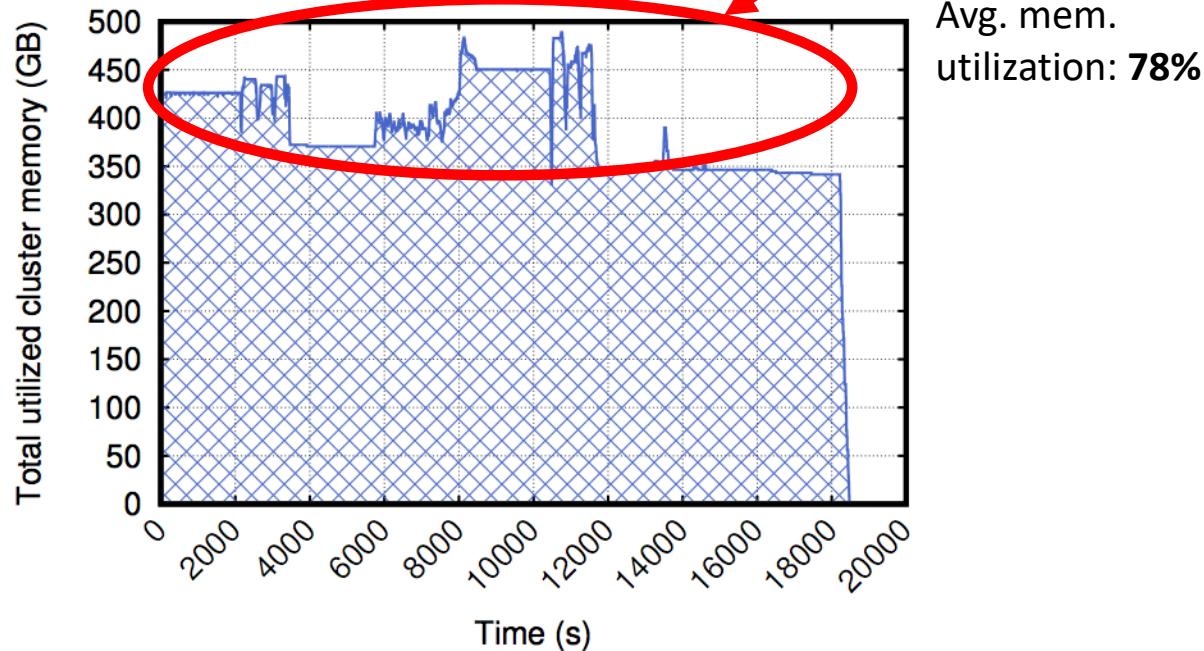


Cluster memory is under-utilized!



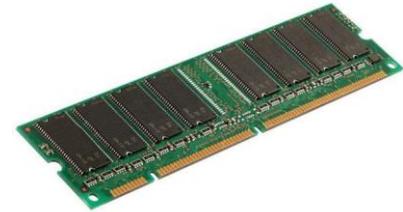
Cluster memory is under-utilized!

Leverage this idle memory!



Impact of memory constraining applications

- Conventional wisdom: do not touch memory!
- Risks:
 - crashes
 - severe performance degradation (e.g., thrashing)



Can we safely, deterministically, and with modest impact constrain memory?

Context: batch jobs and their memory usage

Context: batch jobs and their memory usage

Input

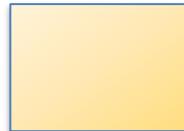


Context: batch jobs and their memory usage

Input



**Ideal
memory**



Context: batch jobs and their memory usage

Input

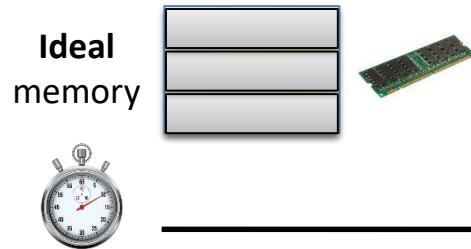


**Ideal
memory**

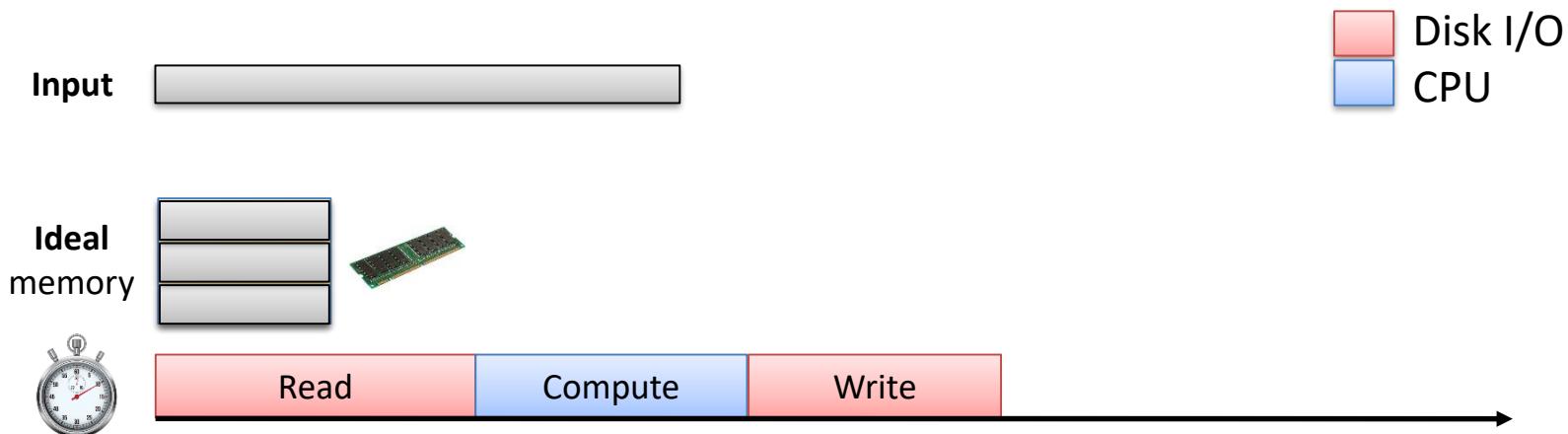


Context: batch jobs and their memory usage

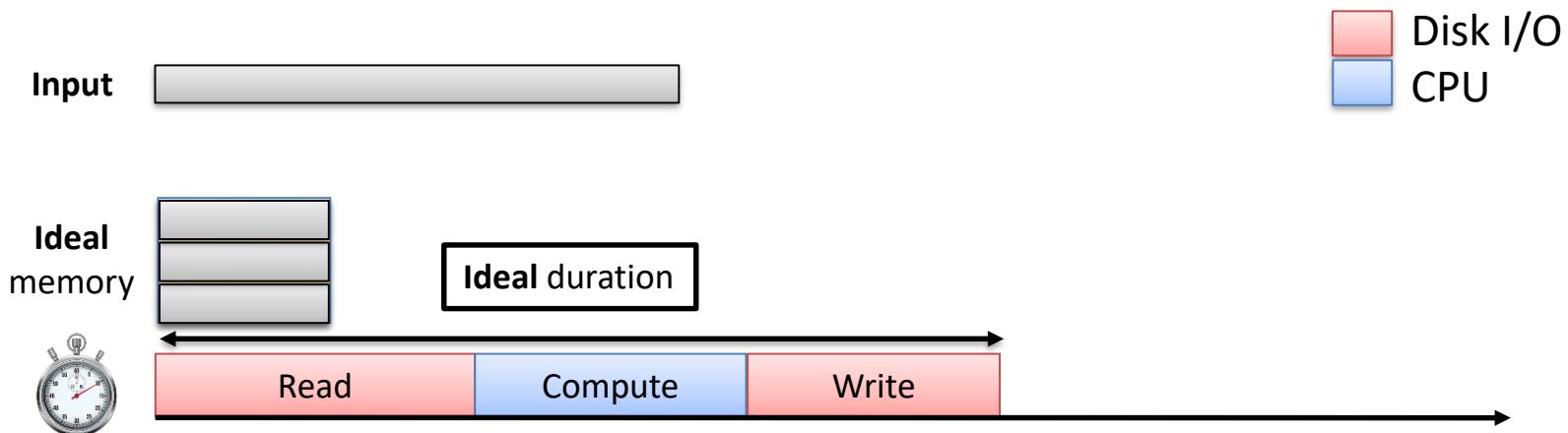
Input 



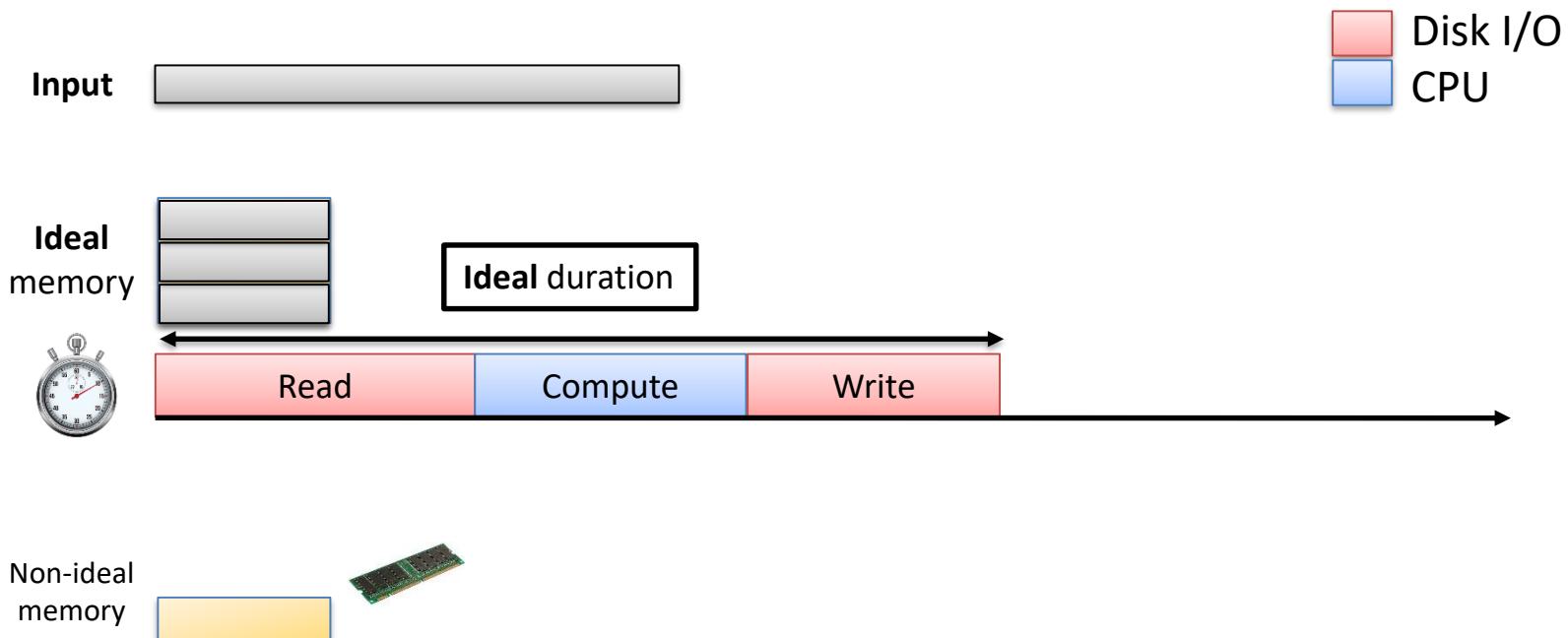
Context: batch jobs and their memory usage



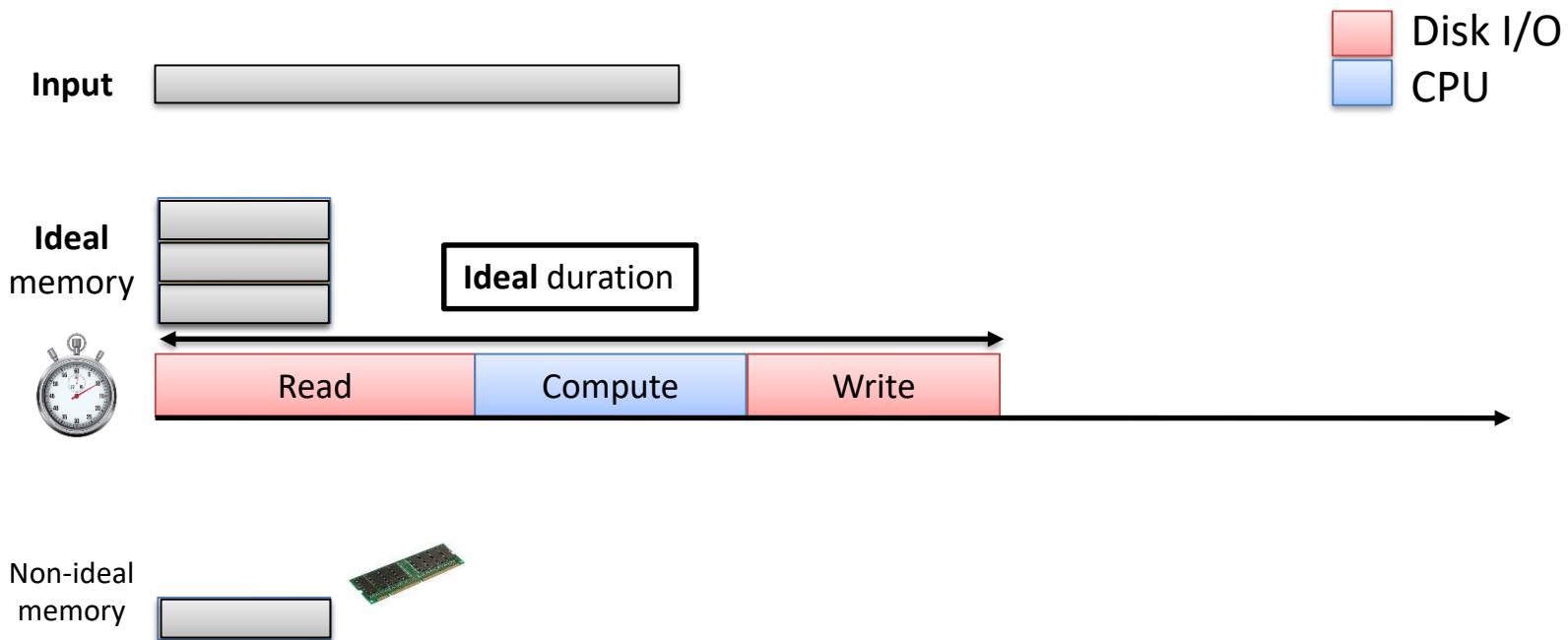
Context: batch jobs and their memory usage



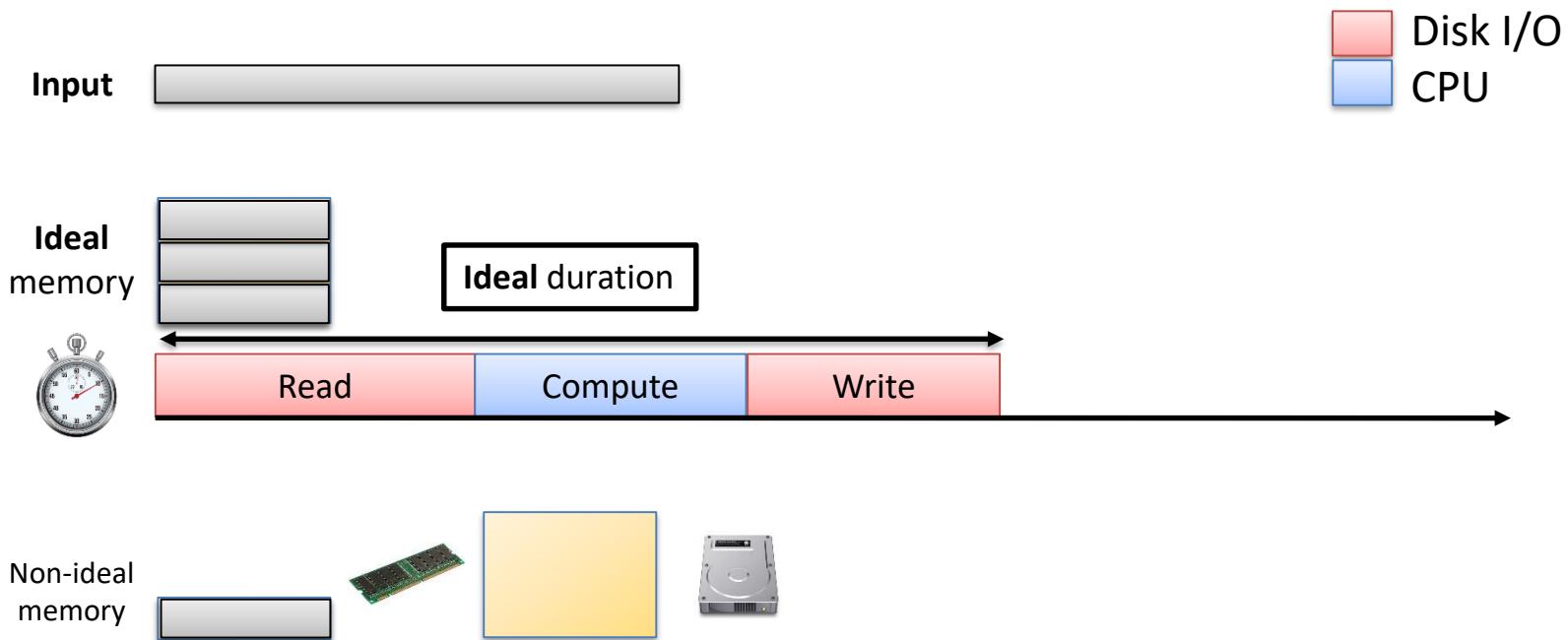
Context: batch jobs and their memory usage



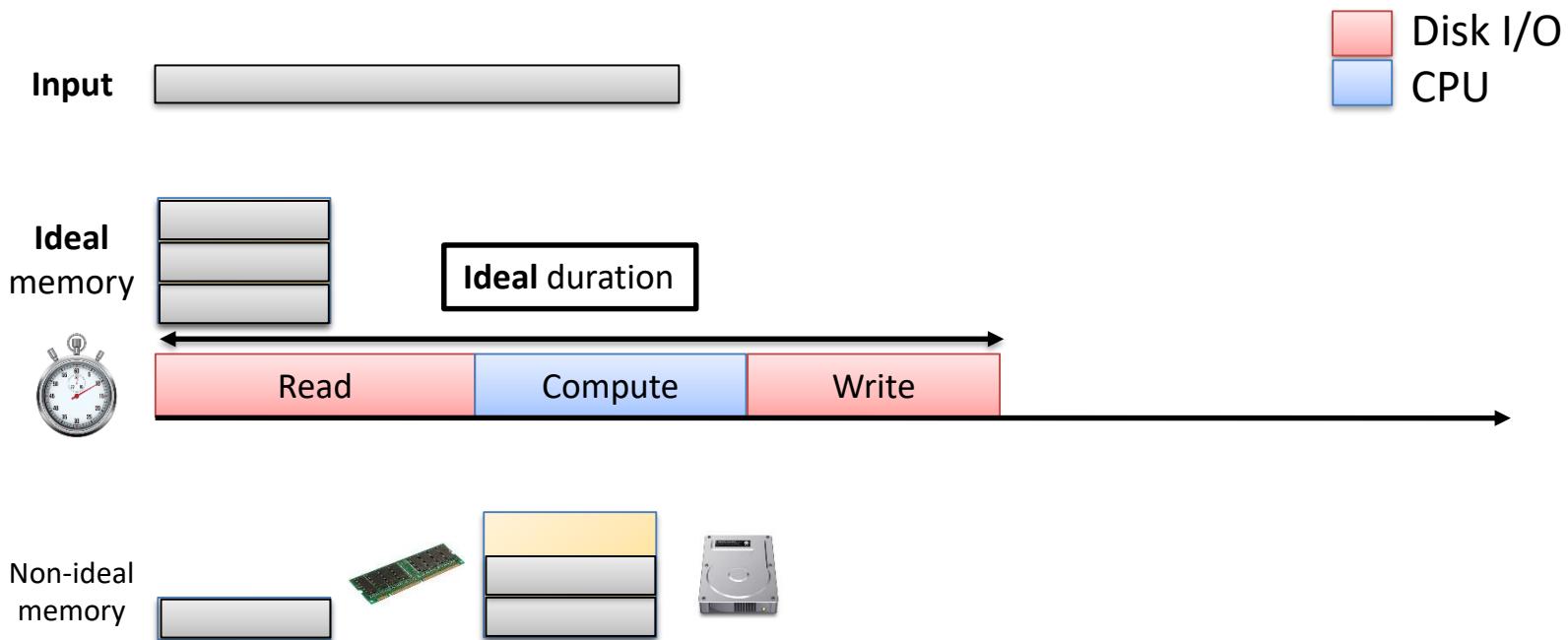
Context: batch jobs and their memory usage



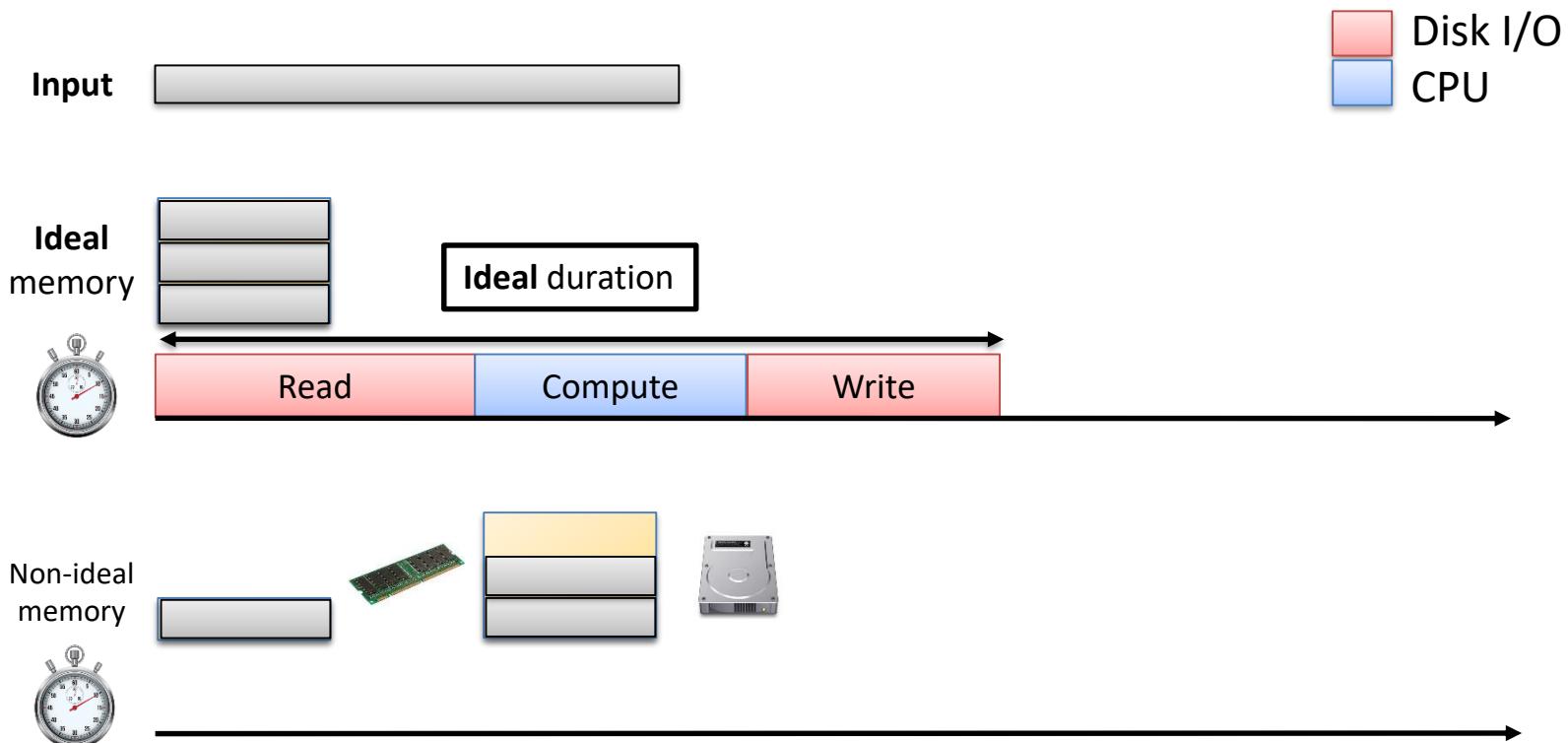
Context: batch jobs and their memory usage



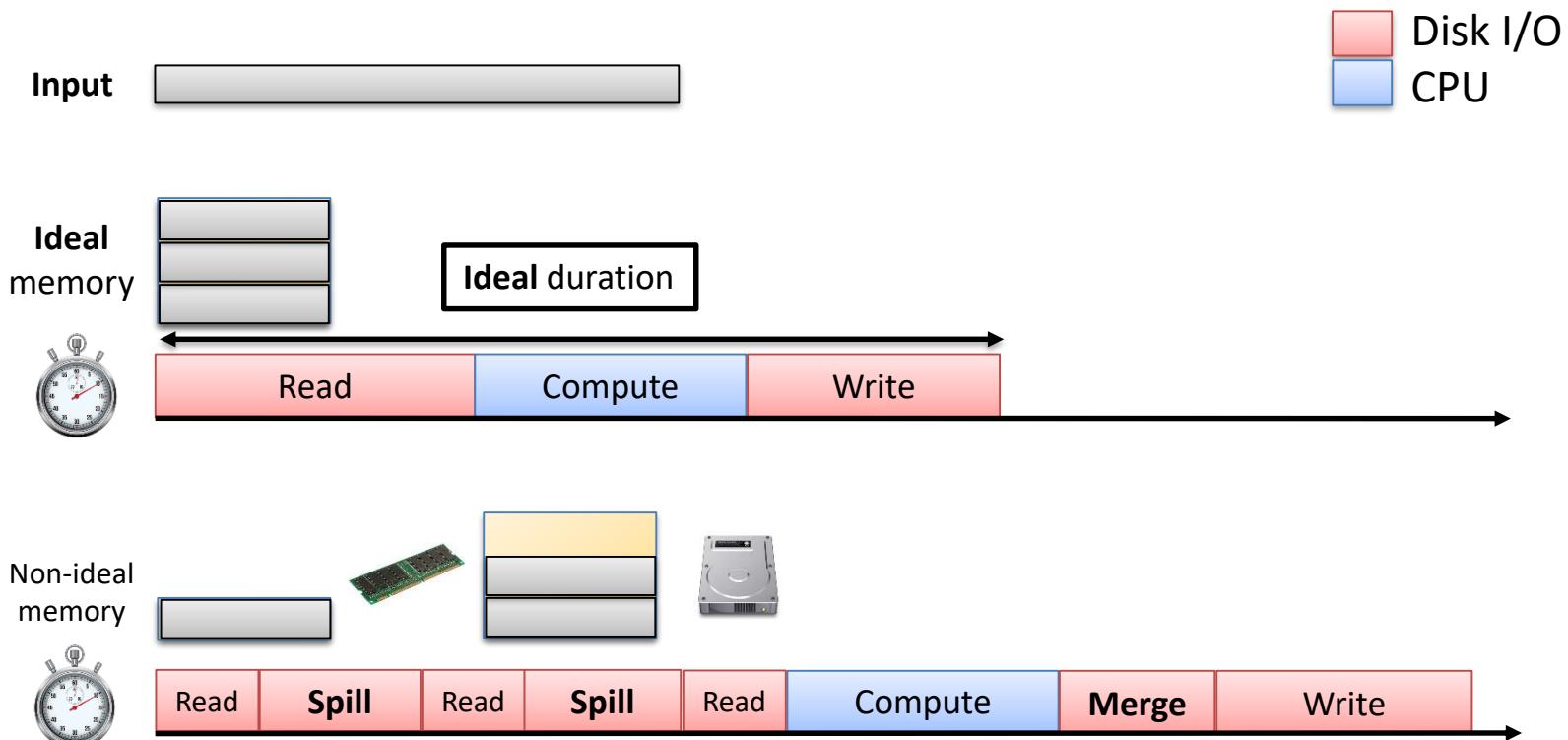
Context: batch jobs and their memory usage



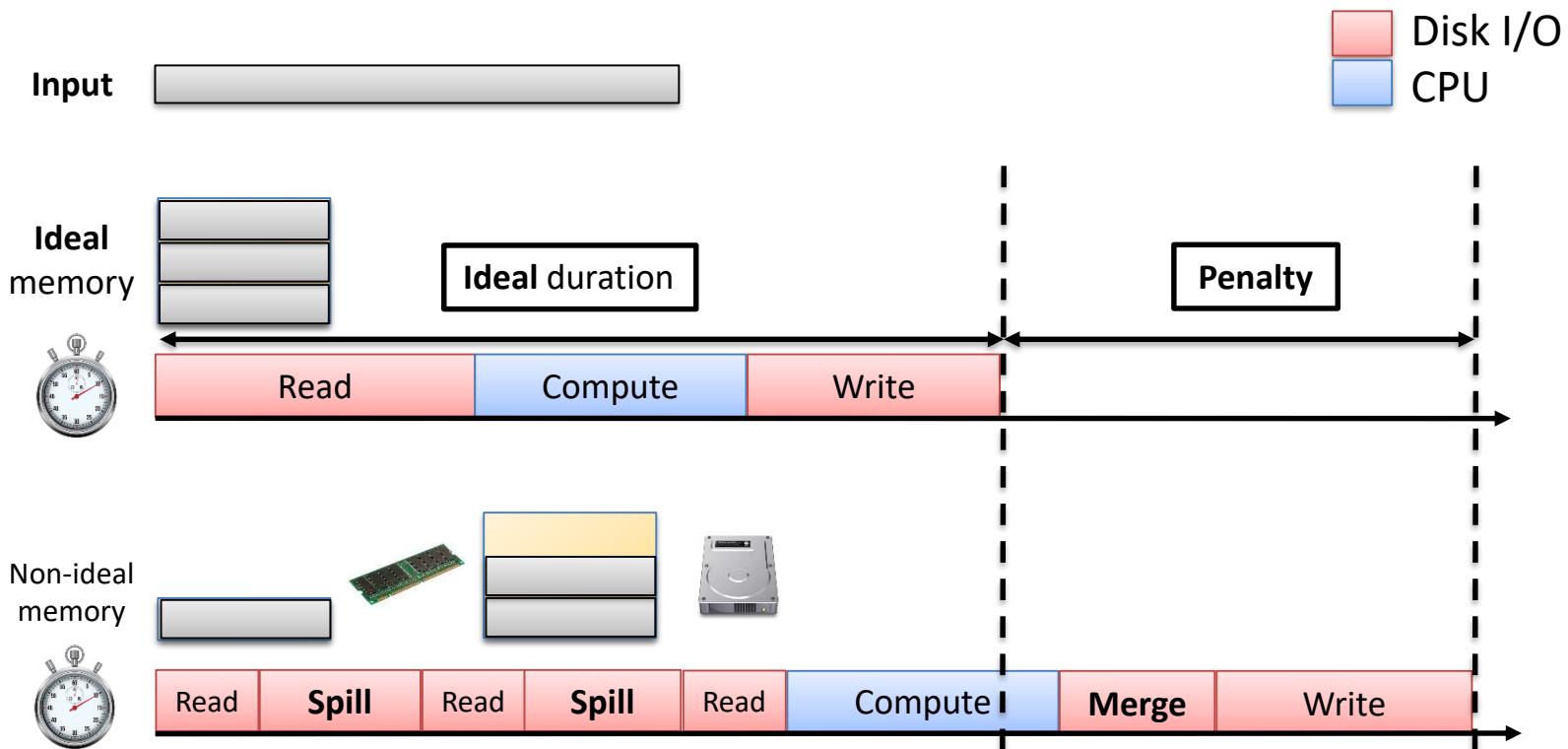
Context: batch jobs and their memory usage



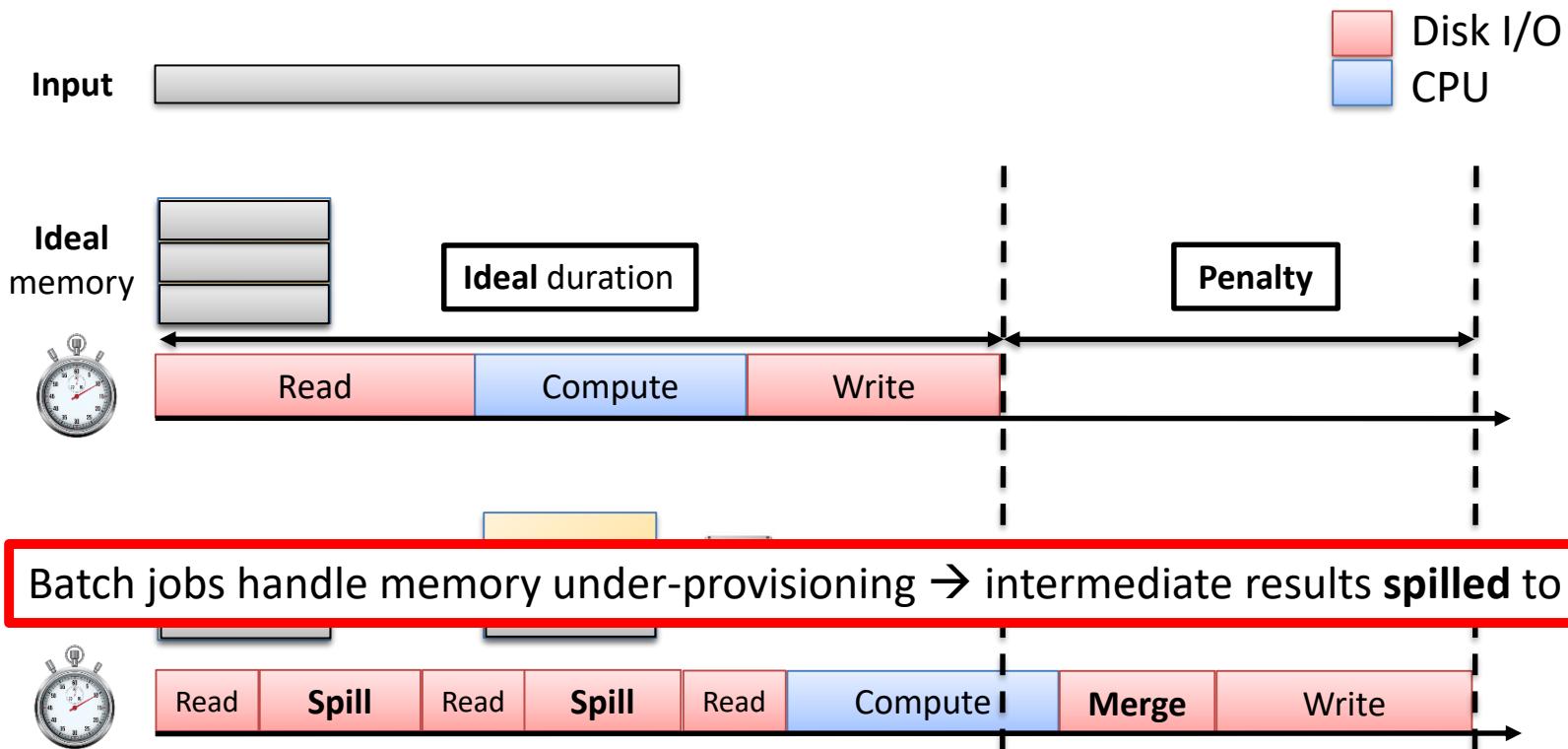
Context: batch jobs and their memory usage



Context: batch jobs and their memory usage

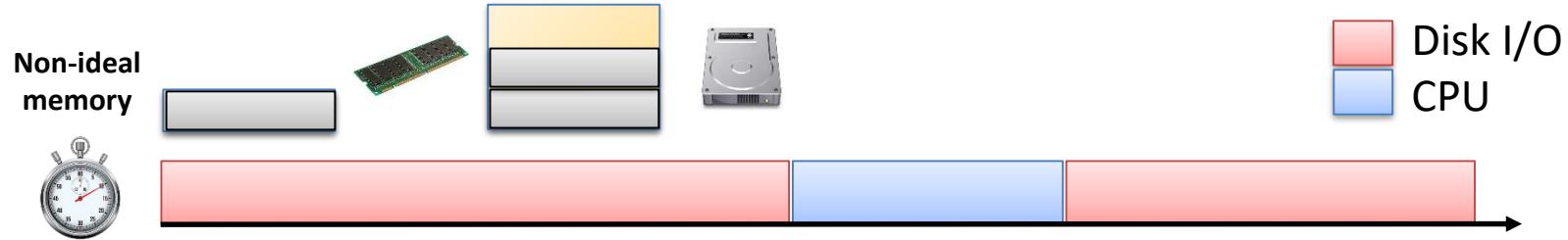


Context: batch jobs and their memory usage

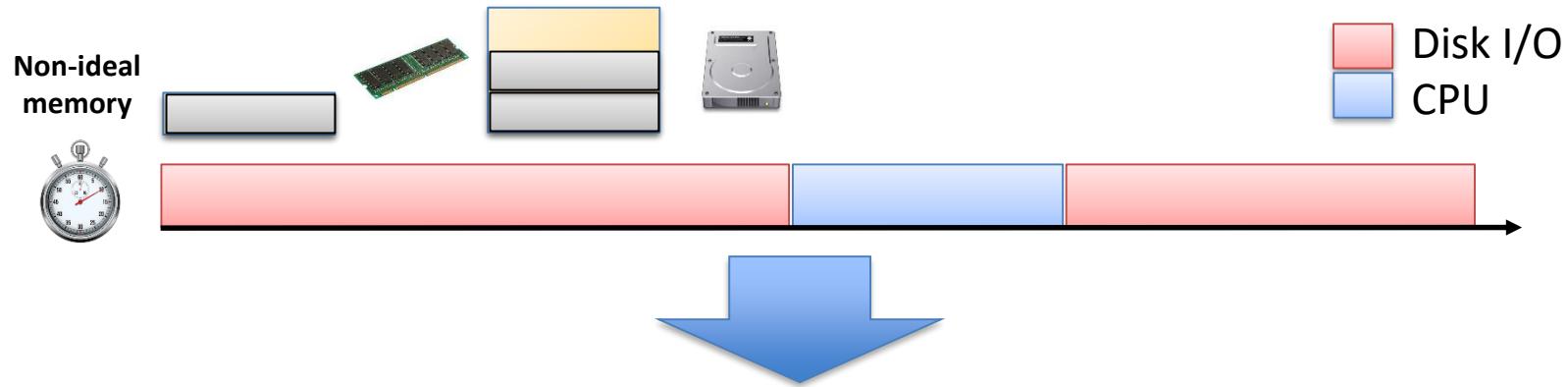


MEMORY ELASTICITY

What is Memory Elasticity?

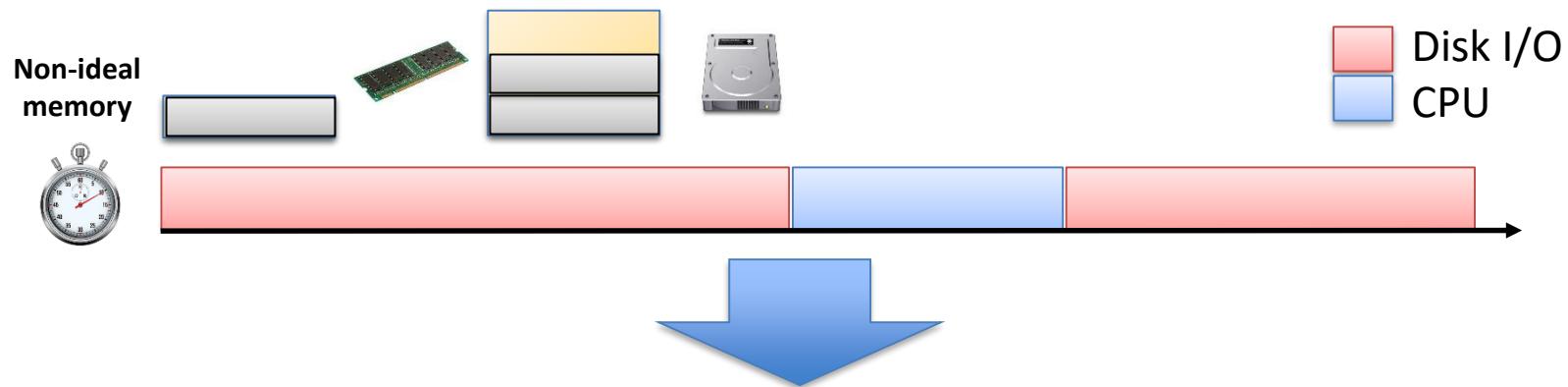


What is Memory Elasticity?



✓ Safely constrain memory

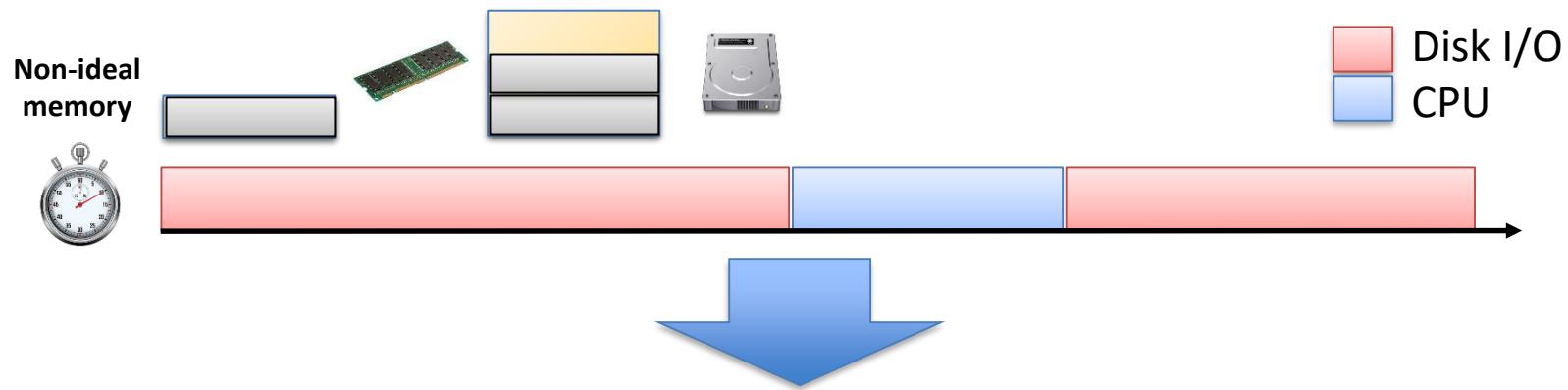
What is Memory Elasticity?



✓ Safely constrain memory

✓ Moderate penalties

What is Memory Elasticity?

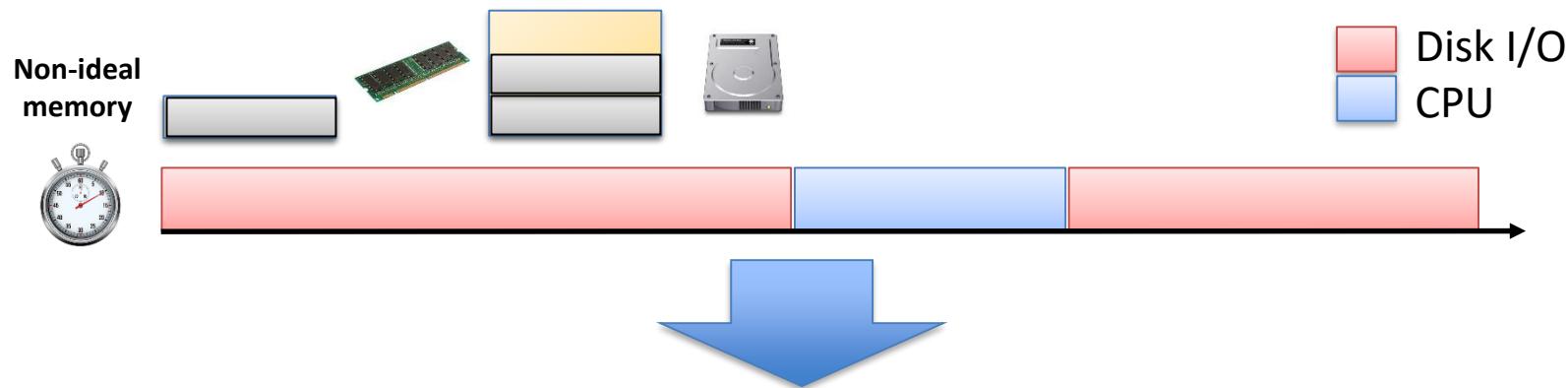


✓ Safely constrain memory

✓ Moderate penalties

✓ Highly predictable

What is Memory Elasticity?



✓ **Safely** constrain memory

✓ **Moderate penalties**

✓ **Highly predictable**

✓ **Ubiquitous** for most
data-parallel apps

An empirical study of Memory Elasticity

- Analysis of 18 jobs across 8 different applications
- Constrain tasks' memory → measure **penalty**
- Bypass disk buffer cache (to not mask impact of spilling to disk)

Questions about Memory Elasticity

Questions about Memory Elasticity

- Are the penalties large?

Questions about Memory Elasticity

- Are the penalties large?
- Do penalties vary considerably w.r.t given memory?

Questions about Memory Elasticity

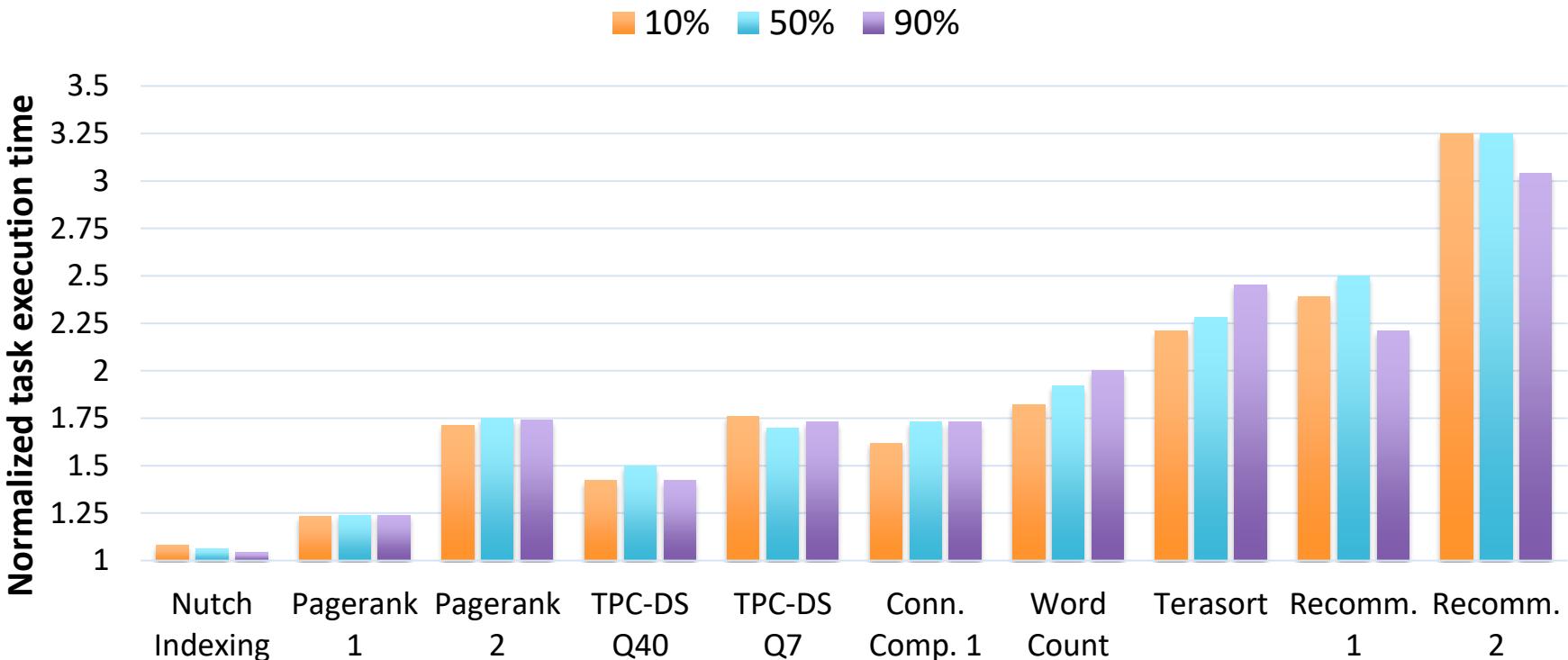
- Are the penalties large?
- Do penalties vary considerably w.r.t given memory?
- Does the additional I/O cause disk contention?

Questions about Memory Elasticity

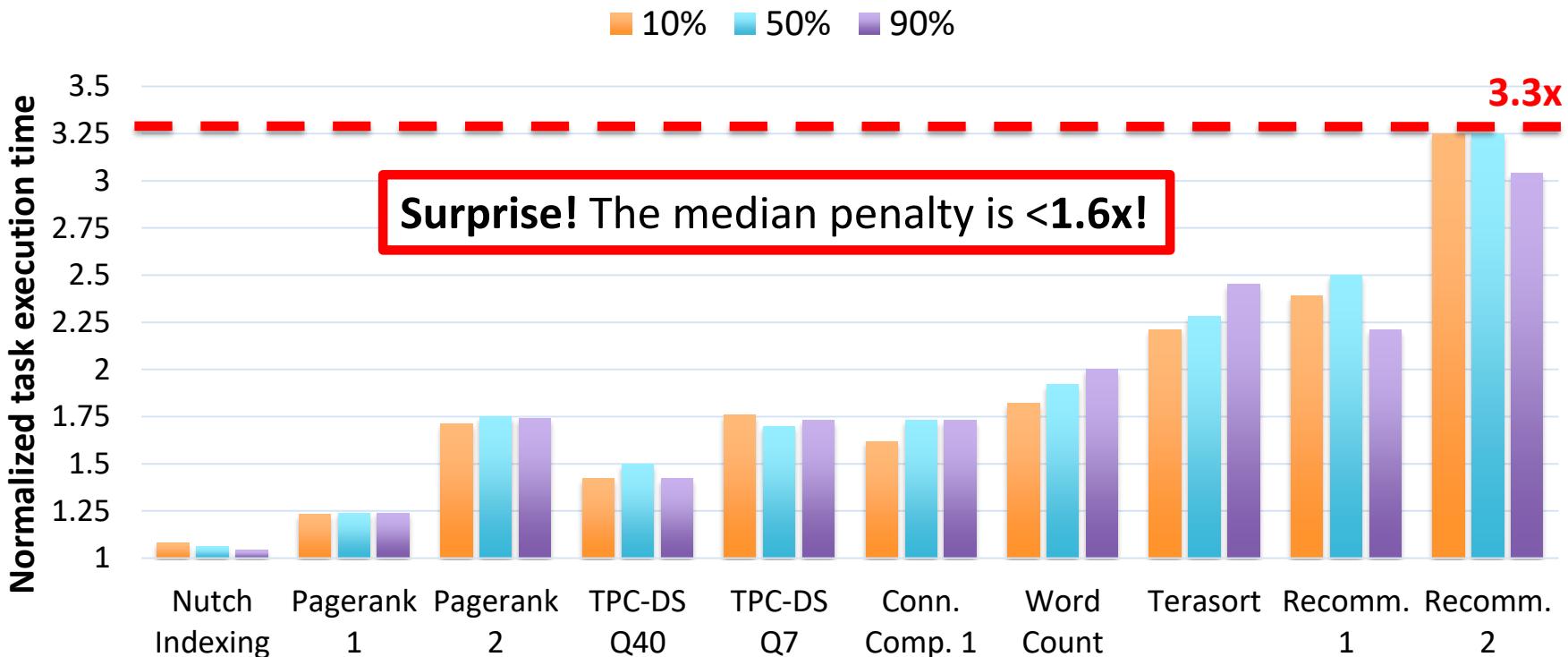
- Are the penalties large?
- Do penalties vary considerably w.r.t given memory?
- Does the additional I/O cause disk contention?

NOT
SO
MUCH!

Elasticity of Hadoop workloads: Reducers



Elasticity of Hadoop workloads: Reducers



Why are the penalties so modest?

Data buffer – most of app. memory

- Memory pressure absorbed by data buffer

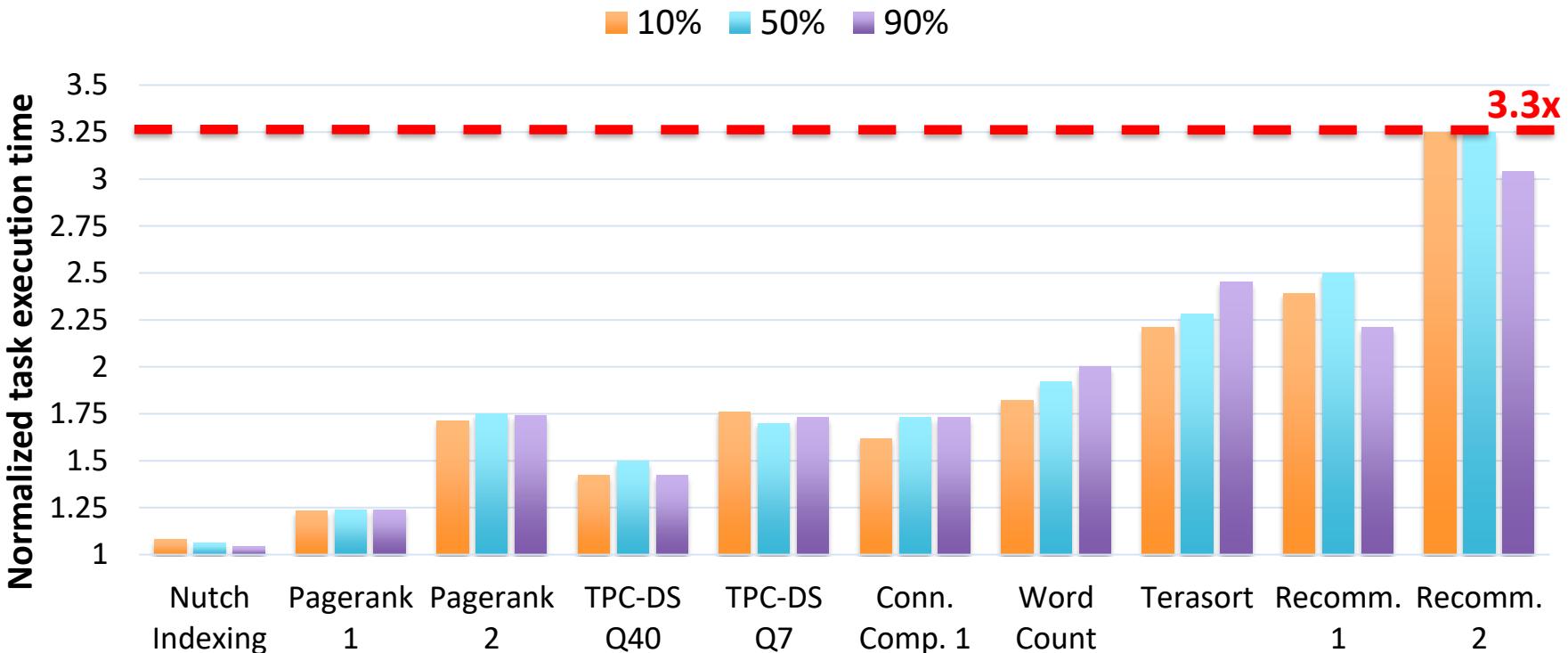
Sequential disk access

- Spilling records to disk is faster than OS paging

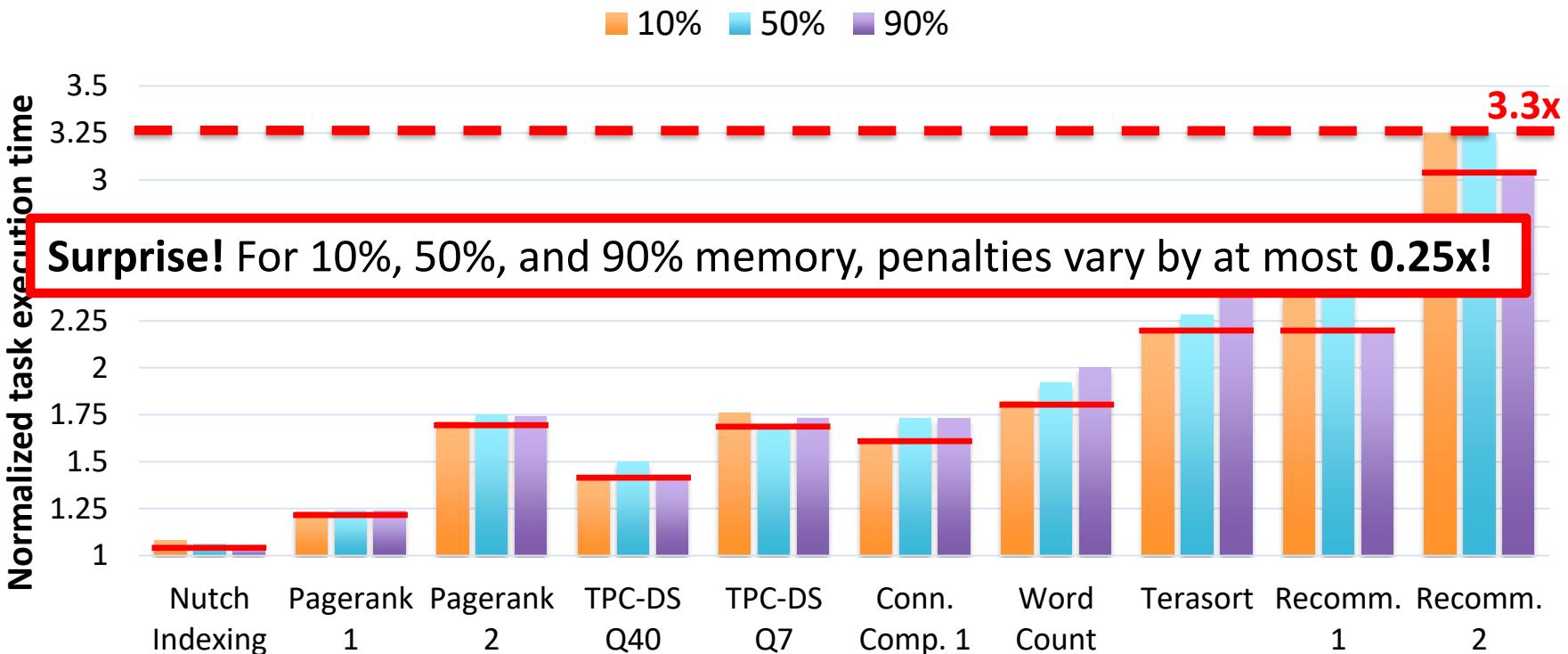
Logarithmic external merge algorithms

- Merge steps required << disk spills

Elasticity of Hadoop workloads: Reducers



Elasticity of Hadoop workloads: Reducers

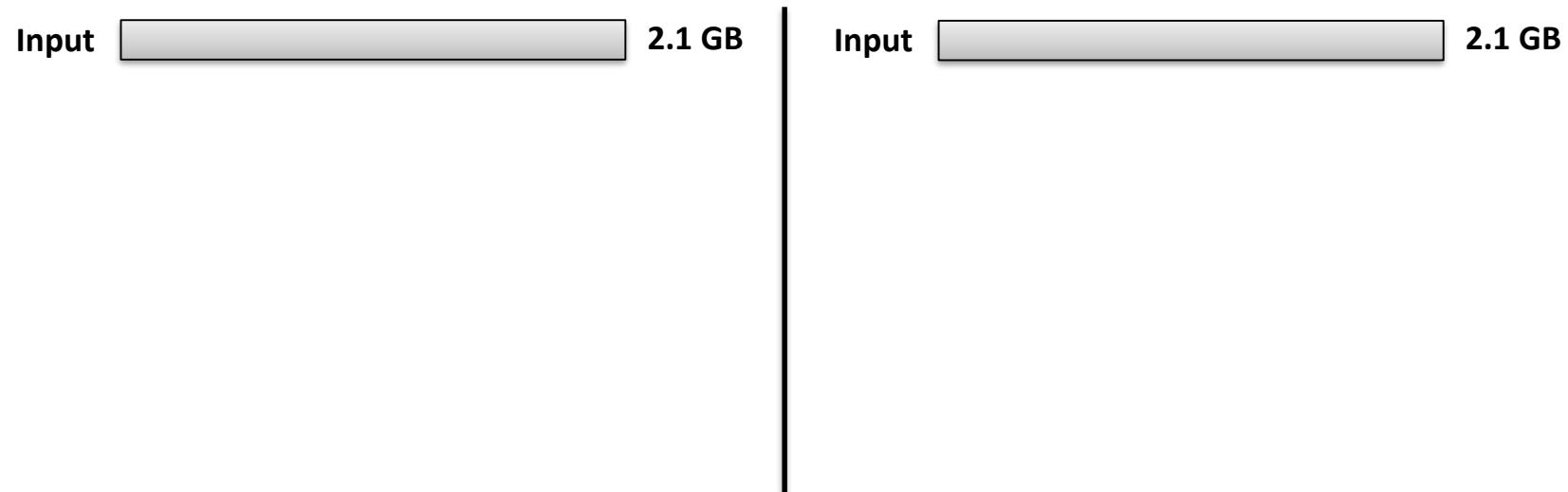


Why do penalties vary so little w/ memory?

- Static spilling threshold → comparable data spilling for 90% and 10% of memory

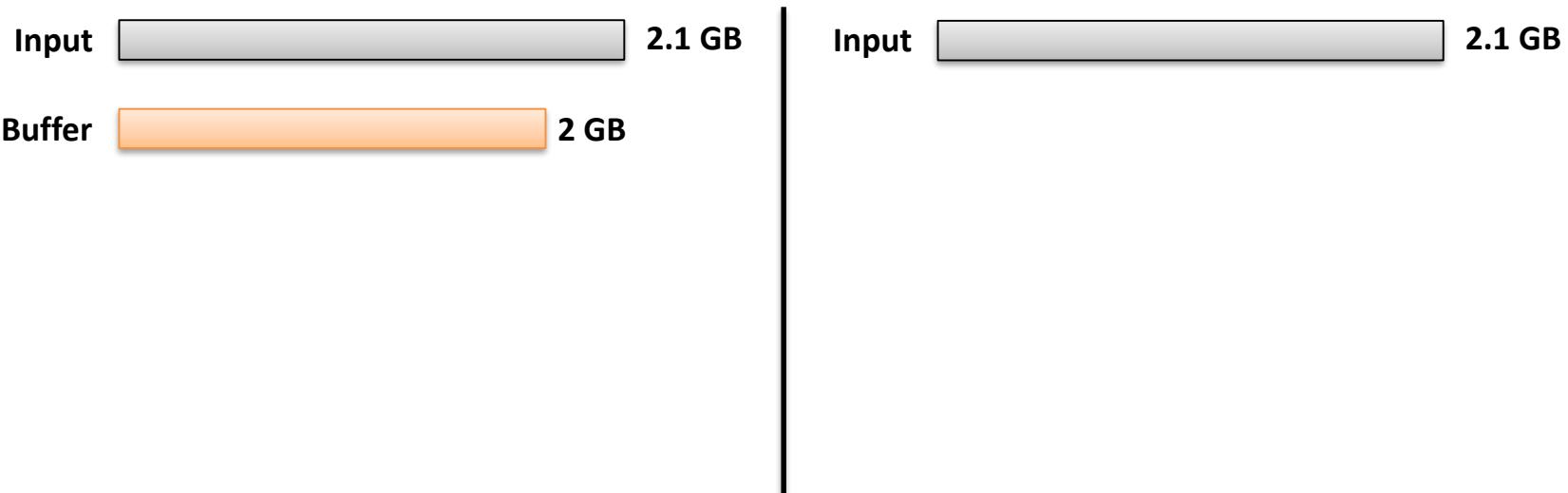
Why do penalties vary so little w/ memory?

- Static spilling threshold → comparable data spilling for 90% and 10% of memory



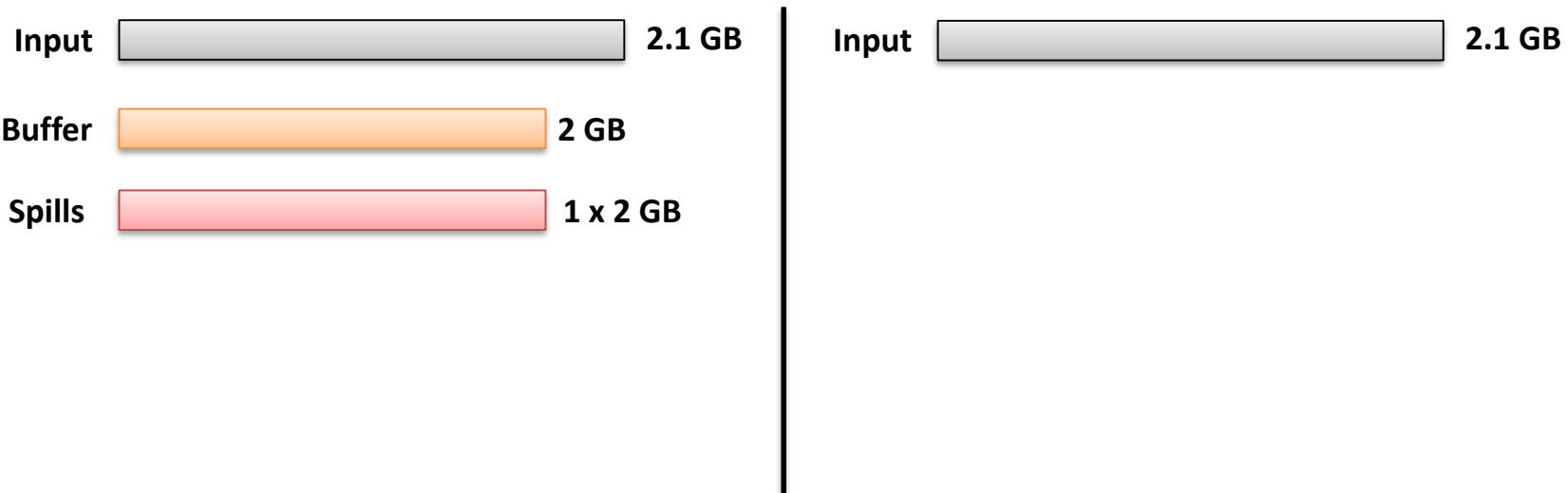
Why do penalties vary so little w/ memory?

- Static spilling threshold → comparable data spilling for 90% and 10% of memory



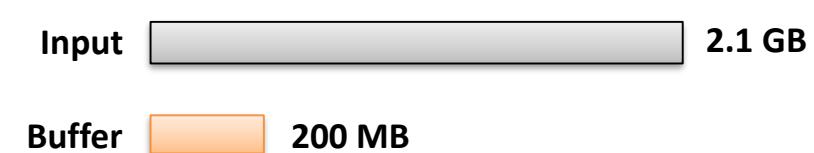
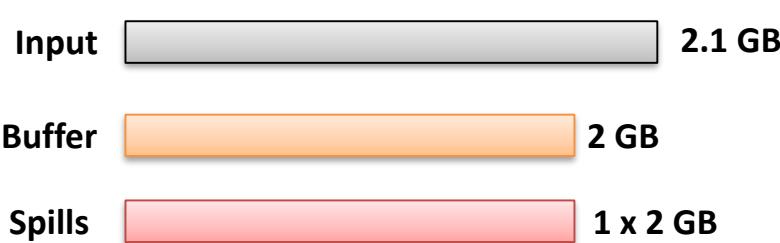
Why do penalties vary so little w/ memory?

- Static spilling threshold → comparable data spilling for 90% and 10% of memory



Why do penalties vary so little w/ memory?

- Static spilling threshold → comparable data spilling for 90% and 10% of memory



Why do penalties vary so little w/ memory?

- Static spilling threshold → comparable data spilling for 90% and 10% of memory

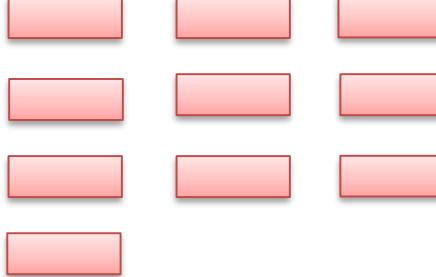
Input  **2.1 GB**

Buffer  **2 GB**

Spills  **1 x 2 GB**

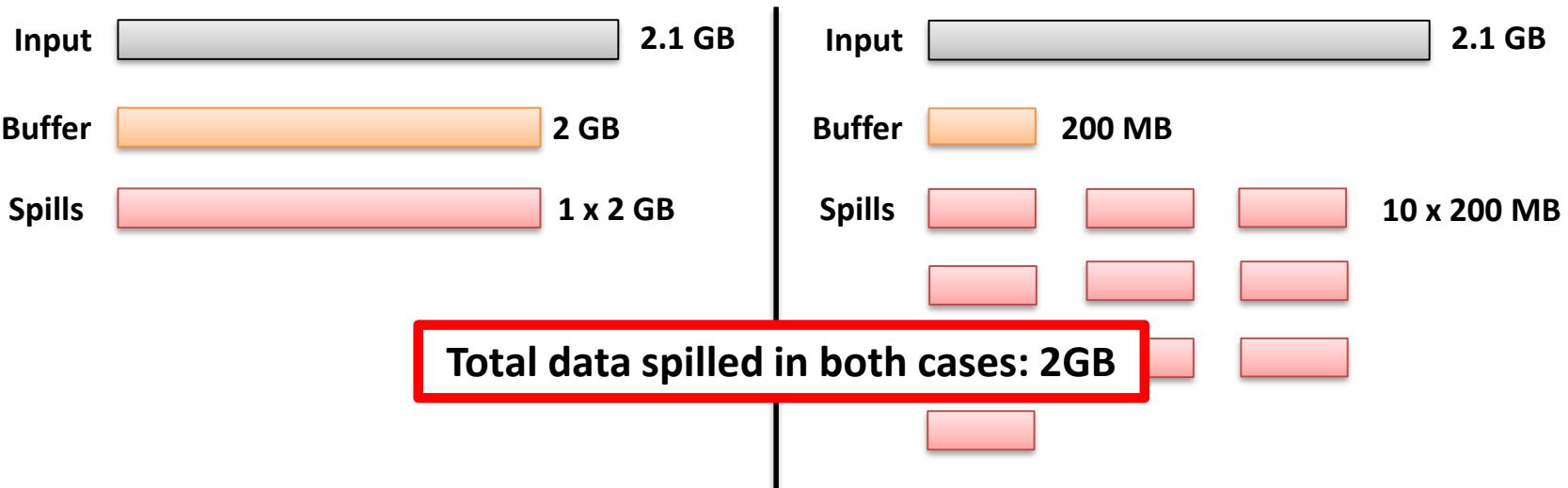
Input  **2.1 GB**

Buffer  **200 MB**

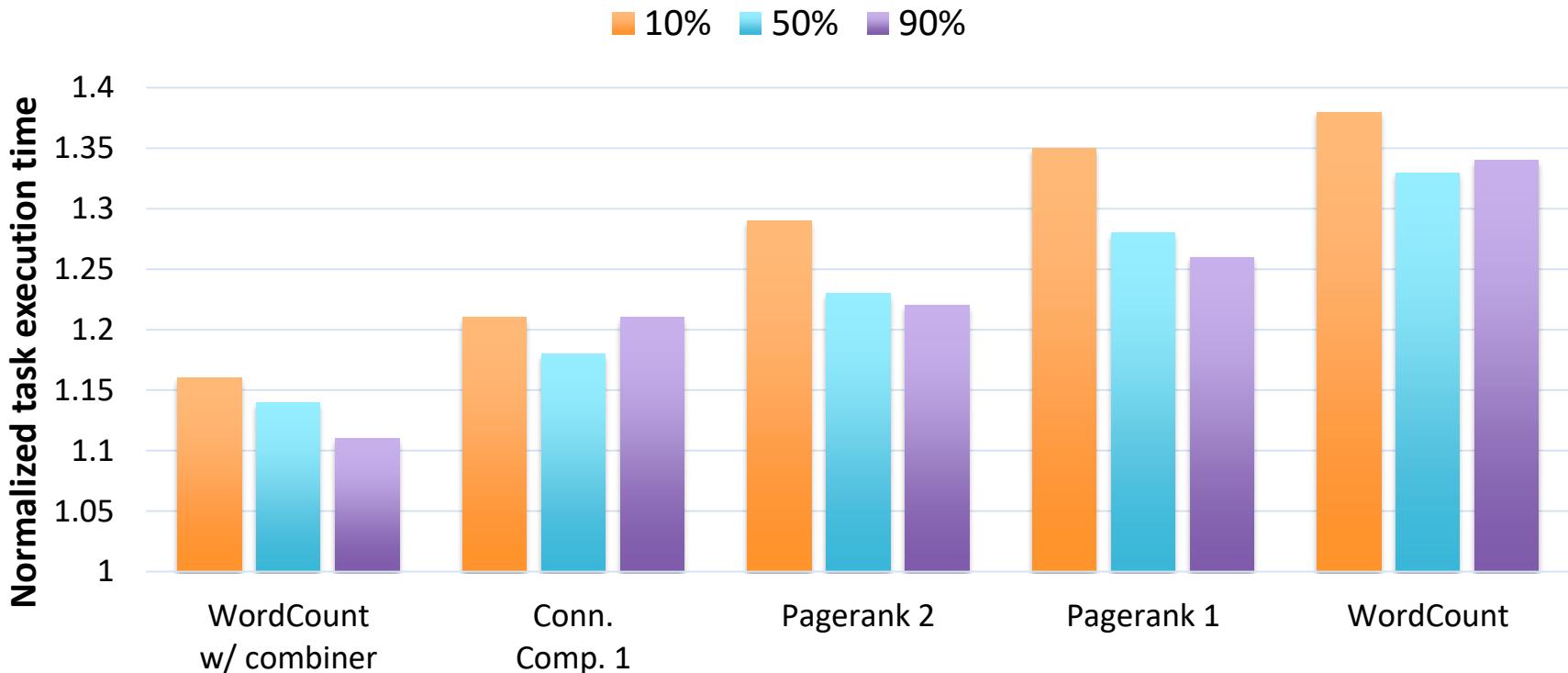
Spills  **10 x 200 MB**

Why do penalties vary so little w/ memory?

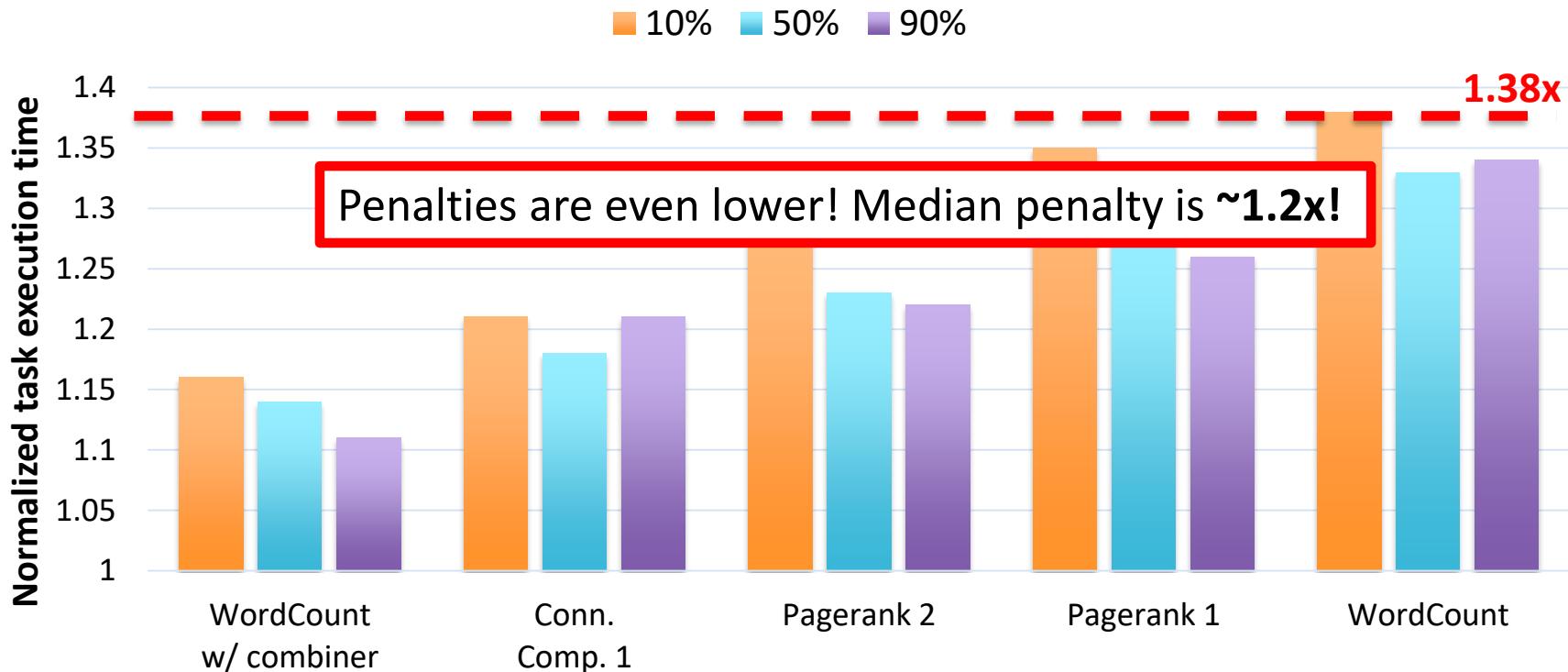
- Static spilling threshold → comparable data spilling for 90% and 10% of memory



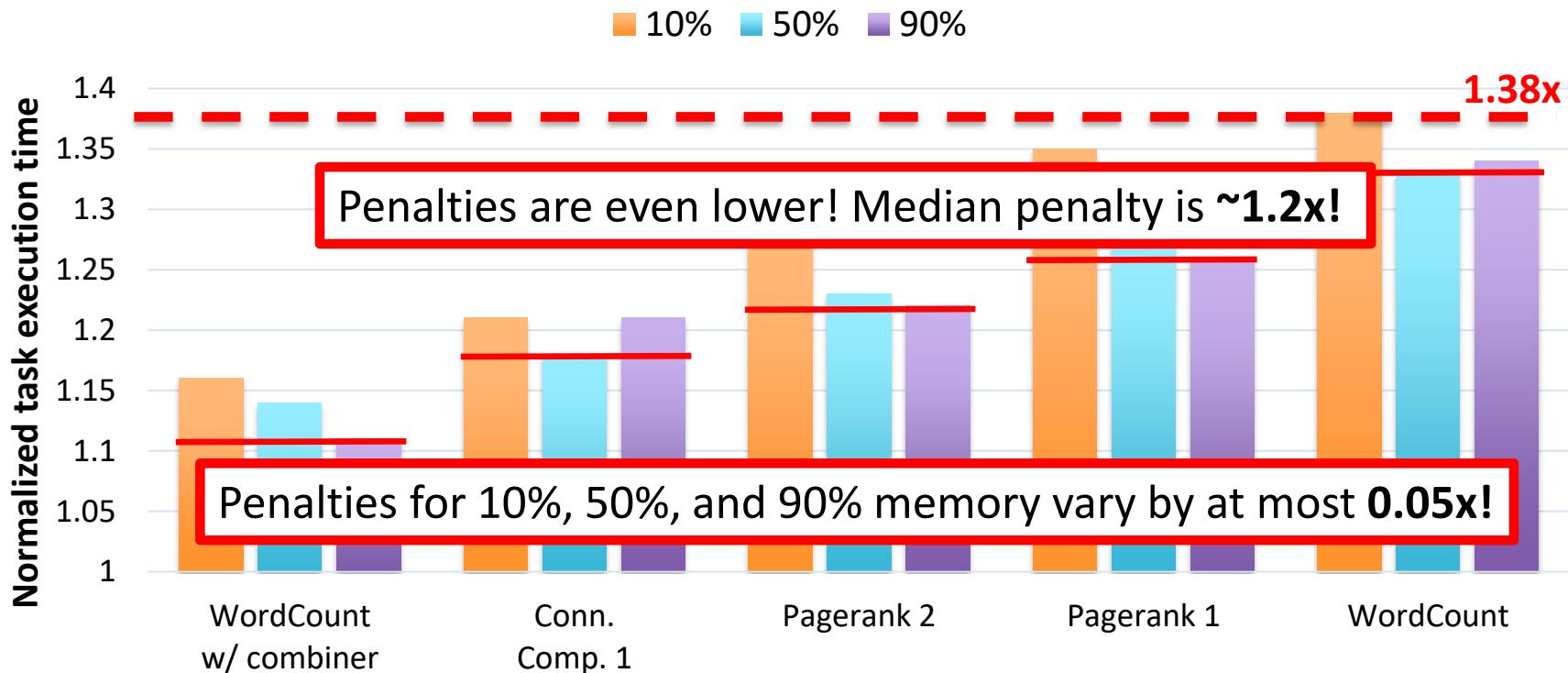
Elasticity of Hadoop workloads: Mappers



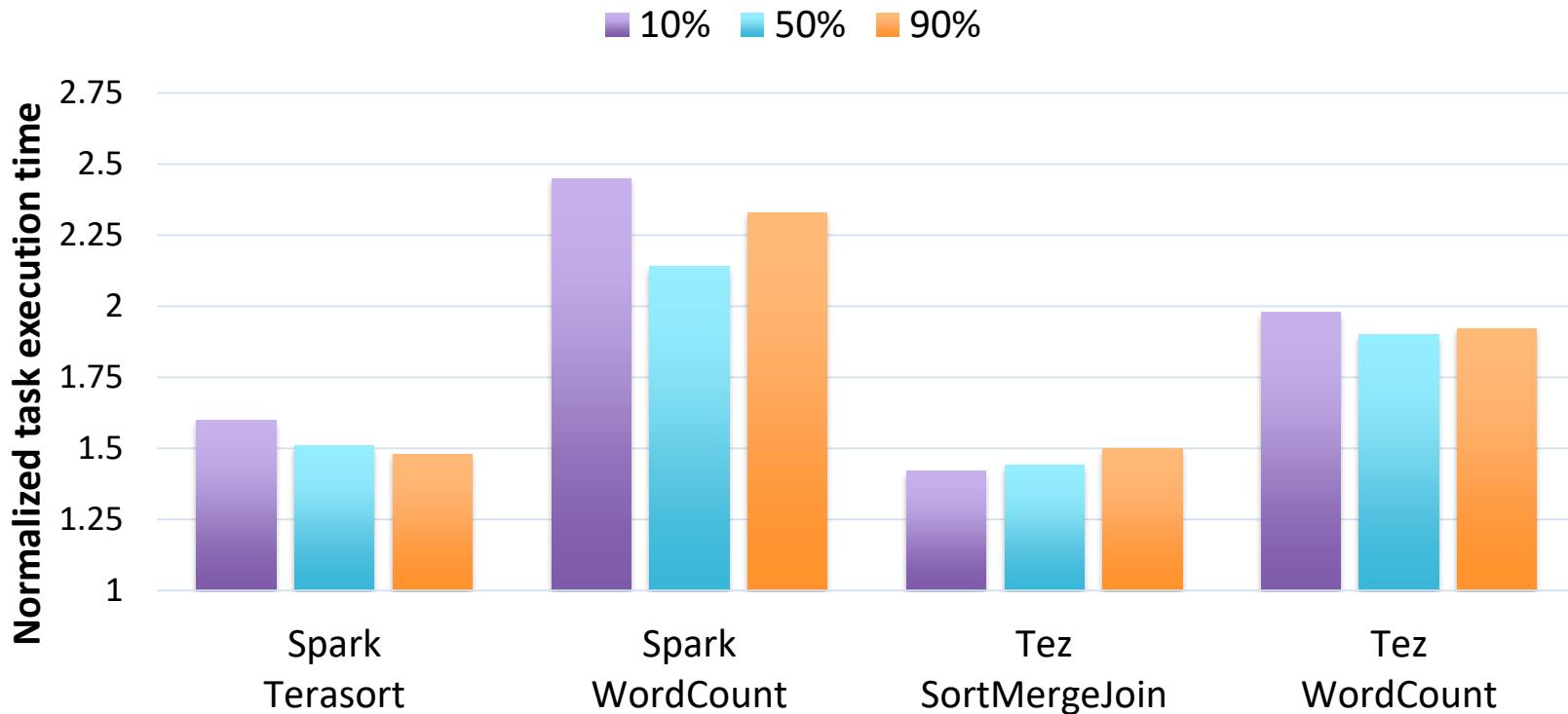
Elasticity of Hadoop workloads: Mappers



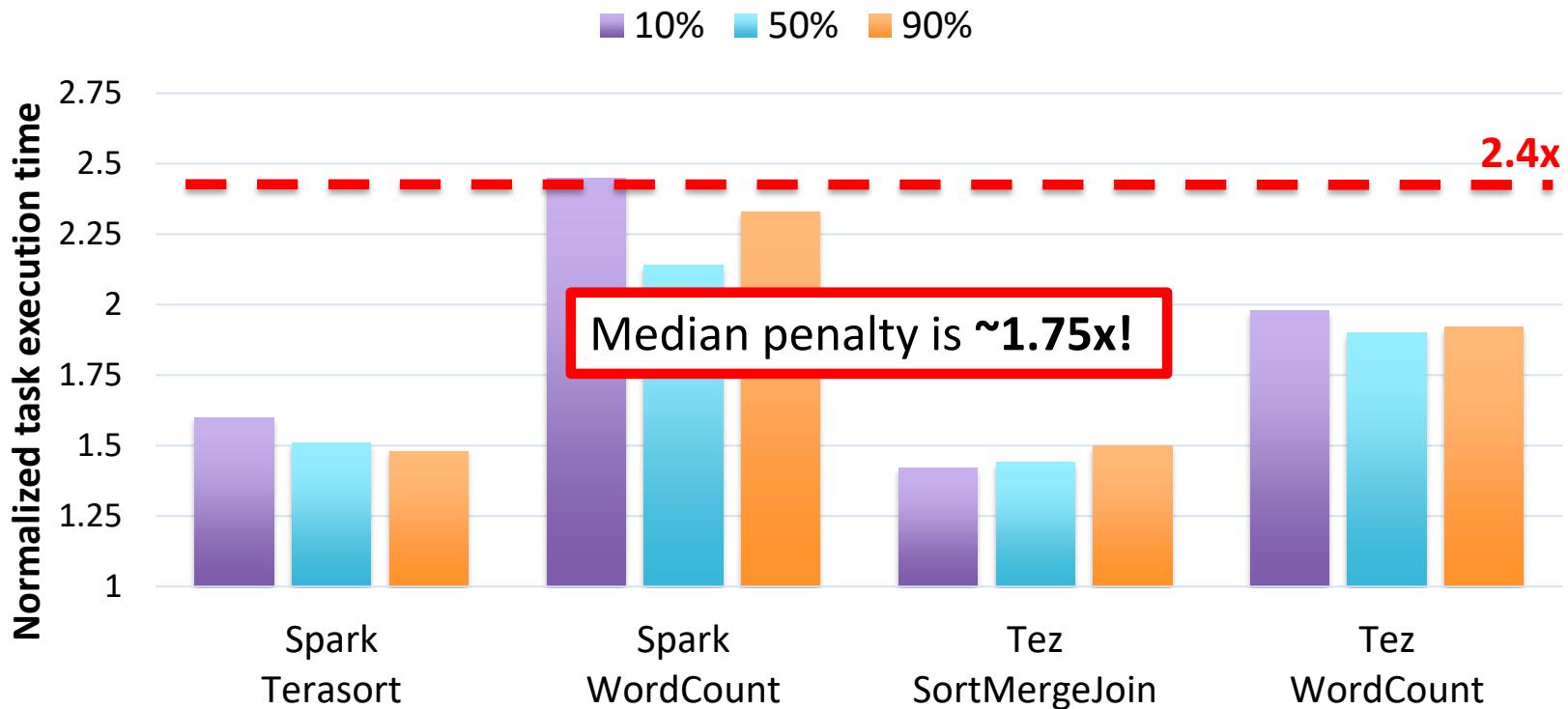
Elasticity of Hadoop workloads: Mappers



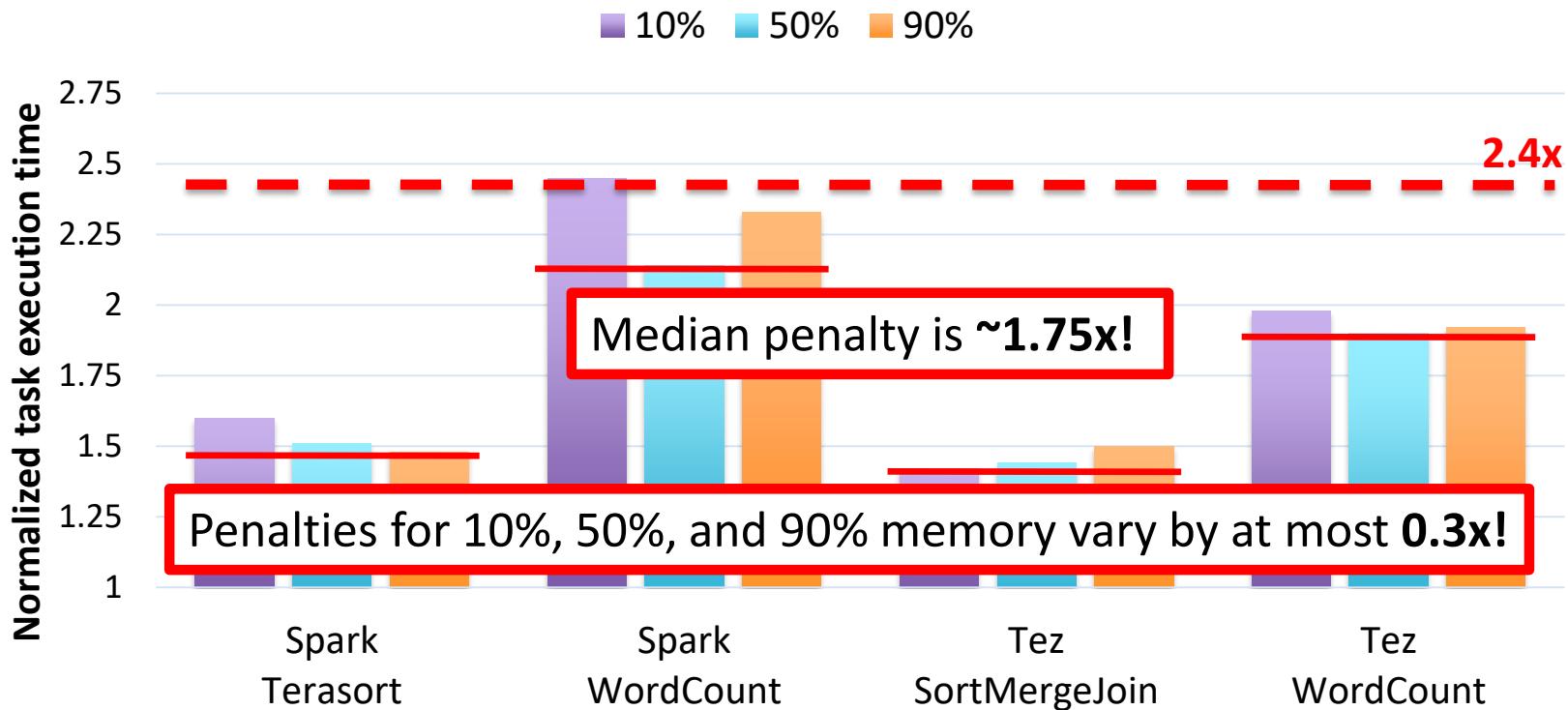
Elasticity of Spark and Tez workloads



Elasticity of Spark and Tez workloads

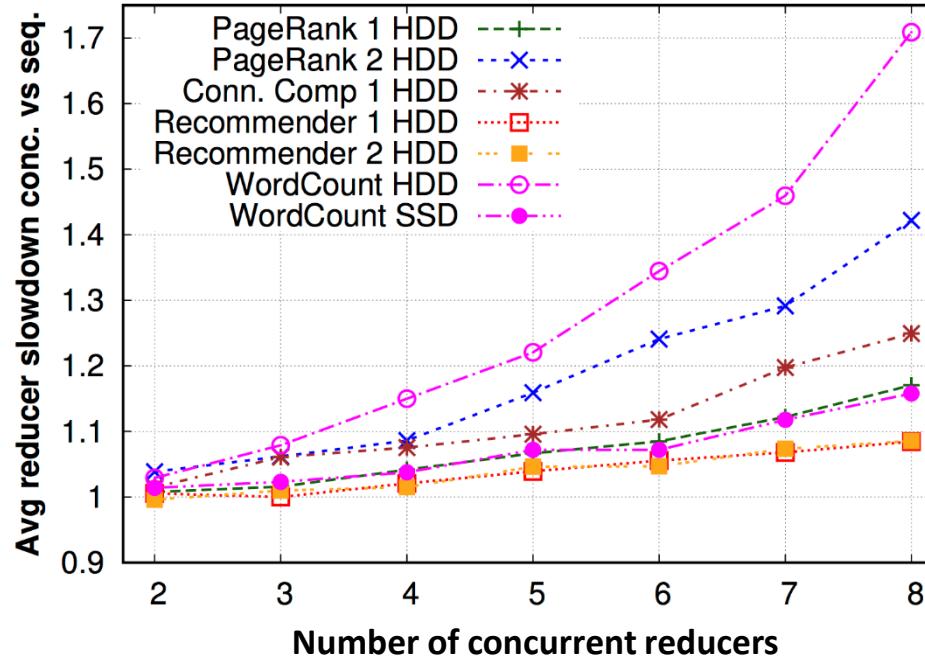


Elasticity of Spark and Tez workloads



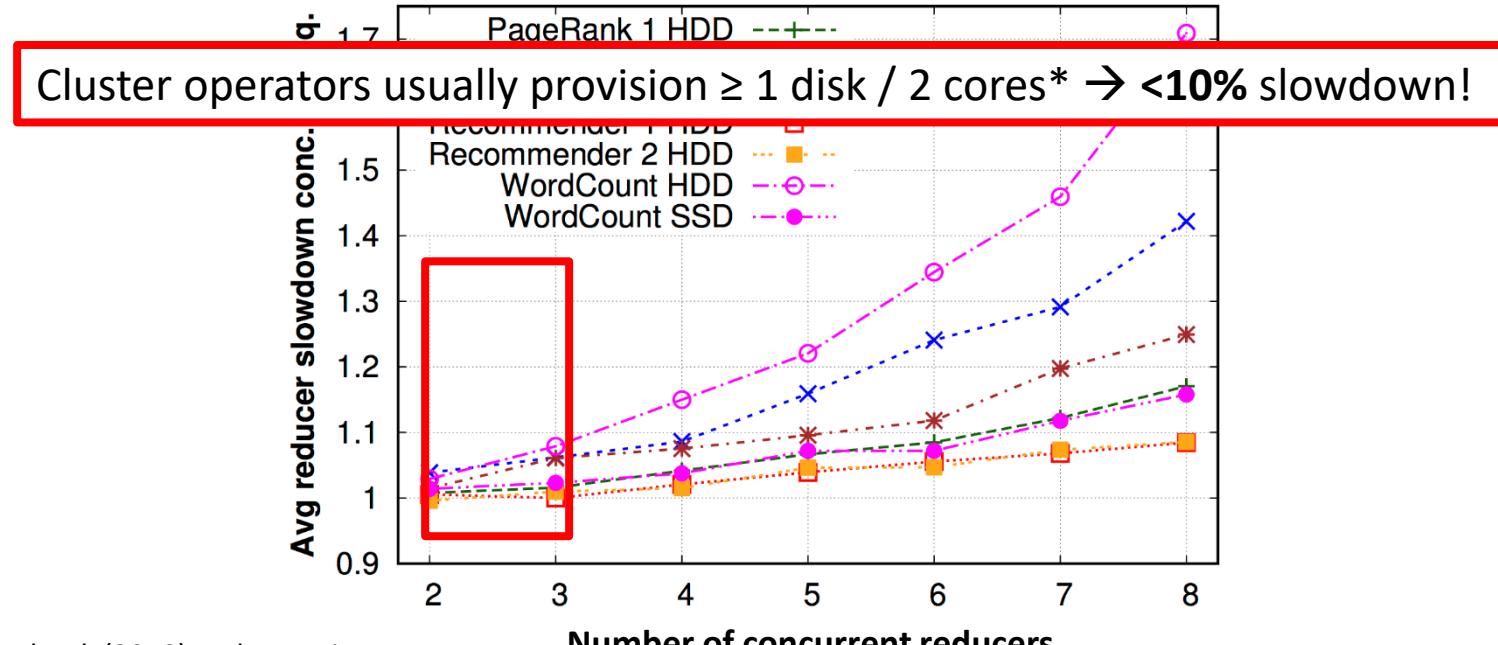
Does the additional I/O cause disk contention?

- Measure slowdown of elastic tasks on same machine spilling to the same disk



Does the additional I/O cause disk contention?

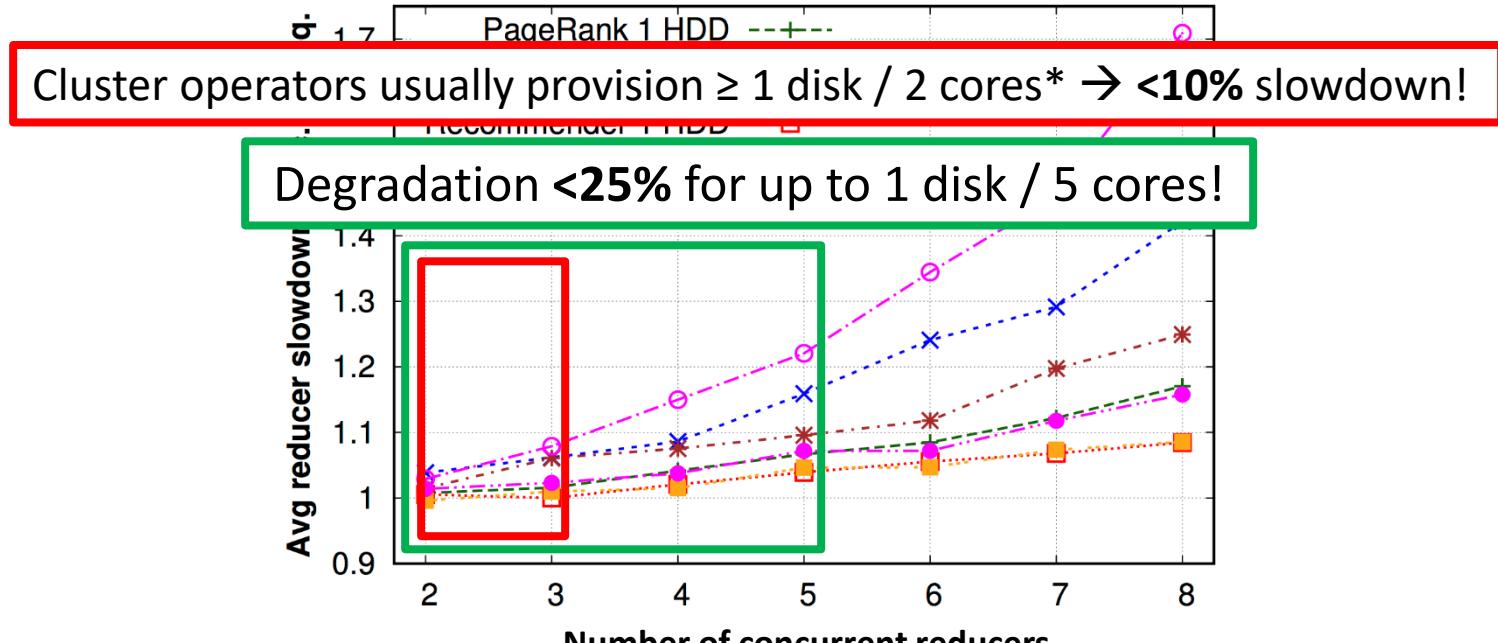
- Measure slowdown of elastic tasks on same machine spilling to the same disk



* Facebook (2010) and Nutanix

Does the additional I/O cause disk contention?

- Measure slowdown of elastic tasks on same machine spilling to the same disk



* Facebook (2010) and Nutanix

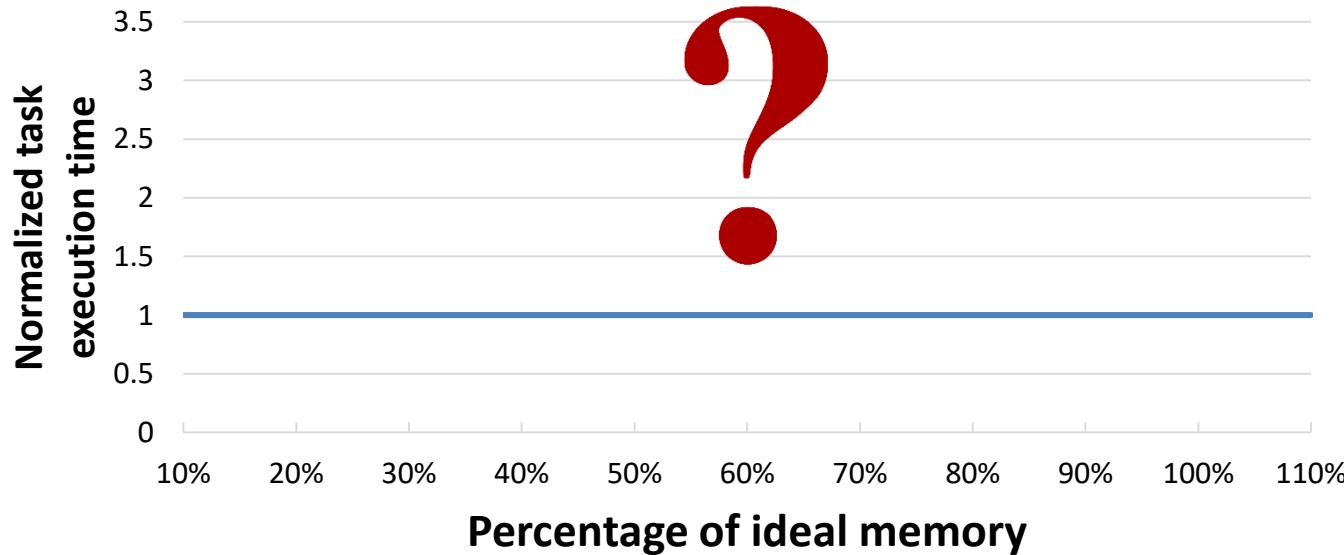
Summary: Memory Elasticity of real workloads

- ✓ Modest performance penalties (<1.6x median)
- ✓ Similar penalties for 10% and 90% of ideal memory
- ✓ Disk contention negligible for existing clusters' setup (<10%)

MODELING MEMORY ELASTICITY

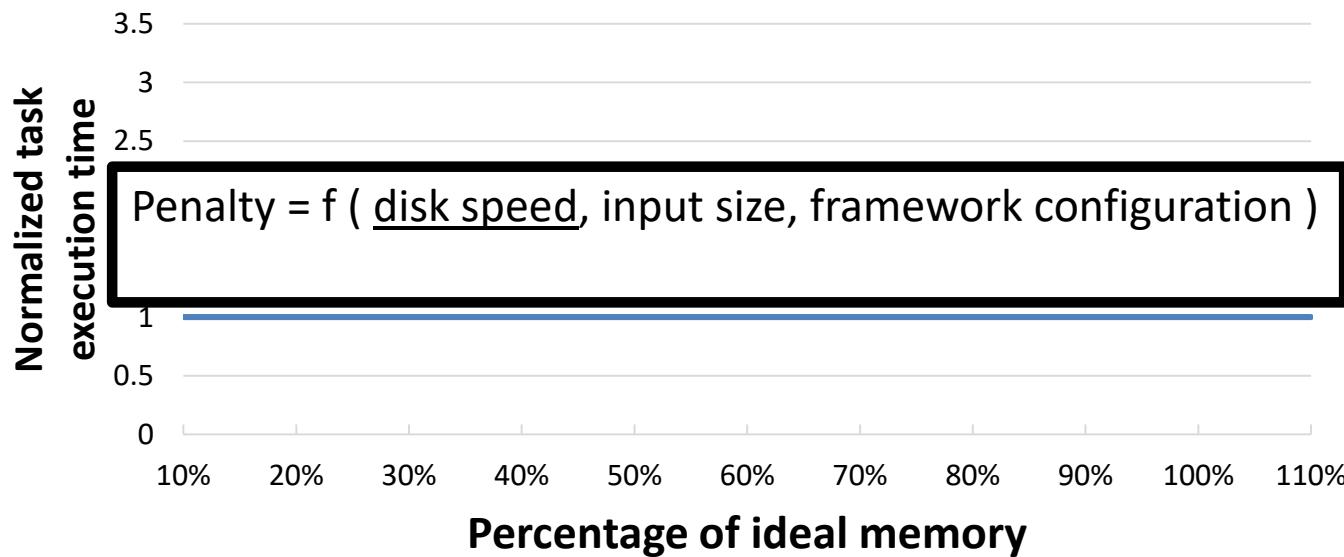
Modeling Memory Elasticity

- How does penalty vary for a task?



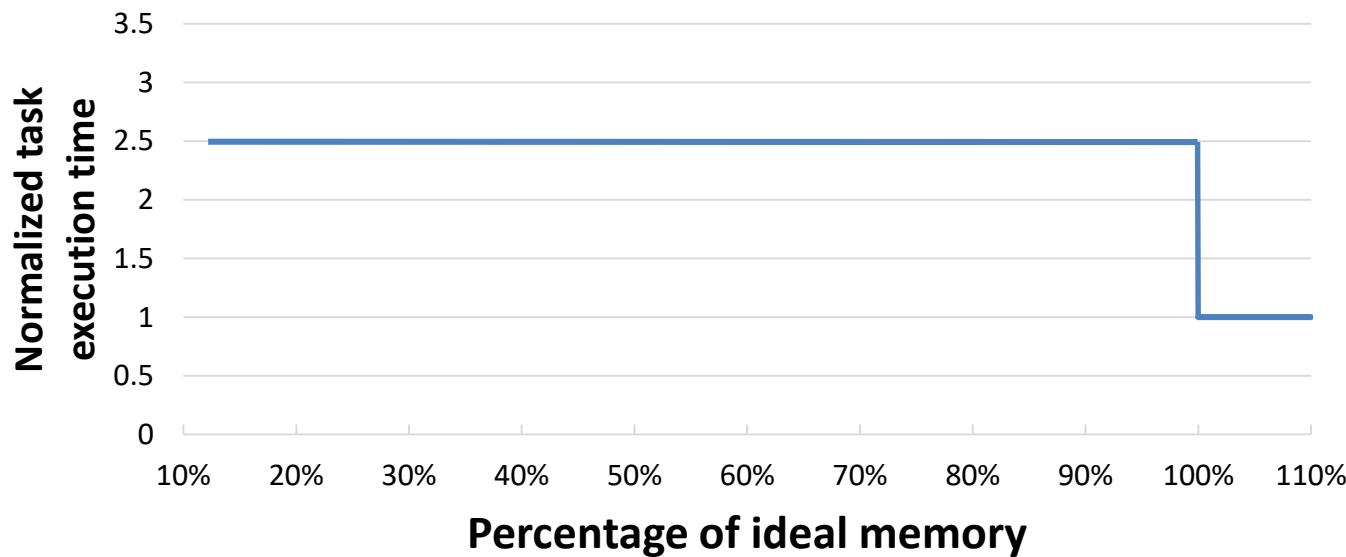
Modeling Memory Elasticity

- How does penalty vary for a task?



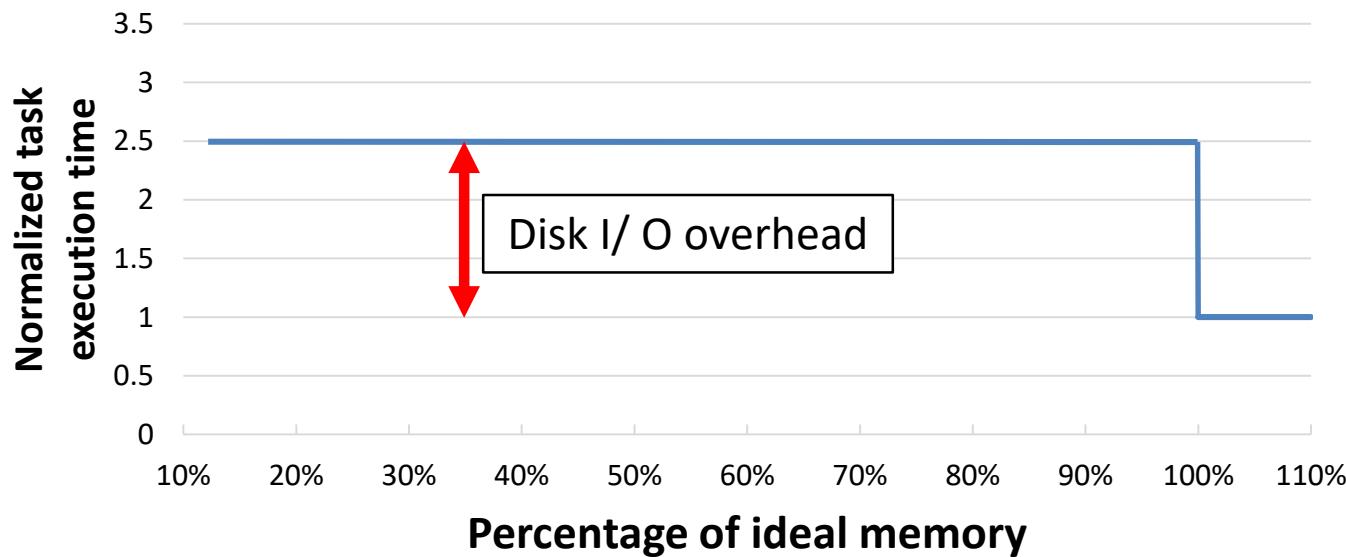
Modeling Memory Elasticity

- Penalties vary little between percentages → step model



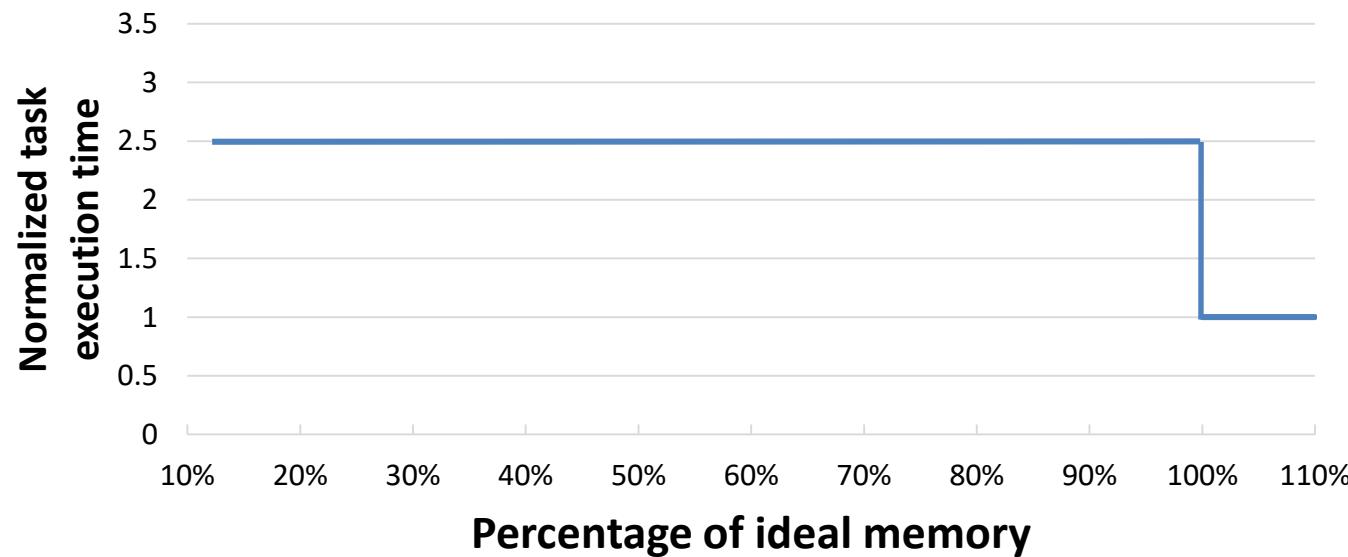
Modeling Memory Elasticity

- Penalties vary little between percentages → step model



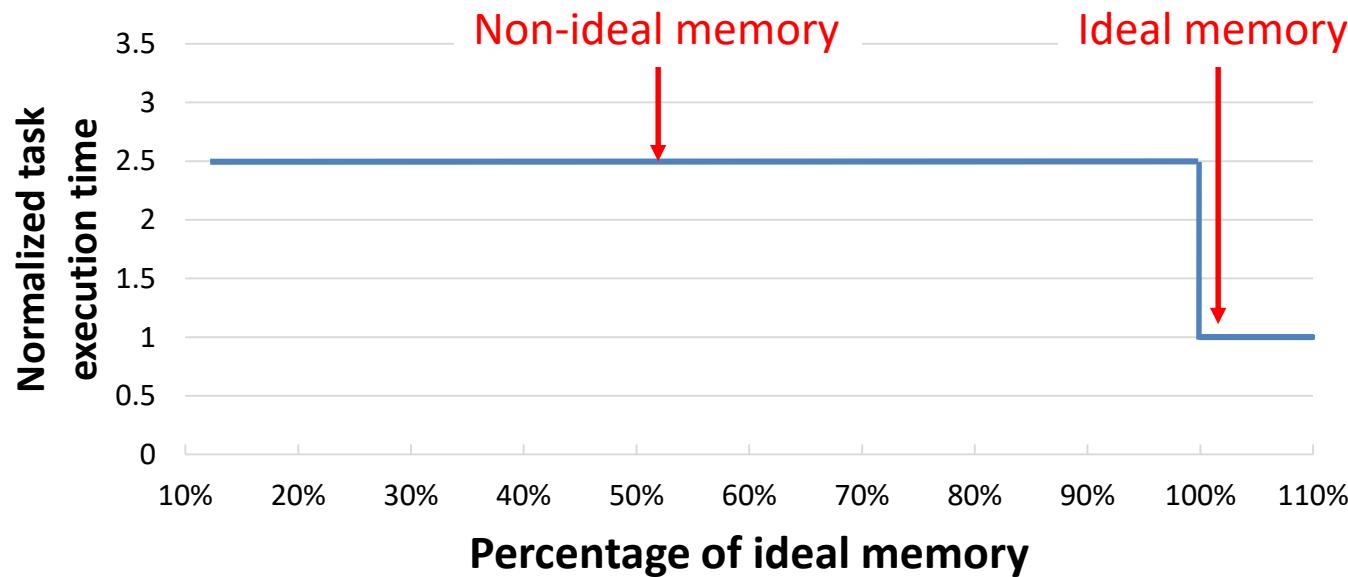
Modeling Memory Elasticity

- Requires 2 profiling runs → infers all other points



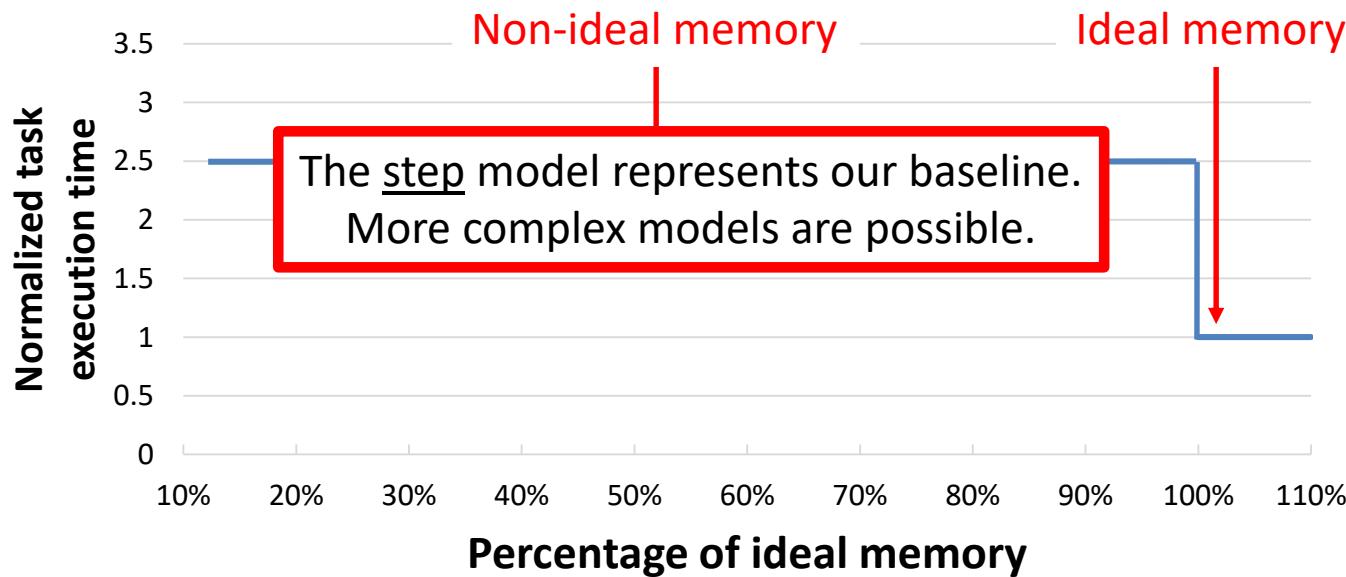
Modeling Memory Elasticity

- Requires 2 profiling runs → infers all other points

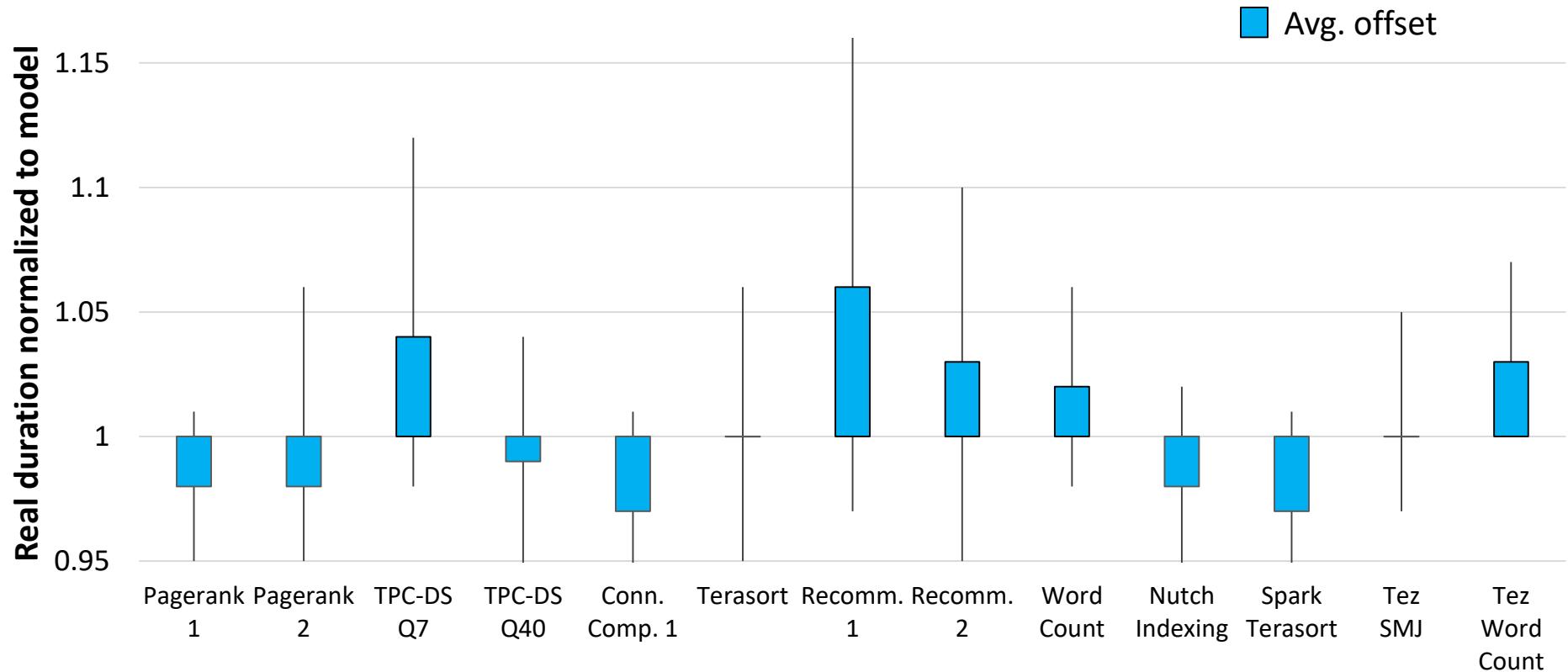


Modeling Memory Elasticity

- Requires 2 profiling runs → infers all other points



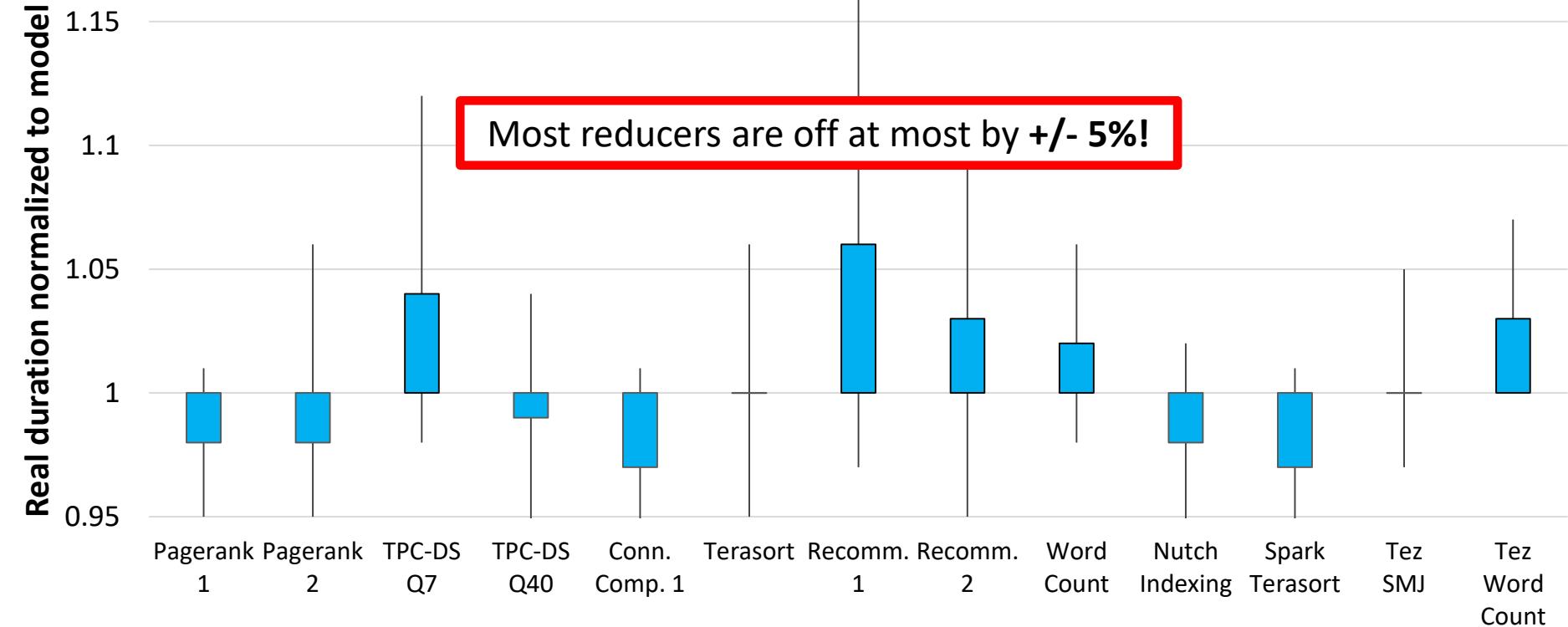
Accuracy of our reducer model



Accuracy of our reducer model

Avg. offset

Most reducers are off at most by +/- 5%!



Summary: Modeling memory elasticity

- ✓ Step model is adequate (more complex models available)
- ✓ Only 2 profiling points → full model
- ✓ Models are very robust (+/- 5% error for most reducers)

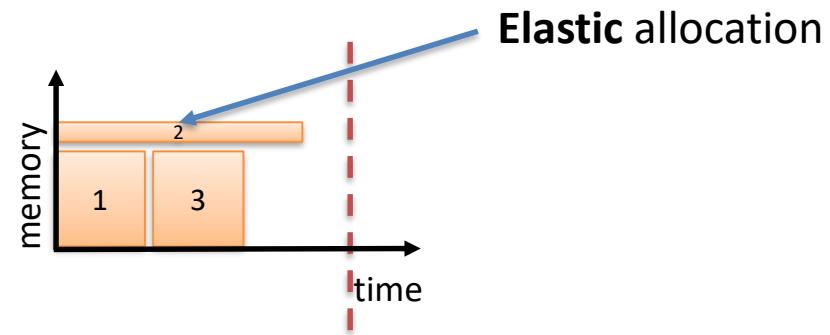
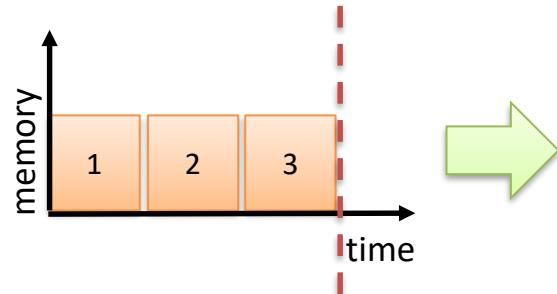
LEVERAGING MEMORY ELASTICITY IN CLUSTER SCHEDULING

How can a scheduler reason about Memory Elasticity?

Trade-off between

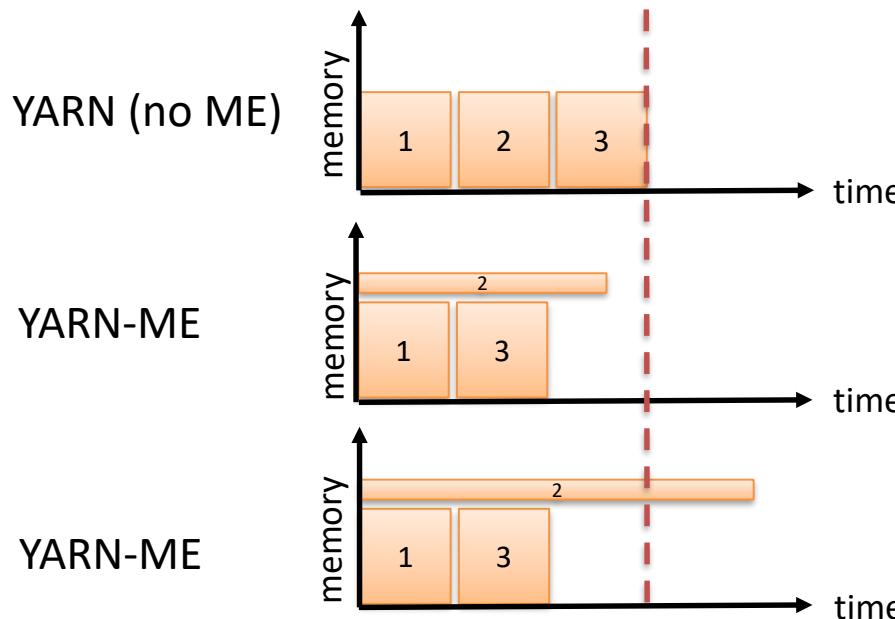
↓ task queueing time

↑ task execution time

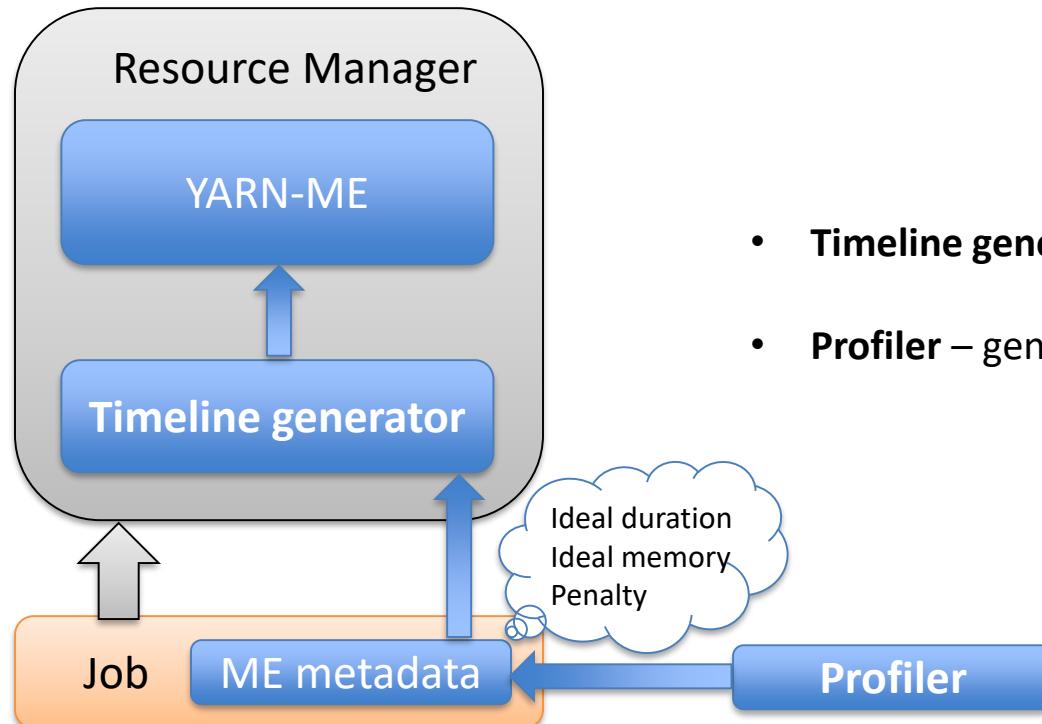


YARN-ME: Decision process

- Make an elastic allocation *iff* it does not exceed the expected job completion time



YARN-ME: Design and components

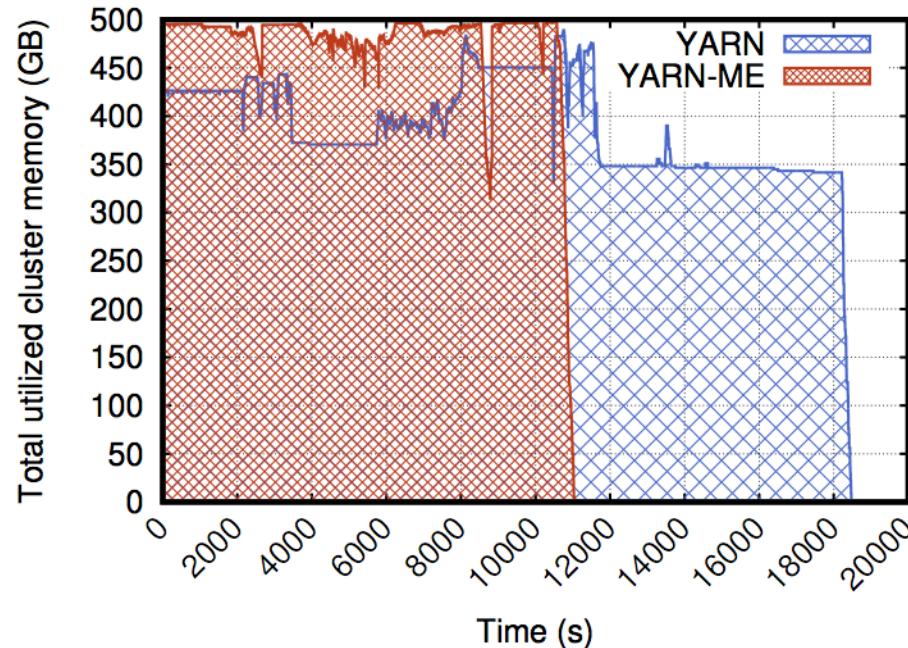


Components

- **Timeline generator** – computes expected JCTs
- **Profiler** – generates the model metadata

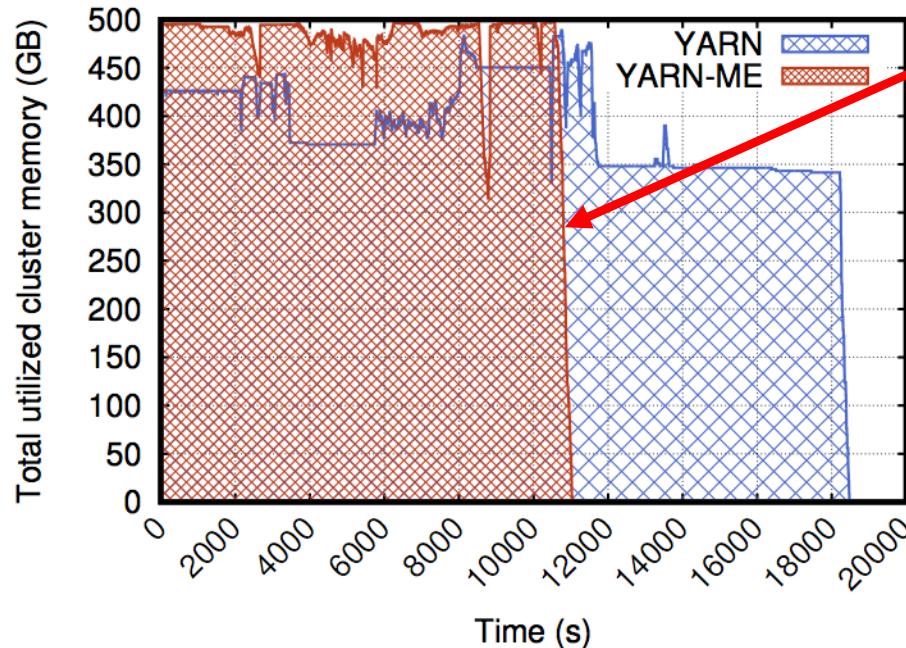
Memory utilization analysis for YARN-ME

- 50 node cluster
 - Homogeneous trace: 5x Pagerank jobs



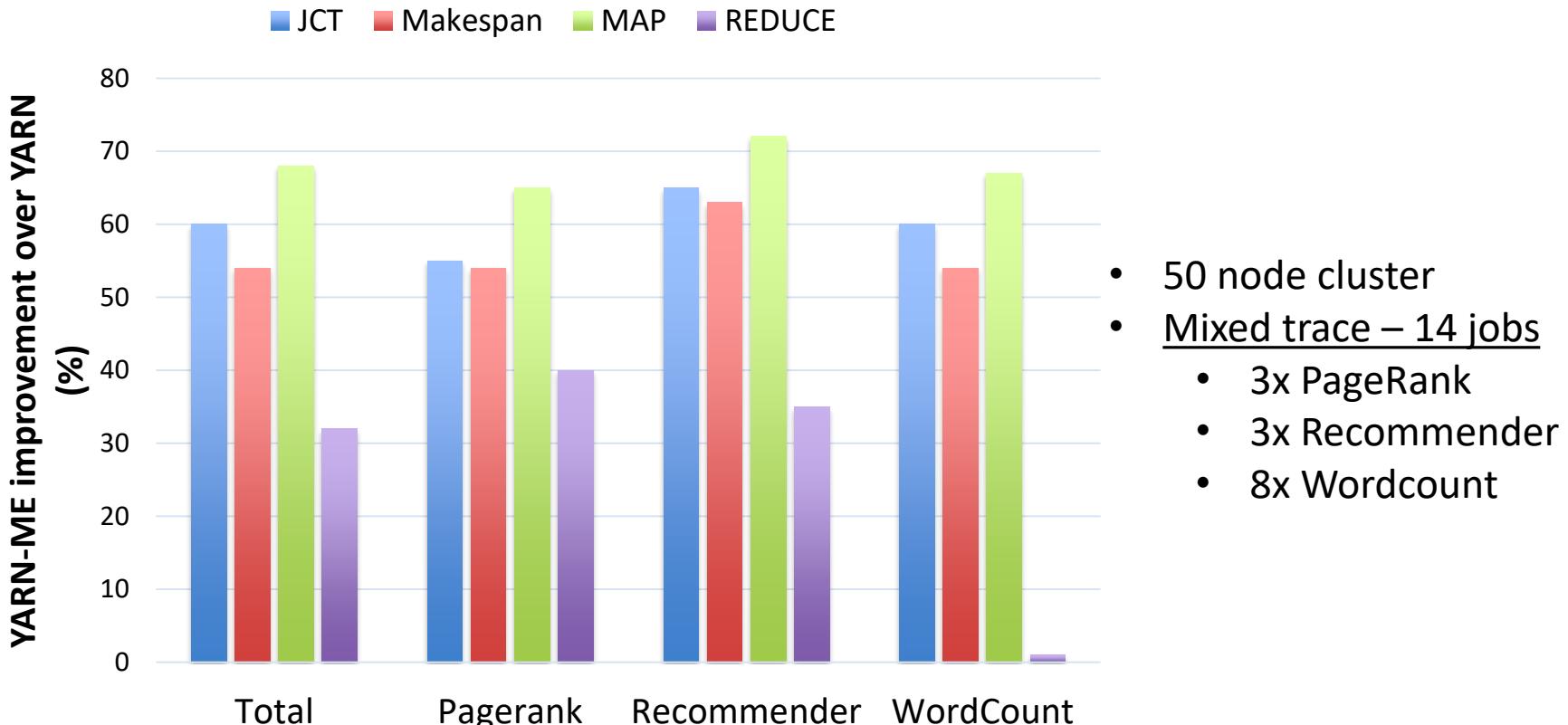
Memory utilization analysis for YARN-ME

- 50 node cluster
- Homogeneous trace: 5x Pagerank jobs

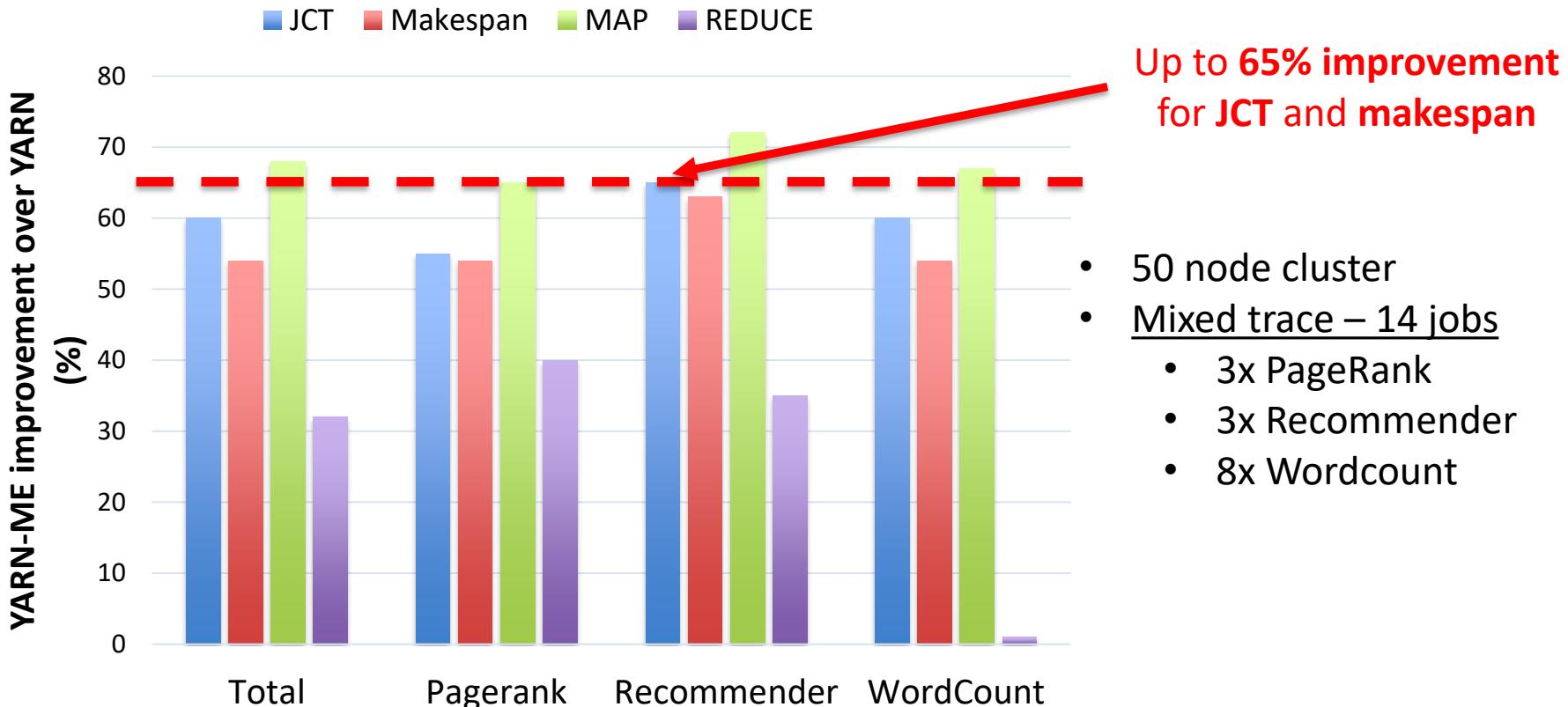


Memory utilization increased to 95%

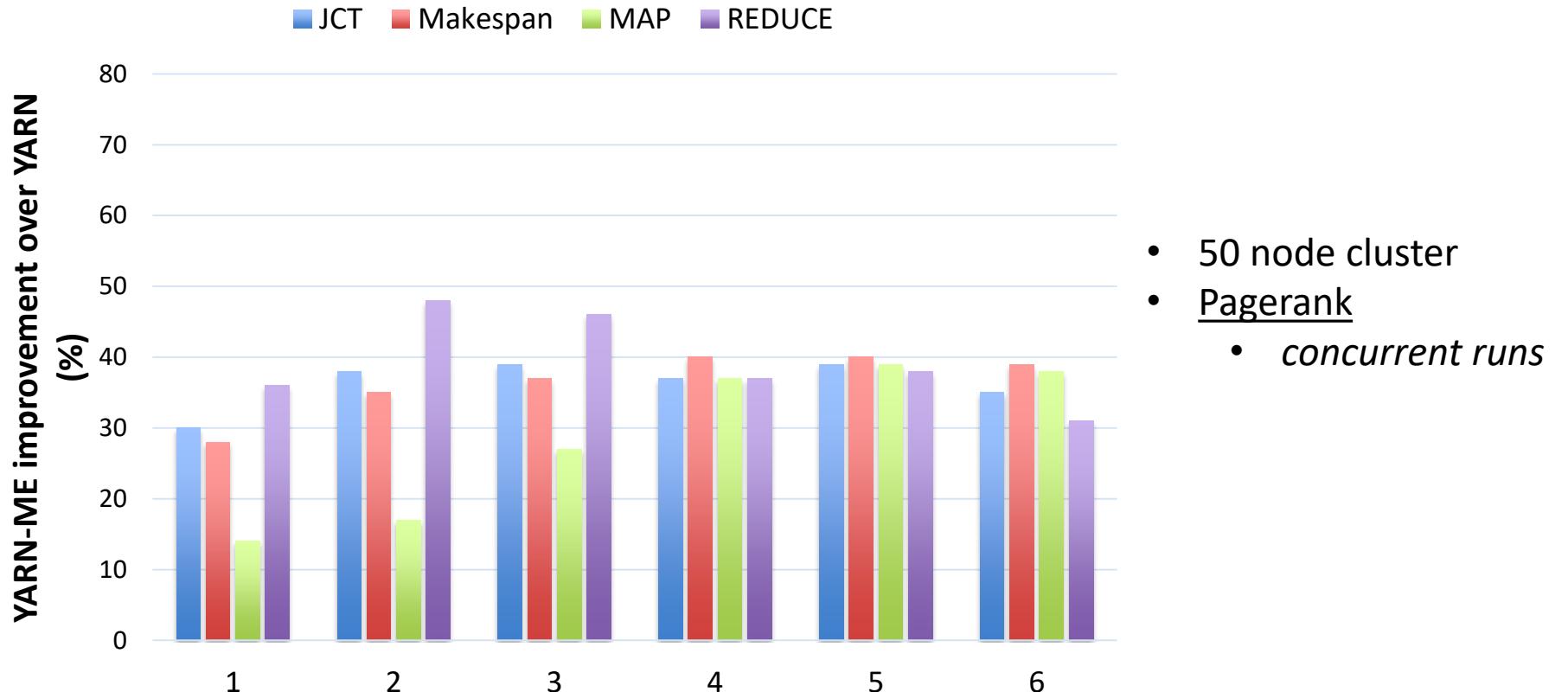
What gains can YARN-ME achieve for heterogeneous workloads?



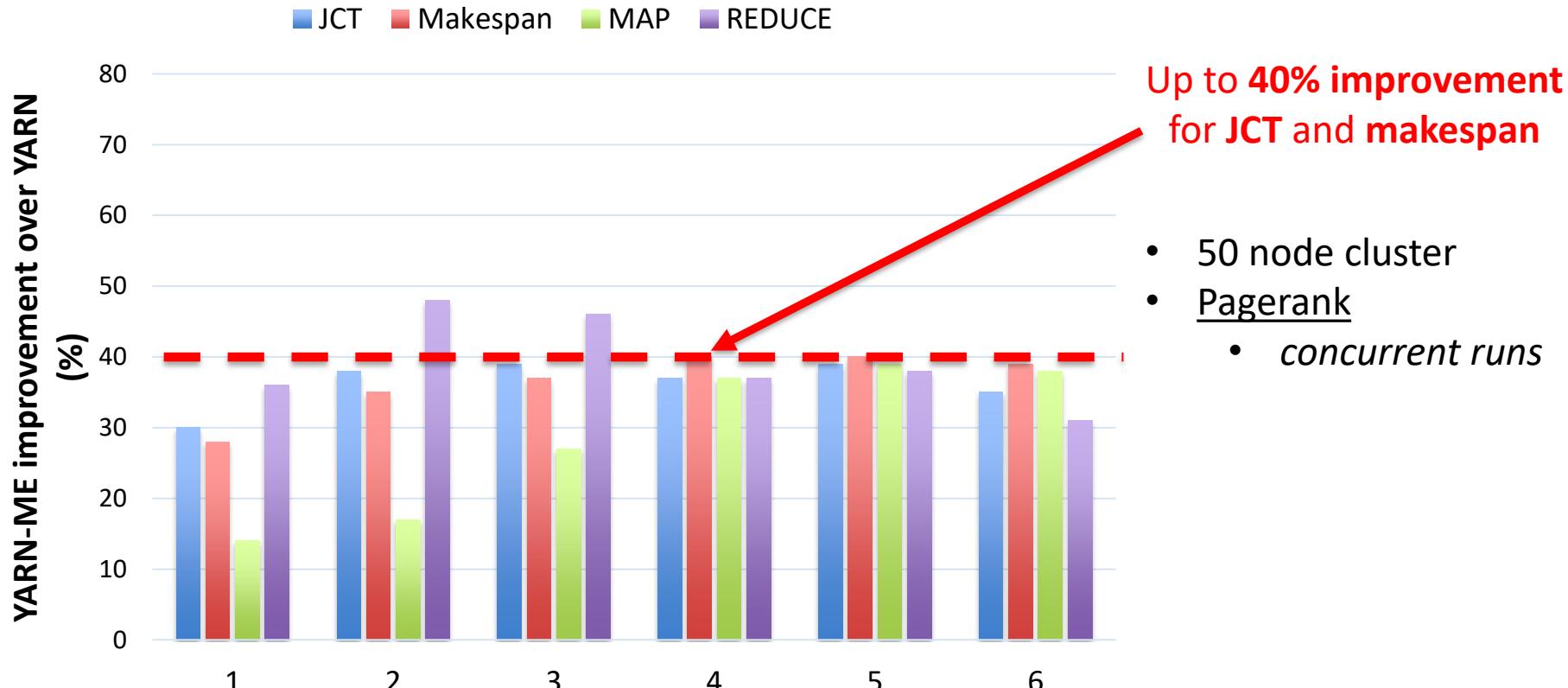
What gains can YARN-ME achieve for heterogeneous workloads?



What gains can YARN-ME achieve for homogeneous workloads?



What gains can YARN-ME achieve for homogeneous workloads?



Trace-driven simulation of YARN-ME

- We built DSS (the Discrete Scheduler Simulator)
 - and it is open-source!

Trace parameter sweep

- > 8,000 traces
- up to 3,000 nodes
- results comparable to real workloads

Robustness analysis

- > 20,000 traces
- YARN-ME is robust to model mis-estimations

Related work

- Efficient packing → better resource utilization
 - Tetris [**SIGCOMM '14**], GRAPHENE [**OSDI '16**]
- Collocate batch-jobs with latency-critical services
 - Heracles [**ISCA '15**]
- Resource over-committing
 - Apollo [**OSDI '14**], Borg [**EuroSys '15**]
- Suspend tasks under memory pressure
 - ITask [**SOSP '15**]

Conclusion: Don't cry over spilled records!

- ✓ Memory Elasticity → highly **predictable**, low penalty
- ✓ Memory Elasticity in scheduling → trade task **queueing-time** for **running-time**
- ✓ YARN-ME → up to 60% improvement in average JCT
- ✓ DSS code available: <https://github.com/epfl-labos/dss>