# PAL: A VARIABILITY-AWARE POLICY FOR SCHEDULING ML WORKLOADS IN GPU CLUSTERS

Rutwik Jain, Brandon Tran, Keting Chen, Matthew D. Sinclair, Shivaram Venkataraman

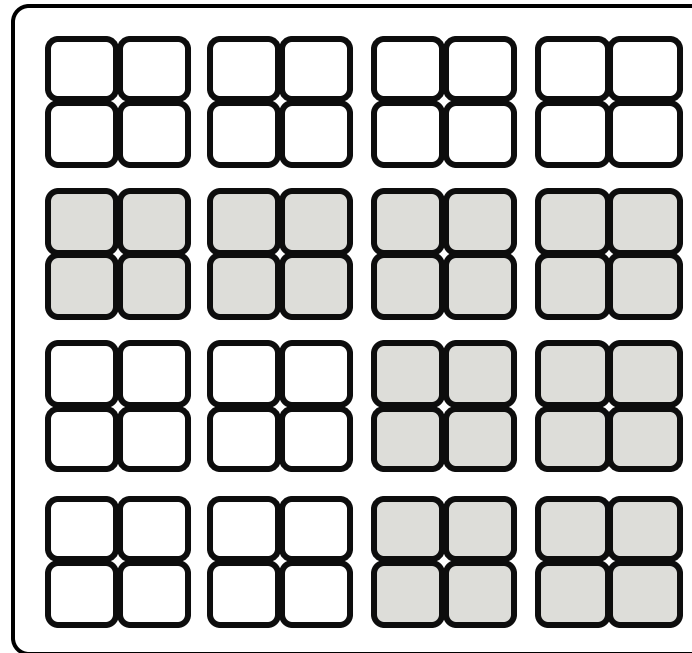Computer Sciences Department, University of Wisconsin-Madison
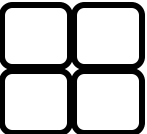
**Computer Sciences**
SCHOOL OF COMPUTER, DATA & INFORMATION SCIENCES
UNIVERSITY OF WISCONSIN–MADISON

SC24
Atlanta, GA | hpc creates.

GPU Cluster

= 1 node (4 GPUs)

System Administrator

GPU Cluster

Job Instance 1

Job Instance 2

User

= 1 node (4 GPUs)

System Administrator
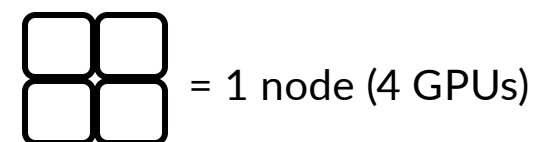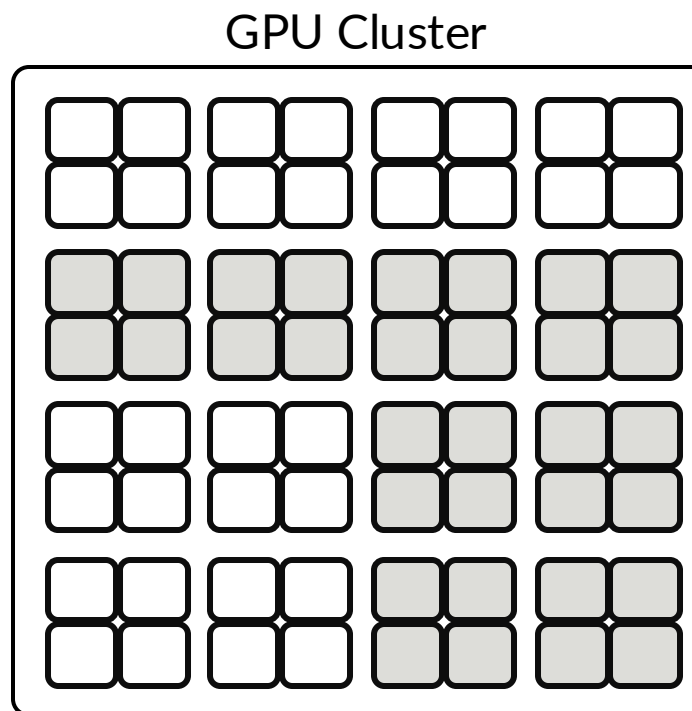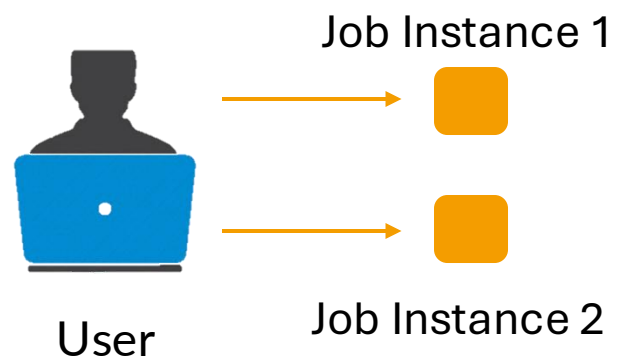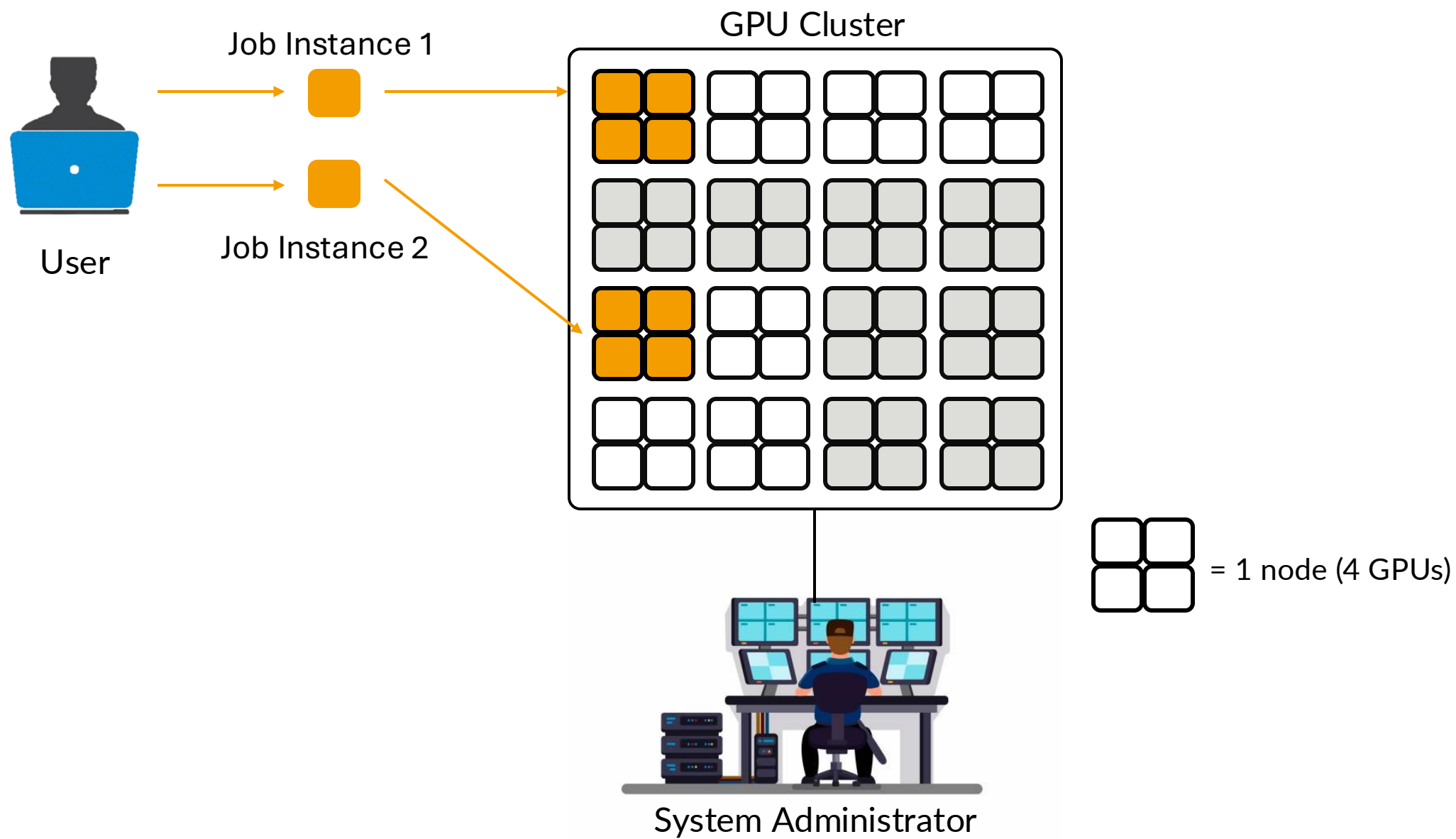
GPU Cluster

Job Instance 1

Job Instance 2

User

Node 0

22% slower!

Node 8

= 1 node (4 GPUs)

System Administrator
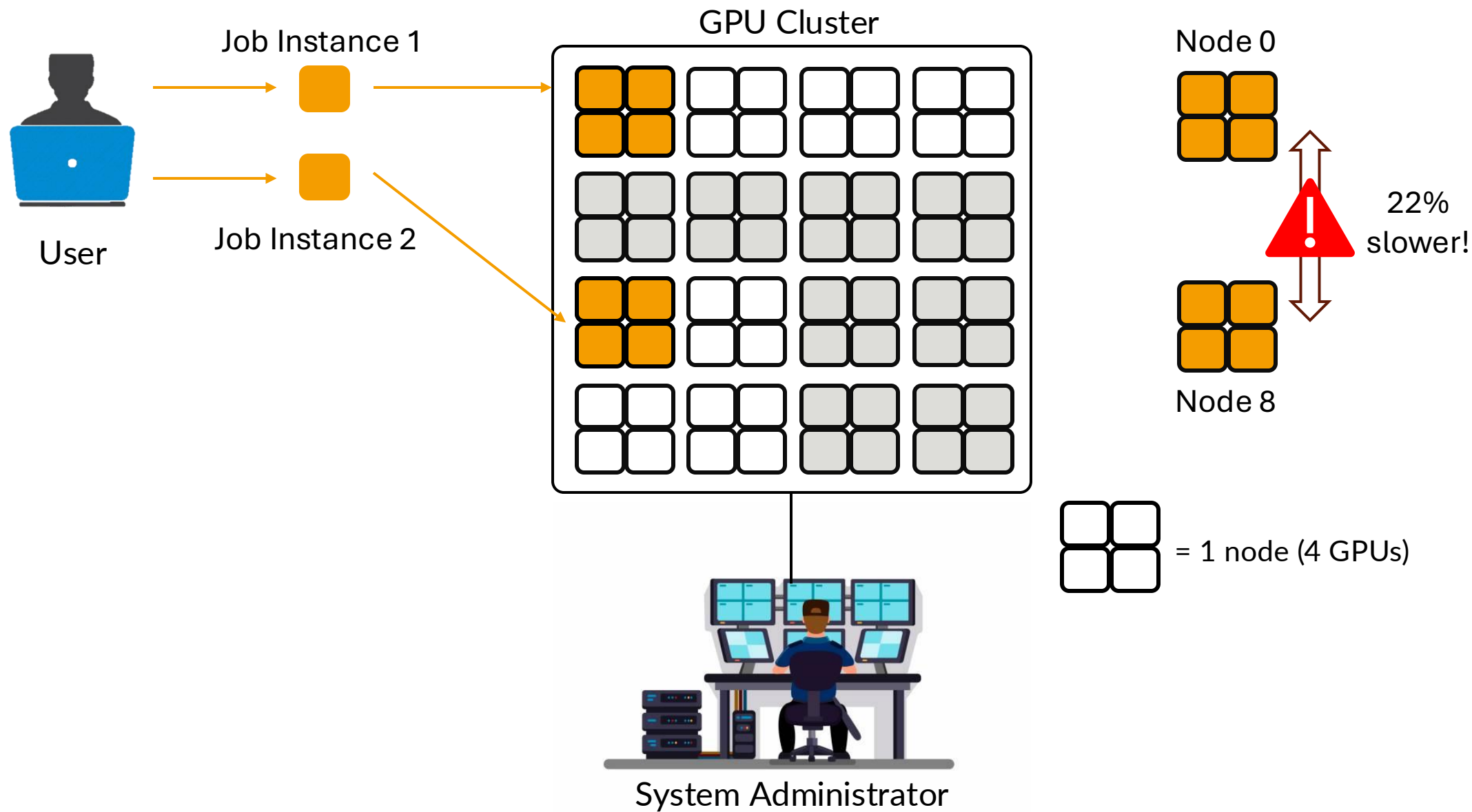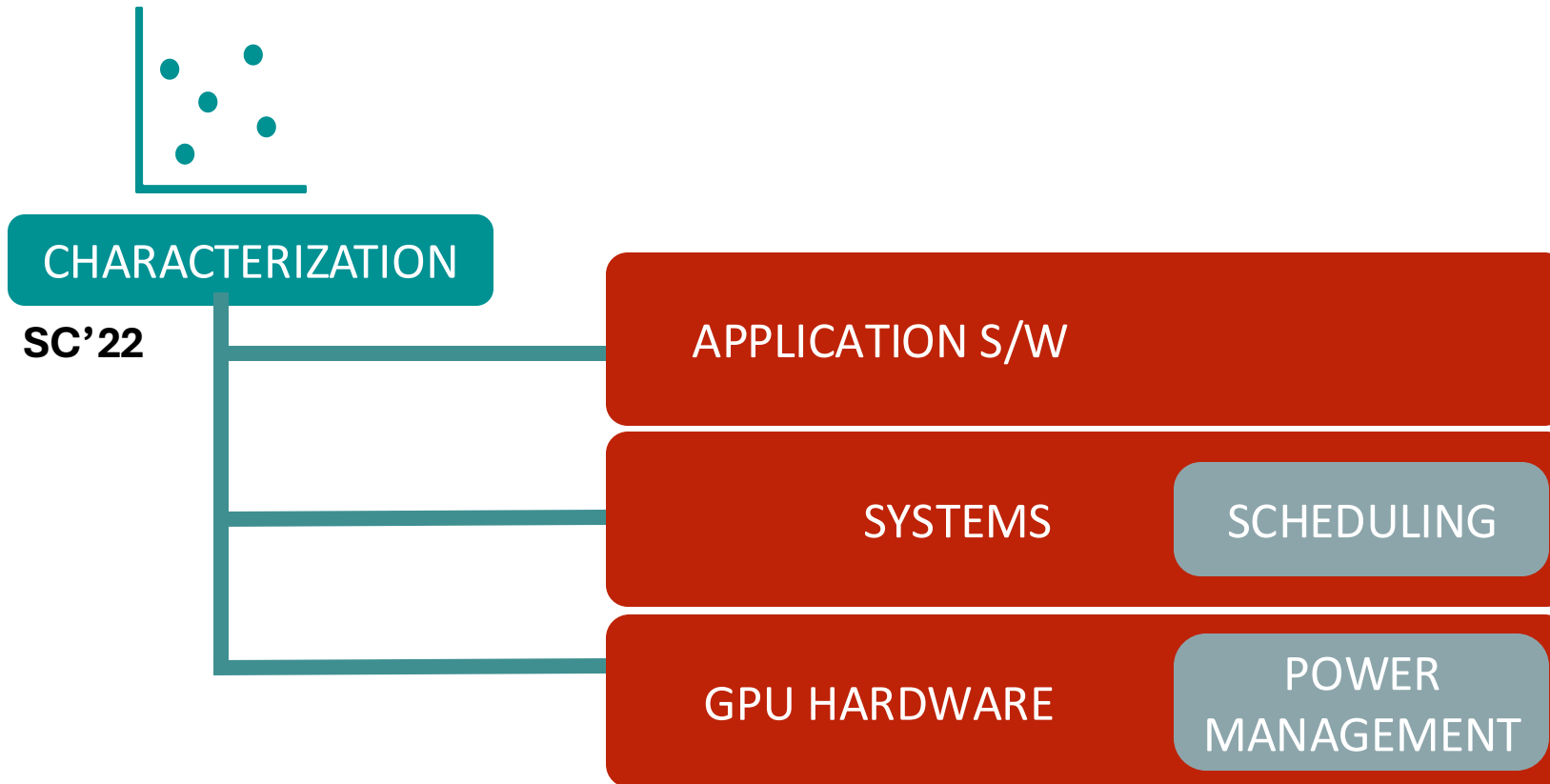
# GPU PERFORMANCE VARIABILITY

- Culprit: **Variability**!
- GPUs exhibit variability because of:
    - Local Power Management (PM)
    - Thermal/non-uniform cooling effects
    - Process variation
- **Impact**
    - Users: Hard to get **repeatable, consistent** performance for applications.
    - Sysadmins: **Resource underutilization** for multi-GPU jobs – faster GPUs wait for stragglers

# CHARACTERIZING VARIABILITY



**CHARACTERIZATION**
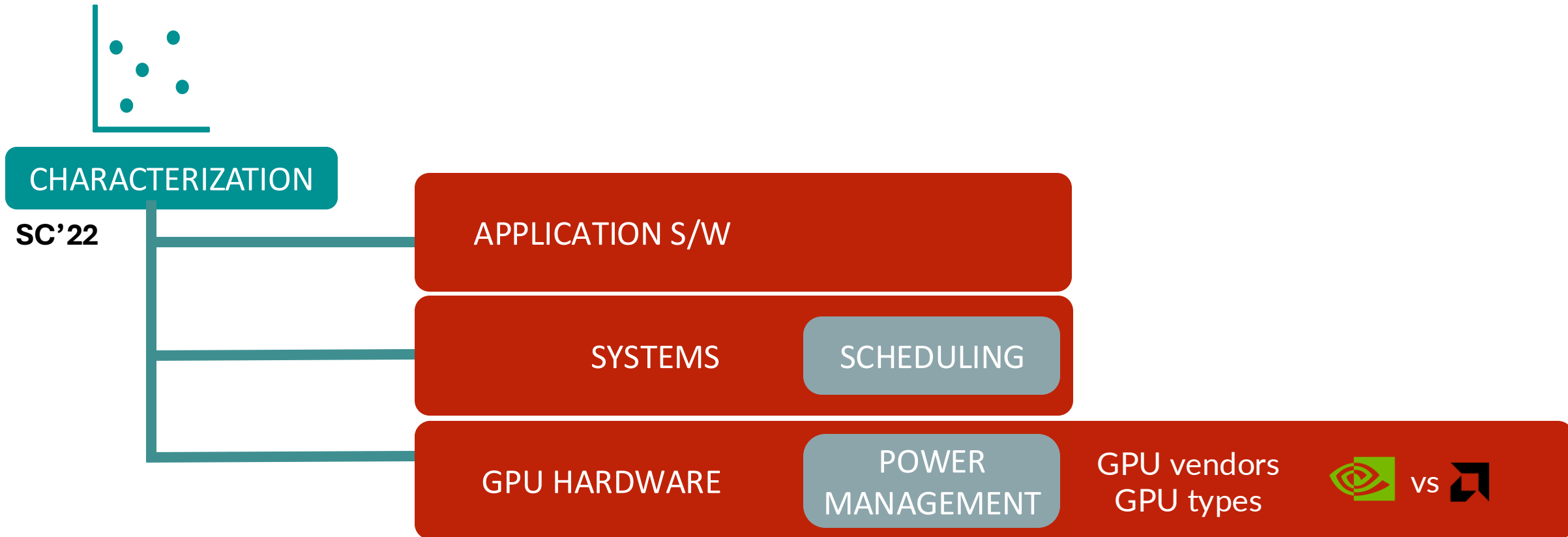
**SC'22**

APPLICATION S/W

SYSTEMS — SCHEDULING

GPU HARDWARE — POWER MANAGEMENT

Not All GPUs Are Created Equal: Characterizing Variability in Large-Scale, Accelerator-Rich Systems (SC22)
https://dl.acm.org/doi/abs/10.5555/3571885.3571971

# CHARACTERIZING VARIABILITY

**CHARACTERIZATION**

**SC'22**

**APPLICATION S/W**

**SYSTEMS** | **SCHEDULING**

**GPU HARDWARE** | **POWER MANAGEMENT** | GPU vendors GPU types | vs

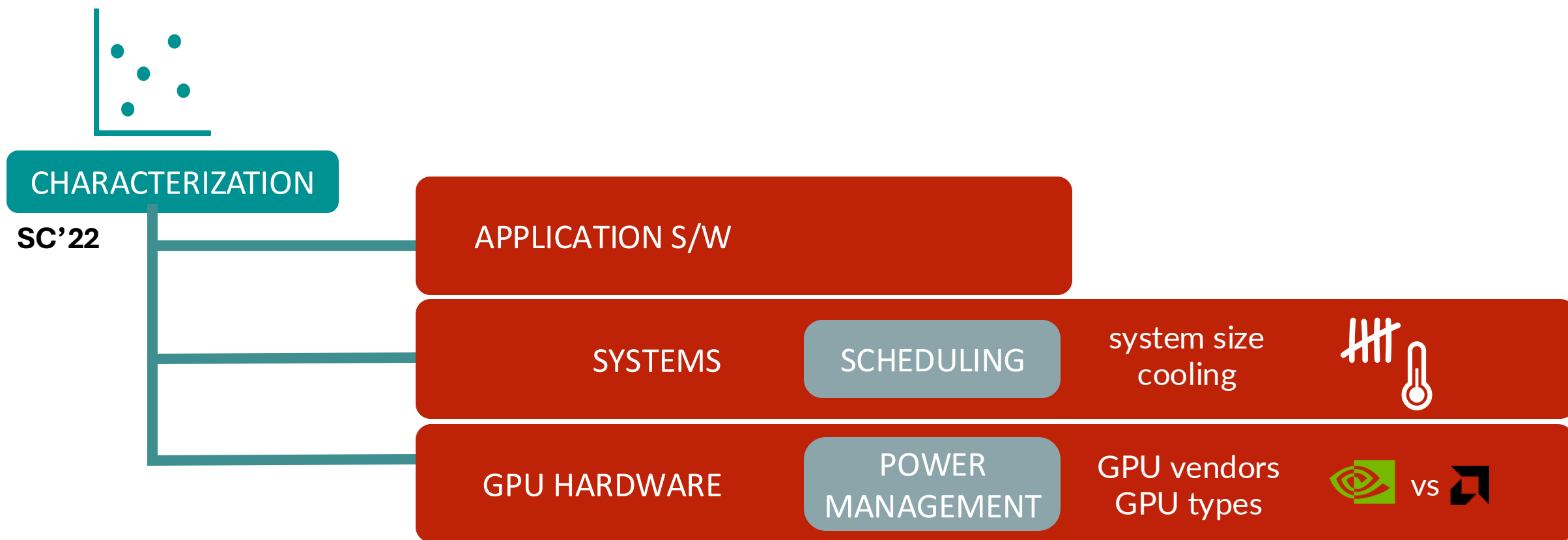Not All GPUs Are Created Equal: Characterizing Variability in Large-Scale, Accelerator-Rich Systems (SC22)
https://dl.acm.org/doi/abs/10.5555/3571885.3571971

# CHARACTERIZING VARIABILITY



CHARACTERIZATION

**SC'22**

APPLICATION S/W

SYSTEMS    SCHEDULING    system size cooling

GPU HARDWARE    POWER MANAGEMENT    GPU vendors GPU types    vs

Not All GPUs Are Created Equal: Characterizing Variability in Large-Scale, Accelerator-Rich Systems (SC22)
https://dl.acm.org/doi/abs/10.5555/3571885.3571971

# CHARACTERIZING VARIABILITY

**CHARACTERIZATION**

**SC'22**

**APPLICATION S/W** — applications

IMAGE RECOGNITION
PHYSICS SIMULATION
MACHINE TRANSLATION
GRAPH ANALYTICS
MOLECULAR DYNAMICS

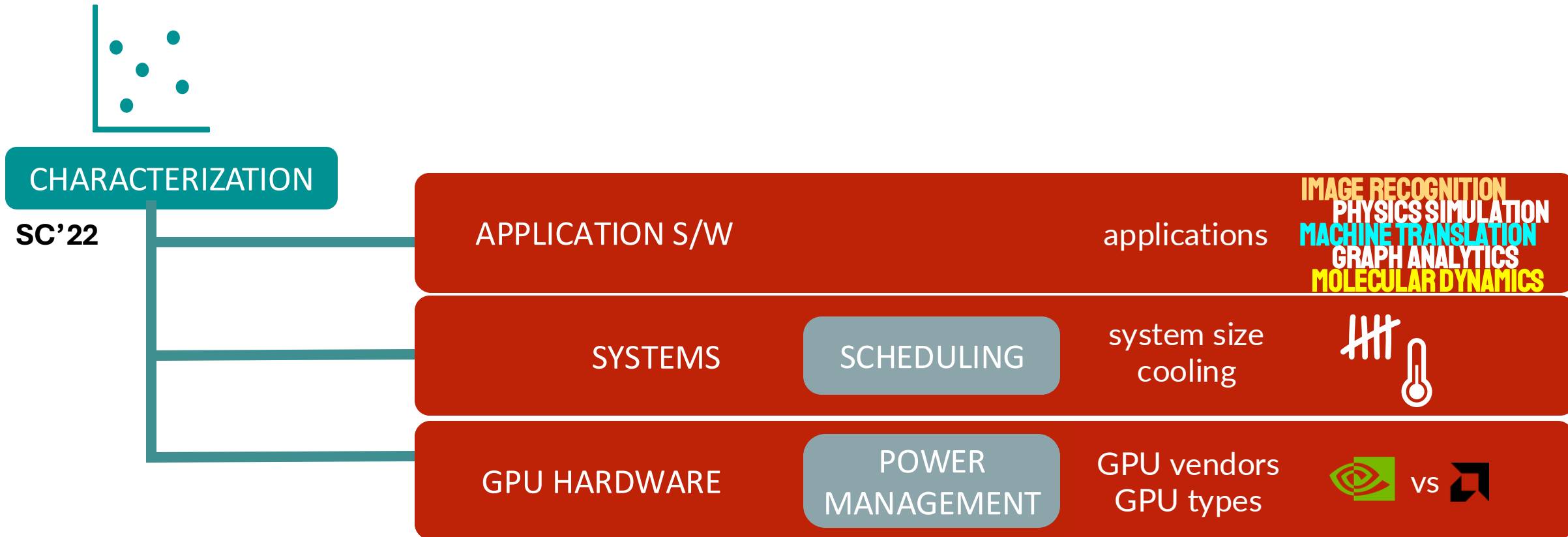**SYSTEMS** — SCHEDULING — system size cooling

**GPU HARDWARE** — POWER MANAGEMENT — GPU vendors GPU types — vs

Not All GPUs Are Created Equal: Characterizing Variability in Large-Scale, Accelerator-Rich Systems (SC22)
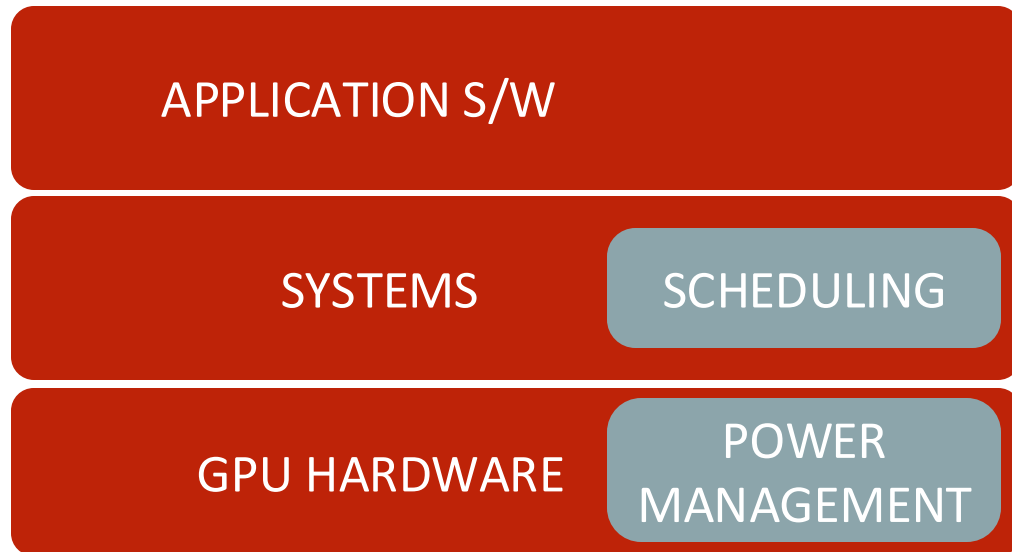https://dl.acm.org/doi/abs/10.5555/3571885.3571971

# CHARACTERIZATION STUDY: TAKEAWAYs

**1** Consistent performance variability across clusters, GPU vendors and cooling methods
(8% for SGEMM, outliers up to 1.5x slower)

**2** Variability is not a transient effect, ill-performing GPUs are consistently ill-performing

**3** **Variability is application-specific** – compute-intensive workloads are more variability sensitive
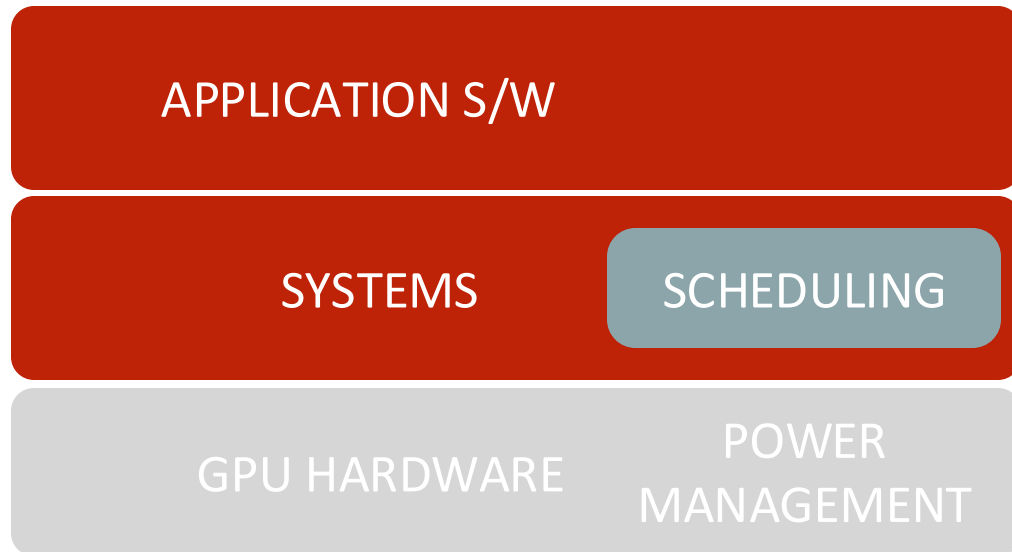(ResNet-50 shows 22% variability while PageRank has <1%)

Not All GPUs Are Created Equal: Characterizing Variability in Large-Scale, Accelerator-Rich Systems
https://dl.acm.org/doi/abs/10.5555/3571885.3571971

# KEY INSIGHT: EMBRACE VARIABILITY

APPLICATION S/W

SYSTEMS      SCHEDULING

GPU HARDWARE      POWER MANAGEMENT

Not All GPUs Are Created Equal: Characterizing Variability in Large-Scale, Accelerator-Rich Systems (SC22)
https://dl.acm.org/doi/abs/10.5555/3571885.3571971
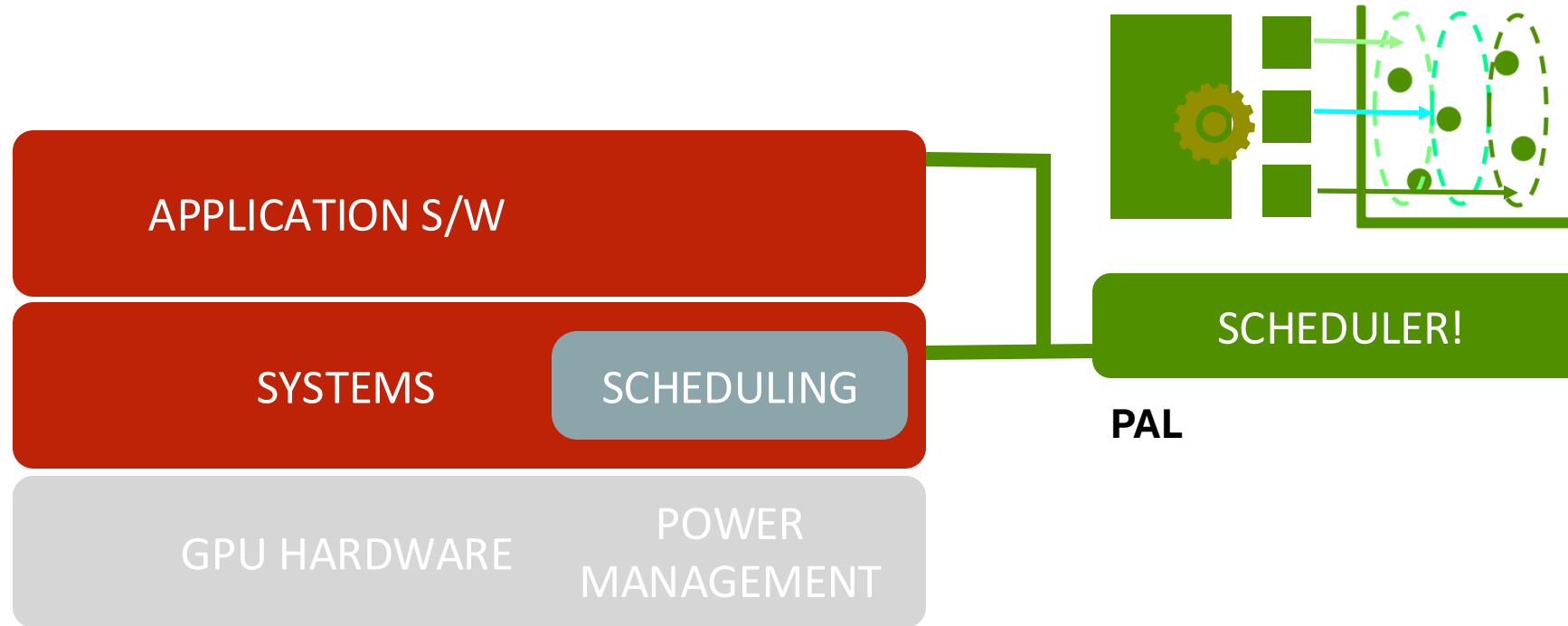
# KEY INSIGHT: EMBRACE VARIABILITY

**APPLICATION S/W**

**SYSTEMS**    **SCHEDULING**

**GPU HARDWARE**    **POWER MANAGEMENT**

Not All GPUs Are Created Equal: Characterizing Variability in Large-Scale, Accelerator-Rich Systems (SC22)
https://dl.acm.org/doi/abs/10.5555/3571885.3571971

# KEY INSIGHT: EMBRACE VARIABILITY

**APPLICATION S/W**

**SYSTEMS** | **SCHEDULING**

**GPU HARDWARE** | **POWER MANAGEMENT**

**SCHEDULER!**

**PAL**

Not All GPUs Are Created Equal: Characterizing Variability in Large-Scale, Accelerator-Rich Systems (SC22)
https://dl.acm.org/doi/abs/10.5555/3571885.3571971

# OUR APPROACH

- We characterize which applications are more likely to suffer from performance **variability** & take that into account **while placing jobs on the cluster.**

- Novel placement policies: **PM-First** and **PAL**
  - **PM-First** uses **application-specific variability profiles** to improve performance and utilization.
  - **PAL** further improves scheduling by balancing **variability** with **locality**.

- Overall, PAL improves geomean Job Completion Time (JCT) by **42%** and cluster utilization by **28%** over state-of-the-art schedulers like Tiresias
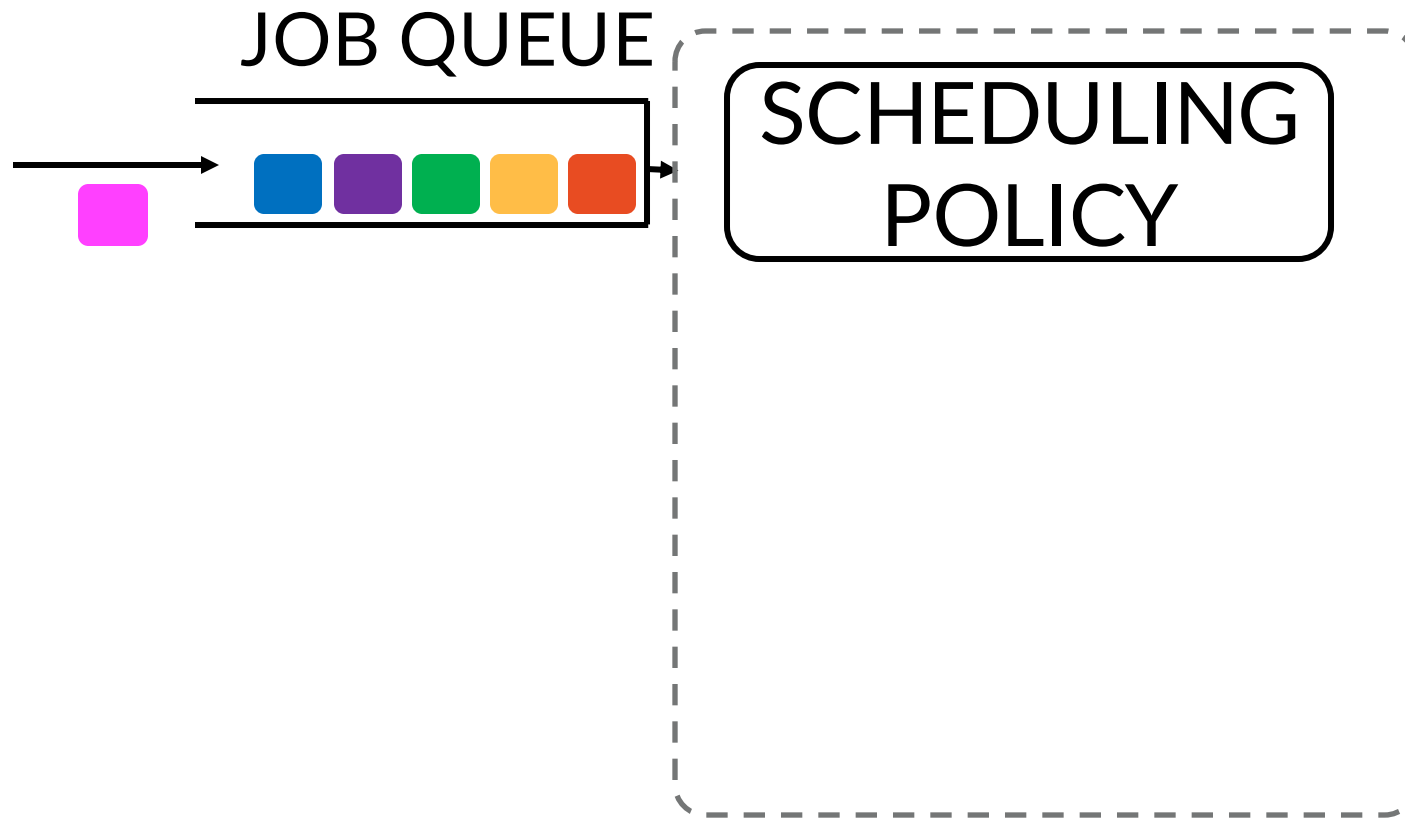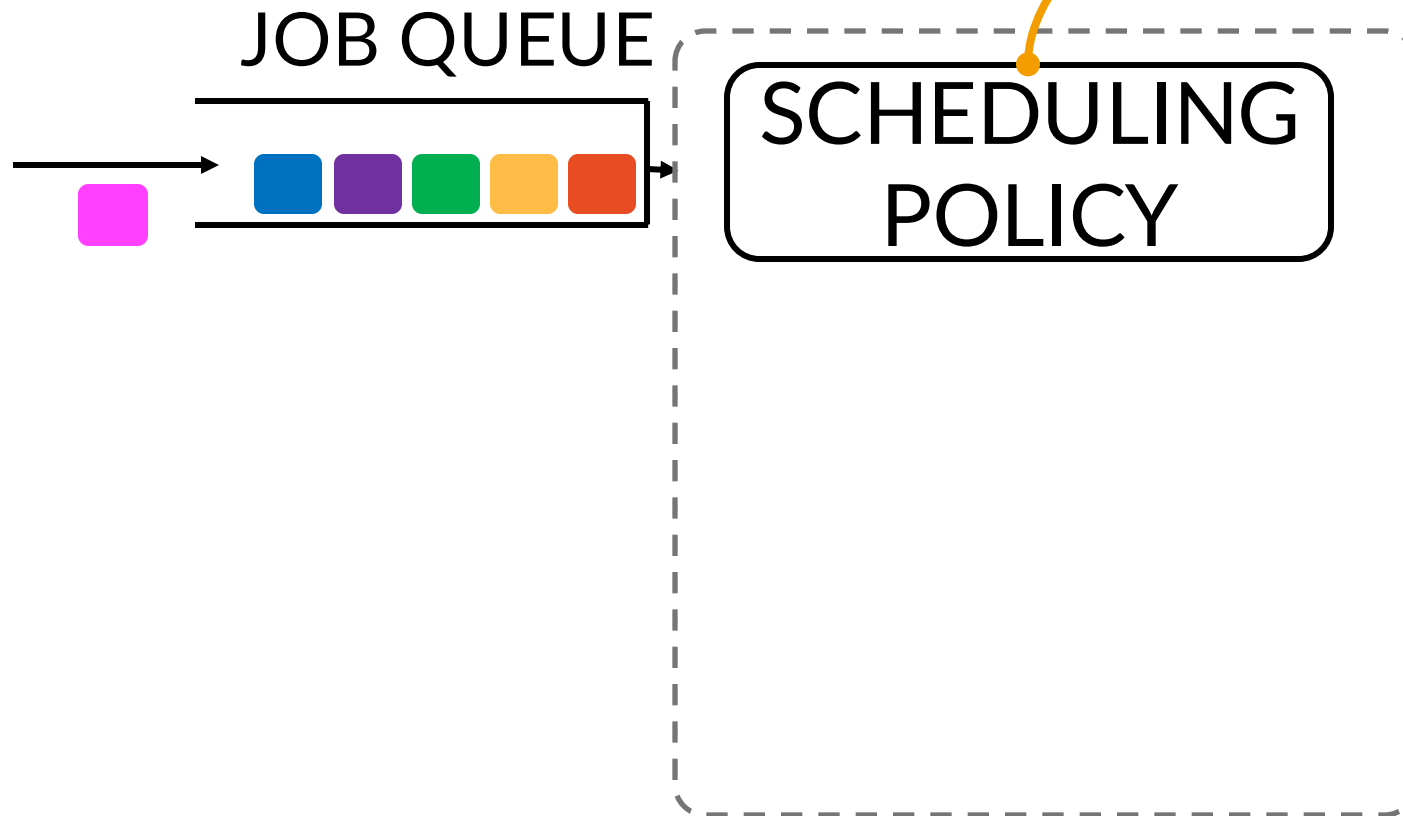
# OUTLINE

- Introduction
- **Cluster Scheduling**
- **Design**
- **Methodology**
- **Evaluation**
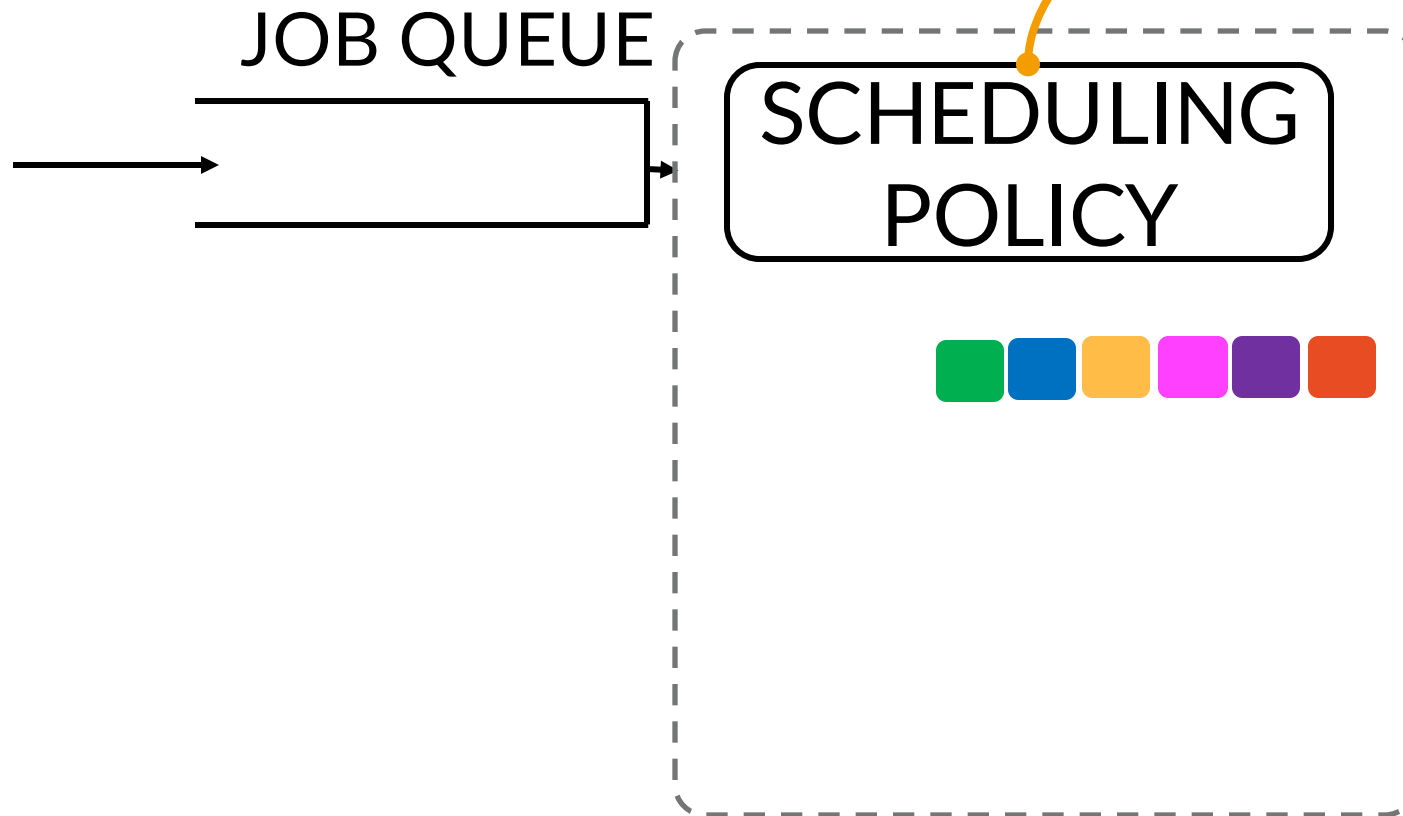- **Conclusion**

# CLUSTER SCHEDULING

JOB QUEUE



SCHEDULING POLICY

# CLUSTER SCHEDULING

## JOB QUEUE

First In First Out (FIFO)
Least Attained Service (LAS)
Shortest Job First (SJF)

SCHEDULING POLICY

CLUSTER SCHEDULING

JOB QUEUE

SCHEDULING POLICY

First In First Out (FIFO)
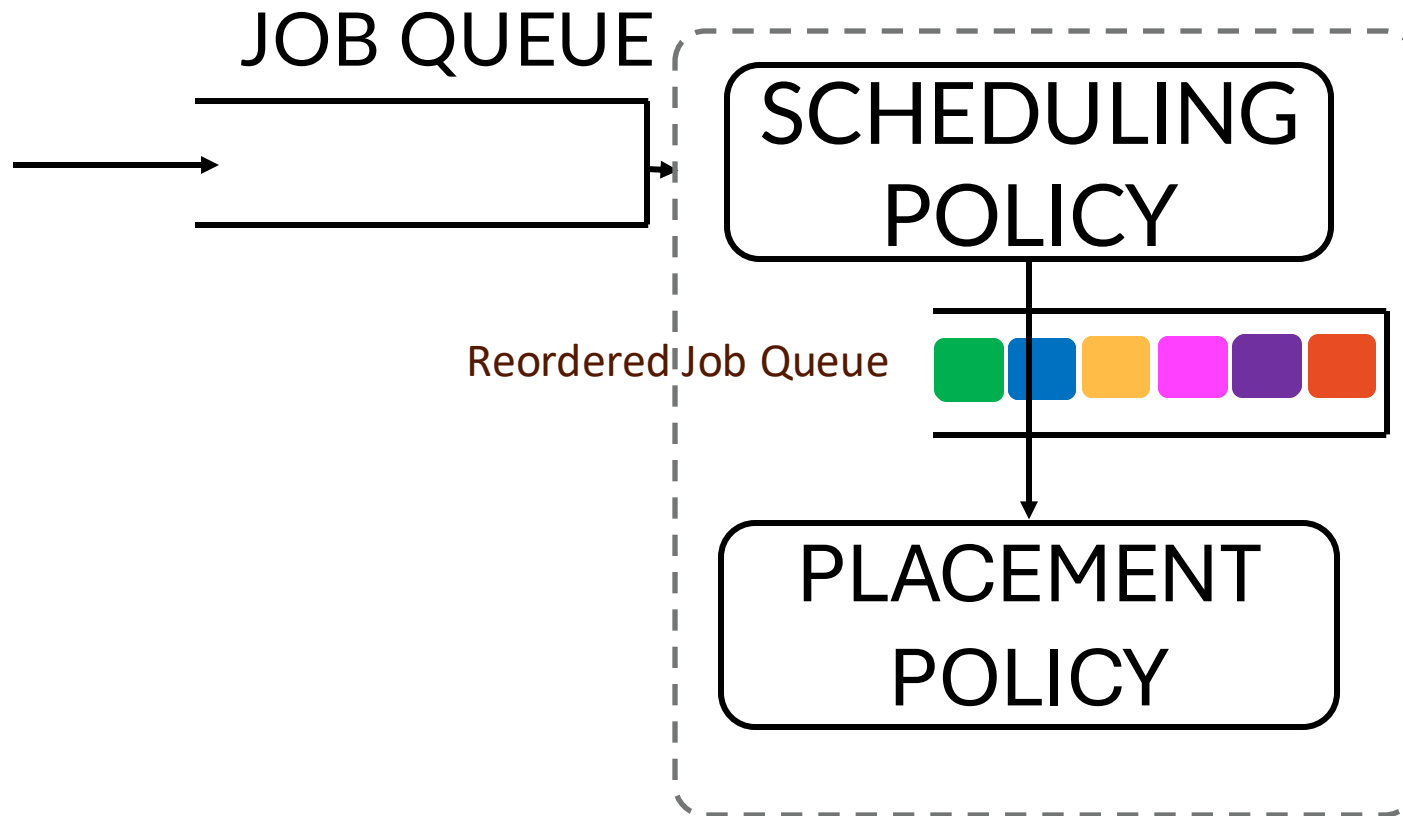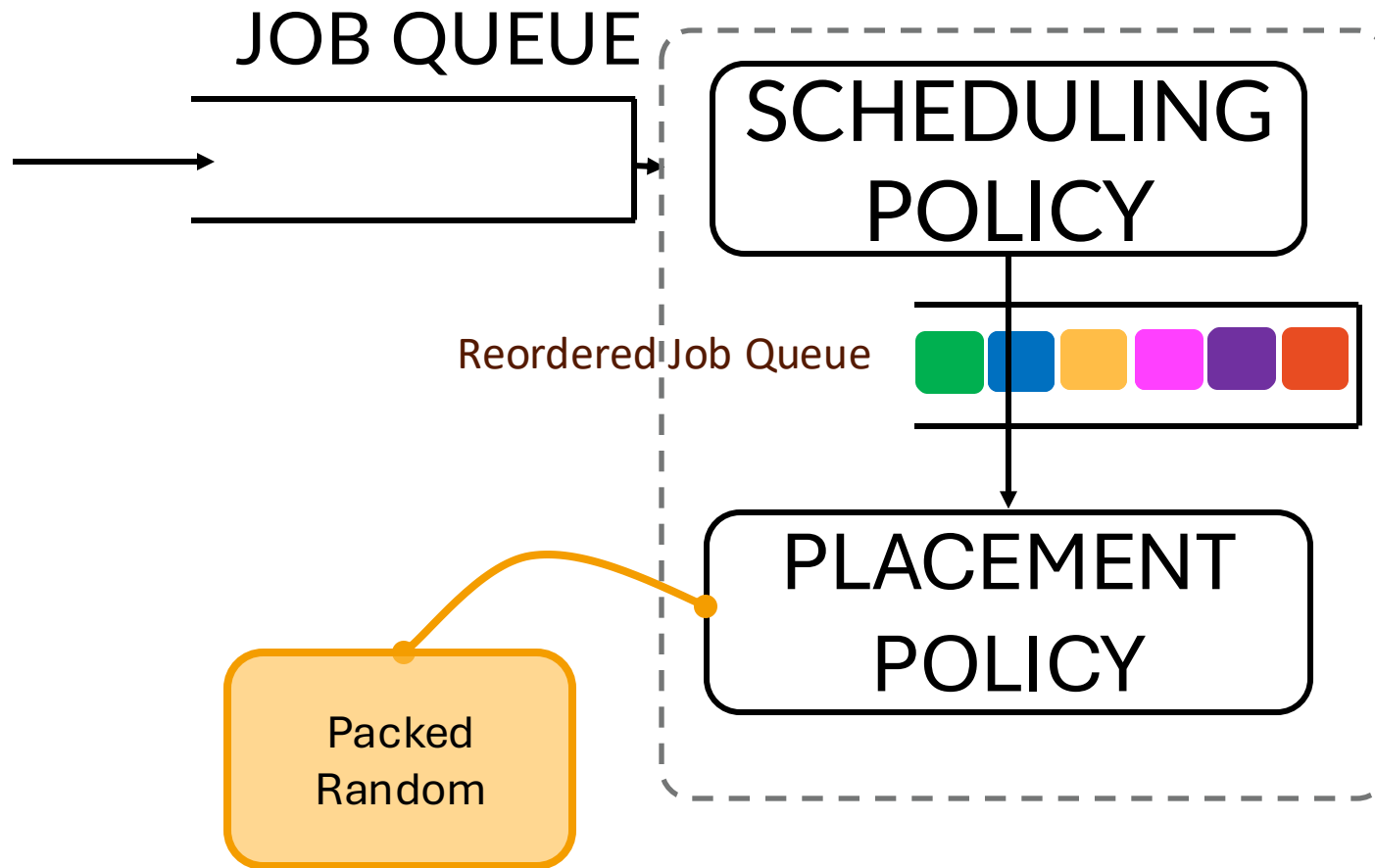Least Attained Service (LAS)
Shortest Job First (SJF)

Reordered Job Queue

# CLUSTER SCHEDULING



JOB QUEUE

SCHEDULING POLICY

Reordered Job Queue

PLACEMENT POLICY

# CLUSTER SCHEDULING

JOB QUEUE

SCHEDULING POLICY

Reordered Job Queue

PLACEMENT POLICY

Packed Random

# CLUSTER SCHEDULING

JOB QUEUE

SCHEDULING POLICY

Reordered Job Queue

PLACEMENT POLICY

CLUSTER

GPU allocation

# CLUSTER SCHEDULING



JOB QUEUE

SCHEDULING POLICY

Reordered Job Queue

PLACEMENT POLICY

VARIABILITY AGNOSTIC

CLUSTER

GPU allocation

# VARIABILITY-INFORMED PLACEMENT

# VARIABILITY-INFORMED PLACEMENT



application-specific **profiling** to generate variability data

JOB QUEUE

SCHEDULING POLICY

Reordered Job Queue

PLACEMENT POLICY

CLUSTER

**Problem: Variability is application-specific**, and systems run a wide variety of workloads

# VARIABILITY-INFORMED PLACEMENT



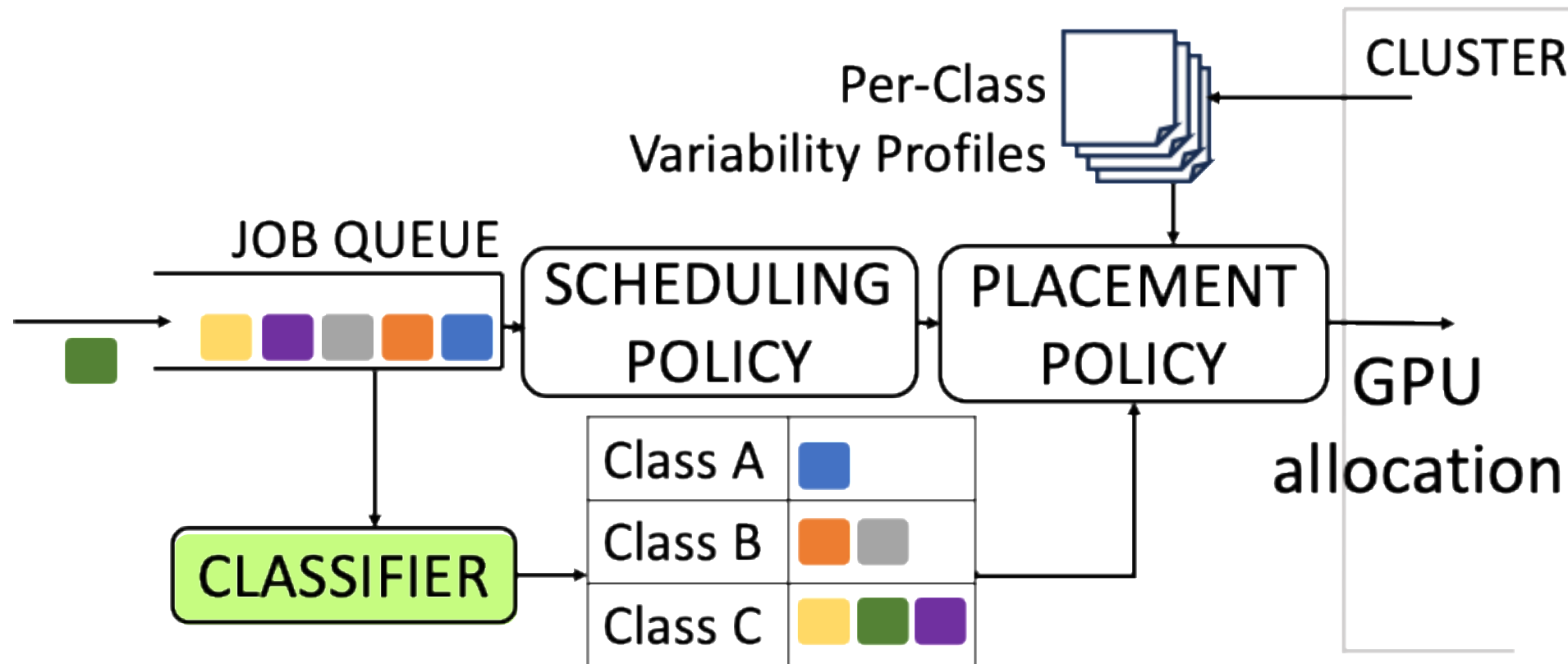**Solution**: have a finite number of classes
(compute bound ↔ memory bound)

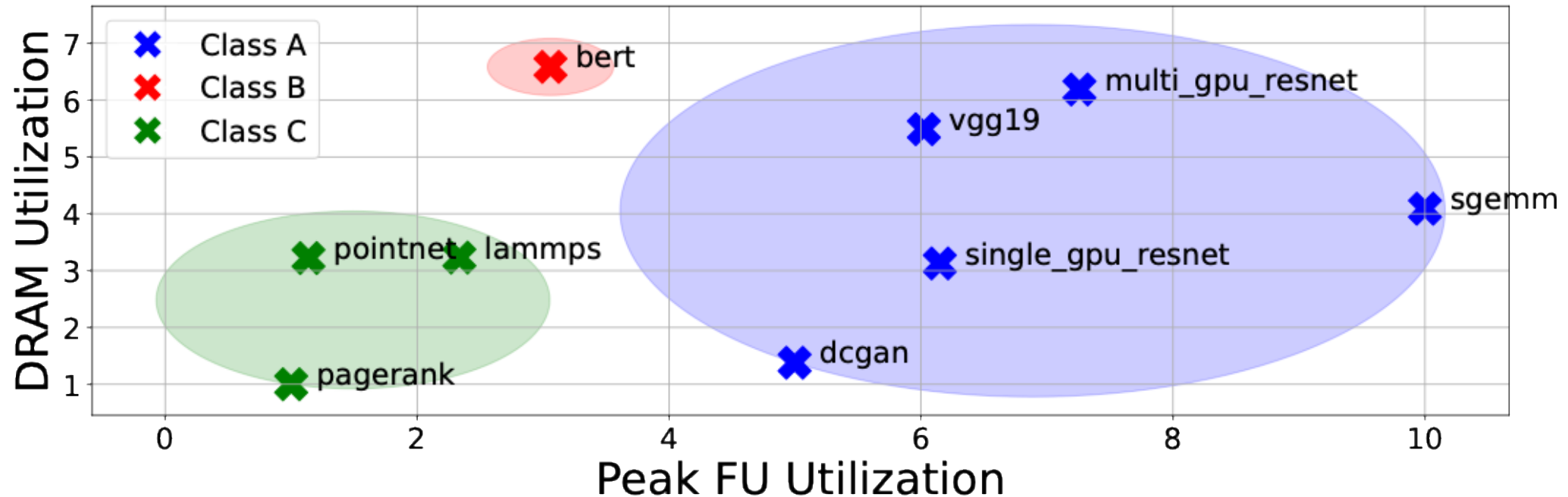# VARIABILITY-INFORMED PLACEMENT

use one variability profile per "class" as proxy for variability behavior of all applications in that class



**Solution**: have a finite number of classes (compute bound ↔ memory bound)

# CLASSIFICATION LAYER



- Similar to Guerreiro et al. [Parallel Computing 2019], we used nsight compute to measure workloads' DRAM Throughput and Peak SM Throughput.
- **2D K-Means clustering** to produce ordered classes.
- Any new application can fall into one of these classes based on its DRAM and SM throughput.

# SOLUTION 1: PM-FIRST PLACEMENT

# SOLUTION 1: PM-FIRST PLACEMENT

- Give variability first-order precedence when assigning GPUs to jobs.

# SOLUTION 1: PM-FIRST PLACEMENT

- Give variability first-order precedence when assigning GPUs to jobs.

- We associate a **PM-score** with each GPU
  - How fast/slow the GPU is relative to median GPU in the cluster (normalized performance)

$$P_i = \frac{t_i}{t_{median}}$$

  - PM-score is class-specific.

# SOLUTION 1: PM-FIRST PLACEMENT

- Give variability first-order precedence when assigning GPUs to jobs.

- We associate a **PM-score** with each GPU
  - How fast/slow the GPU is relative to median GPU in the cluster (normalized performance)

$$P_i = \frac{t_i}{t_{median}}$$

  - PM-score is class-specific.

- We also associate a **placement priority** for each job
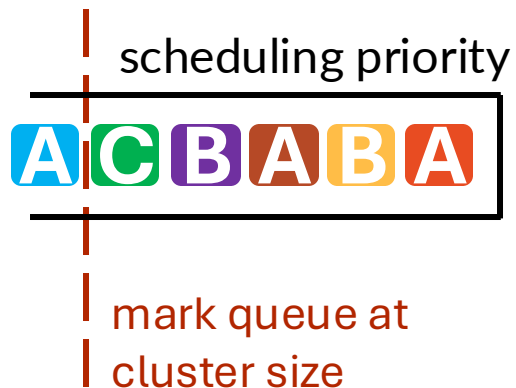
scheduling priority

# SOLUTION 1: PM-FIRST PLACEMENT

- Give variability first-order precedence when assigning GPUs to jobs.

- We associate a **PM-score** with each GPU
  - How fast/slow the GPU is relative to median GPU in the cluster (normalized performance)

$$P_i = \frac{t_i}{t_{median}}$$

  - PM-score is class-specific.

- We also associate a **placement priority** for each job



scheduling priority
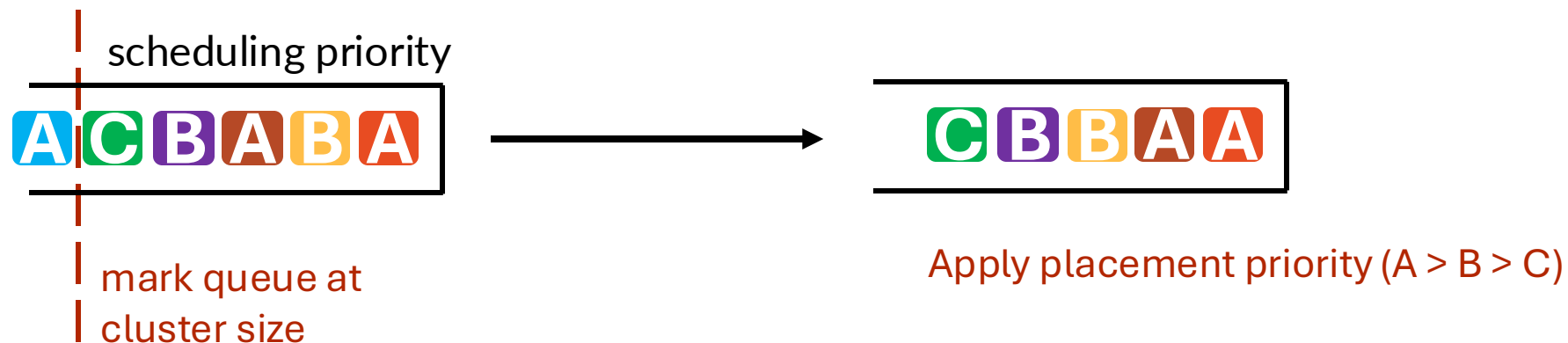
A C B A B A

mark queue at
cluster size

# SOLUTION 1: PM-FIRST PLACEMENT

- Give variability first-order precedence when assigning GPUs to jobs.

- We associate a **PM-score** with each GPU
  - How fast/slow the GPU is relative to median GPU in the cluster (normalized performance)
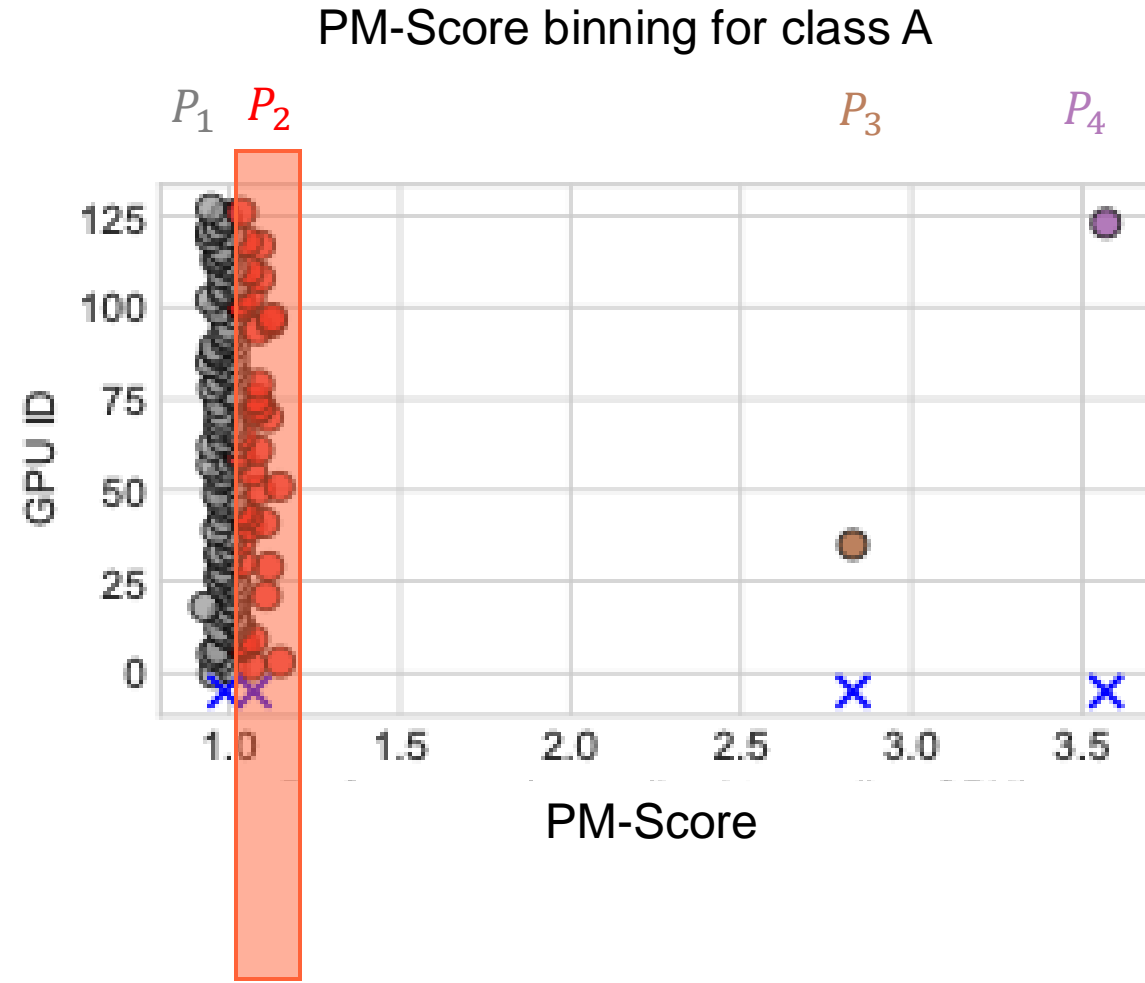
$$P_i = \frac{t_i}{t_{median}}$$

  - PM-score is class-specific.

- We also associate a **placement priority** for each job

scheduling priority



mark queue at
cluster size

Apply placement priority (A > B > C)

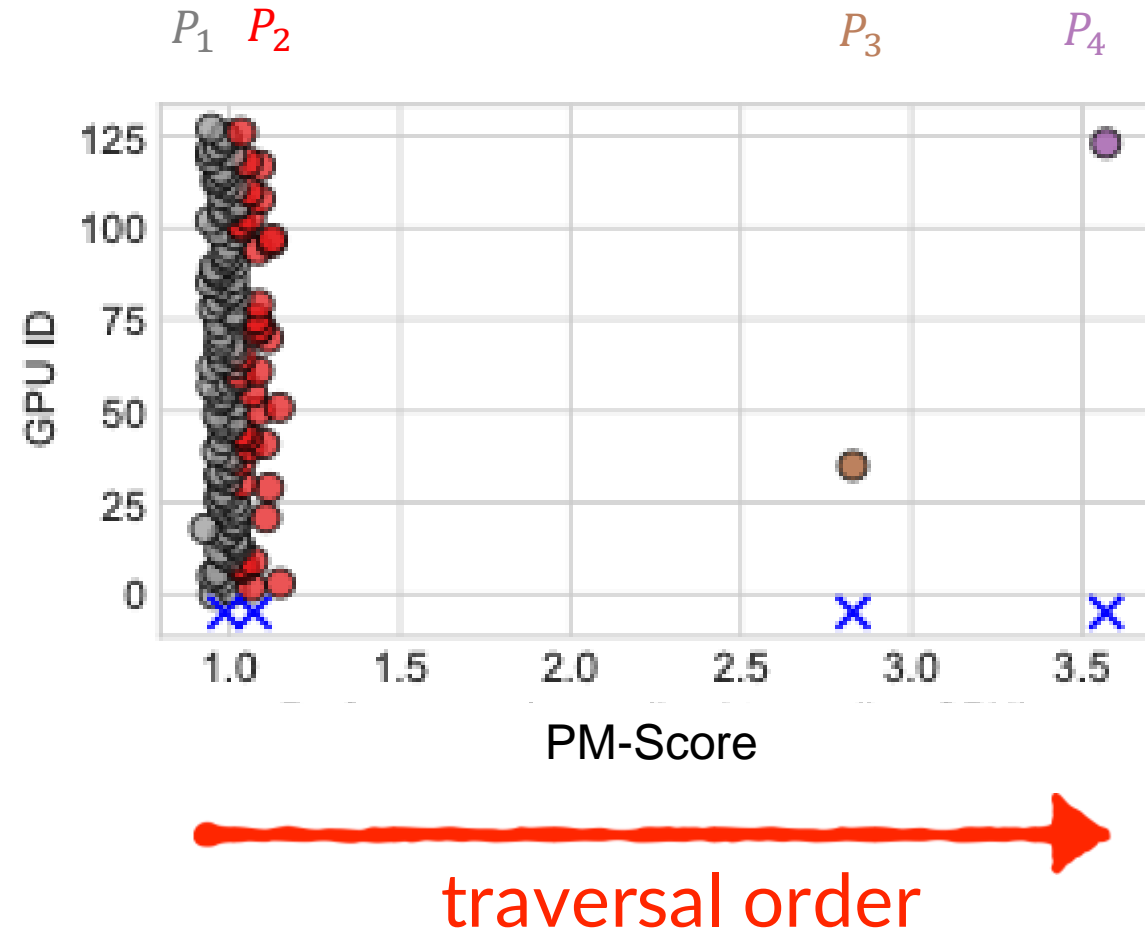# SOLUTION 1: PM-FIRST PLACEMENT

PM-Score binning for class A

- For each class, we make **bins** of GPUs with similar **PM-Scores**.
- Why?
  - Fine-grained variability for large-scale systems is expensive
  - Memory-bound applications have little variability – one large pool of GPUs for them.
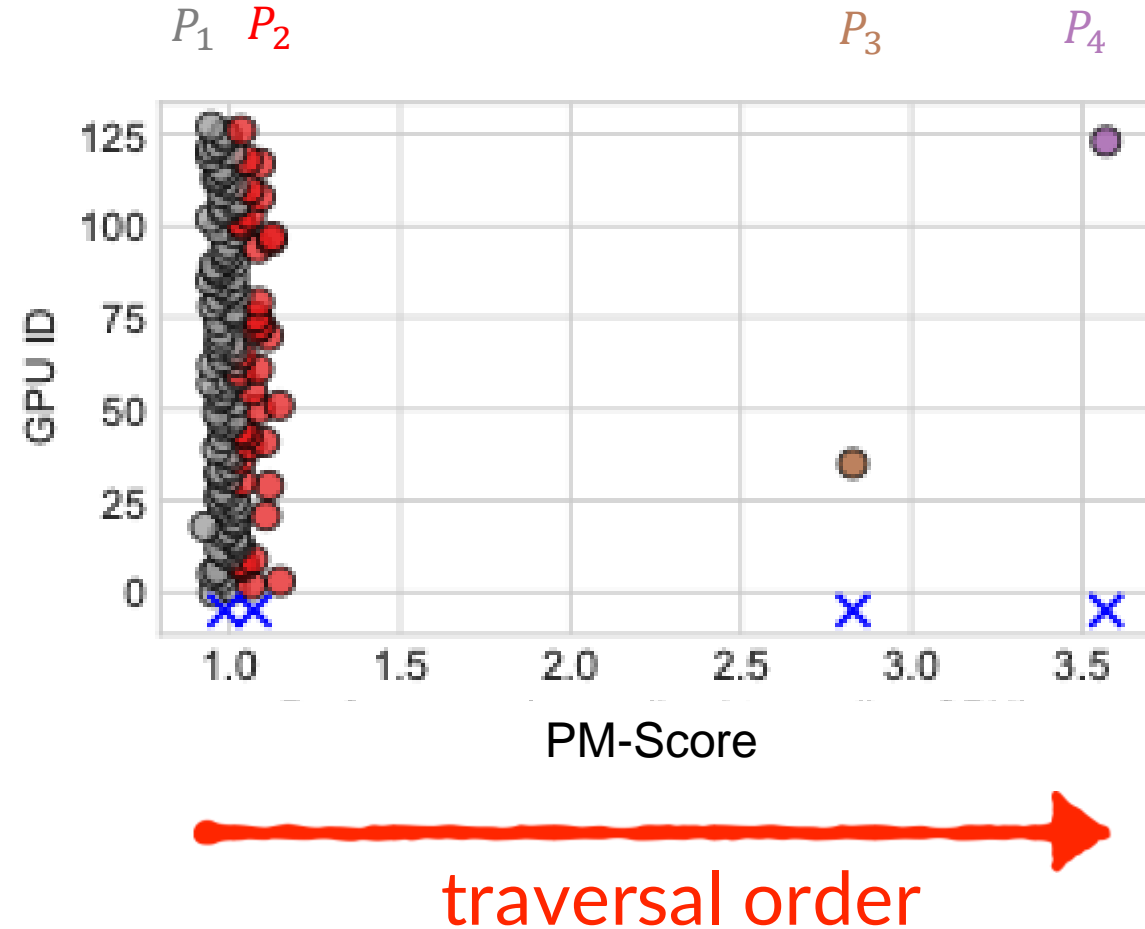


all GPUs of the red bin have similar variability

# PM-FIRST PLACEMENT POLICY

# PM-FIRST PLACEMENT POLICY

In each scheduling round,
- Get job queue from scheduling policy
- Mark queue at cluster size, and assign placement priority



traversal order
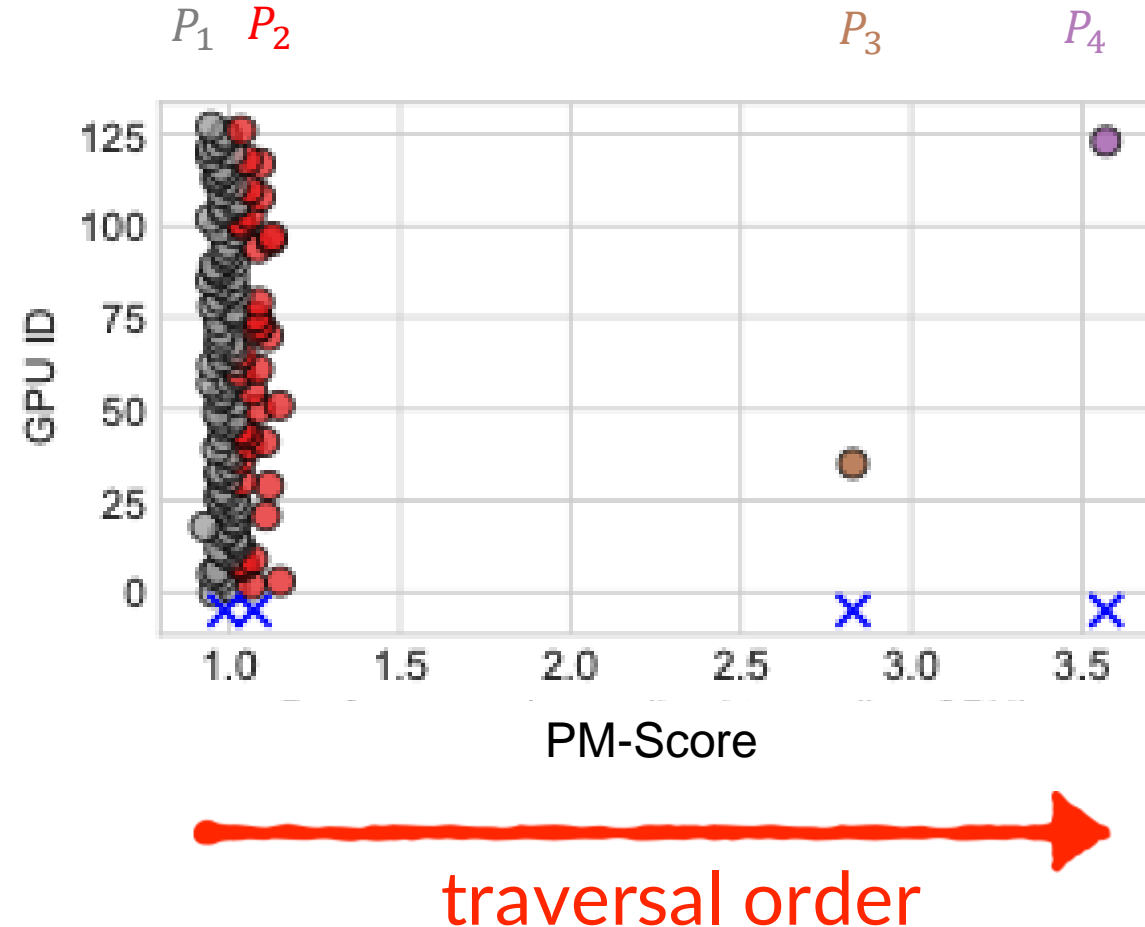
# PM-FIRST PLACEMENT POLICY

In each scheduling round,
- Get job queue from scheduling policy
- Mark queue at cluster size, and assign placement priority

for each job $j$ in order of placement priority:
- get job class and GPU demand $N$
- get freelist of GPUs available and **sort** by their **PM-Score** value (for job class)
- allocate first $N$ GPUs from this sorted list.
- Mark GPUs in use (remove from freelist)



traversal order

# WHAT ABOUT LOCALITY?

- PM-First ignores communication overheads due to ineffective packing
- Some jobs are more sensitive to locality of GPUs than variability

# WHAT ABOUT LOCALITY?

- PM-First ignores communication overheads due to ineffective packing

- Some jobs are more sensitive to locality of GPUs than variability

- Assume a simplified locality model:
  - Multi-GPU job incurs penalty $\boldsymbol{L_{across}}$ if it's allocation spills across nodes
  - No performance degradation if allocation is within a node

- Modified iteration time for a job running with set of GPUs $G$ spread across > 1 node:

$$t_{iter} = L_{across} \times max_{g \in G}(P_g) \times t_{iter}^{orig}$$

where $P_g$ is the PM-Score of GPU $g$

# WHAT ABOUT LOCALITY?

- PM-First ignores communication overheads due to ineffective packing

- Some jobs are more sensitive to locality of GPUs than variability

- Assume a simplified locality model:
  - Multi-GPU job incurs penalty $\boldsymbol{L_{across}}$ if it's allocation spills across nodes
  - No performance degradation if allocation is within a node

- Modified iteration time for a job running with set of GPUs $G$ spread across > 1 node:

$$t_{iter} = L_{across} \times max_{g \in G}(P_g) \times t_{iter}^{orig}$$

where $P_g$ is the PM-Score of GPU $g$

$$t_{iter} = (L_{across} \times V) \times t_{iter}^{orig}$$

# LV-PRODUCT

$$\min LV{-}\text{product} = \min\left(L_{across} \times V\right)$$

Slowdown due to locality          Slowdown due to variability

Lower V values are better

| Example ClassA | $V_1 = 0.89$ | $V_2 = 0.94$ | $V_3 = 1.06$ | $V_4 = 2.55$ |
|---|---|---|---|---|
| Within Node $L_{within} = 1$ | 0.89 | 0.94 | 1.06 | 2.55 |
| Across node $L_{across} = 1.5$ | 1.35 | 1.41 | 1.59 | 3.82 |

# LV-PRODUCT

$$\min LV{-}\text{product} = \min\left(L_{across} \times V\right)$$

Slowdown due to locality

Slowdown due to variability

Lower V values are better

| Example ClassA | $V_1 = 0.89$ | $V_2 = 0.94$ | $V_3 = 1.06$ | $V_4 = 2.55$ |
|---|---|---|---|---|
| Within Node $L_{within} = 1$ | 0.89 | 0.94 | 1.06 | 2.55 |
| Across node $L_{across} = 1.5$ | 1.35 | 1.41 | 1.59 | 3.82 |

PAL traversal order

# LV-PRODUCT

$$\min LV-\text{product} = \min (L_{across} \times V)$$

Slowdown due to locality

Slowdown due to variability

| Example ClassA | $V_1 = 0.89$ | $V_2 = 0.94$ | $V_3 = 1.06$ | $V_4 = 2.55$ |
|---|---|---|---|---|
| Within Node $L_{within} = 1$ | 0.89 | 0.94 | 1.06 | 2.55 |
| Across node $L_{across} = 1.5$ | 1.35 | 1.41 | 1.59 | 3.82 |

PM-First traversal order

# OUTLINE

- Introduction
- Cluster Scheduling
- Design
- **Methodology**
- Evaluation
- Conclusion

# SYSTEM

- **Blox (open-source modular toolkit)**
- **Physical Cluster**

  TACC Frontera
  **360** GPUs
  mineral oil cooled
  NVIDIA Quadro **RTX5000** GPUs
  4x GPUs per node
  16GB memory per GPU

- **Testbed**
  - 16 node (64 GPUs)

# SYSTEM

- **Blox (open-source modular toolkit)**
- **Physical Cluster**

  TACC Frontera

  **360** GPUs

  mineral oil cooled

  NVIDIA Quadro **RTX5000** GPUs

  4x GPUs per node

  16GB memory per GPU

- **Testbed**

  - 16 node (64 GPUs)

| Experiment(s) | Cluster Size | Workload Trace | Paper Section |
|---|---|---|---|
| Cluster Testbed Evaluation | 64 | Sia-Cluster* | V-A |
| Baseline Simulation Varying Locality Penalty | 64 | Sia-Philly* | V-B |
| Varying Job Load Varying Scheduling Policy | 256 | Synergy† | V-C |

\*  [Subramanya et al. SOSP'23]
† [Mohan et al.  OSDI'22 ]

# BASELINE PLACEMENT POLICIES

# BASELINE PLACEMENT POLICIES

Locality

Packed ✖        Random ✖

# BASELINE PLACEMENT POLICIES

Sticky

Non-Sticky

Reallocation

# METHODOLOGY: PLACEMENT POLICIES

Locality

Packed-Sticky     Random-Sticky

Packed Non-Sticky     Random-Non-Sticky

Reallocation

# METHODOLOGY: PLACEMENT POLICIES

Locality →

Packed-Sticky     Random-Sticky

❌       ❌

Packed Non-Sticky    Random-Non-Sticky

❌       ❌

↓ Reallocation

Baseline placement policies

- **Tiresias** [Gu et al. NSDI'19]
  - Packed Sticky
- **Gandiva** [Xia et al. OSDI'18]
  - Packed Non-sticky
- **Random-Sticky**
- **Random-Non-Sticky**

HotGauge
[Hankin et al. IISWC'21]

# METHODOLOGY: PLACEMENT POLICIES

Locality →

Packed-Sticky     Random-Sticky

Packed Non-Sticky     Random-Non-Sticky

Reallocation

Baseline placement policies

- **Tiresias** [Gu et al. NSDI'19]
  - Packed Sticky

- **Gandiva** [Xia et al. OSDI'18]
  - Packed Non-sticky

- **Random-Sticky**

- **Random-Non-Sticky**    HotGauge [Hankin et al. IISWC'21]

- **PM-First and PAL**
  - Both non-sticky

# APPLICATIONS: PHYSICAL CLUSTER EXPERIMENTS

| Task | Model | Dataset | Batch Size | Class |
|------|-------|---------|------------|-------|
| Image | PointNet | ShapeNet | 32 | C |
| Image | vgg19 | ImageNet2012 | 32 | A |
| Vision | DCGAN | LSUN | 128 | A |
| Language | BERT | WikiText | 64 | B |
| Image | ResNet-50 | Imagenet2012 | 32 | A |
| Language | GPT2 | Wikitext | 128 | B |

# PM-SCORES: VARIABILITY PROFILES

Class A          Class B          Class C



TACC Frontera

$$P_i = \frac{t_i}{t_{median}}$$

# LOCALITY PENALTY



8 GPU ResNet-50 job
Batch Size 128



4 GPU ResNet-50 job
Batch Size 64

# LOCALITY PENALTY

8 GPU ResNet-50 job
Batch Size 128 $\longrightarrow$ $t_{iter,avg}$

$$\frac{t_{iter,avg}}{t_{iter,avg}} = L_{across}$$

4 GPU ResNet-50 job
Batch Size 64 $\longrightarrow$ $t_{iter,avg}$

# OUTLINE

- Introduction
- Cluster Scheduling
- Design
- Methodology
- **Evaluation**
- Conclusion

# PHYSICAL CLUSTER EXPERIMENT: 64 GPUs

# PHYSICAL CLUSTER EXPERIMENT: 64 GPUs

# PHYSICAL CLUSTER EXPERIMENT: 64 GPUs

# PHYSICAL CLUSTER EXPERIMENT: 64 GPUs



First-order trends are accurate and correlate well in simulation

# SIA-PHILLY SIMULATIONS

Scheduler: **FIFO**
# of GPUs: **64**
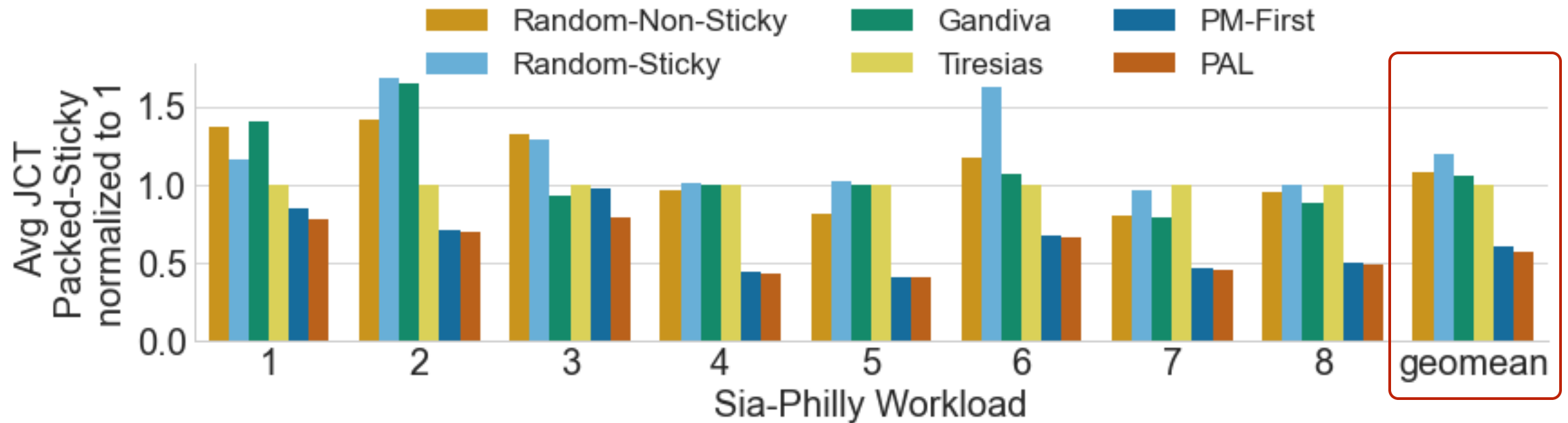Job trace: **Sia-Philly trace (160 jobs, requesting up to 48 GPUs)**

# SIA-PHILLY SIMULATIONS

Scheduler: **FIFO**

# of GPUs: **64**

Job trace: **Sia-Philly trace (160 jobs, requesting up to 48 GPUs)**



PM-First improves geomean JCT by 40%, PAL further improves JCTs by considering locality in addition to variability
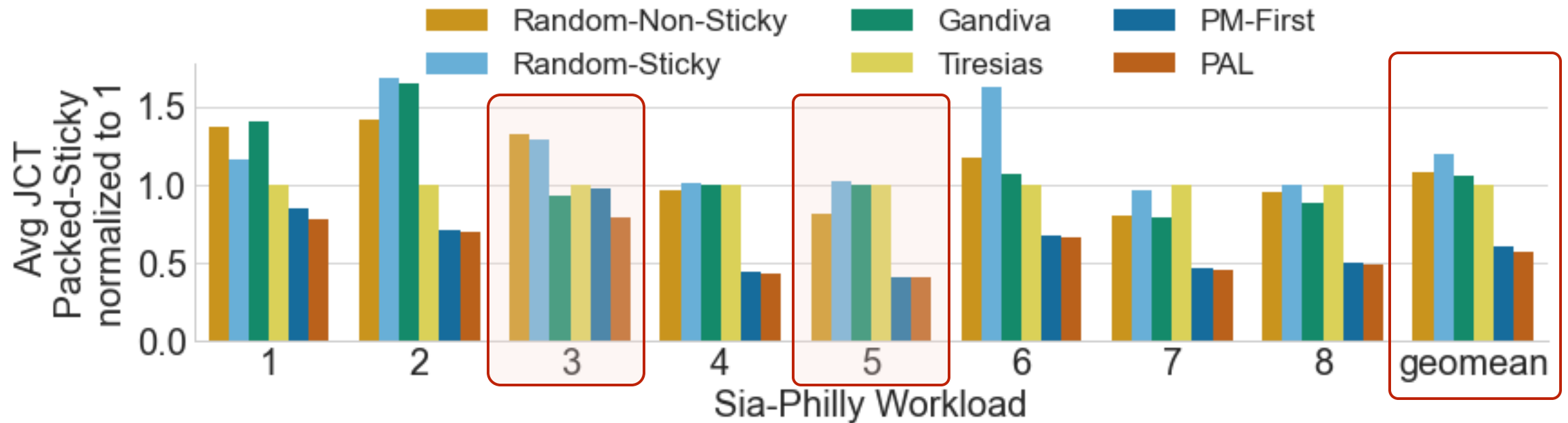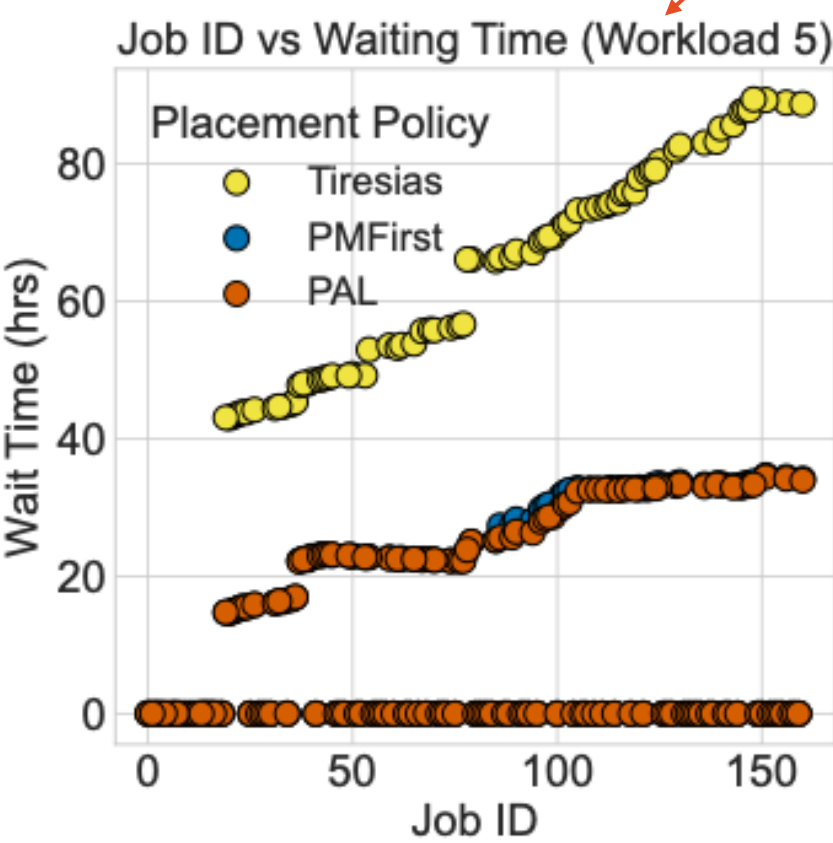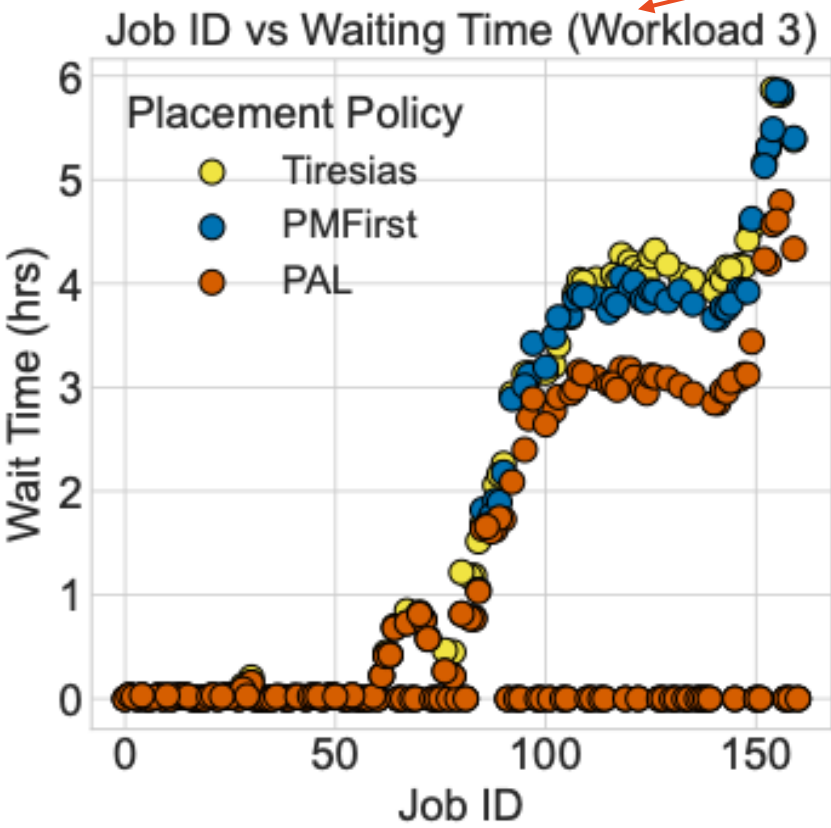
# SIA-PHILLY SIMULATIONS

Scheduler: **FIFO**
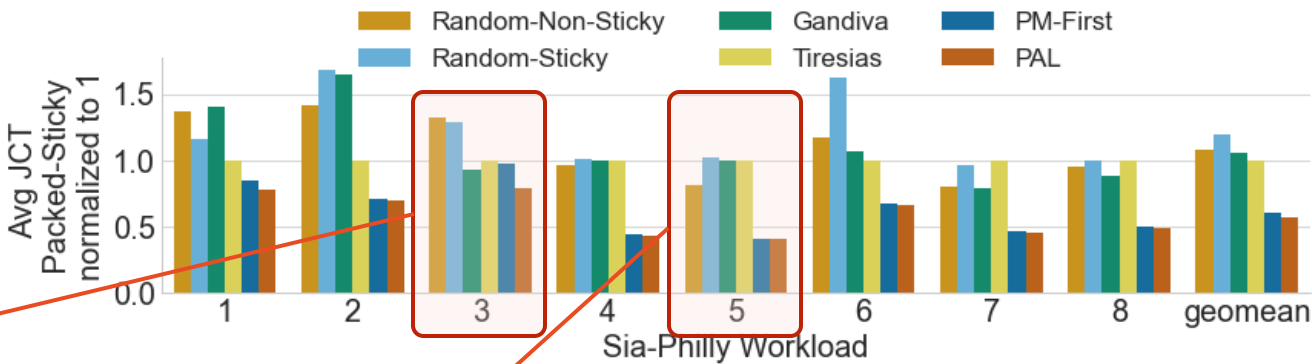
# of GPUs: **64**

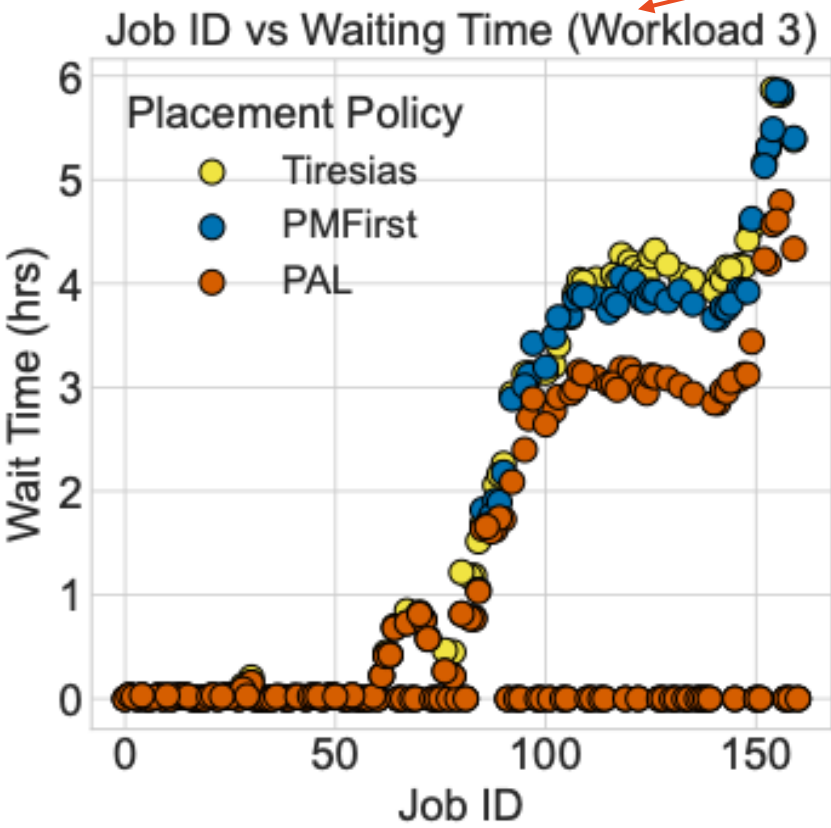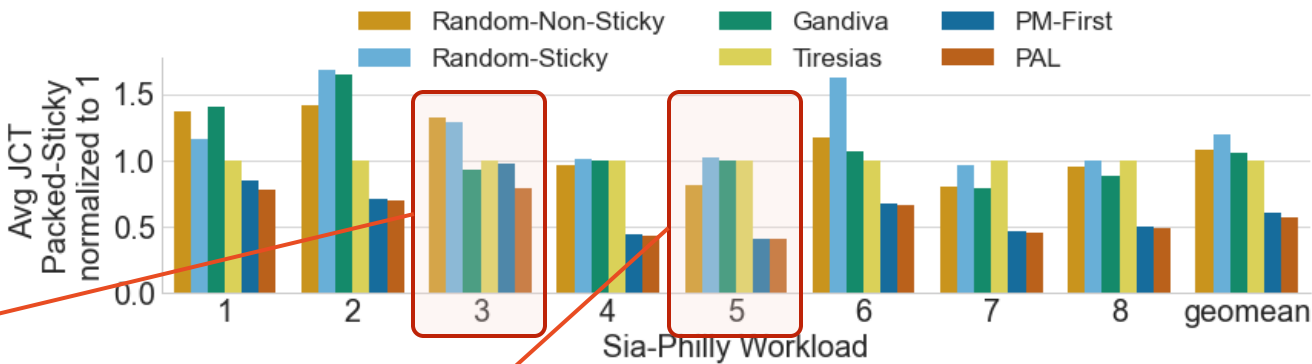Job trace: **Sia-Philly trace (160 jobs, requesting up to 48 GPUs)**



PM-First improves geomean JCT by 40%, PAL further improves JCTs by considering locality in addition to variability
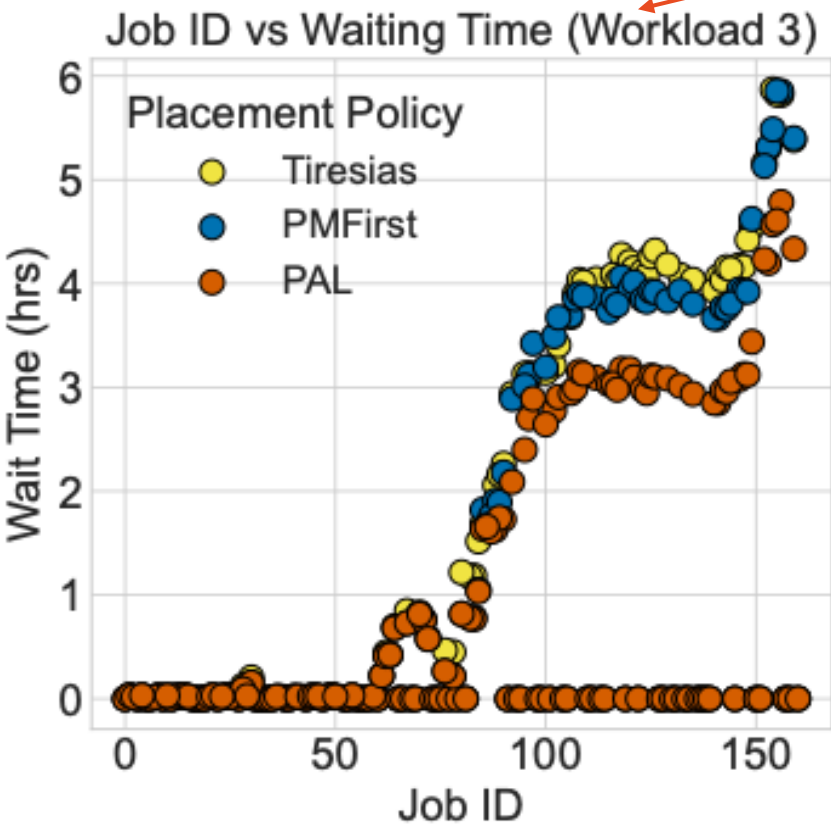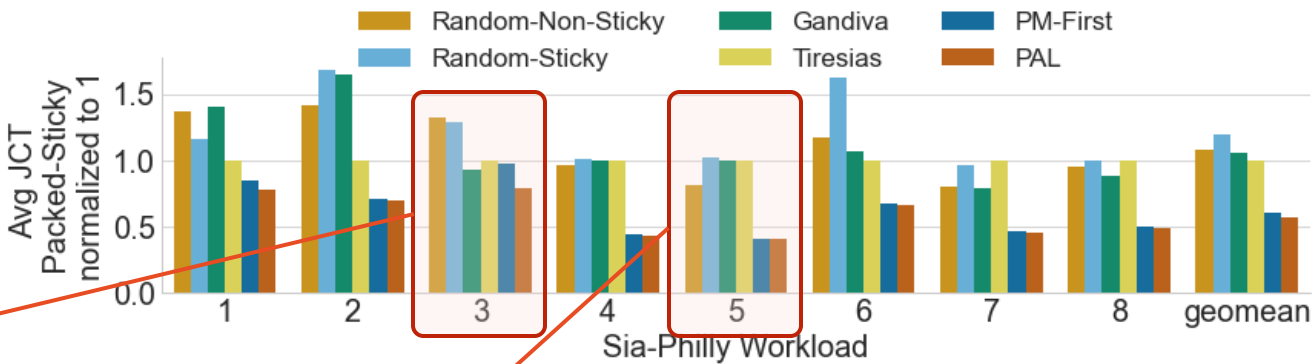
# SIA-PHILLY SIMULATIONS

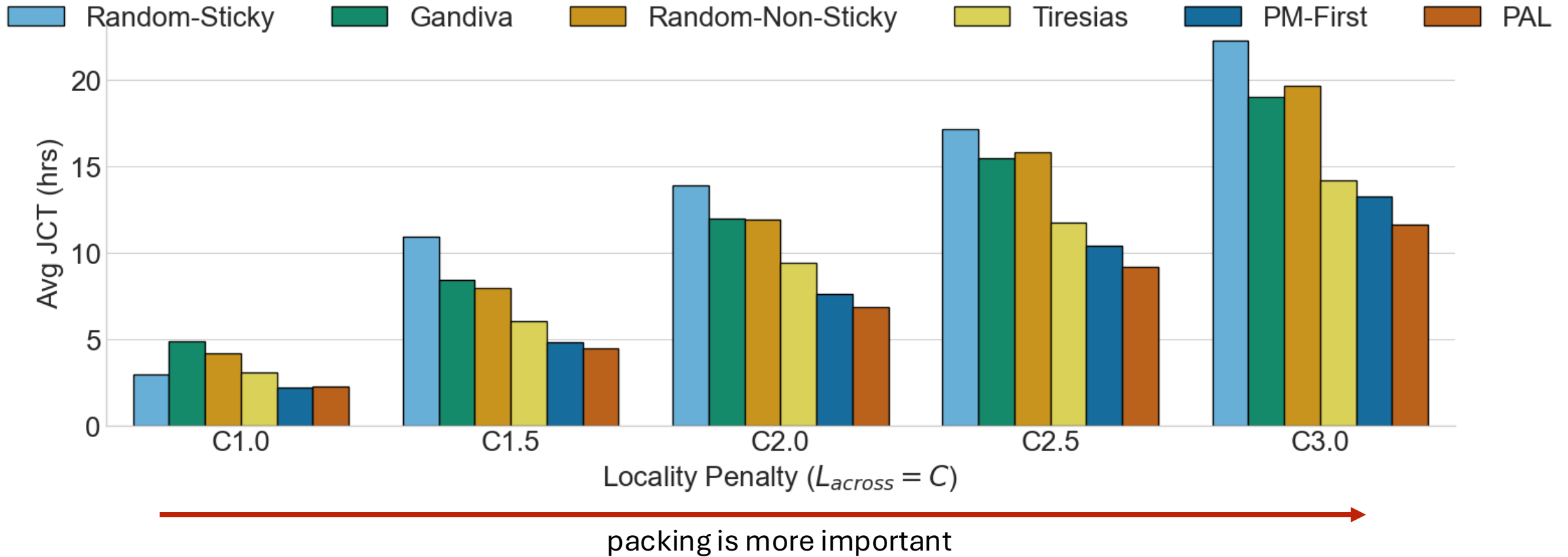# SIA-PHILLY SIMULATIONS

# SIA-PHILLY SIMULATIONS



Large jobs hog cluster resources, causing cascading wait time impact!

Jobs finish faster with PAL and subsequent wait times consequently reduce

# Varying Locality Penalty

# job arrival rate, cluster utilization, placement overheads...

- Varying cluster size (64 vs 256) – **Section V-B vs V-C**

- Varying locality penalty – **Section V-B (1) and V-C (2)**

- Varying job load (jobs/hr) – **Section V-C (1)**
  - Average JCT
  - Utilization

- Different proportions of single GPU jobs - **Section V-B vs V-C**

- Varying scheduling policy - **Section V-C (1)**
  - FIFO
  - SRTF
  - LAS

- Placement policy overheads  - **Section V-C (2)**

# CONCLUSION

# CONCLUSION

Variability can affect cluster performance, utilization and load balancing on GPU clusters. Likely to get worse as **ML algorithms** and c**luster size** scales up, while **transistors** shrink.

# CONCLUSION

Variability can affect cluster performance, utilization and load balancing on GPU clusters. Likely to get worse as **ML algorithms** and c**luster size** scales up, while **transistors** shrink.

Key Insight: Schedulers must **embrace** and harness this variability.

# CONCLUSION

Variability can affect cluster performance, utilization and load balancing on GPU clusters. Likely to get worse as **ML algorithms** and c**luster size** scales up, while **transistors** shrink.

Key Insight: Schedulers must **embrace** and harness this variability.

We proposed **PAL** which uses application-specific variability to optimize job placement. Across a mix of ML workloads, PAL improves on SOTA ML cluster schedulers across metrics.

# CONCLUSION

Variability can affect cluster performance, utilization and load balancing on GPU clusters. Likely to get worse as **ML algorithms** and c**luster size** scales up, while **transistors** shrink.
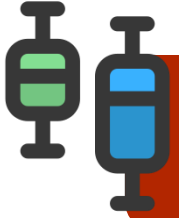
Key Insight: Schedulers must **embrace** and harness this variability.

We proposed **PAL** which uses application-specific variability to optimize job placement. Across a mix of ML workloads, PAL improves on SOTA ML cluster schedulers across metrics.

We are working to extend this to HPC and HPC+ML workloads. Looking at HW-SW codesign with variability as a first-order constraint.

# SUMMARY

Variability can affect cluster performance, utilization and load balancing on GPU clusters
Likely to get worse as ML algorithms and cluster size scales up, while transistors shrink

Key Insight: Schedulers must embrace and harness this variability

We proposed **PAL** which uses application-specific variability to optimize job placement
Across a mix of ML workloads, PAL improves on SOTA ML cluster schedulers across metrics.

We are working to extend this to HPC and HPC+ML workloads.
Looking at HW-SW codesign with variability as a first-order constraint.
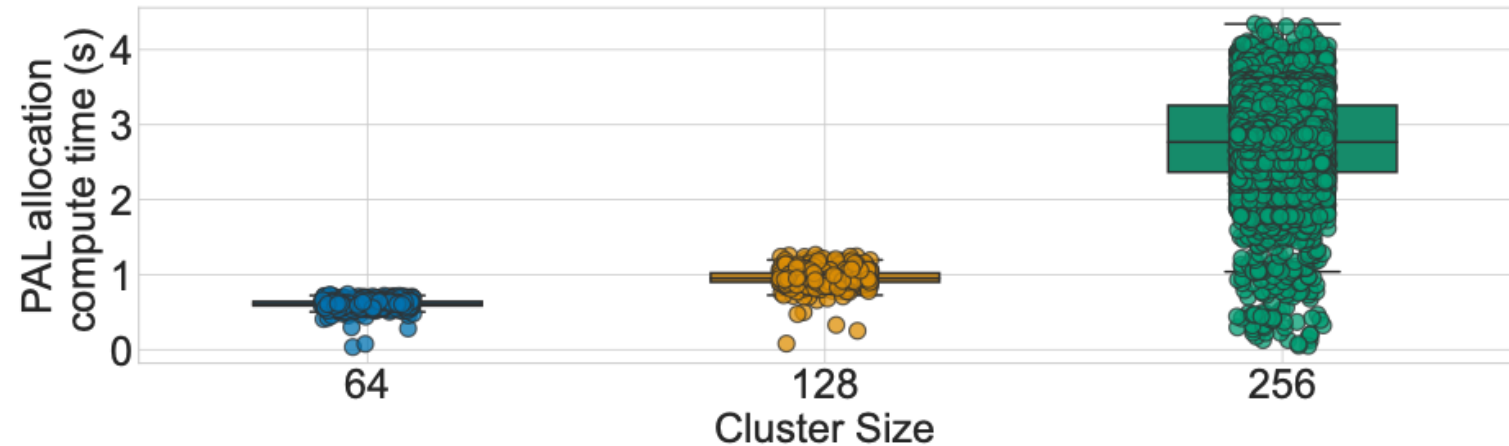
Contact

rnjain@wisc.edu

https://pages.cs.wisc.edu/~rnjain/

Artifact

THANK YOU
*Questions?*

# Placement Policy Overheads

- Scheduling round duration is 300s (5 mins)

- PAL
  - worst-case 4 seconds
  - median of 2.8 seconds
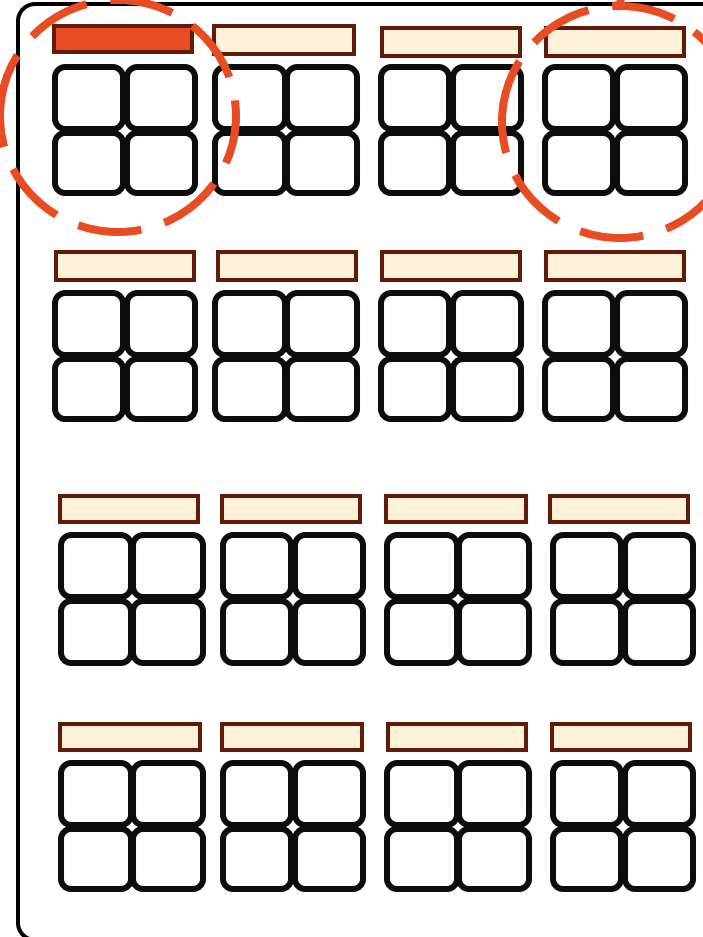
- PM-First
  - worst-case 2 seconds

# Blox Cluster Setup

## GPU Cluster



**Scheduler Node: CPU Host**
- gRPC Server that registers all GPU UUIDs.
- Another CPU threads sends job requests at specified arrival times.
- Scheduler thread runs placement policy algorithm to determine set of GPUs to run on.
- Server dispatches job by sending a gRPC packet to the clients running on relevant node(s).
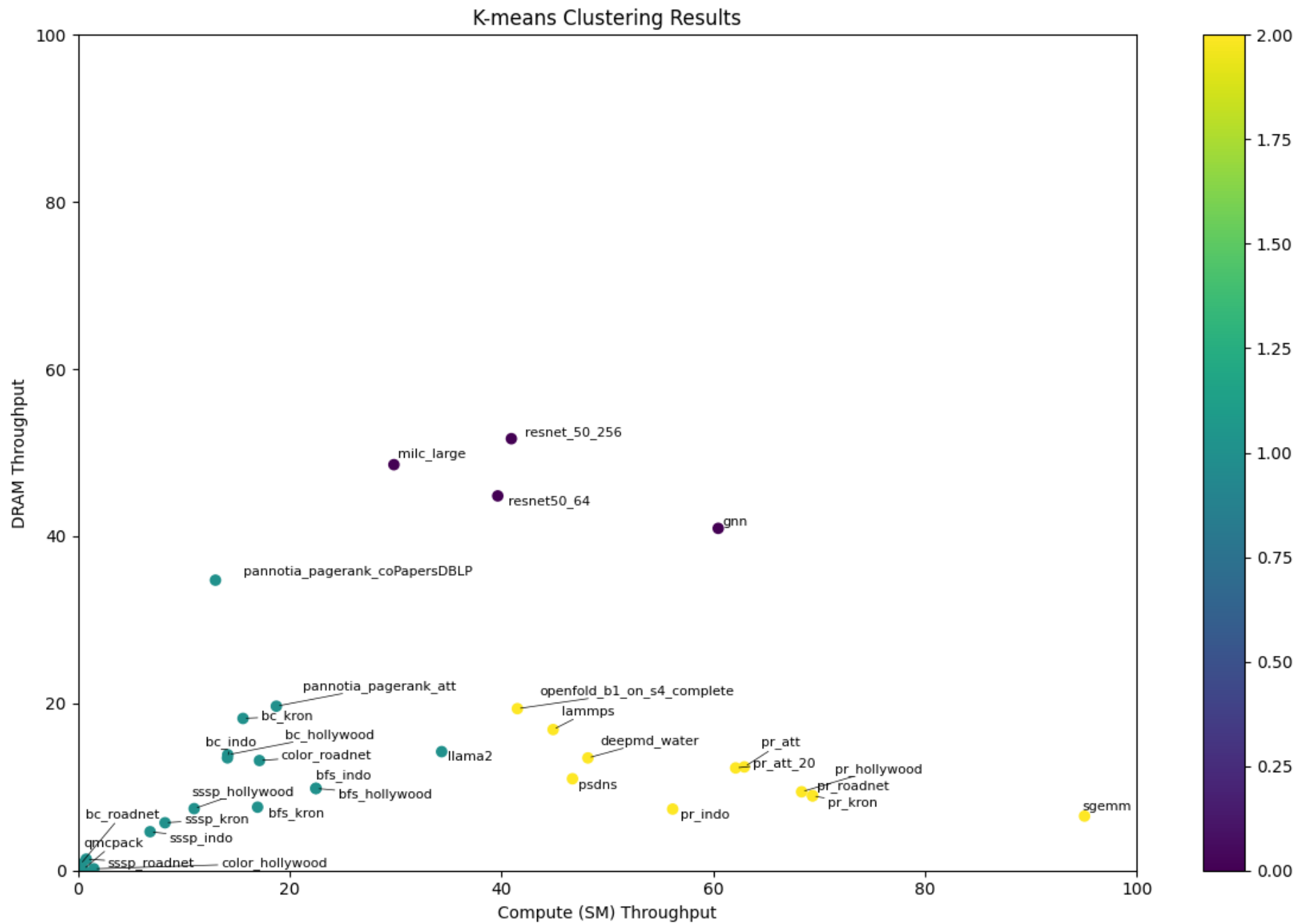
**Execution Node**
- Host on each execution node acts as gRPC client
- At startup, client registers it's GPUs with the scheduler node and receives acknowledgement.
- It then waits for job request packets from server
- Whenever it receives a job request, it runs corresponding job command on specified list of GPUs.

K-means Clustering Results

# How do we define variability?
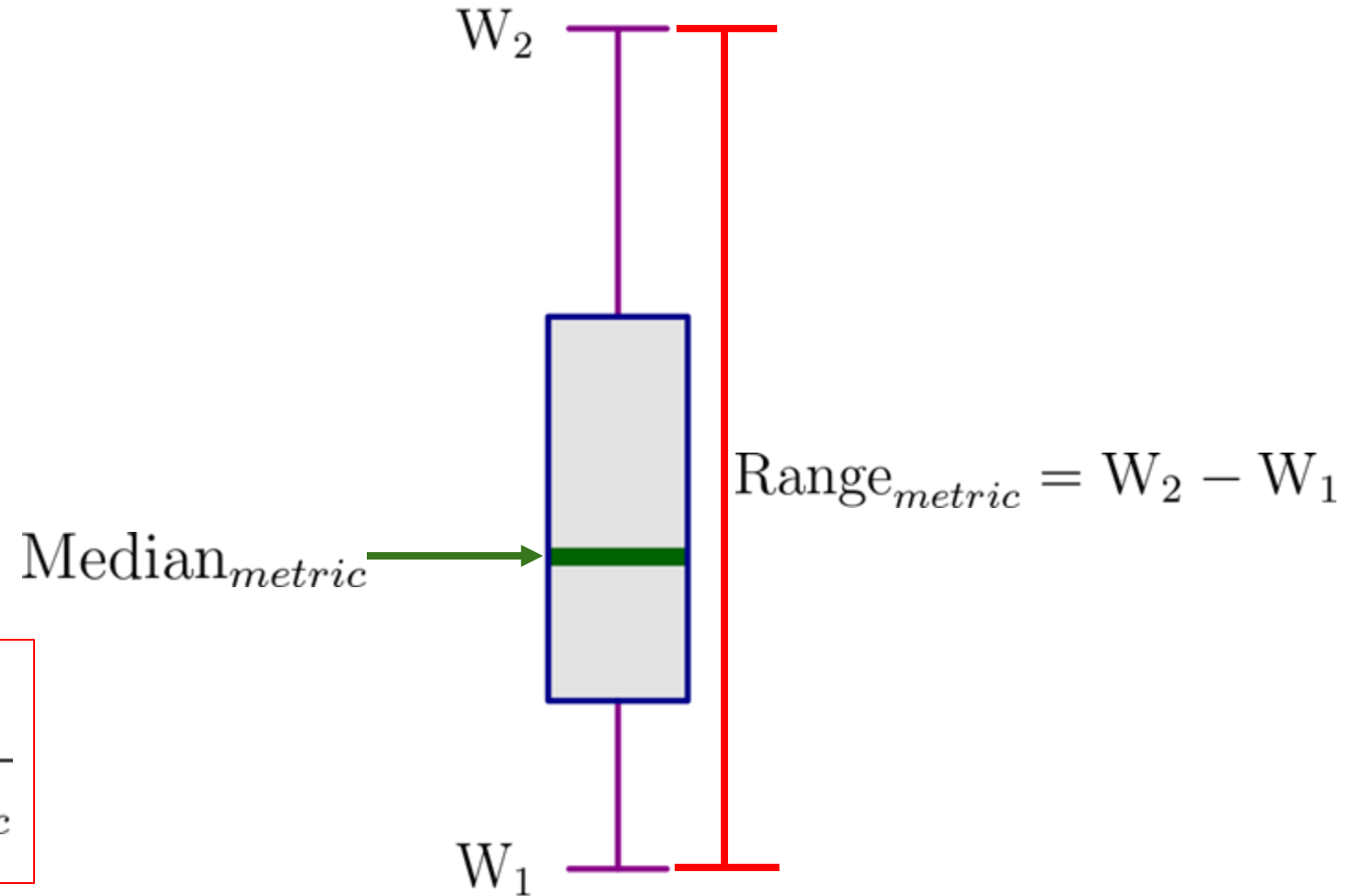
Interquartile Range

$$\mathrm{IQR} = Q_3 - Q_2$$

Whiskers and Range

$$W_2 = Q_3 + 1.5 \times \mathrm{IQR}$$
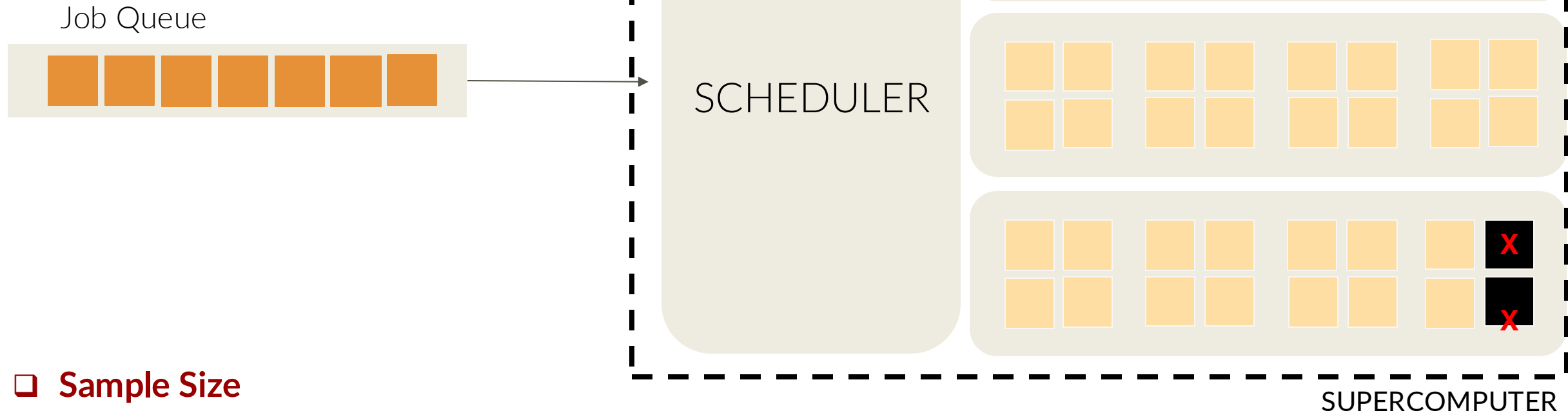$$W_1 = Q_1 - 1.5 \times \mathrm{IQR}$$
$$\mathrm{Range}_{metric} = W_2 - W_1$$

$$\mathrm{Variability}_{metric} = \frac{\mathrm{Range}_{metric}}{\mathrm{Median}_{metric}}$$



$$W_2$$

$$\mathrm{Range}_{metric} = W_2 - W_1$$

$$\mathrm{Median}_{metric}$$

$$W_1$$

# Methodology



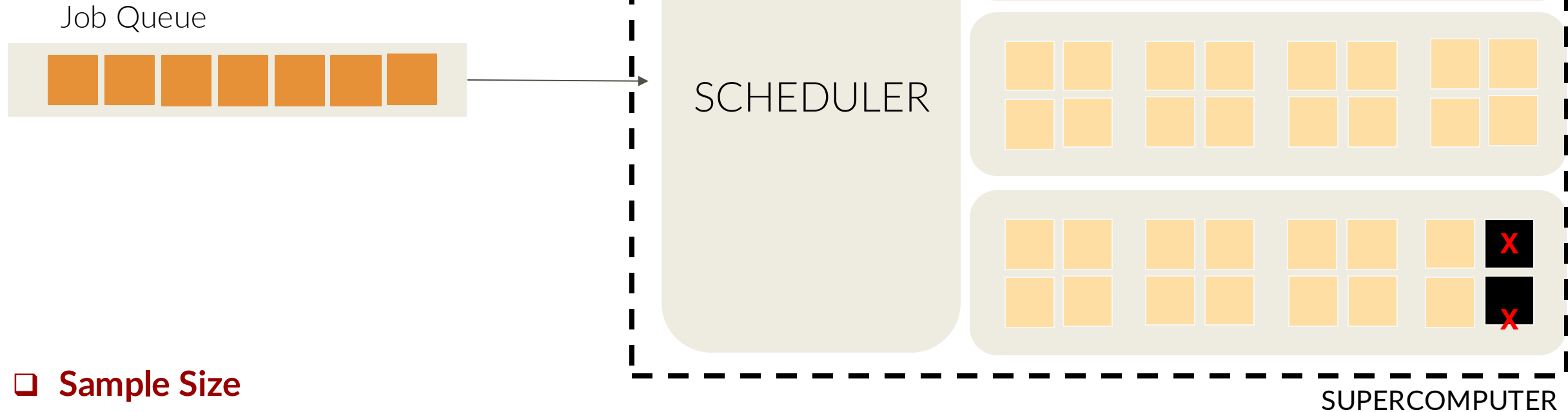Job Queue

SCHEDULER

SUPERCOMPUTER

❑ **Sample Size**
  ❑ Sample measurements from almost all GPUs in each cluster
  ❑ Profiled 2.9x more GPUs than worst-case recommendations for statistical significance [Scogland, et al. SC'15]

# Methodology

**Legend:** Available · Running SGEMM · X Unavailable · Completed SGEMM
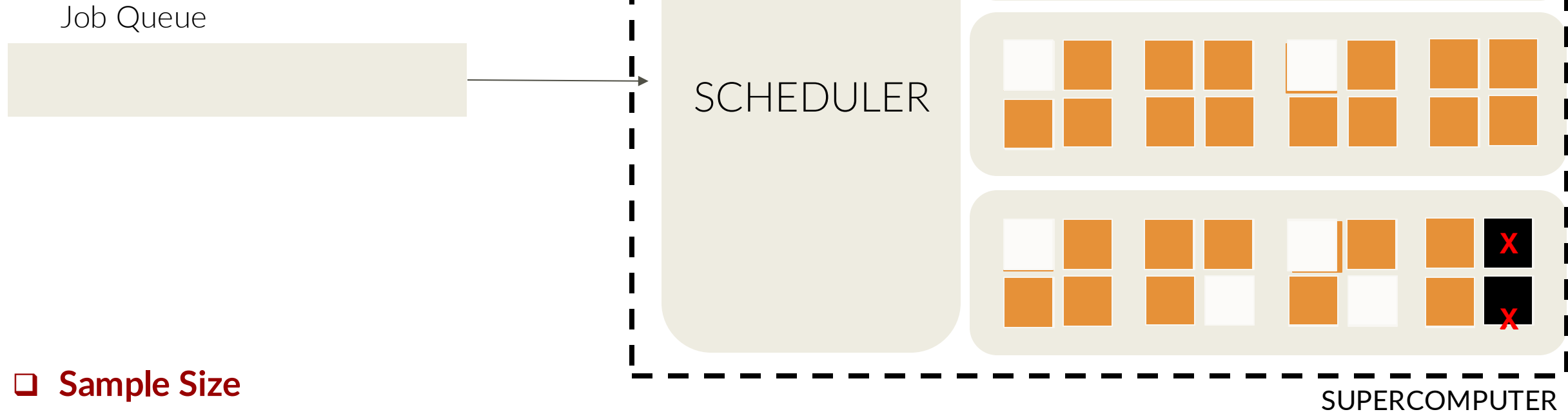
Job Queue

SCHEDULER

SUPERCOMPUTER

- ❑ **Sample Size**
  - ❑ Sample measurements from almost all GPUs in each cluster
  - ❑ Profiled 2.9x more GPUs than worst-case recommendations for statistical significance [Scogland, et al. SC'15]

# Methodology

**Available** **Running SGEMM** **X Unavailable** **Completed SGEMM**

Job Queue

SCHEDULER

SUPERCOMPUTER

❑ **Sample Size**
  ❑ Sample measurements from almost all GPUs in each cluster
  ❑ Profiled 2.9x more GPUs than worst-case recommendations for statistical significance [Scogland, et al. SC'15]

# CLUSTER PARAMETERS

| | Cluster | GPU Architecture | Number of GPUs | Cooling Method |
|---|---|---|---|---|
| Lawrence Livermore National Laboratory | **LLNL Corona** | AMD MI60 | 328 | air cooled |
| TACC Texas Advanced Computing Center | **TACC Longhorn** | NVIDIA V100 | 416 | air cooled |
| TACC Texas Advanced Computing Center | **TACC Frontera** | NVIDIA Quadro RTX 5000 | 360 | mineral oil cooled |
| Sandia National Laboratories | **SNL Vortex** | NVIDIA V100 | 216 | water cooled |
| Oak Ridge National Laboratory | **ORNL Summit** | NVIDIA V100 | 27648 | air cooled |

Not All GPUs Are Created Equal: Characterizing Variability in Large-Scale, Accelerator-Rich Systems
https://dl.acm.org/doi/abs/10.5555/3571885.3571971

# VARIABILITY ACROSS APPLICATIONS



Performance Variability per application

22%    8%    7%    <1%    <1%

Compute bound for the input size and configurations they were run in

Memory bound