

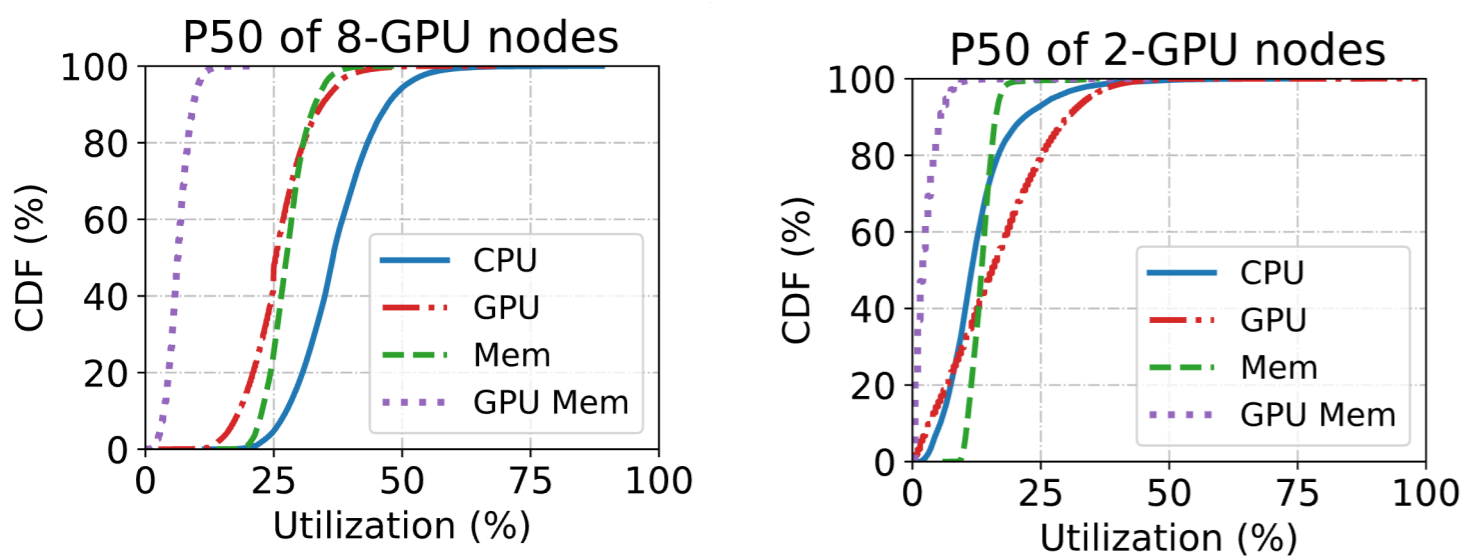
## Scheduling GPU-Sharing Workloads with Fragmentation Gradient Descent

Qizhen Weng<sup>†\*</sup> Lingyun Yang<sup>†\*</sup> Yinghao Yu<sup>^†</sup> Wei Wang<sup>†</sup> Xiaochuan Tang<sup>^</sup> Guodong Yang<sup>^</sup> Liping Zhang<sup>^</sup>  
<sup>†</sup>: Hong Kong University of Science and Technology    <sup>^</sup>: Alibaba Group    (\*: equal contribution)

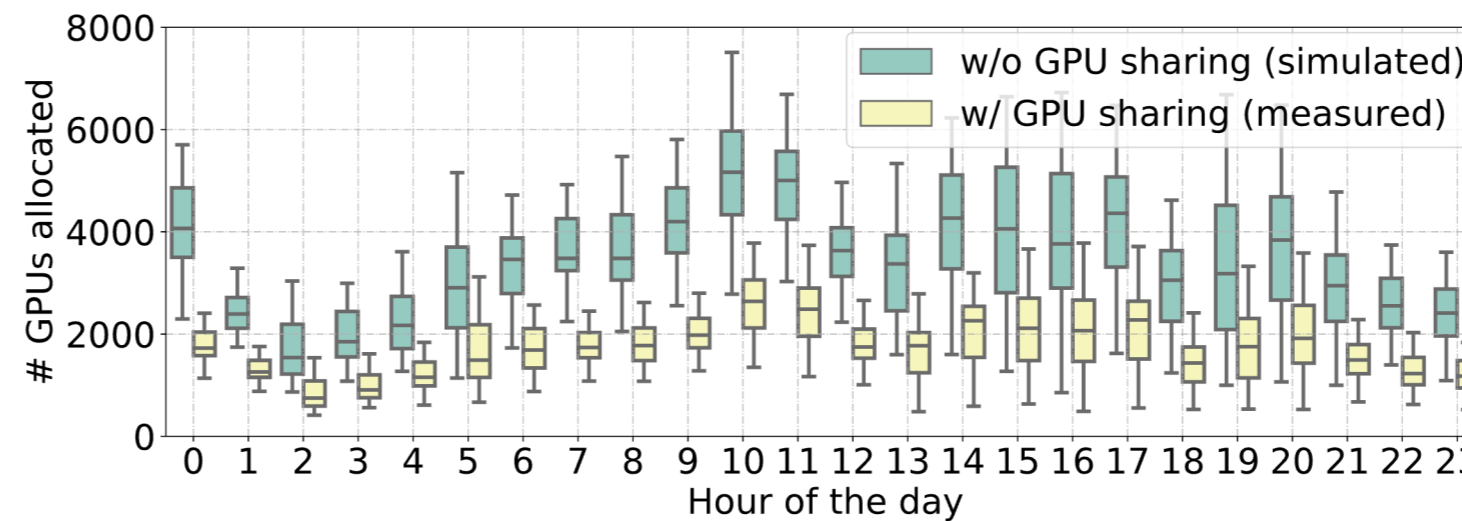
TL;DR: We propose a novel measure of fragmentation to statistically quantify the degree of GPU fragmentation caused by different sources. Based on this measure, we invent a scheduling policy FGD that packs tasks to minimize the growth of fragmentation and maximize GPU allocation.

### ML-as-a-Service clouds suffer low GPU utilization

Avg. 25-50% GPU utilization in production MLaaS clouds [1-3].



[1] Weng et al., "MLaaS in the Wild: Workload analysis and scheduling in large-scale heterogeneous GPU clusters," in NSDI 2022.  
 [2] Hu et al., "Characterization and prediction of deep learning workloads in large-scale GPU datacenters," in SC 2021.  
 [3] Narayanan et al., "Heterogeneity-aware cluster scheduling policies for deep learning workloads," in OSDI 2020



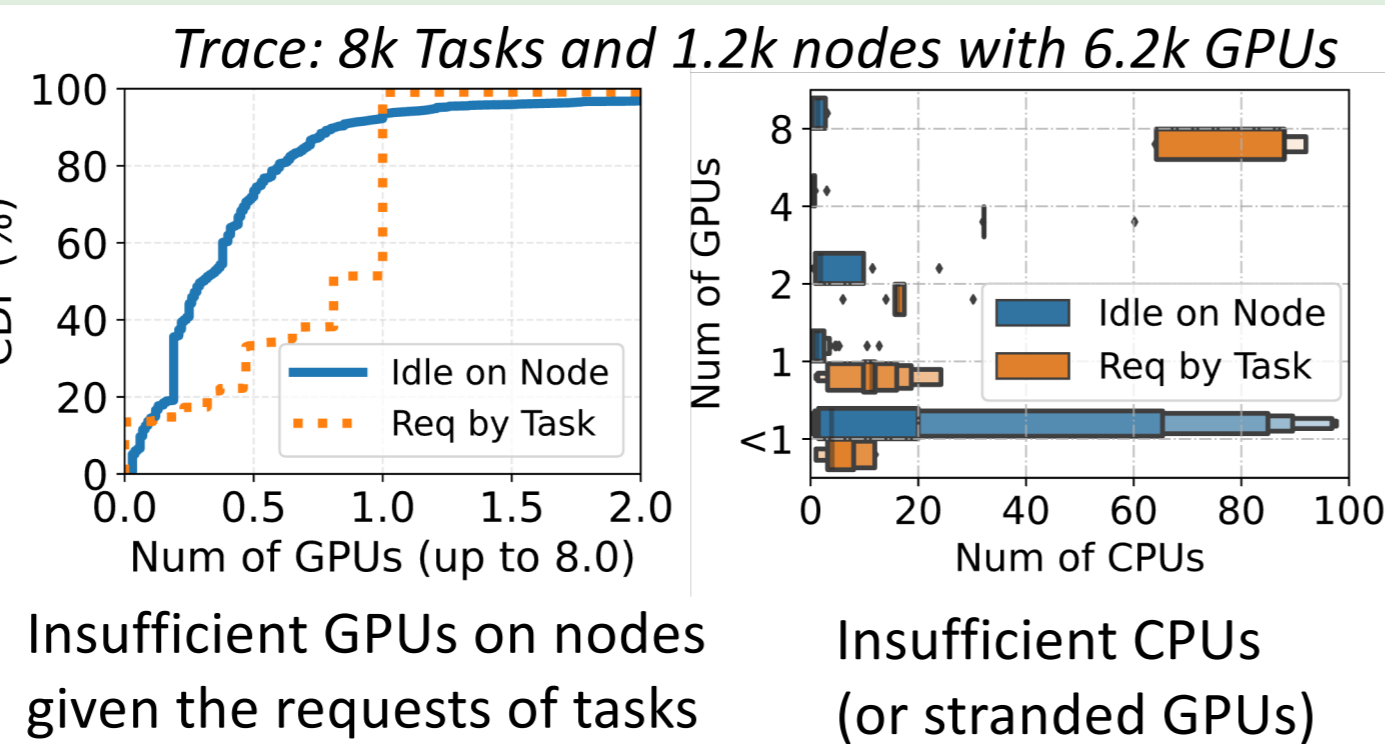
GPU sharing lets multiple tasks run on a single GPU, via DL framework manipulation, or CUDA API interception, or hardware-assisted methods (e.g., MIG).

← Sharing saves 50% GPUs in Alibaba [1].

### Yet, GPU sharing doesn't always improve allocation. Often, allocating partial GPUs results in fragmentation

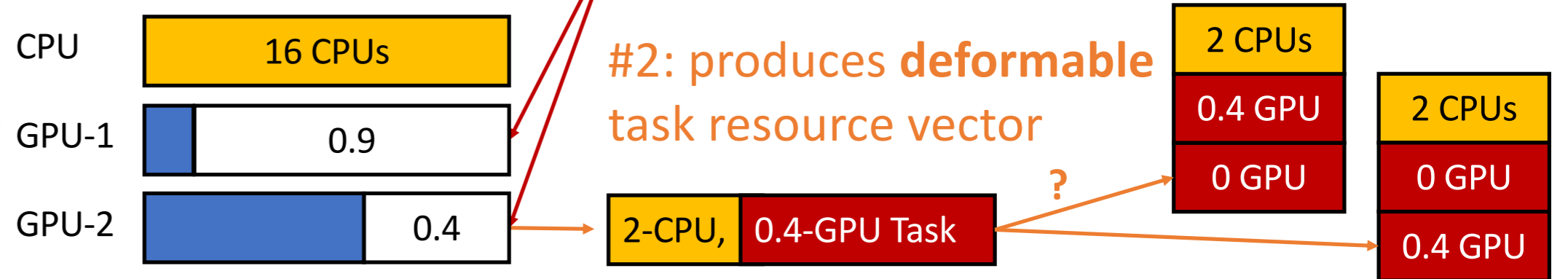
In many clusters, the GPU allocation rate can reach 85-90% maximum, leaving hundreds of GPUs unable to allocate!

Many users experienced scheduling failures even with sufficient GPU allocation quotas.



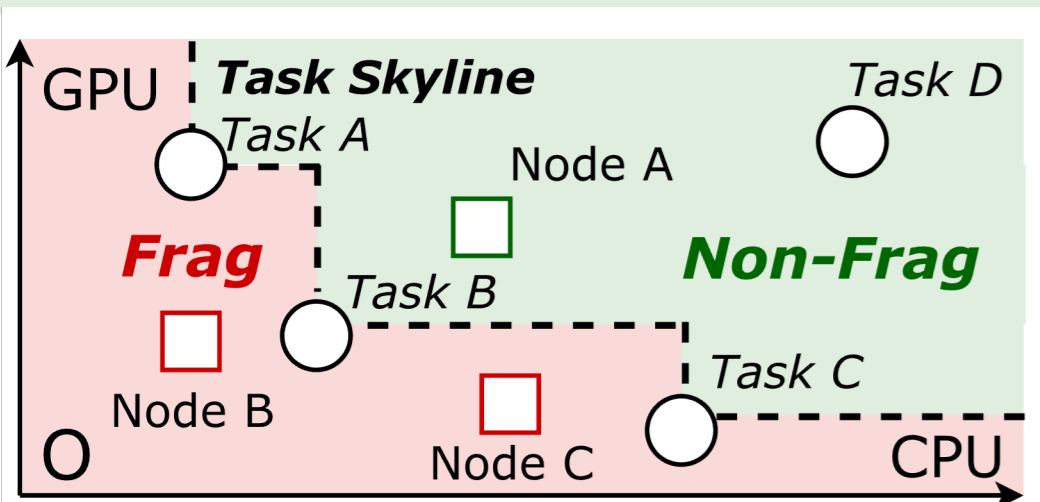
### Classical multi-resource bin-packing cannot work effectively on GPUs due to formulation mismatch

Neither of these formulation attempts works: (1) treating multiple GPUs as a unified logical device; (2) treating each GPU as an independent resource dimension.



### Definition of GPU Fragmentation:

The absolute measure is defective. Be statistical



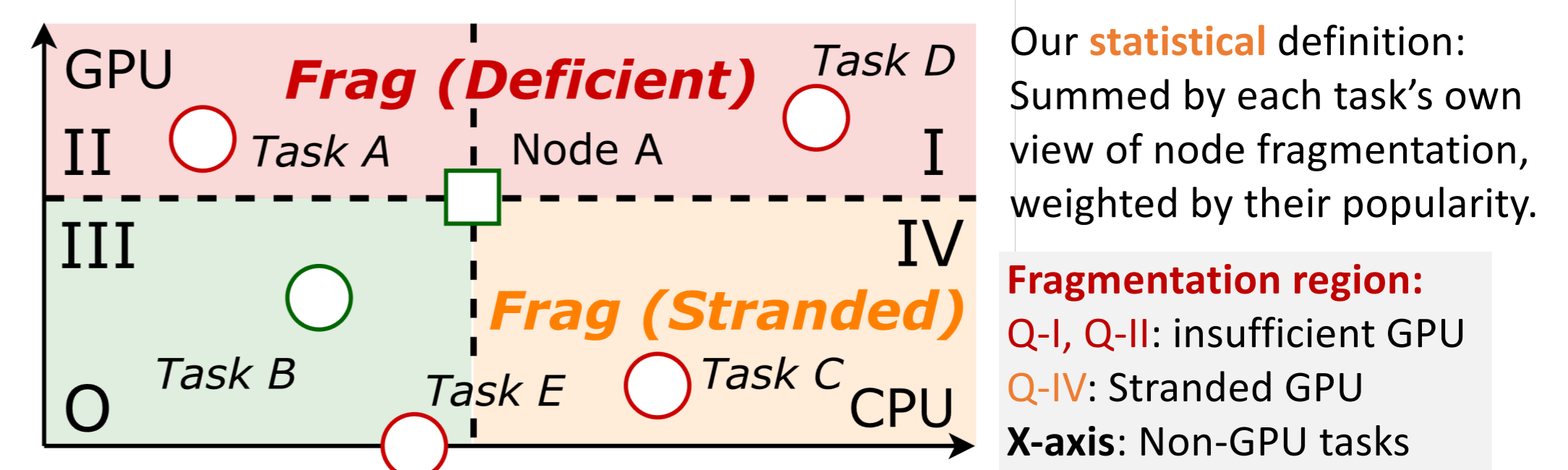
A defective definition of fragmentation in absolute terms — "a node is fragmented if and only if it cannot run any task". Task skyline determines the frag / non-frag boundary, yet, only 0.06% task instances belong to the skyline 😞



Absolute fragmentation stays low (<5%) throughout scheduling simulation (8k tasks to 6.2k GPUs) — 😞 fail to provide useful feedback to the scheduling quality

$$F_n(M) = \sum_{m \in M} p_m F_n(m) \quad (p_m: \text{task popularity})$$

For each task  $m$  in task set  $M$ , Sum the fragmentation viewed by task  $m$



Our statistical definition: Summed by each task's own view of node fragmentation, weighted by their popularity.

Fragmentation region:  
 Q-I, Q-II: insufficient GPU  
 Q-III, Q-IV: Stranded GPU  
 X-axis: Non-GPU tasks

Fragmentation rate: the likelihood of tasks in fragmentation regions:

- 😊 Aware of workload distribution while stable to small changes.
- 😞 Break down fragmentation into Deficient and Stranded.
- 😊 Independent of scheduling policy and node distribution.

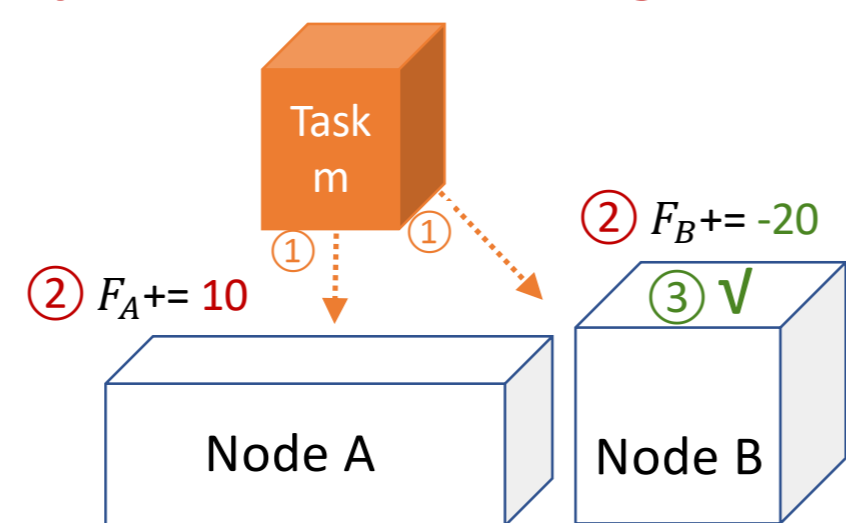
### Schedule Alg.: Fragmentation Gradient Descent

Algorithm 1: Node selecting process of FGD

Input : Node set  $N$ , incoming task  $m$ , workload set  $M$   
 Output : Assigned node  $n^*$

- Initialize node score set  $S \leftarrow \emptyset$ , and output  $n^* \leftarrow \emptyset$ .
- parallel for node  $n \in N$  do
- if Insufficient resources || constraints not met then
- Return ▷ Filter out unavailable nodes
- $n^- \leftarrow \text{AssignTaskToNode}(m, n)$  ▷ Hypothetically
- $\Delta \leftarrow F_n^{\text{GPU}}(M) - F_n^{\text{GPU}}(M)$  ▷ Fragmentation increment
- $S \leftarrow S \cup (n, \Delta)$
- if  $S \neq \emptyset$  then
- $n^* \leftarrow \arg \min_{n \in S} \Delta$  ▷ pick the node with the least  $\Delta$ .

Schedule Tasks towards the Steepest Descent of Fragmentation



### Formal Description of Computation $F_n(m)$

- Case 1: All Residuals are Frag. (Q-I, Q-II, Q-IV, x-axis):

$$F_n(m) = \sum_{1 \leq g \leq G_n} R_{n,g}^{\text{GPU}} \quad \text{Residual resource on GPU } g \text{ Node } n$$

$G_n$ : GPU set on node  $n$

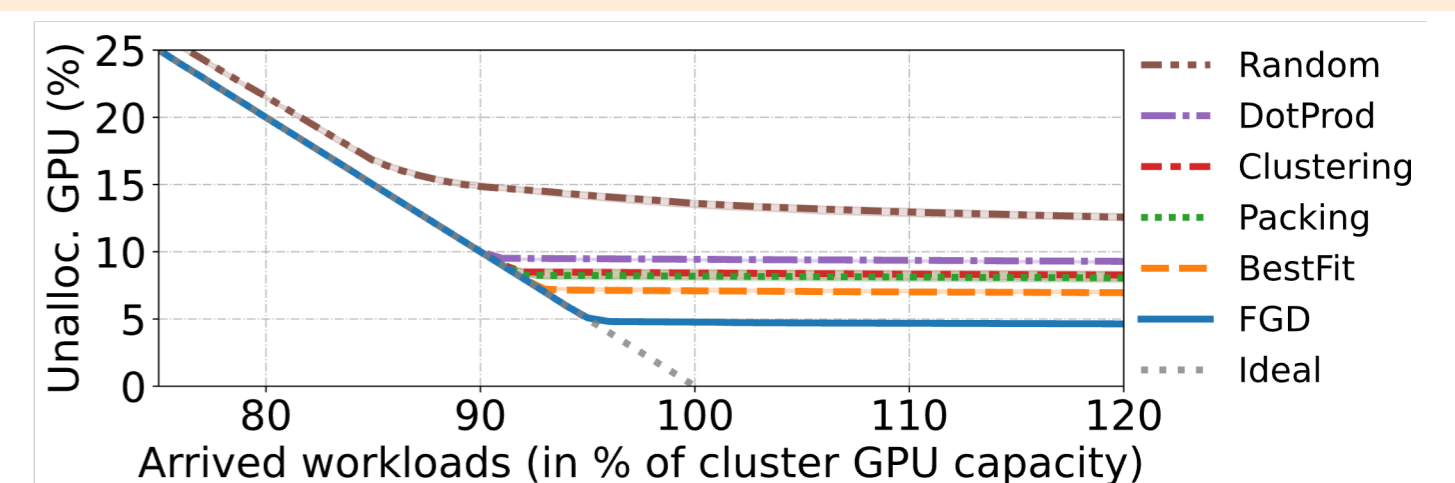
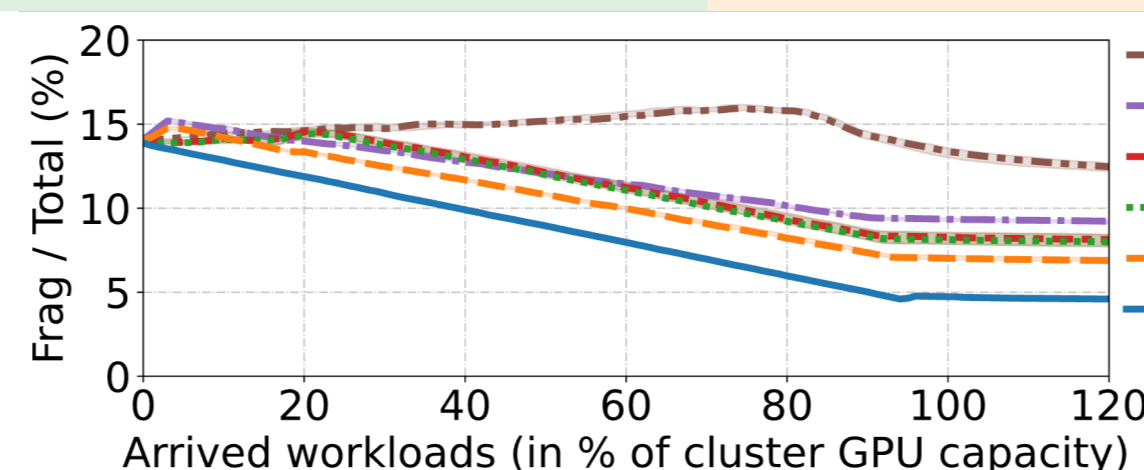
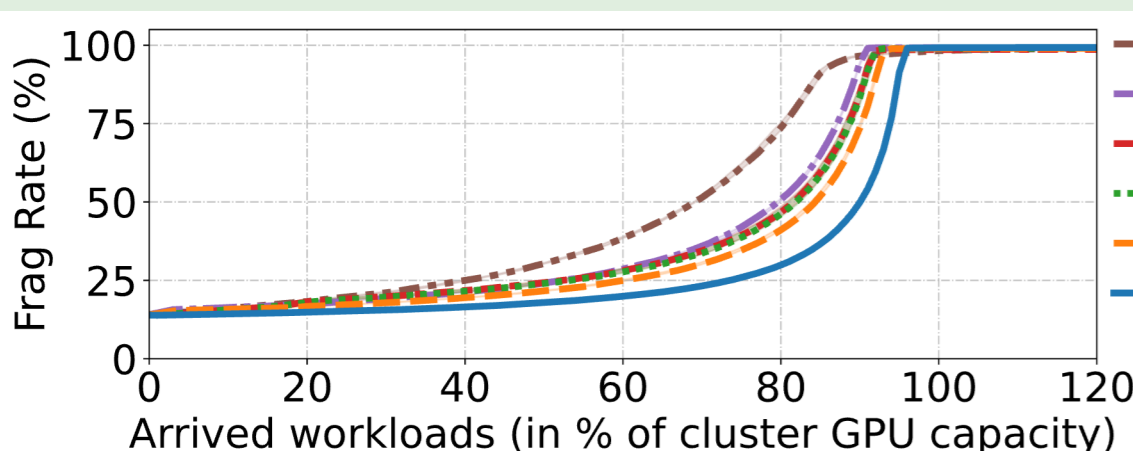
- Case 2: Partial or No Residuals are Frag. (Q-III):

$$F_n(m) = \sum_{1 \leq g \leq G_n} R_{n,g}^{\text{GPU}} \mathbb{1}(R_{n,g}^{\text{GPU}} < \min\{D_m^{\text{GPU}}, 1\})$$

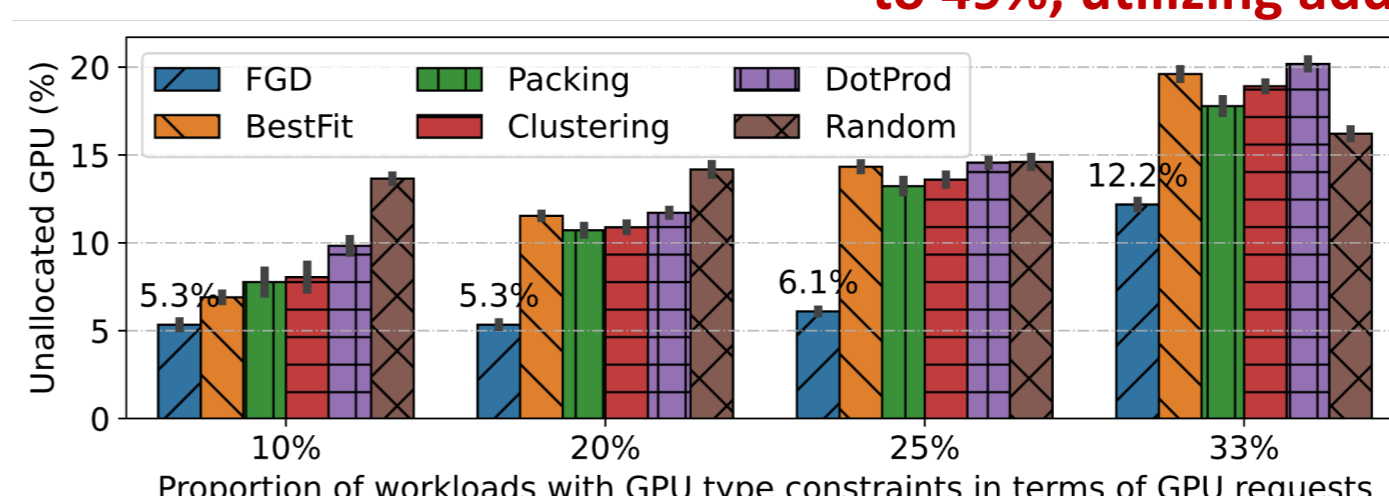
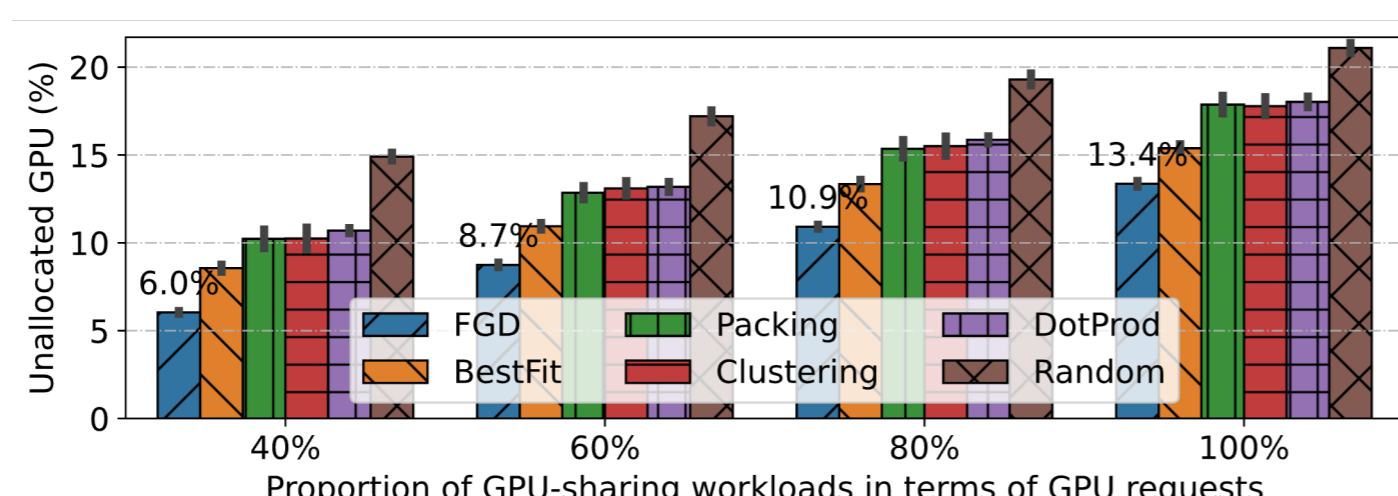
1, if remaining resource is smaller than the demand of task  $m$ , else 0.

### Evaluation: Schedule 8k tasks to 6.2k GPUs (1.2k nodes)

### FGD: Lowest Frag. Rate & Fewest GPUs Unallocated



Compared to existing packing-based schedulers, FGD pursues the lowest fragmentation and fewest GPUs unallocated: FGD reduces unallocated GPUs by up to 49%, utilizing additional 290 GPUs.



FGD allocates more GPUs across a variety of settings. See more results and task distributions in paper and code.

