



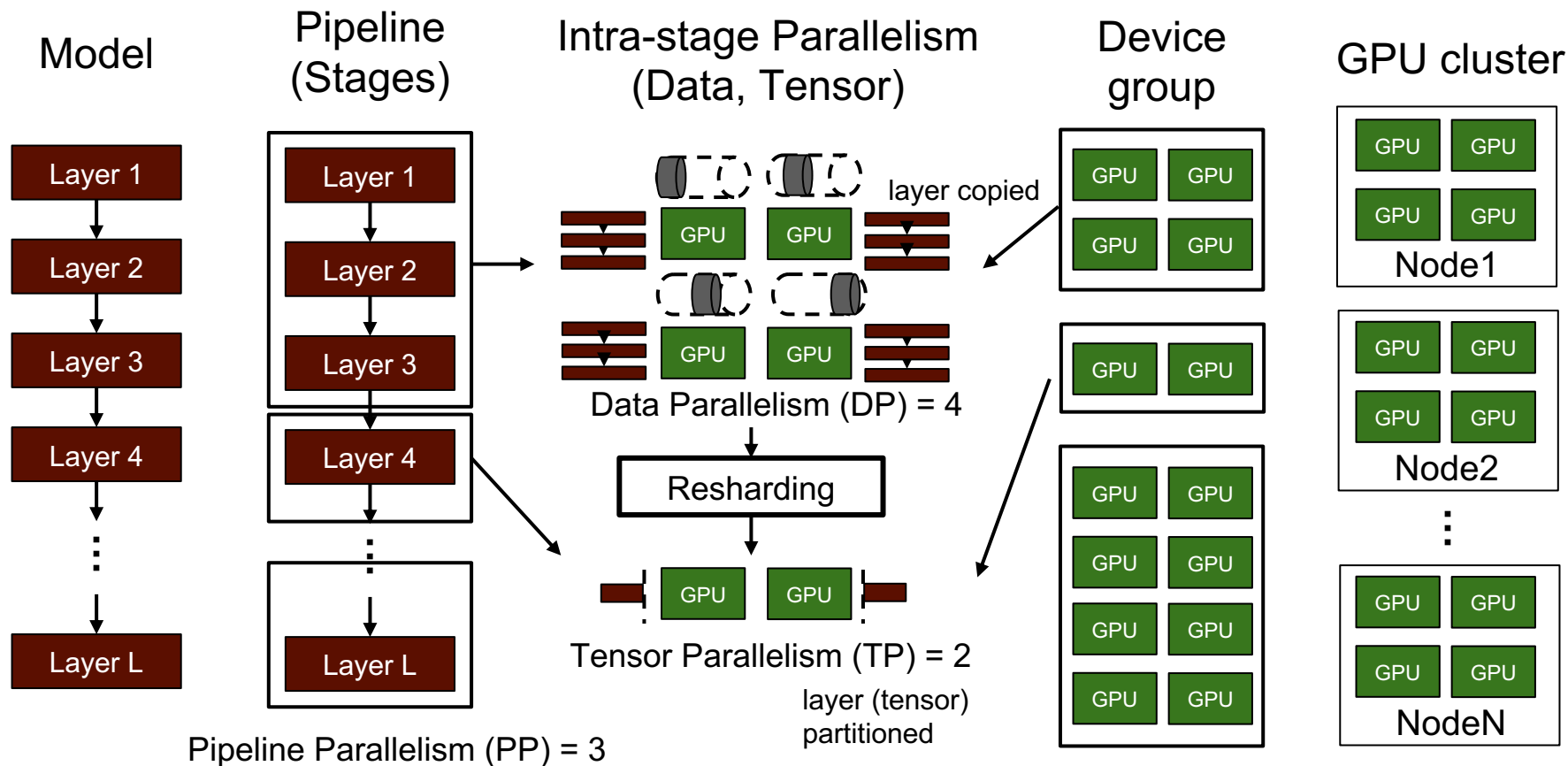
Metis: Fast Automatic Distributed Training on Heterogeneous GPUs

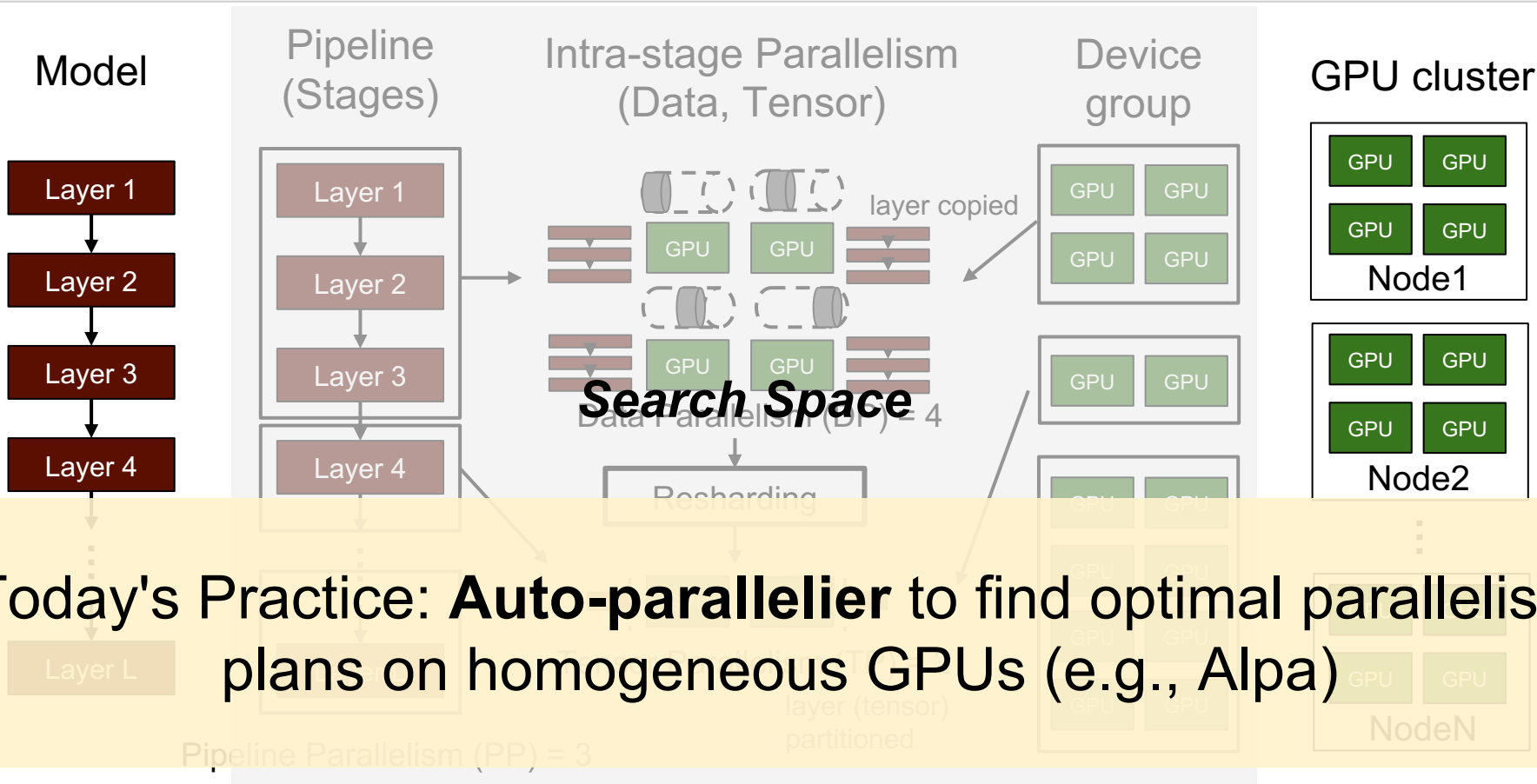
**Taegeon Um§, *Byungsoo Oh§, *Minyoung Kang§, Woo-Yeon Lee§, Goeun Kim§
Dongseob Kim§, Youngtaek Kim§, Mohd Muzzammil§ and Myeongjae Jeon¶*

** Equal contribution*

§ Samsung Research ¶ UNIST

Automatic Distributed Training of Large Model





Today's Practice: **Auto-parallelizer** to find optimal parallelism plans on homogeneous GPUs (e.g., Alpa)

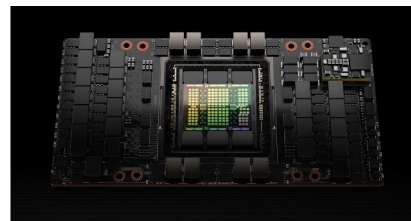
Heterogeneous GPUs in GPU Clusters



V100, Volta ('17)



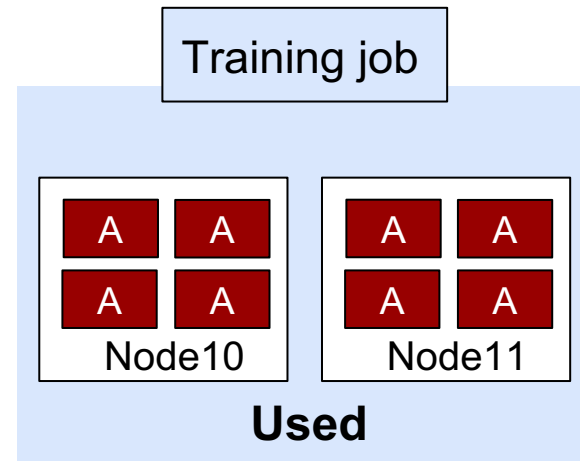
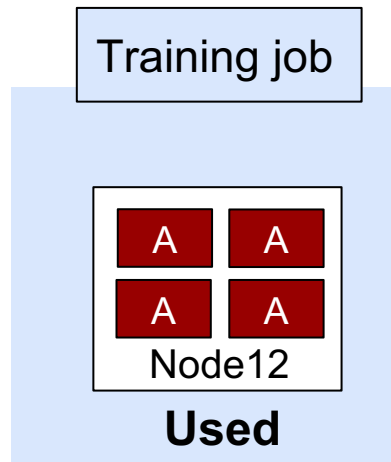
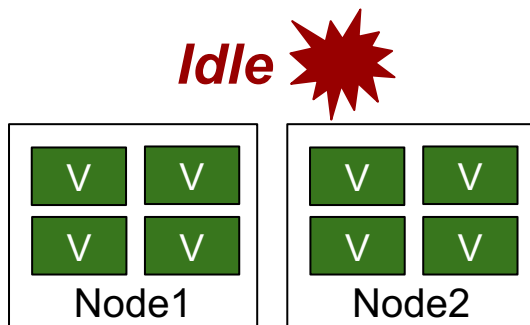
A100, Ampere ('20)

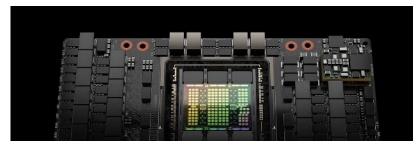


H100, Hopper ('23)



B100, Blackwell ('24?)





**Our focus: Auto-parallelizer
on heterogeneous GPUs**

V100, Volta ('17)

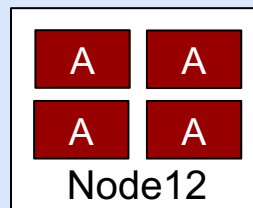
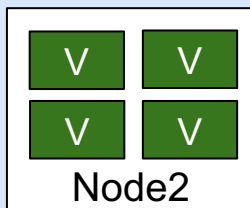
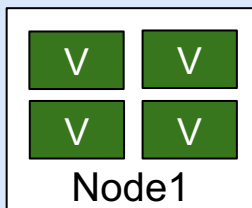
A100, Ampere ('20)

H100, Hopper ('23)

B100, Blackwell ('24?)

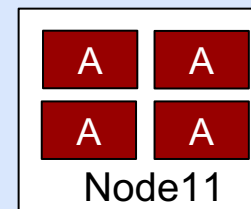
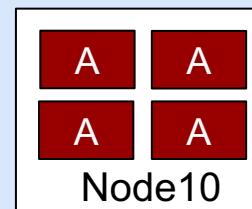
Training job

Distributed training on heterogeneous GPUs!



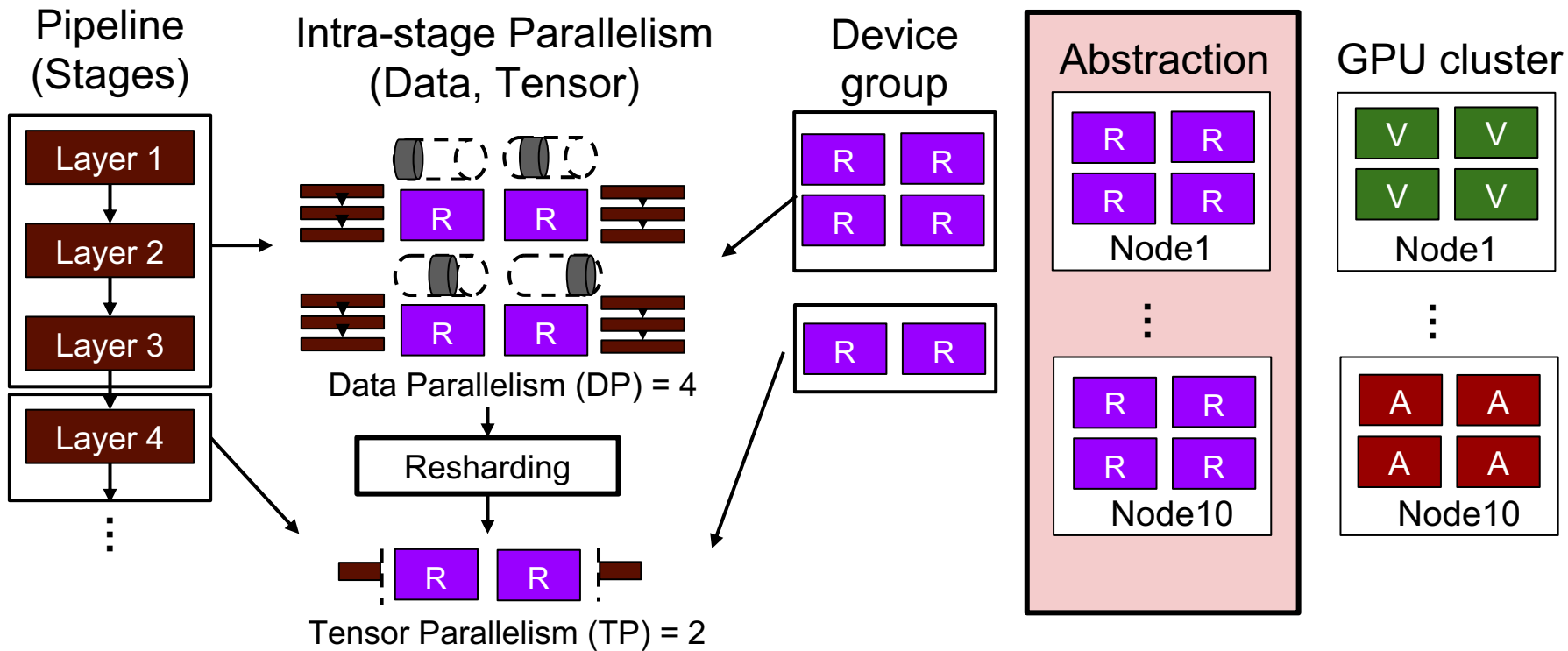
Good performance, Good utilization 😊

Training job

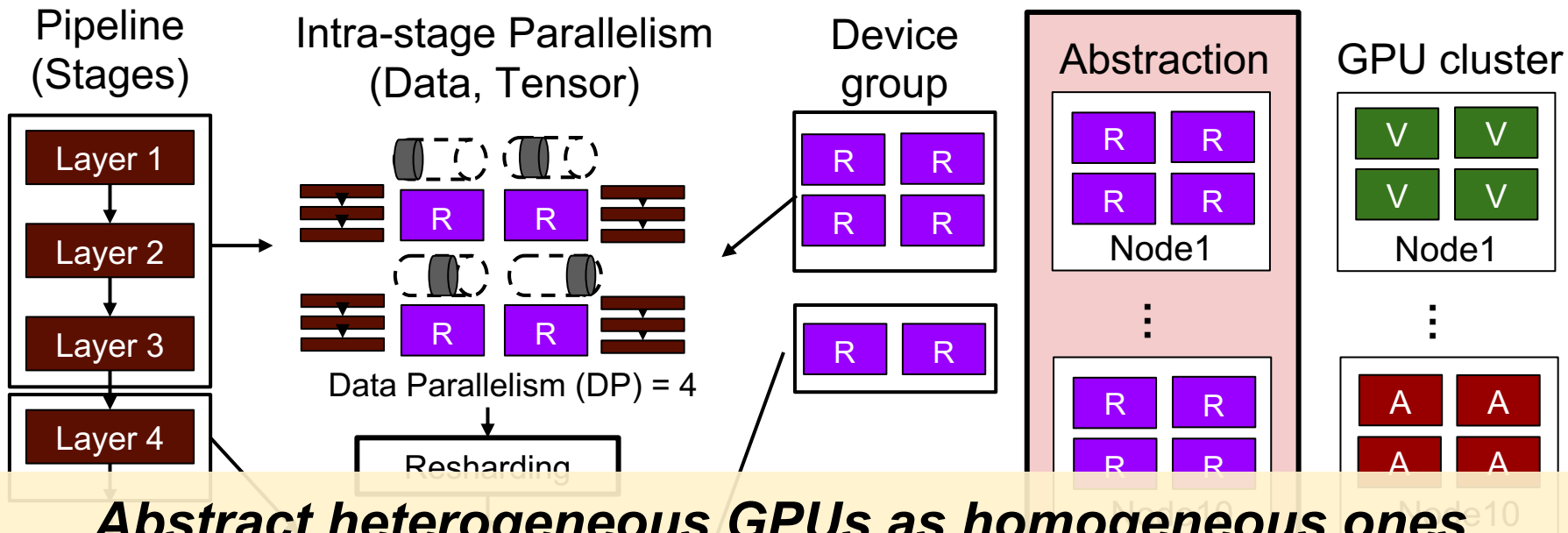


Used

Existing Auto-Parallelizer on Heterogeneous GPUs Samsung Research

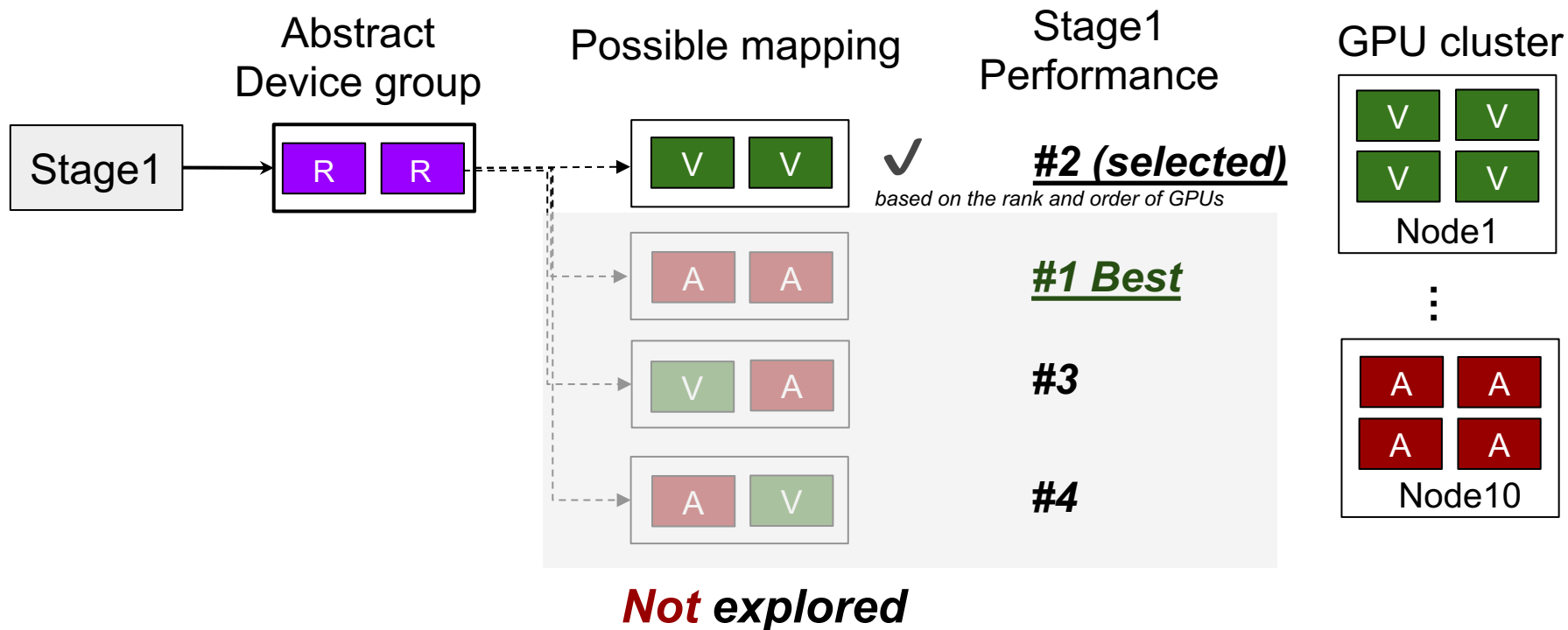


Existing Auto-Parallelizer on Heterogeneous GPUs Samsung Research

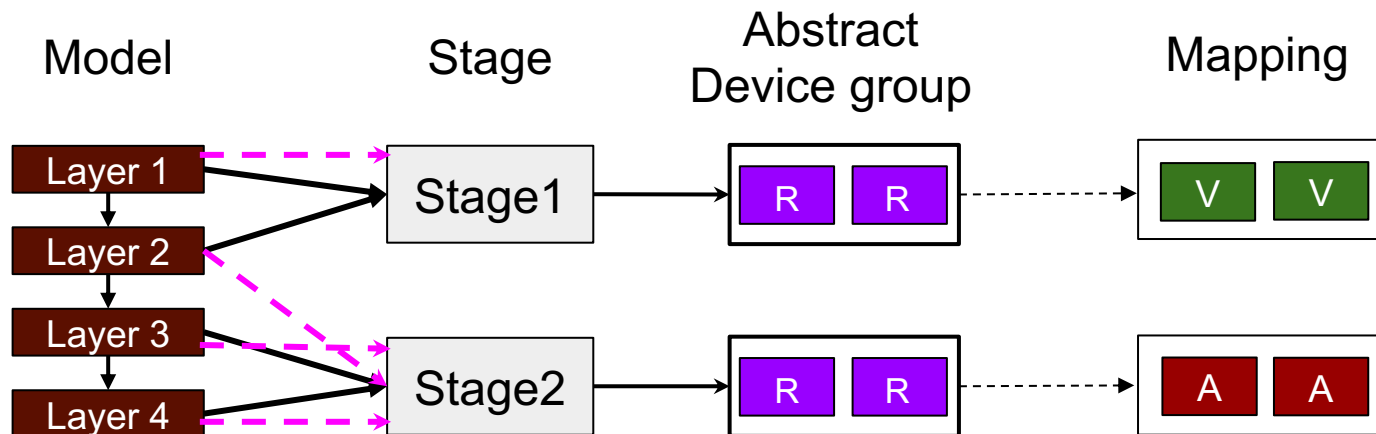


Tensor Parallelism (TP) = 2

1) *Unexplored device groups*



2) Load balancing based on the number of GPUs

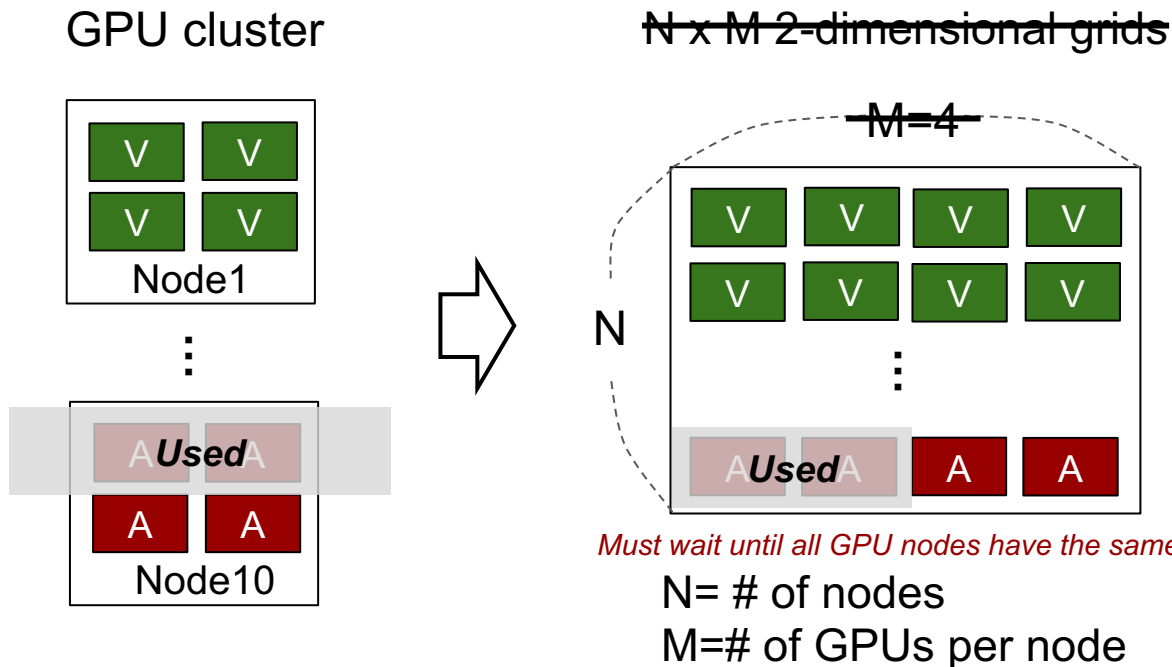


*  : Optimal layer load balancing

*  : Selected layer load balancing

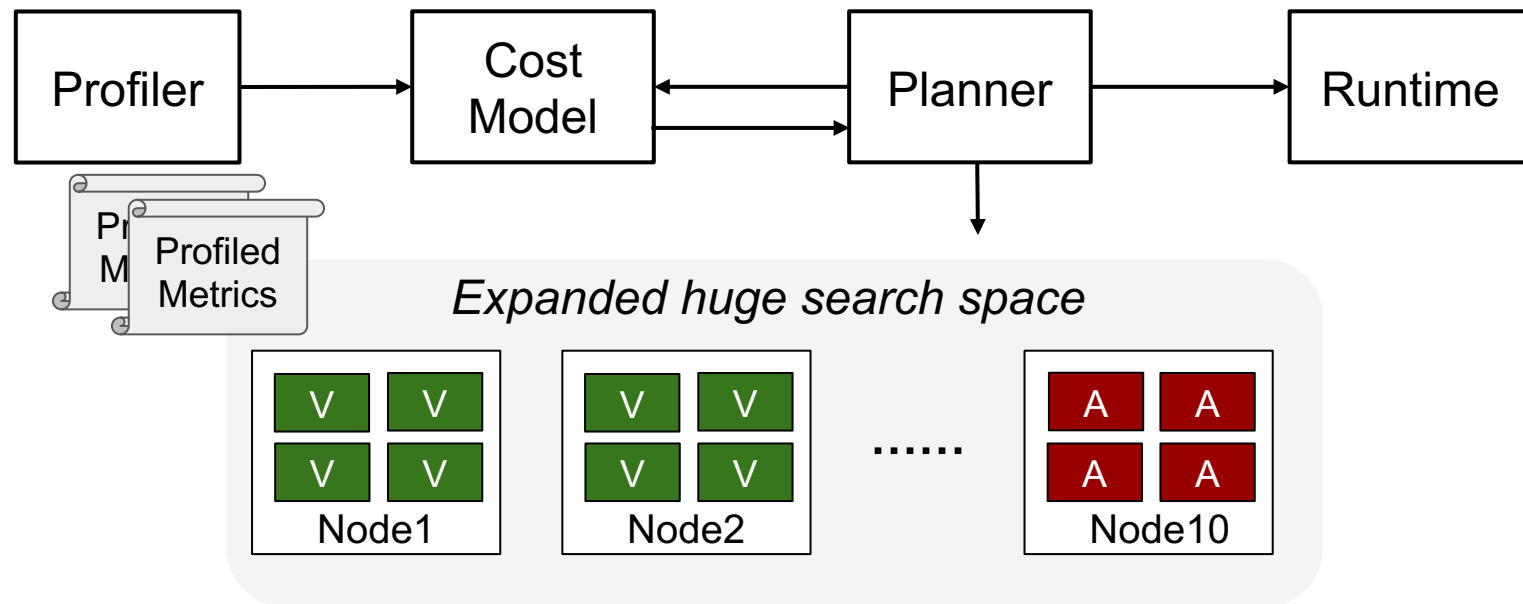
* GPU performance:  > 

3) Cluster shape constrained by 2-D grids

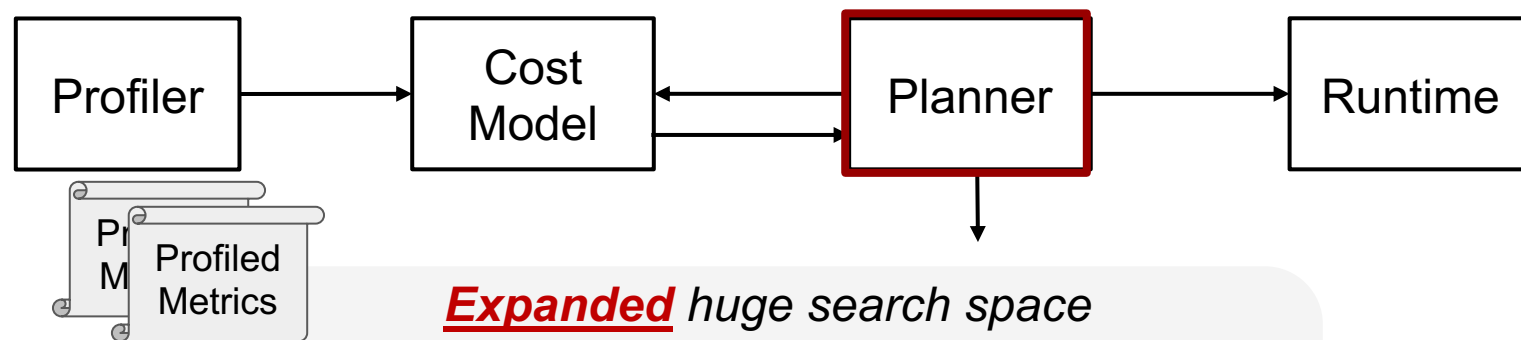


Breaks 2-D abstraction \Rightarrow Not supported

Expands the search space of plans by being aware of heterogeneous computing, memory, and number of GPUs



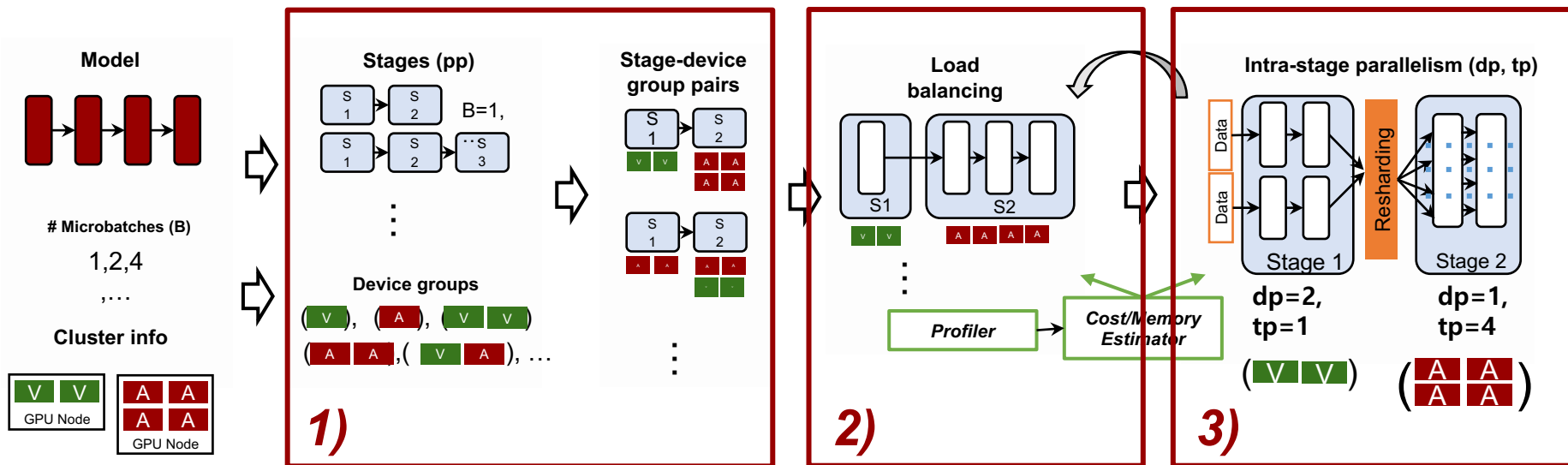
Expands the search space of plans by being aware of heterogeneous computing, memory, and number of GPUs



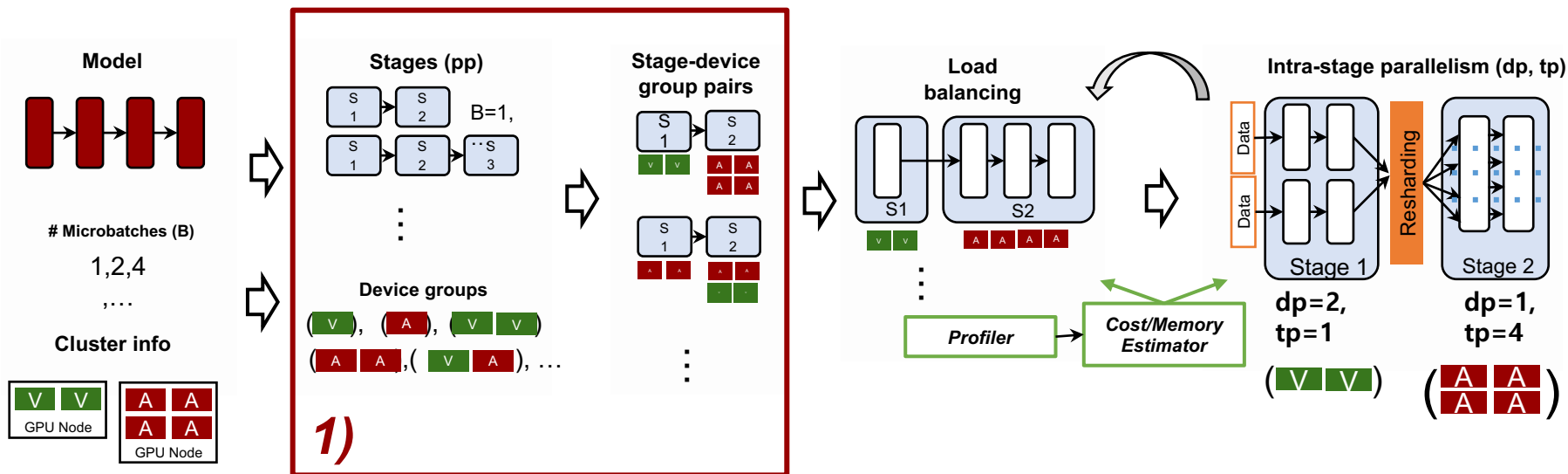
Challenge: profiling and search overheads that may take several days or weeks

** Please see the paper for the details of profiler and cost model*

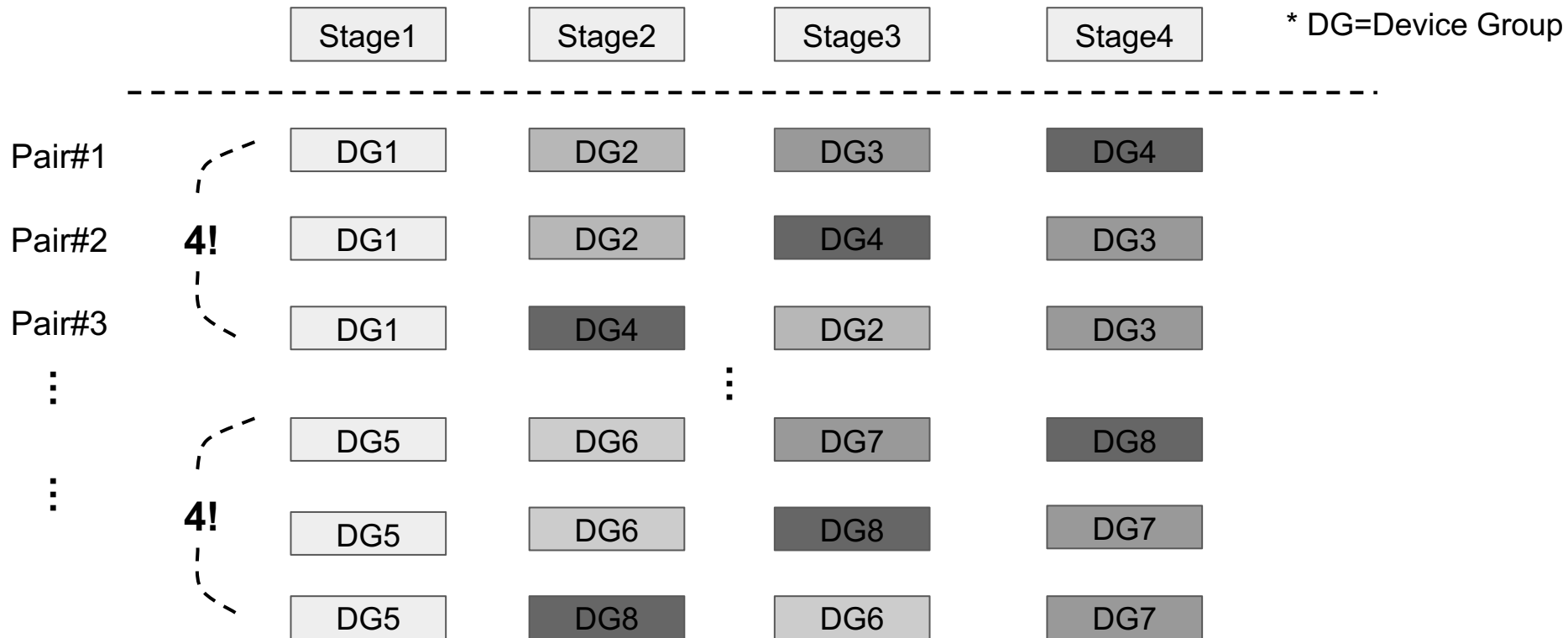
- 1) *Pruning inefficient/similar combinations of stage-device group pairs*
- 2) *Balancing layers across stages with capacity-aware allocation*
- 3) *Navigating efficient intra-stage plans based on DP/TP characteristics*



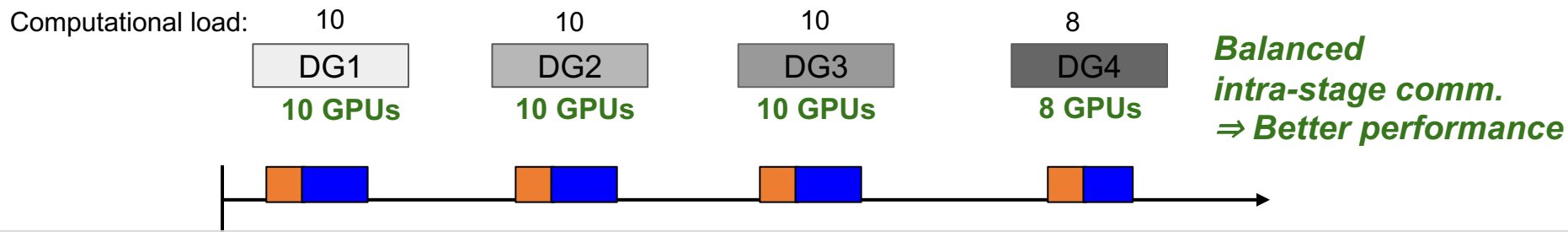
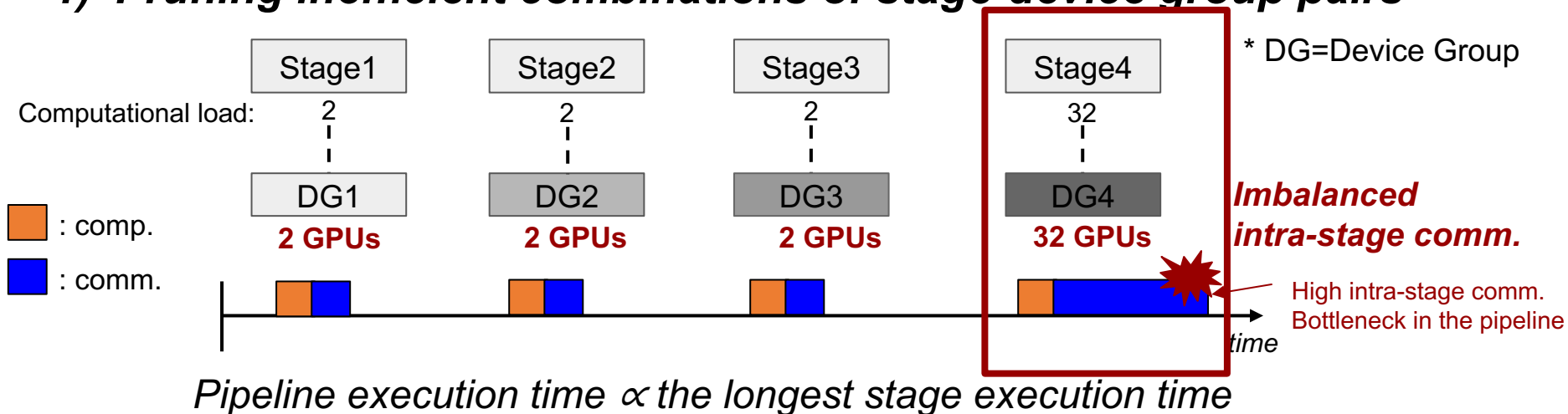
- 1) **Pruning inefficient/similar combinations of stage-device group pairs**
- 2) *Balancing layers across stages with capacity-aware allocation*
- 3) *Navigating efficient intra-stage plans based on DP/TP characteristics*



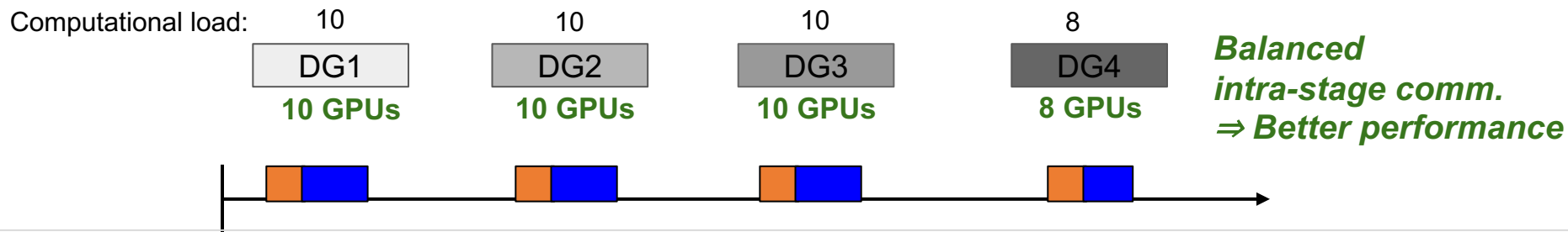
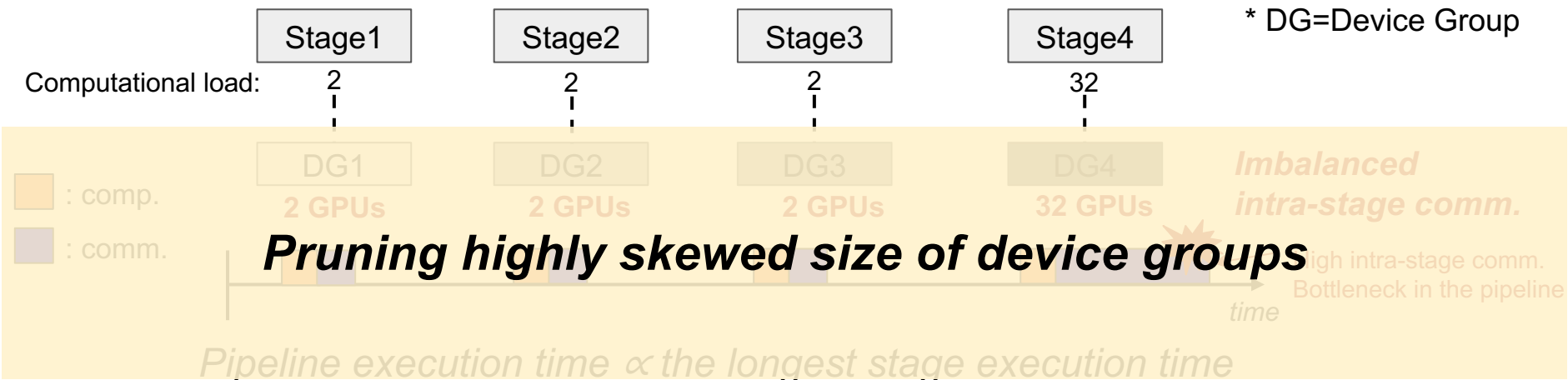
1) Pruning inefficient combinations of stage-device group pairs



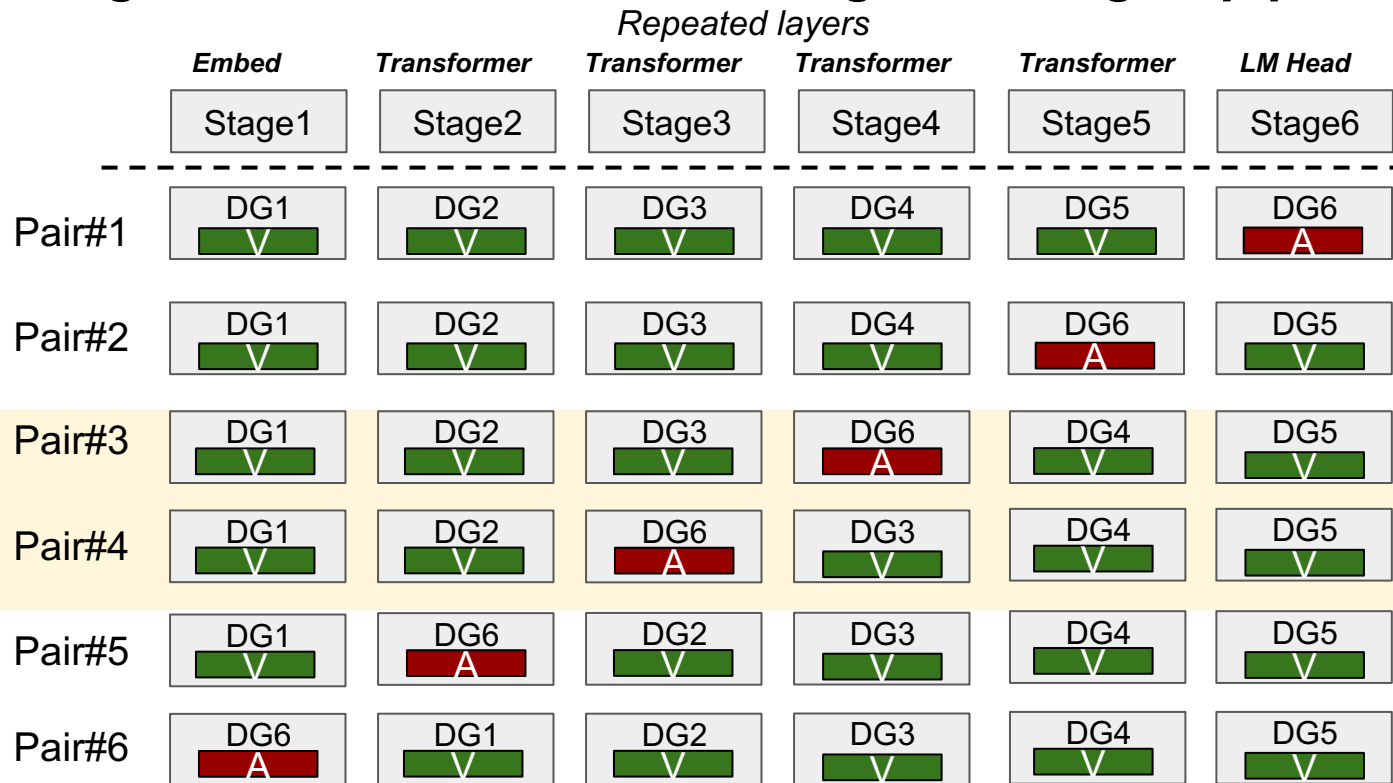
1) Pruning inefficient combinations of stage-device group pairs



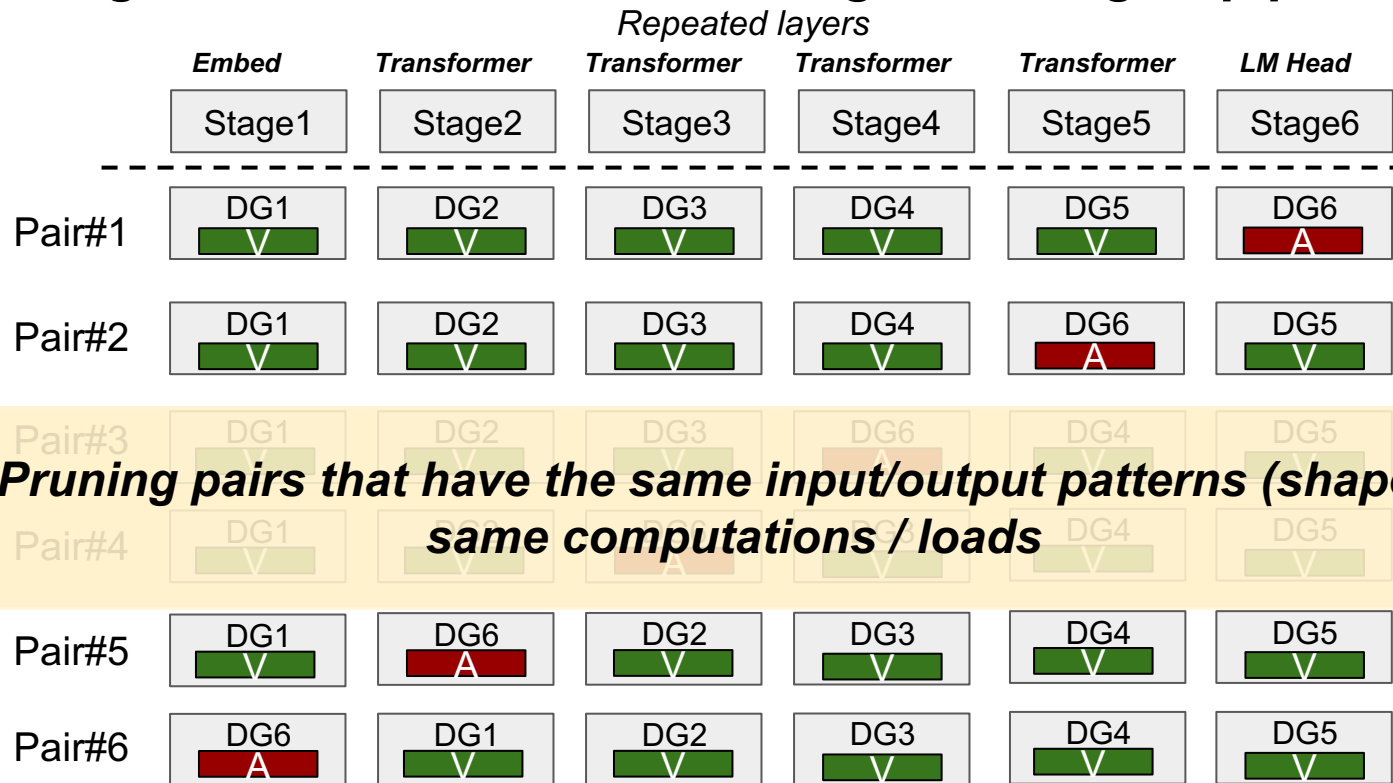
1) Pruning inefficient combinations of stage-device group pairs



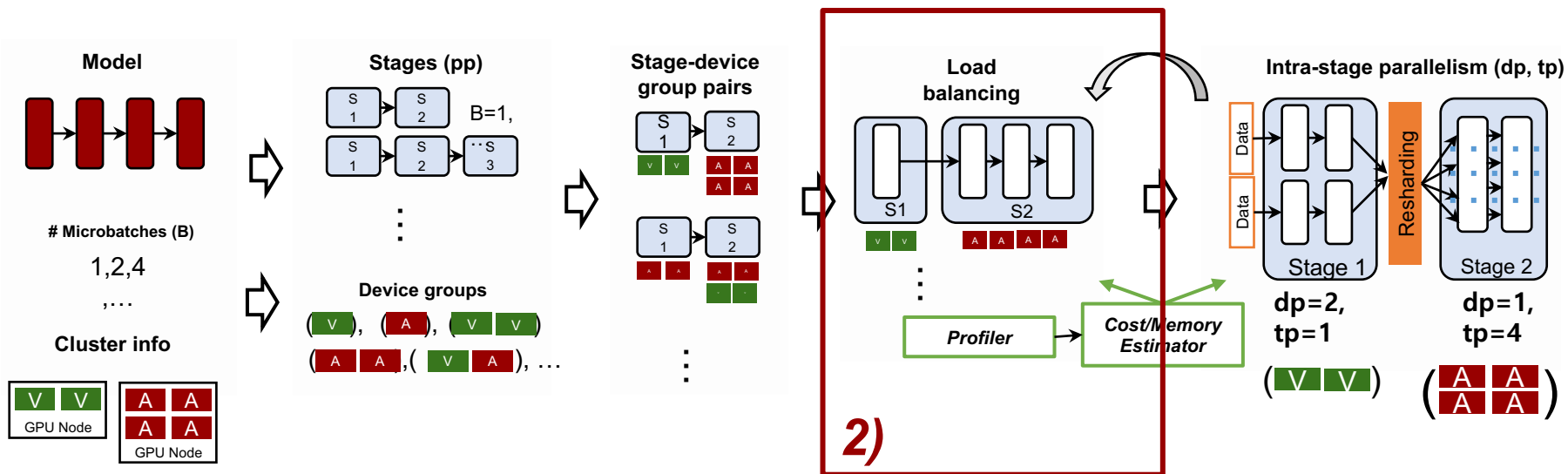
1) Pruning similar combinations of stage-device group pairs



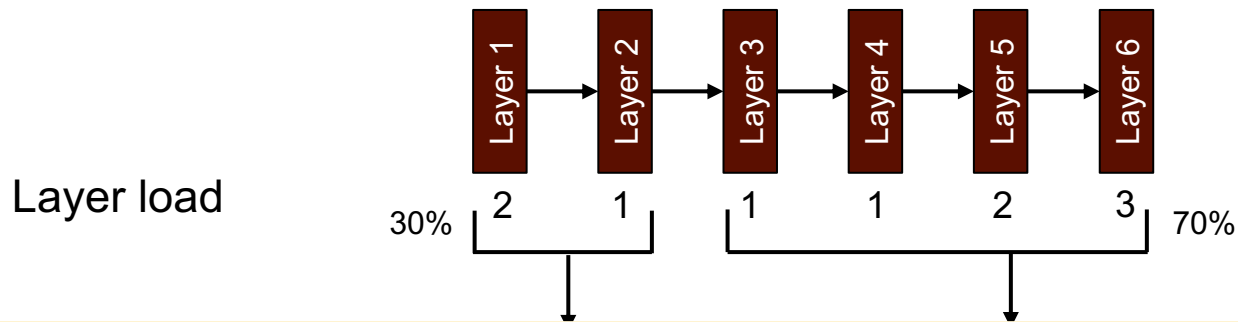
1) Pruning similar combinations of stage-device group pairs



- 1) *Pruning inefficient/similar combinations of stage-device group pairs*
- 2) **Balancing layers across stages with capacity-aware allocation**
- 3) *Navigating efficient intra-stage plans based on DP/TP characteristics*



2) *Balancing layers across stages with capacity-aware allocation*



Relative performance

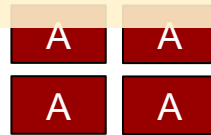
$O(L)$ layer load balancing based on the relative load of layers and performance of device groups

Estimated execution time of a model: 10s



(DP=2, TP=2)

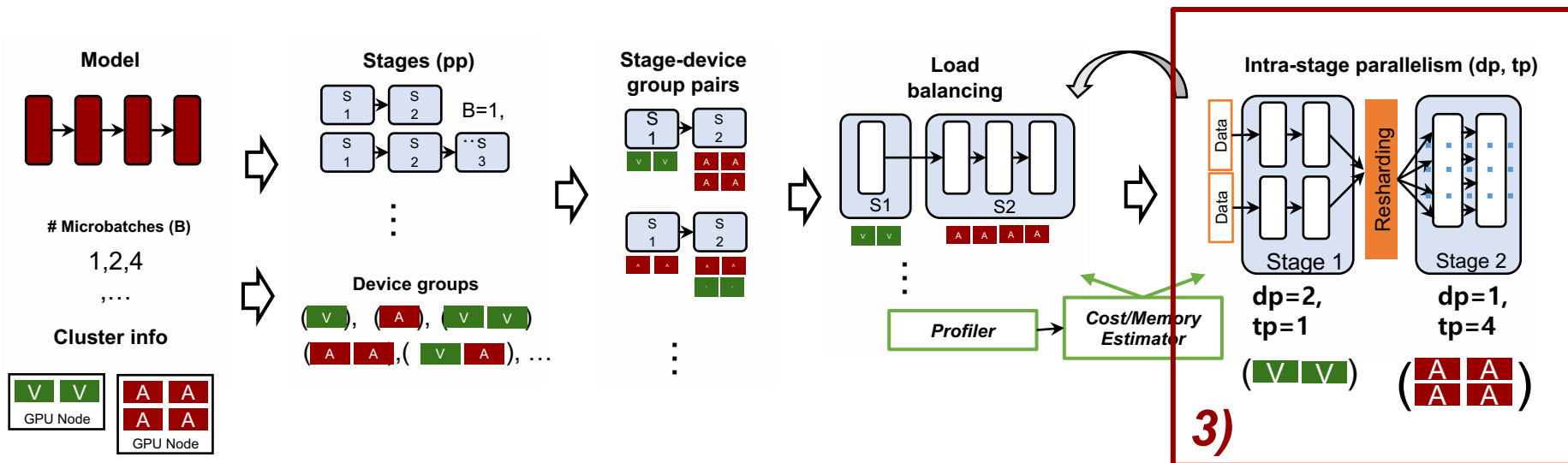
Assigned device groups



(DP=4, TP=1)

Estimated execution time of a model: 5s

- 1) *Pruning inefficient/similar combinations of stage-device group pairs*
- 2) *Balancing layers across stages with capacity-aware allocation*
- 3) ***Navigating efficient intra-stage plans based on DP/TP characteristics***

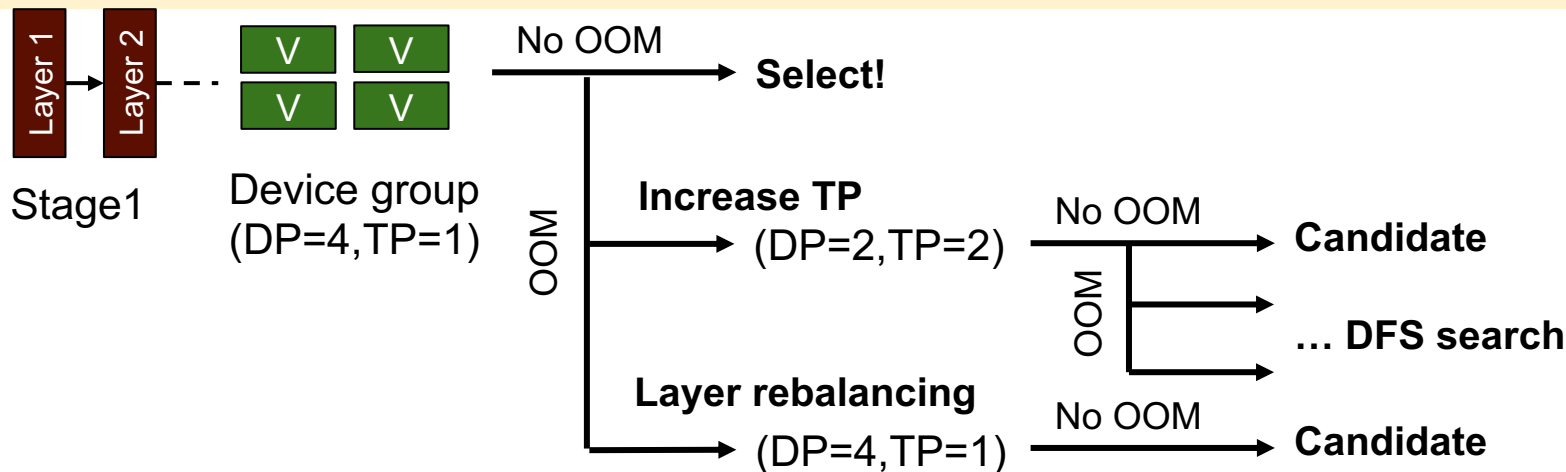


3) Navigating efficient intra-stage plans based on DP/TP characteristics

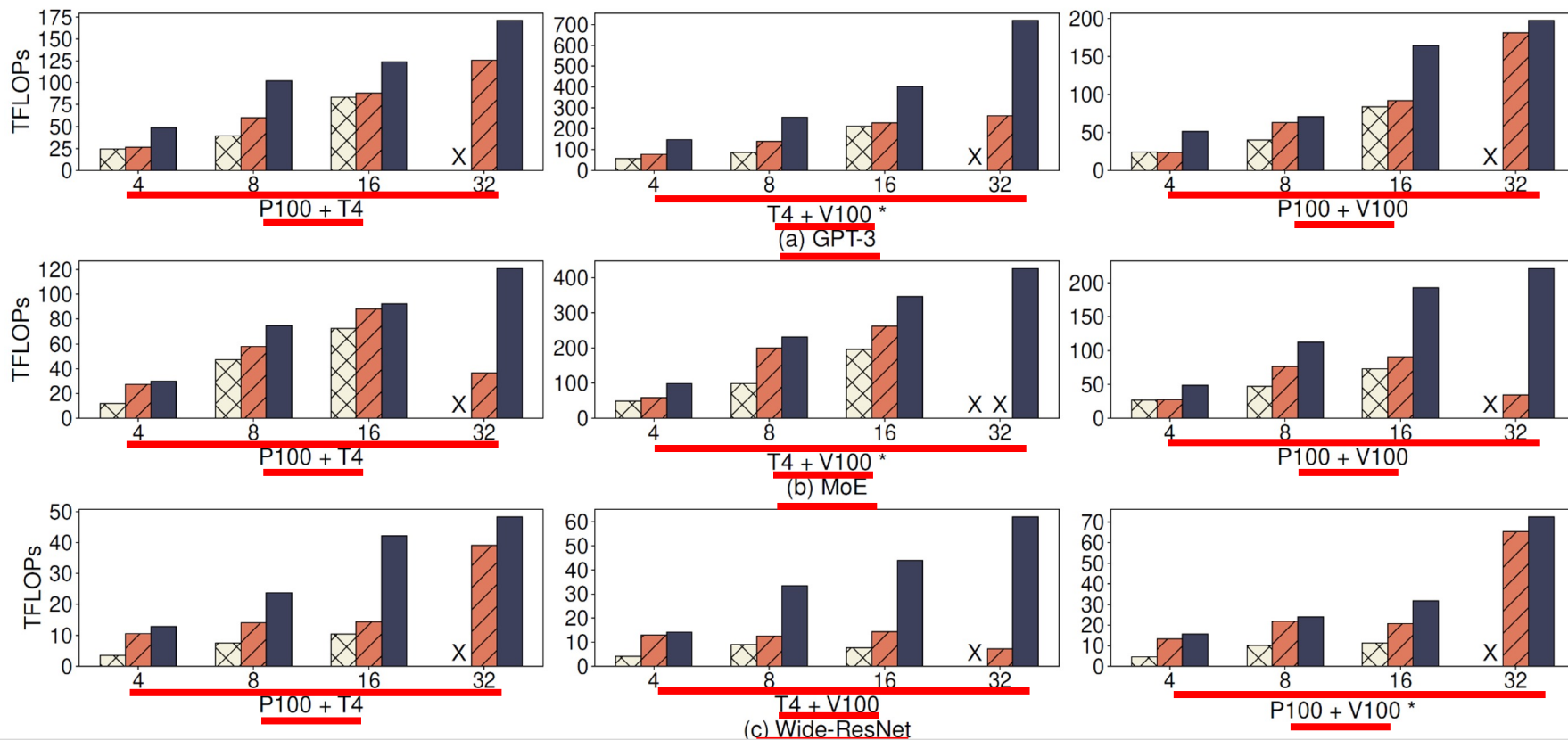
- DP => All-reduce for gradient after mini-batch

Efficient DFS search that prioritizes DP over TP with OOM detection

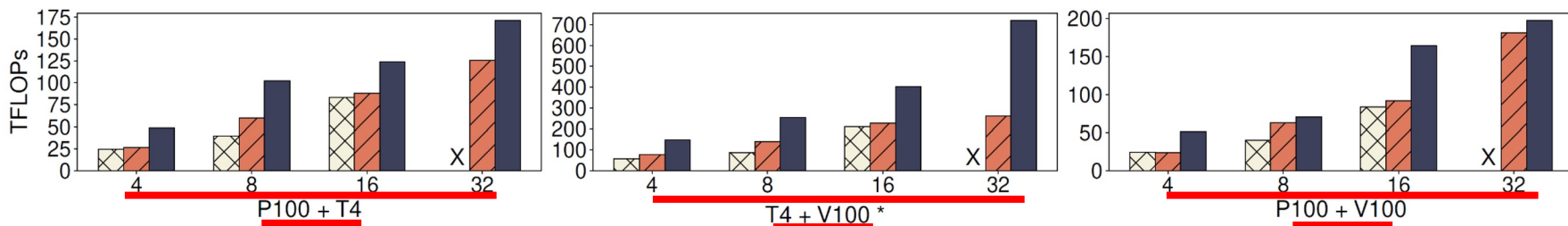
- TP comm. overhead > DP comm. overhead



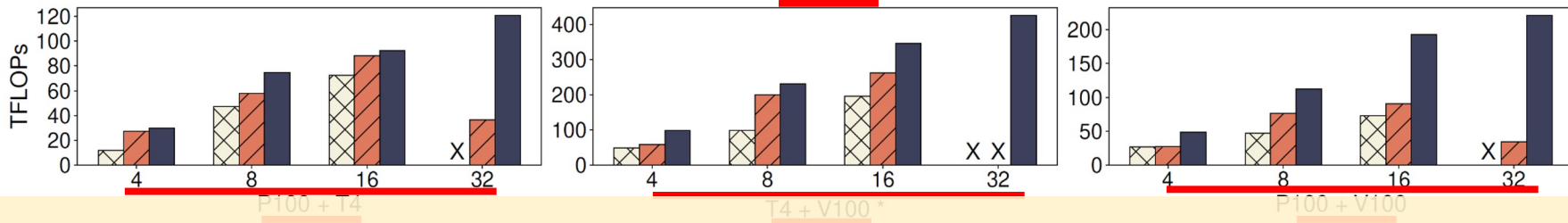
Evaluation: Training Speed



Evaluation: Training Speed

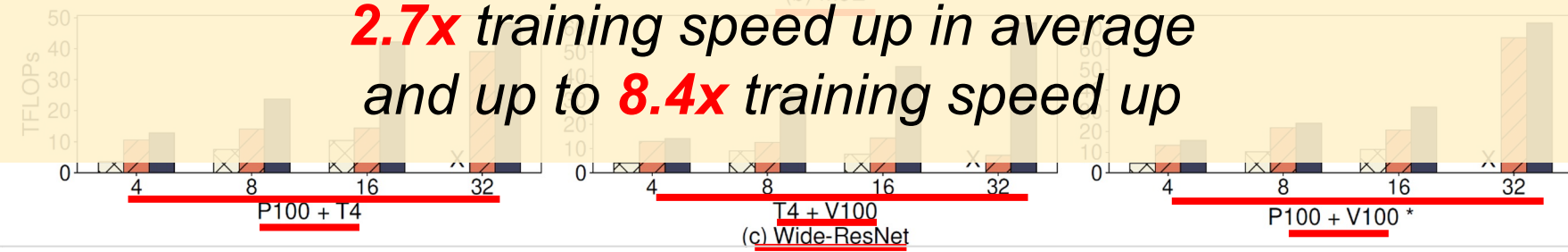


(a) GPT-3



(b) MoE

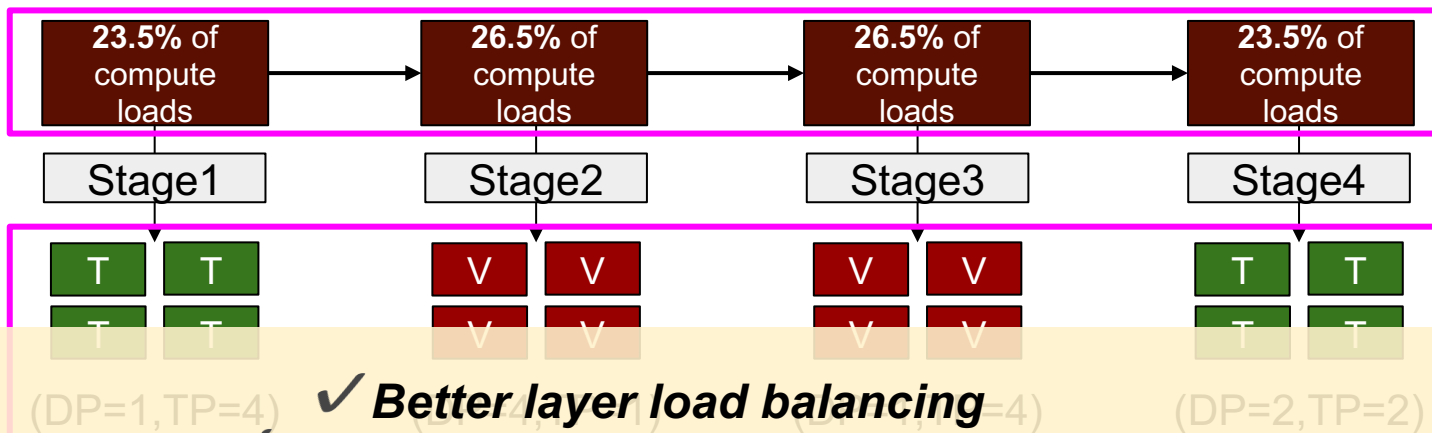
2.7x training speed up in average
and up to **8.4x** training speed up



(c) Wide-ResNet

Evaluation: Plan Comparison * GPT-3, 16 GPUs, T4+V100

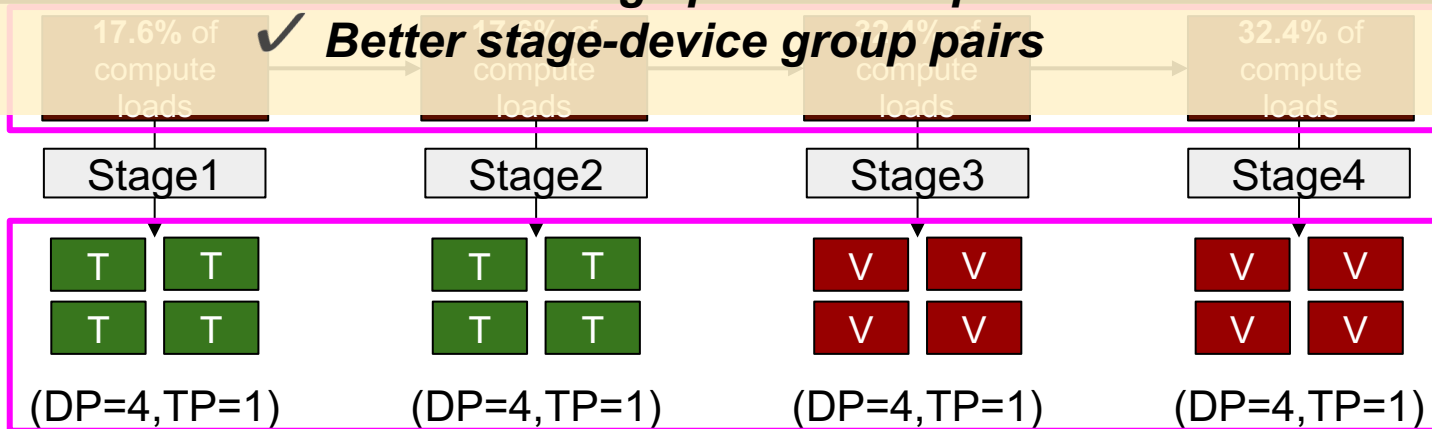
Alpa



V > T
V100 T4

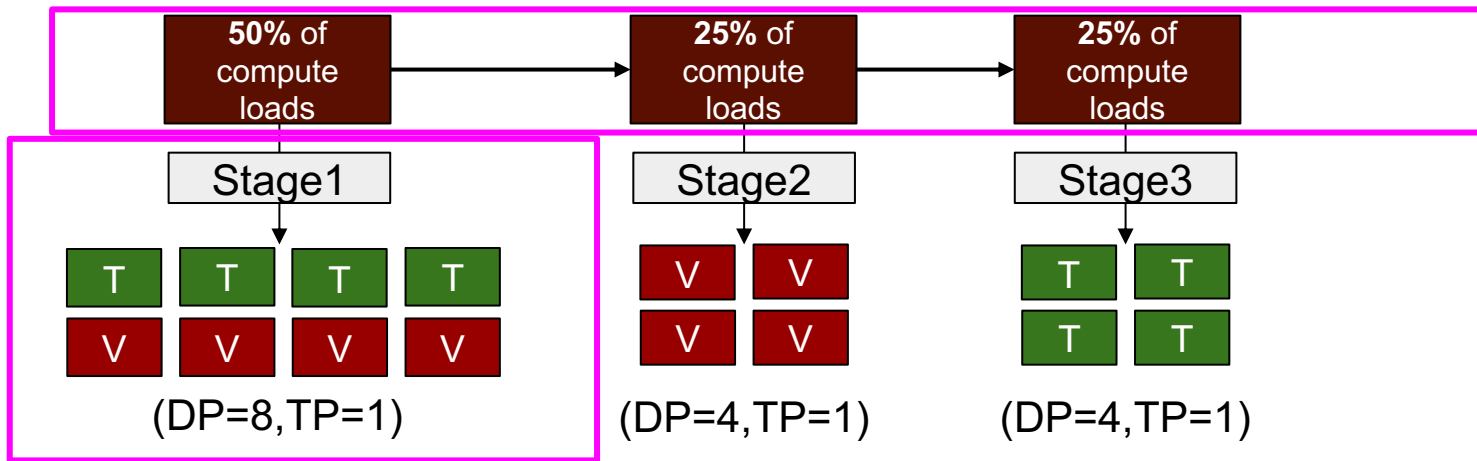
✓ Better stage-device group pairs

Metis

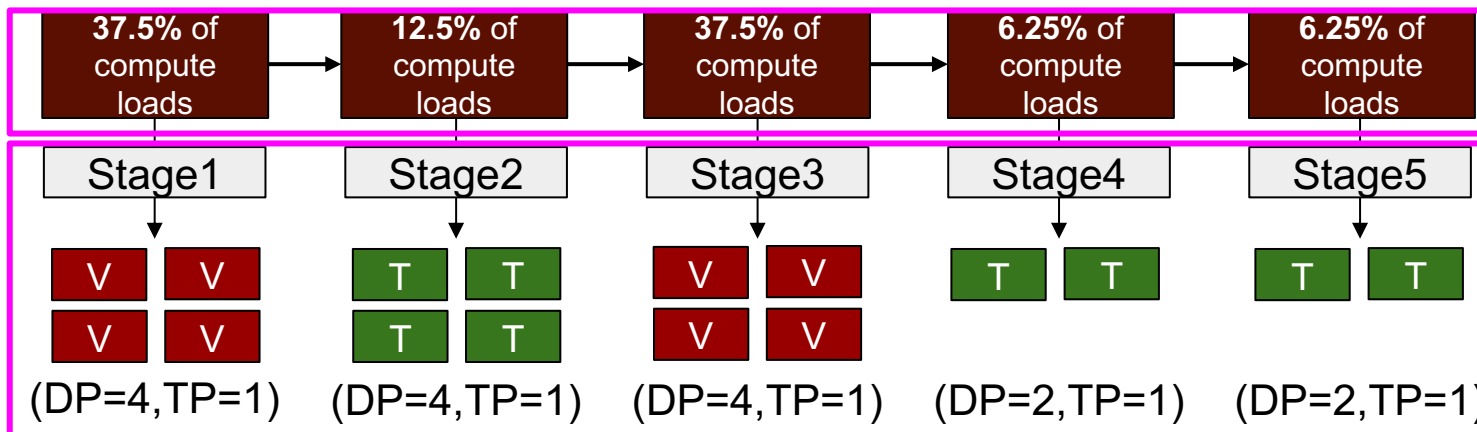


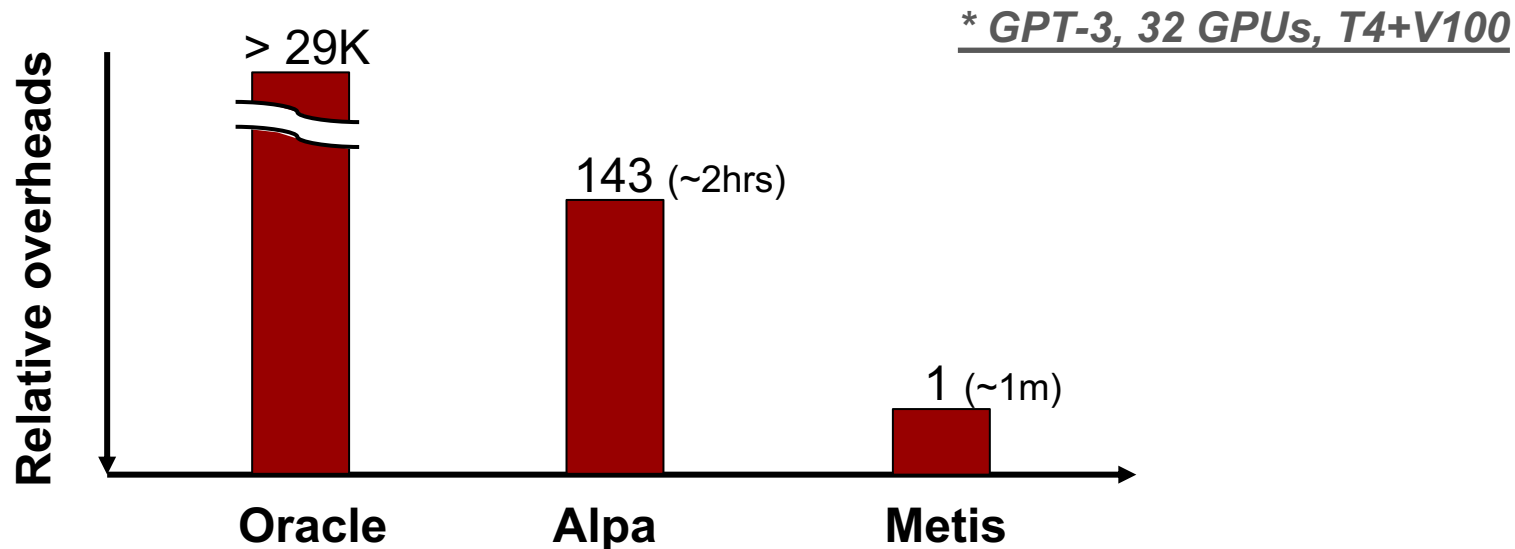
Evaluation: Plan Comparison * MoE, 16 GPUs, T4+V100

Alpa



Metis





143x (vs Alpa) and **29000x** (vs Oracle) search speed-up while finding near-optimal plans (5% diff compared to Oracle)

Conclusion

- *Show* that the **existing auto-parallelizer is not optimized for distributed training on heterogeneous GPUs**
- *Design* **Metis**, a system that automatically finds good parallelism plans on heterogeneous GPUs **by expanding the search space**
- *Develop* a **new hetero-aware search algorithm** that prunes inefficient plans, balances layers based on capacity-awareness, and finds efficient intra-stage parallelism while prioritizing DP over TP
- *Evaluate* that Metis **improves training speed by up to 8x with less search overheads compared to SoTA**
- *Believe* that Metis can be adopted to **improve training speed and GPU utilization** and used for **cost-efficient distributed training** on heterogeneous GPU clusters

Thank you

taegeon.um@samsung.com

<https://taegeonum.github.io>