

# HexiScale: Facilitating Large Language Model Training over Heterogeneous Hardware

---

Ran Yan<sup>1\*</sup>, Youhe Jiang<sup>1\*</sup>, Xiaonan Nie<sup>2</sup>, Fangcheng Fu<sup>3</sup>, Bin Cui<sup>2</sup>, Binhang Yuan<sup>1</sup>

1. HKUST, 2. Peking University, 3. Shanghai Jiao Tong University



# Introduction & Motivation

---

## The homogeneous GPUs are expensive:

- Training LLMs (e.g., Qwen, Llama) demands thousands of GPUs, while homogeneous high-end GPUs are scarce and expensive.

## The existing resources are underutilized:

- GPU generations update fast: Turing (2018) → Ampere (2020) → Hopper (2022) → Blackwell (2024).
- Older GPUs remain in service for years (e.g., K80 still on AWS p2 instances).

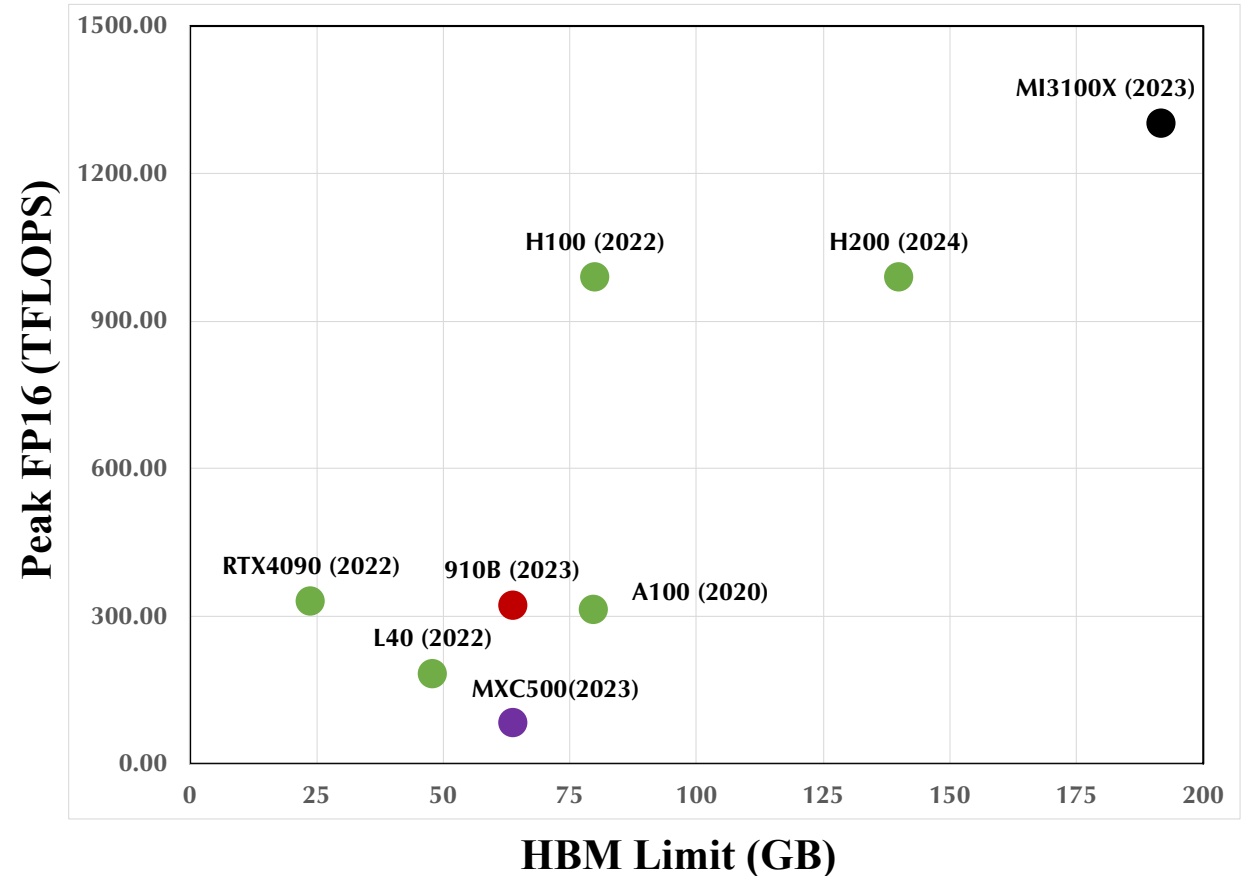
## Key question:

How can we achieve high performance for LLM training using heterogeneous GPUs?

# Challenges

## Challenge 1: Heterogeneous GPUs have varied hardware specs:

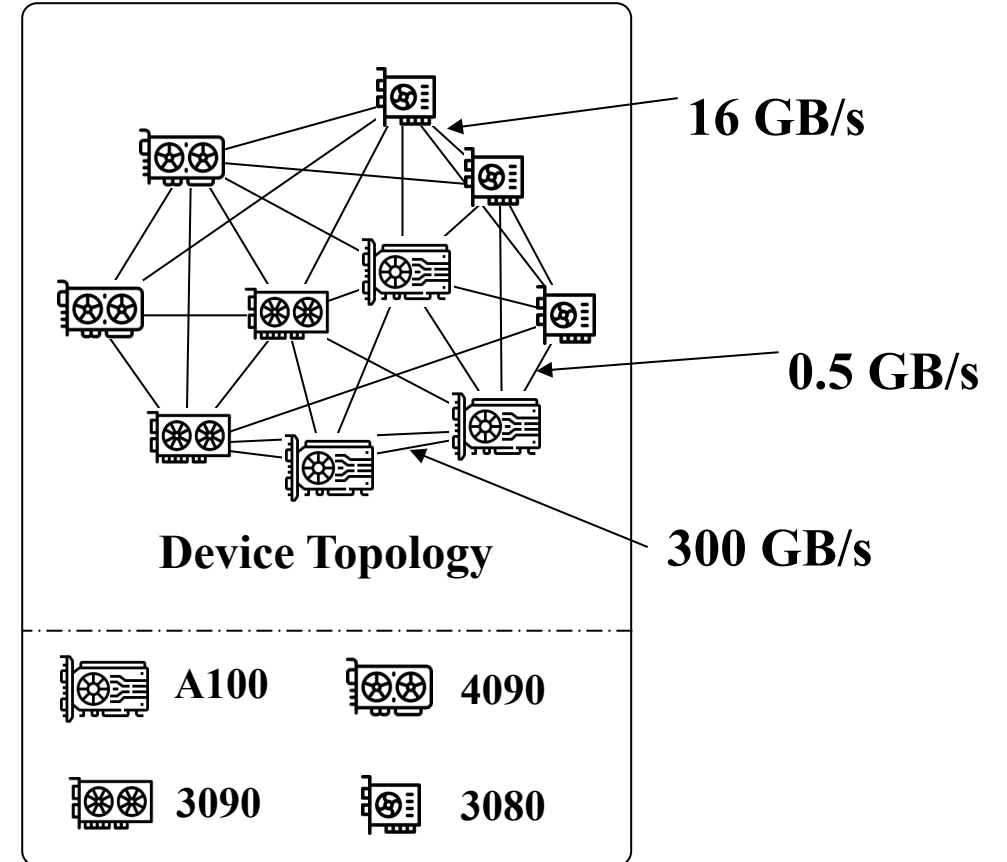
- Compute and memory resources demonstrate high level of heterogeneity.
- More capable GPUs may be underutilized without proper management.



# Challenges

## Challenge 2: Heterogeneous GPUs have varied network bandwidth:

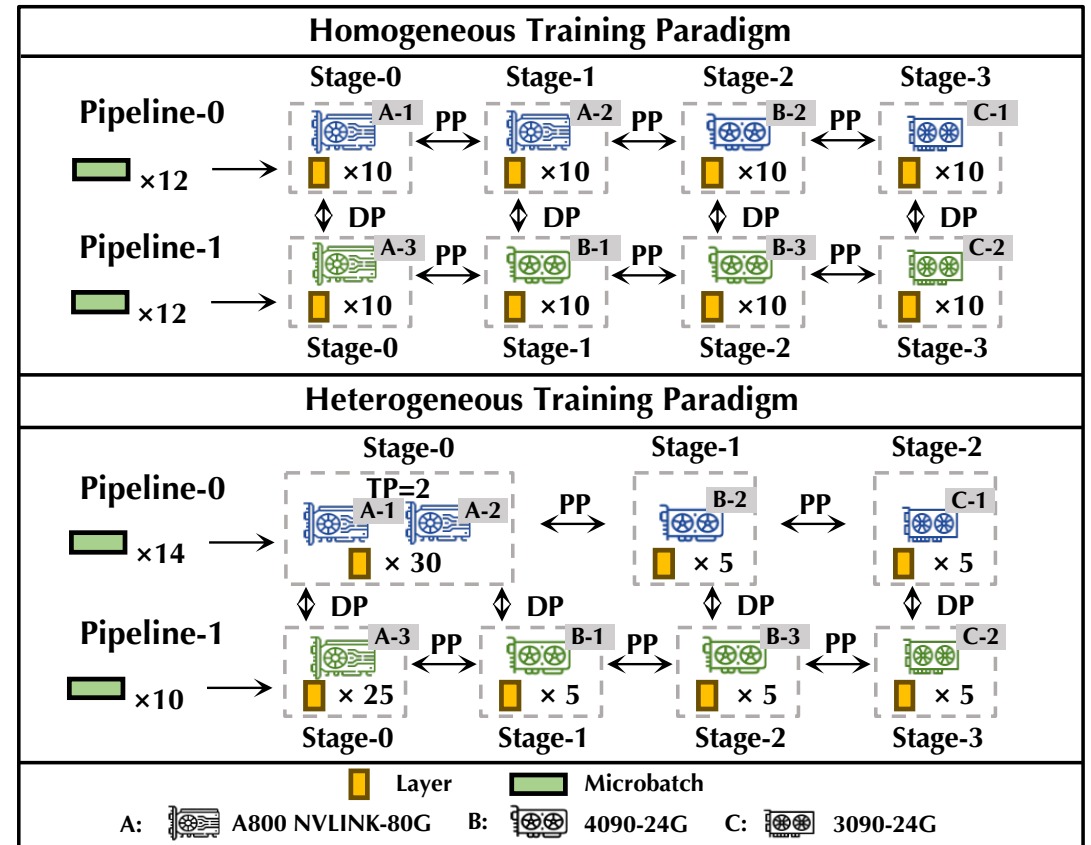
- GPU-GPU network bandwidth can vary significantly.
- Effective management of communication overhead is essential to prevent faster connections from stalling.



# System Design

## Asymmetric parallel support of HexiScale:

- Within one pipeline, PP stages can have varied TP degrees.
- PP stages from different pipelines can have distinct TP degrees.
- Each PP stage can handle arbitrary number of layers.
- Each pipeline can handle uneven batch-size.



# Scheduling Algorithm

## Problem Formulation:

- Objective: Given a heterogeneous cluster, we find the optimal placement that have the highest training throughput.
- Formally:

$$\begin{aligned} \sigma^* &= \min_{\sigma} \text{Comm-Cost}(\sigma) + \text{Comp-Cost}(\sigma) \\ \text{s. t. } &\text{Mem-Cumsum}(d) \leq m_d \quad \forall d \in \mathbf{D} \end{aligned}$$

- $d$  is the GPU device,  $m_d$  is the memory limit,  $\mathbf{D}$  is the heterogeneous devices set.
- $\sigma$  is the parallel configuration,  $\sigma^*$  is the optimal parallel configuration.
- $\text{Comm-Cost}(\sigma)$  is the communication cost,  $\text{Comp-Cost}(\sigma)$  is the computation cost,  $\text{Mem-Cumsum}(d)$  is the memory cost of device  $d$ .

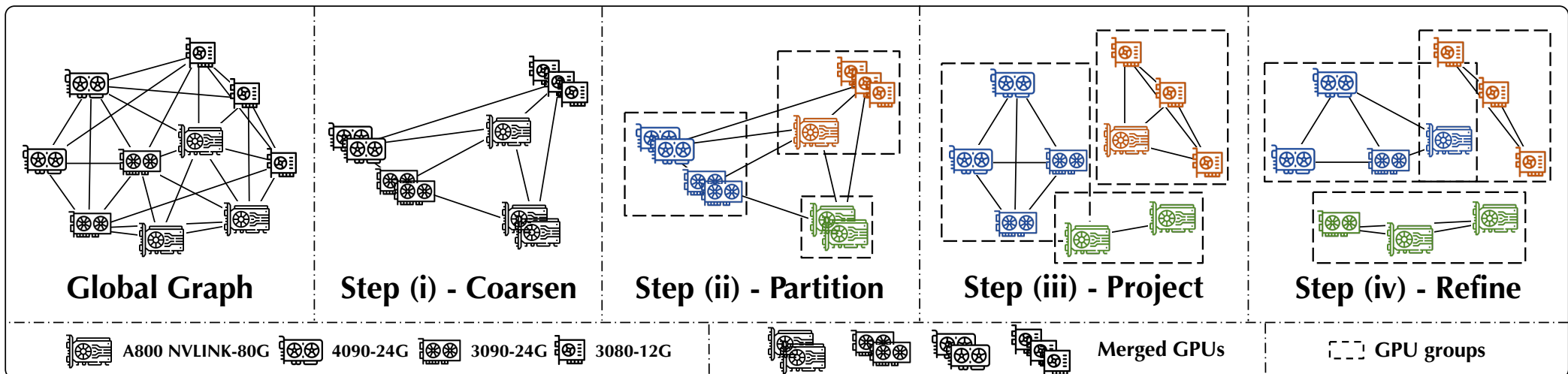
# Scheduling Algorithm

## Phase-1, global graph partitioning:

- Model GPUs as graph: vertices weighted by FLOPS, edges by pairwise bandwidth.
- Partition into  $D_{dp}$  groups (each for one pipeline) via 4-step graph partition algorithm.

## Iterative optimization from phase-1:

- Enumerate  $D_{dp}$  to optimize the number of pipelines.
- Adaptively allocate bandwidth for DP and PP.



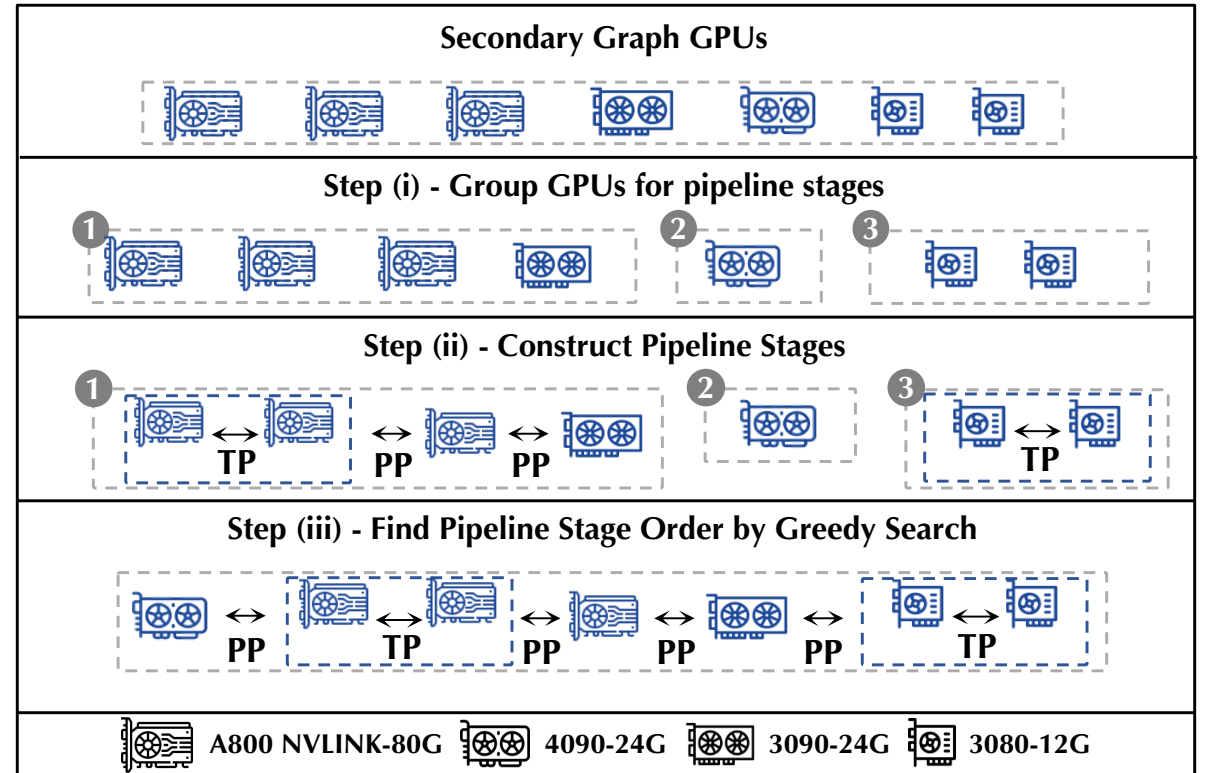
# Scheduling Algorithm

## Phase-2, pipeline construction:

- (i) Group high-bandwidth GPUs via secondary graph partition.
- (ii) Search intra-group TP/PP/layer strategy per machine.
- (iii) Top- $\tau$  greedy search for optimal stage ordering.

## Iterative optimization from phase-2:

- Construct different GPU groups through different secondary graph partition.



# Evaluations

## Compare with baselines over homogeneous clusters:

- Baseline frameworks: Megatron, Galvatron and FSDP.
- Comparison scheme: Under the same total FLOPS, we compare the PFLOPS of HexiScale over *heterogeneous clusters* and baselines over *homogeneous clusters*.

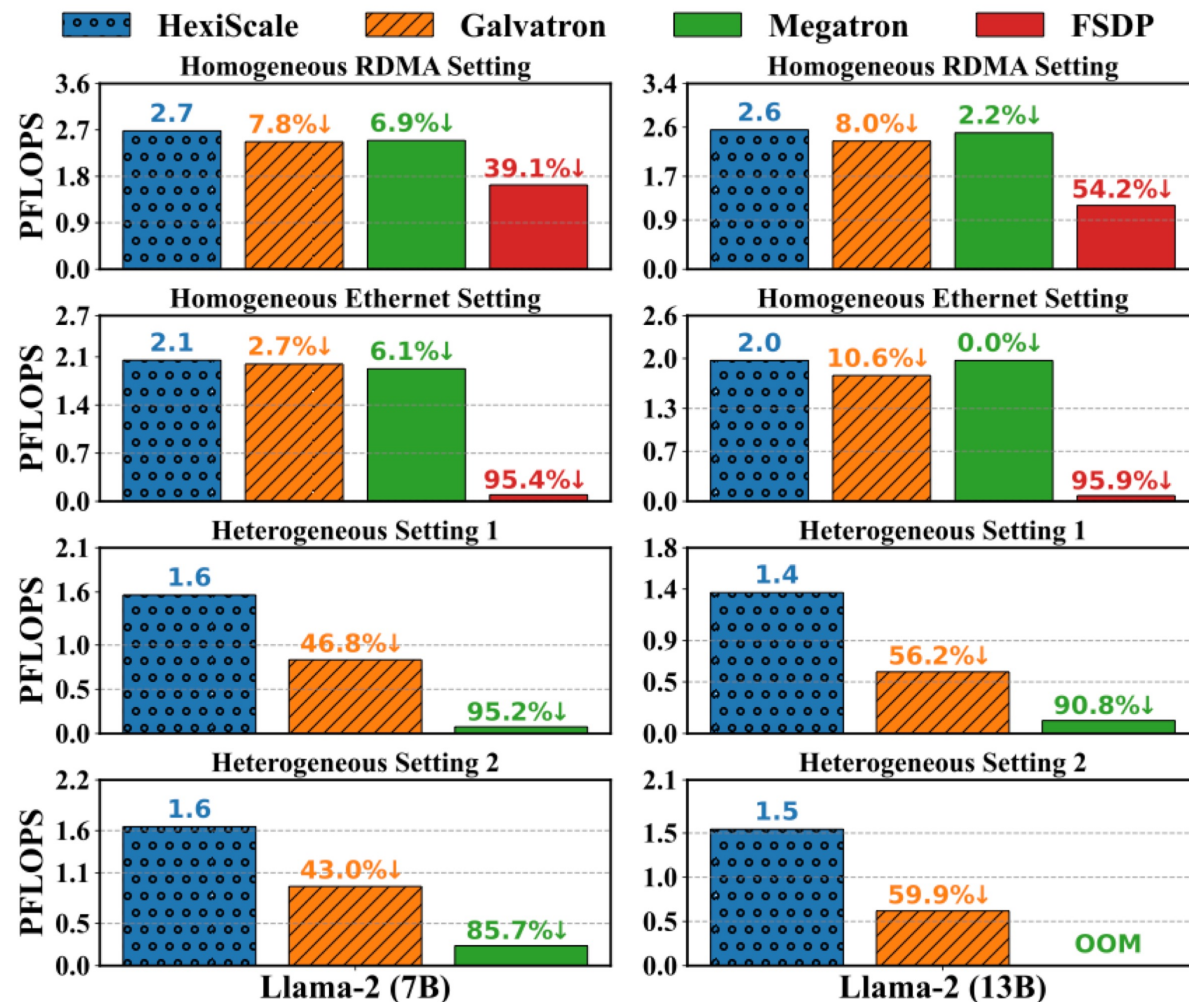
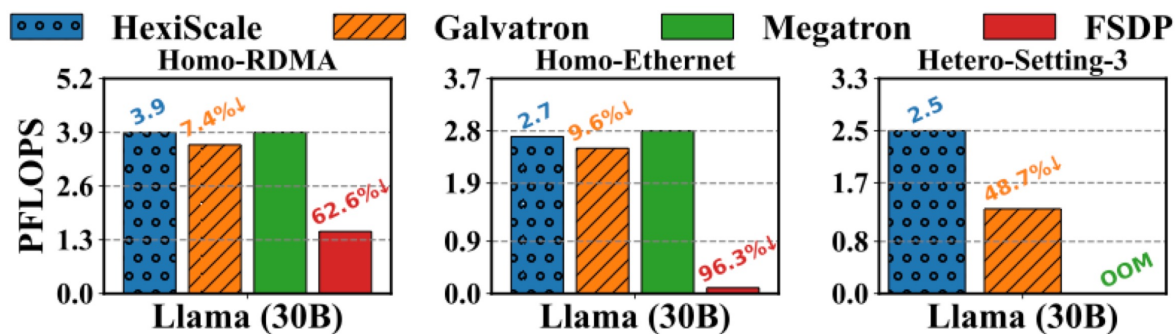
## Experimental settings:

Setting	GPU Configuration	FLOPS vs Baseline	Target Homogeneous Baseline Cluster	Target Model(s)
<b>Hetero-Setting 1</b>	1×(8×3080Ti) + 1×(8×3090) + 3×(8×4090)	+1.36% higher total FLOPS	2×(8×A800 PCIe-80G)	Llama-2 7B / 13B
<b>Hetero-Setting 2</b>	1×(8×3080Ti) + 1×(8×3090) + 1×(8×4090) + 1×(8×A800 NVLINK-80G)	-1.59% lower total FLOPS	2×(8×A800 PCIe-80G)	Llama-2 7B / 13B
<b>Hetero-Setting 3</b>	1×(8×3090) + 2×(4×3090) + 4×(8×4090) + 1×(8×A800 NVLINK-80G)	-4.67% lower total FLOPS	4×(8×A800 PCIe-80G)	Llama 30B

# Evaluations

## Evaluation results against homogeneous baselines:

- HexiScale exhibits up to  $1.01 \times$  throughput and an average throughput gap of  $0.83 \times$  compared to homogeneous scenarios.

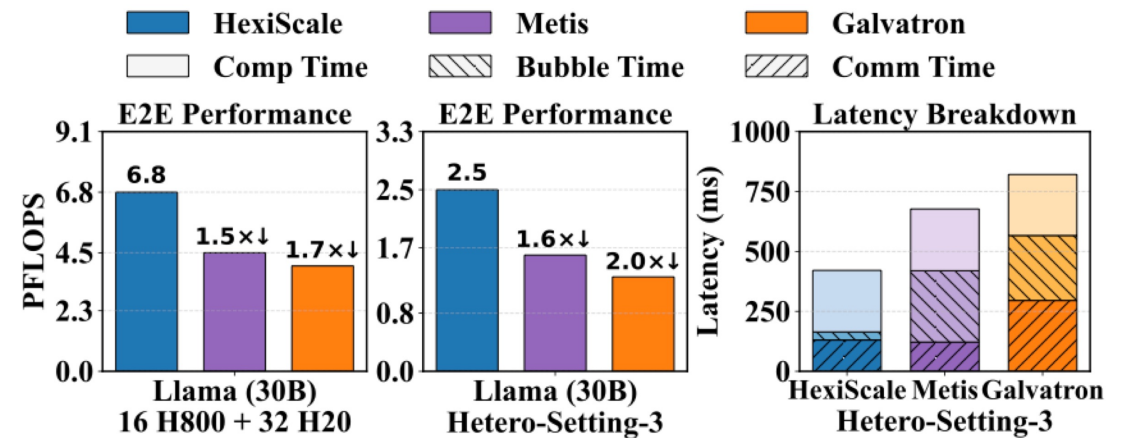
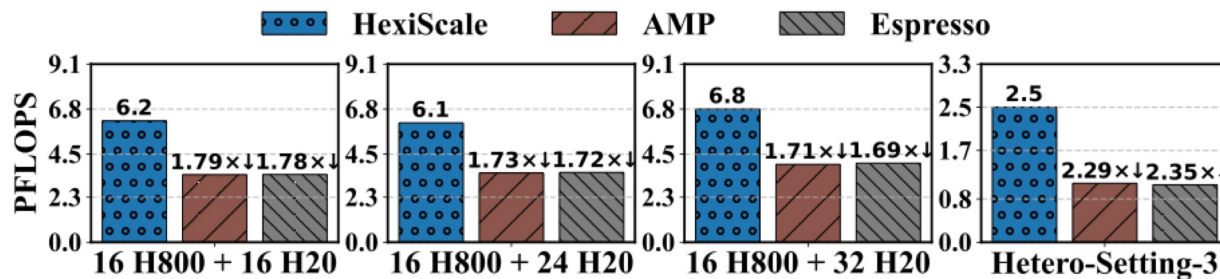


# Evaluations

Compare with heterogeneity-aware baselines over heterogeneous clusters:

- Baseline frameworks: Metis, AMP, and Espresso.
- Comparison scheme: Under the same heterogeneous clusters, we compare the PFLOPS of HexiScale against the heterogeneity-aware baselines.

Experimental results:



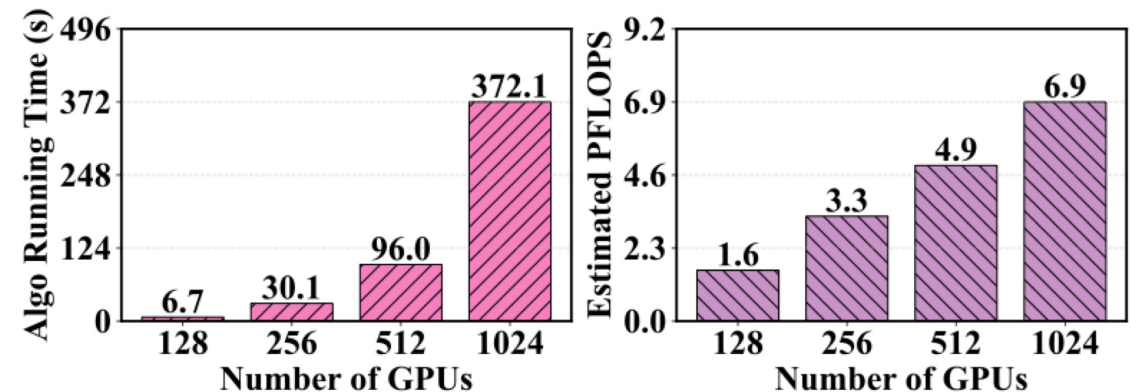
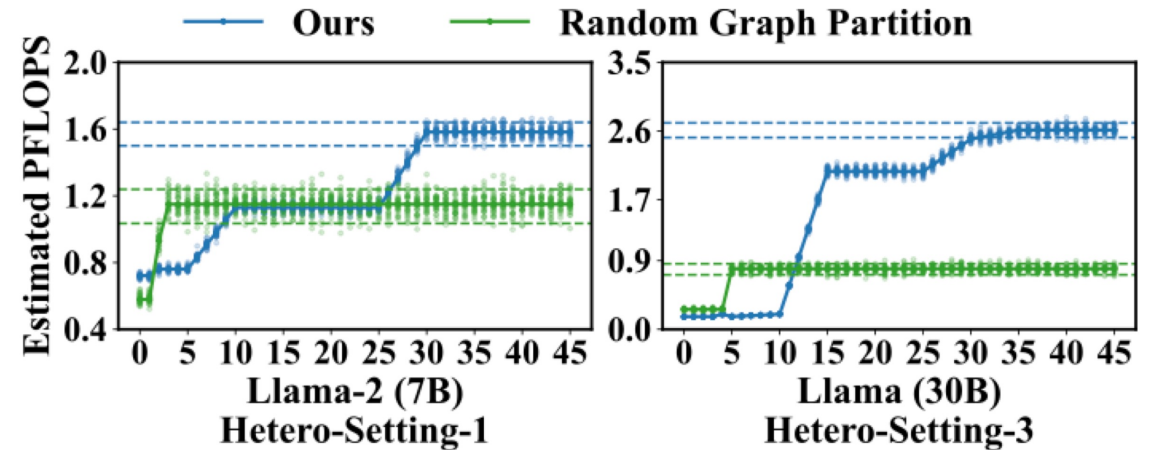
# Evaluations

## Algorithm effectiveness evaluation:

- Compared to Random Graph Partition algorithm, HexiScale's algorithm demonstrates earlier convergence and 1.3x – 3.3x higher estimated throughput.

## Algorithm efficiency evaluation:

- HexiScale's scheduling time for large clusters remain manageable.
- HexiScale demonstrates strong scalability across large GPU clusters.



# Summary

---

## **Benefit of using heterogeneous GPUs for LLM training:**

- Integrate computational resources from different GPU generations.

## **Key challenges:**

- Heterogeneity in hardware specs and network connections may lead to underutilization of resources.

## **HexiScale solutions:**

- Design and implement flexible asymmetric parallelism for heterogeneous GPUs.
- Design a novel graph partition algorithm to identify the efficient parallel strategies.

# **Thank You**

## Questions & Discussion

---

Contact: [ryanaf@connect.ust.hk](mailto:ryanaf@connect.ust.hk)