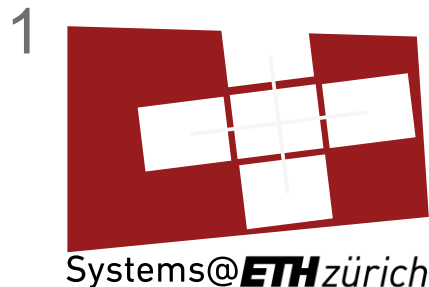


HetRL: Efficient Reinforcement Learning for LLMs in Heterogeneous Environments

Yongjun He¹ (yongjun.he@inf.ethz.ch), Shuai Zhang², Jiading Gai², Xiyuan Zhang², Boran Han², Bernie Wang², Huzefa Rangwala², George Karypis²



The Rise of RL for LLMs

RL has become a key component of post-training pipelines for LLMs.

- Human Alignment with RLHF
- Reasoning Tasks (e.g., math and coding)
- Agentic Tasks (multi-turn, tool-use, environment interactions)

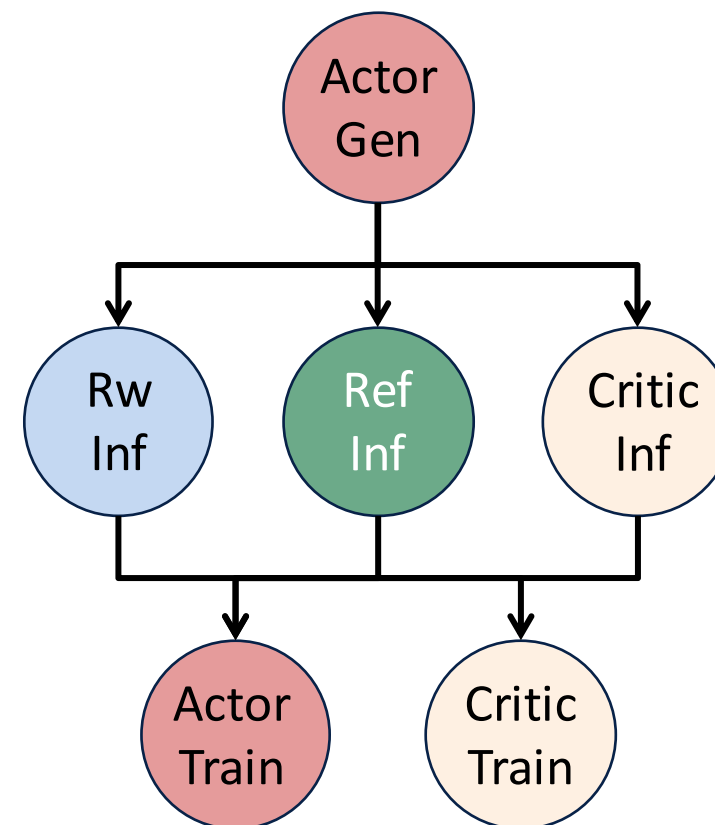


...

Heterogeneity in RL Workflows

Observations

- Multiple models and tasks with complex computational and data dependencies
 - RL models: actor model, reference model, reward model, and critic model
 - RL tasks: actor generation, reference inference, reward inference, critic inference, actor training, and critic training.

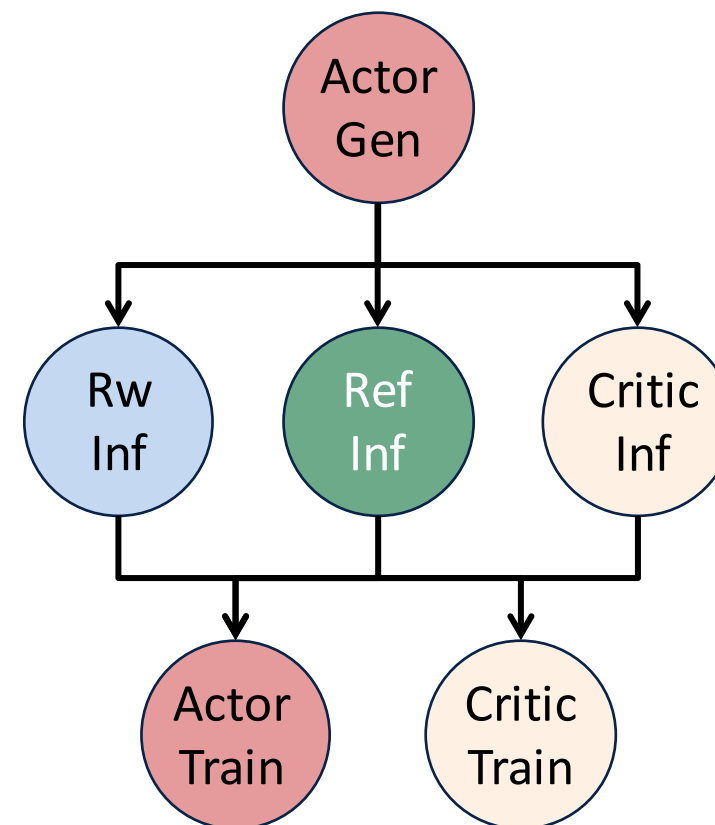


Proximal policy optimization (PPO)

Heterogeneity in RL Workflows

Observations

- Multiple models and tasks with complex computational and data dependencies
- Heterogeneous computational characteristics across tasks
 - Generation task: memory-bound
 - Inference task: compute-bound
 - Training task: compute-bound & inter-node communication heavy

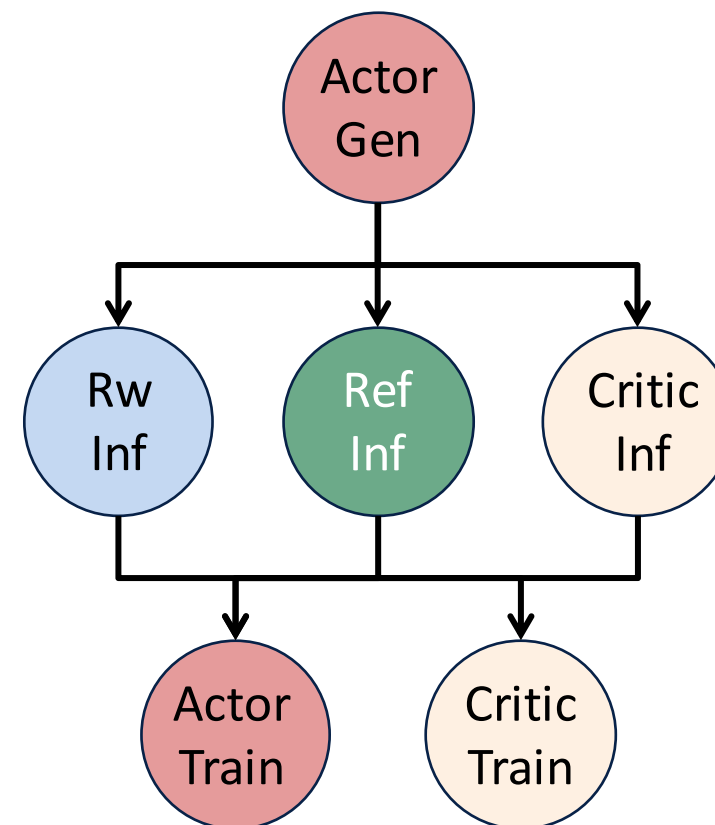


Proximal policy optimization (PPO)

Heterogeneity in RL Workflows

Observations

- Multiple models and tasks with complex computational and data dependencies
- Heterogeneous computational characteristics across tasks
- Heterogeneous parallelization strategies across tasks



Proximal policy optimization (PPO)

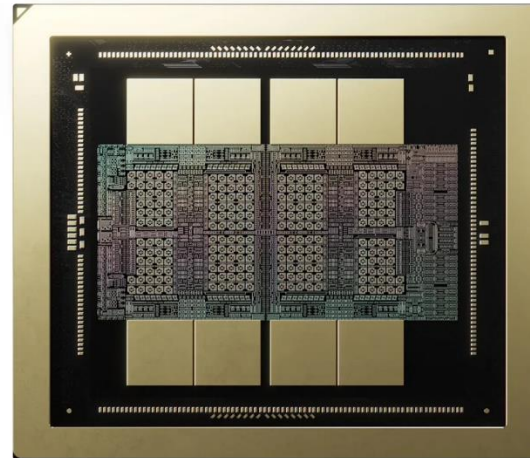
Heterogeneity in Computing Environments

Observations

- Heterogeneous device capabilities (e.g., compute, memory capacity, and memory bandwidth).
- Heterogeneous interconnects (e.g., network bandwidth and latency).

Uniting Processors of Extreme FLOPS and Bandwidth

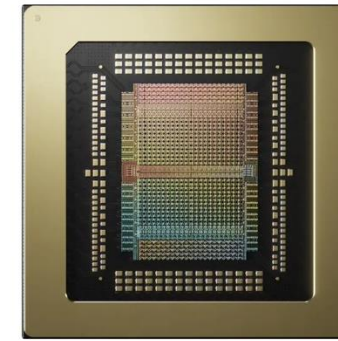
NVIDIA Rubin GPU + Groq 3 LPU



288 GB HBM4

22 TB/s

50 PFLOPs (NVFP4)



500 MB SRAM

150 TB/s SRAM Bandwidth

1.2 PFLOPs (FP8)

Heterogeneity in Computing Environments

Observations

- Heterogeneous device capabilities (e.g., compute, memory capacity, and memory bandwidth).
- Heterogeneous interconnects (e.g., network bandwidth and latency).

Uniting Processors of Extreme FLOPS and Bandwidth

NVIDIA Rubin GPU + Groq 3 LPU



Cerebras is coming to AWS

aws cerebras

AWS TRAINIUM 3 CEREBRAS WSE 3

Heterogeneity in Computing Environments

Observations

- Heterogeneous device capabilities (e.g., compute, memory capacity, and memory bandwidth).
- Heterogeneous interconnects (e.g., network bandwidth and latency).
- Heterogeneous resource availability across regions (e.g., device type and quantity).

Uniting Processors of Extreme FLOPS and Bandwidth

NVIDIA Rubin GPU + Groq 3 LPU



Can we deploy RL training across a set of heterogeneous GPUs connected via heterogeneous networks?

Can we deploy RL training across a set of heterogeneous GPUs connected via heterogeneous networks?

Existing solutions

- Lack support for heterogeneity-aware deployment
- Limited search space
- Time-consuming scheduling algorithms

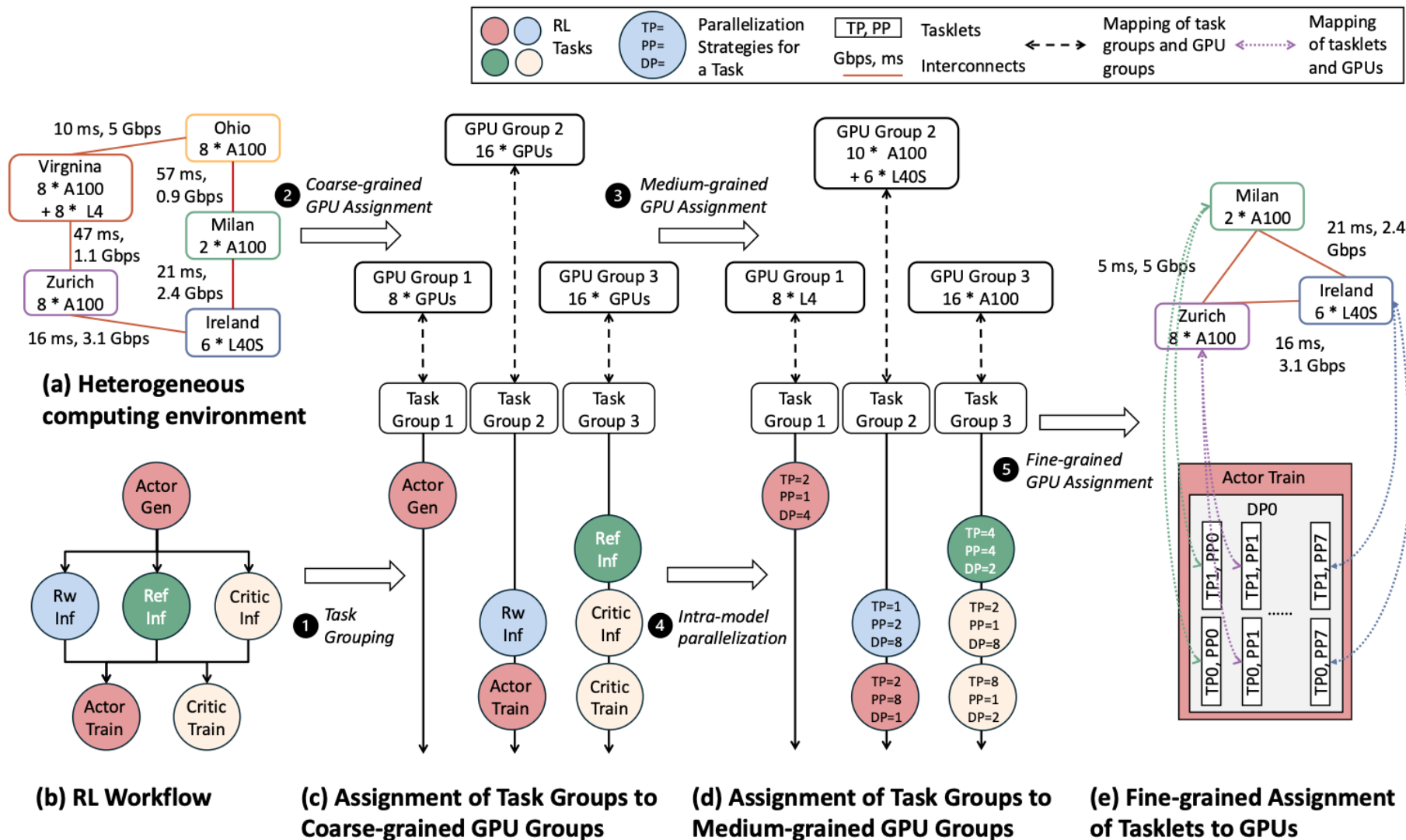
HetRL Methodology

- **Problem formulation:** constrained joint optimization problem
 - **Setting:** given a computational graph G for an RL workflow and a device topology graph G_D for a heterogeneous environment
 - **Goal:** determine an optimal scheduling strategy, which consists of a partitioning strategy ρ and an assignment strategy σ , such that the execution time of the RL workflow is minimized and resource constraints are satisfied.
 - **Hardness:** NP-hard

HetRL Methodology

- Problem formulation: constrained joint optimization problem
- Multi-Level search framework: coarse-to-fine constructive approach that operationalizes the joint optimization over (ρ, σ)
 - Level 1: task grouping
 - Level 2: coarse-grained GPU assignment
 - Level 3: medium-grained GPU assignment
 - Level 4: intra-model parallelization
 - Level 5: fine-grained GPU assignment

HetRL Overview



HetRL Methodology

- Problem formulation: constrained joint optimization problem
- Multi-Level search framework: coarse-to-fine constructive approach that operationalizes the joint optimization over (ρ, σ)
- Hybrid scheduling algorithm efficiently identifies near-optimal solutions.
 - Successive halving algorithm for search budget allocation at higher levels (Level 1 and 2)
 - Evolutionary algorithm for low-level plan generation

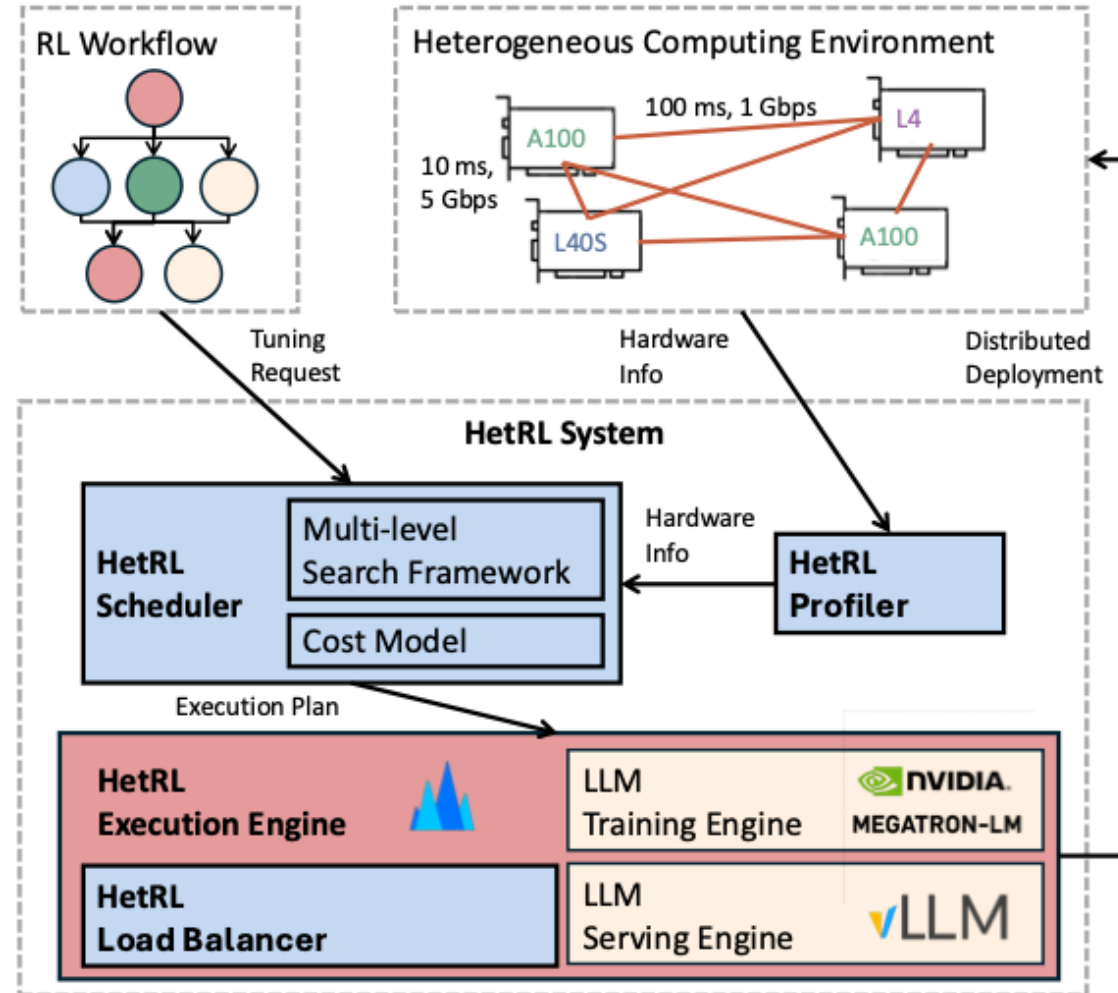
HetRL Methodology

- Problem formulation: constrained joint optimization problem
- Multi-Level search framework: coarse-to-fine constructive approach that operationalizes the joint optimization over (ρ, σ)
- Hybrid scheduling algorithm efficiently identifies near-optimal solutions.
- ILP-based scheduling algorithm provides optimal solutions when budgets allows.

HetRL System

Principal components

- Scheduler
 - Search framework
 - Cost model
- Profiler
- Execution engine
 - Training engine
 - Serving engine
 - Load Balancer



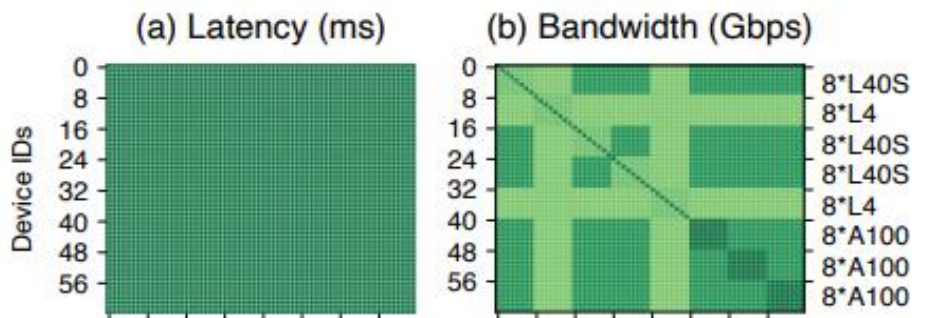
Evaluation Setup

- Models, datasets, and RL algorithms
 - Qwen 2B, 4B, 8B, and 14B
 - GSM8K and Math-500
 - PPO and GRPO (sync and async)
- Baselines: verl and StreamRL.
- Scenarios: single-region, multi-region, multi-country, and multi-continent.
- GPU specs:

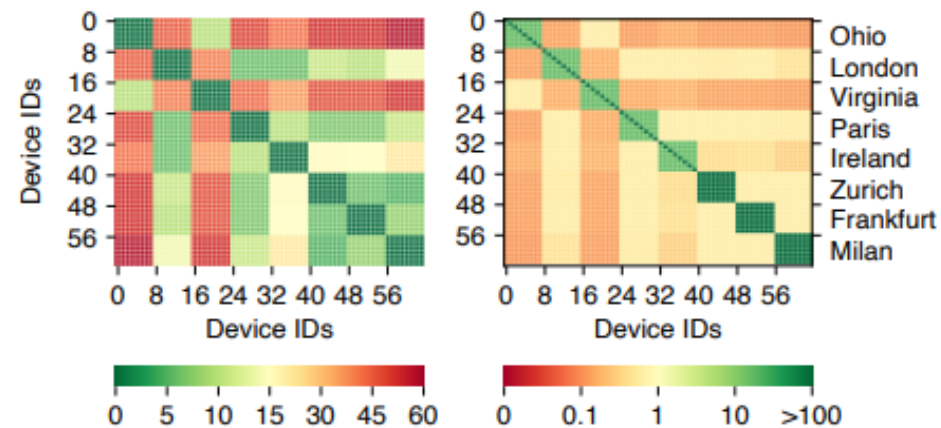
Model	Arch	Size (GB)	FP16 Perf (TFLOPS)	HBM (GB/s)	Intra (GB/s)
A100 (24)	Ampere	40	312	2039	600
L40S (24)	Ada	48	366	864	64
L4 (16)	Ada	24	121	300	64

End-to-End Comparison

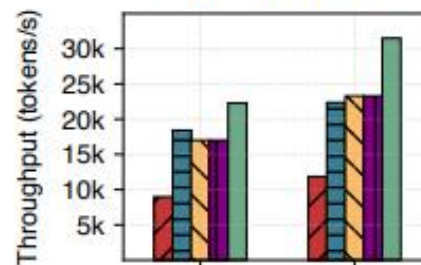
Scenario 1
Single-Region



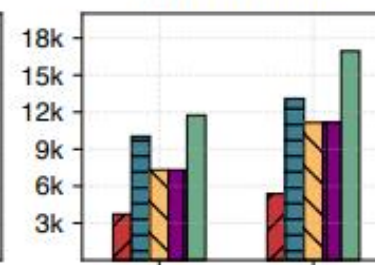
Scenario 4
Multi-Continent



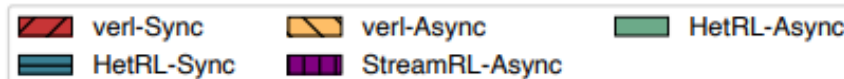
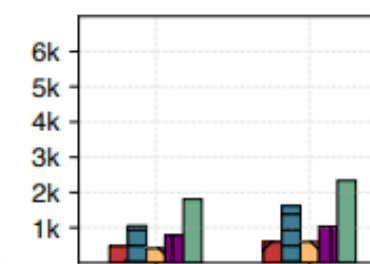
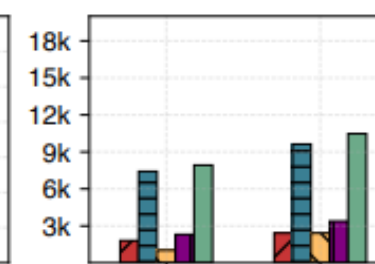
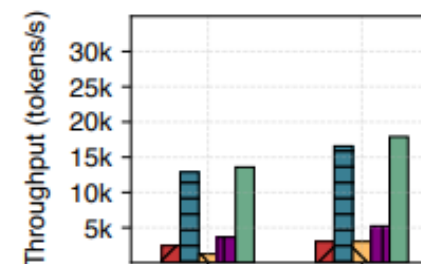
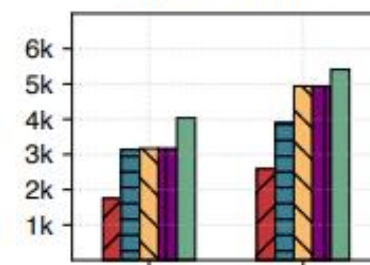
(c) Qwen-4B



(d) Qwen-8B



(e) Qwen-14B



Search Efficiency Comparison

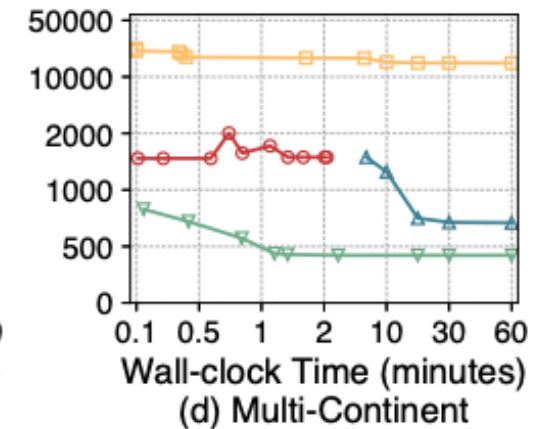
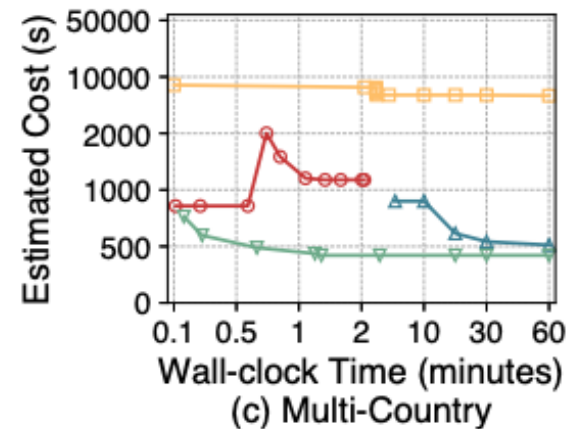
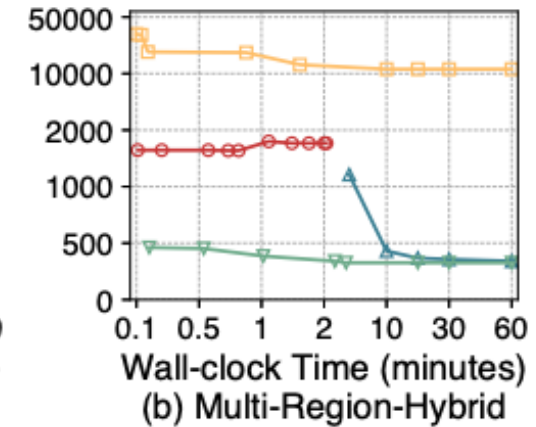
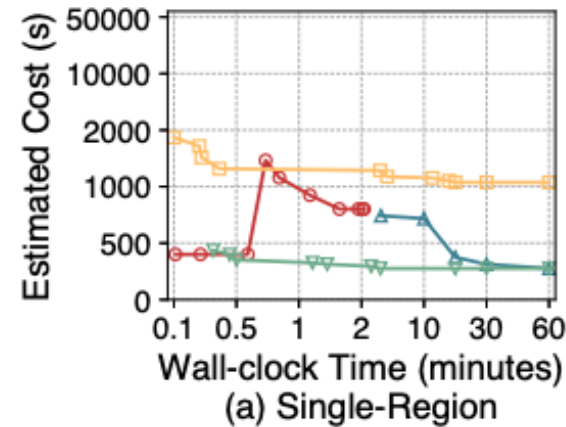
Hybrid Scheduling Algorithm

- Efficiently identifies near-optimal solutions

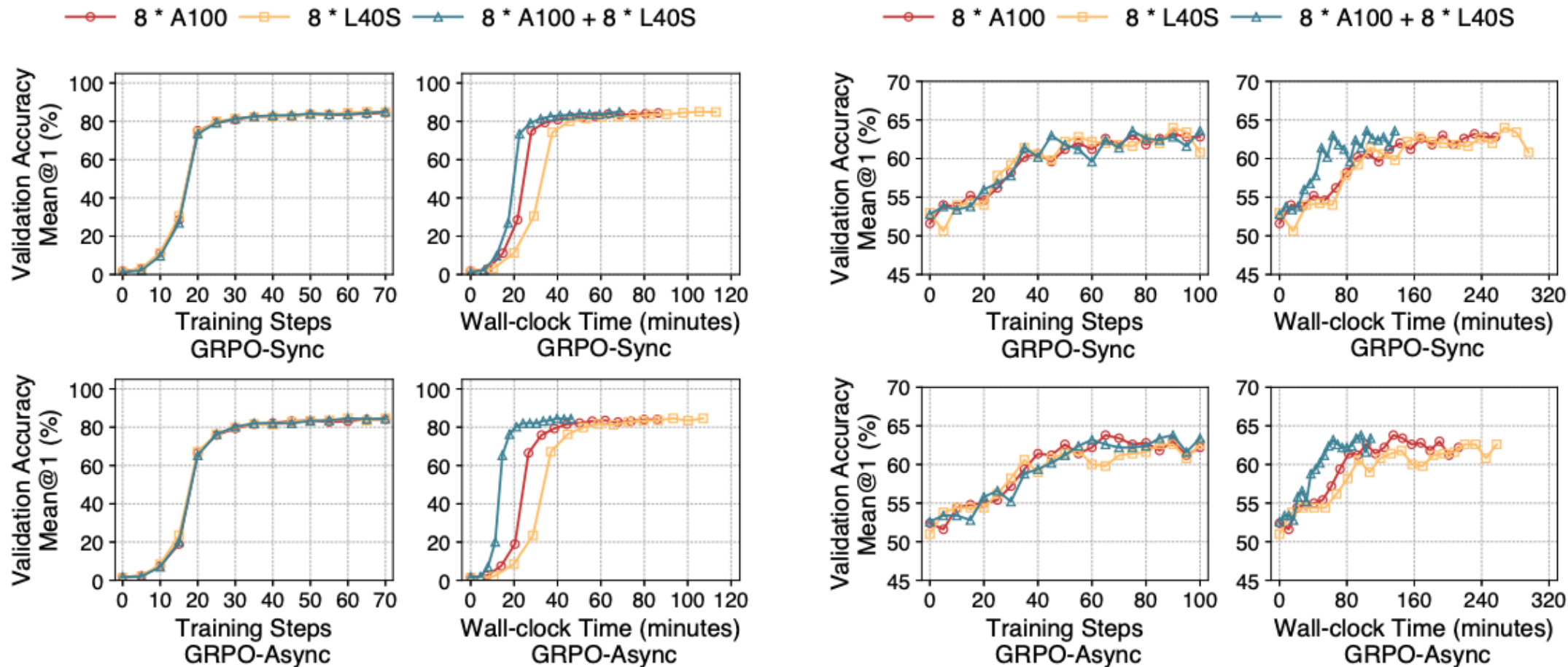
ILP-based scheduling algorithm

- Provides optimal solutions when budgets allows

—○— verl —□— DEAP —△— HetRL (ILP) —▽— HetRL (SHA-EA)



Impact on Training Quality



Training dynamics of Qwen3-1.7B-Base on GSM8K (left) and MATH-500 (right) with GRPO

Summary



HetRL: Efficient Reinforcement Learning for LLMs in heterogeneous environments

- Support for heterogeneity-aware deployment
 - Alleviating the shortage of homogeneous high-end GPUs within a single AZ
 - Leveraging underutilized mid-range or previous-generation GPUs
- Fast and heterogeneity-aware scheduling
 - Hybrid scheduling algorithm efficiently identifies near-optimal solutions
 - ILP-based scheduling algorithm provides optimal solutions when budgets allows
- Outperforms SoTA systems by $3.17\times$ on average