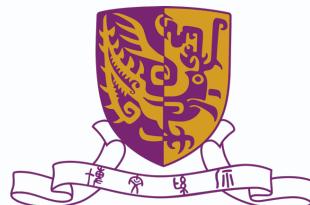


LLMPrism: Black-box Performance Diagnosis for Production LLM Training Platforms

Zhihan Jiang¹, Rui Ren², Guangba Yu¹, Yulun Wu¹, Wenwei Gu¹, Yichen Li¹, Yujie Huang¹
Cong Feng², Yongqiang Yang², Zengyin Yang² and Michael R. Lyu¹

¹The Chinese University of Hong Kong,

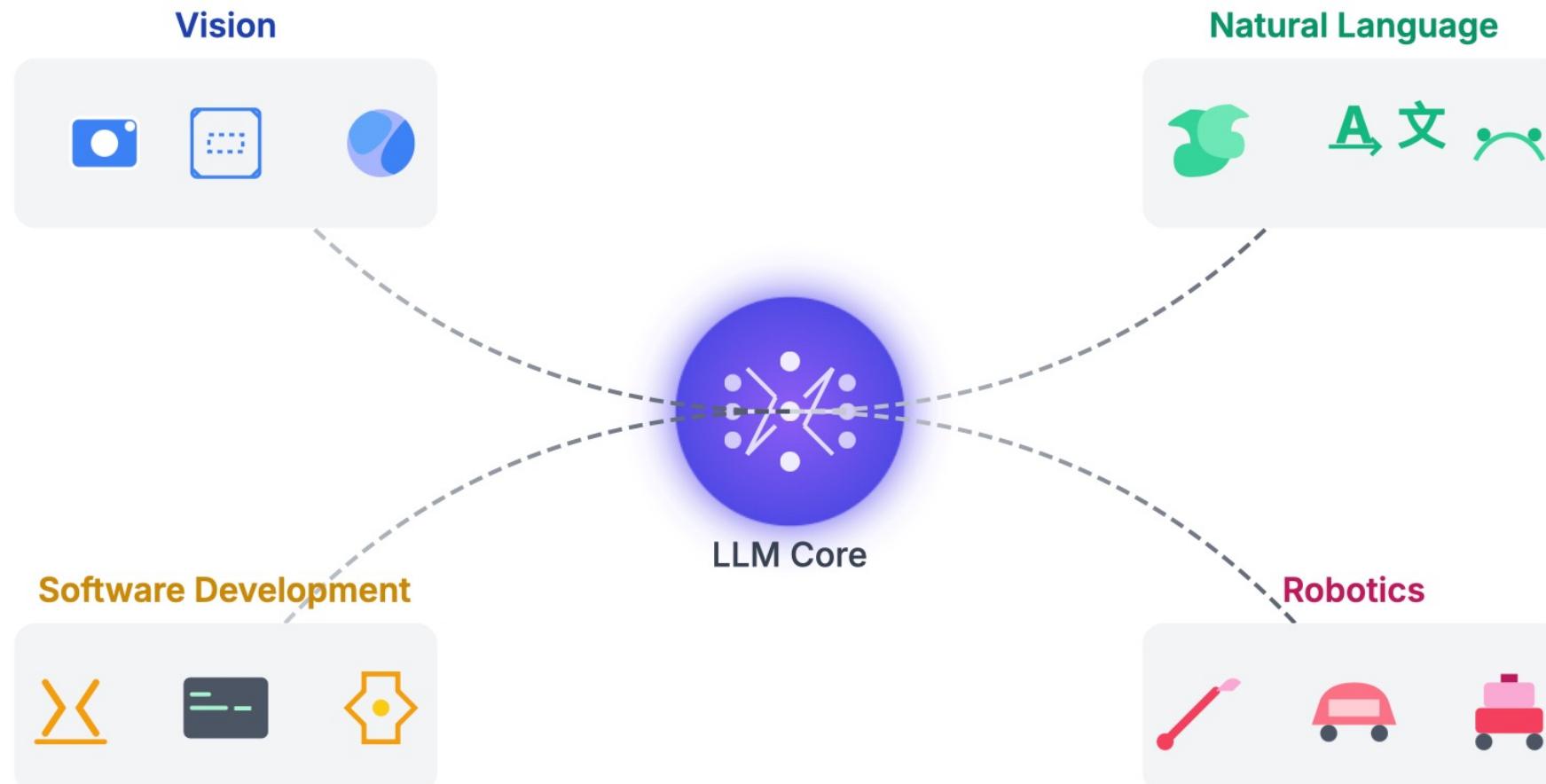
²Huawei Cloud





Background

Large Language Models (LLMs) are everywhere for everyone





Scaling Law in LLMs

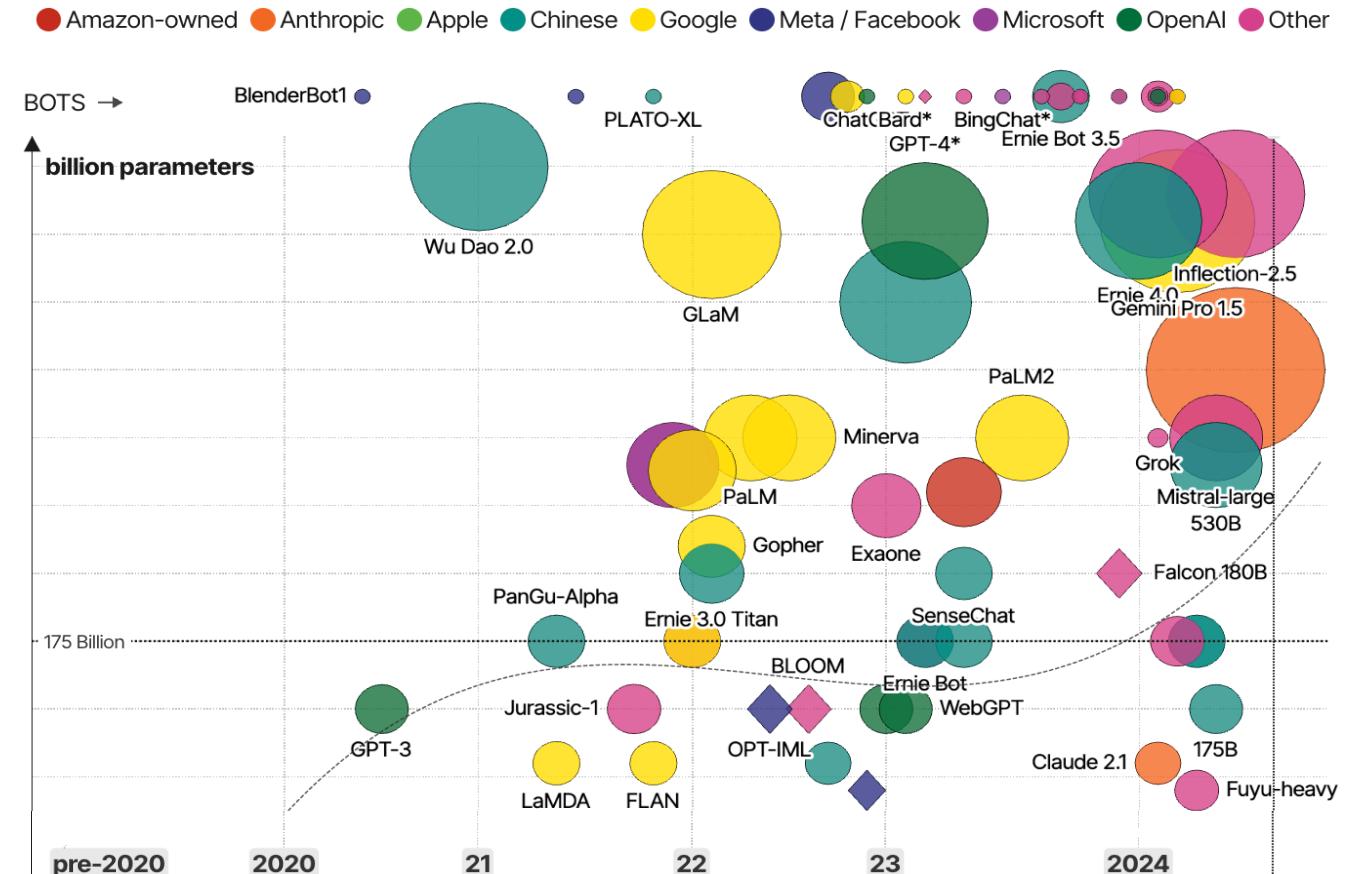
LLMs have grown rapidly in scale, e.g., >100B parameters

 Mistral Large 123B

 OpenAI GPT-3 175B

 Meta LLaMA3.1 405B

 DeepSeek-R1 671B

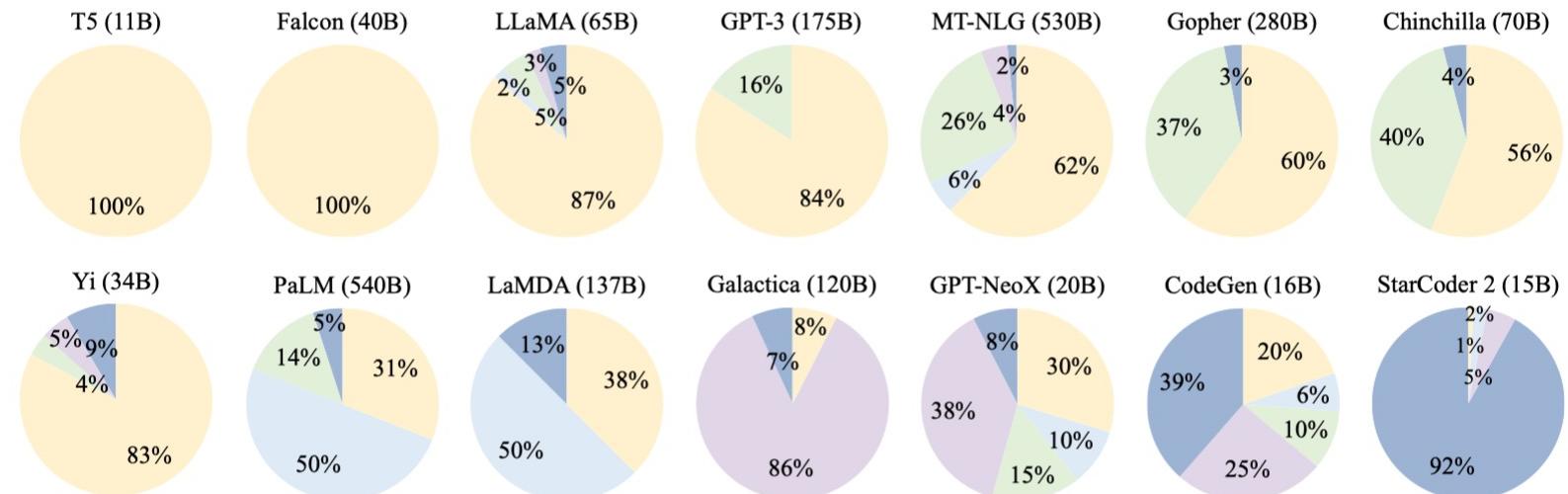




Scaling Law in LLMs

The training datasets are also increasingly diverse and massive

- LAION-2B-en datasets:
over 360B tokens



- RedPajama datasets:
over 30T tokens

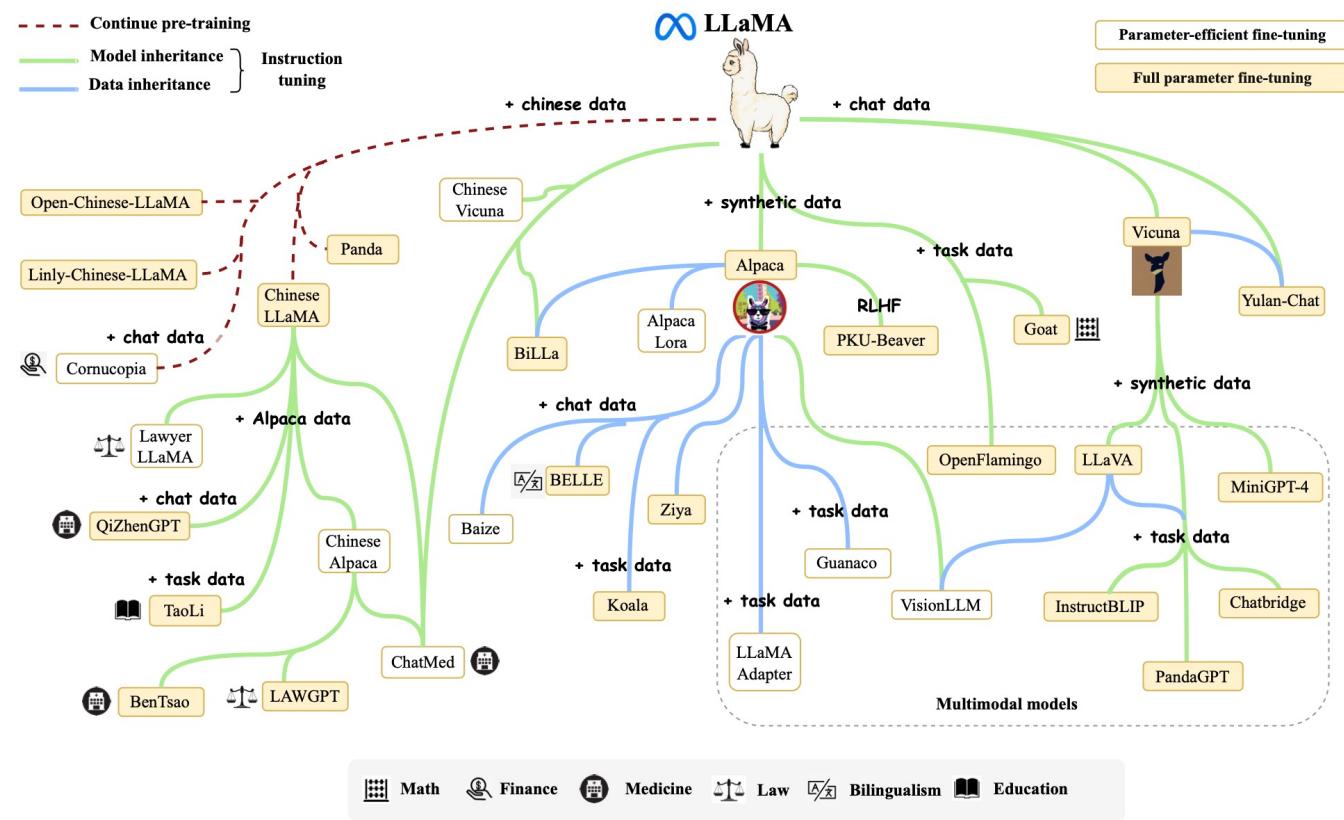
Webpages
Conversation Data
Books & News
Scientific Data
Code

C4 (800G, 2019), OpenWebText (38G, 2023), Wikipedia (21G, 2023)
the Pile - StackExchange (41G, 2020)
BookCorpus (5G, 2015), Gutenberg (-, 2021), CC-Stories-R (31G, 2019), CC-NEWES (78G, 2019), REALNEWS (120G, 2019)
the Pile - ArXiv (72G, 2020), the Pile - PubMed Abstracts (25G, 2020)
BigQuery (-, 2023), the Pile - GitHub (61G, 2020)



LLM Training and Tuning

Fine-tuning LLMs on domain-specific datasets is essential for achieving SOTA performance in specialized fields



A survey of large language models[J]. arXiv preprint arXiv:2303.18223, 2023.



LLM Training and Tuning

Training or tuning LLMs require substantial computing resources

Full-training



BLOOM-176B: 384 A100 GPUs for 3.5 months



LLaMA 3.1 405B: 16,384 H100 GPUs for 54 days

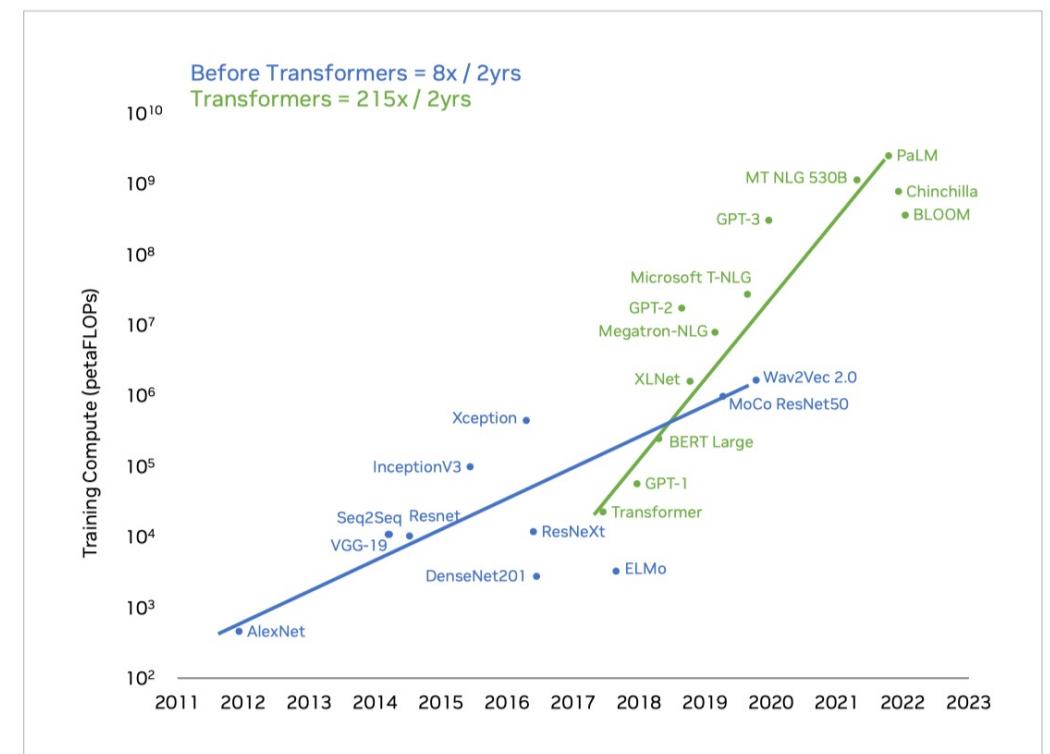


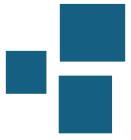
DeepSeek-V3-671B: 2048 H800 GPUs for 56+ days
(2.788M GPU Hours)

Post-training (Fine-tuning)

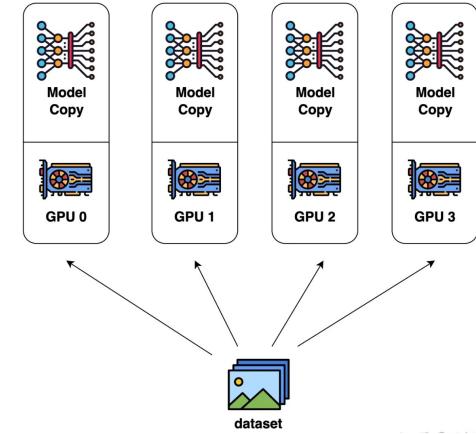
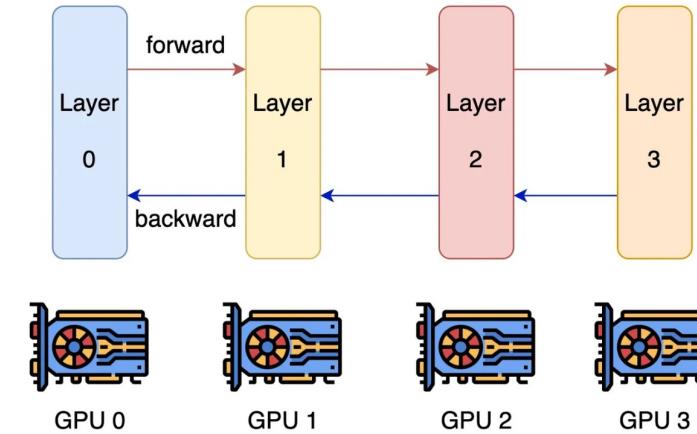
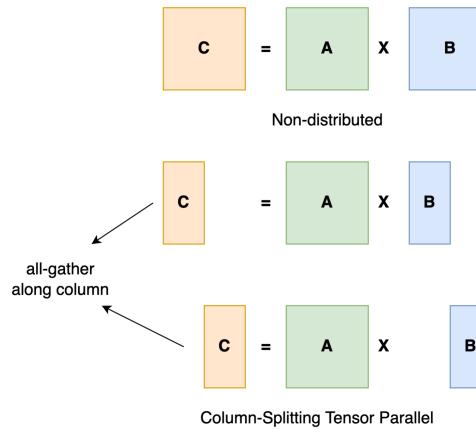


DeepSeek-V3-671B: 64 H800 GPUs for 3+ days
(5K GPU Hours)





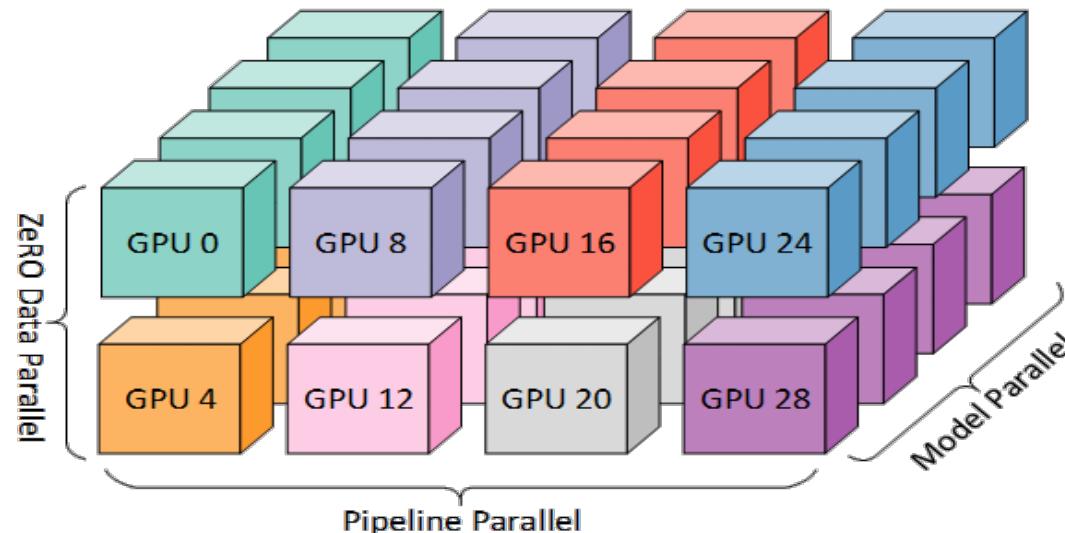
LLM Training Parallelism



Tensor Parallelism

Pipeline Parallelism

Data Parallelism



3D Parallelism



LLM Training Platforms

LLM training platforms: simplify the LLM training and tuning process



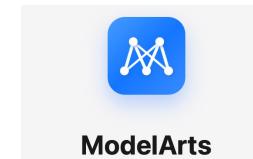
Amazon SageMaker



Google Vertex AI



Databricks



Huawei ModelArts



LLM Training Job Slowdown

Performance degradation is frequent due to:

- The scale and complexity of LLM training jobs
- Strong synchronization requirements at training-step boundaries

Potential Reasons

network congestion

Resource contention

garbage collection

thermal throttling



LLM Training Profiling

LLM training profiling and performance monitoring is important!

Existing LLM training profiler:

PyTorch Profiler

xpu-timer

Meta HTA

Meta Strobelight

Limitation:

Intrusiveness and Compatibility

Performance Overhead

Visibility Constraints



LLM Training Network Monitoring

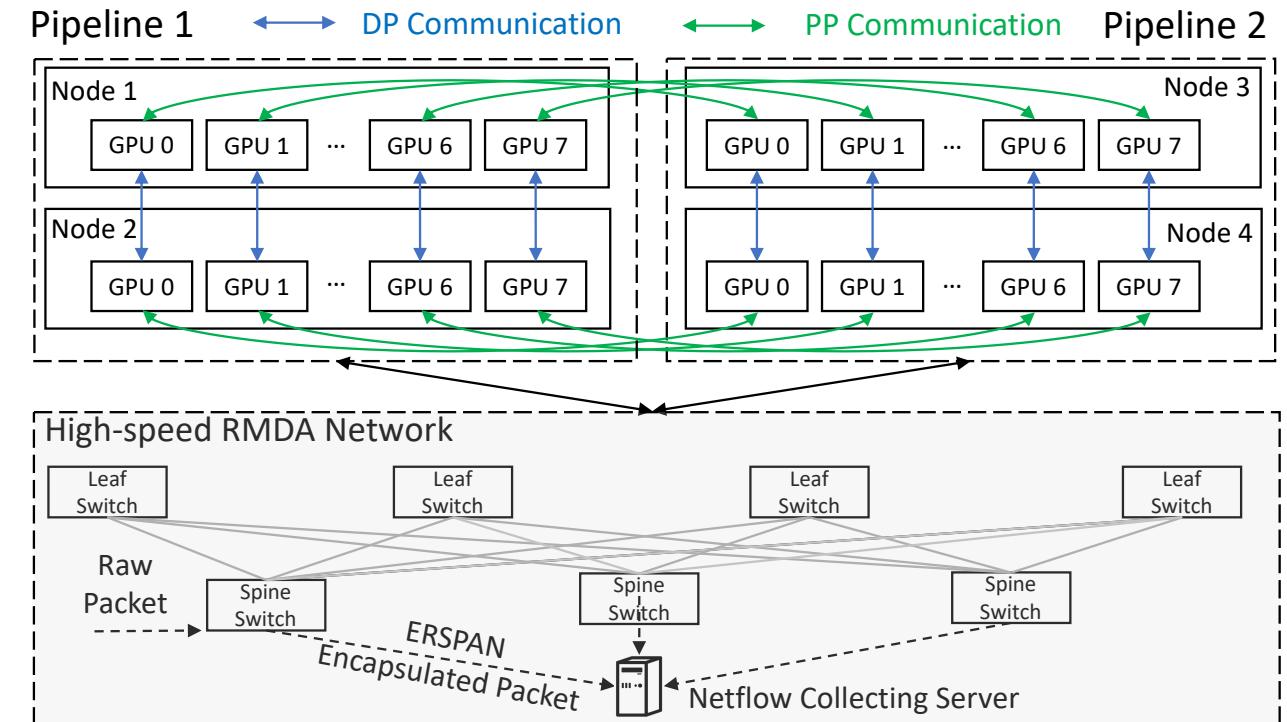
What can we leverage to design a non-intrusive profiling solution with minimal performance impact?

Network traffic monitoring

- non-intrusive and independent approach
- minimal interference with the network and associated workloads
- e.g., Encapsulated Remote Switch Port Analyzer (ERSPAN)

Data format

flow start time, source address, destination address, involved switches, flow size, and flow durations.





LLM Training Network Monitoring

Three observed distinct patterns inherent in LLM training procedure:

Spatial Communication Patterns

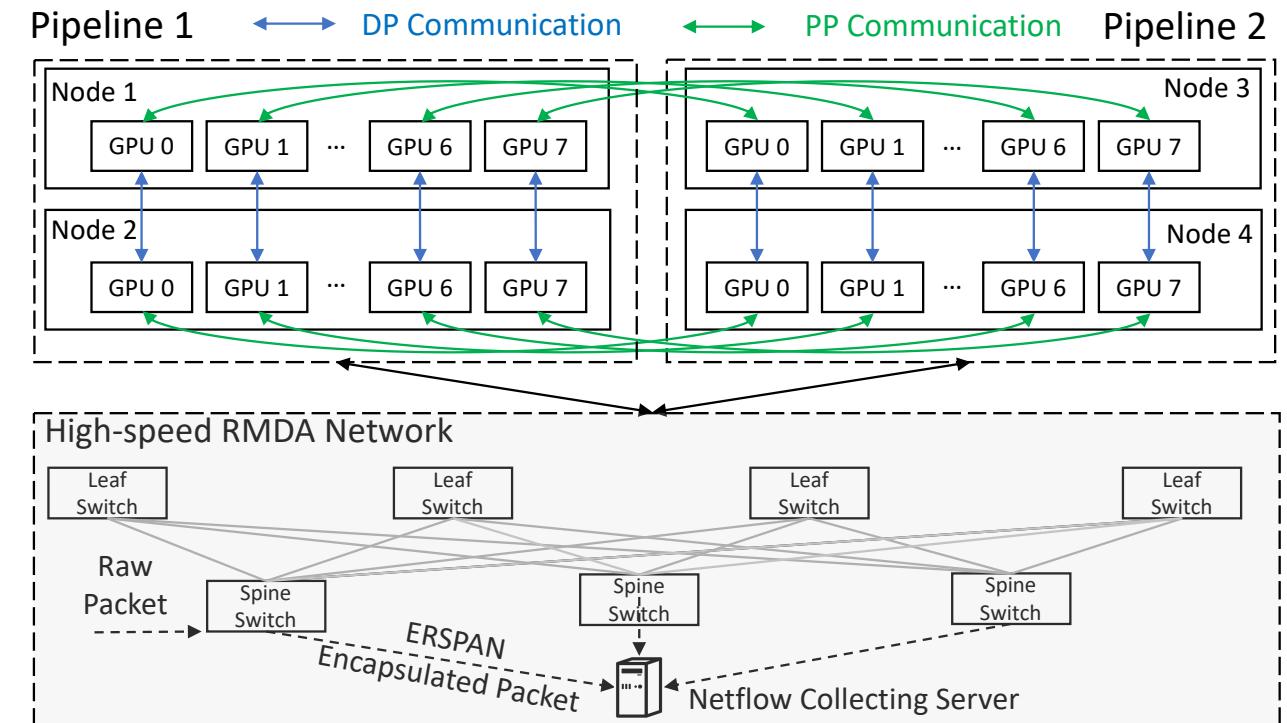
- Communications during LLM training exhibit spatial stability, with interactions strictly confined to GPUs within the same DP or PP groups.

Temporal Communication Patterns

- Communications traffics demonstrate clear temporal periodicity, as the training steps keep cycling and the workload remains stable.

Distinctive Parallelism Communication

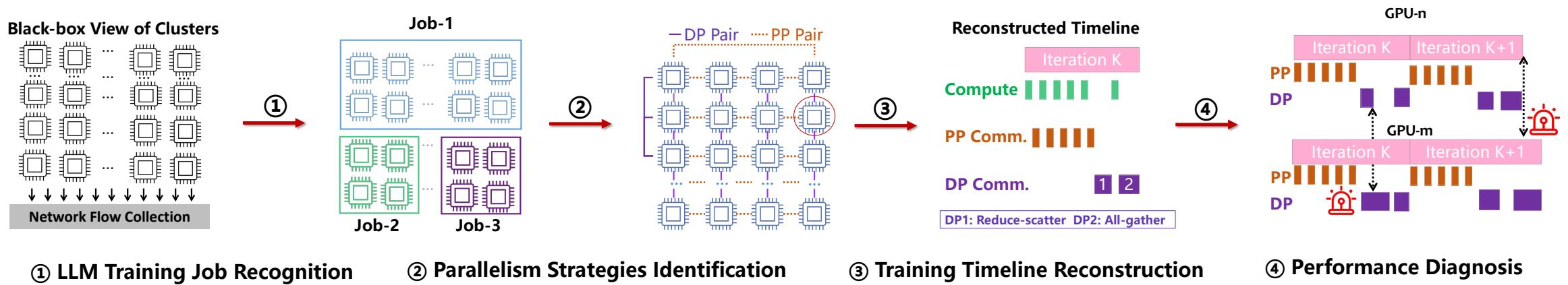
- Different communication modes (DP and PP) display uniquely identifiable characteristics, enabling differentiation between them.





Method: LLMPrism

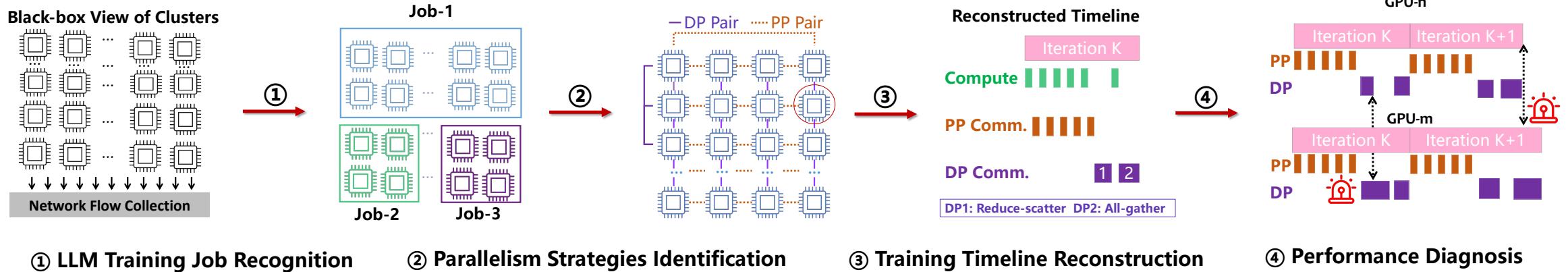
Goal: non-intrusively monitor and diagnose the performance of LLM training jobs



- The first work using underlying network data to reconstruct detailed LLM training timelines
- Transforming the original black-box view into white-box, enabling precise diagnosis of performance issues within large-scale training platforms



Method: LLMPrism



Step-1 LLM Training Job Recognition

Spatial Communication Patterns

Algorithm 1: LLM Training Jobs Recognition

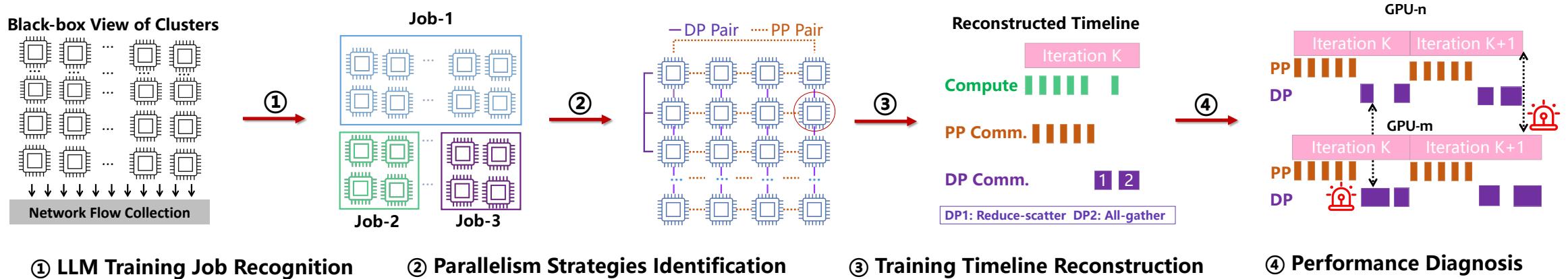
```

Input: RoCE Network Flows:  $\mathcal{F} = \{f_1, f_2, \dots\}$ ; List of
       GPUs:  $\mathcal{X} = \{x_1, x_2, \dots\}$ ; Physical Topology:  $T$ 
Output: Recognized LLM task jobs:  $\mathcal{C} = \{C_1, C_2, \dots\}$ ;
1  $U \leftarrow$  Disjoint-set data structure
2 // (1) Search neighbors and build cross-machine clusters
3 for each network flow  $f_i \in \mathcal{F}$  do
4   | srcAddr, dstAddr  $\leftarrow f_i$ 
5   | if  $U.findSet(srcAddr) \neq U.findSet(dstAddr)$  then
6   |   |  $U.unionSet(srcAddr, dstAddr)$  // merge two GPUs
7  $\mathcal{CR} \leftarrow U.getAllSets()$  // cross-machine clusters
8 // (2) Build task-level clusters based on physical topology  $T$ 
9 for each cross-machine clusters  $c_i \in \mathcal{CR}$  do
10  | get all server list  $s_i$  of  $c_i$  from physical topology
11 for each pair  $c_i$  and  $c_j \in \mathcal{CR}$  do
12  |   | if Jaccard Similarity ( $s_i, s_j$ ) = 1 then
13  |   |   | merge  $c_i$  and  $c_j$ 
14 return  $\mathcal{C} \leftarrow \text{getMergedSet}()$  // job-level clusters

```



Method: LLMP Prism



Step-2 Parallelism Strategies Identification

Distinctive Parallelism Communication

Algorithm 2: Communication Type Identification

```

Input: RoCE Network Flows:  $\mathcal{F} = \{f_1, f_2, \dots\}$ ; Comm. pairs within a job:  $\mathcal{P} = \{(u_1, v_1), \dots, (u_n, v_n)\}$ 
Output: Identified communication types:  $\mathcal{T} = \{t_1, \dots, t_n\}$ ;
1  $\mathcal{G} \leftarrow$  unordered graph of all GPUs within the job
2 for each pair  $(u, v)$  in  $\mathcal{P}$  do
3   // (1) Calculate communication interval
4   extract related network flows:  $\mathcal{F}_{(u,v)} = \{f_1, \dots, f_t\}$ 
5   compute comm. interval:  $\Delta t_i = f_{i+1}(\text{time}) - f_i(\text{time})$ 
6   // (2) Step division
7   employ Bayesian CPD:  $\mathcal{F}_{(u,v)} = \{F_{step_1}, F_{step_2}, \dots\}$ 
8   // (3) Type distinction
9   for each  $step_k$  do
10    | count the distinct number of flow sizes  $N_k$  in  $F_{step_k}$ 
11    type  $t(u, v) \leftarrow$  PP if  $\text{Mode}(N_k) = 1$  else DP
12    add edge  $(u, v)$  to  $\mathcal{G}$  if  $t(u, v)$  is DP
13 // (4) Noise refinement
14 DFS to find all connected component  $CC_i$  in  $\mathcal{G}$ 
15 for all pairs  $(u, v)$  in all  $CC_i$  do
16  | refine  $t(u, v)$  to DP if it is PP
17 return  $\mathcal{T} = \{t_1, \dots, t_n\}$  // final communication types

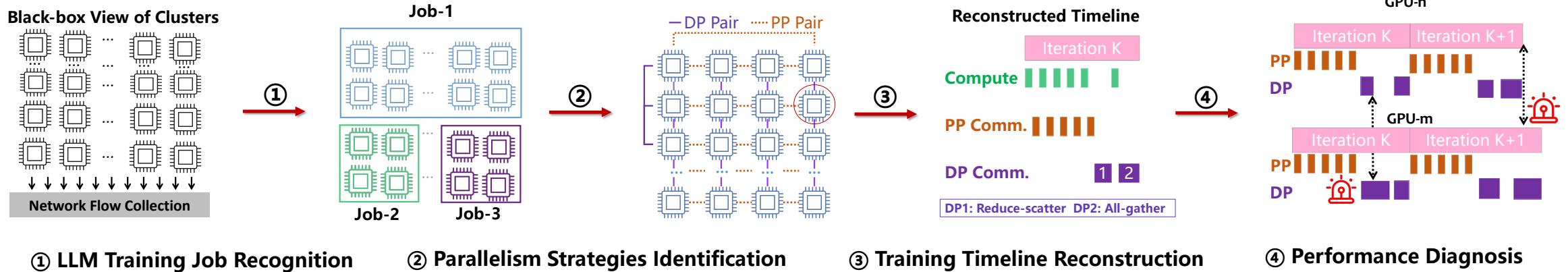
```

Bayesian Online Change-point Detection (BOCD) algorithm

$$r_t = \begin{cases} 0, & \text{if a change-point occurs at } t \\ r_{t-1} + 1, & \text{otherwise} \end{cases}$$



Method: LLMPrism



Step-3 Training Timeline Reconstruction Temporal Communication Patterns

Key Findings:

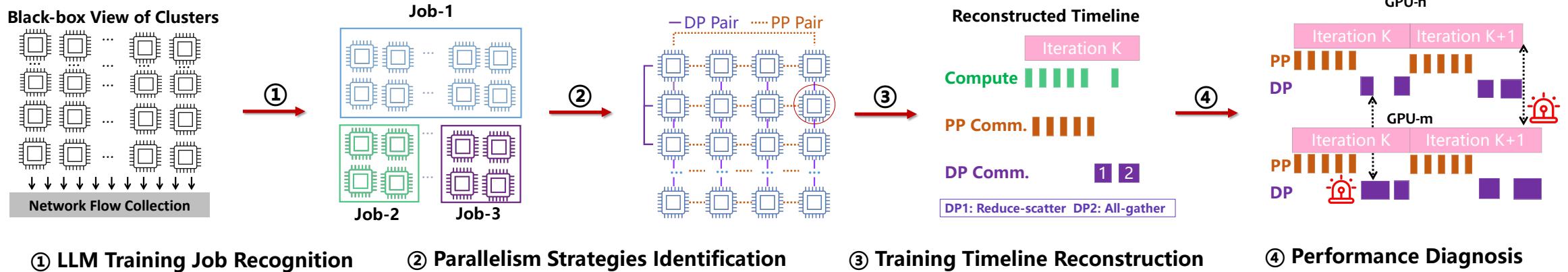
- each training step concludes with a segment of DP collective communication traffic

Bayesian Online Change-point Detection (BOCD) algorithm

$$r_t = \begin{cases} 0, & \text{if a change-point occurs at } t \\ r_{t-1} + 1, & \text{otherwise} \end{cases}$$



Method: LLMPrism



Step-4
Performance Diagnosis

Cross-step Diagnosis

Cross-group Diagnosis

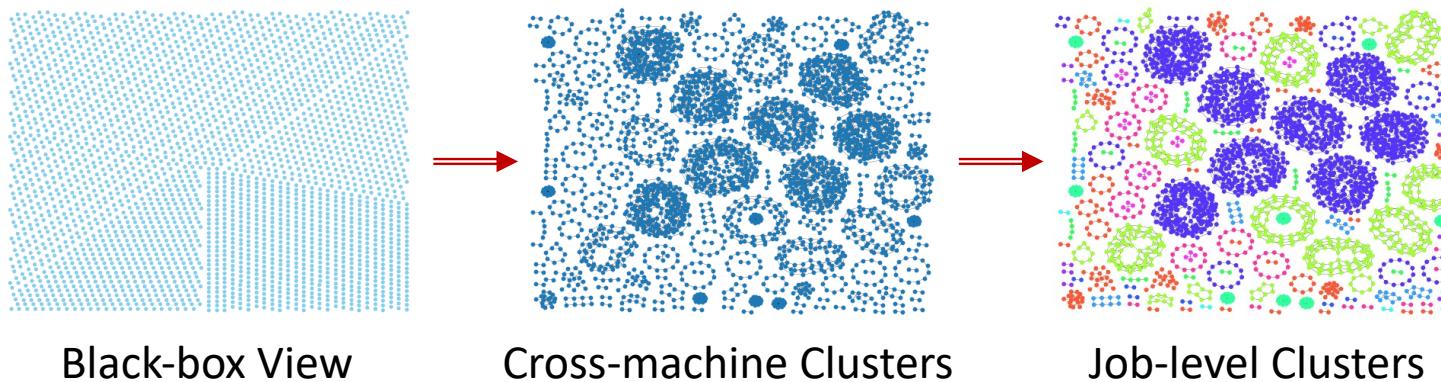
Switch-level Diagnosis



Evaluation

Part-1: Effectiveness of LLM Training Jobs Recognition

successfully identifying 19 training jobs



Internal Training Clusters with 2,880 XPU

Evaluation

Part-2: Effectiveness of Parallelism Strategies Identification

TABLE I: Accuracy in identifying parallelism strategies via flows of varying-length periods for jobs with 1024 GPUs.

Methods	1 min Acc.	3 min Acc.	5 min Acc.	10min Acc.
LLMPrism w/o refinement	96.00%	97.93%	98.03%	99.61%
LLMPrism	100%	100%	100%	100%

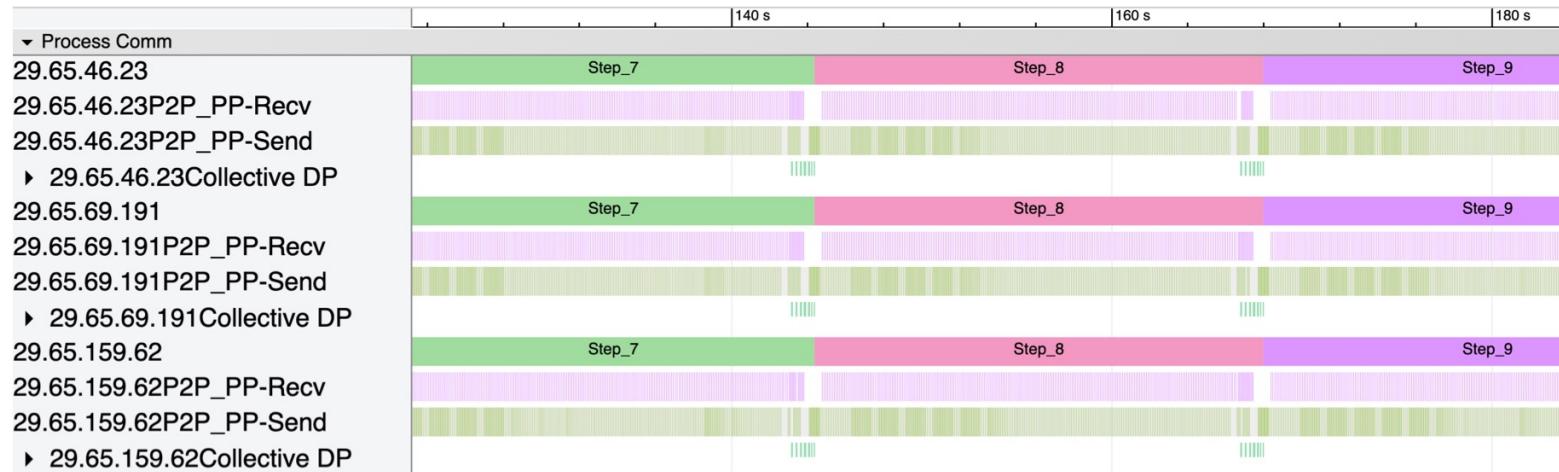
5 LLM training jobs with 1024 XPUs with manual ground truths

- Without the noise refinement process, errors may occur in parallelism strategy classification, where some DP pairs are misclassified as PP pairs.
- This issue becomes more pronounced when network flow data is limited to shorter durations.
- Incorporating the noise refinement process allows LLMPrism to consistently achieve 100% accuracy across all training jobs



Evaluation

Part-3: Effectiveness of Training Timeline Reconstruction



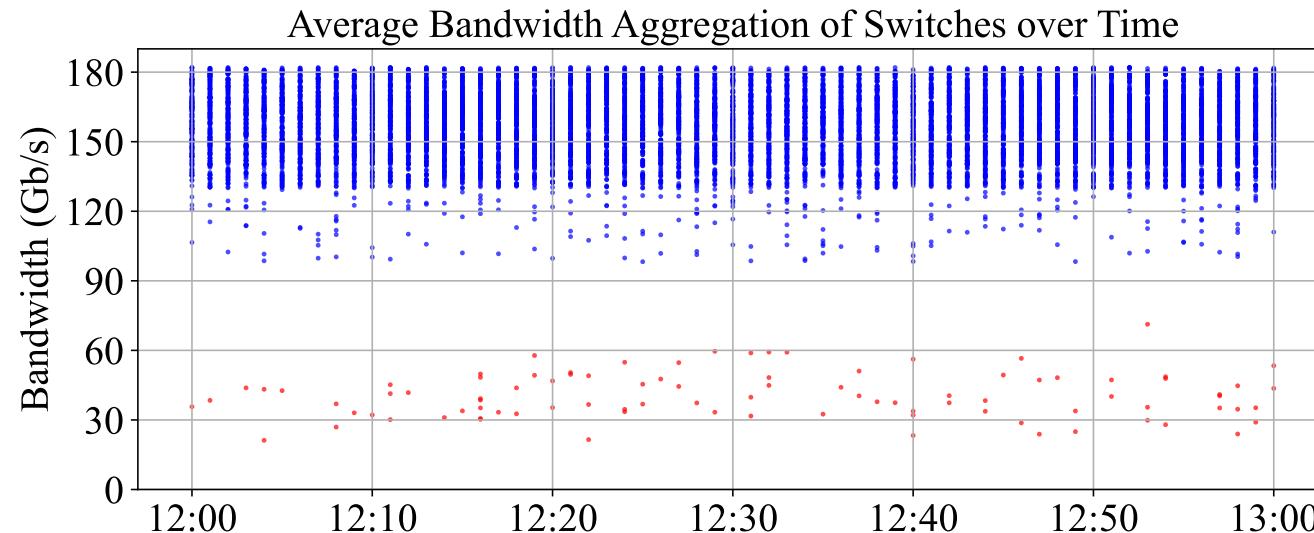
The reconstructed training timeline of a training job using 1,024 XPU

reconstruction error of LLMPrism remained within 0.3%



Evaluation

Part-4: Effectiveness of Performance Diagnosis



Switch-level diagnosis of LLMPrism using one-hour aggregated average bandwidths

successfully identified network issues and triggered alerts

Conclusion

LLM Training Profiling

LLM training profiling and performance monitoring is important!

Existing LLM training profiler:

PyTorch Profiler xpu-timer

Meta HTA Meta Strobelight

Limitation:

Intrusiveness and Compatibility

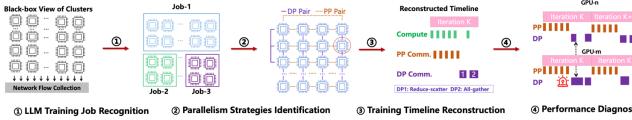
Performance Overhead

Visibility Constraints

Motivation

Method: LLMP Prism

Goal: non-intrusively monitor and diagnose the performance of LLM training jobs



- The first work using underlying network data to reconstruct detailed LLM training timelines
- Transforming the original black-box view into white-box, enabling precise diagnosis of performance issues within large-scale training platforms.

Evaluation

Part-1: Effectiveness of LLM Training Jobs Recognition



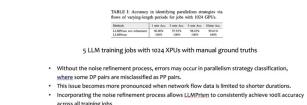
Part-3: Effectiveness of Training Timeline Reconstruction



The reconstructed training timeline of a training job using 1,024 XPU

reconstruction error of LLMPrism remained within 6.3%

Part-2: Effectiveness of Parallelism Strategies Identification



Without the noise refinement process, errors occur in parallelism classification, where some DP pairs are misclassified as PP pairs.

The issue becomes more pronounced when network flow data is limited to shorter durations.

However, the refinement process allows LLMPrism to consistently achieve 100% accuracy across all training jobs.

TABLE 1: Issues in existing parallelism analysis: the flow of a typical single point-to-point job with 1024 GPUs.

5 LLM training jobs with 1024 XPU are used for manual ground truths

Switch-level diagnosis of LLMPrism using one-hour aggregated average bandwidths

successfully identified network issues and triggered alerts

Approach

Experiment

Thank you for listening

Q & A

Contact: zhjiang@link.cuhk.edu.hk