



# L4: Diagnosing Large-scale LLM Training Failures via Automated Log Analysis

Zhihan Jiang<sup>1</sup>, Junjie Huang<sup>1</sup>, Guangba Yu<sup>1</sup>, Zhuangbin Chen<sup>2</sup>, Yichen Li<sup>1</sup>, Renyi Zhong<sup>1</sup>  
Cong Feng<sup>3</sup>, Yongqiang Yang<sup>3</sup>, Zengyin Yang<sup>3</sup> and Michael R. Lyu<sup>1</sup>

<sup>1</sup>The Chinese University of Hong Kong,

<sup>2</sup>Sun Yat-sen University,

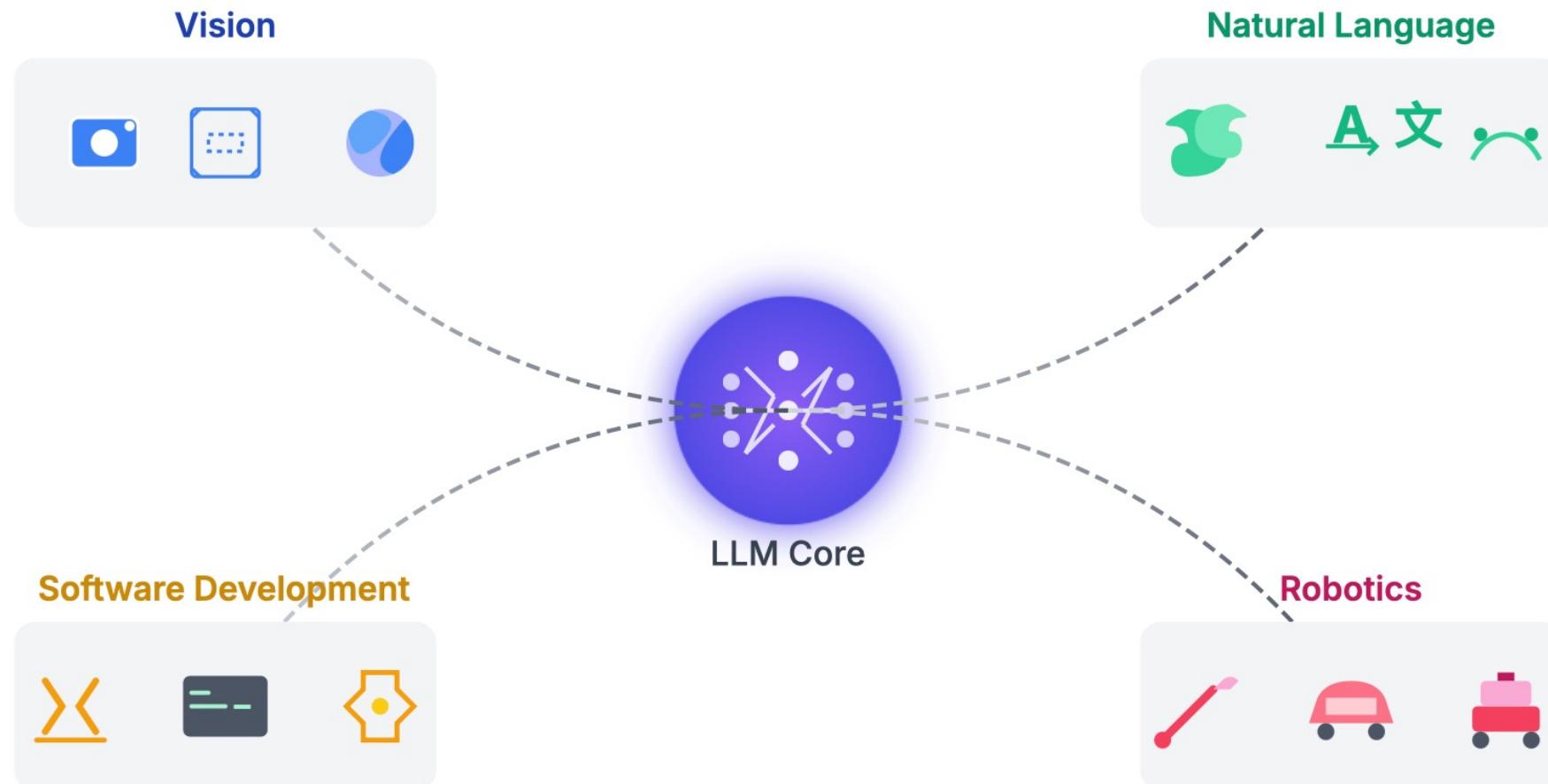
<sup>3</sup>Huawei Cloud





# Background

**Large Language Models (LLMs) are everywhere for everyone**





# Scaling Law in LLMs

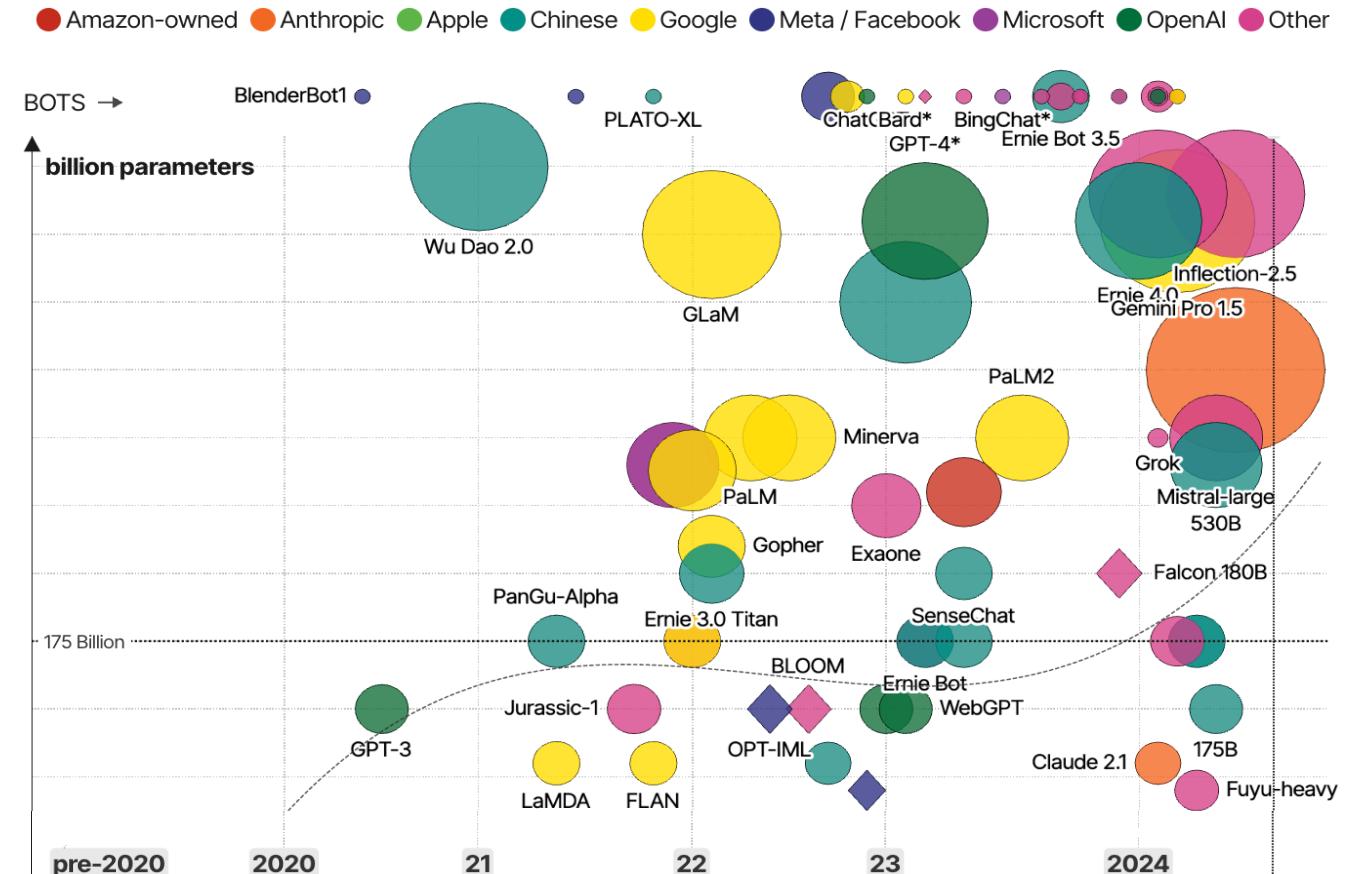
LLMs have grown rapidly in scale, e.g., >100B parameters

 Mistral Large 123B

 OpenAI GPT-3 175B

 Meta LLaMA3.1 405B

 DeepSeek-R1 671B



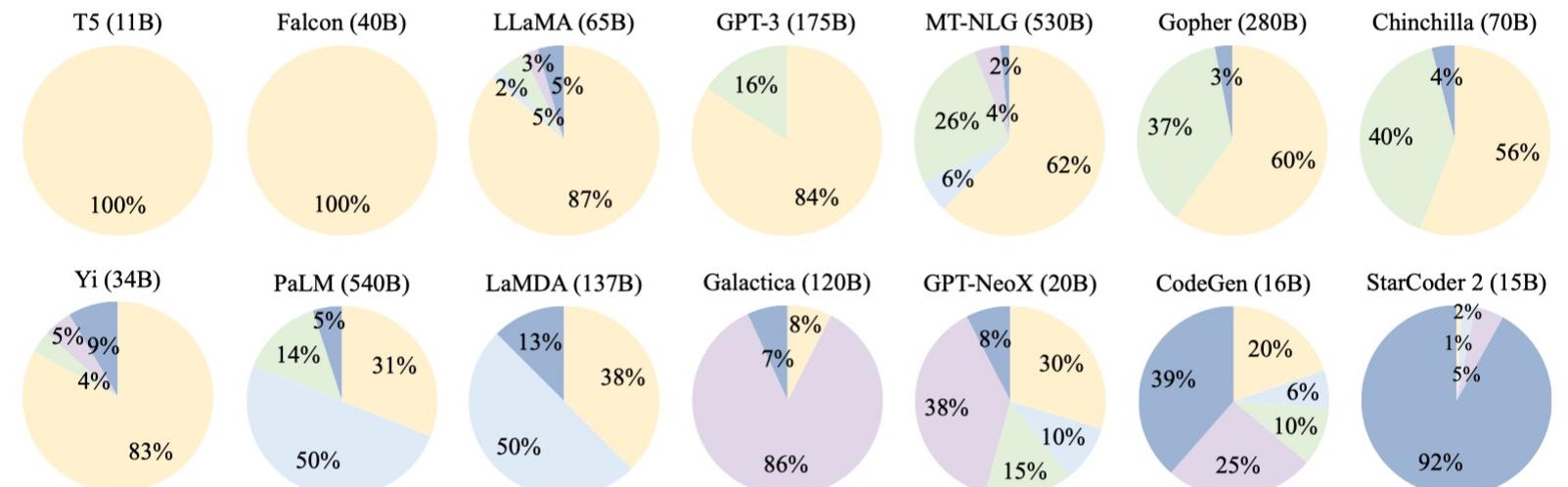


# Scaling Law in LLMs

The training datasets are also increasingly diverse and massive

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

- RedPajama datasets:  
over 30T tokens



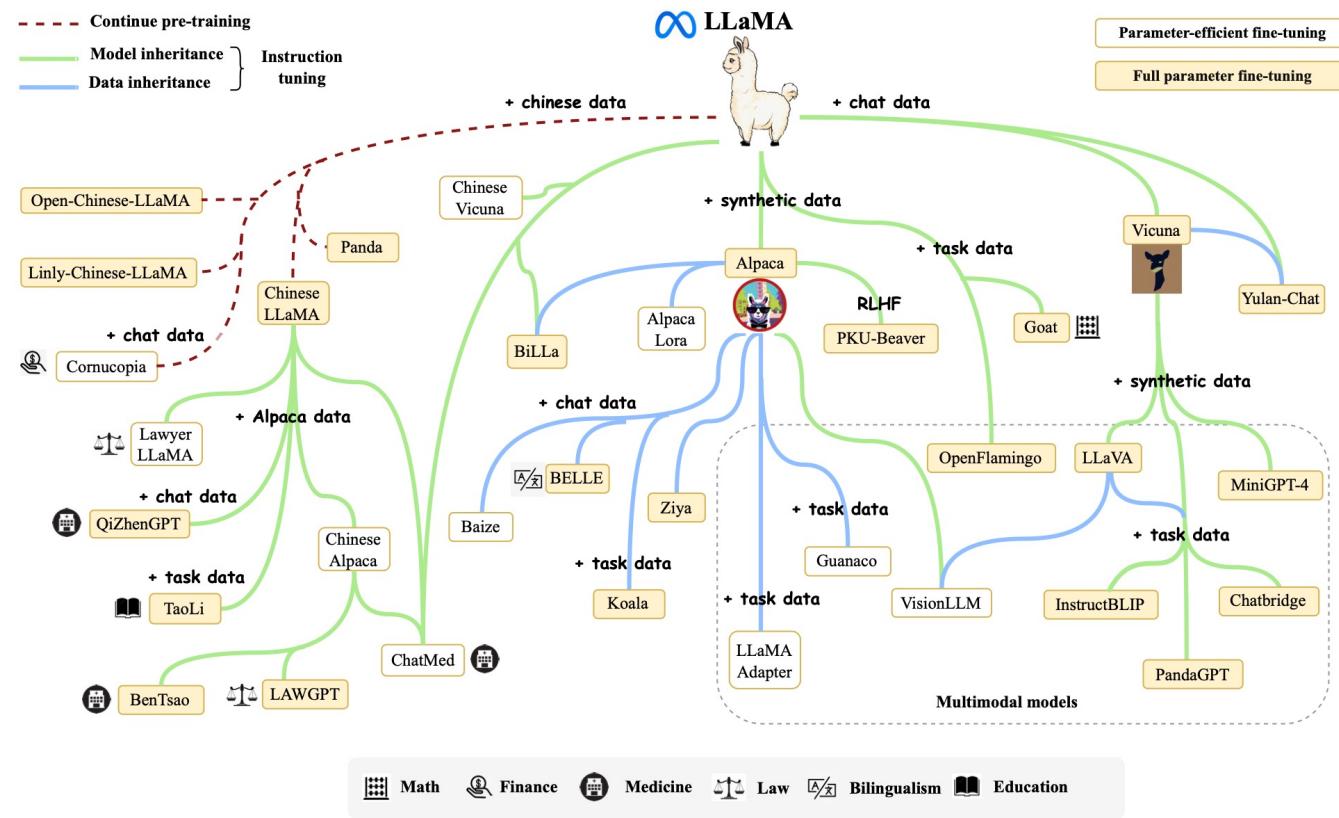
Webpages  
Conversation Data  
Books & News  
Scientific Data  
Code

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





# 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

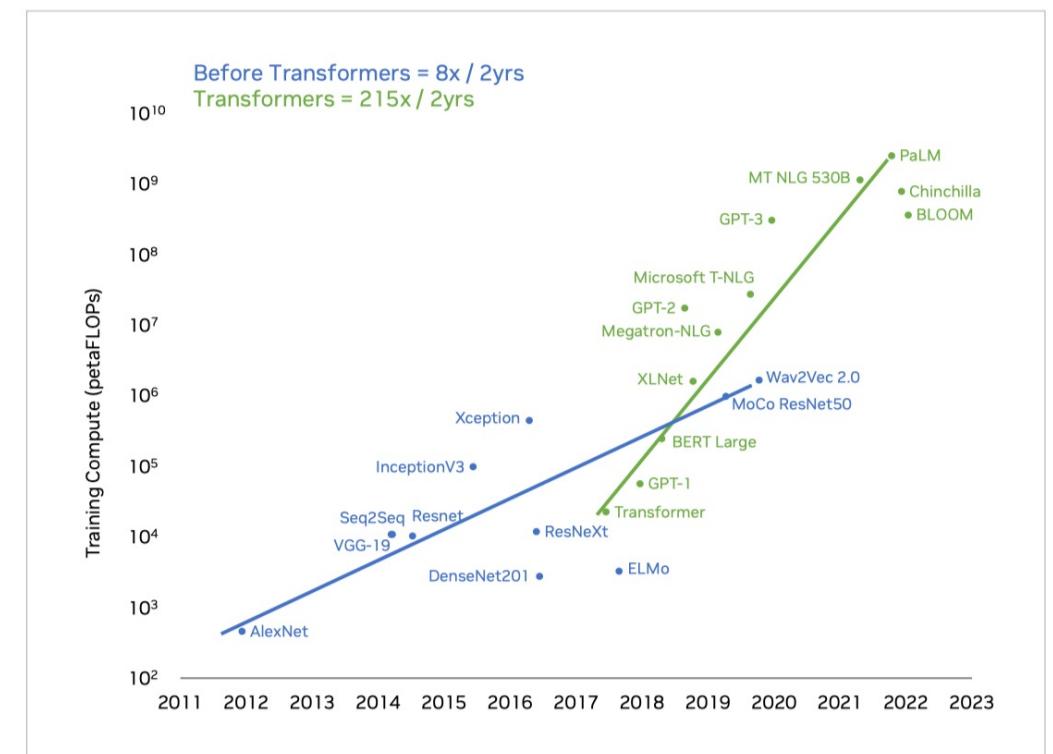


**DeepSeek-V3-671B:** 2048 H800 GPUs for 56+ days  
(2.788M GPU Hours)

## Post-training (Fine-tuning)

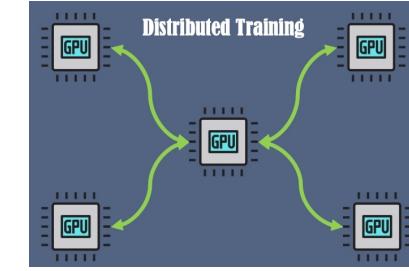


**DeepSeek-V3-671B:** 64 H800 GPUs for 3+ days  
(5K GPU Hours)





# LLM Training Reliability

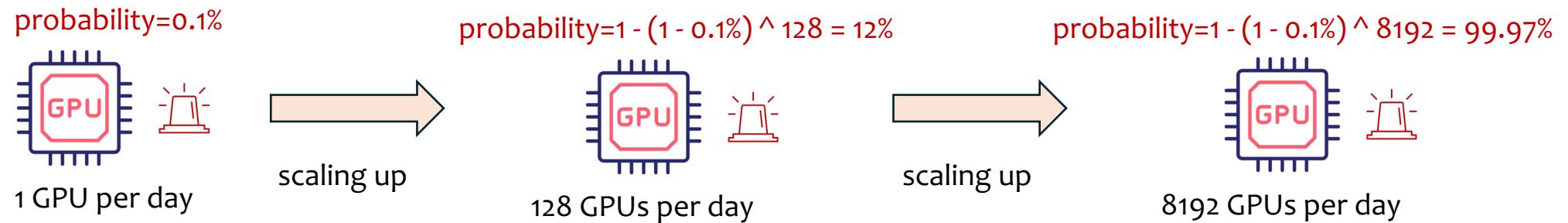


As the scales of LLMs and the training resources grow, the likelihood of failures during the training process increases significantly

The core reason: synchronization properties during the distributed training process



A local fault in a specific component (e.g., a GPU), can disrupt the entire training process





# LLM Training Reliability

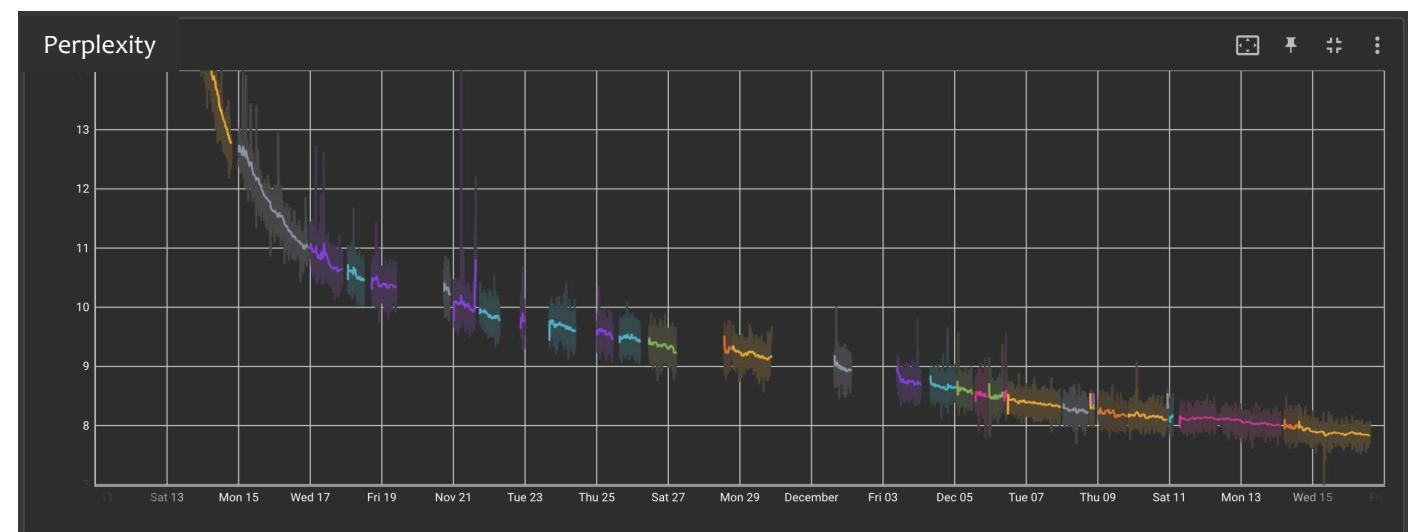
**As the scales of LLMs and the training resources grow, the likelihood of failures during the training process increases significantly.**

Example: Development life-cycle of OPT-175B from Oct. 2021 to Jan. 2022

35+ manual restarts

70+ automatic restarts

100+ cycling hosts



# LLM Training Platforms

LLM training platforms: simplify the LLM training and tuning process



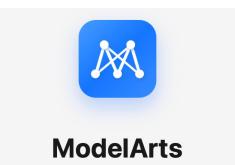
Amazon SageMaker



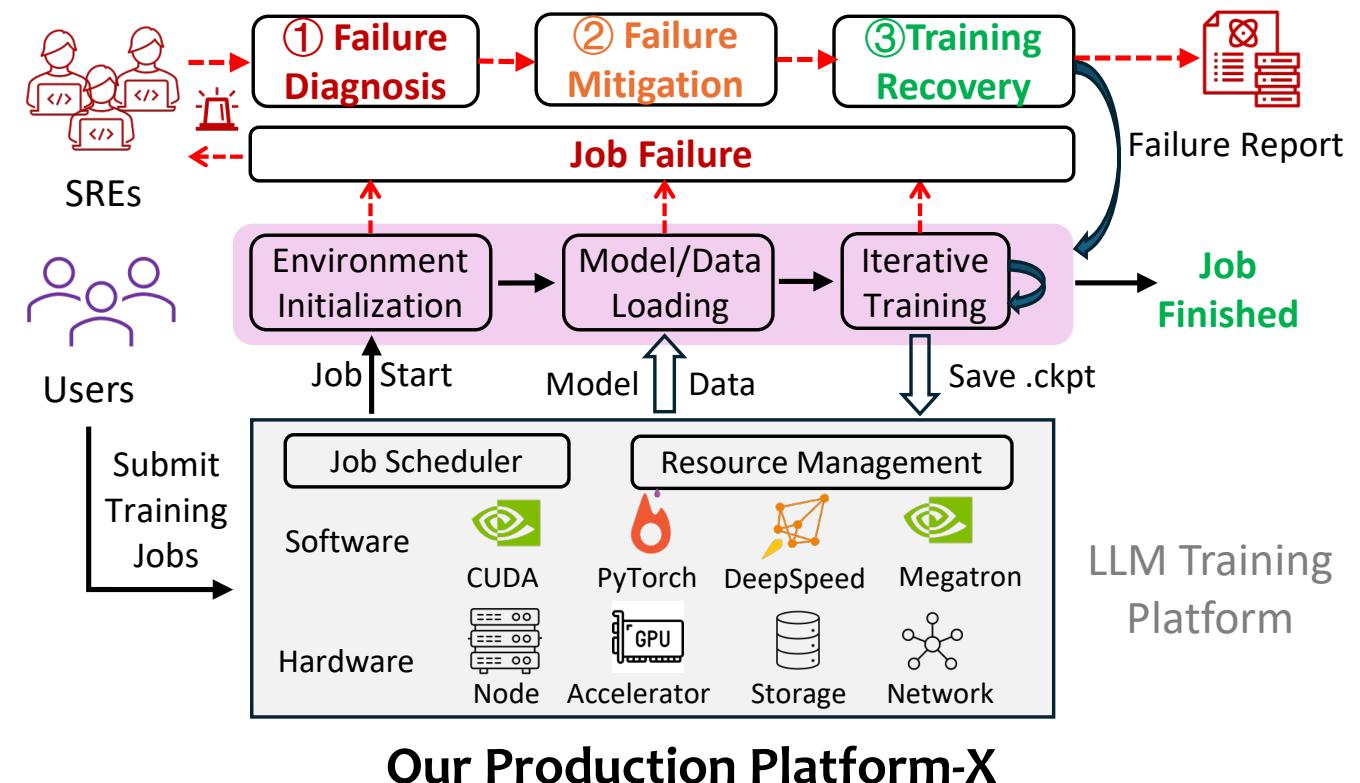
Google Vertex AI



Databricks



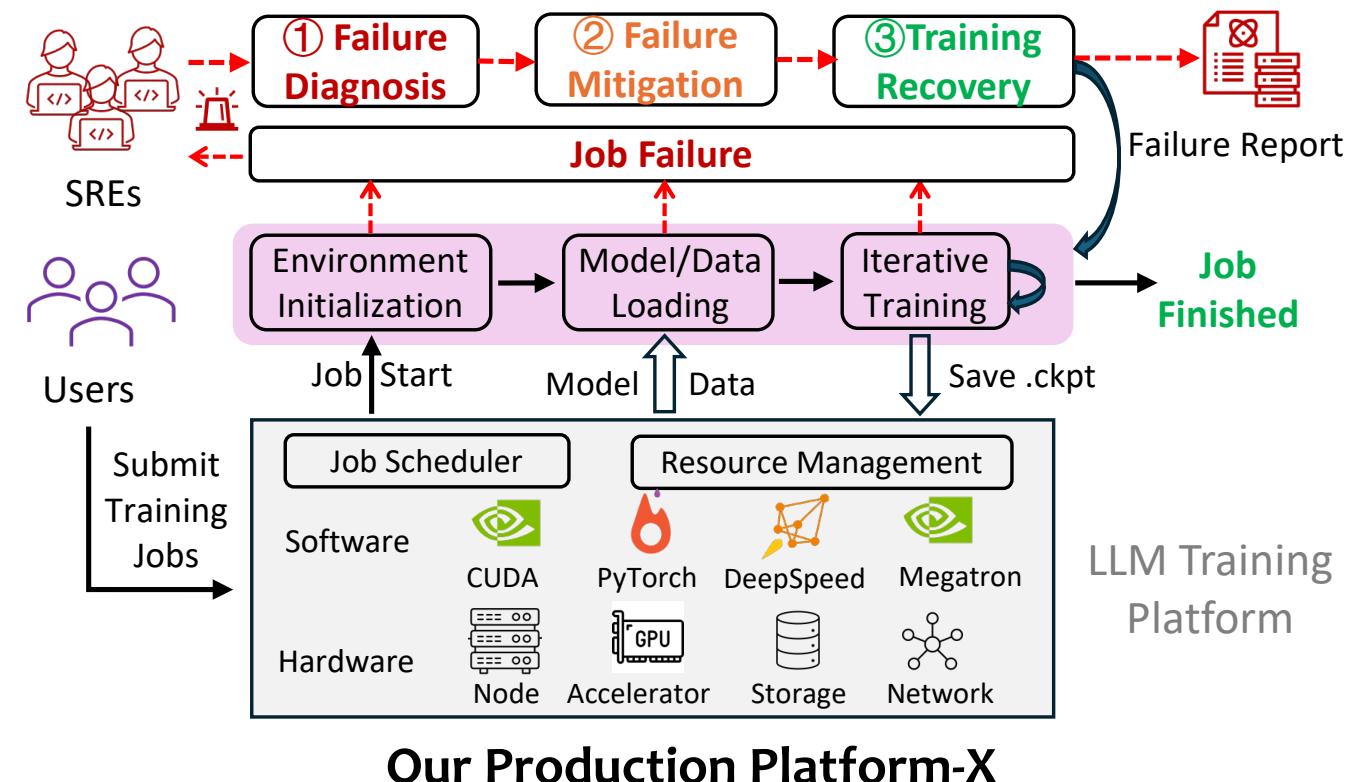
Huawei ModelArts





# LLM Training Failure Diagnosis

Rapid and accurate failure diagnosis is critical for failure management





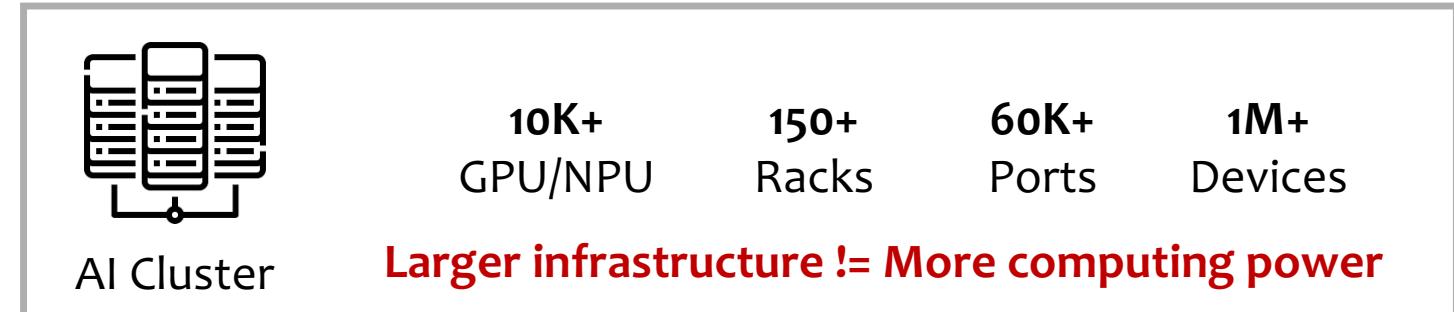
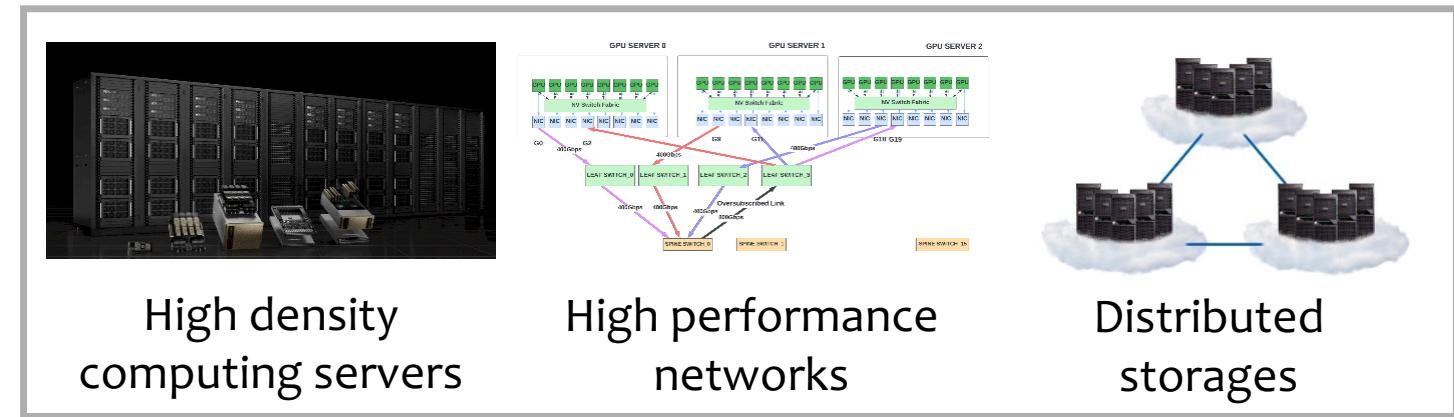
# LLM Training Failure Diagnosis

Rapid and accurate failure diagnosis is critical for failure management

## Challenges:

### Horizontal:

- large-scale clusters and a multitude of components





# LLM Training Failure Diagnosis

Rapid and accurate failure diagnosis is a critical step in failure recovery

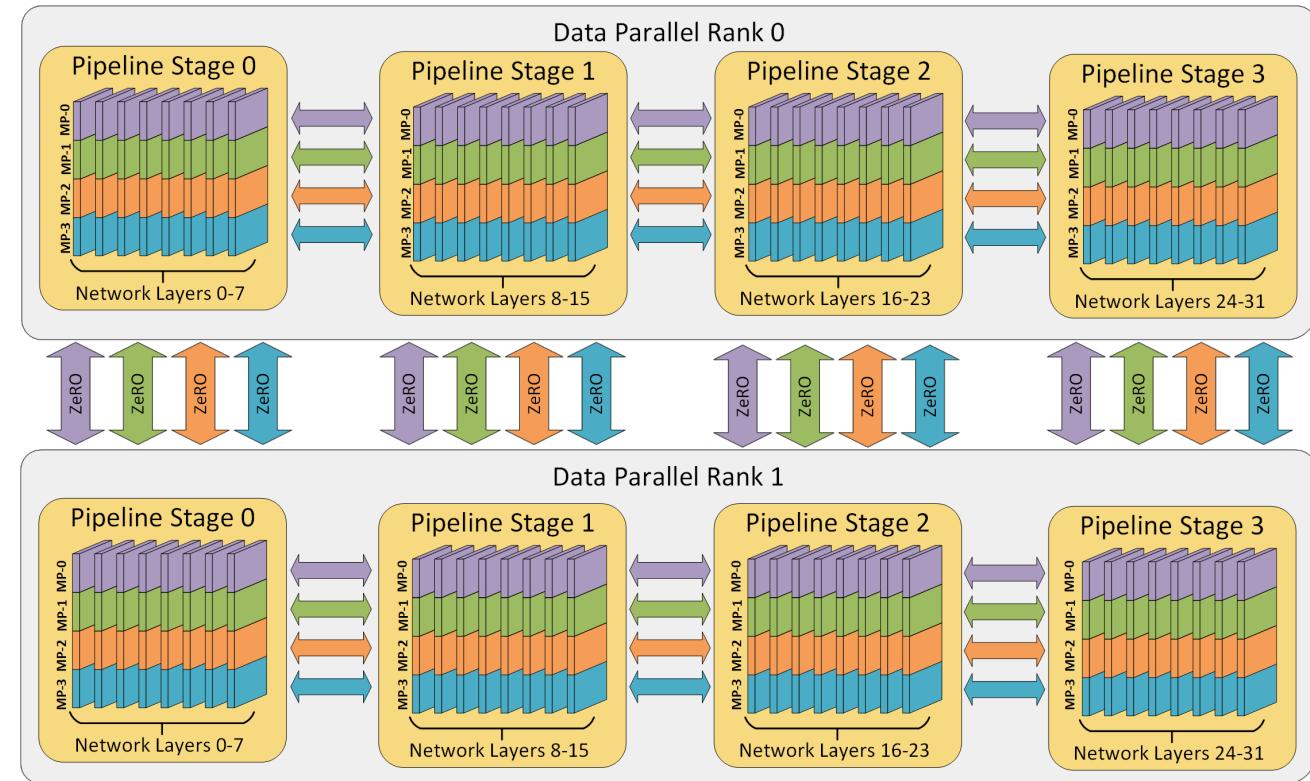
## Challenges:

### Horizontal:

- large-scale clusters and a multitude of components

### Vertical:

- complex configurations and multi-layered architectures





# LLM Training Failure Diagnosis

Rapid and accurate failure diagnosis is a critical step in failure recovery

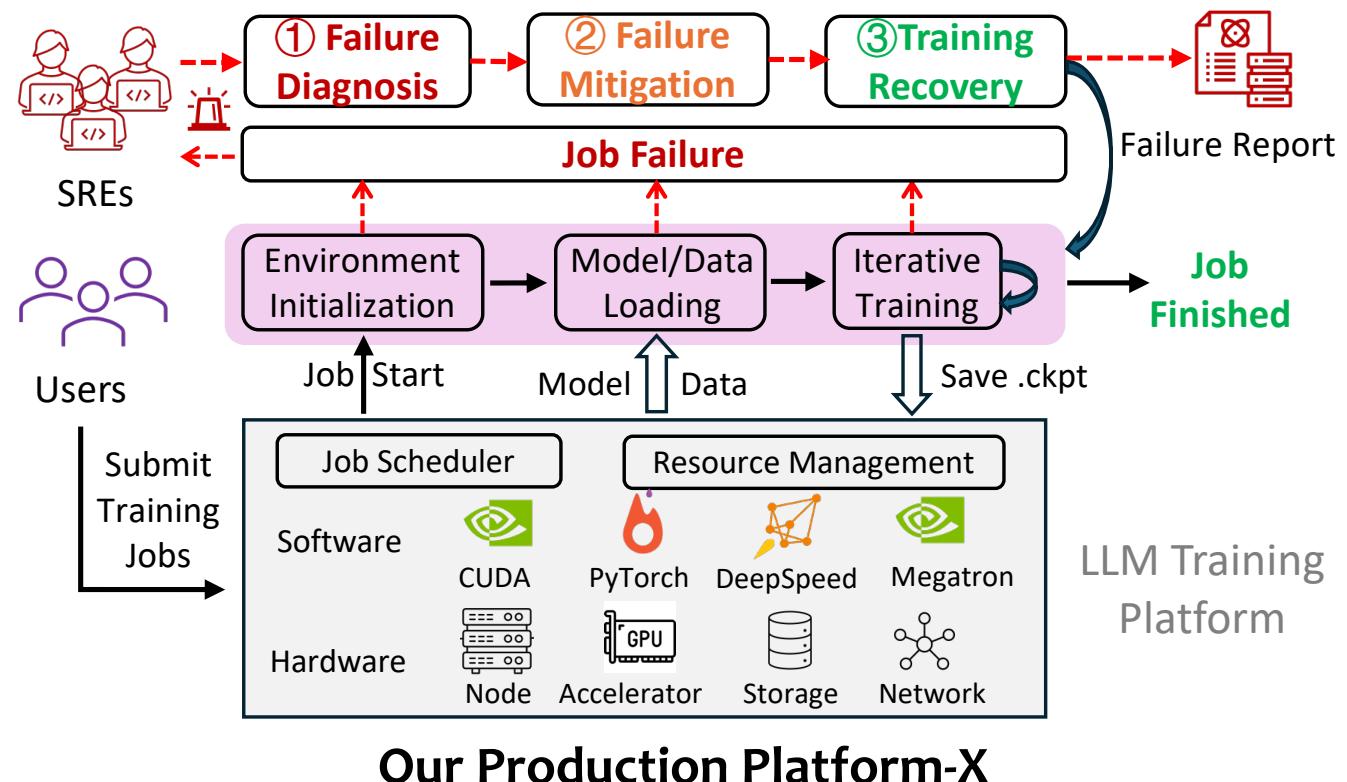
## Challenges:

### Horizontal:

- large-scale clusters and a multitude of components

### Vertical:

- complex configurations and multi-layered architectures



Our Production Platform-X

Understanding failures in LLM training and exploring opportunities for automated diagnosis are important!



# Empirical Study on LLM Training Failures

## Study Platform

Our multi-tenant production LLM training Platform-X supporting hundreds of internal users and partner companies

## Study Subject

Collected 428 failure reports of failed large-scale LLM training jobs in Platform-X from May 2023 to April 2024

- average model size: **72.8B parameters**
- average accelerators: **941 GPUs or NPUs**

|   |   |  |
|---|---|--|
| Job ID: 1437825    Report submission time: 2024/05/22 13:24:42<br>Training start time: 2024/05/21 16:27:14<br>Training end time: 2024/05/22 13:14:52  |   | <b>Job metadata</b><br><br><b>Report status:</b> Resolved  |
| <b>Hardware resource</b><br><br>Resource region: DC region-1<br>Storage position:<br>S3://job_1437825/<br>Comp. resource: 256 x Node-type2  | <b>Software environment</b><br><br>OS image: Ubuntu 20.04<br>Driver: NVIDIA 525.60.13<br>Framework: PyTorch 2.1.0<br>Library: Transformers 4.33.0 ....  | <b>Monitoring data</b><br><br>Training log: <br><br>Other data:  |
| <b>Failure description</b><br><br>Model info: Llama-2-70b-chat-hf<br>Symptom: Job failed after 1120 steps running.<br>Comment: I've attempted to relaunch the job; however, the launch still wasn't successful. | <b>Diagnostic Information</b><br><br>SRE leader: Jackie<br>Diagnosis root cause: Node-17 GPU-3 detached<br>Diagnosis procedure: there are error logs like "rank-131 connection lost", stress test failed..... |  |
|   |   |  |

## Study RQs

- RQ1: What are the common symptoms of LLM training failures?
- RQ2: What are the common root causes of LLM training failures?
- RQ3: What monitor data sources are typically used to diagnose LLM training failures?



# RQ1: LLM training failures symptoms

## Launching Failure:

- environment initialization
- model and dataset downloading

21.3%

## Training Crash:

- iterative training failed

57.5%

## Abnormal Behavior:

- training hanging
- training slowdown

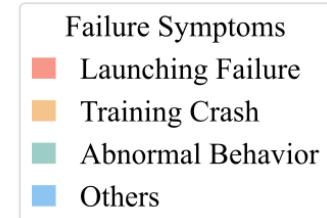
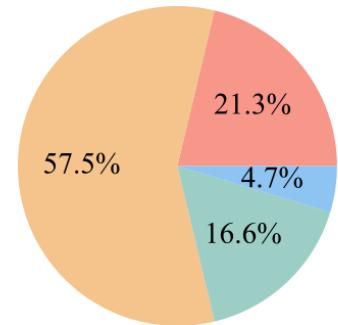
16.6%

## Others:

- resource unavailability

4.7%

74.1% during  
iterative training



Classification of LLM training failure symptoms



## Main Finding

Most LLM training failures (74.1%) occur during the iterative training stage, which can waste significant computational resources and training time.



# RQ2: LLM training failure root causes

## 1<sup>st</sup> : Hardware Fault:

- network fault
- accelerator fault
- node fault
- storage fault

## 2<sup>nd</sup>: User Fault

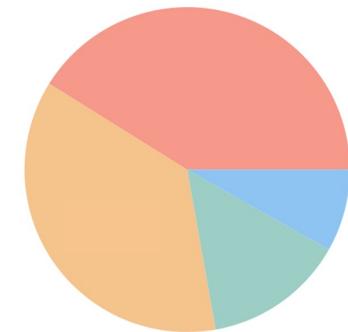
- configuration error
- program/script bug
- software incompatibility
- mis-operation

## 3<sup>rd</sup> : Framework Fault

- various LLM training frameworks and libraries:
  - PyTorch, DeepSpeed, MindSpore, etc.

## 4<sup>th</sup> : Platform Fault

- platform resource management
- platform job scheduling and configurations



| Failure Root Causes |
|---------------------|
| Hardware Fault      |
| User Fault          |
| Framework Fault     |
| Platform Fault      |

Classification of LLM training failure root causes



## Main Finding

**LLM training failures stem from diverse causes, with hardware faults being significantly more prevalent than in traditional computing or DL workloads.**



# RQ3: LLM training failure diagnosis data

- **Manual diagnosis is challenging**

average diagnosis time: 34.7 hours

diagnosis > 24 hour: 41.9%

average training log size: 16.92GB

- **Monitor data for failure diagnosis**

## Log-only Diagnosable:

- only training logs are involved

53.9%  
36.0%  
89.9% failures

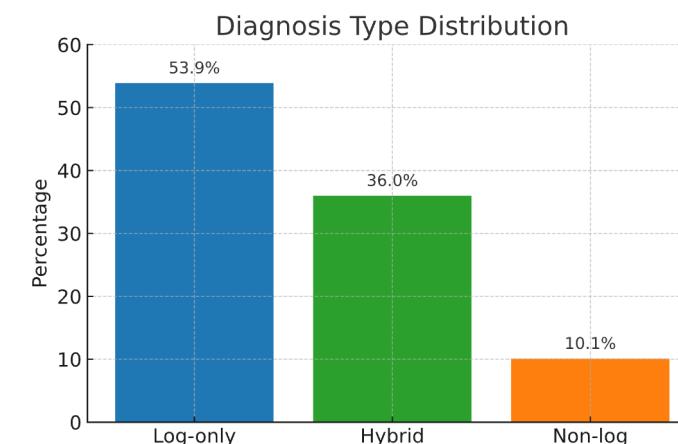
## Hybrid Diagnosable:

- Both training logs and other monitoring data are jointly used

## Non-log Diagnosable:

- Training logs do not provide clues

10.1%



## Main Finding

Training logs are crucial for diagnosing most (89.9%) LLM failures, but their sheer volume highlights the need for **automated failure-indicating log identification**.

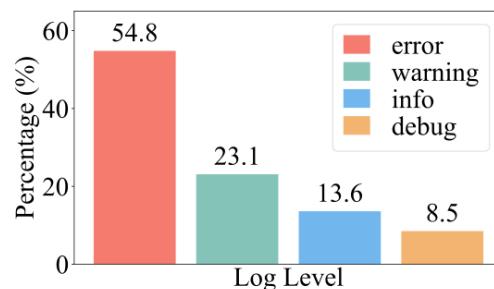
# Limitation for Existing Approaches

**Goal:** Automated identification of *failure-indicating logs* from substantial LLM training logs

**Features used in existing log-based anomaly detector:**

## Logging Level:

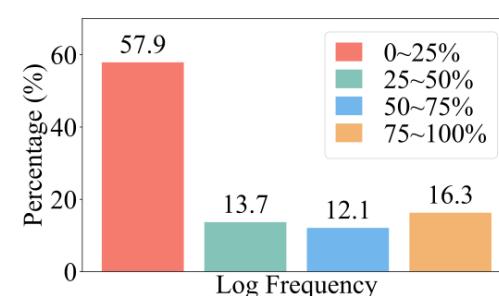
- more serious logs (error)



(a) Different logging levels.

## Event Frequency:

- infrequent log events



(b) Different event frequencies.

## Error Semantic:

- logs with error semantics

## Our observations:

- Error-semantic logs from specific components/stages may not impact training.
- Successful training jobs often contain various error-semantic logs.
- Not all failure-indicating logs show explicit error semantics.



## Main Finding

Existing Log AD methods struggle to effectively identify failure-indicating logs, as traditional anomalous log features are not suitable for LLM training log scenarios.



# Distinct Patterns of LLM Training Logs

---

We have observed **three distinct patterns** of LLM training logs that can be used to automatically pinpoint *failure-indicating log*

## ① Cross-job Pattern

- In practice, each failed training job is usually associated with a series of successful jobs with identical settings (e.g., models and frameworks)
- The log events from both successful and failed jobs are typically noisy

## ② Spatial Pattern

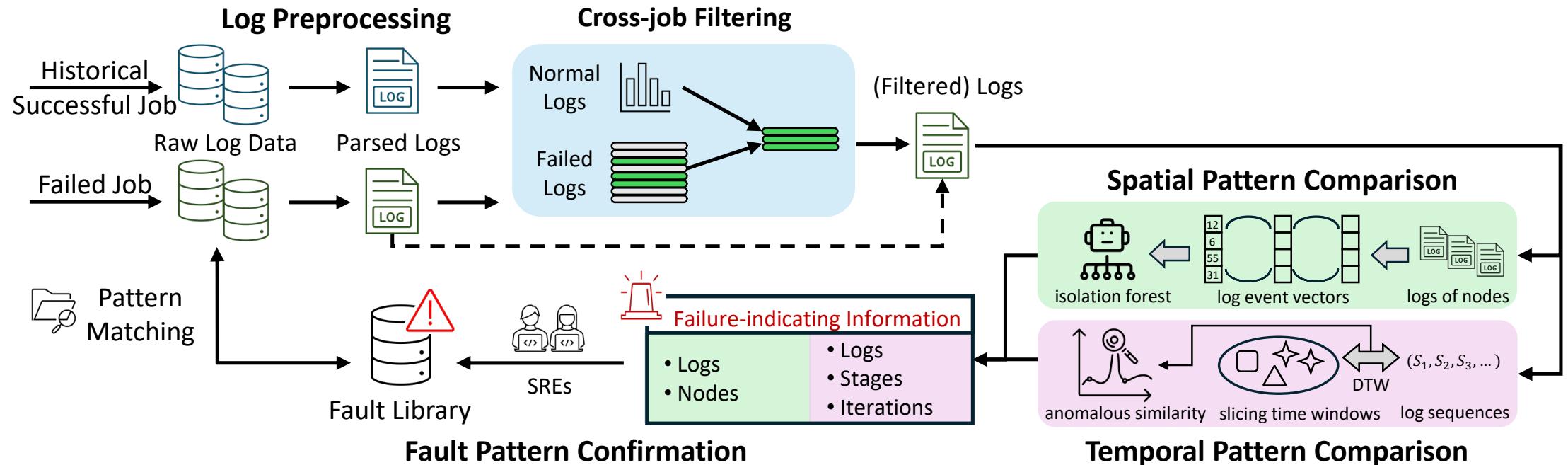
- Unlike traditional distributed systems, LLM training systems exhibit highly synchronized and nearly identical node workflows
- The distributed log sequences generated by different nodes are very similar

## ③ Temporal Pattern

- During the iterative training procedure, all nodes follow an identical workflow per iteration
- The deviation in log sequences across iterations can signal potential anomalies



# Our Framework: L4



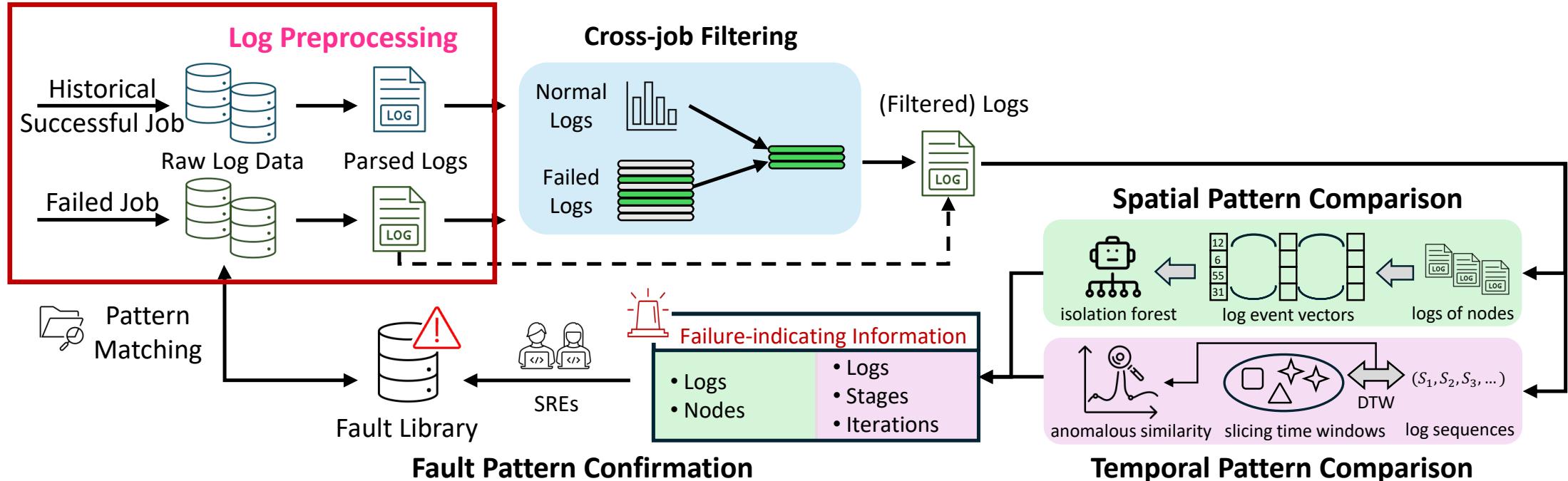
## L4: a Log-based Large-scale LLM training failure diagnosis framework



- automatically extract failure-indicating information from extensive training logs
- improve LLM training failure diagnostic efficiency



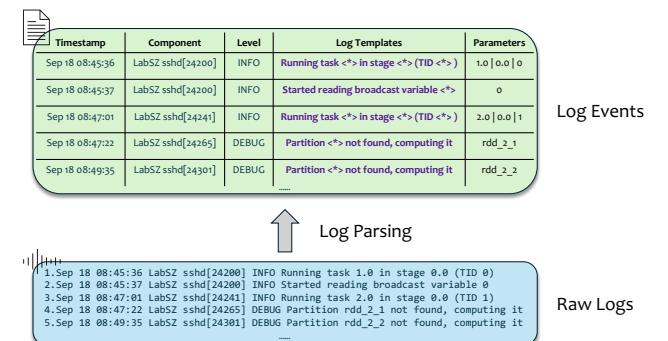
# Our Framework: L4



## Log Preprocessing

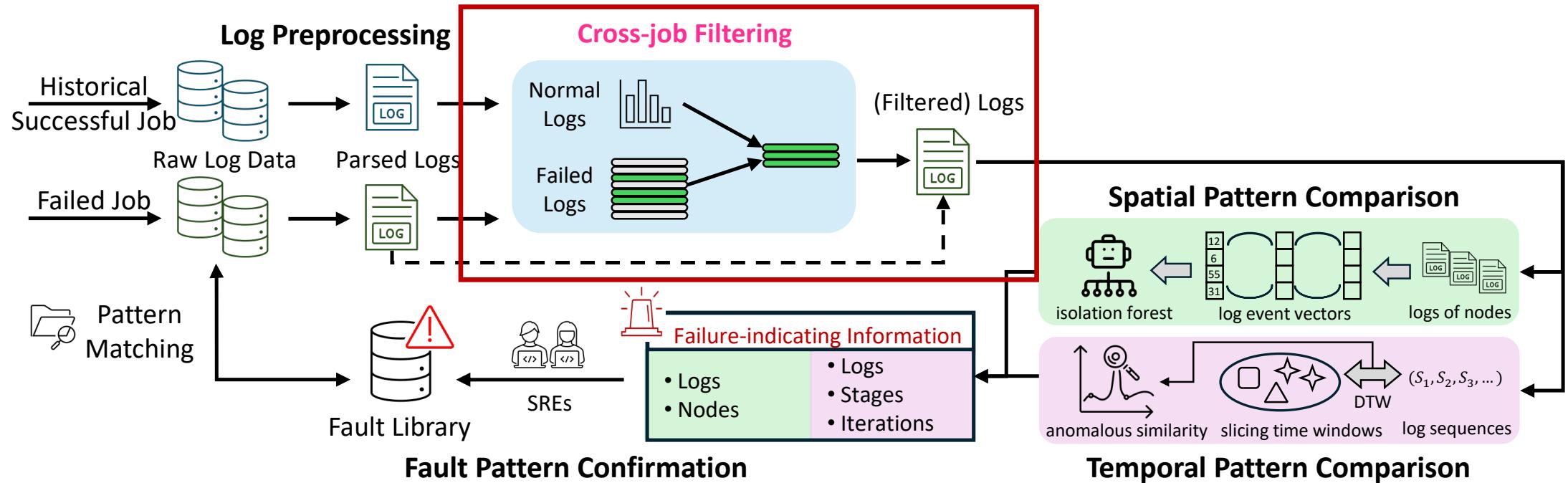
Convert unstructured raw training logs into structured log events:

- adopt the widely-used and most efficient log parser, Drain





# Our Framework: L4



## Cross-job Filtering

If there are related historical successful jobs:

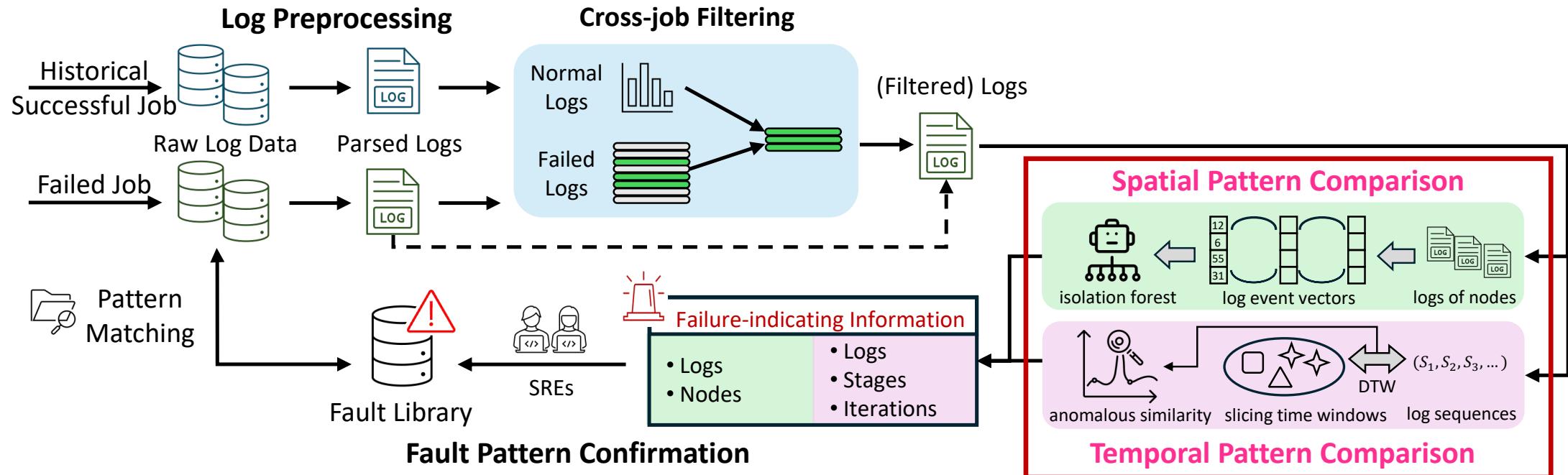
- constructing a normal log event pool, denoted as  $N = \{e_{n_1}, e_{n_2}, \dots\}$
- filtering frequent log events within  $N$

If there are no related historical successful jobs:

- using the parsed logs from failed jobs for analysis



# Our Framework: L4



## Spatial Pattern Comparision

Identify suspicious logs and nodes

- transforming the parsed log events of each node into log event vector  $V_i$
- employing *Isolation Forest* to detect deviant log event vectors across these nodes
  - unsupervised and interpretability

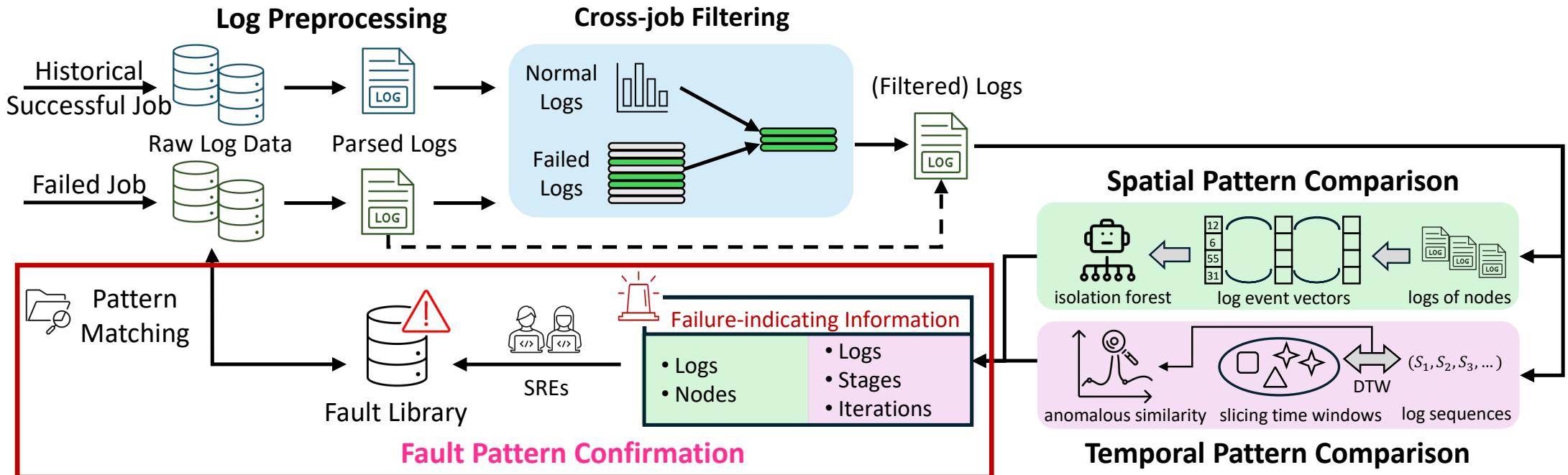
## Temporal Pattern Comparision

Identify suspicious logs and iterations

- converting the log events of each iteration into event sequence  $S_i$
- employing dynamic time warping (DTW) distance for similarity measurement
- using slicing time window with 10 iterations for detect anomalous iterations



# Our Framework: L4



## Fault Pattern Confirmation

Summarized fault patterns:

- failure-indicating logs, nodes, iterations and stages
- manually confirmed by engineers for future matching
- improving the efficiency of handling similar failures in the future

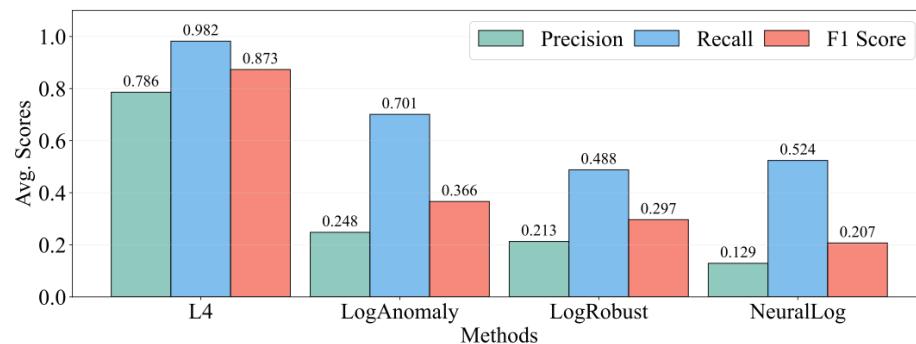


# Evaluation

## Research Questions

- RQ4: How effective is L4 in identifying failure-indicating logs
- RQ5: How effective is L4 in locating faulty nodes?

## RQ4: failure-indicating log identification

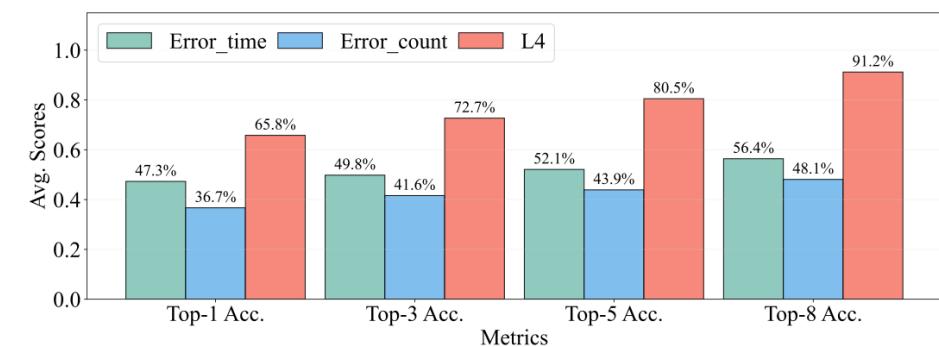


**L4 achieved the best accuracy in failure-indicating log identification, with an F1 score of 0.873.**

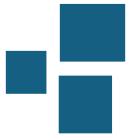
## Evaluation Datasets

- 100 randomly sampled failed LLM training jobs
  - averaging 632 accelerators and 12.3GB of training logs
- 42 cases were caused by hardware faults related to specific nodes

## RQ5: Faulty node localization



**L4 achieved the best accuracy in faulty node localization, with a Top-8 Accuracy of 91.2.**



# Case Study

## fine-grained localization of hardware fault

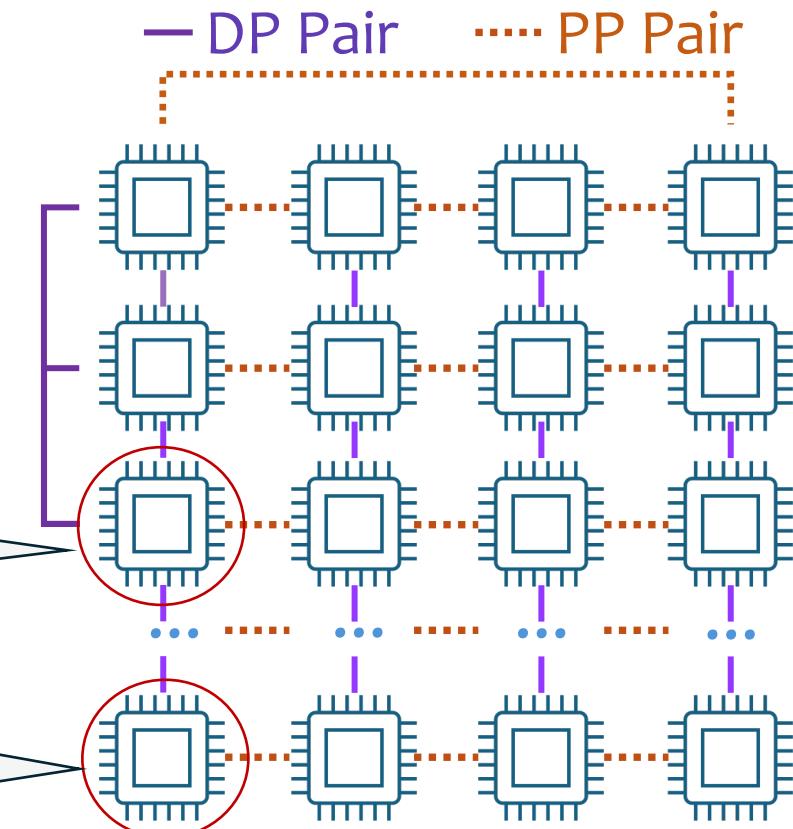
an LLM training job involving 1024 nodes with 4096 accelerators

failed after 20 hours of training,  
producing 71 GBs of latest training logs

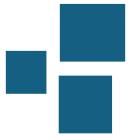
deployed for one year  
(since Jun. 2024)

log event n1.....  
log event n2.....  
ROCE(hccp\_service.bin):error cqe status.  
log event n3.....  
log event n4.....

log event n1.....  
notify wait from ... timeout  
log event n3.....



More case studies can be found in our paper



# Discussion

---

- **Future Research Directions**

- **LLM Training Monitoring**

- The observability of LLM training framework is still evolving
  - the logging quality should be improved
  - not only logs, but also efficient and effective profiling approach

- **Multi-modal Failure Diagnosis**

- *36.0% of failures require hybrid monitoring data for accurate diagnosis*
- Integrating multi-modal monitoring data for automated failure diagnosis and root cause analysis is essential for LLM training reliability

# Conclusion

## Empirical Study on LLM Training Failures

### Study Platform

Our multi-tenant production LLM training Platform-X supporting hundreds of internal users and partner companies

### Study Subject

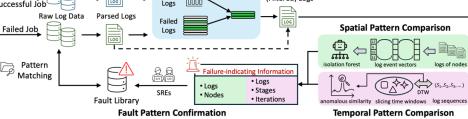
Collected 428 failure reports of failed large-scale LLM training jobs in Platform-X from May 2023 to April 2024  
• average model size: 72.8B parameters  
• average accelerators: 941 GPUs or NPUs

### Study RQs

- RQ1: What are the common symptoms of LLM training failures?
- RQ2: What are the common root causes of LLM training failures?
- RQ3: What monitor data sources are typically used to diagnose LLM training failures?

## Empirical Study

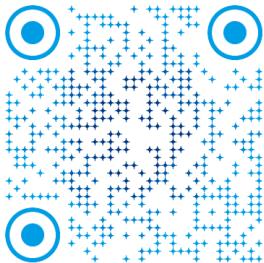
## Our Framework: L4



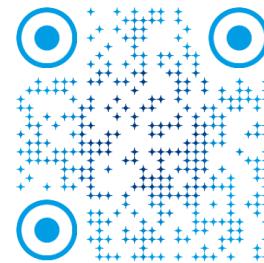
### L4: a Log-based Large-scale LLM training failure diagnosis framework

- automatically extract failure-indicating information from extensive training logs
- improve LLM training failure diagnostic efficiency

## Approach



Pre-print paper



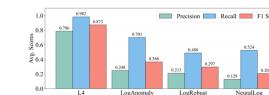
LogPAI website

## Evaluation

### Research Questions

- RQ4: How effective is L4 in identifying failure-indicating logs?
- RQ5: How effective is L4 in locating faulty nodes?

### RQ4: failure-indicating log identification

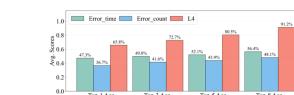


L4 achieved the best accuracy in failure-indicating log identification, with an F1 score of 0.873.

### Evaluation Datasets

- 100 randomly sampled failed LLM training jobs  
• averaging 632 accelerators and 12.3GB of training logs
- 42 cases were caused by hardware faults related to specific nodes

### RQ5: Faulty node localization



L4 achieved the best accuracy in faulty node localization, with a Top-8 Accuracy of 91.2.

## Experiment

Thank you for listening.

Q & A

Contact: [zhjiang@link.cuhk.edu.hk](mailto:zhjiang@link.cuhk.edu.hk)