

MegaScale: Scaling Large Language Model Training to More Than 10,000 GPUs

Ziheng Jiang^{1,*} Haibin Lin^{1,*} Yinmin Zhong^{2,*}

Qi Huang¹ Yangrui Chen¹ Zhi Zhang¹ Yanghua Peng¹ Xiang Li¹ Cong Xie¹ Shibiao Nong¹
Yulu Jia¹ Sun He¹ Hongmin Chen¹ Zhihao Bai¹ Qi Hou¹ Shipeng Yan¹ Ding Zhou¹ Yiyao
Sheng¹ Zhuo Jiang¹ Haohan Xu¹ Haoran Wei¹ Zhang Zhang¹ Pengfei Nie¹ Leqi Zou¹ Sida
Zhao¹ Liang Xiang¹ Zherui Liu¹ Zhe Li¹ Xiaoying Jia¹ Jianxi Ye¹ Xin Jin², Xin Liu¹

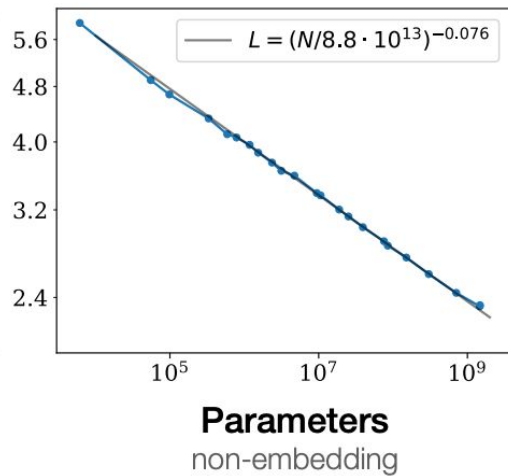
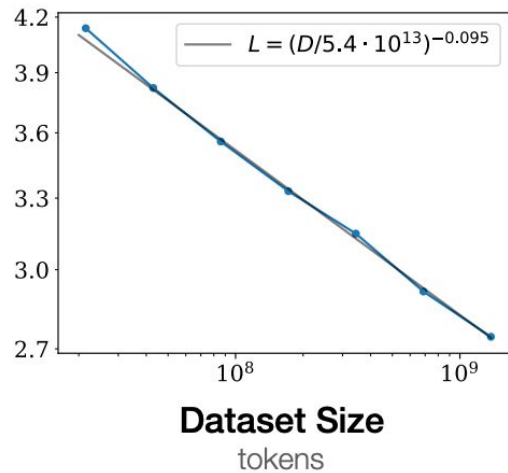
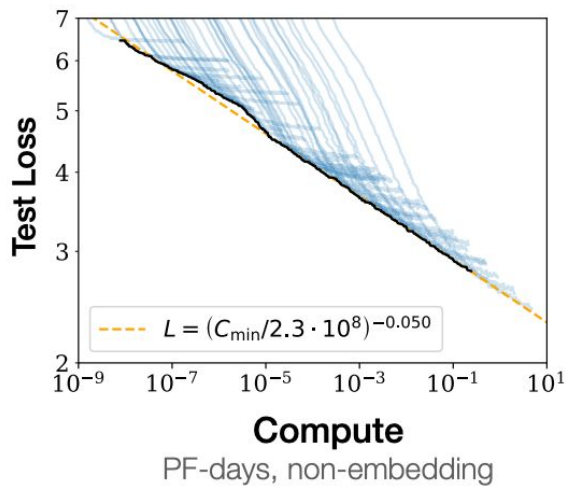
¹ByteDance ²Peking University



Large Language Model (LLM) Applications

- conversational agents
- code development
- content creation
- ...

The power behind LLMs is scaling



Source: Jared Kaplan, et al. "Scaling Laws for Neural Language Models". In preprint 2020.

Compute required to train an LLM

$$C = \tau T \approx 6ND$$

τ : cluster's effective throughput

T : training time

N : model size (#parameters)

D : dataset size (#tokens)

Source: Jared Kaplan, et al. "Scaling Laws for Neural Language Models". In preprint 2020.

Challenges

- Achieve high training efficiency at scale
 - LLM training is not embarrassingly parallel
 - communication and many other factors contribute significantly to the efficiency

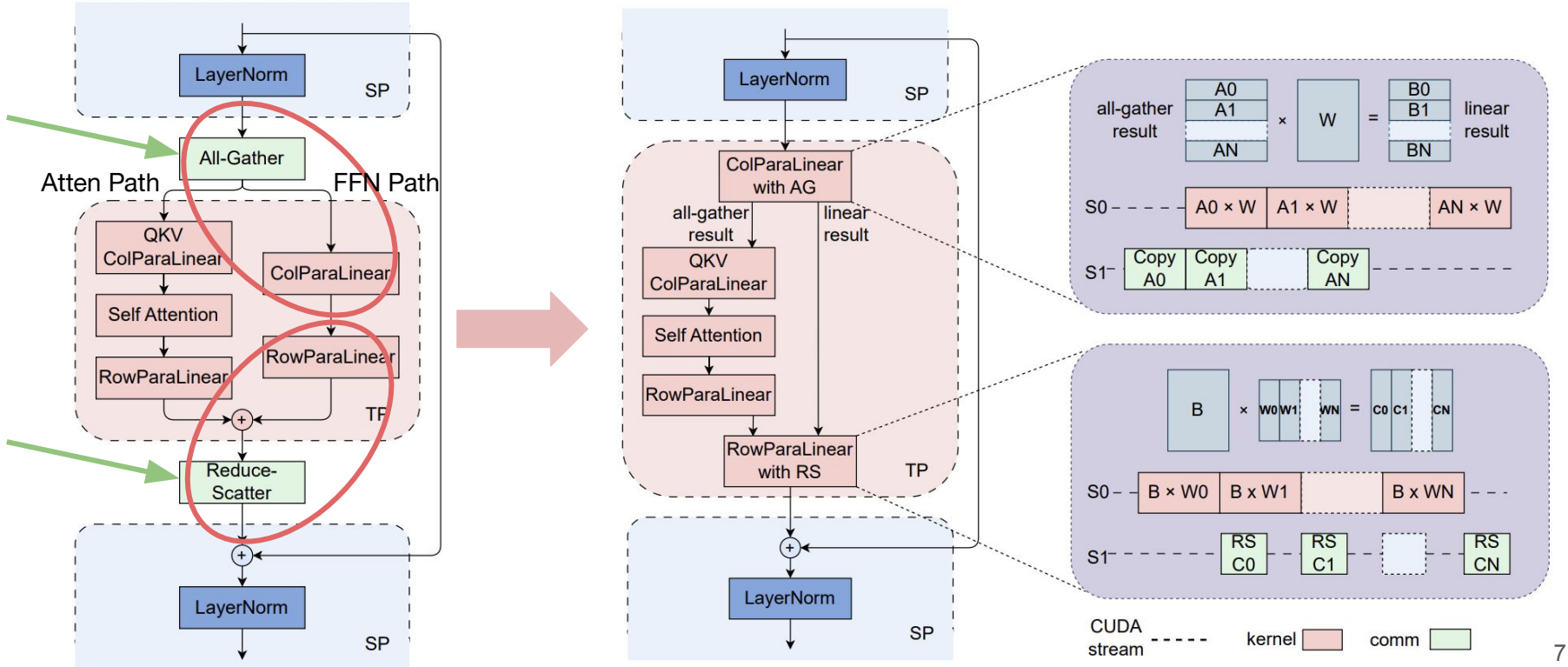
- Achieve high training stability at scale
 - the training duration can extend beyond months
 - failures and stragglers are the norm for LLM training

Our solution – MegaScale

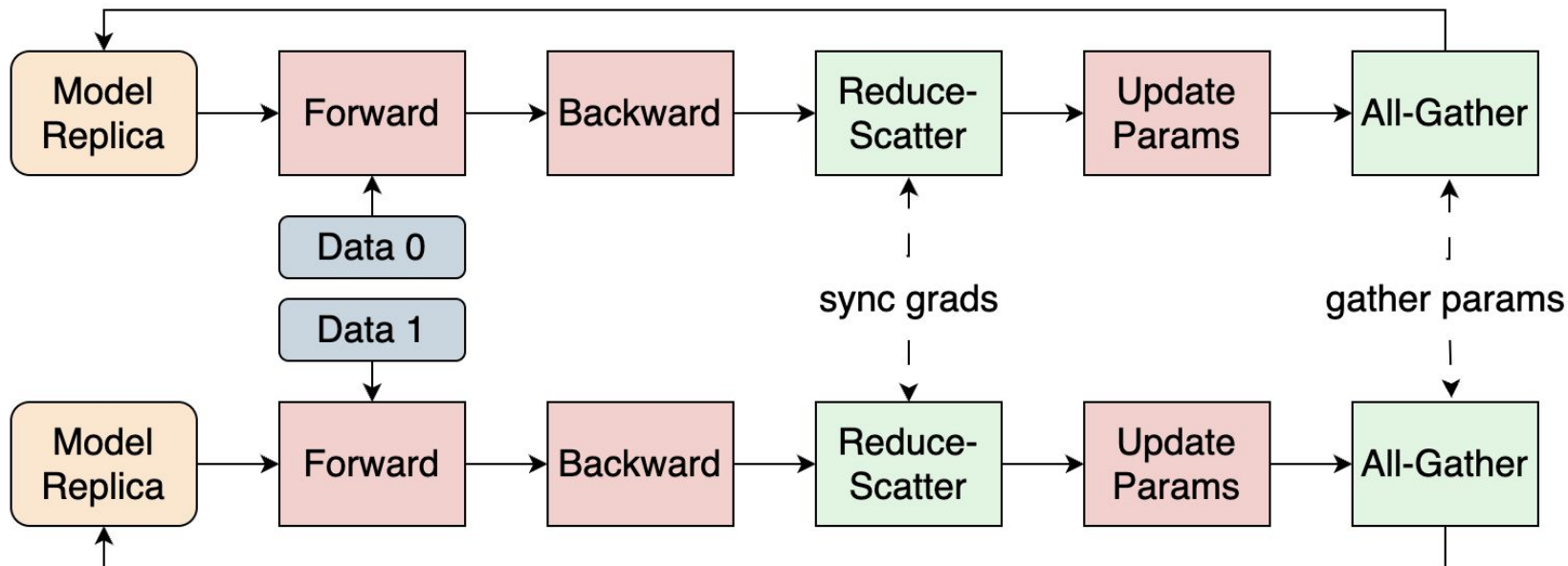
- Achieve high training efficiency at scale
 - Algorithmic optimizations
 - Communication overlapping in 3D parallelism
 - Training data preprocessing and loading
 - Collective communication group initialization
 - Network performance tuning
- Achieve high training stability at scale



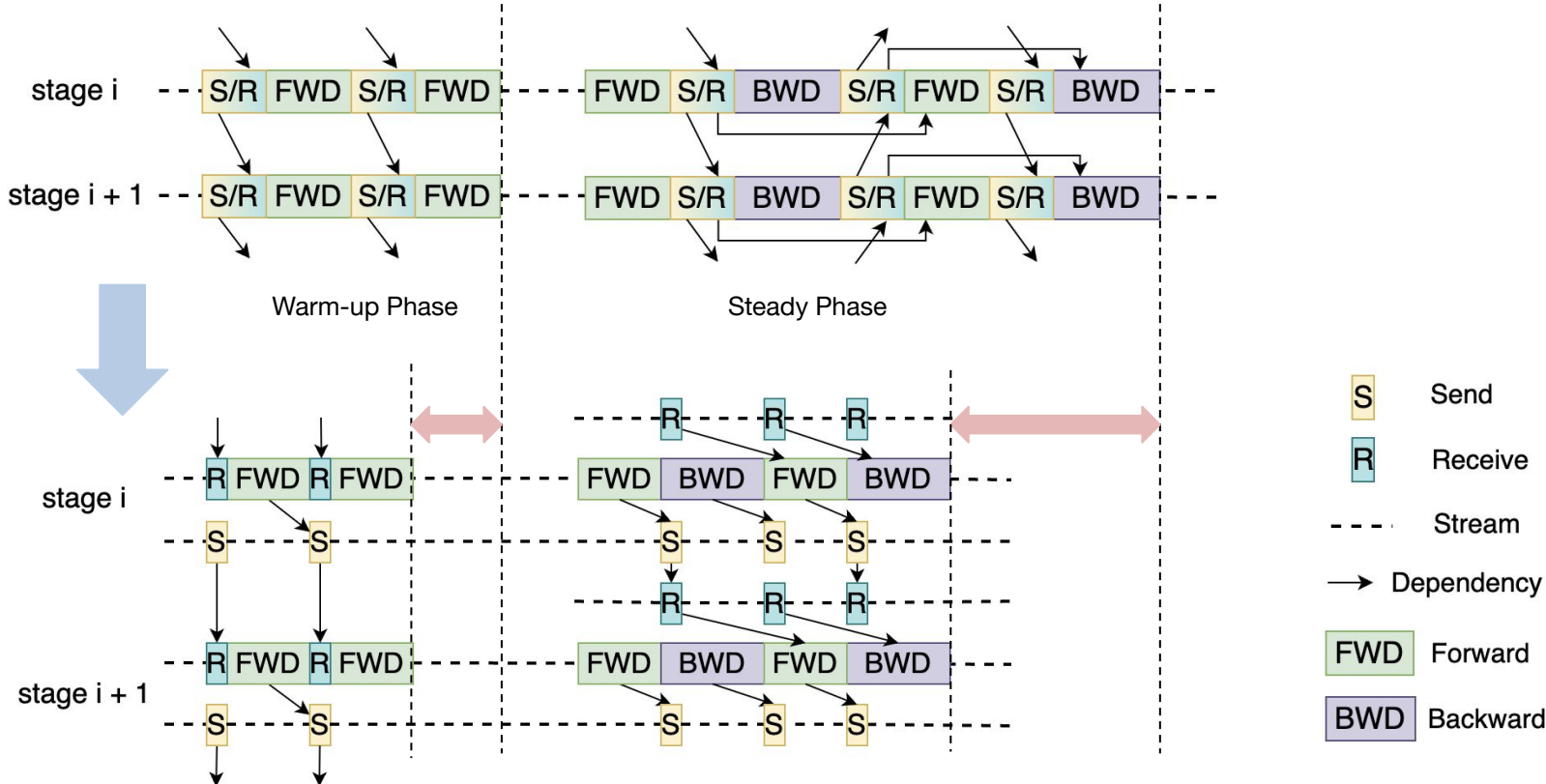
Overlapping in tensor/sequence parallelism



Overlapping in data parallelism (with ZeRO)



Overlapping in pipeline parallelism

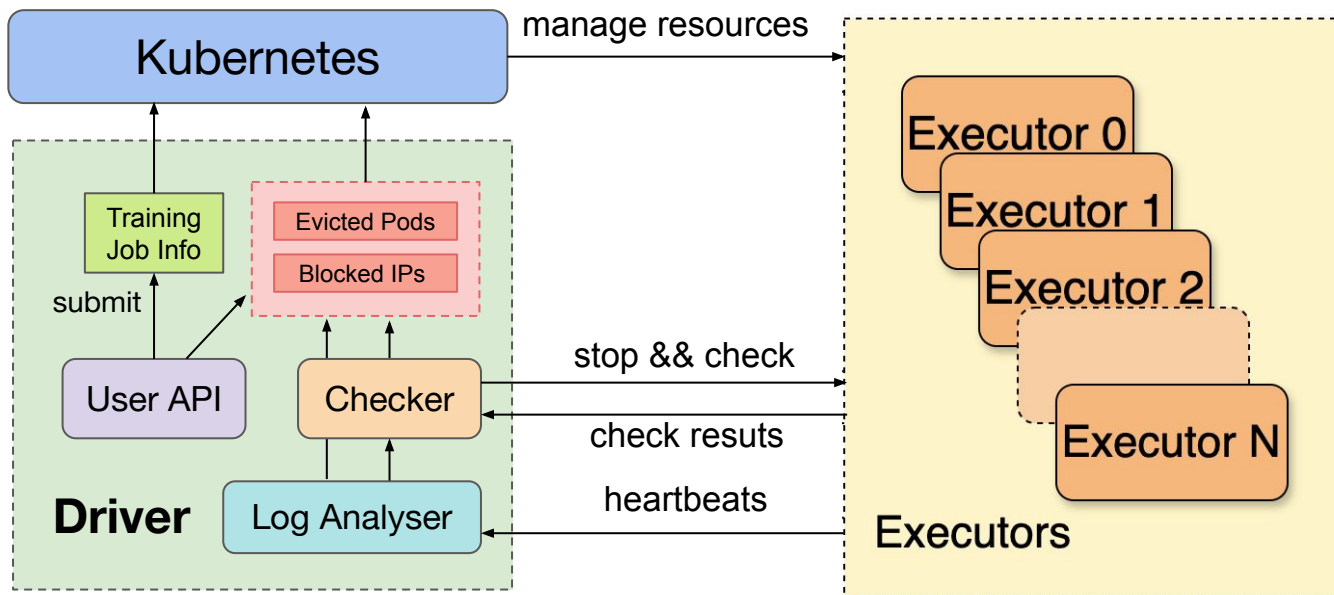


Our solution – MegaScale

- Achieve high training efficiency at scale
 - Algorithmic optimizations
 - Communication overlapping in 3D parallelism
 - Training data preprocessing and loading
 - Collective communication group initialization
 - Network performance tuning
- Achieve high training stability at scale
 - Robust training framework
 - Fast checkpointing and recovery
 - Monitoring and analysis tools



Robust Training Workflow



Deployment Experience

- Training performance
- Training stability
- Problems discovered and fixed

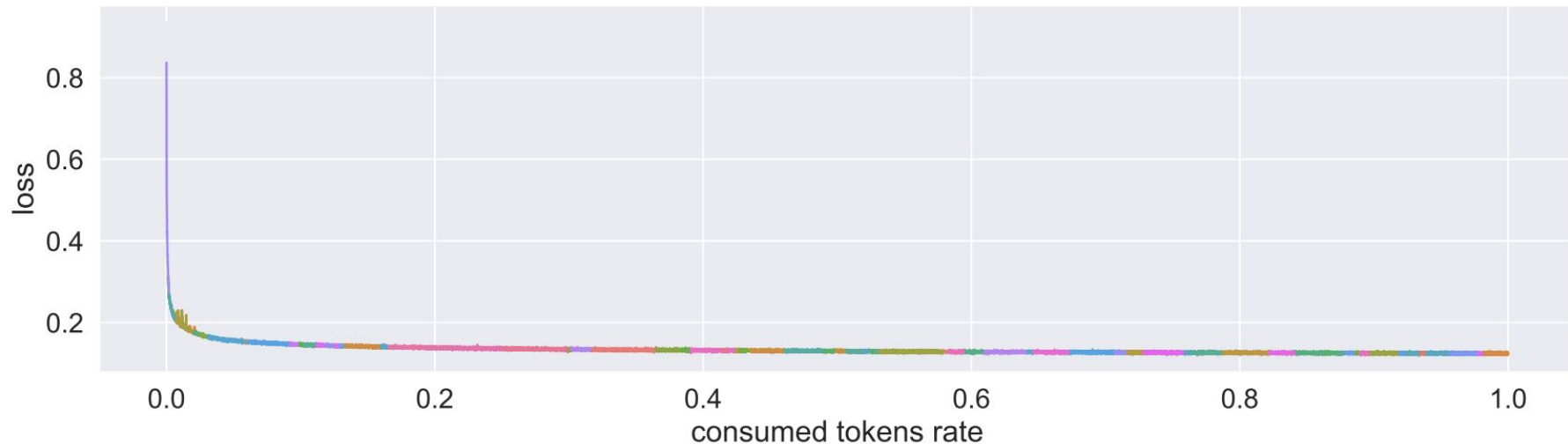
Training Performance

Strong-scaling training performance for the 175B model over 300B tokens compared to Megatron-LM.

Batch Size	Method	GPUs	Iteration Time (s)	Throughput (tokens/s)	Training Time (days)	MFU	Aggregate PFlops/s
768	Megatron-LM	256	40.0	39.3k	88.35	53.0%	43.3
		512	21.2	74.1k	46.86	49.9%	77.6
		768	15.2	103.8k	33.45	46.7%	111.9
		1024	11.9	132.7k	26.17	44.7%	131.9
	MegaScale	256	32.0	49.0k	70.86	65.3%(1.23 ×)	52.2
		512	16.5	95.1k	36.51	63.5%(1.27 ×)	101.4
		768	11.5	136.7k	25.40	61.3%(1.31 ×)	146.9
		1024	8.9	176.9k	19.62	59.0%(1.32 ×)	188.5
6144	Megatron-LM	3072	29.02	433.6k	8.01	48.7%	466.8
		6144	14.78	851.6k	4.08	47.8%	916.3
		8192	12.24	1027.9k	3.38	43.3%	1106.7
		12288	8.57	1466.8k	2.37	41.2%	1579.5
	MegaScale	3072	23.66	531.9k	6.53	59.1%(1.21 ×)	566.5
		6144	12.21	1030.9k	3.37	57.3%(1.19 ×)	1098.4
		8192	9.56	1315.6k	2.64	54.9%(1.26 ×)	1400.6
		12288	6.34	1984.0k	1.75	55.2%(1.34 ×)	2166.3

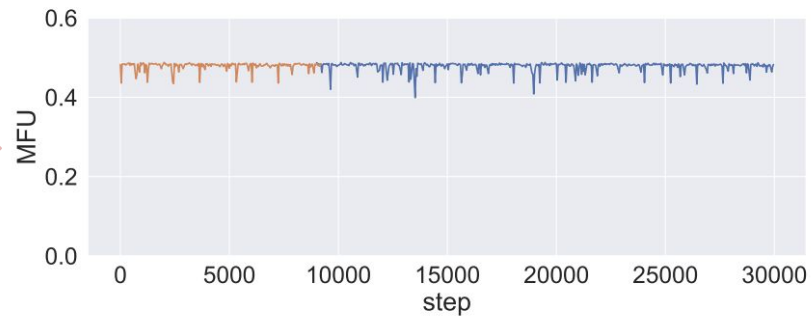
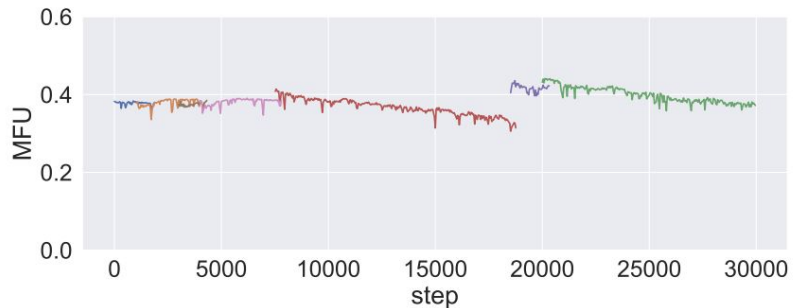
Training stability

The loss curve of a real production run that trains a proprietary model with hundreds of billions of parameters on multi-trillion tokens.



Problems discovered and fixed

- Computational stragglers
- MFU decreasing



Summary of MegaScale

- Achieve high training efficiency at scale
 - Algorithmic optimizations
 - Communication overlapping in 3D parallelism
 - Training data preprocessing and loading
 - Collective communication group initialization
 - Network performance tuning
- Achieve high training stability at scale
 - Robust training framework
 - Fast checkpointing and recovery
 - Monitoring and analysis tools



github.com/volcengine/veScale



Large Language Models (LLMs)