

OLCF Training

# AI on Frontier

**Junqi Yin**

Analytics and AI Methods at Scale



ORNL is managed by UT-Battelle, LLC for the US Department of Energy

# Outline

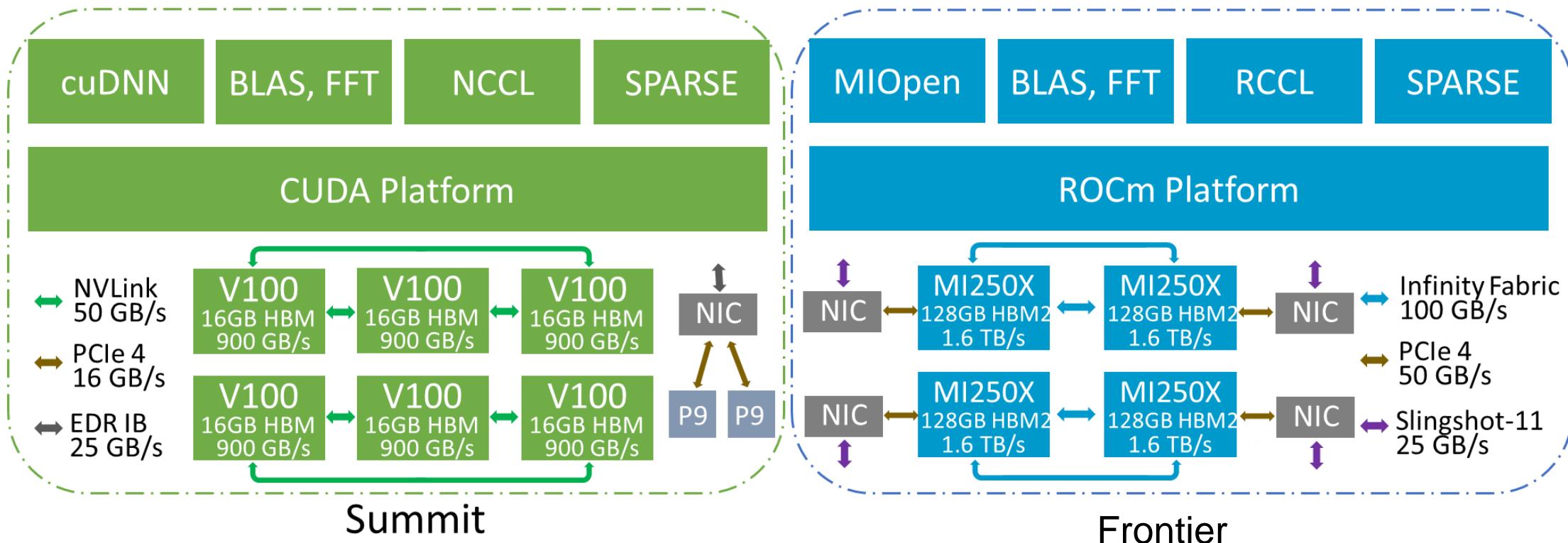
- Frontier DL Environment
- Distributed Training
- Examples
  - PyTorch
  - TensorFlow
- Training Large Language Models (LLMs)

# Deep Learning Stacks

Framework

TensorFlow, PyTorch, ...

DL Stack



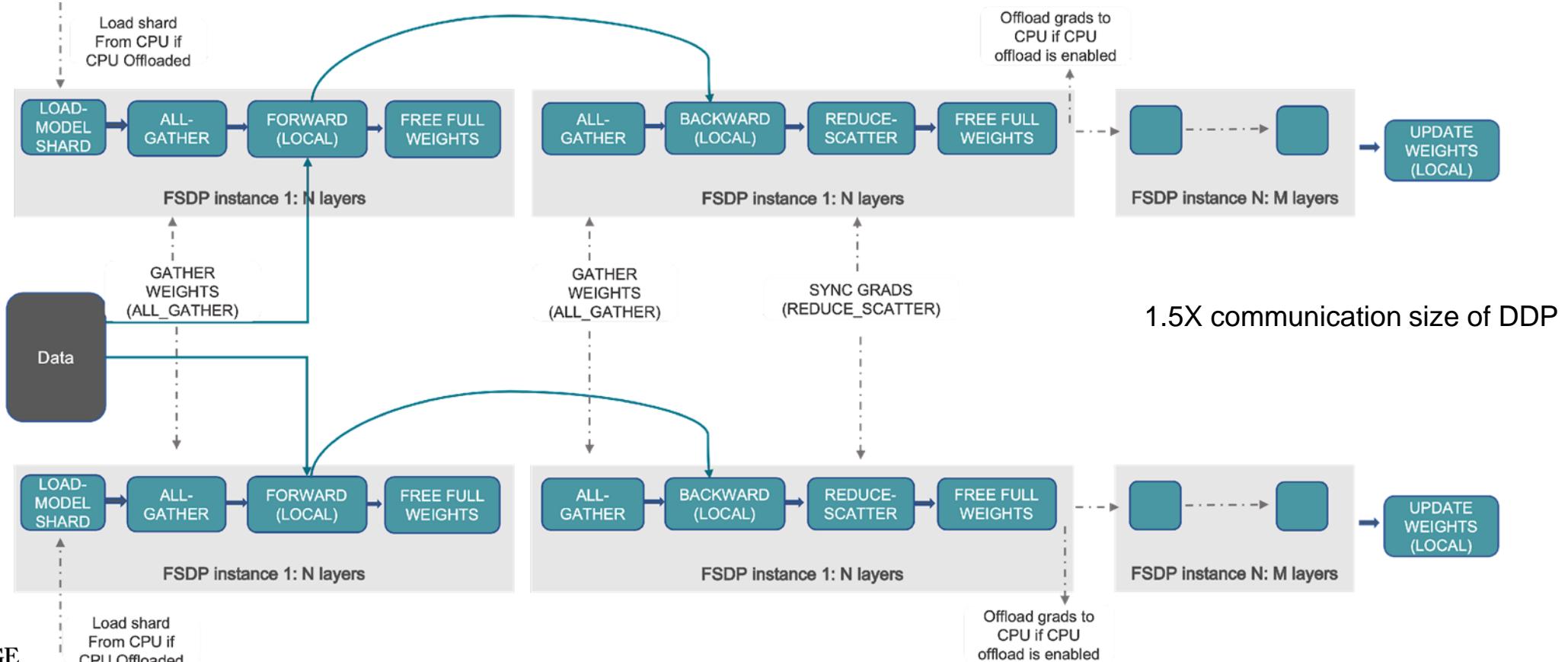
- Most DL codes work on Frontier as is

# Distributed Training Strategies

- Data Parallel
  - PyTorch DDP, TensorFlow Multiworker Mirrored strategy
  - Horovod
- Model Parallel
  - PyTorch FSDP
  - DeepSpeed-Megatron 3D parallelism

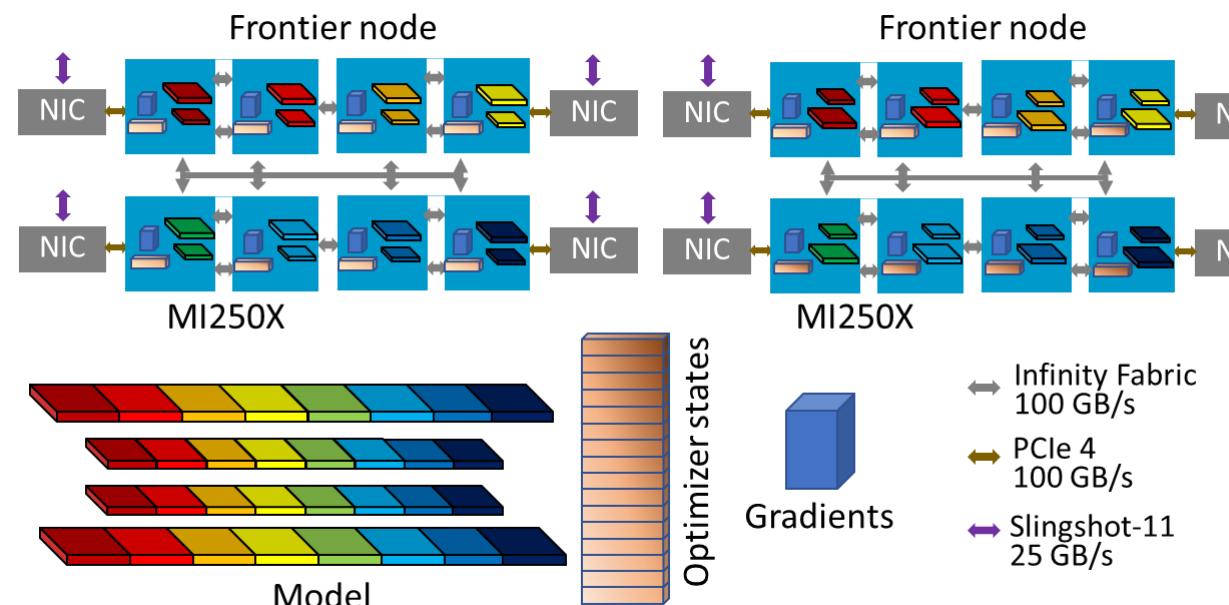
# Distributed Training Strategies: FSDP

- FSDP – fully sharded data parallel, similar to ZeRO



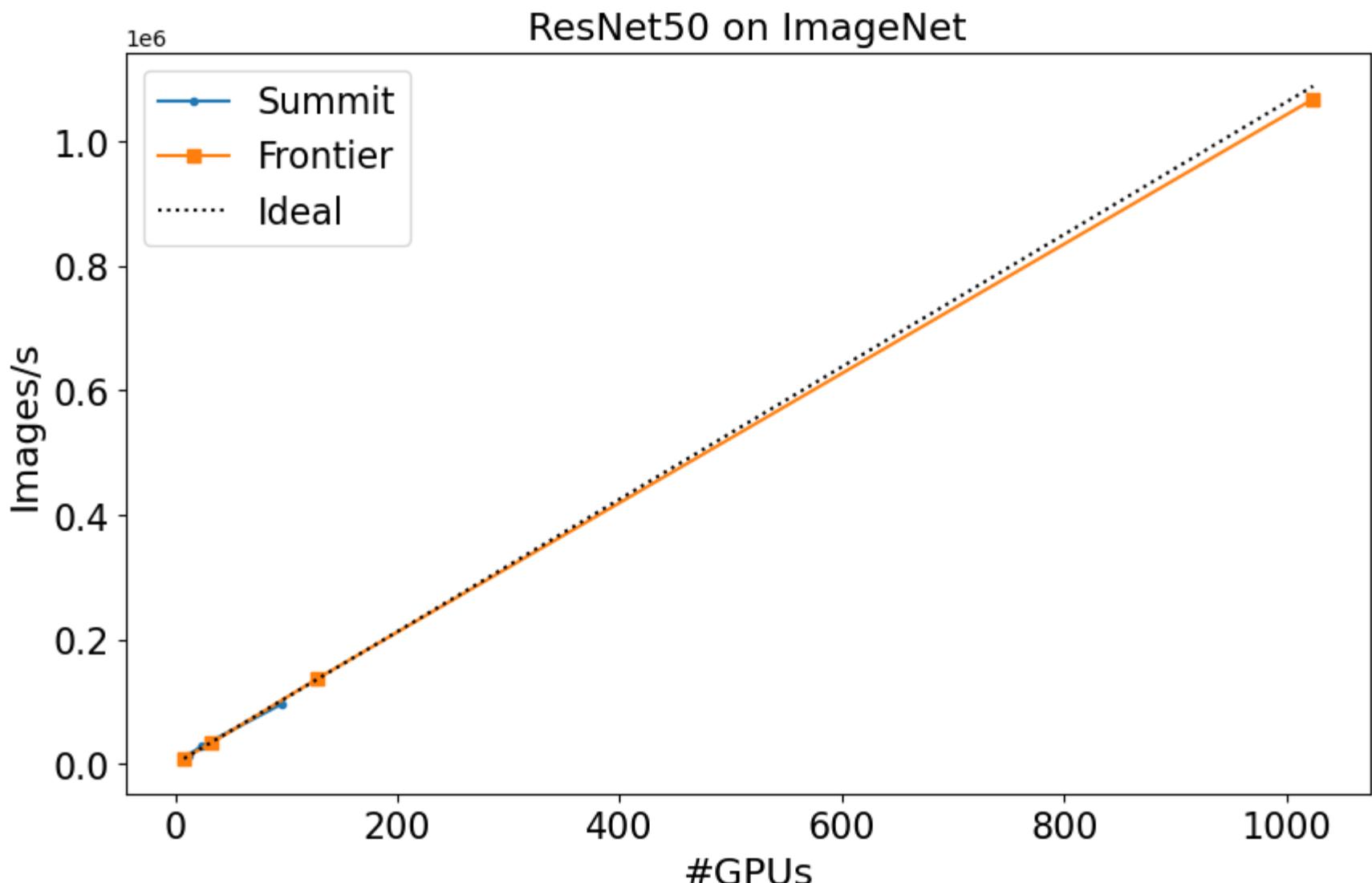
# Distributed Training Strategies: DeepSpeed+Megatron

- DeepSpeed + Megatron
  - Sharded data parallel – ZeRO, shared optimizer states, gradients, and parameters
  - Pipeline parallel – layer placement + micro-batch (Gpipe)
  - Tensor parallel – sharded tensor



# Performance baselines: data-parallel

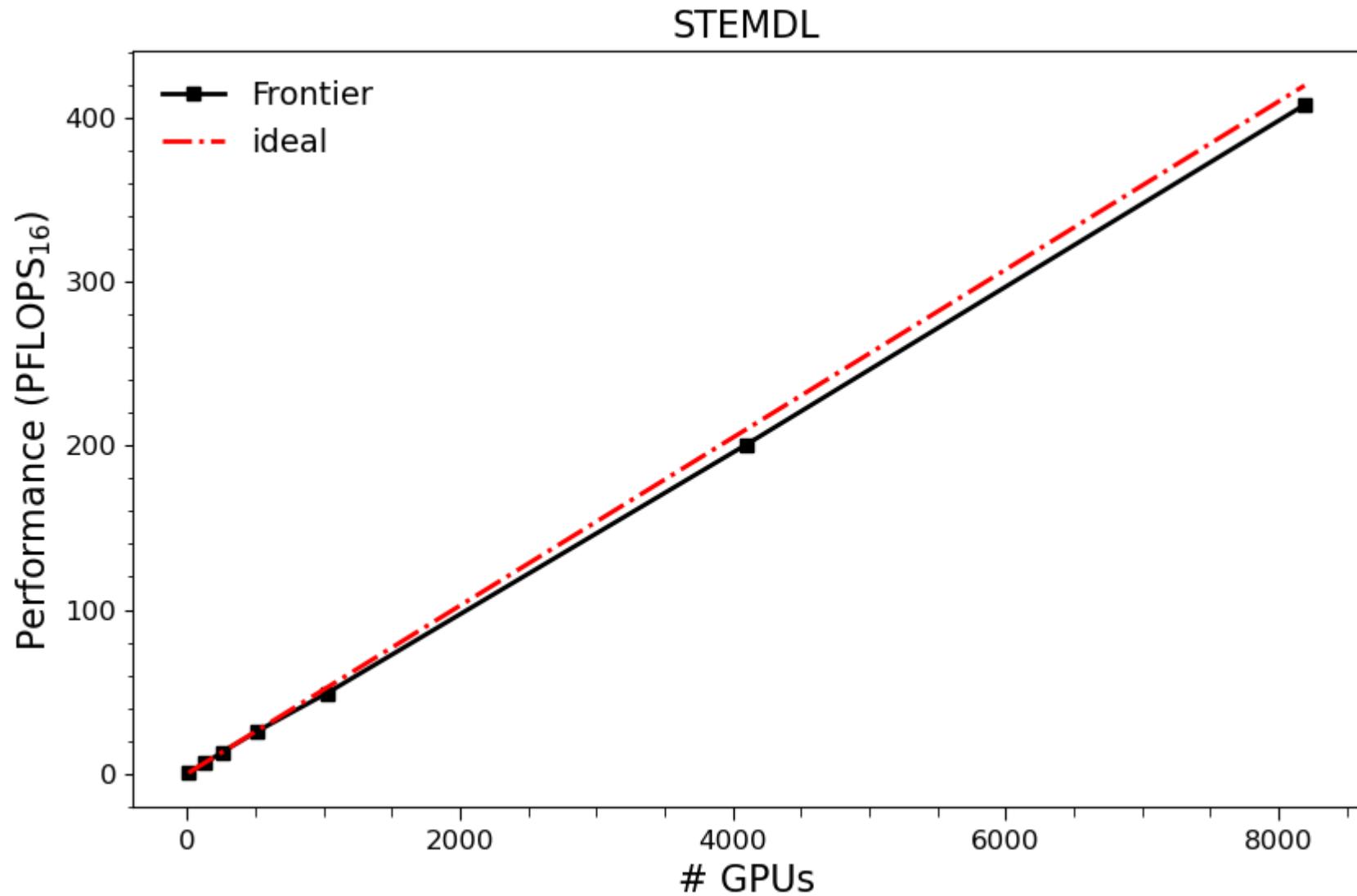
- ResNet50
  - Mixed
  - ~1.0x per GCD
  - 98% at 1024



# Performance baselines: data-parallel

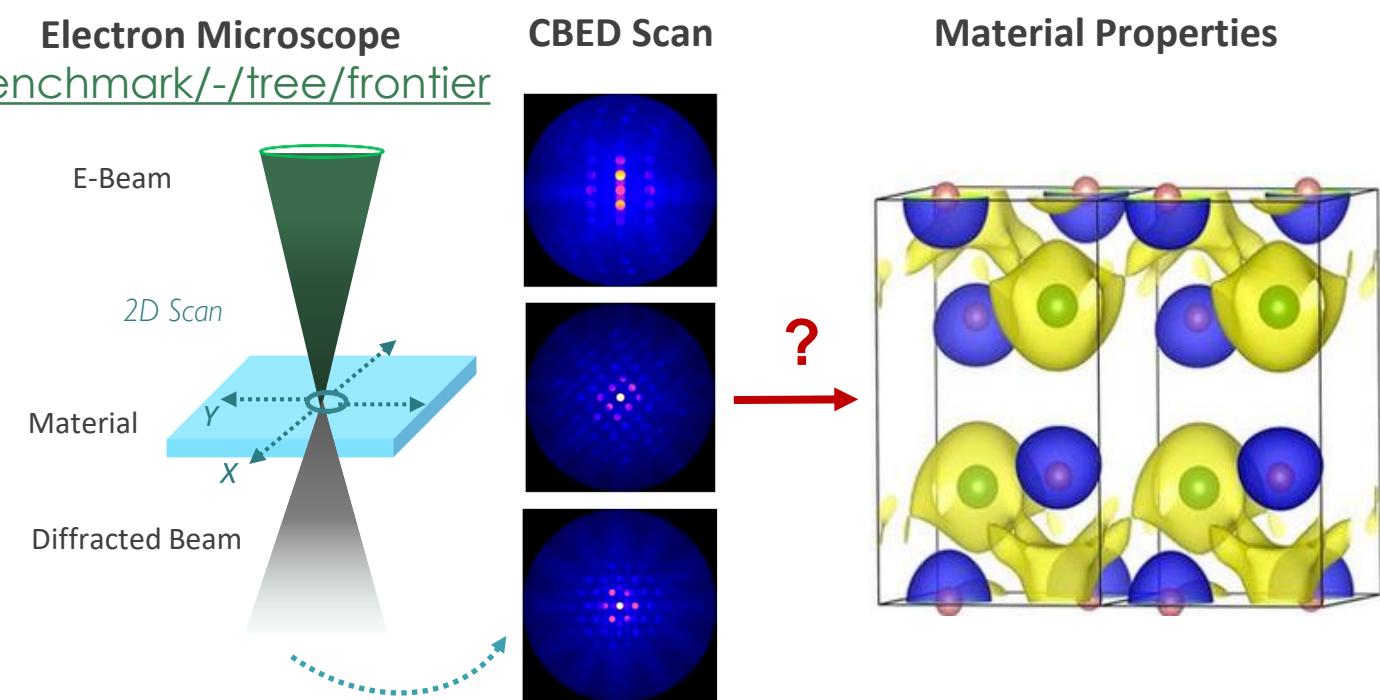
- STEMDL
  - Tiramisu network
  - 220M parameters
  - 97% at 8192 GPUs

Accelerating Collective  
Communication in Data  
Parallel Training across Deep  
Learning Frameworks.  
USENIX NSDI'22, 2022



# Distributed Training Examples: STEM DL Benchmark

- One of MLCommon Science benchmarks
- Code on Frontier
  - <https://code.ornl.gov/jqyin/stemdl-benchmark/-/tree/frontier>
  - Classification – PyTorch
  - Reconstruction – TensorFlow
  - Built/job scripts
  - Conda environment



# Distributed Training Examples: STEM DL Benchmark

code.ornl.gov/jqyin/stemdl-benchmark/-/tree/frontier?ref\_type=heads

## Software Environment

A pre-built environment is provided at [/lustre/orion/world-shared/stf218/junqi/miniconda](#).

```
module load rocm/5.1.0
BUILD=/lustre/orion/world-shared/stf218/jur
export PATH=${BUILD}/bin:$PATH
source ${BUILD}/etc/profile.d/conda.sh
```

To activate PyTorch env:

```
conda activate ${BUILD}/envs/pyt_env
```

To activate TensorFlow env:

```
conda activate ${BUTLD}/envs/tf_env
```

To build from source, please follow the scripts in `utils`

## 🔗 Output

The classification task outputs the mlperf and tensorboard logs to `logs` folder (default)

```
#stemdl_classification.log
:::MLLOG {"namespace": "", "time_ms": 1692106306826, "event_type": "POINT_IN_TIME", "key": "eval_accuracy", "value": 7.01891407370
:::MLLOG {"namespace": "", "time_ms": 1692106306826, "event_type": "POINT_IN_TIME", "key": "f1_score", "value": 2.2504812106490135
:::MLLOG {"namespace": "", "time_ms": 1692106306826, "event_type": "INTERVAL_END", "key": "eval_stop", "value": null, "metadata": {
:::MLLOG {"namespace": "", "time_ms": 1692106307223, "event_type": "INTERVAL_START", "key": "epoch_stop", "value": null, "metadata":
```

The reconstruction task outputs the training performance in FLOPS and loss values to `output_logs` folder (default)

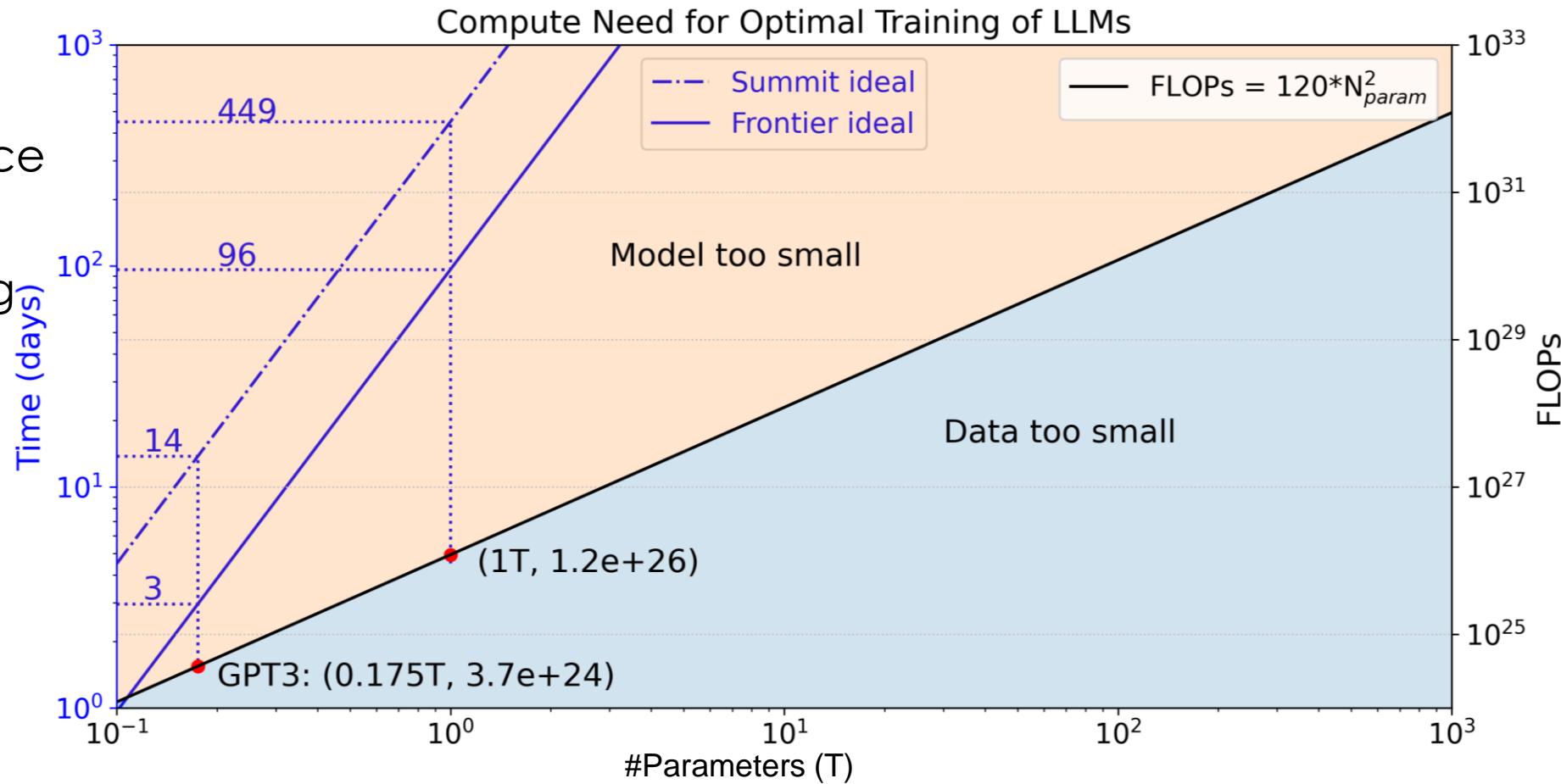
```
558.0, step= 10, epoch= 9.55e-03, loss= 0.52, lr= 1.00e-06, step_time= 15.93 sec, ranks= 8, examples/sec= 4.0, flops = 4.40e+13, ave  
75.1, step= 20, epoch= 1.91e-02, loss= 0.46, lr= 1.00e-06, step_time= 1.71 sec, ranks= 8, examples/sec= 37.5, flops = 4.10e+14, ave  
92.3, step= 30, epoch= 2.87e-02, loss= 0.53, lr= 1.00e-06, step_time= 1.71 sec, ranks= 8, examples/sec= 37.4, flops = 4.09e+14, ave  
99.4, step= 40, epoch= 3.82e-02, loss= 0.52, lr= 1.00e-06, step_time= 1.71 sec, ranks= 8, examples/sec= 37.4, flops = 4.09e+14, ave
```

Data

Sample data (pre-processed) are located at `/lustre/orion/world-shared/stf218/junqi/stemdl-data`

# Training LLMs on Frontier

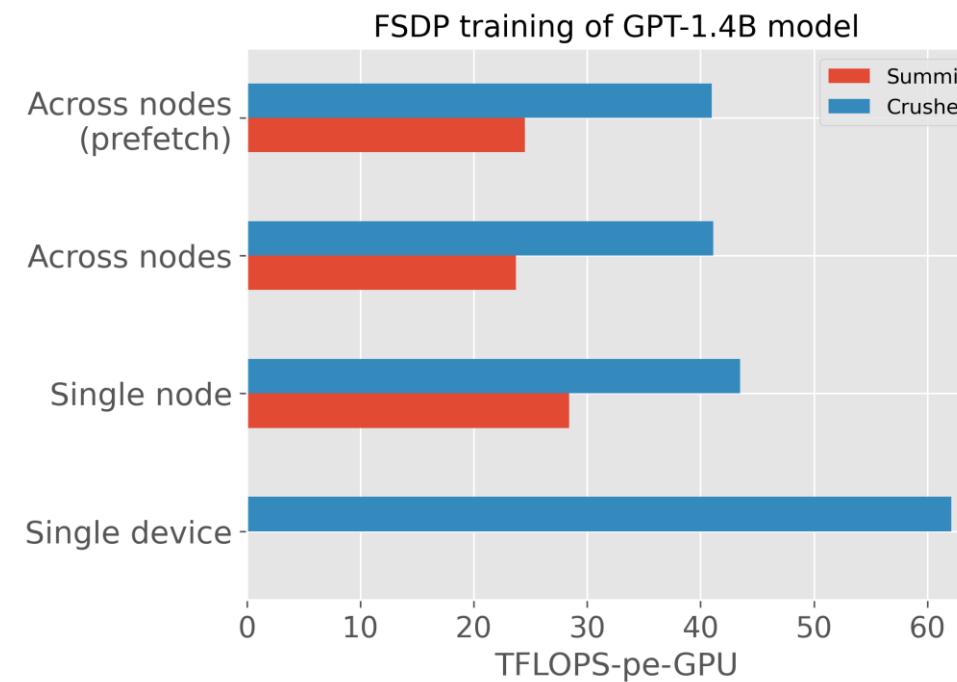
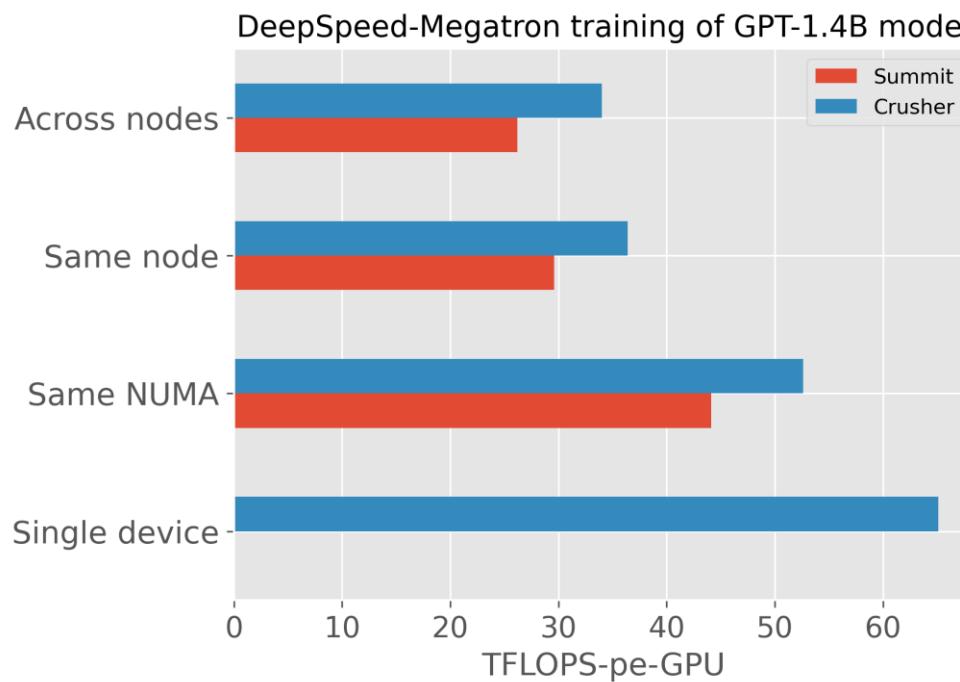
- Ideal case
  - Peak performance
  - Linear scaling
  - Chinchilla scaling



Evaluation of Pre-Training Large Language Models on Leadership Platforms, Journal of Supercomputing, 2023

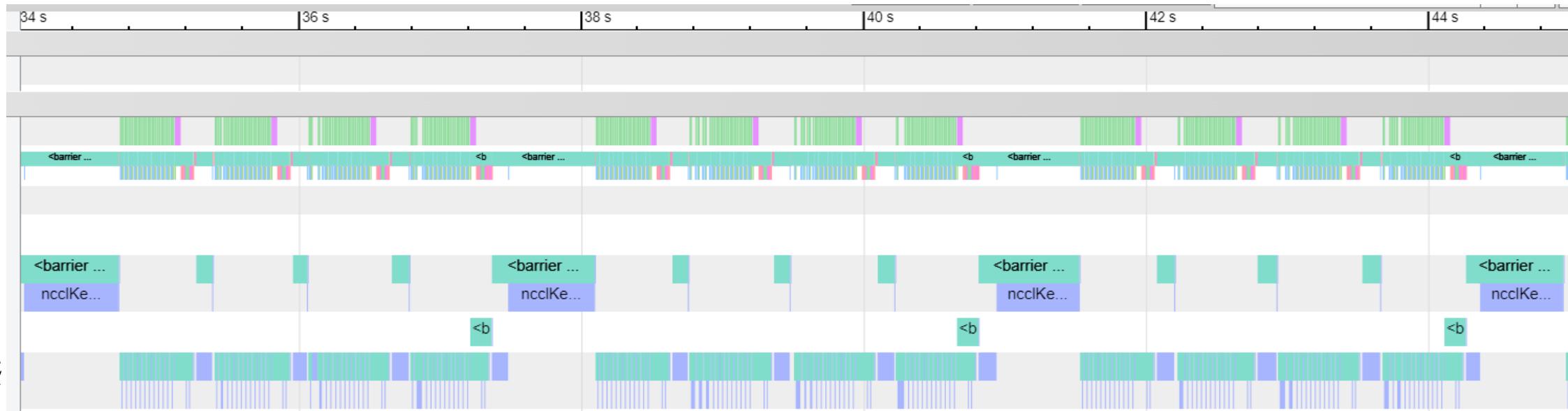
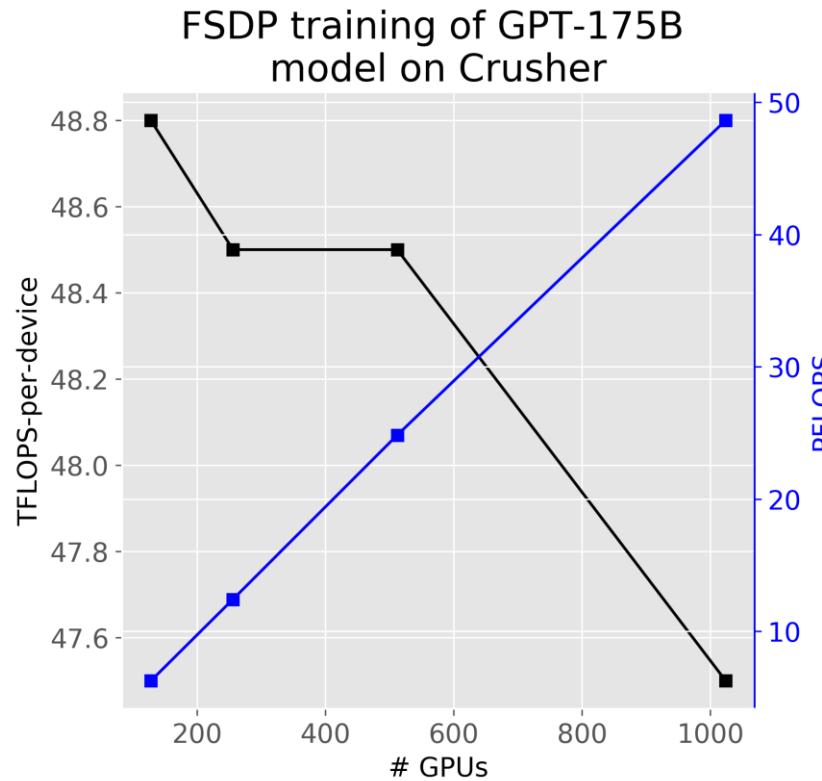
# Training LLMs on Frontier

- Actual measurement
  - 2 frameworks: FSDP, DeepSpeed-Megatron
  - 3D parallelisms: DP, TP, PP



# Training LLMs on Frontier

- Scaling
  - DP is preferable
  - Large batch
  - bf16 more stable



# Training LLMs on Frontier

- Computation and energy needed to pretrain LLMs

Model	Summit			Frontier		
	Time (node hours)	Energy (MWh)	TFLOPS/ Watt	Time (node hours)	Energy (MWh)	TFLOPS/ Watt
GPT-13B	23K	30	0.185	12K	24	0.239
GPT-175B	[0.6M, 12M]	[309, 6176]	0.165	[0.1M, 2.6M]	[217, 4338]	0.235
GPT-1T	[1B, 21B]	[3.7e5, 7.4e6]	0.01	[8M, 160M]	[4727, 94539]	0.343

[ # tokens = # params, # tokens = 20\*# params ]

- Takeaway: pre-train O(13B) ~ O(DD), pre-train O(175B) ~ O(INCITE)

# Simulation-ML Integration

Method	Pro	Con
Framework C++ API (TensorFlow/PyTorch C++)	<ul style="list-style-type: none"><li>Portable</li><li>Better latency</li><li>Easy to deploy</li></ul>	<ul style="list-style-type: none"><li>Not flexible</li></ul>
Framework Server (TensorFlow Serving/TorchServe)	<ul style="list-style-type: none"><li>Flexible</li><li>Better throughput</li><li>Options to deploy</li></ul>	<ul style="list-style-type: none"><li>High maintenance</li></ul>
Third-party API (SmartRedis/RedisAI)	<ul style="list-style-type: none"><li>Easy integration</li><li>More functionality</li><li>Model support</li></ul>	<ul style="list-style-type: none"><li>Portability</li></ul>

Strategies for Integrating Deep Learning Surrogate Models with HPC Simulation Applications, IPDPSW 2022

Questions ?

[yinj@ornl.gov](mailto:yinj@ornl.gov)