



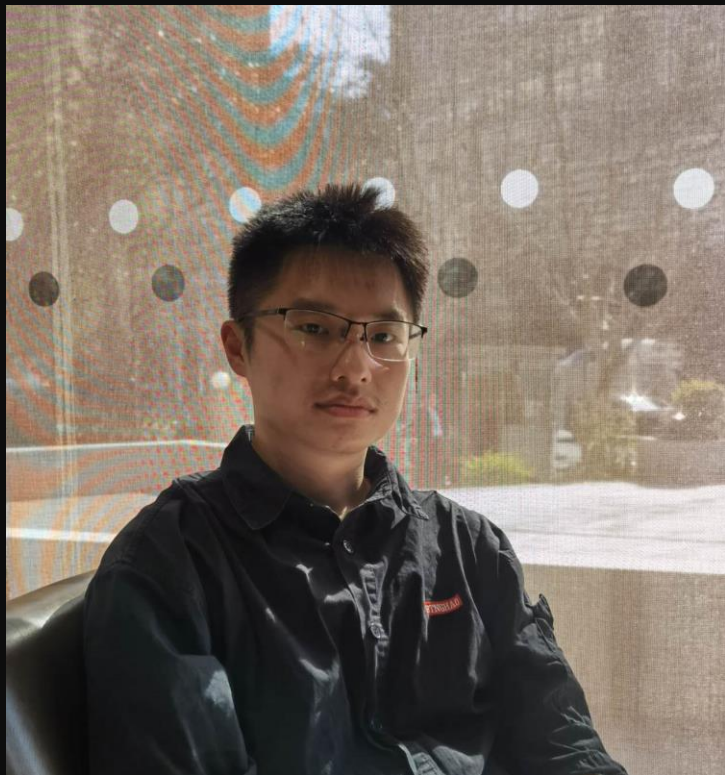
A PyTorch-Native Auto-Parallel Framework for *Ease of Use*

veScale Team

ByteDance

2024-8-22

About Me



Hongyu Zhu

Received my PhD degree from
University of Toronto (advisor:
Gennady Pekhimenko)

Joined ByteDance AML group in
2022.3

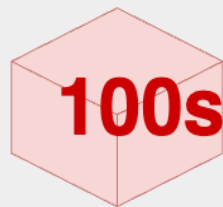
Currently working on LLM training
frameworks

Agenda

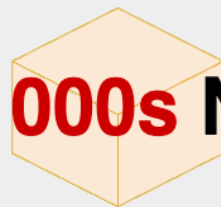
- Why veScale
- What is veScale
- Preliminary Results of veScale
- Future of veScale

Why veScale

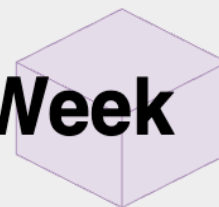
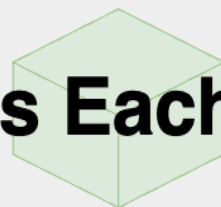
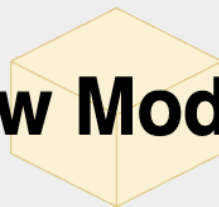
Company:



100s~1000s



New Models Each Week



Industrial Training Framework



Only Performance



Ease of Use



Many
weeks to
write one
model

But Current Frameworks are **Hard to Use**

Not PyTorch

**Intertwined System
and Model Design**

**Not Automated
Enough**

GradBuffer Defrag
AllReduce Overlap

nn.Linear

ColumnParallelLinear

Hard to Debug

**Distributed
Checkpoint**

Intertwined Bugs



Heavy Maintenance Effort

But Current Frameworks are **Hard to Use**

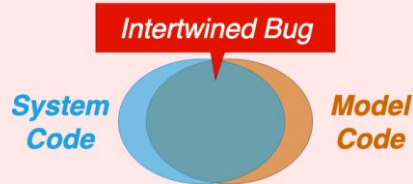
Not PyTorch

Only 8% non-PyTorch



HuggingFace Models

Intertwined System and Model Design



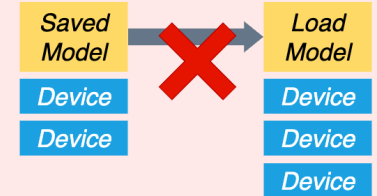
Not Automated Enough



Hard to Debug



No Distributed Checkpoint



A PyTorch-Native Auto-Parallel Framework for *Ease of Use*

PyTorch Native

92% PyTorch



HuggingFace

Decoupled System and Model Design

System Code

Model Code

GradBuffer Defrag
AllReduce Overlap

nn.Linear

No Intertwining

Automatic Parallelism

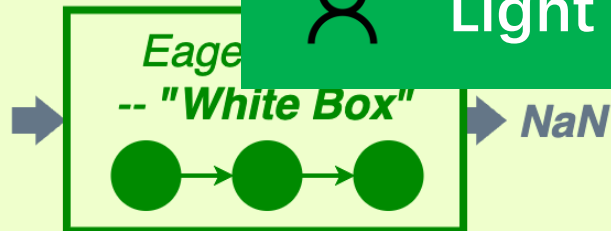
Tensor, Sequence, Data, ZeRO,
Pipeline Parallelism



Minimal Manual Effort

Easy to Debug

Line-by-Line Debug



Easy to Distributed Checkpoint

Load Model

Device

Device

Device

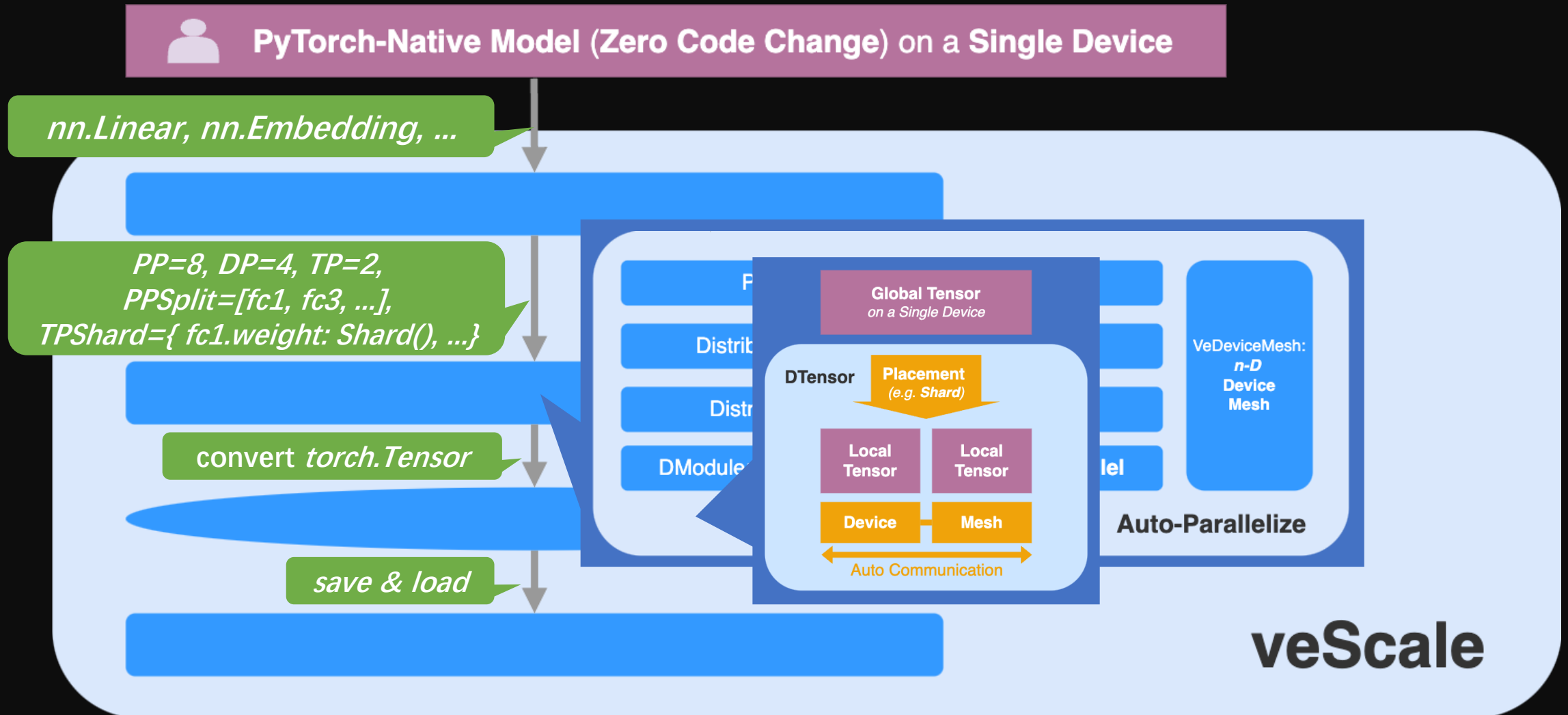
Light Maintenance Effort

Online
Auto
Reshard

Agenda

- Why veScale
- **What is veScale**
- Preliminary Results of veScale
- Future of veScale

What is veScale



Agenda

- Why veScale
- What is veScale
- **Preliminary Results of veScale**
- Future of veScale

Preliminary Results of veScale

Simple API of nD Parallel Training (WIP)

Python ▾

```
1  ### user provides model on single device
2  from internal_model/huggingface.transformers import AutoConfig, AutoModel
3  config = AutoConfig.from_pretrained('/path/to/config')
4  import vescale
5  model = AutoModel.from_config(config)
6
7  ### vescale creates nD parallel plan
8  plan = vescale.generate_plan(model, settings_and_constraints, ...)
9
10 ### vescale creates nD parallel model
11 model, optimizer, ... = vescale.parallelize(plan, model, optimizer_fn, ...)
12
13 ### vescale loads nD parallel model
14 vescale.load("/path", { "plan": plan, "model" : model, "optimizer" : optimizer })
15
16 ### user trains nD parallel model as if on single device
17 for batch in dataloader:
18     loss = model(batch)
19     loss.backward()
20     optimizer.step()
21     optimizer.zero_grad()
22     ...
23
24 ### vescale saves nD parallel model
25 vescale.save("/path", { "plan": plan, "model" : model, "optimizer" : optimizer })
```

Zero Code Change
of Model

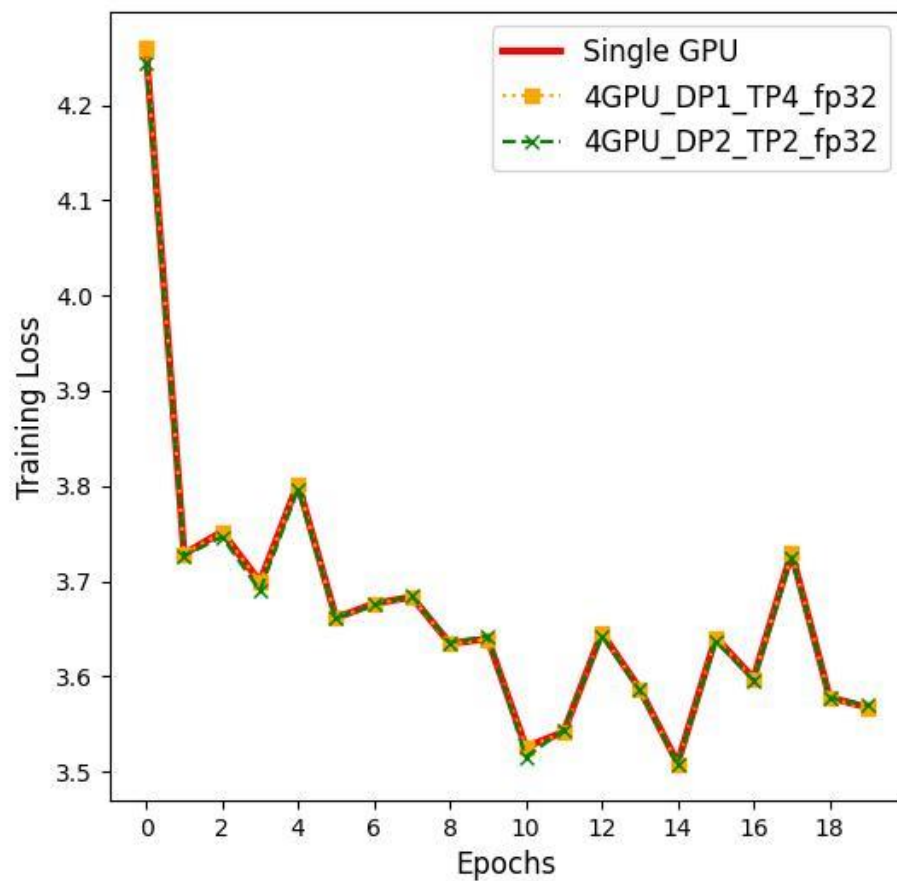
Zero Code Change
of Training Loop

nD Parallel Training
in 5 LoC

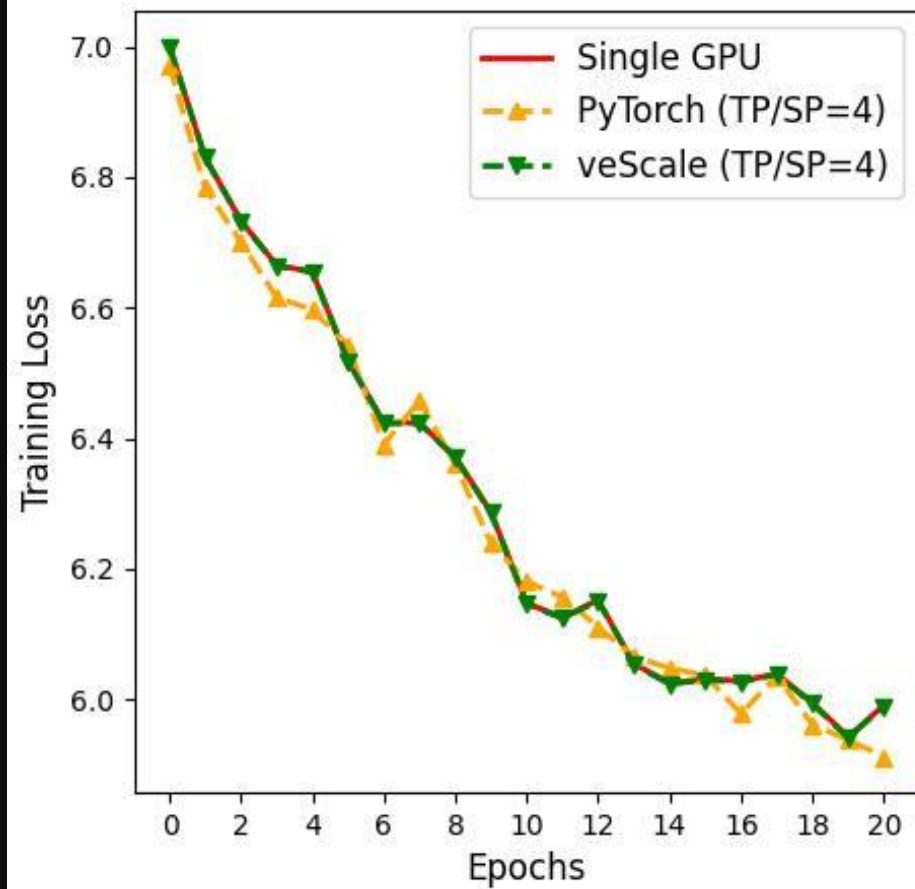
Preliminary Results of veScale

Bitwise Correctness of 4D Parallel Training

nanoGPT



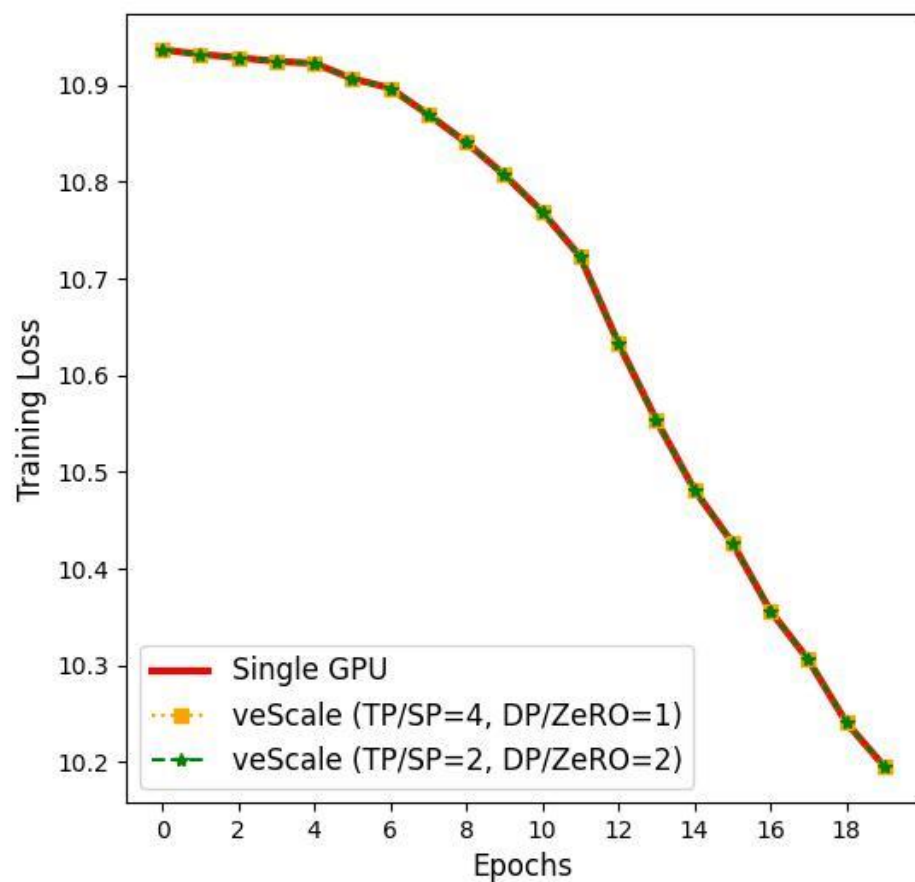
nanoGPT Training



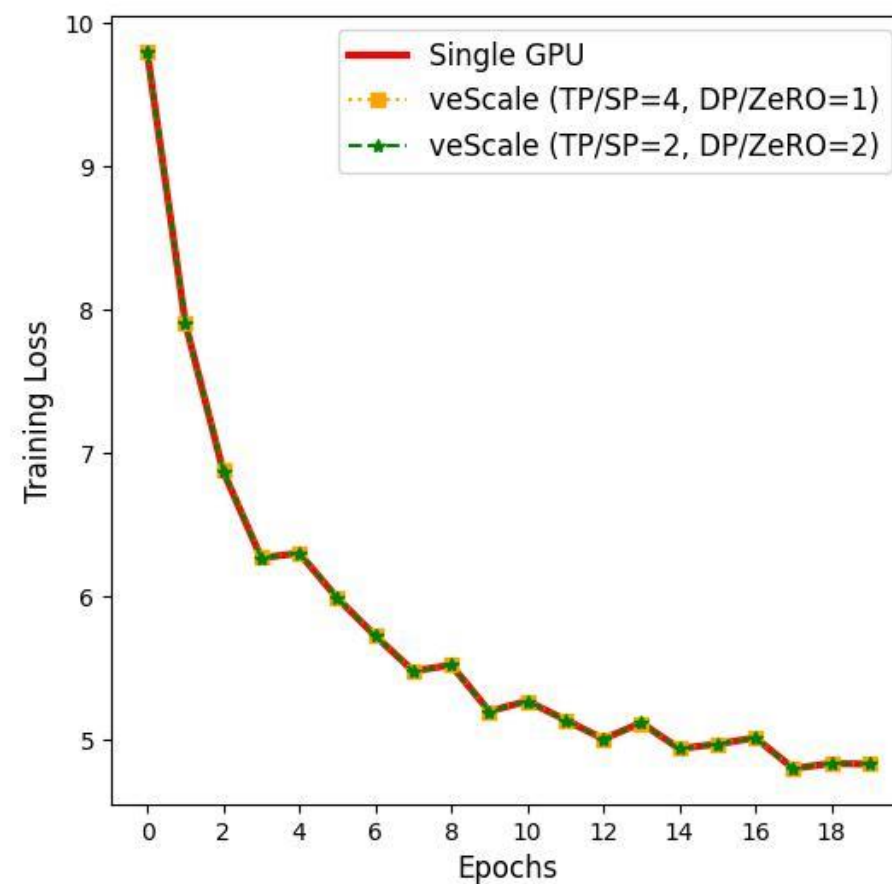
Preliminary Results of veScale

Bitwise Correctness of 4D Parallel Training

Mixtral

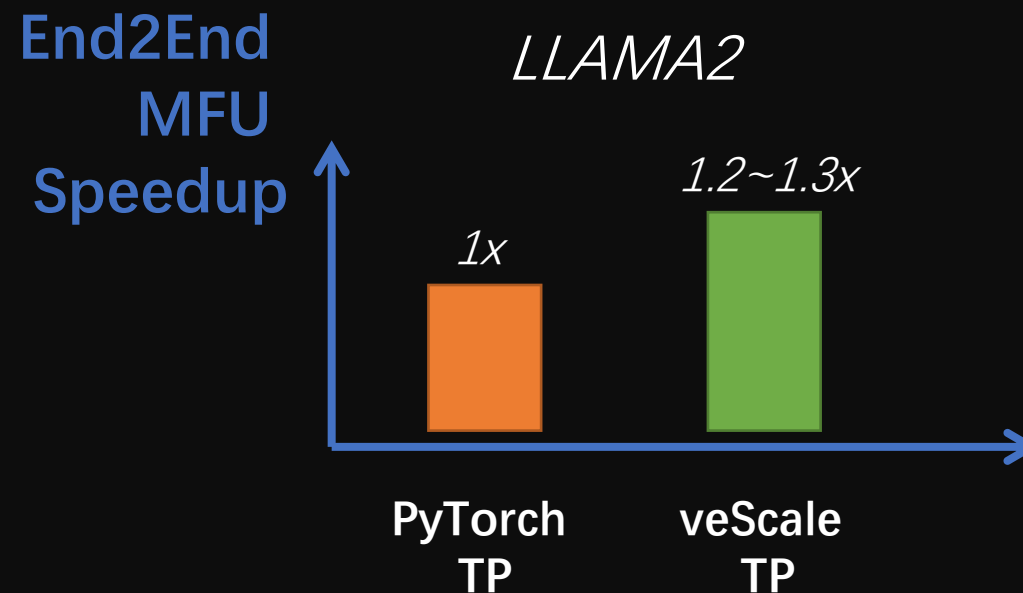
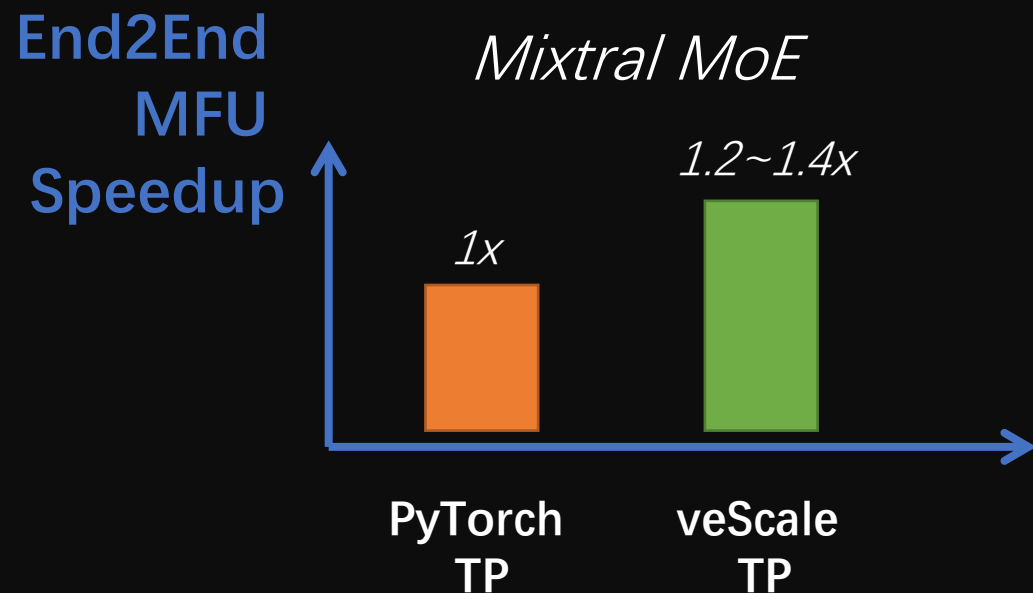


LLama2



Preliminary Results of veScale

Decent Performance of TP (WIP)

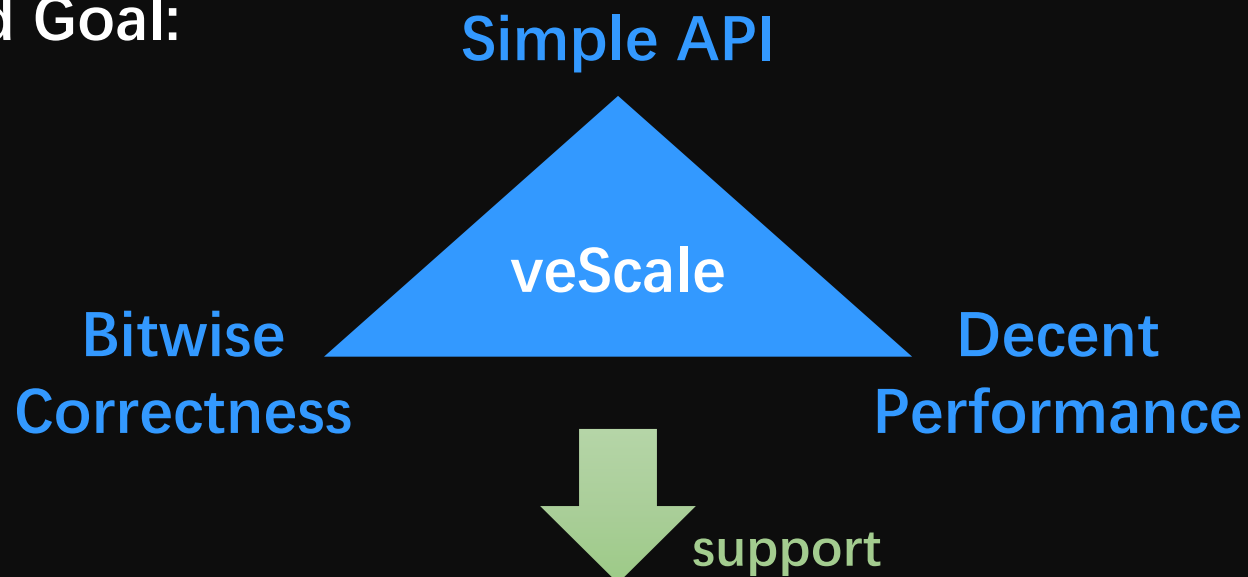


Agenda

- Why veScale
- What is veScale
- Preliminary Results of veScale
- **Future of veScale**

Future of veScale

End Goal:



Company:

100s~1000s

New Models Each Week

Impact!

Open Source Community:

Everyone

"A Promising Work!"

-- AWS AI Lab

-- Octol AI

-- Boson AI

"An Ambitious Work!"

-- Llama Training Lead

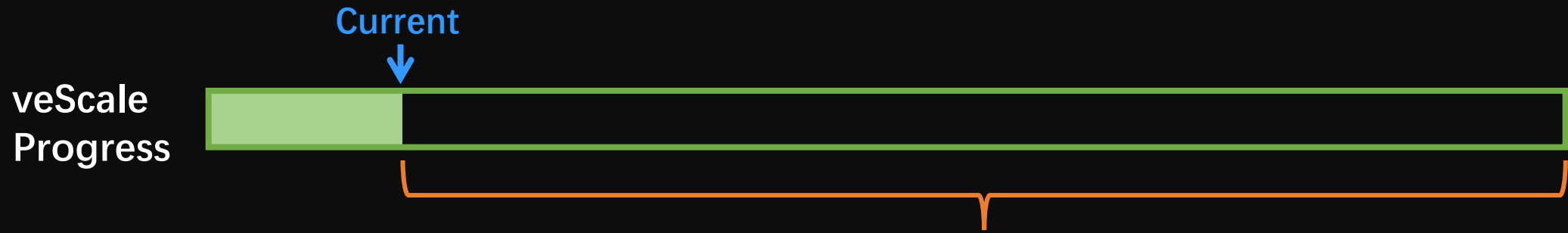
-- PyTorch Training Lead

***"But Many Effort Ahead;
Long-Term Effort Ahead ..."***

-- Llama Training Lead

What's next for veScale

- **Better eager-mode n-D parallelism**
 - Ease of use & performance
- **Better fsdp2**
 - Performance & fsdp2+pp+tp
- **Compile mode for performance**
- **Auto-planner**



Future Challenges

“800” Operator Support for PyTorch

Bitwise Correctness for “nD Parallel”

“Easy to Use” vs “High Performance”

Acknowledgement

(random order)

[Leaders] **Li-wen Chang, Yanghua Peng, Haibin Lin, Xin Liu**

[Contributors] **Xinyi Di, Jiawei Wu, Hongyu Zhu, Ziang Song, Jiacheng Yang, Youjie Li**

[Collaborators] **Minji Han, Chengji Yao, Chenyuan Wang, Yan Xu, Changming Yu, Wenlei Bao, Hao Gong, Ming Zhang, Ningxin Zheng, Xuanrun Zhang**



Open Source for All



vescale.xyz