# OpenBioML PyTorch Lightning workshop

Session 1: Thu 23 Feb, 3pm ET

Session 2: Thu 2 Mar, 3pm ET

https://harvard.zoom.us/j/97375262666

# OpenBioML PyTorch Lightning workshop

## Session 1

- Intro to PyTorch Lightning + Fabric
- Hands-on: raw PyTorch -> Fabric -> PyTorch Lightning Trainer
- A look into OpenFold

By the end of Session 1 we saw how to build a model with PyTorch Lightning and Fabric and train it. Distributed. On a SLURM cluster.

Lightning AI

# OpenBioML PyTorch Lightning workshop

**Session 2**

- Intro to core distributed concepts and what's new in PyTorch 2.0
- Hands-on: how to debug and optimize performance, single node and distributed, running benchmarks
- More on OpenFold

By the end of Session 2 you will know how to make sure you are setting up your training correctly and verify you are leveraging your hardware the best.

Lightning AI

# Join us here

discord     discord.gg/MWAEvnC5fU

forums     lightning.ai/forums

twitter     @LightningAI

https://linktr.ee/lightningai

Lightning AI

# PyTorch Lightning

*You do the science, we do the engineering*

33 million+
DOWNLOADS

780+
CONTRIBUTORS

pip install lightning    21,512

13,000+
PROJECTS USING LIGHTNING

6,000+
SLACK MEMBERS

10,000+ ORGANIZATIONS BUILD WITH LIGHTNING

amazon    NVIDIA    Meta    Microsoft    habana An Intel Company    MARS

## https://lightning.ai

**OpenBioML PyTorch Lightning workshop**

Lightning AI

# Scale your models, without the boilerplate

Lightning's open-source ecosystem is designed for researchers and developers who require flexibility and performance at scale.


pip install lightning

### Build AI without the boilerplate
Lightning simplifies your deep learning code by taking care of engineering boilerplate, so you can focus on the problems that matter to you.

### Unlock deep learning at scale
Work seamlessly with distributed computing environments like multi-GPU and TPU clusters and scale projects to large models and data.

### Create with the community
Join over 100,000 users and companies using Lightning to create their AI future. Tap into cutting-edge research and take it to production.

OpenBioML PyTorch Lightning workshop

⚡ Lightning 2.0 ⚡ Launch Party 🎉

March 15 – Lightning AI Discord

OpenBioML PyTorch Lightning workshop

# Intro: Scaling out and going fast

OpenBioML PyTorch Lightning workshop

Lightning

# TRAIN 1 TRILLION+ PARAMETER MODELS

When training large models, fitting larger batch sizes, or trying to increase throughput using multi-GPU compute, Lightning provides advanced optimized distributed training strategies to support these cases and offer substantial improvements in memory usage.

Note that some of the extreme memory saving configurations will affect the speed of training. This Speed/Memory trade-off in most cases can be adjusted.

Some of these memory-efficient strategies rely on offloading onto other forms of memory, such as CPU RAM or NVMe. This means you can even see memory benefits on a **single GPU**, using a strategy such as DeepSpeed ZeRO Stage 3 Offload.

Check out this amazing video explaining model parallelism and how it works behind the scenes:

## Choosing an Advanced Distributed GPU Strategy

# PyTorch Lightning

## Organized PyTorch

LightningModule

Trainer

```python
import lightning as L
from torch import n, optim

encoder = nn.Sequential(nn.Linear(28 * 28, 64), nn.ReLU(), nn.Linear(64, 3))
decoder = nn.Sequential(nn.Linear(3, 64), n.ReLU() , nn. Linear (64, 28 * 28))

class LitAutoEncoder(L.LightningModule):
    def _init__(self, encoder, decoder):
        super()._init__()
        self.encoder = encoder
        self.decoder = decoder

    def training_step(self, batch, batch_idx):
        X, y = batch
        X = x. view(x.size(0), -1)
        z = self.encoder(x)
        x_hat = self.decoder (z)
        loss = nn. functional.mse_loss(x_hat, x)
        self.log("train_loss", loss) return loss

    def configure_optimizers( self):
        optimizer = optim.Adam(self.parameters(), lr=le-3)
        return optimizer

autoencoder = LitAutoEncoder(encoder, decoder)
dataset = MIST(os.getcwd(), download=True, transform=ToTensor())

train_loader = utils.data.DataLoader (dataset)

trainer = L.Trainer(limit_train_batches=100, max_epochs=1)
trainer.fit(model=autoencoder, train_dataloaders=train_loader)
```

Lightning AI

# PyTorch Lightning

Accelerators

   GPU, TPU, HPU, IPU, MPS

Strategies

   DDP, FSDP, DeepSpeed, Colossal AI

Precision

Callbacks

```python
# train on 4 GPUs
trainer = Trainer(
    devices=4,
    accelerator="gpu"
)


# train 1B+ parameter models with Deepspeed/fsdp
trainer = Trainer(
    devices=4,
    accelerator="gpu",
    strategy="deepspeed_stage_2",
    precision=16
)


# 20+ helpful flags for rapid idea iteration
trainer = Trainer(
    max_epochs=10,
    min_epochs=5,
    overfit_batches=1
)


# access the latest state of the art techniques
trainer = Trainer(callbacks=[StochasticWeightAveraging(...)])
```

OpenBioML PyTorch Lightning workshop

Lightning AI

# Fabric

③ Your code

① Fabric object

② Setup model, optimizer, dataloader

```python
import lightning as L

def train(fabric, model, optimizer, dataloader):
    # Training loop
    model.train()
    for epoch in range(num_epochs):
        for i, batch in enumerate(dataloader):
            ...

def main():
    # (Optional) Parse command line options
    args = parse_args()

    # Configure Fabric
    fabric = L.Fabric(..., strategy="deepspeed")

    # Instantiate objects
    model = ...
    optimizer = ...
    train dataloader = ...

    # Set up objects
    model, optimizer = fabric.setup(model, optimizer)
    train_dataloader = fabric.setup_dataloaders(train_dataloader)

    # Run training loop
    train(fabric, model, optimizer, train_dataloader)

if __name__ == "__main__":
    main()
```

⚡ Lightning AI

# Scaling out

```
fabric = Fabric(accelerator="gpu", strategy="ddp", devices=8, num_nodes=4)
```

```
trainer = Trainer(accelerator="gpu", strategy="ddp", devices=8, num_nodes=4)
```

OpenBioML PyTorch Lightning workshop

Lightning AI

# Scaling out

```
fabric = Fabric(accelerator="gpu", strategy="ddp", devices=8, num_nodes=4)
```

```
trainer = Trainer(accelerator="gpu", strategy="ddp", devices=8, num_nodes=4)
```

Lightning^AI

# Scaling out

```
fabric = Fabric(accelerator="gpu", strategy="ddp", devices=8, num_nodes=4)
```

```
trainer = Trainer(accelerator="gpu", strategy="ddp", devices=8, num_nodes=4)
```

**OpenBioML PyTorch Lightning workshop**

Lightning<sup>AI</sup>

# Scaling out

```
fabric = Fabric(accelerator="gpu", strategy="ddp", devices=8, num_nodes=4)
```

```
trainer = Trainer(accelerator="gpu", strategy="ddp", devices=8, num_nodes=4)
```

Think one process per accelerator

OpenBioML PyTorch Lightning workshop

# Under the hood

```
MASTER_PORT: free port on machine with NODE_RANK 0
MASTER_ADDR: address of NODE_RANK 0 node
WORLD_SIZE: the total number of GPUs/processes
NODE_RANK: id of the node in the cluster
```

```
MASTER_PORT=23456
MASTER_ADDR="10.10.10.25"
WORLD_SIZE=32
```

GPU

NODE_RANK=0
10.10.10.25

NODE_RANK=1
10.10.10.26

NODE_RANK=2
10.10.10.27

NODE_RANK=3
10.10.10.28

OpenBioML PyTorch Lightning workshop

Lightning AI

# Under the hood

```
WORLD_SIZE: the total number of GPUs/processes
NODE_RANK: id of the node
LOCAL_RANK: id of the process in each node
GLOBAL_RANK: unique id of the process
```

```
MASTER_PORT=23456
MASTER_ADDR="10.10.10.25"
WORLD_SIZE=32
```

LOCAL_RANK



NODE_RANK=0
10.10.10.25

NODE_RANK=1
10.10.10.26

NODE_RANK=2
10.10.10.27

NODE_RANK=3
10.10.10.28

**OpenBioML PyTorch Lightning workshop**

# Under the hood

WORLD_SIZE: the total number of GPUs/processes
NODE_RANK: id of the node
LOCAL_RANK: id of the process in each node
GLOBAL_RANK: unique id of the process

```
MASTER_PORT=23456
MASTER_ADDR="10.10.10.25"
WORLD_SIZE=32
```

GLOBAL_RANK

| 0 | 1 |
|---|---|
| 2 | 3 |
| 4 | 5 |
| 6 | 7 |

NODE_RANK=0
10.10.10.25

| 8 | 9 |
|---|---|
| 10 | 11 |
| 12 | 13 |
| 14 | 15 |

NODE_RANK=1
10.10.10.26

| 16 | 17 |
|---|---|
| 18 | 19 |
| 20 | 21 |
| 22 | 23 |

NODE_RANK=2
10.10.10.27

| 24 | 25 |
|---|---|
| 26 | 27 |
| 28 | 29 |
| 30 | 31 |

NODE_RANK=3
10.10.10.28

**OpenBioML PyTorch Lightning workshop**

Lightning<sup>AI</sup>

# In Lightning / Fabric

```python
from lightning.fabric import Fabric

# Devices and num nodes determine how many processes there are
fabric = Fabric(devices=8, num_nodes=4)
fabric.launch()
```

```python
# The total number of processes running across all devices and nodes
fabric.world_ size  # 4 * 8 = 32

# The global index of the current process across all devices and nodes
fabric.global_rank  # -> {0, 1, 2, 3, 4, ..., 31}

# The index of the current process among the processes running on the
local node
fabric.local_rank   # -> {0, 1, 2, 3, 4, 5, 6, 7}

# The index of the current node
fabric.node_rank    # -> {0, 1, 2, 3}

# Do something only on rank 0
if fabric.global_rank == 0:
    ...
```

OpenBioML PyTorch Lightning workshop

Lightning AI

# Scaling out

data

model

GPU

Lightning AI

# Scaling out



data

batch

model

GPU

# Scaling out



data

batch

model

GPU

parameters

activations

gradients

optimizer state

memory

Lightning AI

# Gradient accumulation

data

μ-batches

GPU

model

compute backward each μ-batch, but call
optimizer.step() for the whole batch

parameters

activations

gradients

optimizer state

memory

# Activation/Gradient checkpointing



data

batch

model

GPU

avoid storing activations during forward, only save (checkpoint) a few of them in-between to avoid recomputing too much

parameters

activations

gradients

optimizer state

memory

OpenBioML PyTorch Lightning workshop

# Mixed precision



data

batch

model

GPU

reduce numerical precision (32 -> 16-bit) so that
memory footprint is reduced; normally full precision is
automatically restored at computation time

parameters

activations

gradients

optimizer state

memory

Lightning AI

# Mixed precision



float 32

Sign (1 bit)  Exponent (8 bits)  Fraction (23 bits)

float 16 ("half" precision)

Sign (1 bit)  Exponent (5 bits)  Fraction (10 bits)

bfloat 16 ("brain" floating point, more "dynamic range" like float 32)

Sign (1 bit)  Exponent (8 bits)  Fraction (7 bits)

**Overflow and Underflow**

```
torch.tensor(10**6, dtype=torch.float32)
```
tensor(1000000.)

```
torch.tensor(10**6, dtype=torch.float16)
```
tensor(inf, dtype=torch.float16)

**During training:**

FP32 weights → FP16 weights → FP16 gradients → FP32 gradients → Optimizer → FP32 weights

Sebastian Rashka, https://lightning.ai/pages/courses/deep-learning-fundamentals/

Lightning AI

# Offloading



data

batch

model

GPU

GPU memory

CPU memory

move portions of parameters, gradients, optimizer state and activations to CPU when they are not involved in active computations

# Scaling out



data

model

GPU

Lightning AI

# Scaling out



data



GPU

Lightning AI

# Data parallelism



data

batch

model

GPU

batch

model

GPU

Different processes (devices on the same node, or devices on separate nodes)

# Data parallelism



data

batch

model

GPU

batch

model

GPU

Separate instances of DistributedSampler load distinct batches from within different processes

Gradients need to be synchronized across processes during the backward pass

Lightning AI

# Data parallelism



batch

data

batch

Separate instances of DistributedSampler load distinct batches from within different processes

**Process 0**

param0 grad0
param1 grad1
param2 grad2
param3 grad3

bucket1
bucket0

allreduce

**Process 1**

param0 grad0
param1 grad1
param2 grad2
param3 grad3

bucket1
bucket0

allreduce

model

GPU

https://pytorch.org/docs/master/notes/ddp.html

Lightning AI

# Data parallelism



data

batch

model

GPU

batch

model

GPU

Model still needs to fit in device memory

# Sharded Data parallelism



batch

data

batch

model

GPU

model

GPU

Each node keeps a portion of the parameters, etc in memory, gathers the weights needed to perform a computation and deallocates

Lightning AI

# Model parallelism



data

batch

model

model

GPU

GPU

Distribute model parameters across machines in contiguous blocks

Lightning AI

# Model parallelism



data

batch

model

GPU A

model

GPU B

**OpenBioML PyTorch Lightning workshop**

# Model parallelism



data

batch

GPU B

GPU A

forward

backward

time

Lightning<sup>AI</sup>

# Model parallelism



data

batch

GPU B

GPU A

forward

backward

time

Bubbles! Bubbles!
Bubbles! Bubbles!

Lightning AI

# Naive model parallelism



https://siboehm.com/articles/22/pipeline-parallel-training

OpenBioML PyTorch Lightning workshop

# Pipeline parallelism



data

μ-batches

μ-batch

GPU B

GPU A

forward

backward

time

Gradient accumulation across μ-batches

**OpenBioML PyTorch Lightning workshop**

Lightning<sup>AI</sup>

# Interleaved pipeline parallelism



data

μ-batches

μ-batch

GPU B

GPU A

forward          backward

time

Compute backward as soon as possible
Bubble gets better as more nodes participate

**OpenBioML PyTorch Lightning workshop**

Lightning<sup>AI</sup>

# Tensor parallelism



data

batch

model

GPU

model

GPU

Each node keeps a portion of the weights in memory, performs ops column- or row-wise

# Tensor parallelism

https://sebastianraschka.com/blog/2023/pytorch-faster.html



**Column-wise**

GPU 1

GPU 2

**Row-wise**

OpenBioML PyTorch Lightning workshop

# Combining parallelism: 2D, 3D, 4D

E.g. Sharded Data Parallel + Tensor Parallel in separate parallel dimensions:

- Data Parallel across hosts
- Tensor Parallel within each host

Other example: Megatron-LM



FullyShardedDataParallel (FSDP)

| Host 1 | Host 2 | … | Host N |
|---|---|---|---|
| 8 GPUs | 8 GPUs | | 8 GPUs |
| (TP) | (TP) | | (TP) |
| [0, 1, 2, …, 7] | [8, 9, …, 15] | | [8*N - 8, …, 8*N - 1] |

[0, 8, …, 8*N - 8] [1, 9, …, 8*N - 7] [2, 10, …, 8*N - 6] … [7, 15, …, 8*N - 1]

Wanchao Liang, Two Dimensional Parallelism Using Distributed Tensors

OpenBioML PyTorch Lightning workshop

# DDP

DDP - sharded DDP that shards **model parameters**, optimizer state and gradients across DDP ranks. It can optionally offload to CPU.

```
trainer = Trainer(accelerator="gpu", strategy="ddp", devices=8, num_nodes=4)
```

# DeepSpeed

- DeepSpeed ZeRO Stage 1 - Shard **optimizer states**, remains at speed parity with DDP whilst providing memory improvement

- DeepSpeed ZeRO Stage 2 - Shard **optimizer states** and **gradients**, remains at speed parity with DDP whilst providing even more memory improvement

- DeepSpeed ZeRO Stage 2 Offload - **Offload optimizer states** and **gradients** to CPU. Increases distributed communication volume and GPU-CPU device transfer, but provides significant memory improvement

- DeepSpeed ZeRO Stage 3 - **Shard optimizer states**, **gradients**, **parameters** and optionally **activations**. Increases distributed communication volume, but provides even more memory improvement

- DeepSpeed ZeRO Stage 3 Offload - **Offload optimizer states**, **gradients**, **parameters** and optionally **activations** to CPU. Increases distributed communication volume and GPU-CPU device transfer, but even more significant memory improvement.

- DeepSpeed Activation Checkpointing - **Free activations after forward pass**. Increases computation, but provides memory improvement for all stages.

# DeepSpeed

```python
trainer = Trainer(accelerator="gpu", devices=8, num_nodes=4, strategy="deepspeed_stage_1")
trainer = Trainer(accelerator="gpu", devices=8, num_nodes=4, strategy="deepspeed_stage_2")
trainer = Trainer(accelerator="gpu", devices=8, num_nodes=4, strategy="deepspeed_stage_2_offload")
trainer = Trainer(accelerator="gpu", devices=8, num_nodes=4, strategy="deepspeed_stage_3")
trainer = Trainer(accelerator="gpu", devices=8, num_nodes=4, strategy="deepspeed_stage_3_offload")
```

OpenBioML PyTorch Lightning workshop

Lightning^AI

# DeepSpeed

```python
import lightning as L
from lightning.pytorch.strategies import DeepSpeedStrategy
import deepspeed


class MyModel(L.LightningModule):
    ...

    def configure_sharded_model(self):
        self.block_1 = nn.Sequential(nn.Linear(32, 32),
nn.ReLU())
        self.block_2 = torch.nn.Linear(32, 2)

    def forward(self, x):
        # Use the DeepSpeed checkpointing function instead of
calling the module directly
        x = deepspeed.checkpointing.checkpoint(self.block_1, x)
        return self.block_2(x)
```

```python
model = MyModel()

trainer = L.Trainer(accelerator="gpu", devices=4,
strategy="deepspeed_stage_3_offload", precision=16)

# Enable CPU Activation Checkpointing
trainer = Trainer(
    accelerator="gpu",
    devices=4,
    strategy=DeepSpeedStrategy(
        stage=3,
        offload_optimizer=True,   # Enable CPU Offloading
        cpu_checkpointing=True,   # Offload activations to CPU
    ),
    precision=16,
)

trainer.fit(model)
```

**OpenBioML PyTorch Lightning workshop**

Lightning AI

# FSDP (Fully Sharded Data Parallel)

FSDP - sharded DDP that **shards model parameters**, **optimizer state** and **gradients** across DDP ranks.

It optionally **offloads activations** and **optimizer state** to CPU.

```python
trainer = Trainer(accelerator="gpu", strategy="fsdp", devices=8, num_nodes=4)
```

https://engineering.fb.com/2021/07/15/open-source/fsdp/

Lightning^AI

# ColossalAI

ColossalAI - Zero-DP with dynamic chunk-based memory management and other configurable parallelization strategies.

```python
class MyModel(LightningModule):
    def __init__(self):
        super().__init__()
        # don't instantiate layers here
        # move the creation of layers to
`configure_sharded_model`

    def configure_sharded_model(self):
        # create all your layers here
        self.layers = nn.Sequential(...)
```

```python
from lightning_colossalai import ColossalAIStrategy

model = MyModel()
my_strategy = ColossalAIStrategy(placement_policy="auto")

trainer = Trainer(accelerator="gpu", devices=4, precision=16,
strategy=my_strategy)

trainer.fit(model)
```

Monitors the consumption of CUDA memory during the warmup phase and collects CUDA memory usage of all auto-grad operations

Automatically manages the data transmission between GPU and CPU according to collected CUDA memory usage information

Lightning^AI

# Going fast

Eager mode

```python
def forward(self, x):
    B, T, C = x.size()                                              ────────────→  execute
    q, k ,v  = self.c_attn(x).split(self.n_embd, dim=2)             ────────────→  execute
    k = k.view(B, T, self.n_head, C // self.n_head).transpose(1, 2) ───────────→  execute
    q = q.view(B, T, self.n_head, C // self.n_head).transpose(1, 2) ───────────→  execute
    v = v.view(B, T, self.n_head, C // self.n_head).transpose(1, 2) ───────────→  ...

    att = (q @ k.transpose(-2, -1)) * (1.0 / math.sqrt(k.size(-1)))
    att = att.masked_fill(self.bias[:,:,:T,:T] == 0, float('-inf'))
    att = F.softmax(att, dim=-1)
    att = self.attn_dropout(att)

    y = att @ v
    y = y.transpose(1, 2).contiguous().view(B, T, C)
    y = self.resid_dropout(self.c_proj(y))
```

GPU

2: compute

1: fetch data →

← 3: store result

**OpenBioML PyTorch Lightning workshop**

Lightning<sup>AI</sup>

# Going fast: bottleneck

Eager mode

```python
def forward(self, x):
    B, T, C = x.size()
    q, k ,v  = self.c_attn(x).split(self.n_embd, dim=2)
    k = k.view(B, T, self.n_head, C // self.n_head).transpose(1, 2)
    q = q.view(B, T, self.n_head, C // self.n_head).transpose(1, 2)
    v = v.view(B, T, self.n_head, C // self.n_head).transpose(1, 2)

    att = (q @ k.transpose(-2, -1)) * (1.0 / math.sqrt(k.size(-1)))
    att = att.masked_fill(self.bias[:,:,:T,:T] == 0, float('-inf'))
    att = F.softmax(att, dim=-1)
    att = self.attn_dropout(att)

    y = att @ v
    y = y.transpose(1, 2).contiguous().view(B, T, C)
    y = self.resid_dropout(self.c_proj(y))
```

execute

execute

execute

execute

...

GPU

2: compute

1: fetch data

3: store result

Bottleneck on modern hardware

**OpenBioML PyTorch Lightning workshop**

Lightning^AI

# Going fast: bottleneck



```python
def forward(self, x):
    B, T, C = x.size()
    q, k ,v  = self.c_attn(x).split(self.n_embd, dim=2)
    k = k.view(B, T, self.n_head, C // self.n_head).transpose(1, 2)
    q = q.view(B, T, self.n_head, C // self.n_head).transpose(1, 2)
    v = v.view(B, T, self.n_head, C // self.n_head).transpose(1, 2)

    att = (q @ k.transpose(-2, -1)) * (1.0 / math.sqrt(k.size(-1)))
    att = att.masked_fill(self.bias[:,:,:T,:T] == 0, float('-inf'))
    att = F.softmax(att, dim=-1)
    att = self.attn_dropout(att)

    y = att @ v
    y = y.transpose(1, 2).contiguous().view(B, T, C)
    y = self.resid_dropout(self.c_proj(y))
```

Eager mode

execute
execute
execute
execute
...

5x

1.5x

data fetch    compute

GPU

2: compute

1: fetch data

3: store result

Bottleneck on modern hardware

# Going fast: compiled mode

```python
def forward(self, x):
    B, T, C = x.size()
    q, k ,v  = self.c_attn(x).split(self.n_embd, dim=2)
    k = k.view(B, T, self.n_head, C // self.n_head).transpose(1, 2)
    q = q.view(B, T, self.n_head, C // self.n_head).transpose(1, 2)
    v = v.view(B, T, self.n_head, C // self.n_head).transpose(1, 2)

    att = (q @ k.transpose(-2, -1)) * (1.0 / math.sqrt(k.size(-1)))
    att = att.masked_fill(self.bias[:,:,:T,:T] == 0, float('-inf'))
    att = F.softmax(att, dim=-1)
    att = self.attn_dropout(att)

    y = att @ v
    y = y.transpose(1, 2).contiguous().view(B, T, C)
    y = self.resid_dropout(self.c_proj(y))
```

Compiled mode

acquire → intermediate representation → fuse/optimize → execute fused block

GPU

2: compute

1: fetch data

3: store result

Infrequent memory access

Uninterrupted compute

OpenBioML PyTorch Lightning workshop

Lightning^AI

# Going fast: PyTorch 2.0

```
model = NanoGPT(config)
model = torch.compile(model)
model(x)

# Works with PyTorch Lightning (+ Fabric)
model = MyLitModule()
model = torch.compile(model)
trainer.fit(model)
```

**dynamo**

Program acquisition
through tracing
bytecode, optimization

**inductor**

Optimized execution,
based on OpenAI Triton

Lightning™

# Going fast: PyTorch 2.0

```
model = NanoGPT(config)
model = torch.compile(model)
model(x)

# Works with PyTorch Lightning (+ Fabric)
model = MyLitModule()
model = torch.compile(model)
trainer.fit(model)
```

dynamo

Program acquisition
through tracing
bytecode, optimization

inductor

Optimized execution,
based on OpenAI Triton

Status: **DDP supported, FSDP support in the works**

Lightning<sup>AI</sup>

# Demo: Model Parallel

**github.com/Lightning-AI/open-bio-ml-workshop**

Lightning<sup>AI</sup>

# Performance Optimization

Lightning AI

# Minimize framework overhead

**Compile:** Use *torch.compile* whenever possible

**Logging:** Log often for development, log less for long training runs

**Checkpointing:** Minimize frequency, checkpoint often only if training is unstable

**Validation:** Tune frequency based on dataset sizes

**OpenBioML PyTorch Lightning workshop**

# Maximize throughput

**Compile:** Use *torch.compile* whenever possible

**Mixed precision:** Speed + memory, prefer bfloat16 if available

**Batch Size:** Increase until OOM, avoid Malloc retries

**Num Workers:** Increase if GPU is waiting on data, consumes more CPU memory

OpenBioML PyTorch Lightning workshop

Lightning AI

# Best practices

Avoid unnecessary GPU **synchronization**

```
# Avoid these:
output.item()
output.numpy()
output.cpu()

torch.cuda.empty_cache()
```

Create tensors directly on the **device**

```
# bad
t = torch.rand(4, 4).cuda()

# LightningModule:
torch.rand(4, 4, device=self.device)

# Fabric:
torch.rand(4, 4, device=fabric.device)
```

**OpenBioML PyTorch Lightning workshop**

Lightning AI

# Performance flags

### Speed (default)

```
# PyTorch / Fabric
torch.use_deterministic_algorithms(False)
torch.backends.cudnn.benchmark = True

# Trainer
trainer = Trainer(benchmark=True, deterministic=False)
```

### Determinism

```
# PyTorch / Fabric
torch.use_deterministic_algorithms(True)
torch.backends.cudnn.benchmark = False

# Trainer
trainer = Trainer(benchmark=False, deterministic=True)
```

OpenBioML PyTorch Lightning workshop

Lightning

# Comparing Trainer vs. PyTorch

**Run "barebones"**

```python
# Disable logging, checkpointing, etc.
trainer = Trainer(..., barebones=True)
```

Recommended for comparing

implementations and **unit testing**!

OpenBioML PyTorch Lightning workshop

Lightning^AI

# Tensor cores

Perform **costly matrix multiplications** in lower precision (internal)

```python
# Default
torch.set_float32_matmul_precision("high")

# Lower precision matrix multiplication
torch.set_float32_matmul_precision("highest")
torch.set_float32_matmul_precision("medium")
```

**Lightning** informs you if you have tensor cores

```
WARNING: You are using a CUDA device ('NVIDIA GeForce RTX 3090') that
has Tensor Cores. To properly utilize them, you should set
`torch.set_float32_matmul_precision('medium' | 'high')` which will
trade-off precision for performance.
```

OpenBioML PyTorch Lightning workshop

Lightning AI

# Find bottlenecks (framework)

**Configure profiler**

```python
# Trainer
trainer = Trainer(profiler="simple", ...)
trainer.fit(...)
```

**Output after fit**

```
FIT Profiler Report

-----------------------------------------------------------------------------------------------
|  Action                                        |  Mean duration (s)  |  Total time (s)  |
-----------------------------------------------------------------------------------------------
|    [LightningModule]BoringModel.prepare_data    |  10.0001            |  20.00           |
|    run_training_epoch                           |  6.1558             |  6.1558          |
|    run_training_batch                           |  0.0022506          |  0.015754        |
|    [LightningModule]BoringModel.optimizer_step  |  0.0017477          |  0.012234        |
|    [LightningModule]BoringModel.val_dataloader  |  0.00024388         |  0.00024388      |
|    on_train_batch_start                         |  0.00014637         |  0.0010246       |
|    [LightningModule]BoringModel.teardown        |  2.15e-06           |  2.15e-06        |
|    [LightningModule]BoringModel.on_train_start  |  1.644e-06          |  1.644e-06       |
|    [LightningModule]BoringModel.on_train_end    |  1.516e-06          |  1.516e-06       |
|    [LightningModule]BoringModel.on_fit_end      |  1.426e-06          |  1.426e-06       |
|    [LightningModule]BoringModel.setup           |  1.403e-06          |  1.403e-06       |
|    [LightningModule]BoringModel.on_fit_start    |  1.226e-06          |  1.226e-06       |
-----------------------------------------------------------------------------------------------
```

Lightning^AI

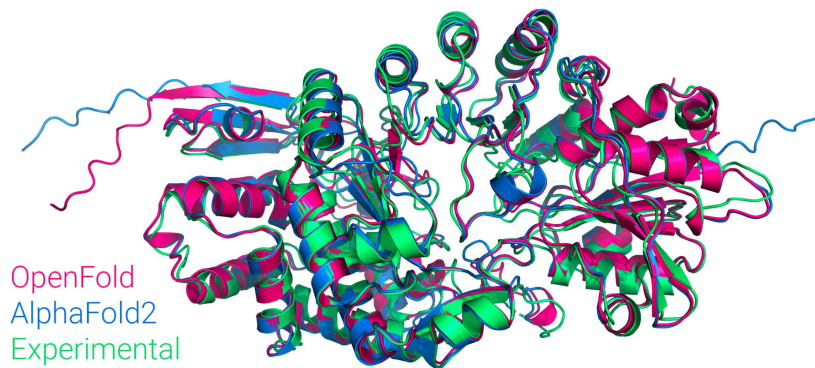# Find bottlenecks (PyTorch)

**Configure profiler**

```python
# Trainer
trainer = Trainer(profiler="pytorch", ...)
trainer.fit(...)
```

**Output after fit**

```
Profiler Report
Profile stats for: records
---------------------------------------------------  ------------  ------------  ------------  ------------  ------------  ------------  ------------  ------------  ---
                                              Name    Self CPU %     Self CPU   CPU total %    CPU total   CPU time avg    Self CUDA   Self CUDA %   CUDA total    CUD
---------------------------------------------------  ------------  ------------  ------------  ------------  ------------  ------------  ------------  ------------  ---
                                     ProfilerStep*         3.48%     47.673ms        39.44%    540.549ms    180.183ms      0.000us         0.00%    113.315ms
         [pl][profile][Strategy]SingleDeviceStrategy.backward...   32.14%    440.445ms        32.18%    441.004ms    147.001ms      0.000us         0.00%      3.000us
autograd::engine::evaluate_function: EmbeddingBackwa...    0.00%     57.000us        27.61%    378.372ms     63.062ms      0.000us         0.00%      4.470us
                                  EmbeddingBackward0         0.00%     21.000us        27.60%    378.277ms     63.046ms      0.000us         0.00%      2.834us
                           aten::embedding_backward         0.00%     15.000us        27.60%    378.256ms     63.043ms      0.000us         0.00%      2.834us
                     aten::embedding_dense_backward         0.03%    360.000us        27.60%    378.241ms     63.040ms    659.000us         0.11%      2.834us
                               cudaStreamSynchronize        27.51%    377.063ms        27.51%    377.063ms     41.896ms      0.000us         0.00%      0.000us
                  [pl][profile]run_training_batch         0.05%    710.000us        26.02%    356.661ms    178.331ms      0.000us         0.00%    117.272ms
       [pl][profile][LightningModule]LitGPT.optimizer_step     0.00%     51.000us        25.97%    355.951ms    177.976ms      0.000us         0.00%    117.272ms
                           Optimizer.step#AdamW.step        23.54%    322.614ms        25.97%    355.900ms    177.950ms      0.000us         0.00%    117.272ms
        [pl][profile][Strategy]SingleDeviceStrategy.training...    0.02%    224.000us         3.64%     49.858ms     16.619ms      0.000us         0.00%    175.918ms
                         [pl][module]gpt.GPT: gpt         0.06%    773.000us         3.62%     49.634ms     16.545ms      0.000us         0.00%    175.918ms
                               cudaDeviceSynchronize         3.26%     44.698ms         3.26%     44.698ms     44.698ms      0.000us         0.00%      0.000us
                                   cudaLaunchKernel         2.01%     27.556ms         2.01%     27.556ms      5.968us     14.669ms         2.55%     14.669ms
                                      aten::linear         0.11%      1.443ms         1.65%     22.615ms     76.922us      0.000us         0.00%    150.840ms
                                          aten::to         0.12%      1.584ms         1.33%     18.255ms     14.318us      0.000us         0.00%     56.977ms
                                     aten::_to_copy         0.27%      3.708ms         1.26%     17.297ms     17.850us      0.000us         0.00%     58.926ms
                                          aten::mm         0.58%      7.996ms         0.85%     11.625ms     26.361us    209.986ms        36.56%    211.492ms
      autograd::engine::evaluate_function: MmBackward0     0.10%      1.334ms         0.82%     11.205ms     76.224us      0.000us         0.00%    138.901ms
                                        MmBackward0         0.10%      1.370ms         0.72%      9.871ms     67.150us      0.000us         0.00%    138.901ms
---------------------------------------------------  ------------  ------------  ------------  ------------  ------------  ------------  ------------  ------------  ---
Self CPU time total: 1.371s
Self CUDA time total: 574.402ms
```

Lightning

OpenFold
AlphaFold2
Experimental

# OpenFold Update: Lightning 2.0!

OpenBioML PyTorch Lightning workshop

Lightning

# Reach out!

discord     discord.gg/MWAEvnC5fU

forums     lightning.ai/forums

twitter     @LightningAI

https://linktr.ee/lightningai

Lightning<sup>AI</sup>

# Thanks

Adrian Wälchli
Akihiro Nitta
Carlos Mocholí
Eden Afek
Ethan Harris
Jirka Borovec
Justus Schock
Thomas Chaton
William Falcon

💜 community

＋

Lightning LEAGUE

https://linktr.ee/lightningai

Lightning AI

# Thanks

OpenBioML PyTorch Lightning workshop