

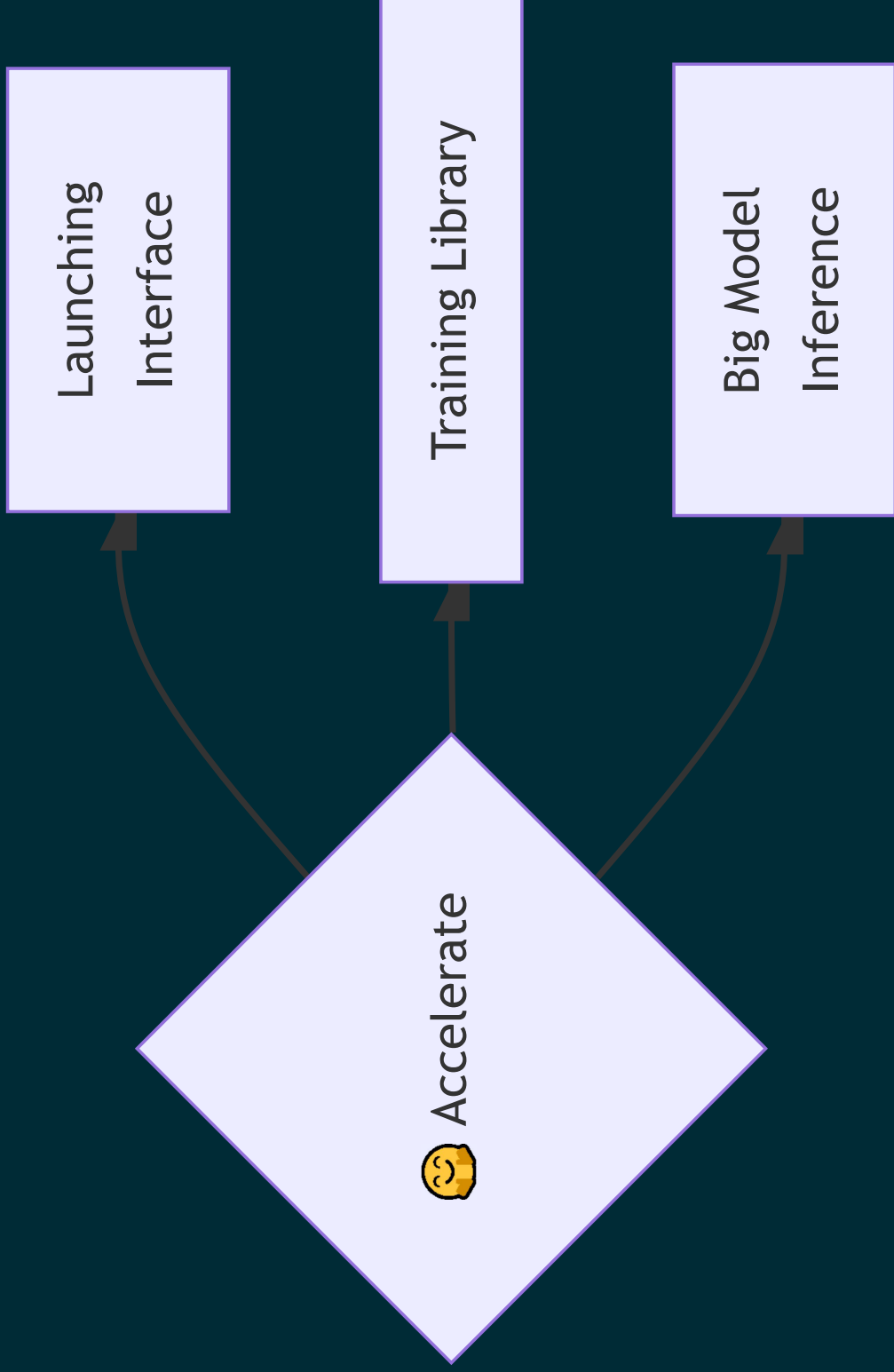
ACCELERATE, THREE POWERFUL SUBLIBRARIES FOR PYTORCH

Zachary Mueller

WHO AM I?

- Zachary Mueller
- Deep Learning Software Engineer at 
- API design geek

WHAT IS 🤖 ACCELERATE?



A LAUNCHING INTERFACE

Can't I just use `python do_the_thing.py`?

A LAUNCHING INTERFACE

Launching scripts in different environments is complicated:

- ```
1 python script.py
```

- ```
1 torchrun --nnodes=1 --nproc_per_node=2 script.py
```

- ```
1 deepspeed --num_gpus=2 script.py
```

And more!

# A LAUNCHING INTERFACE

But it doesn't have to be:

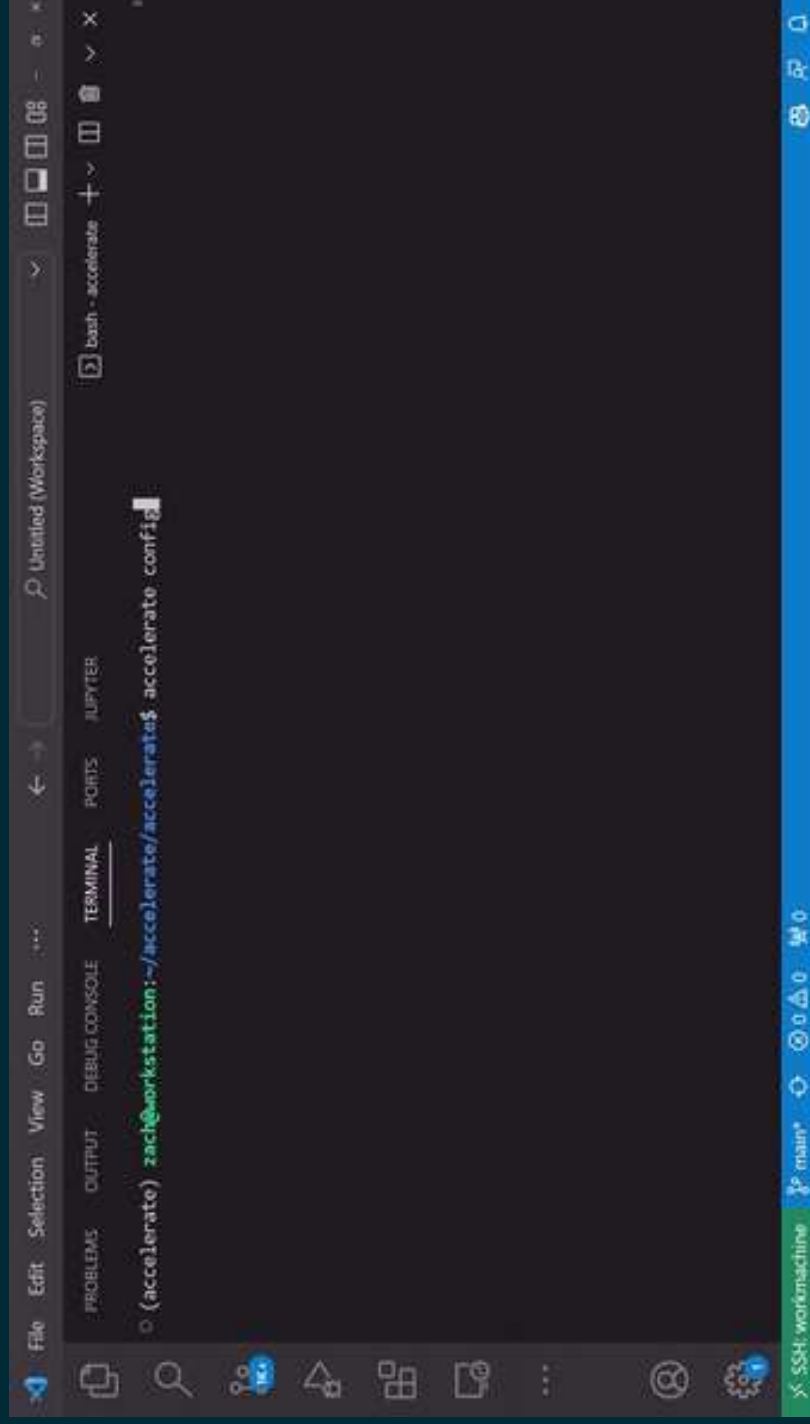
```
1 accelerate launch script.py
```

A single command to launch with **DeepSpeed**, Fully Sharded Data Parallelism, across single and multi CPUs and GPUs, and to train on TPUs<sup>1</sup> too!

1. Without needing to modify your code and create a `mp_fn`

# A LAUNCHING INTERFACE

Generate a device-specific configuration through `accelerate config`



```
bash - accelerate
(Unitled Workspace)
PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS JUPYTER
(accelerate) zach@workstation:~/accelerate/accelerate$ accelerate config
```

# A LAUNCHING INTERFACE

Or don't. **accelerate** **config** doesn't *have* to be done!

```
1 torchrun --nnodes=1 --nproc_per_node=2 script.py
2 accelerate launch --multi_gpu --nproc_per_node=2 script.py
```

A quick default configuration can be made too:

```
1 accelerate config default
```



# A LAUNCHING INTERFACE

With the `notebook_launcher` it's also possible to launch code directly from your Jupyter environment too!

```
1 from accelerate import notebook_launcher
2 notebook_launcher(
3 training_loop_function,
4 args,
5 num_processes=2
6)
```

```
1 Launching training on 2 GPUs.
2 epoch 0: 88.12
3 epoch 1: 91.73
4 epoch 2: 92.58
5 epoch 3: 93.90
6 epoch 4: 94.71
```

# A TRAINING LIBRARY

Okay, will `accelerate` `launch` make `do_the_thing.py`  
use all my GPUs magically?

# A TRAINING LIBRARY

- Just showed that its possible using **accelerate** **launch** to *launch* a python script in various distributed environments
- This does *not* mean that the script will just “use” that code and still run on the new compute efficiently.
- Training on different computes often means *many* lines of code changed for each specific compute.
- 🤖 **accelerate** solves this by ensuring the same code can be ran on a CPU or GPU, multiples, and on TPUs!

# A TRAINING LIBRARY

```
1 for batch in dataloader:
2 optimizer.zero_grad()
3 inputs, targets = batch
4 inputs = inputs.to(device)
5 targets = targets.to(device)
6 outputs = model(inputs)
7 loss = loss_function(outputs, targets)
8 loss.backward()
9 optimizer.step()
10 scheduler.step()
```

# A TRAINING LIBRARY

```
1 # For alignment purposes
2 for batch in dataloader:
3 optimizer.zero_grad()
4 inputs, targets = batch
5 inputs = inputs.to(device)
6 targets = targets.to(device)
7 outputs = model(inputs)
8 loss = loss_function(outputs, targets)
9 loss.backward()
10 optimizer.step()
11 scheduler.step()
```

```
1 from accelerate import Accelerator
2 accelerator = Accelerator()
3 dataloader, model, optimizer scheduler = (
4 accelerator.prepare(
5 dataloader, model, optimizer, scheduler
6)
7)
8 for batch in dataloader:
9 optimizer.zero_grad()
10 inputs, targets = batch
11 # inputs = inputs.to(device)
12 # targets = targets.to(device)
13 outputs = model(inputs)
14 loss = loss_function(outputs, targets)
15 accelerator.backward(loss) # loss.backward()
16 optimizer.step()
17 scheduler.step()
18
```

# A TRAINING LIBRARY

What all happened in `Accelerator.prepare`?

1. `Accelerator` looked at the configuration
2. The `dataLoader` was converted into one that can dispatch each batch onto a separate GPU
3. The `model` was wrapped with the appropriate DDP wrapper from either `torch.distributed` or `torch_xla`
4. The `optimizer` and `scheduler` were both converted into an `AcceleratedOptimizer` and `AcceleratedScheduler` which knows how to handle any distributed scenario

# A TRAINING LIBRARY, MIXED PRECISION

🤪 **accelerate** also supports *automatic mixed precision*.

Through a single flag to the **Accelerator** object when calling **accelerator.backward()** the mixed precision of your choosing (such as **bf16** or **fp16**) will be applied:

```
1 from accelerate import Accelerator
2 accelerator = Accelerator(mixed_precision="fp16")
3 ...
4 for batch in dataloader:
5 optimizer.zero_grad()
6 inputs, targets = batch
7 outputs = model(inputs)
8 loss = loss_function(outputs, targets)
9 accelerator.backward(loss)
10 optimizer.step()
11 scheduler.step()
```

# A TRAINING LIBRARY, GRADIENT ACCUMULATION

Gradient accumulation in distributed setups often need extra care to ensure gradients are aligned when they need to be and the backward pass is computationally efficient.

 **accelerate** can just easily handle this for you:

```
1 from accelerate import Accelerator
2 accelerator = Accelerator(gradient_accumulation_steps=4)
3 ...
4 for batch in dataloader:
5 with accelerator.accumulate(model):
6 optimizer.zero_grad()
7 inputs, targets = batch
8 outputs = model(inputs)
9 loss = loss_function(outputs, targets)
10 accelerator.backward(loss)
11 optimizer.step()
12 scheduler.step()
```



# A TRAINING LIBRARY, GRADIENT ACCUMULATION

```
1 ddp_model, dataloader = accelerator.prepare(model, dataloader)
2
3 for index, batch in enumerate(dataloader):
4 inputs, targets = batch
5 if index != (len(dataloader)-1) or (index % 4) != 0:
6 # Gradients don't sync
7 with accelerator.no_sync(model):
8 outputs = ddp_model(inputs)
9 loss = loss_func(outputs, targets)
10 accelerator.backward(loss)
11
12 else:
13 # Gradients finally sync
14 outputs = ddp_model(inputs)
15 loss = loss_func(outputs)
16 accelerator.backward(loss)
```

# BIG MODEL INFERENCE

Stable Diffusion taking the world by storm

# BIGGER MODELS == HIGHER COMPUTE

As more large models were being released, Hugging Face quickly realized there must be a way to continue our decentralization of Machine Learning and have the day-to-day programmer be able to leverage these big models.

Born out of this effort by Sylvain Gugger:



Accelerate: Big Model Inference.

# THE BASIC PREMISE

- In PyTorch, there exists the **meta** device.
- Super small footprint to load in huge models quickly by not loading in their weights immediatly.
- As an input gets passed through each layer, we can load and unload *parts* of the PyTorch model quickly so that only a small portion of the big model is loaded in at a single time.
- The end result? Stable Diffusion v1 can be ran on < 800mb of vRAM

# THE CODE

Generally you start with something like so:

```
1 import torch
2
3 my_model = ModelClass(...)
4 state_dict = torch.load(checkpoint_file)
5 my_model.load_state_dict(state_dict)
```

But this has issues:

1. The full version of the model is loaded at **3**
2. Another version of the model is loaded into memory at **4**

If a *6 billion* parameter model is being loaded, each model class has a dictionary of 24GB so 48GB of vRAM is needed

# EMPTY MODEL WEIGHTS

We can fix step 1 by loading in an empty model skeleton at first:

```
1 from accelerate import init_empty_weights
2
3 with init_empty_weights():
4 my_model = ModelClass(...)
5 state_dict = torch.load(checkpoint_file)
6 my_model.load_state_dict(state_dict)
```

⚠ This code will not run

It is likely that just calling `my_model(x)` will fail as not all tensor operations are supported on the `meta` device.

# SHARDED CHECKPOINTS - THE CONCEPT

The next step is to have “Sharded Checkpoints” saved for your model.

Basically smaller chunks of your model weights stored that can be brought in at any particular time.

This reduces the amount of memory step 2 takes in since we can just load in a “chunk” of the model at a time, then swap it out for a new chunk through PyTorch hooks

# SHARDED CHECKPOINTS - THE CODE

```
1 from accelerate import init_empty_weights, load_checkpoint_and_dispatch
2
3 with init_empty_weights():
4 my_model = ModelClass(...)
5
6 my_model = load_checkpoint_and_dispatch(
7 my_model, "sharded-weights", device_map="auto"
8)
```

**device\_map="auto"** will tell 🤖 Accelerate that it should determine where to put each layer of the model:

1. Maximum space on the GPU(s)
2. Maximum space on the CPU(s)
3. Utilize disk space through memory-mapped tensors



# BIG MODEL INFERENCE PUT TOGETHER

```
1 from accelerate import init_empty_weights, load_checkpoint_and_dispatch
2
3 with init_empty_weights():
4 my_model = ModelClass(...)
5
6 my_model = load_checkpoint_and_dispatch(
7 my_model, "sharded-weights", device_map="auto"
8)
9 my_model.eval()
10
11 for batch in dataloader:
12 output = my_model(batch)
```

# IS THERE AN EASIER WAY?

The **transformers** library combined with the Hub makes all this code wrapping much easier for you with the **pipeline**

```
1 import torch
2 from transformers import pipeline
3 pipe = pipeline(
4 task="text-generation",
5 model="EleutherAI/gpt-j-6B",
6 device_map="auto",
7 torch_dtype=torch.float16
8)
9
10 text = pipe("This is some generated text, I think")
```

# WHAT ABOUT STABLE DIFFUSION?

A demo with `diffusers`

# SOME HANDY RESOURCES

- 🤪 Accelerate documentation
- Launching distributed code
- Distributed code and Jupyter Notebooks
- Migrating to 🤪 Accelerate easily
- Big Model Inference tutorial
- DeepSpeed and 🤪 Accelerate
- Fully Sharded Data Parallelism and 🤪 Accelerate