# HUGGING FACE ACCELERATE: MAKING DEVICE-AGNOSTIC ML TRAINING AND INFERENCE EASY AT SCALE

# WHO AM I?

- Zachary Mueller

- Technical Lead for the 🤗 Accelerate project

- Maintain the `transformers` Trainer

- API design geek

# WHAT IS 🤗 ACCELERATE?

- A training framework

- An inference framework

- A command-line interface

# A TRAINING FRAMEWORK

- Powered by PyTorch

- Change a few lines of code, gain device *and* hardware-agnostic capabilities

- Low-code, with minimal magic aimed at easy hackability and use without high-level abstractions

- We handle the intracies so you don't have to

# A TRAINING FRAMEWORK

- Support for any hardware-accelerator on the market:
    - CPU, GPU, TPU, XPU, NPU, MLU
- Automatic mixed-precision training *safely* in whatever fashion you may choose:
    - FP16, BF16, FP8 (through either `TransformerEngine` or `MS-AMP`)
- Automatic and efficient gradient accumulation
- Support for quantization through `bitsandbytes`
- Support your favorite experiment trackers (`aim`, `clearml`, `comet_ml`, `dvc-lite`, `ml-flow`, `tensorboard`, `wandb`)
- Easy to configure plugin or YAML-level API for setting up advanced frameworks like `FSDP`, `DeepSpeed`, and `Megatron-LM`

# LOW-CODE

- Biggest friction with "wrapper" libraries is control of your code
- By being minimally intrusive, your code just "works" while still giving you complete control

```
 1    import torch
 2    import torch.nn.functional as F
 3    from datasets import load_dataset
 4  + from accelerate import Accelerator
 5
 6  + accelerator = Accelerator()
 7  - device = 'cpu'
 8  + device = accelerator.device
 9
10    model = torch.nn.Transformer().to(device)
11    optimizer = torch.optim.Adam(model.parameters())
12    dataset = load_dataset('my_dataset')
13    data = torch.utils.data.DataLoader(dataset, shuffle=True)
14
15  + model, optimizer, dataloader = accelerator.prepare(model, optimizer, dataloader
16
17    model.train()
18    for epoch in range(10):
19        for source, targets in dataloader:
20            source, targets = source.to(device), targets.to(device)
21            optimizer.zero_grad()
22            output = model(source)
23            loss = F.cross_entropy(output, targets)
24  -         loss.backward()
25  +         accelerator.backward(loss)
26            optimizer.step()
```

# EASY TO INTEGRATE

- Due to the low-code nature, it's trivial to integrate into existing PyTorch frameworks:

  1. Create an `Accelerator`

```
 1    import torch
 2    import torch.nn.functional as F
 3    from datasets import load_dataset
 4  + from accelerate import Accelerator
 5
 6  + accelerator = Accelerator()
 7    device = 'cpu'
 8
 9    model = torch.nn.Transformer().to(device)
10    optimizer = torch.optim.Adam(model.parameters())
11    dataset = load_dataset('my_dataset')
12    data = torch.utils.data.DataLoader(dataset, shuffle=True)
13
14    model.train()
15    for epoch in range(10):
16        for source, targets in dataloader:
17            source, targets = source.to(device), targets.to(device)
18            optimizer.zero_grad()
19            output = model(source)
20            loss = F.cross_entropy(output, targets)
21            loss.backward()
22            optimizer.step()
```

# EASY TO INTEGRATE

- Due to the low-code nature, it's trivial to integrate into existing PyTorch frameworks:

  2. Wrap your PyTorch objects with `accelerator.prepare` and remove device-placements
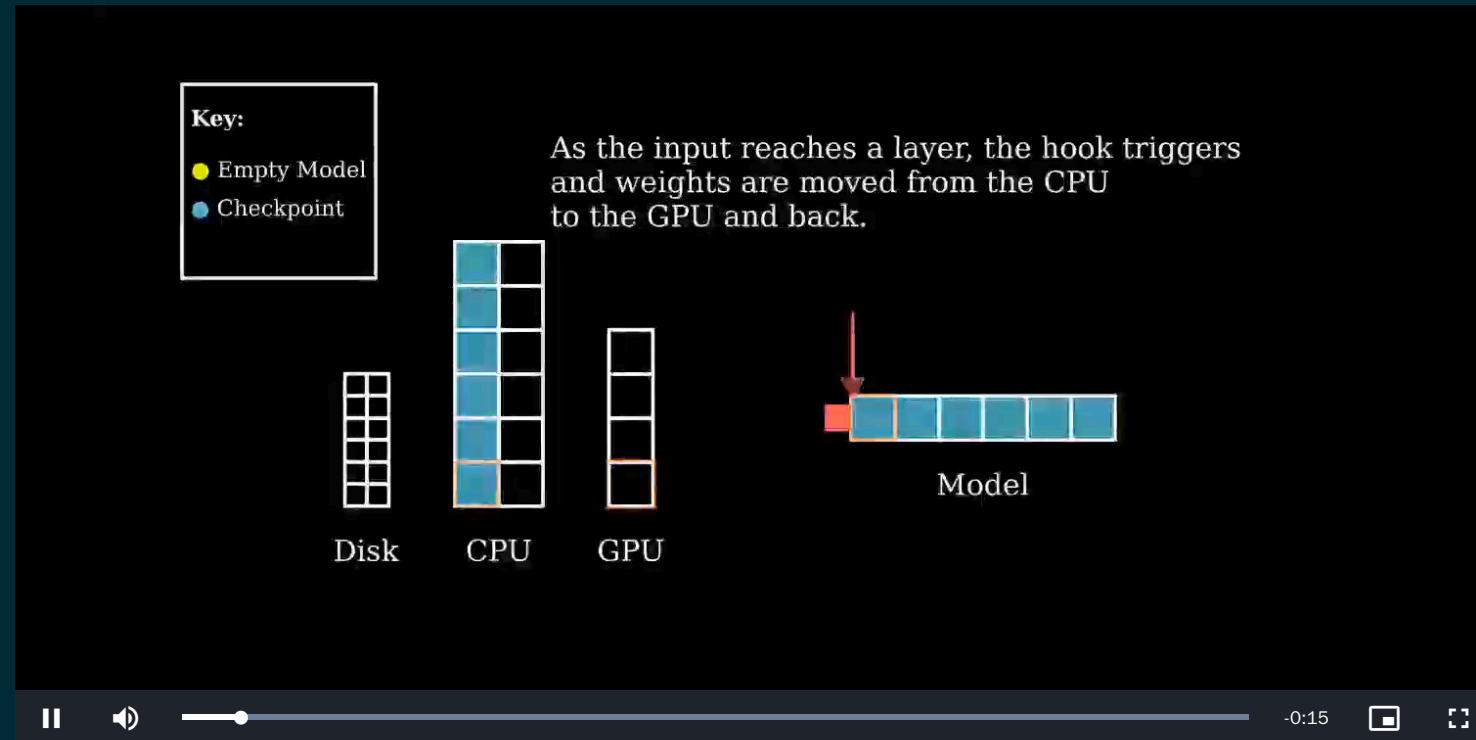
```
 1    import torch
 2    import torch.nn.functional as F
 3    from datasets import load_dataset
 4    from accelerate import Accelerator
 5
 6    accelerator = Accelerator()
 7  - device = 'cpu'
 8
 9    model = torch.nn.Transformer().to(device)
10    optimizer = torch.optim.Adam(model.parameters())
11    dataset = load_dataset('my_dataset')
12    data = torch.utils.data.DataLoader(dataset, shuffle=True)
13
14  + model, optimizer, dataloader = accelerator.prepare(model, optimizer, dataloader
15
16    model.train()
17    for epoch in range(10):
18        for source, targets in dataloader:
19            source, targets = source.to(device), targets.to(device)
20            optimizer.zero_grad()
21            output = model(source)
22            loss = F.cross_entropy(output, targets)
23            loss.backward()
24            optimizer.step()
```

# EASY TO INTEGRATE

- Due to the low-code nature, it's trivial to integrate into existing PyTorch frameworks:

  3. Use `accelerator.backward` for the backward pass

```
 1    import torch
 2    import torch.nn.functional as F
 3    from datasets import load_dataset
 4    from accelerate import Accelerator
 5
 6    accelerator = Accelerator()
 7
 8    model = torch.nn.Transformer().to(device)
 9    optimizer = torch.optim.Adam(model.parameters())
10    dataset = load_dataset('my_dataset')
11    data = torch.utils.data.DataLoader(dataset, shuffle=True)
12
13    model, optimizer, dataloader = accelerator.prepare(model, optimizer, dataloader
14
15    model.train()
16    for epoch in range(10):
17        for source, targets in dataloader:
18            source, targets = source.to(device), targets.to(device)
19            optimizer.zero_grad()
20            output = model(source)
21            loss = F.cross_entropy(output, targets)
22 -          loss.backward()
23 +          accelerator.backward(loss)
24            optimizer.step()
```

# BUT WHAT ABOUT INFERENCE?

- 🤗 Accelerate is not just for training, and has helped make the GPU-Poor take control of the narrative

- Using tools like Big Model Inference, users with *tiny* compute can run large models locally

- Started with the boom of stable diffusion, and now has scaled to having the ability to run huge LLMs locally with a single graphics card

# HOW DOES IT WORK?

- PyTorch introduced `device="meta"`

- 🤗 Accelerate introduced `device_map="auto"`

# A CLI INTERFACE

- `accelerate config`

  - Configure the environment

- `accelerate launch`

  - How to run your script

# LAUNCHING DISTRIBUTED TRAINING IS HARD

```
1  python script.py
```

vs.

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

vs.

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

How can we make this better?

# accelerate launch

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

```
1  accelerate launch --multi_gpu --num_processes 2 script.py
```

```
1  accelerate launch \
2      --multi_gpu \
3      --use_deepspeed \
4      --num_processes 2 \
5      script.py
```

# accelerate config

- Rely on `config.yaml` files

- Choose to either running `accelerate config` or write your own:

```
ddp_config.yaml
1  compute_environment: LOCAL_MACHINE
2  distributed_type: MULTI_GPU
3  main_training_function: main
4  mixed_precision: bf16
5  num_machines: 1
6  num_processes: 8
```

```
fsdp_config.yaml
1   compute_environment: LOCAL_MACHINE
2   distributed_type: FSDP
3   fsdp_config:
4     fsdp_auto_wrap_policy: TRANSFORMER_BASED_WRA
5     fsdp_backward_prefetch: BACKWARD_PRE
6     fsdp_cpu_ram_efficient_loading: true
7     fsdp_forward_prefetch: false
8     fsdp_offload_params: false
9     fsdp_sharding_strategy: FULL_SHARD
10    fsdp_state_dict_type: SHARDED_STATE_DICT
11    fsdp_sync_module_states: true
12    fsdp_use_orig_params: false
13  main_training_function: main
14  mixed_precision: bf16
15  num_machines: 1
16  num_processes: 8
```

# NOW THAT YOU'RE UP TO SPEED, WHAT'S NEW?

WE'VE HAD A BUSY LAST YEAR, AND SO HAS THE ML COMMUNITY!

# NEW TRAINING TECHNIQUES

- Quantization has taken the field by storm

- New ideas such as FSDP + QLoRA to train huge models on tiny compute!

- New precision backends as we train natively on smaller precision

- Optimizing futher how much we can push on a single machine through efficient RAM and timing techniques

# LARGER COMPUTE LANDSCAPE

- As we search for alternatives to NVIDIA, new compilers rise:

  - XPU (Intel)

  - NPU (Intel)

  - MLU (Cambricon)

All of which are supported by 🤗 Accelerate

# LOWER ABSTRACTIONS

- While the `Accelerator` was great, needed better abstractions focused on controlling behaviors

- Introduced the `PartialState`

```
1  from accelerate import PartialState
2
3  if PartialState().is_main_process:
4    # Run on only 1 device
5
6  with PartialState().main_process_first:
7    # Useful for dataset processing
8
9  # Device-agnostic without the bulk of the `Accelerator`
10 device = PartialState().device
```

# FASTER AND BETTER INFERENCE ALTERNATIVES

- `PiPPy` gives us efficient pipeline-parallelism in distributed environments to increase throughput while keeping a simple torch-bound API

- Rather than having to wait for each GPU, every GPU can be busy in parallel

```python
1  import torch
2  from transformers import AutoModelForSequenceClassification
3
4  from accelerate import PartialState, prepare_pippy
5
6  model = AutoModelForSequenceClassification.from_pretrained("gpt2")
7  model.eval()
8
9  input = torch.randint(
10     low=0,
11     high=model.config.vocab_size,
12     size=(2, 1024),   # bs x seq_len
13     device="cpu",
14 )
15
16 model = prepare_pippy(model, split_points="auto", example_args=(input,
17
18 with torch.no_grad():
19     output = model(input)
```

# ADOPTION: ACCELERATE IN THE ECOSYSTEM

# ACCELERATE IN THE ECOSYSTEM

- Many of the frameworks you use daily already rely on 🤗 Accelerate!
  - Nearly all of 🤗
  - `axolotl`
  - `fastai`
  - `FastChat`
  - `lucidrains`
  - `kornia`

# ACCELERATE IN THE ECOSYSTEM

- Started as a way to isolate out distributed code on TPU and `DistributedDataParallelism`

# ACCELERATE IN THE ECOSYSTEM

- Now is the backbone of some of the largest PyTorch training frameworks in the ecosystem

# WHAT'S NEXT?

# ELEVATING THE COMMUNITY

- Now that more advanced training techniques are reachable (FSDP, DeepSpeed, etc), we need to focus on educating the community on how to use it best

- Goes beyond how to use the `Trainer` or `Accelerator`, but how to use *what* where

- Keep Accelerate as a tool for the community to utilize when new techniques come out and play with, to push new ideas to scale quickly

# 1.0.0: SOON!

- Tried and battle-tested by over 7M users/month

- As we've been stable for over a year now, we're near ready to release 1.0.0

# THANKS FOR JOINING!

- 🤗 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
- FSDP vs DeepSpeed In-Depth