

The Alan Turing Institute

Tutorial 3: Multi-node training with Lightning

HPC Training – 25-26 November 2025
David Llewellyn-Jones



Motivation

GPUs are great for accelerated training and inference, but processing is bounded by:

1. Available memory
2. Speed of computation

How to make AI better?

1. Improved algorithms
2. Larger models
3. More training data

The last two require *more compute*

Scaling up

1. Strong motivation for scaling across GPUs
2. Single device: 4 or 8 GPUs maximum
3. Eventually want to scale across nodes
4. Get this right... the sky's the limit

Preparation

1. Open Baskerville docs
`https://docs.baskerville.ac.uk`
2. Select Baserville Portal
3. Login if necessary
4. Select JupyterLab Apps list
5. Configure the session
6. Launch

Session configuration

1. Kernel: Python 3.11.3
2. Show Conda Envs: No
3. Number of hours: 2
4. Number of GPUs: 1
5. Project: vjgo8416-hpc2511
6. Queue: turing

JupyterLab - Baskerville OnDem X

portal.baskerville.ac.uk/pur/sys/dashboard/batch_connect/sys_bc_bask_jupyter/session_contexts/new

Interactive Apps

JupyterLab

GUTs

CST Studio Suite

Fiji

Linaro-Forge

RELION

cisTEM

Servers

TensorBoard

doppio

JupyterLab

This app will launch a [JupyterLab](#) server (supporting [Python](#) and [Julia](#) kernels) on [Baskerville](#).

Kernel to load

Python 3.11.3 (2023a / GCCcore-12.3.0)

This defines the kernel you wish to load. The extra packages for use in Python can be selected from inside Jupyter using the Lmod extension. For further information please refer to the [JupyterLab documentation](#).

☐ Show Conda Environments

Include kernels for compatible Conda environments found in: `~/ .conda/ environments.txt`. For further information please see our [Conda environment documentation](#).

Number of hours

2

Number of GPUs

1

Number of GPUs to request. For each GPU you will get 25% of a node's total cores.

Baskerville Project

vjgo8416-hpc2511

Please select the Baskerville Project to which the job will be attached.

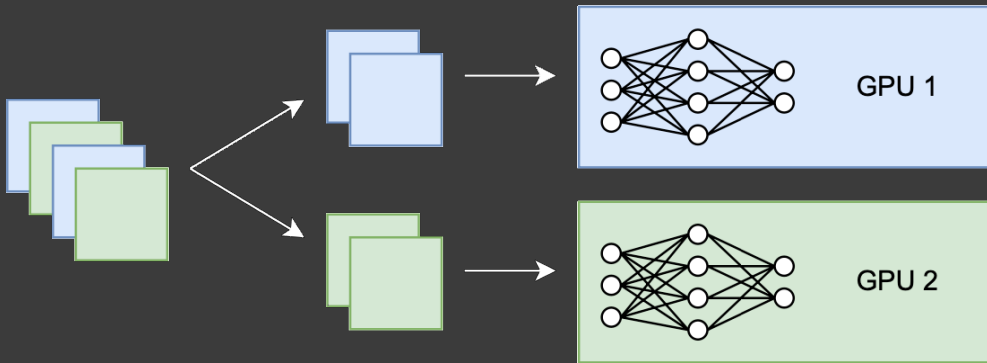
Queue

turing

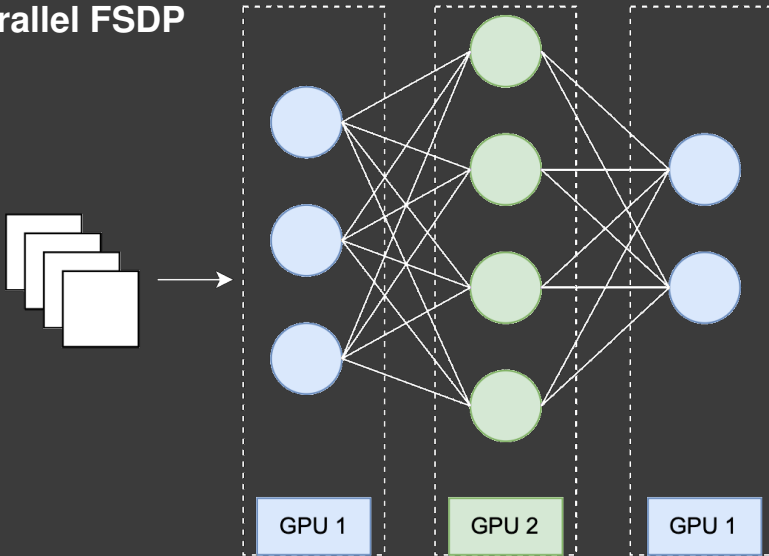
Please select the Queue/QoS on which your job will run.

Launch

Data parallel

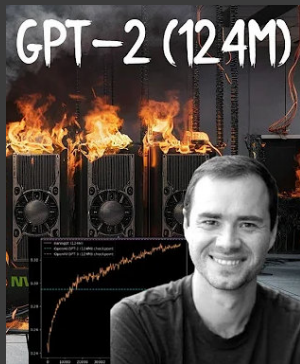


Model parallel FSDP

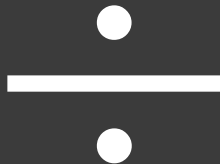


The plan

1. No more theory
2. Simplified GPT2 nano code for one GPU
3. Andrej Karpathy:
<https://youtu.be/l8pRSuU81PU>
4. Extend to support Distributed Data Parallel
5. Extend using PyTorch Lightning



Your most important tools



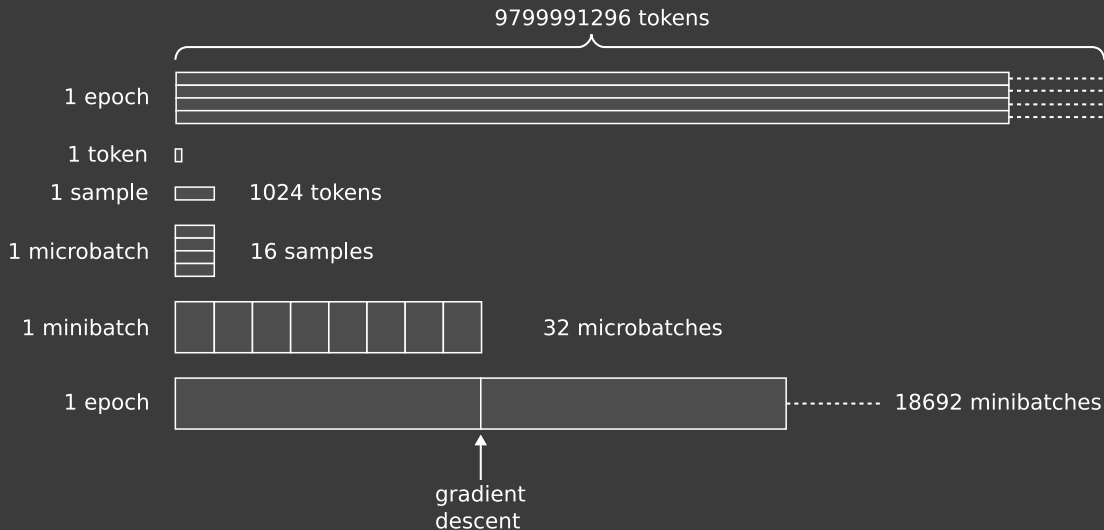
$(*, /)$

Training terminology

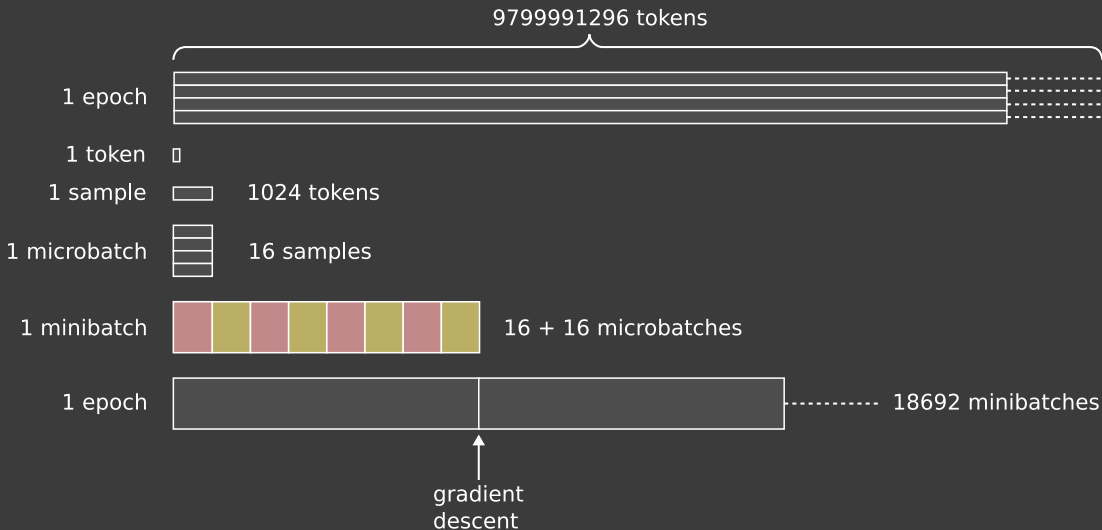
1. Vocabulary: 50257 embeddings
2. Dataset: 9799991296 FineWeb tokens, our training data
3. Sample: a sequence of 1024 tokens from the dataset
4. Microbatch: 16 samples
5. Minibatch (also called a batch): 32 microbatches, gradient accumulation performed afterwards
6. Step: the process for training on one minibatch
7. Epoch: one training sweep of the entire dataset, 18692 steps

$$18692 \times 16 \times 32 \times 1024 = 9799991296$$

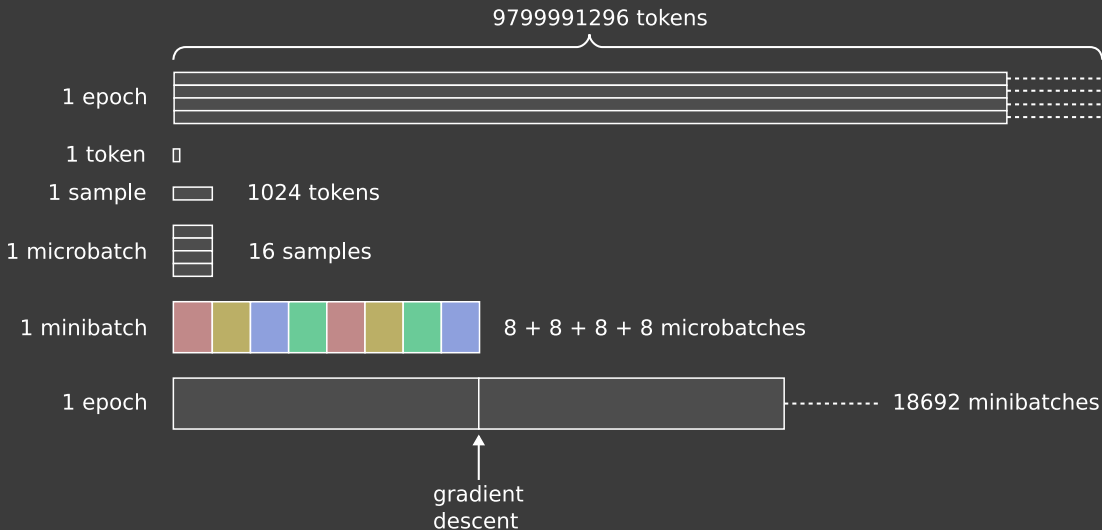
Training with 1 GPU



Training with 2 GPUs



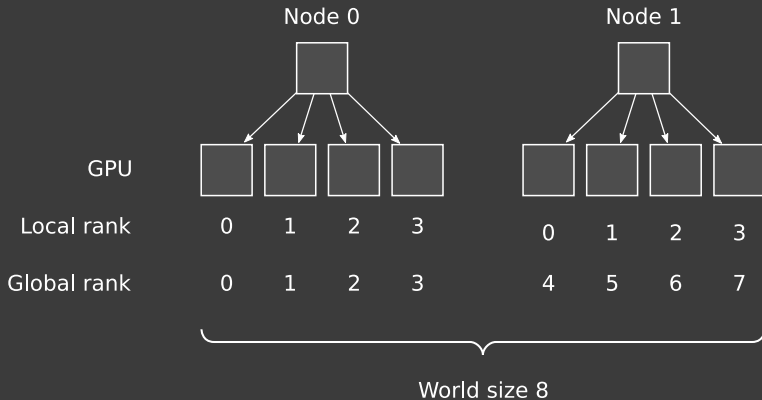
Training with 4 GPUs



Distributed training terminology

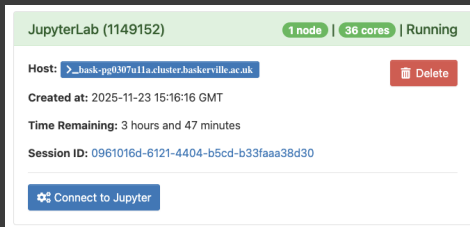
1. Training on n nodes, each with m GPUs
2. Total $n \times m$ GPUs
3. World size: $w = m \times n$
4. Global rank r_G is a unique index $r_G \in \{0, \dots, w - 1\}$
5. Local rank r_L is unique per device $r_L \in \{0, \dots, m - 1\}$
6. No node index, but we do have *hostnames*, e.g. bask-pg0309u05a, bask-pg0309u06a
7. We may also have multiple CPU *workers* for each node

Training with multiple nodes



Time to look at the code

1. Open your JupyterLab page
`https://portal.baskerville.ac.uk`
2. Connect to Jupyter
3. Select the File Browser in the vertical bar on the left
4. Move to the directory:
`/vjgo8416-hpc2511/$USER/hpc-training-nov-2025/3-Training/code`
5. Open the `train_gpt2.py` file by double-clicking on it



GPT2 nano classes and functions

1. CausalSelfAttention, MLP, Block: model components
2. GPTConfig: model configuration
3. GPT: the model
4. generate(): inference
5. configure_optimizers(): training configuration
6. training_step(): one training step
7. load_tokens(): load a single FineWeb shard
8. get_shards(): find the FineWeb shards on disk
9. DataIterator: our dataset and data loader
10. get_lr(): calculate the learning rate

GPT2 nano execution

1. Set hyperparameters
2. Create model
3. Perform training loop

Training time!

1. Right click on the `train_gpt2.py` tab
2. Select **New View for Python File**
3. Open a **GPU Resources** pane from the GPU Dashboard
4. Drag the pane to the tab space on the right
5. Open a terminal in the left tab space
6. `source activate.sh`
7. `python train_gpt2.py`

train_gpt2.py - JupyterLab

portal.baskerville.ac.uk | node/bask-pg0307u29a.cluster.baskerville.ac.uk/29569/lab/tree/vjgo8416-hpc2511/ovau2564/hpc-training-nov-2025/3-Training/code/train_gpt2.py

File Edit View Run Kernel Git Tabs Settings Help

train_gpt2.py

Terminal 2

get_shards

Filter files by name

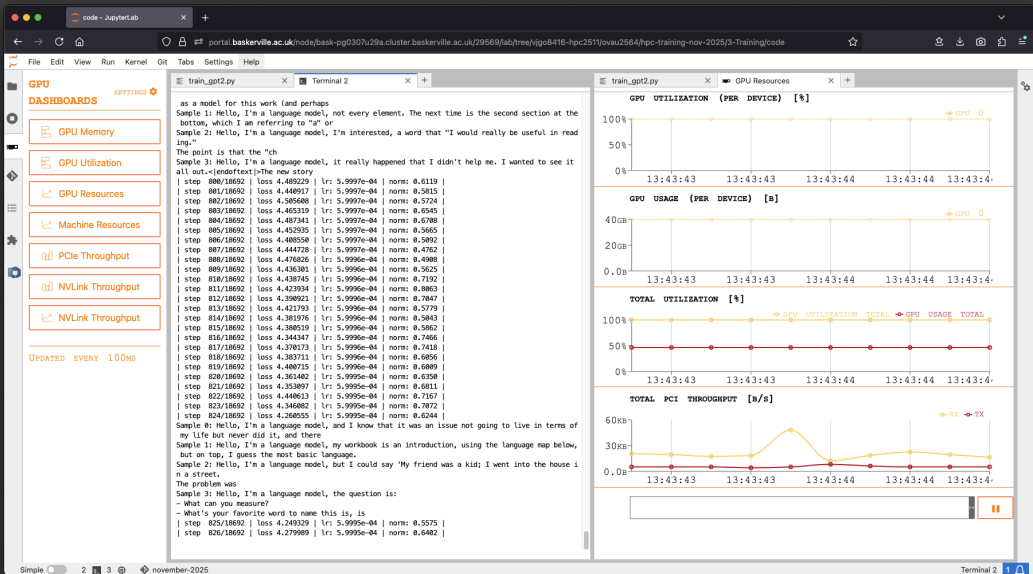
Name Modified

Name	Modified
batchfiles	4h ago
example-logs	4h ago
fineweb	1h ago
notebooks	1h ago
solutions	4h ago
venv	56m ago
activate.sh	57m ago
README.md	4h ago
requirements.txt	4h ago
train_gpt2.py	2h ago

```
163
164
165
166
167 for i
168 to
169 de
170 pr
171
172
173 def config
174 param_
175 param_
176 decay_
177 nodeca
178 optim
179
180 ("params": decay_params, "weight_decay": 0.1},
181 {"params": nodecay_params, "weight_decay": 0.0})
182
183 use_fused = "fused" in inspect.signature(torch.optim.AdamW).parameters
184 # AdamW docs: https://pytorch.org/docs/stable/generated/torch.optim.AdamW.html#torch.optim.AdamW
185 optimizer = torch.optim.AdamW(optimizer_groups, lr=max_lr, betas=(0.9, 0.95), eps=1e-8, fused=use_fused)
186 return optimizer
187
188
189 def training_step(self, batch, batch_idx):
190 x, y = batch
191 x, y = x.to(device), y.to(device)
192 with torch.autocast(device_type="cuda", dtype=torch.bfloat16):
193 logits, loss = self(x, y)
194 loss = loss / grad_accum_steps
195 return loss
196
197
198 def load_tokens(filename):
199 npt = np.load(filename)
200 npt = npt.astype(np.int32)
201 ptt = torch.tensor(npt, dtype=torch.long)
202 return ptt
203
204 def get_shards(split):
205 data_root = "/bask/projects/v/vjgo8416-hpc2511/edu_fineweb10B"
206 shards = os.listdir(data_root)
207 shards = [s for s in shards if split in s]
208 shards = sorted(shards)
209 shards = [os.path.join(data_root, s) for s in shards]
210 return shards
```

Simple 2 Python november-2025

Ln 186, Col 47 Spaces: 4 train_gpt2.py



Upgrade to DDP

1. Open `train_gpt2.py` in the left hand tab space
2. Open `notebooks/diff_ddp.ipynb` in the right hand tab space
3. Execute the first cell of `diff_ddp.ipynb`

diff_ddp.ipynb [auto-s] - Jupyter: X

portal.baskerville.ac.uk/node/bask-pg0308u06a.cluster.baskerville.ac.uk/81025/lab/workspaces/auto-s/tree/vjgo8416-hpc2511/ovau2564/hpc-training-nov-2025/3-1 193%

File Edit View Run Kernel Git Tabs Settings Help

+

Filter files by name

/ ... / code / notebooks /

Name	Modified
clean-output.md	3d ago
diff_ddp.ipynb	28s ago
diff_lit.ipynb	21h ago
README.md	3d ago

Terminal 1

train_gpt2.py

```
1 #!python3
2 # vim: et:ts=4:sts=4:sw=4
3
4 from torch.distributed import init_process_group,
5   destroy_process_group
6 from torch.nn.parallel import DistributedDataParallel as
7   DDP
8 import torch.distributed as dist
9
10 from dataclasses import dataclass
11 import torch
12 import torch.nn as nn
13 from torch.nn import functional as F
14 import math
15 import tiktoken
16 import inspect
17 import os
18 import numpy as np
19
20
21 class CausalSelfAttention(nn.Module):
22
23     def __init__(self, config):
24         super().__init__()
25         assert config.n_embd % config.n_head == 0
26         # Key, query, value projections for all heads,
27         # but in a batch
28         self.c_attn = nn.Linear(config.n_embd, 3 *
29 config.n_embd)
30         # Output projection
31         self.c_proj = nn.Linear(config.n_embd,
```

diff_ddp.ipynb

Code

git

...

```
[2]: !git diff -R diff-ddp -- ../train_gpt2.py

diff --git b/code/train_gpt2.py a/code/train_gpt
2.py
index 44b31a3..47aea41 100644
--- b/code/train_gpt2.py
+++ a/code/train_gpt2.py
@@ -1,6 +1,10 @@
#!python3
# vim: et:ts=4:sts=4:sw=4

+from torch.distributed import init_process_group
+, destroy_process_group
+from torch.nn.parallel import DistributedDataPar
+allel as DDP
+import torch.distributed as dist
+
+
+from dataclasses import dataclass
+import torch
+import torch.nn as nn
@@ -220,11 +224,11 @@ class DataIterator:
     self.process_rank = process_rank
     self.current_shard = 0
     self.tokens = load_tokens(self.shards[se
lf.current_shard])
-
-     self.current_position = 0
-     assert (len(self) * self.B * self.T) <
- (len(self.tokens) * len(self.shards)), f"Not enou
gh data for a complete epoch"
+
+     self.current_position = self.B * self.T
+ * self.process_rank
+
+     assert (len(self) * self.B * self.T * se
```

Simple 1 1 november-2025 Python | Idle Mode: Command Ln 1, Col 1 diff_ddp.ipynb 1

Understanding unified diffs

1. @@ prefix indicates line numbers
@@ -before,len +after,len @@
2. + prefix in green indicates lines added
3. - prefix in red indicates lines removed
4. Replay the cell to update the diff after making changes

```
@@ -145,7 +149,7 @@ class GPT(nn.Module):
    loss = F.cross_entropy(logits.view(
s.view(-1))
    return logits, loss

- def generate(self):
+ def generate(self, rank):
    num_return_sequences = 4
    max_length = 32
    tokens = enc.encode("Hello, I'm a lang
@@ -153,7 +157,7 @@ class GPT(nn.Module):
    tokens = tokens.unsqueeze(0).repeat(nu
xgen = tokens.to(device)
    sample_rng = torch.Generator(device=de
- sample_rng.manual_seed(42)
+ sample_rng.manual_seed(42 + rank)
    while xgen.size(1) < max_length:
        with torch.no_grad():
            with torch.autocast(device_type

6):
@@ -207,22 +211,24 @@ def get_shards(split):

    class DataIterator:

- def __init__(self, B, T):
+ def __init__(self, B, T, num_processes, pr
```

Manually apply the diff

1. Imports
2. Wrap the model in the DDP class
3. Initialise the world parameters
 - 3.1 Harvest **world size**, **rank** and **local rank** from the environment
 - 3.2 Where do these come from?
4. Initialise the process group
5. Add a stride to our data iterator

Observations

1. Most of this is boilerplate
2. The hardest part is sharding the data correctly
3. Because this is *data parallel*
4. Communication between processes is also hard

...but done for us by `torch.distributed` and DDP

Training time!

See `batch-ddp-1n1g.sh`, `batch-ddp-1n2g.sh` and `batch-ddp-2n2g.sh`

```
1 # Single node
2 python -m torch.distributed.launch \
3     --standalone \
4     --nproc-per-node=${SLURM_GPUS_PER_NODE} \
5     train_gpt2.py
6
7 # Multi-node
8 srun bash -c 'python -m torch.distributed.launch \
9     --nproc_per_node=${SLURM_GPUS_PER_NODE} \
10    --nnodes=${SLURM_NNODES} \
11    --master-port=${MASTER_PORT} \
12    --master-addr=${MASTER_ADDR} \
13    --node_rank=${SLURM_PROCID} \
14    train_gpt2.py'
```

Upgrade to Lightning

1. Lightning is conceptually different
2. Code is organised in `LightningModule`: init, train step, validation step, test step, optimisers
3. Code outside `LightningModule` is automated by `Trainer`
4. Remove code moving data to the GPU

Upgrade to Lightning

1. Open `train_gpt2.py` in the left hand tab space
2. Open `diff_lit.ipynb` in the right hand tab space
3. Execute the first cell of `diff_lit.ipynb`

Manually apply the diff

1. Imports
2. Switch from `Module` to `LightningModule`
3. The device is established for us
4. The Learning Rate scheduler is baked in
5. Lightning handles where the data goes
6. Drop periodic example generation
7. Dataloader must be managed by a `DataLoader`
8. World initialisation now handled by Lightning
9. The training loop is all Lightning

Observations

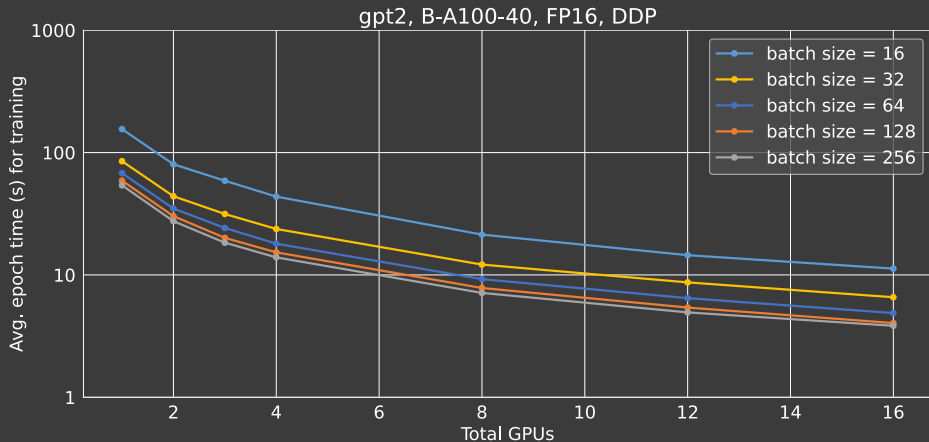
1. The training code is greatly simplified
2. Lightning is SLURM-aware
3. So multi-node training is even easier
4. We can now switch to other strategies

Training time!

See `batch-lit-1n1g.sh` – `batch-lit-2n4g.sh`

```
1 # From inside a SLURM batch script
2 srun python train_gpt2.py
3
4 # Queue execution to run
5 sbatch batch-lit-1n1g.sh
```

Scaling across nodes on A100 (40 GiB) GPUs



Wrapping up

1. Converting code for distributed training is doable
2. Review and try out the batch scripts
3. More examples in the `hpc-landscapes` repository

<https://github.com/alan-turing-institute/hpc-landscape>