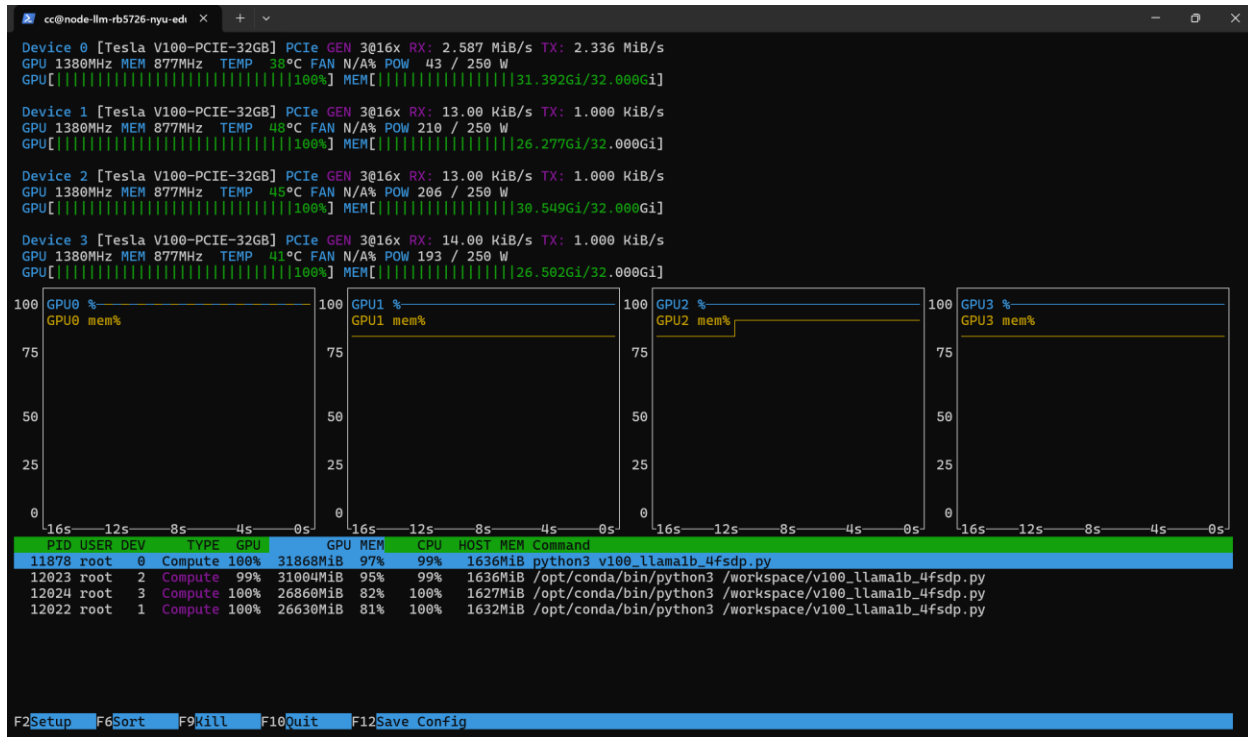


TinyLlama 1.1B model on 4x V100 32GB with FSDP



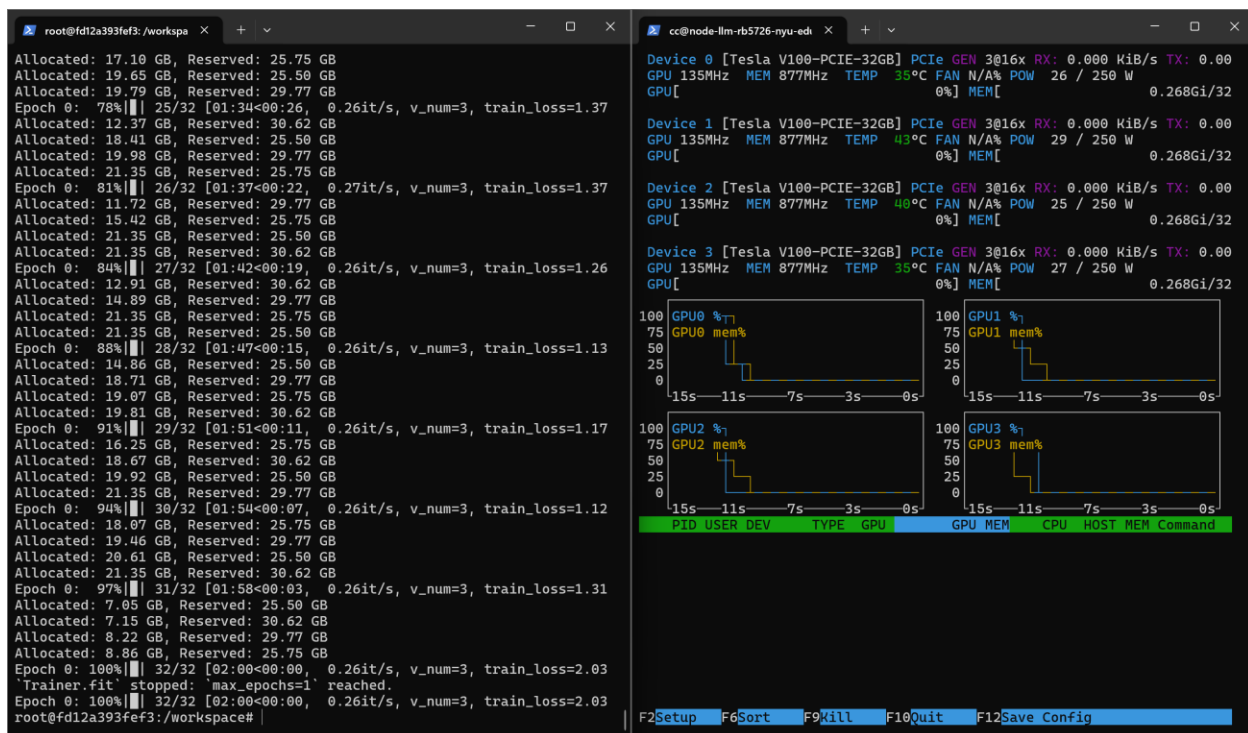
nvidia-smi

```

root@fd12a393fef3:/workspace
Allocated: 17.10 GB, Reserved: 25.75 GB
Allocated: 19.65 GB, Reserved: 25.50 GB
Allocated: 19.79 GB, Reserved: 29.77 GB
Epoch 0: 78%| 25/32 [01:34<00:26, 0.26it/s, v_num=3, train_loss=1.37]Allocated: 12.37 GB, Reserved: 30.62 GB
Allocated: 18.41 GB, Reserved: 25.50 GB
Allocated: 19.98 GB, Reserved: 29.77 GB
Allocated: 21.35 GB, Reserved: 25.75 GB
Epoch 0: 81%| 26/32 [01:37<00:22, 0.27it/s, v_num=3, train_loss=1.37]Allocated: 11.72 GB, Reserved: 29.77 GB
Allocated: 15.42 GB, Reserved: 25.75 GB
Allocated: 21.35 GB, Reserved: 25.50 GB
Allocated: 21.35 GB, Reserved: 30.62 GB
Epoch 0: 84%| 27/32 [01:42<00:19, 0.26it/s, v_num=3, train_loss=1.26]Allocated: 12.91 GB, Reserved: 30.62 GB
Allocated: 14.89 GB, Reserved: 29.77 GB
Allocated: 21.35 GB, Reserved: 25.75 GB
Allocated: 21.35 GB, Reserved: 25.50 GB
Epoch 0: 88%| 28/32 [01:47<00:15, 0.26it/s, v_num=3, train_loss=1.13]Allocated: 14.86 GB, Reserved: 25.50 GB
Allocated: 18.71 GB, Reserved: 29.77 GB
Allocated: 19.07 GB, Reserved: 25.75 GB
Allocated: 19.81 GB, Reserved: 30.62 GB
Epoch 0: 91%| 29/32 [01:51<00:11, 0.26it/s, v_num=3, train_loss=1.17]Allocated: 16.25 GB, Reserved: 25.75 GB
Allocated: 18.67 GB, Reserved: 30.62 GB
Allocated: 19.92 GB, Reserved: 25.50 GB
Allocated: 21.35 GB, Reserved: 29.77 GB
Epoch 0: 94%| 30/32 [01:54<00:07, 0.26it/s, v_num=3, train_loss=1.12]Allocated: 18.07 GB, Reserved: 25.75 GB
Allocated: 19.46 GB, Reserved: 29.77 GB
Allocated: 20.61 GB, Reserved: 25.50 GB
Allocated: 21.35 GB, Reserved: 30.62 GB
Epoch 0: 97%| 31/32 [01:58<00:03, 0.26it/s, v_num=3, train_loss=1.31]Allocated: 7.05 GB, Reserved: 25.50 GB
Allocated: 7.15 GB, Reserved: 30.62 GB
Allocated: 8.22 GB, Reserved: 29.77 GB
Allocated: 8.86 GB, Reserved: 25.75 GB
Epoch 0: 100%| 32/32 [02:00<00:00, 0.26it/s, v_num=3, train_loss=2.03]Trainer.fit' stopped: 'max_epochs=1' reached.
Epoch 0: 100%| 32/32 [02:00<00:00, 0.26it/s, v_num=3, train_loss=2.03
root@fd12a393fef3:/workspace#

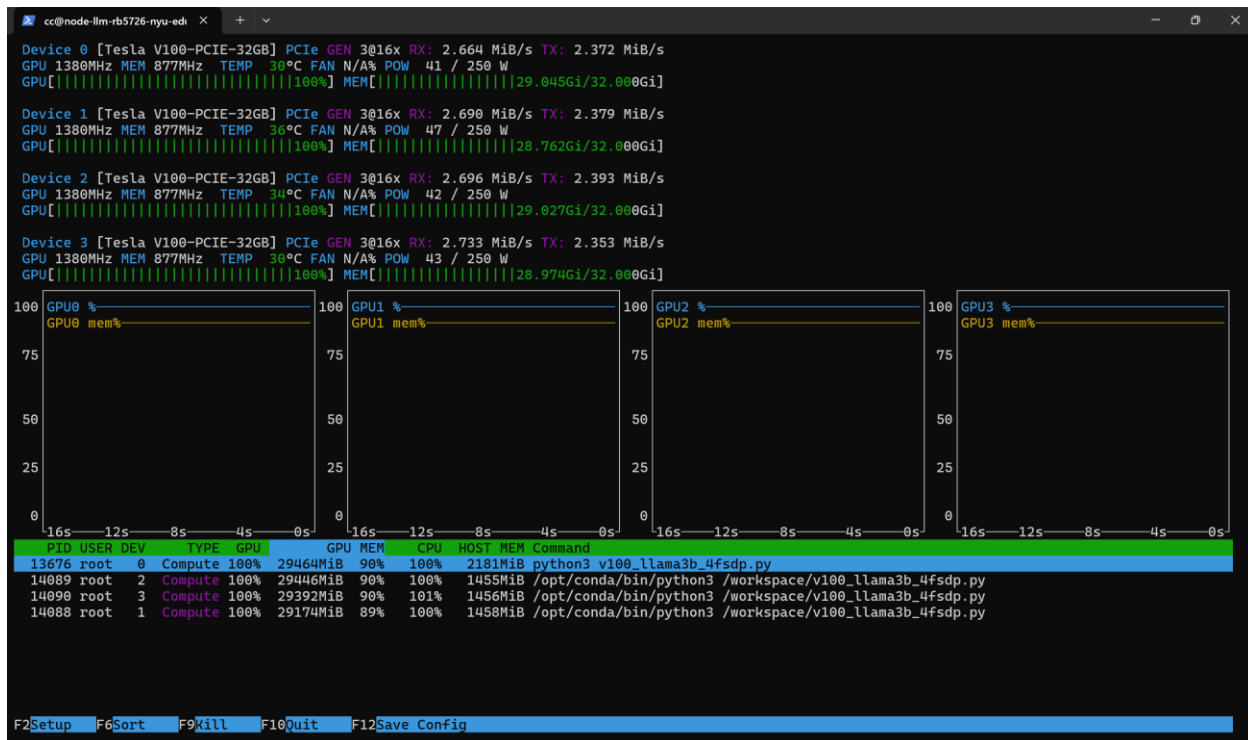
```

Training time



After completion of TinyLlama 1.1B model on 4x V100 32GB with FSDP

OpenLLaMA 3B model on 4x V100 32GB with FSDP



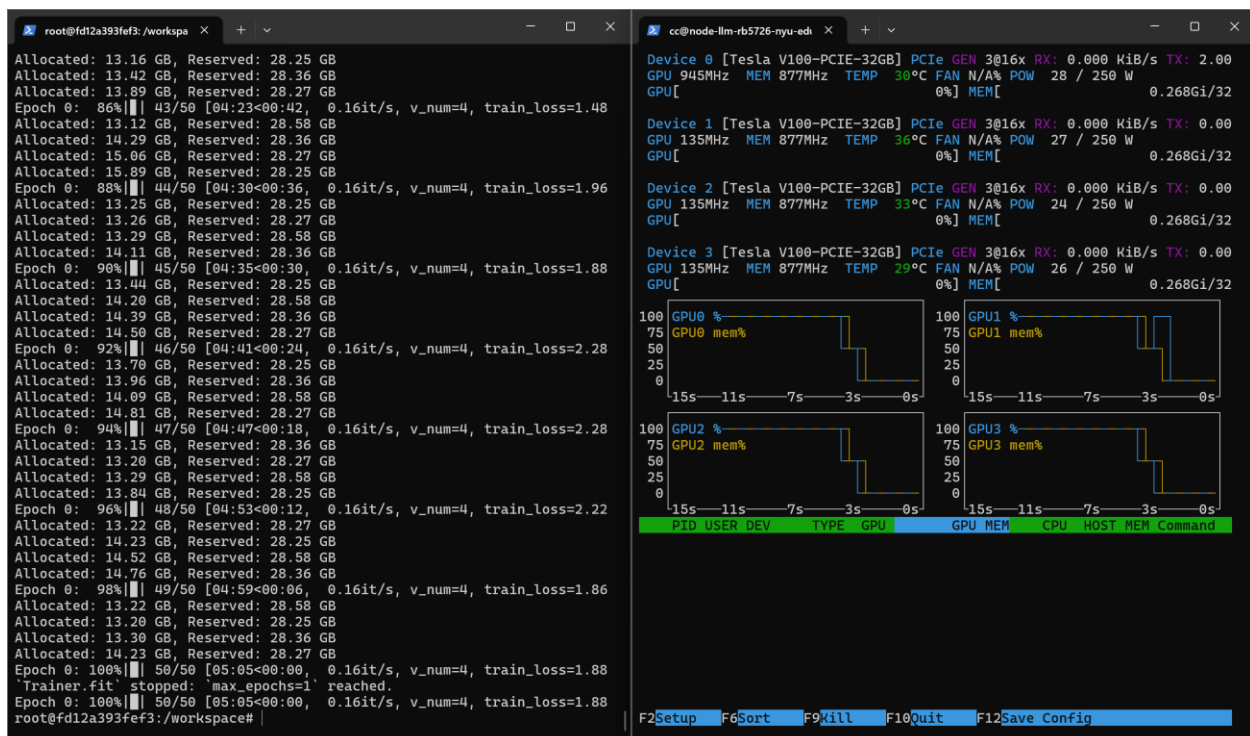
nvidia-smi

```

root@fd12a393fef3:/workspace
Allocated: 13.16 GB, Reserved: 28.25 GB
Allocated: 13.42 GB, Reserved: 28.36 GB
Allocated: 13.89 GB, Reserved: 28.27 GB
Epoch 0: 86% 43/50 [04:23<00:42, 0.16it/s, v_num=4, train_loss=1.48] Allocated: 13.12 GB, Reserved: 28.58 GB
Allocated: 14.29 GB, Reserved: 28.36 GB
Allocated: 15.06 GB, Reserved: 28.27 GB
Allocated: 15.89 GB, Reserved: 28.25 GB
Epoch 0: 88% 44/50 [04:30<00:36, 0.16it/s, v_num=4, train_loss=1.96] Allocated: 13.25 GB, Reserved: 28.25 GB
Allocated: 13.26 GB, Reserved: 28.27 GB
Allocated: 13.29 GB, Reserved: 28.58 GB
Allocated: 14.11 GB, Reserved: 28.36 GB
Epoch 0: 90% 45/50 [04:35<00:30, 0.16it/s, v_num=4, train_loss=1.88] Allocated: 13.44 GB, Reserved: 28.25 GB
Allocated: 14.20 GB, Reserved: 28.58 GB
Allocated: 14.39 GB, Reserved: 28.36 GB
Allocated: 14.50 GB, Reserved: 28.27 GB
Epoch 0: 92% 46/50 [04:41<00:24, 0.16it/s, v_num=4, train_loss=2.28] Allocated: 13.70 GB, Reserved: 28.25 GB
Allocated: 13.96 GB, Reserved: 28.36 GB
Allocated: 14.09 GB, Reserved: 28.58 GB
Allocated: 14.81 GB, Reserved: 28.27 GB
Epoch 0: 94% 47/50 [04:47<00:18, 0.16it/s, v_num=4, train_loss=2.28] Allocated: 13.15 GB, Reserved: 28.36 GB
Allocated: 13.20 GB, Reserved: 28.27 GB
Allocated: 13.29 GB, Reserved: 28.58 GB
Allocated: 13.84 GB, Reserved: 28.25 GB
Epoch 0: 96% 48/50 [04:53<00:12, 0.16it/s, v_num=4, train_loss=2.22] Allocated: 13.22 GB, Reserved: 28.27 GB
Allocated: 14.23 GB, Reserved: 28.25 GB
Allocated: 14.52 GB, Reserved: 28.58 GB
Allocated: 14.76 GB, Reserved: 28.36 GB
Epoch 0: 98% 49/50 [04:59<00:06, 0.16it/s, v_num=4, train_loss=1.86] Allocated: 13.22 GB, Reserved: 28.58 GB
Allocated: 13.20 GB, Reserved: 28.25 GB
Allocated: 13.30 GB, Reserved: 28.36 GB
Allocated: 14.23 GB, Reserved: 28.27 GB
Epoch 0: 100% 50/50 [05:05<00:00, 0.16it/s, v_num=4, train_loss=1.88] 'Trainer.fit' stopped: 'max_epochs=1' reached.
Epoch 0: 100% 50/50 [05:05<00:00, 0.16it/s, v_num=4, train_loss=1.88]
root@fd12a393fef3:/workspace#

```

Training time



After completion of OpenLLaMA 3B model on 4x V100 32GB with FSDP

We cannot train this model on a single V100 without running out of memory - try

run inside pytorch container

python3 v100_llama3b_1device.py

```
root@fd12a393fef3: /workspa x + v
File "/opt/conda/lib/python3.11/site-packages/lightning/pytorch/strategies/strategy.py", line 239, in optimizer_step
    return self.precision_plugin.optimizer_step(optimizer, model=model, closure=closure, **kwargs)
File "/opt/conda/lib/python3.11/site-packages/lightning/pytorch/plugins/precision/precision.py", line 123, in optimizer_step
    return optimizer.step(closure=closure, **kwargs)
File "/opt/conda/lib/python3.11/site-packages/torch/optim/optimizer.py", line 487, in wrapper
    out = func(*args, **kwargs)
File "/opt/conda/lib/python3.11/site-packages/torch/optim/optimizer.py", line 91, in _use_grad
    ret = func(self, *args, **kwargs)
File "/opt/conda/lib/python3.11/site-packages/torch/optim/adamw.py", line 220, in step
    adamw(
File "/opt/conda/lib/python3.11/site-packages/torch/optim/optimizer.py", line 154, in maybe_fallback
    return func(*args, **kwargs)
File "/opt/conda/lib/python3.11/site-packages/torch/optim/adamw.py", line 782, in adamw
    func(
File "/opt/conda/lib/python3.11/site-packages/torch/optim/adamw.py", line 606, in _multi_tensor_adamw
    exp_avg_sq_sqrt = torch._foreach_sqrt(device_exp_avg_sqs)
torch.cuda.OutOfMemoryError: CUDA out of memory. Tried to allocate 54.00 MiB. GPU 0 has a total capacity of 31.73 GiB of which 18.19 MiB is free. Process
15783 has 31.71 GiB memory in use. Of the allocated memory 30.68 GiB is allocated by PyTorch, and 669.94 MiB is reserved by PyTorch but unallocated
. If reserved but unallocated memory is large try setting PYTORCH_CUDA_ALLOC_CONF=expandable_segments:True to avoid fragmentation.  See documentatio
n for Memory Management (https://pytorch.org/docs/stable/notes/cuda.html#environment-variables)
root@fd12a393fef3: /workspace#
```