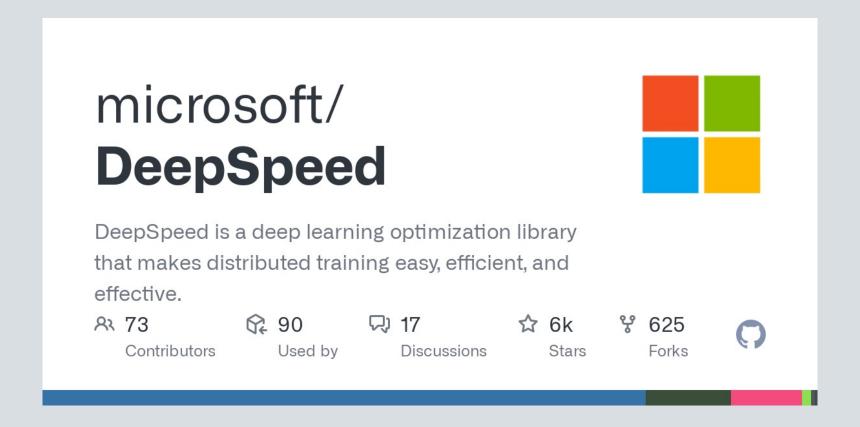# Deploy and Test DeepSpeed on ThetaGPU

**Zhen Xie, SDL Day 1**

# Another Distributed DNN Training Framework



- Microsoft's DeepSpeed is a new open-source framework focused on optimizing the training of massively large deep learning models.

# Another Distributed DNN Training Framework



- **Scale:** DeepSpeed provides system support to run models up to 100 billion parameters.

- **Speed:** DeepSpeed showed 4x-5x higher throughput than other libraries.

- **Cost:** DeepSpeed has three times less cost than other libraries.

- **Usability:** DeepSpeed does not require refactoring PyTorch models and could be used with just a few lines of code.

Argonne
NATIONAL LABORATORY

# Deploy DeepSpeed and Test on ThetaGPU

PATH:
sdl_ai_workshop/01_distributedDeepLearning/DeepSpeed/submissions/thetagpu/deepspeed_cifar10_runner_single_node.sh

```
1    #!/bin/bash
2    #COBALT -n 2
3    #COBALT -t 0:10:00 -q full-node
4    #COBALT -A SDL_Workshop
5    #COBALT --attrs=pubnet
6    #COBALT -O cifar10_1node_4gpus
7
8    #submisstion script for running cifar10 with deepspeed
9
10   echo "Running Cobalt Job $COBALT_JOBID."
11
12   echo "Setting up env"
13
14   conda env create --name deepspeed --file /lus/theta-fs0/projects/datascience/zhen/env_deepspeed.yml
15
16   conda activate deepspeed
17
18   cd /lus/theta-fs0/projects/datascience/zhen/DeepSpeed
19
20   echo "Current directory: "
21   pwd
22
23   echo "Run script: "
24   deepspeed --hostfile=$COBALT_NODEFILE cifar10_deepspeed.py --deepspeed --deepspeed_config ds_config.json $@
```

#Step1: "Setting up env"

#Step2: "Run script"

Argonne
NATIONAL LABORATORY

# Output of cifar-10 example

```
Current directory:
/lus/theta-fs0/projects/datascience/zhen/DeepSpeed
Run script:
[2021-10-05 10:05:46,787] [WARNING] [runner.py:122:fetch_hostfile] Unable to find hostfile, will proceed with training with local resources only.
[2021-10-05 10:05:46,880] [INFO] [runner.py:360:main] cmd = /home/zhen/zhen/anaconda3/envs/deepspeed/bin/python3.9 -u -m deepspeed.launcher.launch --world_info=eyJsb2NhbGhvc3QiOiBbMF19 --maste
r_addr=127.0.0.1 --master_port=29500 cifar10_deepspeed.py --deepspeed --deepspeed_config ds_config.json
[2021-10-05 10:05:47,909] [INFO] [launch.py:80:main] WORLD INFO DICT: {'localhost': [0]}
[2021-10-05 10:05:47,909] [INFO] [launch.py:86:main] nnodes=1, num_local_procs=1, node_rank=0
[2021-10-05 10:05:47,909] [INFO] [launch.py:101:main] global_rank_mapping=defaultdict(<class 'list'>, {'localhost': [0]})
[2021-10-05 10:05:47,909] [INFO] [launch.py:102:main] dist_world_size=1
[2021-10-05 10:05:47,909] [INFO] [launch.py:104:main] Setting CUDA_VISIBLE_DEVICES=0
[2021-10-05 10:05:52,955] [INFO] [distributed.py:46:init_distributed] Initializing torch distributed with backend: nccl
Downloading https://www.cs.toronto.edu/~kriz/cifar-10-python.tar.gz to ./data/cifar-10-python.tar.gz
170499072it [00:01, 111861266.78it/s]
Extracting ./data/cifar-10-python.tar.gz to ./data
[W ProcessGroupNCCL.cpp:1569] Rank 0 using best-guess GPU 0 to perform barrier as devices used by this process are currently unknown. This can potentially cause a hang if this rank to GPU mapp
ing is incorrect.Specify device_ids in barrier() to force use of a particular device.
Files already downloaded and verified
 frog  frog  deer horse
[2021-10-05 10:06:09,109] [INFO] [logging.py:68:log_dist] [Rank 0] DeepSpeed info: version=0.5.3, git-hash=unknown, git-branch=unknown
[2021-10-05 10:06:09,126] [INFO] [logging.py:68:log_dist] [Rank 0] initializing deepspeed groups
[2021-10-05 10:06:09,126] [INFO] [logging.py:68:log_dist] [Rank 0] initializing deepspeed model parallel group with size 1
[2021-10-05 10:06:09,128] [INFO] [logging.py:68:log_dist] [Rank 0] initializing deepspeed expert parallel group with size 1
[2021-10-05 10:06:09,128] [INFO] [logging.py:68:log_dist] [Rank 0] creating expert data parallel process group with ranks: [0]
[2021-10-05 10:06:09,128] [INFO] [logging.py:68:log_dist] [Rank 0] creating expert parallel process group with ranks: [0]
[2021-10-05 10:06:09,236] [INFO] [engine.py:197:__init__] DeepSpeed Flops Profiler Enabled: False
Installed CUDA version 11.2 does not match the version torch was compiled with 11.1 but since the APIs are compatible, accepting this combination
Using /gpfs/mira-home/zhen/.cache/torch_extensions as PyTorch extensions root...
Detected CUDA files, patching ldflags
Emitting ninja build file /gpfs/mira-home/zhen/.cache/torch_extensions/fused_adam/build.ninja...
Building extension module fused_adam...
Allowing ninja to set a default number of workers... (overridable by setting the environment variable MAX_JOBS=N)
ninja: no work to do.
Loading extension module fused_adam...
Time to load fused_adam op: 0.48449182510375977 seconds
[2021-10-05 10:06:11,552] [INFO] [engine.py:821:_configure_optimizer] Using DeepSpeed Optimizer param name adam as basic optimizer
```

Argonne
NATIONAL LABORATORY

# Future plans on DeepSpeed

- **More test cases: Megatron-LM, 1Cycle, and BERT model.**

- **More test cases on multiple nodes.**

- **More performance comparison with Horovod**

# Thanks

**https://github.com/Argonne-lcf/sdl_ai_workshop/tree/master/01_distributedDeepLearning/DeepSpeed**