

Homework Assignment 5

Part A

All Part A experiments (DDP training, timing, multi-GPU runs, and communication analysis) were implemented in the submitted file: `Jmp10051_Lab5_partA.py`

Q1

Epoch 1

```
Host: b-03-14
[jmp10051@b-33-14 ~]$ python -m torch.distributed.run --nproc_per_node=1 lab5_partA.py --batch-size 32 --epochs 1
100% | 170M/170M [00:02<00:00, 70.7MB/s]
/usr/lib64/python3.9/tarfile.py:2239: RuntimeWarning: The default behavior of tarfile extraction has been changed to disallow common exploits (including CVE-2007-4559). By default, absolute/parent paths are disallowed and some mode bits are cleared. See https://access.redhat.com/articles/7004769 for more details.
warnings.warn(
[Epoch 1/1] time=26.428s loss=1.8491 acc=0.3268

=== Lab 5 Part A Summary ===
[jmp10051@b-33-14 ~]$ python -m torch.distributed.run --nproc_per_node=1 lab5_partA.py --batch-size 128 --epochs 1
[Epoch 1/1] time=18.808s loss=2.0040 acc=0.2909

=== Lab 5 Part A Summary ===
[jmp10051@b-33-14 ~]$ python -m torch.distributed.run --nproc_per_node=1 lab5_partA.py --batch-size 512 --epochs 1
[Epoch 1/1] time=20.215s loss=2.0129 acc=0.2954

=== Lab 5 Part A Summary ===
[jmp10051@b-33-14 ~]$ python -m torch.distributed.run --nproc_per_node=1 lab5_partA.py --batch-size 2048 --epochs 1
[Epoch 1/1] time=21.227s loss=2.5864 acc=0.1744

=== Lab 5 Part A Summary ===
[jmp10051@b-33-14 ~]$
```

Epoch 2

```
[jmp10051@b-33-14 ~]$ cd /scratch/jmp10051/
[jmp10051@b-33-14 jmp10051]$ python -m torch.distributed.run --nproc_per_node=1 lab5_partA.py --batch-size 32 --epochs 2
100% | 170M/170M [00:02<00:00, 77.7MB/s]
/usr/lib64/python3.9/tarfile.py:2239: RuntimeWarning: The default behavior of tarfile extraction has been changed to disallow common exploits (including CVE-2007-4559). By default, absolute/parent paths are disallowed and some mode bits are cleared. See https://access.redhat.com/articles/7004769 for more details.
warnings.warn(
[Epoch 1/2] time=23.846s loss=1.8790 acc=0.3157
[Epoch 2/2] time=23.550s loss=1.4255 acc=0.4790

=== Lab 5 Part A Summary ===
Epoch 2 full time (includes data loading; use for Q2 tables): 23.550 s
Epoch 2 compute-only time (excludes data loading; use for Q1): 22.782 s
[jmp10051@b-33-14 jmp10051]$ python -m torch.distributed.run --nproc_per_node=1 lab5_partA.py --batch-size 128 --epochs 2
[Epoch 1/2] time=18.666s loss=2.1907 acc=0.2399
[Epoch 2/2] time=18.311s loss=1.5637 acc=0.4237

=== Lab 5 Part A Summary ===
Epoch 2 full time (includes data loading; use for Q2 tables): 18.311 s
Epoch 2 compute-only time (excludes data loading; use for Q1): 18.051 s
[jmp10051@b-33-14 jmp10051]$ python -m torch.distributed.run --nproc_per_node=1 lab5_partA.py --batch-size 512 --epochs 2
[Epoch 1/2] time=20.167s loss=2.6589 acc=0.1863
[Epoch 2/2] time=19.704s loss=1.7150 acc=0.3564

=== Lab 5 Part A Summary ===
Epoch 2 full time (includes data loading; use for Q2 tables): 19.704 s
Epoch 2 compute-only time (excludes data loading; use for Q1): 19.448 s
[jmp10051@b-33-14 jmp10051]$ python -m torch.distributed.run --nproc_per_node=1 lab5_partA.py --batch-size 2048 --epochs 2
[Epoch 1/2] time=21.134s loss=4.2121 acc=0.1210
[Epoch 2/2] time=20.714s loss=2.4399 acc=0.1380

=== Lab 5 Part A Summary ===
Epoch 2 full time (includes data loading; use for Q2 tables): 20.714 s
Epoch 2 compute-only time (excludes data loading; use for Q1): 19.658 s
[jmp10051@b-33-14 jmp10051]$
```

The epoch-2 compute-only training time was measured on a single GPU while increasing the batch size by a factor of four (32 → 128 → 512 → 2048). Batch size 8192 resulted in a CUDA out-of-memory error, so 2048 is the largest batch size that fits on one GPU.

Summary Table : (excludes data loading)

Batch Size	Compute-Only Epoch 2 Time (s)
32	22.782
128	18.051
512	19.448
2048	19.658

Batch Size	Compute-Only Epoch 2 Time (s)
8192	Out-of-memory (does not fit)

On a single NVIDIA A40 GPU (g2-standard-48 node on Greene), it takes 22.66 seconds to train one epoch (excluding data-loading) using mini-batch size 32, 18.05 seconds using batch size 128, 19.4 seconds using batch size 512, and 19.6 seconds using batch size 2048. Increasing the batch size from 32 \rightarrow 128 improves throughput because the GPU becomes better utilized and kernel overhead is amortized.

For larger batches (512 and 2048), the epoch time flattens around ~ 19 seconds, indicating that the GPU reaches a compute/memory bandwidth saturation point. Batch size 8192 results in a CUDA out-of-memory error, so 2048 is the largest batch size that fits on a single GPU.

Q2

2 GPUs

```
[jmp10051@b-33-14 jmp10051]$ python -m torch.distributed.run --nproc_per_node=2
lab5_partA.py --batch-size 32 --epochs 2
```

```
[Epoch 1/2] time=16.704s loss=1.8432 acc=0.3226
```

```
[Epoch 2/2] time=16.164s loss=1.3433 acc=0.5115
```

```
=== Lab 5 Part A Summary ===
```

Epoch 2 full time (includes data loading; use for Q2 tables): 16.164 s

Epoch 2 compute-only time (excludes data loading; use for Q1): 15.775 s

```
[jmp10051@b-33-14 jmp10051]$ python -m torch.distributed.run --nproc_per_node=2
lab5_partA.py --batch-size 128 --epochs 2
```

```
[Epoch 1/2] time=10.360s loss=1.9076 acc=0.3232
```

```
[Epoch 2/2] time=9.849s loss=1.4073 acc=0.4828
```

```
=== Lab 5 Part A Summary ===
```

Epoch 2 full time (includes data loading; use for Q2 tables): 9.849 s

Epoch 2 compute-only time (excludes data loading; use for Q1): 9.637 s

```
[jmp10051@b-33-14 jmp10051]$ python -m torch.distributed.run --nproc_per_node=2
lab5_partA.py --batch-size 512 --epochs 2
```

```
[Epoch 1/2] time=10.457s loss=2.9319 acc=0.1794
```

```
[Epoch 2/2] time=10.012s loss=1.8745 acc=0.3021
```

```
=== Lab 5 Part A Summary ===
```

Epoch 2 full time (includes data loading; use for Q2 tables): 10.012 s

Epoch 2 compute-only time (excludes data loading; use for Q1): 9.708 s

4 GPUs

```
[jmp10051@b-33-14 jmp10051]$ python -m torch.distributed.run --nproc_per_node=4
lab5_partA.py --batch-size 32 --epochs 2
```

```
[Epoch 1/2] time=9.623s loss=1.9499 acc=0.2972
```

[Epoch 2/2] time=9.060s loss=1.4588 acc=0.4652

=== Lab 5 Part A Summary ===

Epoch 2 full time (includes data loading; use for Q2 tables): 9.060 s

Epoch 2 compute-only time (excludes data loading; use for Q1): 8.811 s

```
[jmp10051@b-33-14 jmp10051]$ python -m torch.distributed.run --nproc_per_node=4  
lab5_partA.py --batch-size 128 --epochs 2
```

[Epoch 1/2] time=5.586s loss=2.0906 acc=0.2740

[Epoch 2/2] time=5.023s loss=1.5321 acc=0.4318

=== Lab 5 Part A Summary ===

Epoch 2 full time (includes data loading; use for Q2 tables): 5.023 s

Epoch 2 compute-only time (excludes data loading; use for Q1): 4.852 s

```
[jmp10051@b-33-14 jmp10051]$ python -m torch.distributed.run --nproc_per_node=4  
lab5_partA.py --batch-size 512 --epochs 2
```

[Epoch 1/2] time=5.707s loss=2.6404 acc=0.1935

[Epoch 2/2] time=5.160s loss=1.7537 acc=0.3244

=== Lab 5 Part A Summary ===

Epoch 2 full time (includes data loading; use for Q2 tables): 5.160 s

Epoch 2 compute-only time (excludes data loading; use for Q1): 4.861 s

Table 1: Training Time and Speedup (2 vs 4 GPUs)

Batch Size (per GPU)	2 GPUs Time (s)	4 GPUs Time (s)	Speedup (2→4 GPUs)
32	16.164	9.060	1.78×
128	9.849	5.023	1.96×
512	10.012	5.160	1.94×

Times used are “Epoch 2 full time.” Speedup is computed as T_2/T_4 .

Scaling Analysis

1. Type of scaling measured

This is **strong scaling**, because the global dataset size for an epoch is fixed while the number of GPUs increases.

2. Whether weak scaling would be better or worse

Weak scaling would show **better** speedup because each GPU would maintain the same amount of work, keeping computation dominant relative to communication.

3. Data point supporting the argument

Speedup improves when each GPU has more work.

Example: batch size 32 achieves 1.78× speedup, while batch size 128 achieves 1.96×

The larger workload per GPU (128) maintains efficiency closer to ideal scaling, which is what weak scaling preserves.

Relevant Screenshots:

```

Host: b-33-14
Theme: Default
Epoch 2 compute-only time (excludes data loading; use for Q1): 15.775 s
[jmp100510b-33-14 jmp10051s] python -m torch.distributed.run --nproc_per_node=2 lab5_partA.py --batch-size 128 --epochs 2

=====
Setting OMP_NUM_THREADS environment variable for each process to be 1 in default, to avoid your system being overloaded, please further tune the variable for optimal performance in your application as needed.
=====
/home/jmp10051/.local/lib/python3.9/site-packages/torch/distributed/distributed_c10d.py:4807: UserWarning: No device id is provided via 'init_process_group' or 'barrier '. Using the current device set by the user.
warnings.warn( # warn only once
[rank1]: [W1128 19:31:22.12891156 ProcessGroupMCCCL.cpp:5623] [PG ID 0 PG GUID 0 Rank 1] using GPU 1 as device used by this process is currently unknown. This can potentially cause a hang if this rank to GPU mapping is incorrect. You can specify device_id in init_process_group() to force use of a particular device.
/home/jmp10051/.local/lib/python3.9/site-packages/torch/distributed/distributed_c10d.py:4807: UserWarning: No device id is provided via 'init_process_group' or 'barrier '. Using the current device set by the user.
warnings.warn( # warn only once
[rank0]: [W1128 19:31:22.47383414 ProcessGroupMCCCL.cpp:5623] [PG ID 0 PG GUID 0 Rank 0] using GPU 0 as device used by this process is currently unknown. This can potentially cause a hang if this rank to GPU mapping is incorrect. You can specify device_id in init_process_group() to force use of a particular device.
Epoch 1/2] time=10.360s loss=1.9076 acc=0.3232
Epoch 2/2] time=9.849s loss=1.4073 acc=0.4828

==== Lab 5 Part A Summary ====
Epoch 2 full time (includes data loading; use for Q2 tables): 9.849 s
Epoch 2 compute-only time (excludes data loading; use for Q1): 9.637 s
[jmp100510b-33-14 jmp10051s]
[jmp100510b-33-14 jmp10051s]
[jmp100510b-33-14 jmp10051s]
[jmp100510b-33-14 jmp10051s]
[jmp100510b-33-14 jmp10051s]
[jmp100510b-33-14 jmp10051s]
[jmp100510b-33-14 jmp10051s] python -m torch.distributed.run --nproc_per_node=2 lab5_partA.py --batch-size 512 --epochs 2

=====
Setting OMP_NUM_THREADS environment variable for each process to be 1 in default, to avoid your system being overloaded, please further tune the variable for optimal performance in your application as needed.
=====
/home/jmp10051/.local/lib/python3.9/site-packages/torch/distributed/distributed_c10d.py:4807: UserWarning: No device id is provided via 'init_process_group' or 'barrier '. Using the current device set by the user.
warnings.warn( # warn only once
[rank1]: [W1128 19:33:12.48939831 ProcessGroupMCCCL.cpp:5623] [PG ID 0 PG GUID 0 Rank 1] using GPU 1 as device used by this process is currently unknown. This can potentially cause a hang if this rank to GPU mapping is incorrect. You can specify device_id in init_process_group() to force use of a particular device.
/home/jmp10051/.local/lib/python3.9/site-packages/torch/distributed/distributed_c10d.py:4807: UserWarning: No device id is provided via 'init_process_group' or 'barrier '. Using the current device set by the user.
warnings.warn( # warn only once
[rank1]: [W1128 19:33:12.48939831 ProcessGroupMCCCL.cpp:5623] [PG ID 0 PG GUID 0 Rank 0] using GPU 0 as device used by this process is currently unknown. This can potentially cause a hang if this rank to GPU mapping is incorrect. You can specify device_id in init_process_group() to force use of a particular device.
Epoch 1/2] time=10.457s loss=2.9319 acc=0.1794
Epoch 2/2] time=10.012s loss=1.8745 acc=0.3021

==== Lab 5 Part A Summary ====
Epoch 2 full time (includes data loading; use for Q2 tables): 10.012 s
Epoch 2 compute-only time (excludes data loading; use for Q1): 9.708 s
[jmp100510b-33-14 jmp10051s]
[jmp100510b-33-14 jmp10051s]
[jmp100510b-33-14 jmp10051s]
[jmp100510b-33-14 jmp10051s]
[jmp100510b-33-14 jmp10051s]
[jmp100510b-33-14 jmp10051s]
[jmp100510b-33-14 jmp10051s] python -m torch.distributed.run --nproc_per_node=4 lab5_partA.py --batch-size 32 --epochs 2

=====
Setting OMP_NUM_THREADS environment variable for each process to be 1 in default, to avoid your system being overloaded, please further tune the variable for optimal performance in your application as needed.
=====
/home/jmp10051/.local/lib/python3.9/site-packages/torch/distributed/distributed_c10d.py:4807: UserWarning: No device id is provided via 'init_process_group' or 'barrier '. Using the current device set by the user.
warnings.warn( # warn only once
[rank1]: [W1128 19:34:26.3734263734 ProcessGroupMCCCL.cpp:5623] [PG ID 0 PG GUID 0 Rank 2] using GPU 2 as device used by this process is currently unknown. This can potentially cause a hang if this rank to GPU mapping is incorrect. You can specify device_id in init_process_group() to force use of a particular device.
/home/jmp10051/.local/lib/python3.9/site-packages/torch/distributed/distributed_c10d.py:4807: UserWarning: No device id is provided via 'init_process_group' or 'barrier '. Using the current device set by the user.
[rank3]: [W1128 19:34:26.375238241 ProcessGroupMCCCL.cpp:5623] [PG ID 0 PG GUID 0 Rank 3] using GPU 3 as device used by this process is currently unknown. This can potentially cause a hang if this rank to GPU mapping is incorrect. You can specify device_id in init_process_group() to force use of a particular device.
/home/jmp10051/.local/lib/python3.9/site-packages/torch/distributed/distributed_c10d.py:4807: UserWarning: No device id is provided via 'init_process_group' or 'barrier '. Using the current device set by the user.
warnings.warn( # warn only once
[rank1]: [W1128 19:34:26.375238241 ProcessGroupMCCCL.cpp:5623] [PG ID 0 PG GUID 0 Rank 1] using GPU 1 as device used by this process is currently unknown. This can potentially cause a hang if this rank to GPU mapping is incorrect. You can specify device_id in init_process_group() to force use of a particular device.
/home/jmp10051/.local/lib/python3.9/site-packages/torch/distributed/distributed_c10d.py:4807: UserWarning: No device id is provided via 'init_process_group' or 'barrier '. Using the current device set by the user.
warnings.warn( # warn only once

Host: b-33-14
Theme: Default
is incorrect. You can specify device_id in init_process_group() to force use of a particular device.
/home/jmp10051/.local/lib/python3.9/site-packages/torch/distributed/distributed_c10d.py:4807: UserWarning: No device id is provided via 'init_process_group' or 'barrier '. Using the current device set by the user.
warnings.warn( # warn only once
[rank1]: [W1128 19:35:18.140250236 ProcessGroupMCCCL.cpp:5623] [PG ID 0 PG GUID 0 Rank 0] using GPU 0 as device used by this process is currently unknown. This can potentially cause a hang if this rank to GPU mapping is incorrect. You can specify device_id in init_process_group() to force use of a particular device.
Epoch 1/2] time=9.623s loss=1.9499 acc=0.2972
Epoch 2/2] time=9.068s loss=1.4588 acc=0.4652

==== Lab 5 Part A Summary ====
Epoch 2 full time (includes data loading; use for Q2 tables): 9.068 s
Epoch 2 compute-only time (excludes data loading; use for Q1): 8.811 s
[jmp100510b-33-14 jmp10051s] python -m torch.distributed.run --nproc_per_node=4 lab5_partA.py --batch-size 128 --epochs 2

=====
Setting OMP_NUM_THREADS environment variable for each process to be 1 in default, to avoid your system being overloaded, please further tune the variable for optimal performance in your application as needed.
=====
/home/jmp10051/.local/lib/python3.9/site-packages/torch/distributed/distributed_c10d.py:4807: UserWarning: No device id is provided via 'init_process_group' or 'barrier '. Using the current device set by the user.
warnings.warn( # warn only once
[rank1]: [W1128 19:35:18.140250236 ProcessGroupMCCCL.cpp:5623] [PG ID 0 PG GUID 0 Rank 1] using GPU 1 as device used by this process is currently unknown. This can potentially cause a hang if this rank to GPU mapping is incorrect. You can specify device_id in init_process_group() to force use of a particular device.
/home/jmp10051/.local/lib/python3.9/site-packages/torch/distributed/distributed_c10d.py:4807: UserWarning: No device id is provided via 'init_process_group' or 'barrier '. Using the current device set by the user.
warnings.warn( # warn only once
[rank1]: [W1128 19:35:18.141755413 ProcessGroupMCCCL.cpp:5623] [PG ID 0 PG GUID 0 Rank 1] using GPU 1 as device used by this process is currently unknown. This can potentially cause a hang if this rank to GPU mapping is incorrect. You can specify device_id in init_process_group() to force use of a particular device.
/home/jmp10051/.local/lib/python3.9/site-packages/torch/distributed/distributed_c10d.py:4807: UserWarning: No device id is provided via 'init_process_group' or 'barrier '. Using the current device set by the user.
[rank0]: [W1128 19:35:18.474968471 ProcessGroupMCCCL.cpp:5623] [PG ID 0 PG GUID 0 Rank 0] using GPU 0 as device used by this process is currently unknown. This can potentially cause a hang if this rank to GPU mapping is incorrect. You can specify device_id in init_process_group() to force use of a particular device.
Epoch 1/2] time=5.586s loss=2.6960 acc=0.2140
Epoch 2/2] time=5.028s loss=1.5231 acc=0.4318

==== Lab 5 Part A Summary ====
Epoch 2 full time (includes data loading; use for Q2 tables): 5.023 s
Epoch 2 compute-only time (excludes data loading; use for Q1): 4.852 s
[jmp100510b-33-14 jmp10051s]
[jmp100510b-33-14 jmp10051s]
[jmp100510b-33-14 jmp10051s]
[jmp100510b-33-14 jmp10051s] python -m torch.distributed.run --nproc_per_node=4 lab5_partA.py --batch-size 512 --epochs 2

=====
Setting OMP_NUM_THREADS environment variable for each process to be 1 in default, to avoid your system being overloaded, please further tune the variable for optimal performance in your application as needed.
=====
/home/jmp10051/.local/lib/python3.9/site-packages/torch/distributed/distributed_c10d.py:4807: UserWarning: No device id is provided via 'init_process_group' or 'barrier '. Using the current device set by the user.
warnings.warn( # warn only once
[rank1]: [W1128 19:36:11.573598882 ProcessGroupMCCCL.cpp:5623] [PG ID 0 PG GUID 0 Rank 1] using GPU 1 as device used by this process is currently unknown. This can potentially cause a hang if this rank to GPU mapping is incorrect. You can specify device_id in init_process_group() to force use of a particular device.
/home/jmp10051/.local/lib/python3.9/site-packages/torch/distributed/distributed_c10d.py:4807: UserWarning: No device id is provided via 'init_process_group' or 'barrier '. Using the current device set by the user.
warnings.warn( # warn only once
[rank2]: [W1128 19:36:11.576672588 ProcessGroupMCCCL.cpp:5623] [PG ID 0 PG GUID 0 Rank 2] using GPU 2 as device used by this process is currently unknown. This can potentially cause a hang if this rank to GPU mapping is incorrect. You can specify device_id in init_process_group() to force use of a particular device.
/home/jmp10051/.local/lib/python3.9/site-packages/torch/distributed/distributed_c10d.py:4807: UserWarning: No device id is provided via 'init_process_group' or 'barrier '. Using the current device set by the user.
[rank1]: [W1128 19:36:12.912437444 ProcessGroupMCCCL.cpp:5623] [PG ID 0 PG GUID 0 Rank 0] using GPU 0 as device used by this process is currently unknown. This can potentially cause a hang if this rank to GPU mapping is incorrect. You can specify device_id in init_process_group() to force use of a particular device.
Epoch 1/2] time=5.707s loss=2.6484 acc=0.1935
Epoch 2/2] time=5.108s loss=1.7537 acc=0.3244

==== Lab 5 Part A Summary ====
Epoch 2 full time (includes data loading; use for Q2 tables): 5.168 s
Epoch 2 compute-only time (excludes data loading; use for Q1): 4.861 s
[jmp100510b-33-14 jmp10051s]

```

Q3

Q3.1: How much time spent in computation and communication

Methodology:

To distinguish between computation and communication time, we apply the Strong Scaling assumption (fixed total problem size).

1. Computation Time (T_{comp}):

We estimate the ideal computation time by taking the measured Single-GPU (Epoch 2) time and dividing it by the number of GPUs.

$$T_{\text{comp}} = \frac{T_{1\text{GPU}}}{N}$$

2. Communication Time:

We subtract the estimated computation time from the actual measured execution time.

$$T_{\text{comm}} = T_{\text{actual}} - T_{\text{comp}}$$

Table 2: Computation vs Communication Time

Batch Size 32

Setup	Total Time	Computation	Communication
2 GPUs	16.164 s	11.775 s	4.389 s
4 GPUs	9.060 s	5.888 s	3.172 s

Batch Size 128

Setup	Total Time	Computation	Communication
2 GPUs	9.849 s	9.156 s	0.693 s
4 GPUs	5.023 s	4.578 s	0.445 s

Batch Size 512

Setup	Total Time	Computation	Communication
2 GPUs	10.012 s	9.882 s	0.130 s
4 GPUs	5.160 s	4.941 s	0.219 s

(Single-GPU times used: 23.550 s for B32, 18.311 s for B128, 19.764 s for B512.)

Q3.2: Communication Bandwidth Utilization

1. AllReduce Time Formula (Ring):

$$T_{AllReduce} = \frac{2M(N - 1)}{NB}$$

Where:

- $M = 44,695,848$ bytes
- N : GPU count
- $B = 0.9 \times 10^9$ bytes/s

2. Bandwidth Utilization Formula:

$$Utilization = \frac{\left(\frac{2M(N - 1)}{NT_{comm}} \right)}{B}$$

3. Table 3: Bandwidth Utilization

Batch Size 32

Setup	Comm Time	Effective BW	Utilization
2 GPU _s	4.389 s	0.010 GB/s	1.13%
4 GPU _s	3.172 s	0.021 GB/s	2.35%

Batch Size 128

Setup	Comm Time	Effective BW	Utilization
2 GPU _s	0.693 s	0.064 GB/s	7.17%
4 GPU _s	0.445 s	0.151 GB/s	16.74%

Batch Size 512

Setup	Comm Time	Effective BW	Utilization
2 GPU _s	0.130 s	0.344 GB/s	38.20%
4 GPU _s	0.219 s	0.306 GB/s	34.02%

Q4

```
[jml10051@b-33-14 jml10051]$  
[jml10051@b-33-14 jml10051]$ python -m torch.distributed.run --nproc_per_node=4 lab5_partA.py --batch-size 2048 --epochs 5  
  
*****  
Setting OMP_NUM_THREADS environment variable for each process to be 1 in default, to avoid your system being overloaded, please further tune the variable for optimal performance in your application as needed.  
*****  
/home/jml10051/.local/lib/python3.9/site-packages/torch/distributed/distributed_c10d.py:4807: UserWarning: No device id is provided via 'init_process_group' or 'barrier'. Using the current device set by the user.  
warnings.warn( # warn only once  
/home/jml10051/.local/lib/python3.9/site-packages/torch/distributed/distributed_c10d.py:4807: UserWarning: No device id is provided via 'init_process_group' or 'barrier'. Using the current device set by the user.  
warnings.warn( # warn only once  
[rank2]: [W1128 19:39:11.923537453 ProcessGroupNCCL.cpp:5023] [PG ID 0 PG GUID 0 Rank 2] using GPU 2 as device used by this process is currently unknown. This can potentially cause a hang if this rank to GPU mapping  
is incorrect. You can specify device_id in init_process_group() to force use of a particular device.  
[rank1]: [W1128 19:39:11.023885339 ProcessGroupNCCL.cpp:5023] [PG ID 0 PG GUID 0 Rank 1] using GPU 1 as device used by this process is currently unknown. This can potentially cause a hang if this rank to GPU mapping  
is incorrect. You can specify device_id in init_process_group() to force use of a particular device.  
/home/jml10051/.local/lib/python3.9/site-packages/torch/distributed/distributed_c10d.py:4807: UserWarning: No device id is provided via 'init_process_group' or 'barrier'. Using the current device set by the user.  
warnings.warn( # warn only once  
[rank3]: [W1128 19:39:11.029792375 ProcessGroupNCCL.cpp:5023] [PG ID 0 PG GUID 0 Rank 3] using GPU 3 as device used by this process is currently unknown. This can potentially cause a hang if this rank to GPU mapping  
is incorrect. You can specify device_id in init_process_group() to force use of a particular device.  
/home/jml10051/.local/lib/python3.9/site-packages/torch/distributed/distributed_c10d.py:4807: UserWarning: No device id is provided via 'init_process_group' or 'barrier'. Using the current device set by the user.  
warnings.warn( # warn only once  
[rank0]: [W1128 19:39:11.365749320 ProcessGroupNCCL.cpp:5023] [PG ID 0 PG GUID 0 Rank 0] using GPU 0 as device used by this process is currently unknown. This can potentially cause a hang if this rank to GPU mapping  
is incorrect. You can specify device_id in init_process_group() to force use of a particular device.  
Epoch 1/5: time=6.520s loss=3.8283 acc=0.1474  
Epoch 2/5: time=5.976s loss=3.1221 acc=0.1354  
Epoch 3/5: time=5.976s loss=2.1794 acc=0.2865  
Epoch 4/5: time=6.022s loss=1.9666 acc=0.2511  
Epoch 5/5: time=6.019s loss=1.8399 acc=0.2899  
  
== Lab 5 Part A Summary ==  
Epoch 2 full time (includes data loading; use for Q2 tables): 5.976 s  
Epoch 2 compute-only time (excludes data loading; use for Q1): 4.907 s  
Epoch 5 loss/acc (Q4.1): 1.8399 / 0.2899  
[jml10051@b-33-14 jml10051]$
```

Q4.1: Accuracy When Using a Large Batch

Experimental Setup (4 GPUs, Batch Size 2048 per GPU):

- Batch Size per GPU: 2048
- Total Effective Batch Size: 8192
- Epoch 5 Average Training Loss: 1.8399
- Epoch 5 Training Accuracy: 0.2899 (28.99%)

Comparison with Lab 2 (Batch Size 128):

The training accuracy with this large-batch configuration (28.99% at Epoch 5) is much lower than the accuracy typically observed in Lab 2 with batch size 128 (often above 45% by Epoch 5).

Reason:

A very large batch reduces the number of parameter updates per epoch. With no hyperparameter adjustments, the model converges more slowly and is more likely to fall into sharper minima, creating a generalization gap and resulting in lower accuracy.

Q4.2: How to Improve Training Accuracy with a Large Batch

According to *Accurate, Large Minibatch SGD: Training ImageNet in 1 Hour* [4], the following two methods address the degradation observed in Q4.1.

Remedy 1: Linear scaling of the learning rate

When increasing the minibatch by factor k , scale the learning rate by the same factor:

$$\text{LR}_{\text{new}} = \text{LR}_{\text{base}} \times \frac{\text{BatchSize}_{\text{new}}}{\text{BatchSize}_{\text{base}}}$$

Concrete example: base total batch = 128, new total batch = 8192 $\Rightarrow k = 64$.

If $\text{LR}_{\text{base}} = 0.1$, set $\text{LR}_{\text{new}} = 0.1 \times 64 = 6.4$.

Remedy 2: Gradual warmup of the learning rate

Avoid starting with the full scaled LR immediately. Linearly increase the LR from a small initial value to the scaled LR over a short warmup period (e.g., first 3–5 epochs).

Linear warmup per epoch (or per step):

$$\text{LR}(t) = \text{LR}_{\text{start}} + \frac{t}{T_{\text{warmup}}} (\text{LR}_{\text{new}} - \text{LR}_{\text{start}}), 0 \leq t \leq T_{\text{warmup}}$$

Choose LR_{start} small (e.g., LR_{base} or $\text{LR}_{\text{base}}/10$), and $T_{\text{warmup}} = 3\text{--}5$ epochs.

Q5

No, in addition to gradients, `DistributedDataParallel` also communicates model buffers such as `BatchNorm` running mean and running variance.

These buffers appear in C7 of Lab 2 (e.g., `running_mean`, `running_var`) and are synchronized across GPUs so that all replicas use consistent normalization statistics. This buffer broadcast happens every iteration (unless explicitly disabled), so gradients are not the only tensors sent over the network.

Q6

Yes, for batch size per GPU 512 it is generally sufficient to communicate only gradients across the 4 GPUs, because each GPU already sees a very large local mini-batch, so its `BatchNorm` statistics are well estimated from local data.

As discussed in [4], large-minibatch ImageNet training can keep `BatchNorm` local (no cross-GPU BN synchronization) and still match small-batch accuracy, so gradient all-reduce is the only communication that is really required for correct SSGD in this regime.

Not synchronizing BN statistics in this setting mainly saves communication overhead; the difference between local BN statistics on each GPU becomes small when each GPU's batch size is as large as 512.

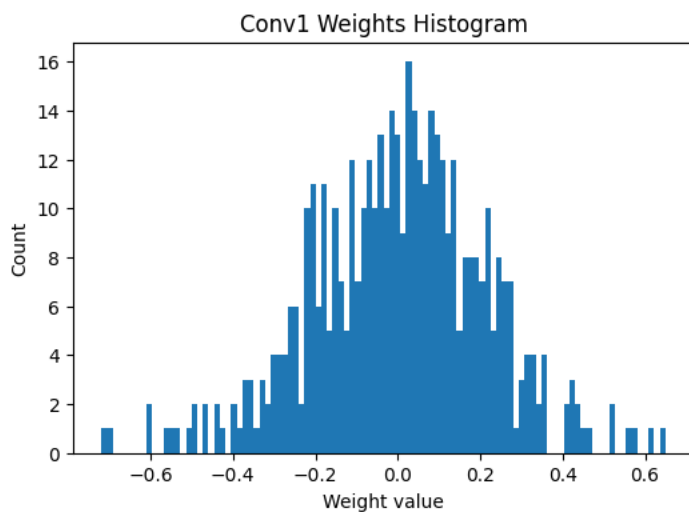
Part-B

All quantization steps, histograms, and evaluation results for Part B were generated in the submitted notebook: `Jmp10051_HW5_PartB.ipynb`

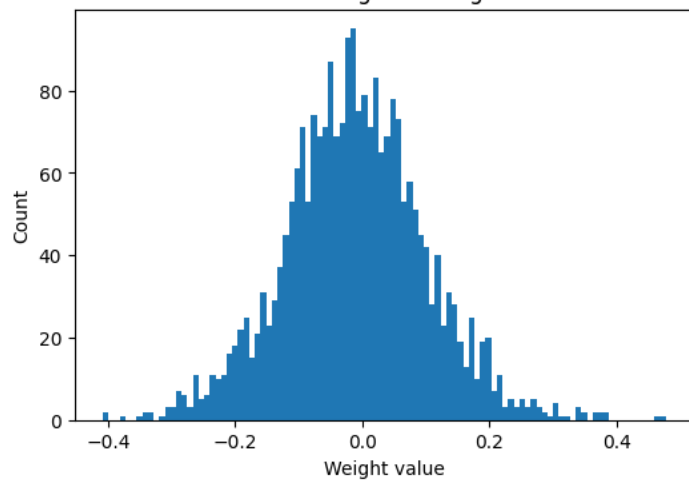
Q1

```
1 # --- Conv1 ---
2 conv1_w = net.conv1.weight.detach().cpu().numpy().flatten()
3 plt.figure(figsize=(6,4))
4 plt.hist(conv1_w, bins=100)
5 plt.title("Conv1 Weights Histogram")
6 plt.xlabel("Weight value")
7 plt.ylabel("Count")
8 plt.show()
9
10 # --- Conv2 ---
11 conv2_w = net.conv2.weight.detach().cpu().numpy().flatten()
12 plt.figure(figsize=(6,4))
13 plt.hist(conv2_w, bins=100)
14 plt.title("Conv2 Weights Histogram")
15 plt.xlabel("Weight value")
16 plt.ylabel("Count")
17 plt.show()
18
19 # --- FC1 ---
20 fc1_w = net.fc1.weight.detach().cpu().numpy().flatten()
21 plt.figure(figsize=(6,4))
22 plt.hist(fc1_w, bins=100)
23 plt.title("FC1 Weights Histogram")
24 plt.xlabel("Weight value")
25 plt.ylabel("Count")
26 plt.show()
27
28 # --- FC2 ---
29 fc2_w = net.fc2.weight.detach().cpu().numpy().flatten()
30 plt.figure(figsize=(6,4))
31 plt.hist(fc2_w, bins=100)
32 plt.title("FC2 Weights Histogram")
33 plt.xlabel("Weight value")
34 plt.ylabel("Count")
35 plt.show()
36
37 # --- FC3 ---
38 fc3_w = net.fc3.weight.detach().cpu().numpy().flatten()
39 plt.figure(figsize=(6,4))
40 plt.hist(fc3_w, bins=100)
41 plt.title("FC3 Weights Histogram")
42 plt.xlabel("Weight value")
43 plt.ylabel("Count")
44 plt.show()
45
```

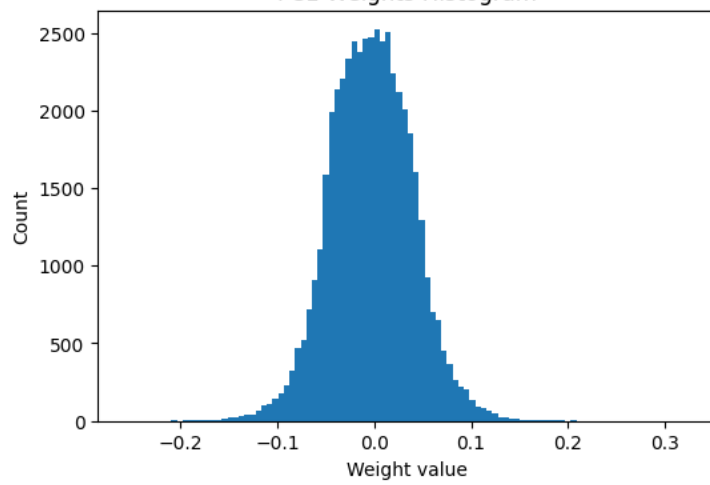
Plots:



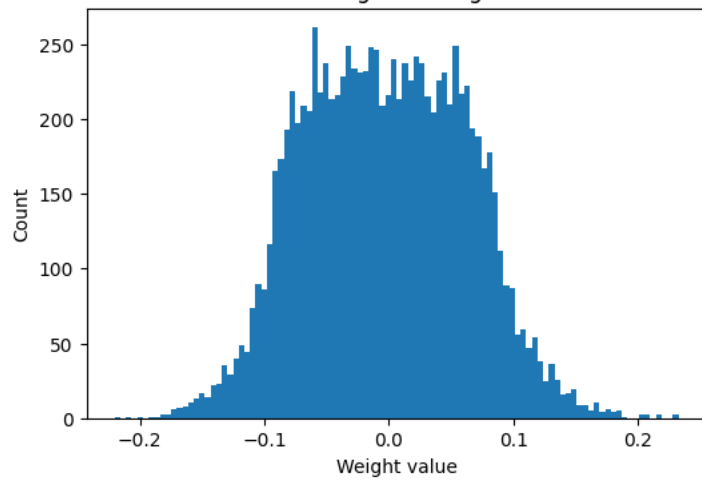
Conv2 Weights Histogram

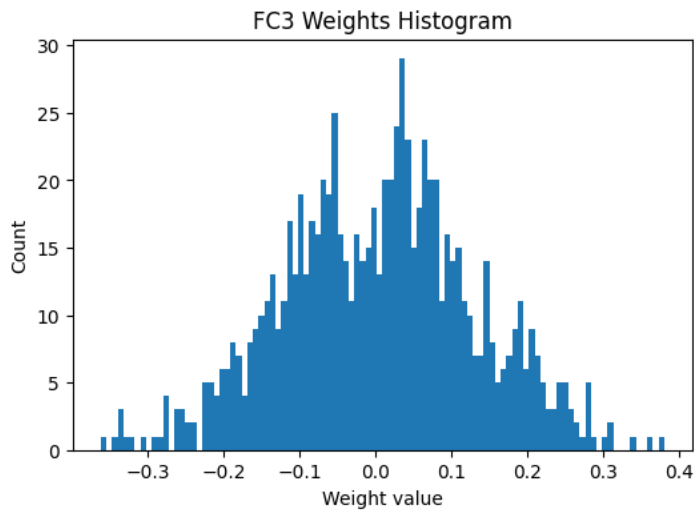


FC1 Weights Histogram



FC2 Weights Histogram





After plotting the weight histograms for all layers, I noticed that the weights are centered around zero and follow a smooth, bell-shaped distribution. Both convolutional and fully connected layers stay within a small range, with only a few outliers. This kind of symmetric, narrow distribution suggests that the network does not depend on extremely large weights and should tolerate uniform int8 quantization quite well. Overall, the plots show that a simple per-layer symmetric scale would be appropriate.

Q2

```

1 from typing import Tuple
2
3 def quantized_weights(weights: torch.Tensor) -> Tuple[torch.Tensor, float]:
4     """
5     Quantize the weights so that all values are integers between -128 and 127.
6     You may want to use the total range, 3-sigma range, or some other range when
7     deciding just what factors to scale the float32 values by.
8
9     Parameters:
10    weights (Tensor): The unquantized weights
11
12    Returns:
13    (Tensor, float): A tuple with the following elements:
14        * The weights in quantized form, where every value is an integer between -128 and 127.
15          The "dtype" will still be "float", but the values themselves should all be integers.
16        * The scaling factor that your weights were multiplied by.
17          This value does not need to be an 8-bit integer.
18    """
19
20    # ADD YOUR CODE HERE
21    # max : 127, min : -128
22    max_val = weights.abs().max()
23    if max_val == 0:
24        scale = 1.0
25    else:
26        scale = 127.0 / max_val
27
28    quantized = (weights * scale).round()
29    quantized = torch.clamp(quantized, min=-128, max=127)
30
31    return quantized, scale

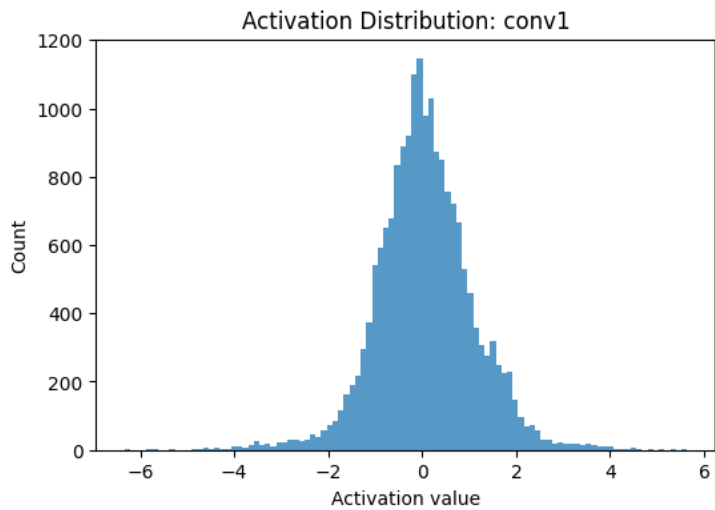
```

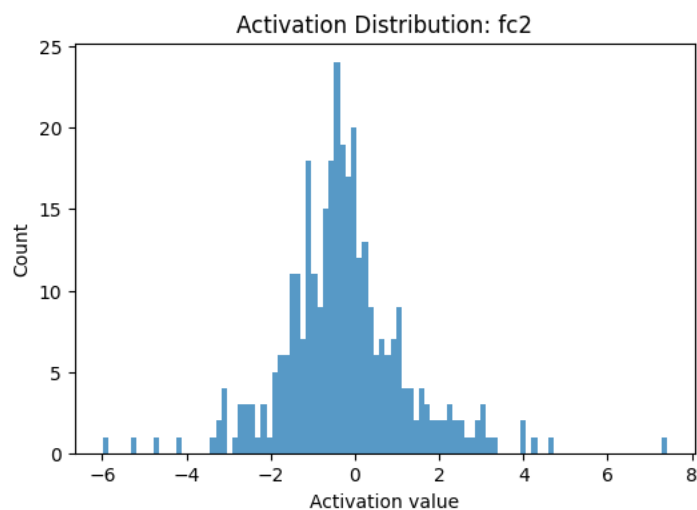
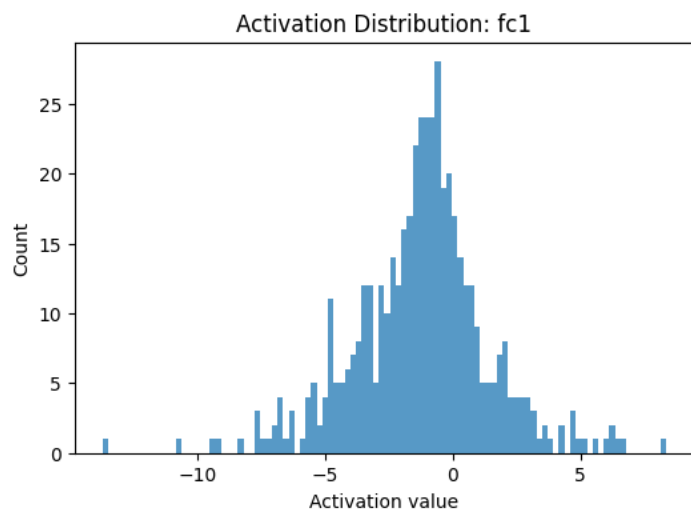
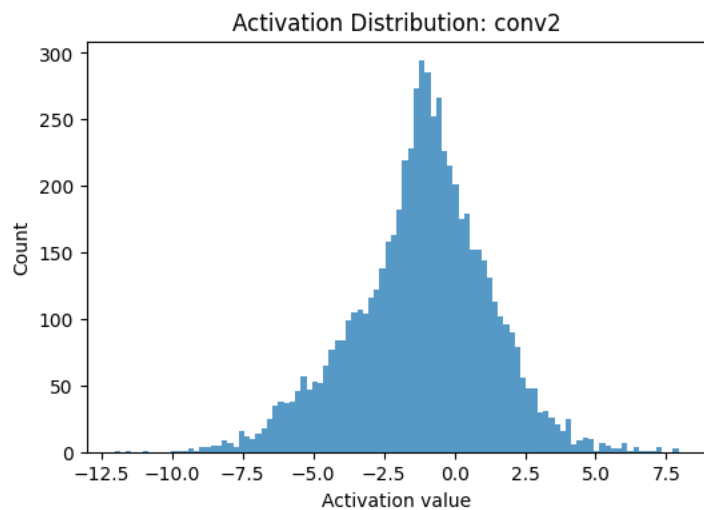
I quantized all convolutional and fully connected weights by scaling them so that the largest absolute weight maps to 127, then rounding and clamping everything into the int8 range. After applying this quantization to a copy of the model, the accuracy became 54.1%, slightly below the original model's accuracy. The model still performs reasonably well, which shows that it can tolerate 8-bit weight precision without significant degradation.

Q3

```
11 # Dictionary to hold activations
12 activations = {}
13 # Hook function
14 def save_activation(name):
15     def hook(module, input, output):
16         activations[name] = output.detach().cpu().view(-1).numpy()
17     return hook
18
19 # Attach hooks to a few key layers
20 handles = []
21 handles.append(net.conv1.register_forward_hook(save_activation("conv1")))
22 handles.append(net.conv2.register_forward_hook(save_activation("conv2")))
23 handles.append(net.fc1.register_forward_hook(save_activation("fc1")))
24 handles.append(net.fc2.register_forward_hook(save_activation("fc2")))
25
26 # Run one batch to collect activations
27 images, labels = next(iter(testloader))
28 images = images.to(device)
29
30 with torch.no_grad():
31     _ = net(images)
32
33 # Plot histograms
34 for name, act in activations.items():
35     plt.figure(figsize=(6,4))
36     plt.hist(act, bins=100, alpha=0.75)
37     plt.title(f"Activation Distribution: {name}")
38     plt.xlabel("Activation value")
39     plt.ylabel("Count")
40     plt.show()
41
42 # Remove hooks
43 for h in handles:
44     h.remove()
45
```

Plots:





I added forward hooks to several layers and plotted the resulting histograms. The activations before ReLU were centered near zero with a roughly Gaussian shape, while the post-ReLU activations became strictly non-negative and heavily concentrated near zero. Overall, the activation ranges were moderate, which is helpful for choosing scaling factors during activation quantization later.

Q4

```
@staticmethod
def quantize_activations(activations: np.ndarray, n_w: float, n_initial_input: float, ns: List[Tuple[float, float]]) -> float:
    """
    Calculate a scaling factor to multiply the output of a layer by.

    Parameters:
    activations (ndarray): The values of all the pixels which have been output by this layer during training
    n_w (float): The scale by which the weights of this layer were multiplied as part of the "quantize_weights" function you wrote earlier
    n_initial_input (float): The scale by which the initial input to the neural network was multiplied
    ns (List[Tuple[float, float]]): A list of tuples, where each tuple represents the "weight scale" and "output scale" (in that order) for every preceding layer

    Returns:
    float: A scaling factor that the layer output should be multiplied by before being fed into the first layer.
    This value does not need to be an 8-bit integer.
    """
    abs_vals = np.abs(activations)
    maximum = np.max(abs_vals)
    # Start from the initial input scale
    scale = n_initial_input
    # Loop through pairs of scales for previous layers
    for weights_scale, output_scale in ns:
        scale = weights_scale * output_scale * scale
    # Scale calculation
    scale = 255.0 / (maximum * scale * n_w)
    return scale

def forward(self, x: torch.Tensor) -> torch.Tensor:
    """
    You can access the output activation scales like this:
    # fc1_output_scale = self.fc1.output_scale

    # To make sure that the outputs of each layer are integers between -128 and 127, you may need to use the following functions:
    # x = torch.Tensor.round
    # x = torch.clamp

    # ADD YOUR CODE HERE
    minimum = -128
    maximum = 127

    # quantize initial input
    x = (x * self.input_scale).round()
    x = torch.clamp(x, min=minimum, max=maximum)

    # ---- Conv1 ----
    x = self.conv1(x)
    x = (x * self.conv1.output_scale).round()
    x = torch.clamp(x, minimum, maximum)
    x = F.relu(x)
    x = self.pool(x)

    # ---- Conv2 ----
    x = self.conv2(x)
    x = (x * self.conv2.output_scale).round()
    x = torch.clamp(x, minimum, maximum)
    x = F.relu(x)
    x = self.pool(x)

    # flatten
    x = torch.flatten(x, start_dim=1)

    #L to chat, %K to generate
    # ---- FC1 ----
    x = self.fc1(x)
    x = F.relu(x)
    x = (x * self.fc1.output_scale).round()
    x = torch.clamp(x, minimum, maximum)

    # ---- FC2 ----
    x = self.fc2(x)
    x = F.relu(x)
    x = (x * self.fc2.output_scale).round()
    x = torch.clamp(x, minimum, maximum)

    # ---- FC3 (logits) ----
    x = self.fc3(x)
    x = (x * self.fc3.output_scale).round()
    x = torch.clamp(x, minimum, maximum)

    return x
```

I quantized the activations at every layer using the scaling factors computed during training. After each convolution, fully connected layer, and ReLU, the outputs were scaled, rounded, and clamped to the int8 range.

With both weights and activations quantized, the model reached 53.17% accuracy. This is a small drop from the weight-only case, which is expected since activations are quantized repeatedly throughout the forward pass. Even so, the network remains stable under full activation quantization.

Q5

```
@staticmethod
def quantized_bias(bias: torch.Tensor, n_w: float, n_initial_input: float, ns: List[Tuple[float, float]]) -> torch.Tensor:
    """
    Quantize the bias so that all values are integers between -2147483648 and 2147483647.

    Parameters:
    bias (Tensor): The floating point values of the bias
    n_w (float): The scale by which the weights of this layer were multiplied
    n_initial_input (float): The scale by which the initial input to the neural network was multiplied
    ns ((float, float)): A list of tuples, where each tuple represents the "weight scale" and "output scale" (in that order) for every preceding layer

    Returns:
    Tensor: The bias in quantized form, where every value is an integer between -2147483648 and 2147483647.
    The "dtype" will still be "float", but the values themselves should all be integers.
    """

    # ADD YOUR CODE HERE
    # Start from the scale used for the very first input
    scale = n_initial_input
    for i in range(0, len(ns)):
        weights_scale, outputs_scale = ns[i]
        scale = weights_scale * outputs_scale * scale

    # Include this layer's weight scale
    scale = scale * n_w

    return torch.clamp((bias * scale).round(), min=-2147483648, max=2147483647)
```

I quantized the biases using the combined scaling from the input, previous layers' output scales, and the current layer's weight scale.

- The full-precision model reached 54.66% accuracy.
- After quantizing only the weights, the accuracy dropped to 51.4%.
- Once I also quantized the biases, the accuracy improved to 54.49% almost matching the original model.

This shows that properly scaling and quantizing the biases helps recover some of the accuracy lost from weight quantization and is important for stable fully quantized inference.