

# Project 2: DNN Quantization via NNI

Oliver Fowler and Brian Park

North Carolina State University, Computer Science 591/791-025

October 2022

## 1 Quantization Methods

We experimented with 5 quantization methods, since they are mostly rinse and repeat with NNI [?]. They are `NaiveQuantizer`, `QAT_Quantizer`, `DoReFaQuantizer`, `BNNQuantizer`, and `ObserverQuantizer` [?, ?, ?]. We applied them to a simple CNN for the MNIST dataset from the PyTorch example repository. Then we stepped it up a notch and applied it to state of the art model ResNet-101 on the CIFAR-10 dataset [?].

## 2 Quantization and Hardware Configurations

### 2.1 Hardware Configuration

We did all of our workload in the NVIDIA A100 on the ARC cluster [?]. Also, we could not do quantization properly on the M1 Mac due to the limitations of NNI not supporting the datatypes supported by the M1 Mac GPU. Thus, all of our experiments was done on the GPU.

### 2.2 Quantization Configuration

First we began with only quantizing the convolutional layers. Then we also added on quantizing the linear layers for a full effect of quantization and speedup. These are the only two layers that could be quantized by the NNI API as they are the core part of the network with weights. Fortunately, there is hardware support for quantization on NVIDIA A100, such as the support for IEEE 16 bit floating point, 16 bit Brain Floating Point (`bf16`), 8 bit integer, etc [?]. Unfortunately, there was trouble with getting NNI setup to fully see the speedups of quantization, as the quantization steps are simulated until they are compiled with TensorRT. There was trouble getting TensorRT to be properly installed on the ARC, thus the speedups will not be shown.

## 3 Experimental Results

The `print(model)` of the models before and after quantization of Simple CNN and ResNet-101 are available in Appendix A.

`NaiveQuantizer`, `QAT_Quantizer`, `DoReFaQuantizer`, and `BNNQuantizer` are quantization methods that require training from scratch. We showcase the loss and accuracies of the training and validation datasets for training and compared it to the baseline training process. We setup the code to train for 180 epochs and exit training early if we hit a training accuracy of 99%, which is why MNIST will often exit early. For a

The results for the four quantization methods are shown in Figure 1 for the simple CNN trained on the MNIST dataset. `BNNQuantizer` seemed to be the most difficult to train, as accuracy and loss wouldn't budge after many epochs. This is probably because there is not much complexity in the weights for MNIST to deal with and the model is not over-parameterized. Plus the weights are either  $-1$  or  $1$ , making quantization and training for a small model difficult. It's also important to point out that Courbariaux only experimented binary quantization with multi layer perceptrons (MLP) on the MNIST dataset, which has no convolutional layers like our model [?]. We wonder if something is wrong with our weight initialization or training methodology, and we would explore this further if given more time. Figure 2 shows the training data for ResNet-101.

We compare these quantization training schemes with their baseline schemes in Figures 3 and 4.

An interesting analysis of the binary neural network shows that there are 503213 1s and 696435 -1s for the weights of the simple CNN trained on MNIST dataset. For ResNet-101 trained on CIFAR-10 dataset, that's 22052154 1s and 22390662 -1s.

## 4 Lessons Learned

We learned that there are varying levels of quantizations that can come into play. The most convenient one is quantization for post-training. This is `ObserverQuantizer`. This simply changes the weights into lower precision. There are tradeoffs

Some difficulties that arise when working with NNI is that it's not clear if the weights are quantized or not. It was not clear that we had to *train from scratch* for most of the quantization methods. We found it really inefficient to train if we already had the weights, especially for very deep neural networks like ResNet-101. Only `ObserverQuantizer` will need some slight fine tuning after post training to stabilize the weights.

Another difficulty was that training with quantized networks caused loss functions to become unstable. We noticed that gradients often exploded and made the loss function output NaN or got stuck at a certain local minima.

We did not experiment with any quantizations for the M1 Mac, as the functionality is limited. Quantization seems to be much more optimized and supported for NVIDIA CUDA devices.

## 5 GitHub Repository

The GitHub repository for this report is publicly accessible here [?]. To reproduce our findings, please read the `README.md` under the `proj2` directory. If there are any setup issues on ARC cluster, please contact `bcpark@ncsu.edu`.

## References

- [1] Microsoft, "Neural Network Intelligence," 1 2021.
- [2] S. Zhou, Z. Ni, X. Zhou, H. Wen, Y. Wu, and Y. Zou, "Dorefa-net: Training low bitwidth convolutional neural networks with low bitwidth gradients," *CoRR*, vol. abs/1606.06160, 2016.
- [3] M. Courbariaux and Y. Bengio, "Binarynet: Training deep neural networks with weights and activations constrained to +1 or -1," *CoRR*, vol. abs/1602.02830, 2016.
- [4] B. Jacob, S. Kligys, B. Chen, M. Zhu, M. Tang, A. G. Howard, H. Adam, and D. Kalenichenko, "Quantization and training of neural networks for efficient integer-arithmetic-only inference," *CoRR*, vol. abs/1712.05877, 2017.
- [5] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," 2015.
- [6] "Arc: A root cluster for research into scalable computer systems." <https://arcb.csc.ncsu.edu/~mueller/cluster/arc/>.
- [7] "Nvidia a100 datasheet." <https://www.nvidia.com/content/dam/en-zz/Solutions/Data-Center/a100/pdf/nvidia-a100-datasheet-nvidia-us-2188504-web.pdf>.
- [8] "Csc 791-025 github repository." <https://github.com/briancpark/csc791-025/tree/main/proj2>.

## 6 Appendix

### 6.0.1 A1: Layers of Simple CNN for MNIST

Original Simple CNN

```
Net(  
  (conv1): Conv2d(1, 32, kernel_size=(3, 3), stride=(1, 1))  
  (conv2): Conv2d(32, 64, kernel_size=(3, 3), stride=(1, 1))  
  (dropout1): Dropout(p=0.25, inplace=False)  
  (dropout2): Dropout(p=0.5, inplace=False)  
  (fc1): Linear(in_features=9216, out_features=128, bias=True)  
  (fc2): Linear(in_features=128, out_features=10, bias=True)  
)
```

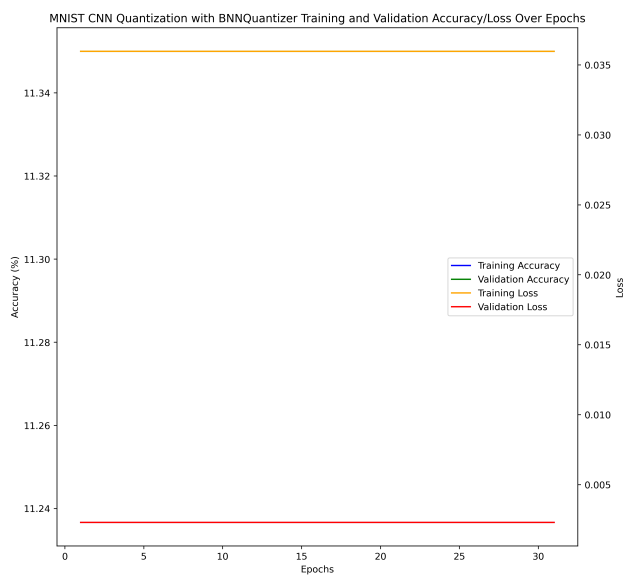
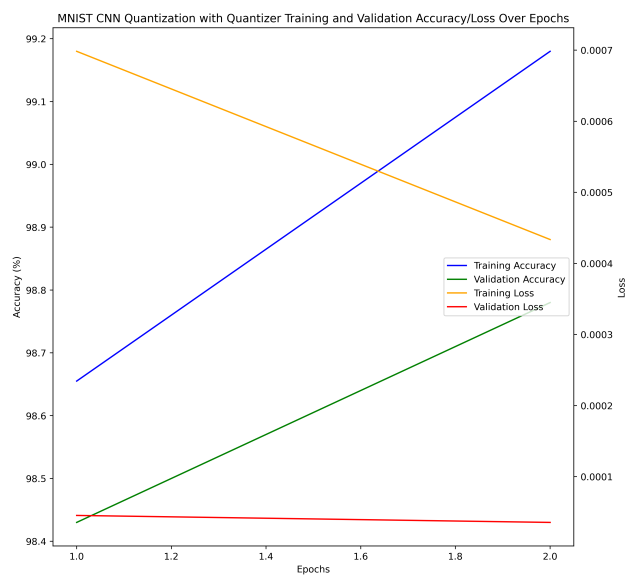
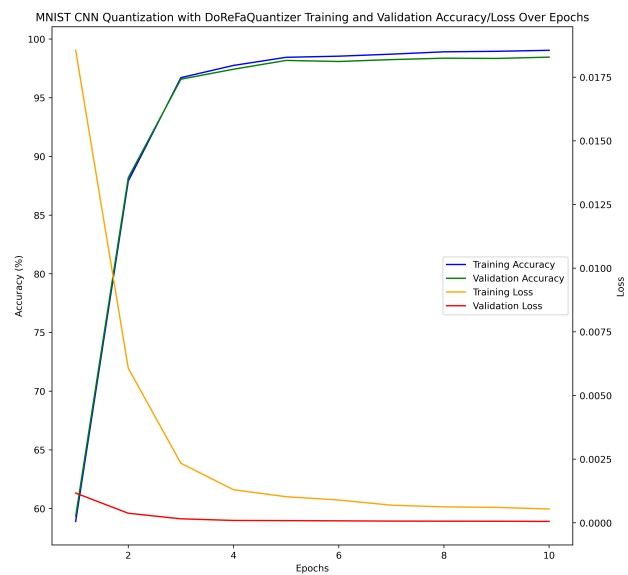
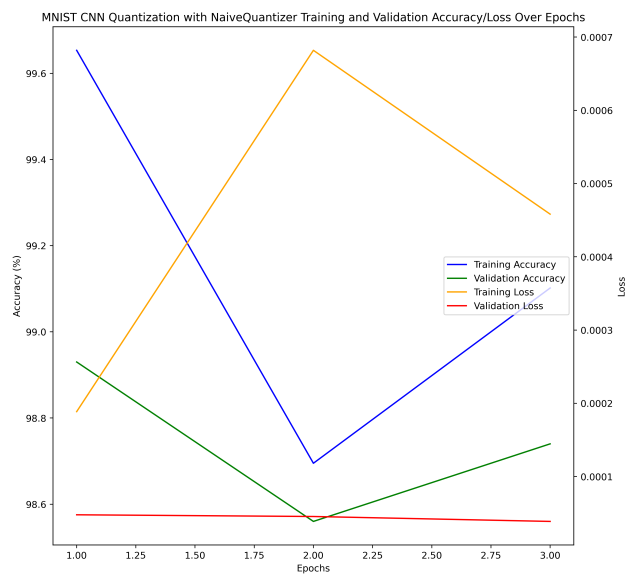
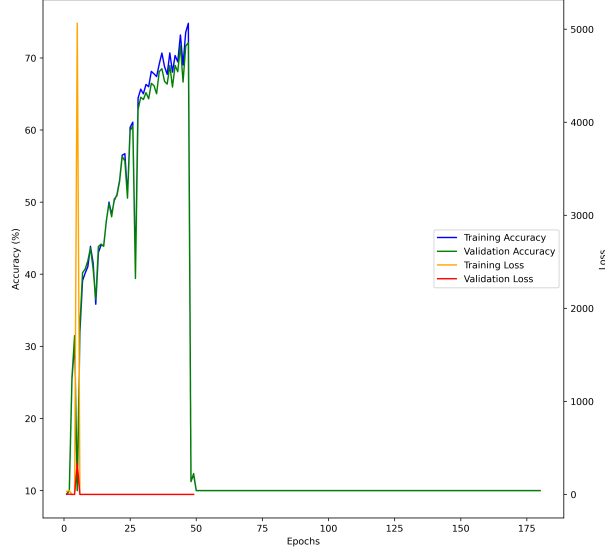
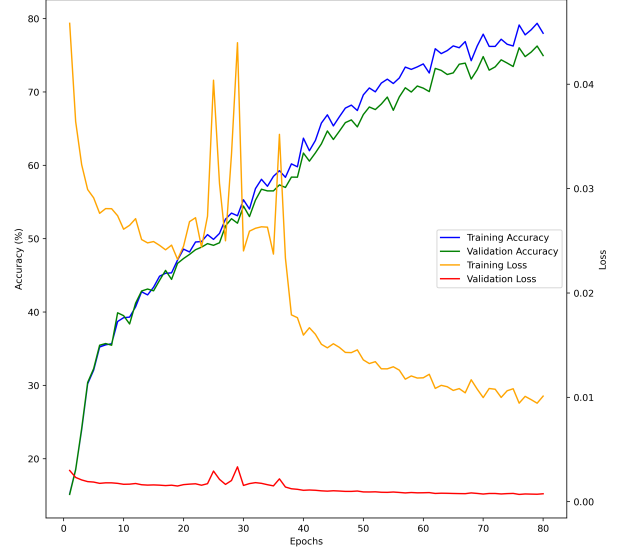


Figure 1: MNIST CNN Training and Validation Accuracy/Loss

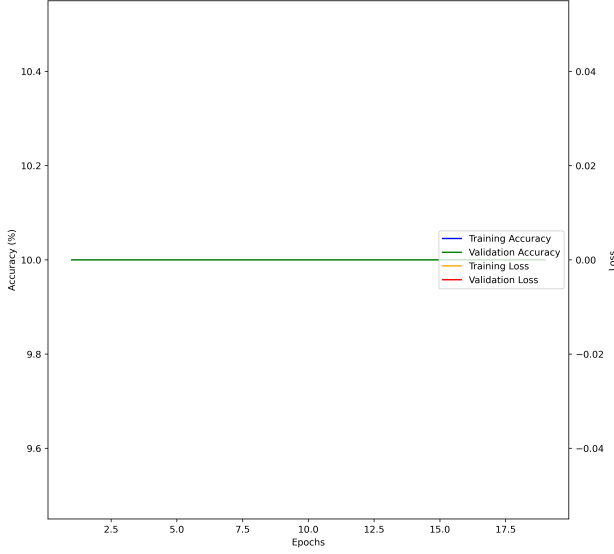
CIFAR-10 ResNet-101 Quantization with NaiveQuantizer Training and Validation Accuracy/Loss Over Epochs



CIFAR-10 ResNet-101 Quantization with DoReFaQuantizer Training and Validation Accuracy/Loss Over Epochs



CIFAR-10 ResNet-101 Quantization with Quantizer Training and Validation Accuracy/Loss Over Epochs



CIFAR-10 ResNet-101 Quantization with BNNQuantizer Training and Validation Accuracy/Loss Over Epochs

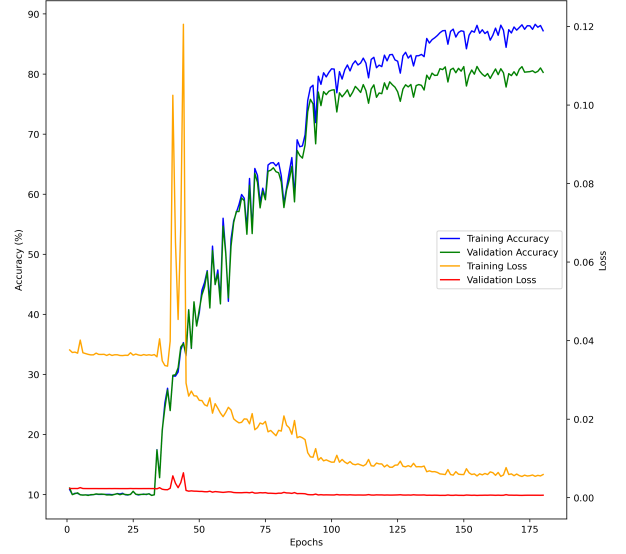


Figure 2: CIFAR-10 ResNet-101 Training and Validation Accuracy/Loss

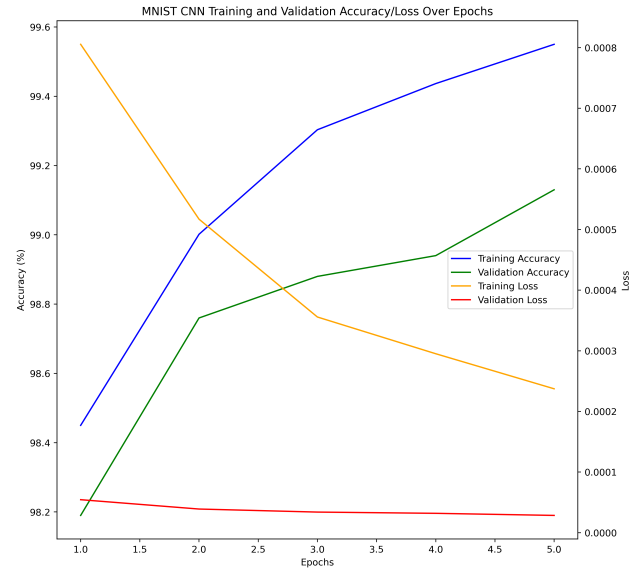


Figure 3: MNIST CNN Training and Validation Accuracy/Loss

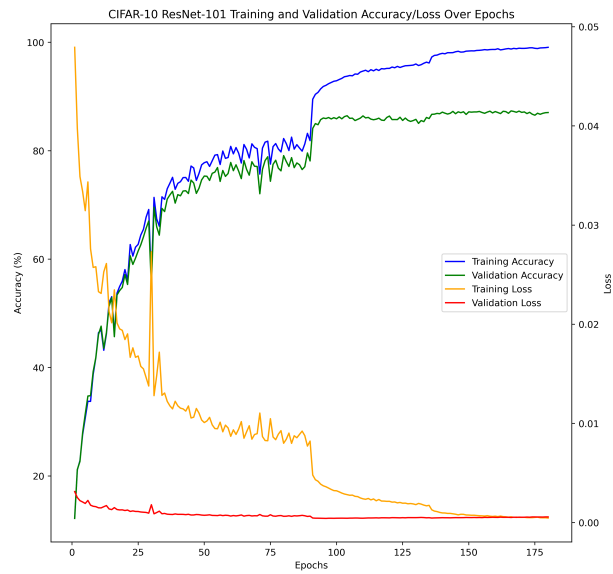


Figure 4: CIFAR-10 ResNet-101 Training and Validation Accuracy/Loss

Quantized with DoReFaQuantizer

```
Net(
  (conv1): QuantizerModuleWrapper(
    (module): Conv2d(1, 32, kernel_size=(3, 3), stride=(1, 1))
  )
  (conv2): QuantizerModuleWrapper(
    (module): Conv2d(32, 64, kernel_size=(3, 3), stride=(1, 1))
  )
  (dropout1): Dropout(p=0.25, inplace=False)
  (dropout2): Dropout(p=0.5, inplace=False)
  (fc1): QuantizerModuleWrapper(
    (module): Linear(in_features=9216, out_features=128, bias=True)
  )
  (fc2): QuantizerModuleWrapper(
    (module): Linear(in_features=128, out_features=10, bias=True)
  )
)
```

### 6.0.2 A1: Layers of ResNet-101 for CIFAR-10

Original ResNet-101 model

```
ResNet(
  (conv1): Conv2d(3, 64, kernel_size=(7, 7), stride=(2, 2), padding=(3, 3), bias=False)
  (bn1): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (relu): ReLU(inplace=True)
  (maxpool): MaxPool2d(kernel_size=3, stride=2, padding=1, dilation=1, ceil_mode=False)
  (layer1): Sequential(
    (0): Bottleneck(
      (conv1): Conv2d(64, 64, kernel_size=(1, 1), stride=(1, 1), bias=False)
      (bn1): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
      (conv2): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
      (bn2): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
      (conv3): Conv2d(64, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
      (bn3): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
      (relu): ReLU(inplace=True)
      (downsample): Sequential(
        (0): Conv2d(64, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
        (1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
      )
    )
  )
  (1): Bottleneck(
    (conv1): Conv2d(256, 64, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn1): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv2): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    (bn2): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv3): Conv2d(64, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn3): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (relu): ReLU(inplace=True)
  )
  (2): Bottleneck(
    (conv1): Conv2d(256, 64, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn1): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv2): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    (bn2): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv3): Conv2d(64, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn3): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (relu): ReLU(inplace=True)
  )
)
```

```

    )
)
(layer2): Sequential(
  (0): Bottleneck(
    (conv1): Conv2d(256, 128, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn1): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv2): Conv2d(128, 128, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1), bias=False)
    (bn2): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv3): Conv2d(128, 512, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn3): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (relu): ReLU(inplace=True)
    (downsample): Sequential(
      (0): Conv2d(256, 512, kernel_size=(1, 1), stride=(2, 2), bias=False)
      (1): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    )
  )
)
  (1): Bottleneck(
    (conv1): Conv2d(512, 128, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn1): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv2): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    (bn2): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv3): Conv2d(128, 512, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn3): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (relu): ReLU(inplace=True)
  )
)
  (2): Bottleneck(
    (conv1): Conv2d(512, 128, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn1): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv2): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    (bn2): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv3): Conv2d(128, 512, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn3): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (relu): ReLU(inplace=True)
  )
)
  (3): Bottleneck(
    (conv1): Conv2d(512, 128, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn1): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv2): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    (bn2): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv3): Conv2d(128, 512, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn3): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (relu): ReLU(inplace=True)
  )
)
)
(layer3): Sequential(
  (0): Bottleneck(
    (conv1): Conv2d(512, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv2): Conv2d(256, 256, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1), bias=False)
    (bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv3): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (relu): ReLU(inplace=True)
    (downsample): Sequential(
      (0): Conv2d(512, 1024, kernel_size=(1, 1), stride=(2, 2), bias=False)
      (1): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    )
  )
)

```

[illegible]







```

)
(20): Bottleneck(
  (conv1): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
  (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv2): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
  (bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv3): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
  (bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (relu): ReLU(inplace=True)
)
(21): Bottleneck(
  (conv1): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
  (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv2): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
  (bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv3): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
  (bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (relu): ReLU(inplace=True)
)
(22): Bottleneck(
  (conv1): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
  (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv2): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
  (bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv3): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
  (bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (relu): ReLU(inplace=True)
)
)
(layer4): Sequential(
  (0): Bottleneck(
    (conv1): Conv2d(1024, 512, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn1): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv2): Conv2d(512, 512, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1), bias=False)
    (bn2): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv3): Conv2d(512, 2048, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn3): BatchNorm2d(2048, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (relu): ReLU(inplace=True)
    (downsample): Sequential(
      (0): Conv2d(1024, 2048, kernel_size=(1, 1), stride=(2, 2), bias=False)
      (1): BatchNorm2d(2048, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    )
  )
  (1): Bottleneck(
    (conv1): Conv2d(2048, 512, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn1): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv2): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    (bn2): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv3): Conv2d(512, 2048, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn3): BatchNorm2d(2048, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (relu): ReLU(inplace=True)
  )
  (2): Bottleneck(
    (conv1): Conv2d(2048, 512, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn1): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv2): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    (bn2): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)

```

```

        (conv3): Conv2d(512, 2048, kernel_size=(1, 1), stride=(1, 1), bias=False)
        (bn3): BatchNorm2d(2048, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
        (relu): ReLU(inplace=True)
    )
)
(avgpool): AdaptiveAvgPool2d(output_size=(1, 1))
(fc): Linear(in_features=2048, out_features=1000, bias=True)
)

Quantized with DoReFaQuantizer

ResNet(
  (conv1): QuantizerModuleWrapper(
    (module): Conv2d(3, 64, kernel_size=(7, 7), stride=(2, 2), padding=(3, 3), bias=False)
  )
  (bn1): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (relu): ReLU(inplace=True)
  (maxpool): MaxPool2d(kernel_size=3, stride=2, padding=1, dilation=1, ceil_mode=False)
  (layer1): Sequential(
    (0): Bottleneck(
      (conv1): QuantizerModuleWrapper(
        (module): Conv2d(64, 64, kernel_size=(1, 1), stride=(1, 1), bias=False)
      )
      (bn1): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
      (conv2): QuantizerModuleWrapper(
        (module): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
      )
      (bn2): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
      (conv3): QuantizerModuleWrapper(
        (module): Conv2d(64, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
      )
      (bn3): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
      (relu): ReLU(inplace=True)
      (downsample): Sequential(
        (0): QuantizerModuleWrapper(
          (module): Conv2d(64, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
        )
        (1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
      )
    )
  )
  (1): Bottleneck(
    (conv1): QuantizerModuleWrapper(
      (module): Conv2d(256, 64, kernel_size=(1, 1), stride=(1, 1), bias=False)
    )
    (bn1): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv2): QuantizerModuleWrapper(
      (module): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    )
    (bn2): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv3): QuantizerModuleWrapper(
      (module): Conv2d(64, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
    )
    (bn3): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (relu): ReLU(inplace=True)
  )
  (2): Bottleneck(
    (conv1): QuantizerModuleWrapper(
      (module): Conv2d(256, 64, kernel_size=(1, 1), stride=(1, 1), bias=False)
    )

```

```

        (bn1): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
        (conv2): QuantizerModuleWrapper(
          (module): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
        )
        (bn2): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
        (conv3): QuantizerModuleWrapper(
          (module): Conv2d(64, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
        )
        (bn3): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
        (relu): ReLU(inplace=True)
      )
    )
(layer2): Sequential(
  (0): Bottleneck(
    (conv1): QuantizerModuleWrapper(
      (module): Conv2d(256, 128, kernel_size=(1, 1), stride=(1, 1), bias=False)
    )
    (bn1): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv2): QuantizerModuleWrapper(
      (module): Conv2d(128, 128, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1), bias=False)
    )
    (bn2): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv3): QuantizerModuleWrapper(
      (module): Conv2d(128, 512, kernel_size=(1, 1), stride=(1, 1), bias=False)
    )
    (bn3): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (relu): ReLU(inplace=True)
    (downsample): Sequential(
      (0): QuantizerModuleWrapper(
        (module): Conv2d(256, 512, kernel_size=(1, 1), stride=(2, 2), bias=False)
      )
      (1): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    )
  )
)
  (1): Bottleneck(
    (conv1): QuantizerModuleWrapper(
      (module): Conv2d(512, 128, kernel_size=(1, 1), stride=(1, 1), bias=False)
    )
    (bn1): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv2): QuantizerModuleWrapper(
      (module): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    )
    (bn2): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv3): QuantizerModuleWrapper(
      (module): Conv2d(128, 512, kernel_size=(1, 1), stride=(1, 1), bias=False)
    )
    (bn3): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (relu): ReLU(inplace=True)
  )
)
  (2): Bottleneck(
    (conv1): QuantizerModuleWrapper(
      (module): Conv2d(512, 128, kernel_size=(1, 1), stride=(1, 1), bias=False)
    )
    (bn1): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv2): QuantizerModuleWrapper(
      (module): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    )

```

```

        (bn2): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
        (conv3): QuantizerModuleWrapper(
          (module): Conv2d(128, 512, kernel_size=(1, 1), stride=(1, 1), bias=False)
        )
        (bn3): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
        (relu): ReLU(inplace=True)
      )
    (3): Bottleneck(
      (conv1): QuantizerModuleWrapper(
        (module): Conv2d(512, 128, kernel_size=(1, 1), stride=(1, 1), bias=False)
      )
      (bn1): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
      (conv2): QuantizerModuleWrapper(
        (module): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
      )
      (bn2): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
      (conv3): QuantizerModuleWrapper(
        (module): Conv2d(128, 512, kernel_size=(1, 1), stride=(1, 1), bias=False)
      )
      (bn3): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
      (relu): ReLU(inplace=True)
    )
  )
(layer3): Sequential(
  (0): Bottleneck(
    (conv1): QuantizerModuleWrapper(
      (module): Conv2d(512, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
    )
    (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv2): QuantizerModuleWrapper(
      (module): Conv2d(256, 256, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1), bias=False)
    )
    (bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv3): QuantizerModuleWrapper(
      (module): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
    )
    (bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (relu): ReLU(inplace=True)
    (downsample): Sequential(
      (0): QuantizerModuleWrapper(
        (module): Conv2d(512, 1024, kernel_size=(1, 1), stride=(2, 2), bias=False)
      )
      (1): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    )
  )
)
(1): Bottleneck(
  (conv1): QuantizerModuleWrapper(
    (module): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv2): QuantizerModuleWrapper(
    (module): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
  )
  (bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv3): QuantizerModuleWrapper(
    (module): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
)

```

```

        (bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
        (relu): ReLU(inplace=True)
    )
(2): Bottleneck(
  (conv1): QuantizerModuleWrapper(
    (module): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv2): QuantizerModuleWrapper(
    (module): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
  )
  (bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv3): QuantizerModuleWrapper(
    (module): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (relu): ReLU(inplace=True)
)
(3): Bottleneck(
  (conv1): QuantizerModuleWrapper(
    (module): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv2): QuantizerModuleWrapper(
    (module): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
  )
  (bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv3): QuantizerModuleWrapper(
    (module): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (relu): ReLU(inplace=True)
)
(4): Bottleneck(
  (conv1): QuantizerModuleWrapper(
    (module): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv2): QuantizerModuleWrapper(
    (module): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
  )
  (bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv3): QuantizerModuleWrapper(
    (module): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (relu): ReLU(inplace=True)
)
(5): Bottleneck(
  (conv1): QuantizerModuleWrapper(
    (module): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv2): QuantizerModuleWrapper(
    (module): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
  )
  (bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)

```

```

(conv3): QuantizerModuleWrapper(
  (module): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
)
(bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(relu): ReLU(inplace=True)
)
(6): Bottleneck(
  (conv1): QuantizerModuleWrapper(
    (module): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv2): QuantizerModuleWrapper(
    (module): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
  )
  (bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv3): QuantizerModuleWrapper(
    (module): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (relu): ReLU(inplace=True)
)
(7): Bottleneck(
  (conv1): QuantizerModuleWrapper(
    (module): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv2): QuantizerModuleWrapper(
    (module): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
  )
  (bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv3): QuantizerModuleWrapper(
    (module): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (relu): ReLU(inplace=True)
)
(8): Bottleneck(
  (conv1): QuantizerModuleWrapper(
    (module): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv2): QuantizerModuleWrapper(
    (module): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
  )
  (bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv3): QuantizerModuleWrapper(
    (module): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (relu): ReLU(inplace=True)
)
(9): Bottleneck(
  (conv1): QuantizerModuleWrapper(
    (module): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv2): QuantizerModuleWrapper(

```



```

    (module): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
)
(bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(conv3): QuantizerModuleWrapper(
  (module): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
)
(bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(relu): ReLU(inplace=True)
)
(10): Bottleneck(
  (conv1): QuantizerModuleWrapper(
    (module): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv2): QuantizerModuleWrapper(
    (module): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
  )
  (bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv3): QuantizerModuleWrapper(
    (module): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (relu): ReLU(inplace=True)
)
(11): Bottleneck(
  (conv1): QuantizerModuleWrapper(
    (module): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv2): QuantizerModuleWrapper(
    (module): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
  )
  (bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv3): QuantizerModuleWrapper(
    (module): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (relu): ReLU(inplace=True)
)
(12): Bottleneck(
  (conv1): QuantizerModuleWrapper(
    (module): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv2): QuantizerModuleWrapper(
    (module): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
  )
  (bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv3): QuantizerModuleWrapper(
    (module): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (relu): ReLU(inplace=True)
)
(13): Bottleneck(
  (conv1): QuantizerModuleWrapper(
    (module): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)

```

```

)
(bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(conv2): QuantizerModuleWrapper(
  (module): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
)
(bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(conv3): QuantizerModuleWrapper(
  (module): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
)
(bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(relu): ReLU(inplace=True)
)
(14): Bottleneck(
  (conv1): QuantizerModuleWrapper(
    (module): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv2): QuantizerModuleWrapper(
    (module): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
  )
  (bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv3): QuantizerModuleWrapper(
    (module): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (relu): ReLU(inplace=True)
)
(15): Bottleneck(
  (conv1): QuantizerModuleWrapper(
    (module): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv2): QuantizerModuleWrapper(
    (module): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
  )
  (bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv3): QuantizerModuleWrapper(
    (module): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (relu): ReLU(inplace=True)
)
(16): Bottleneck(
  (conv1): QuantizerModuleWrapper(
    (module): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv2): QuantizerModuleWrapper(
    (module): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
  )
  (bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv3): QuantizerModuleWrapper(
    (module): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (relu): ReLU(inplace=True)
)
)

```

```

(17): Bottleneck(
  (conv1): QuantizerModuleWrapper(
    (module): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv2): QuantizerModuleWrapper(
    (module): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
  )
  (bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv3): QuantizerModuleWrapper(
    (module): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (relu): ReLU(inplace=True)
)
(18): Bottleneck(
  (conv1): QuantizerModuleWrapper(
    (module): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv2): QuantizerModuleWrapper(
    (module): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
  )
  (bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv3): QuantizerModuleWrapper(
    (module): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (relu): ReLU(inplace=True)
)
(19): Bottleneck(
  (conv1): QuantizerModuleWrapper(
    (module): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv2): QuantizerModuleWrapper(
    (module): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
  )
  (bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv3): QuantizerModuleWrapper(
    (module): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (relu): ReLU(inplace=True)
)
(20): Bottleneck(
  (conv1): QuantizerModuleWrapper(
    (module): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv2): QuantizerModuleWrapper(
    (module): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
  )
  (bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv3): QuantizerModuleWrapper(
    (module): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
)

```

```

        (bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
        (relu): ReLU(inplace=True)
    )
(21): Bottleneck(
  (conv1): QuantizerModuleWrapper(
    (module): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv2): QuantizerModuleWrapper(
    (module): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
  )
  (bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv3): QuantizerModuleWrapper(
    (module): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (relu): ReLU(inplace=True)
)
(22): Bottleneck(
  (conv1): QuantizerModuleWrapper(
    (module): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv2): QuantizerModuleWrapper(
    (module): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
  )
  (bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv3): QuantizerModuleWrapper(
    (module): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (relu): ReLU(inplace=True)
)
)
(layer4): Sequential(
  (0): Bottleneck(
    (conv1): QuantizerModuleWrapper(
      (module): Conv2d(1024, 512, kernel_size=(1, 1), stride=(1, 1), bias=False)
    )
    (bn1): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv2): QuantizerModuleWrapper(
      (module): Conv2d(512, 512, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1), bias=False)
    )
    (bn2): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv3): QuantizerModuleWrapper(
      (module): Conv2d(512, 2048, kernel_size=(1, 1), stride=(1, 1), bias=False)
    )
    (bn3): BatchNorm2d(2048, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (relu): ReLU(inplace=True)
    (downsample): Sequential(
      (0): QuantizerModuleWrapper(
        (module): Conv2d(1024, 2048, kernel_size=(1, 1), stride=(2, 2), bias=False)
      )
      (1): BatchNorm2d(2048, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    )
  )
)
(1): Bottleneck(

```

```

(conv1): QuantizerModuleWrapper(
  (module): Conv2d(2048, 512, kernel_size=(1, 1), stride=(1, 1), bias=False)
)
(bn1): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(conv2): QuantizerModuleWrapper(
  (module): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
)
(bn2): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(conv3): QuantizerModuleWrapper(
  (module): Conv2d(512, 2048, kernel_size=(1, 1), stride=(1, 1), bias=False)
)
(bn3): BatchNorm2d(2048, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(rel): ReLU(inplace=True)
)
(2): Bottleneck(
  (conv1): QuantizerModuleWrapper(
    (module): Conv2d(2048, 512, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn1): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv2): QuantizerModuleWrapper(
    (module): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
  )
  (bn2): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (conv3): QuantizerModuleWrapper(
    (module): Conv2d(512, 2048, kernel_size=(1, 1), stride=(1, 1), bias=False)
  )
  (bn3): BatchNorm2d(2048, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (relu): ReLU(inplace=True)
)
)
(avgpool): AdaptiveAvgPool2d(output_size=(1, 1))
(fc): QuantizerModuleWrapper(
  (module): Linear(in_features=2048, out_features=1000, bias=True)
)
)

```