

Computer Vision for Embedded Systems

Yung-Hsiang Lu
Purdue University
yunglu@purdue.edu

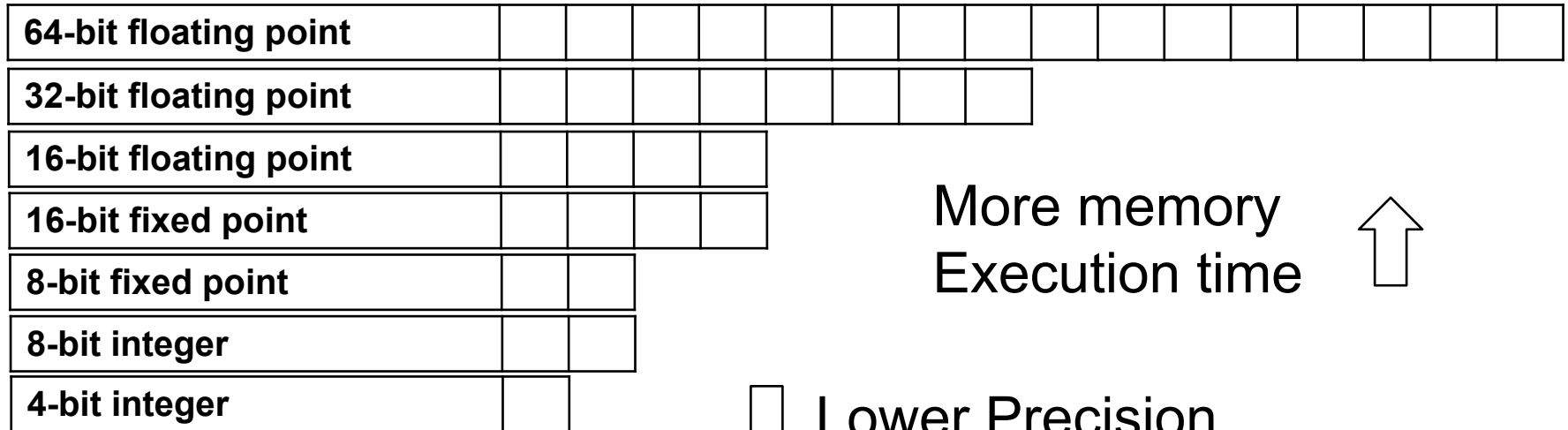


Yung-Hsiang Lu, Purdue University



Reduce Network Sizes

- Reduce #bits for each parameter
- Remove unused connections between layers
- Remove inactive neurons



Benefits of Quantization

- Reduced memory footprint
- Sparse weights
- Faster inference time
- Regularization
- Better cache performance

Tailin Liang, John Glossner, Lei Wang, Shaobo Shi, Xiaotong Zhang, Pruning and quantization for deep neural network acceleration: A survey, Neurocomputing, Volume 461, 2021, Pages 370-403, ISSN 0925-2312, <https://doi.org/10.1016/j.neucom.2021.07.04>

Sparsity

- **Definition:** % of weights that are equal to 0
- **Benefit:** Specialized hardware/software can make computations faster given sparse data

Fast Sparse Matrix Multiplication

<https://www.cs.tau.ac.il/~zwick/papers/sparse.pdf>



Adding Two Floating Point Numbers

How to add 3.75 and 5.125 to get 8.875?

$$3.75 = 1.875 \times 2 = (1 + 0.875) \times 2$$

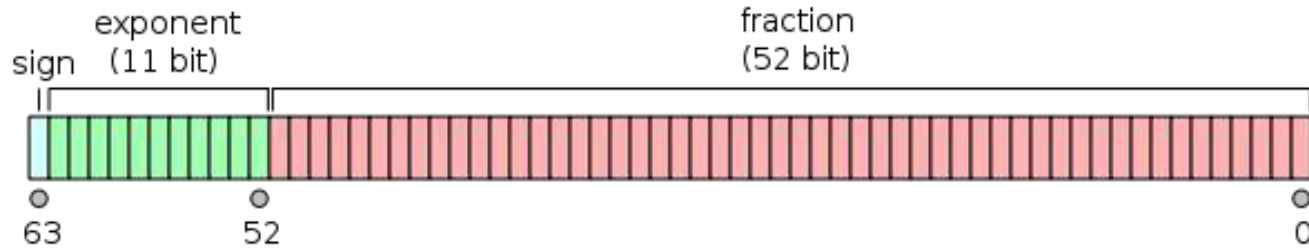
$$0.875 = 0.5 + 0.25 + 0.125 = 2^{-1} + 2^{-2} + 2^{-3}$$

$$5.125 = 1.28125 \times 2^2 = (1 + 0.28125) \times 2^2$$

$$0.28125 = 0.25 + 0.03125 = 2^{-2} + 2^{-5}$$

- A. make them have the same exponent
- B. add the fraction
- C. convert back to the correct format

64-bit Floating Point

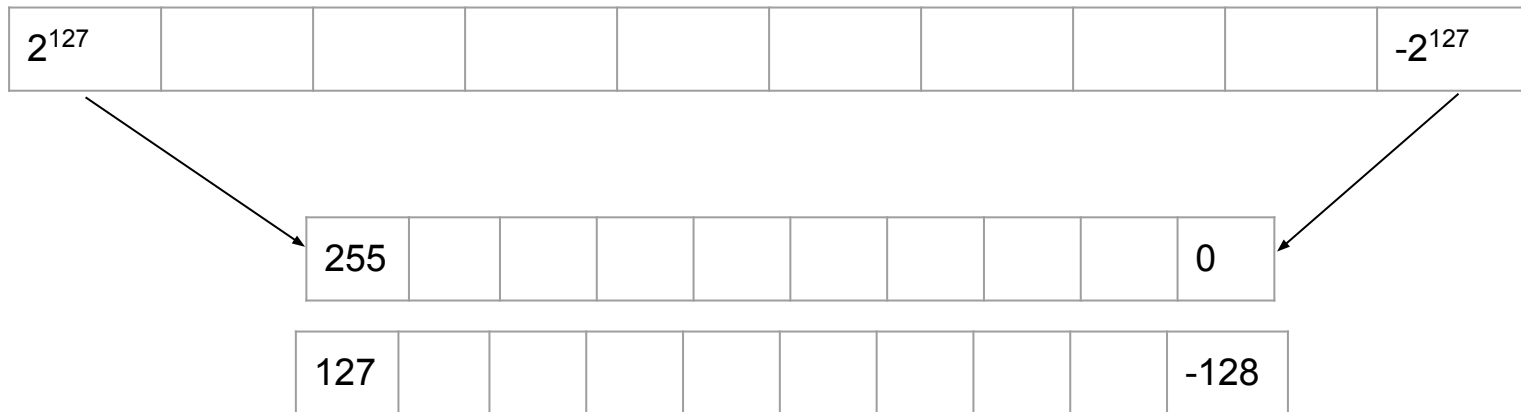


- sign: 0 positive, 1 negative
- 11-bit exponent: $e - 1023$.
- $1 + \text{fraction}$

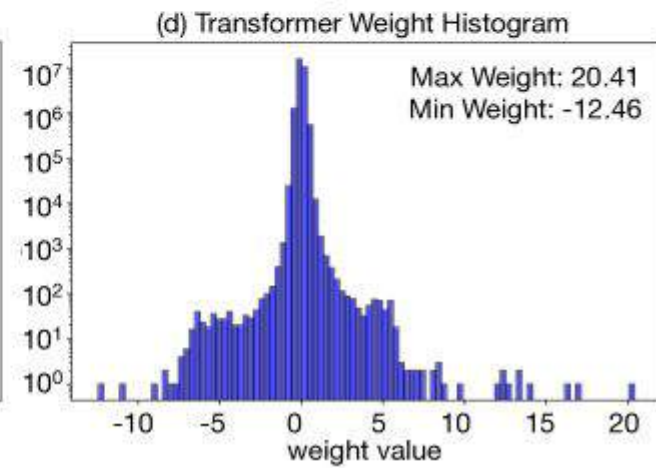
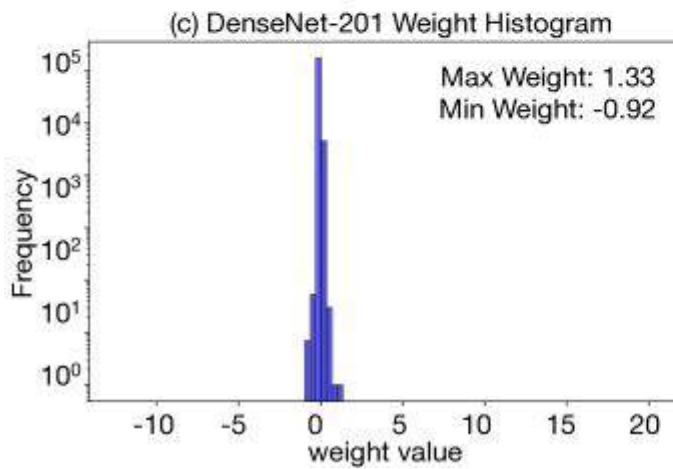
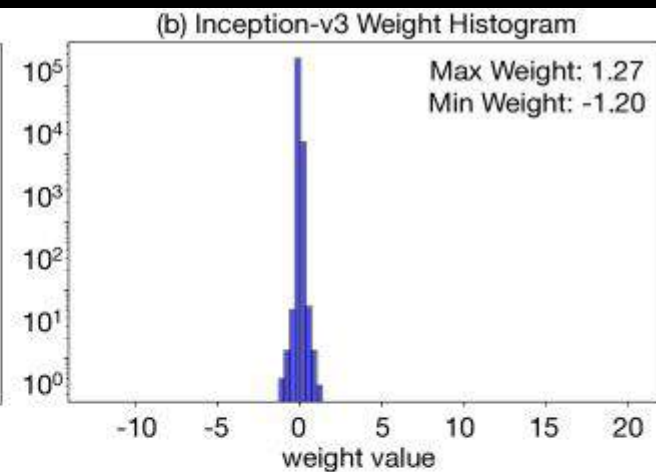
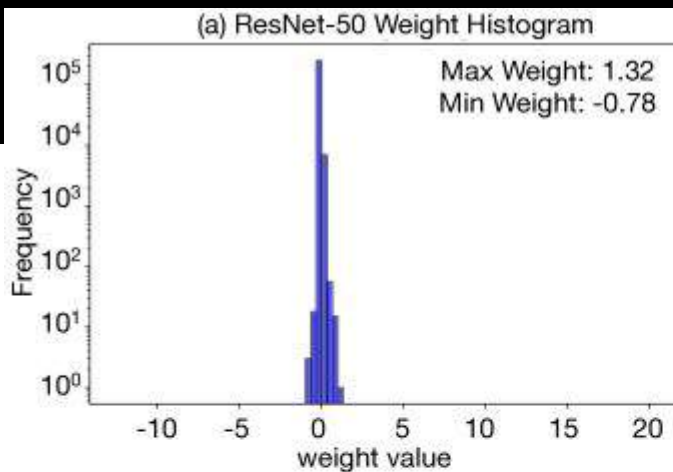
wikipedia.org

Quantization is lossy mapping

32-bit floating point as large as $2^{127} \gg 255$ in 8-bit integer



In reality, neural networks' weights are rarely too large



Quantization as Linear Mapping

floating point $x \in [\alpha, \beta]$ integer $x_q \in [\alpha_q, \beta_q]$

$$x = c(x_q + d) \quad x_q = \left[\frac{x}{c} - d \right]$$

$$\alpha = c(\alpha_q + d) \quad \beta = c(\beta_q + d)$$

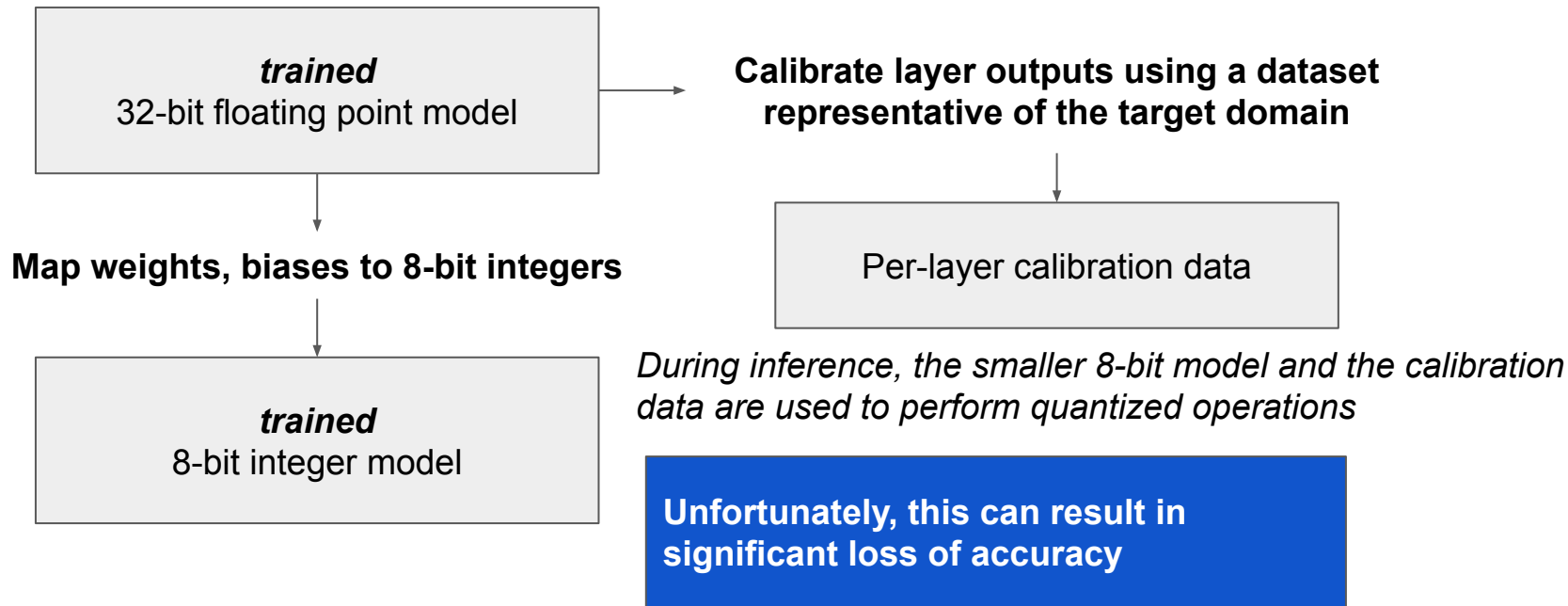
$$c = \frac{\beta - \alpha}{\beta_q - \alpha_q} \quad d = \frac{\alpha\beta_q - \beta\alpha_q}{\beta - \alpha}$$

scale *zero point*

<https://leimao.github.io/article/Neural-Networks-Quantization/>

Method	Inference Latency	Accuracy Lose	Training Data
Dynamic Quantization	Usually faster	Small	No need
Static Quantization	Faster	Larger	Unlabeled
Quantization Aware Training	Faster	Small	Labeled

Post-training Static Quantization Overview



Vincent Vanhoucke, Andrew Senior, & Mark Z. Mao (2011). *Improving the speed of neural networks on CPUs*. In *Deep Learning and Unsupervised Feature Learning Workshop, NIPS 2011*.

Quantization-Aware Training Overview

Forward pass

Round all floating-point weights, biases, and activations to the nearest quantized integer

Backwards propagation

Use floating-point arithmetic as usual

This way, weights can still be adjusted incrementally, while influenced by the quantized forward pass

This still produces a 32-bit floating-point model. However, the weights and biases should now be much more amenable to static quantization, resulting in higher-accuracy than post-training static quantization.

B. Jacob et al., "Quantization and Training of Neural Networks for Efficient Integer-Arithmetic-Only Inference," 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2018, pp. 2704-2713, doi: 10.1109/CVPR.2018.00286.

32 bits to 8 bits: Example MatMul Layer

Let's say we have a **matrix multiplication layer** for fp32 tensors input **X** and weight **W**. It'll give us output **O** (note, SF: scale factor, ZP: zero point). You could do similar stuff for Conv2D and Linear layers.

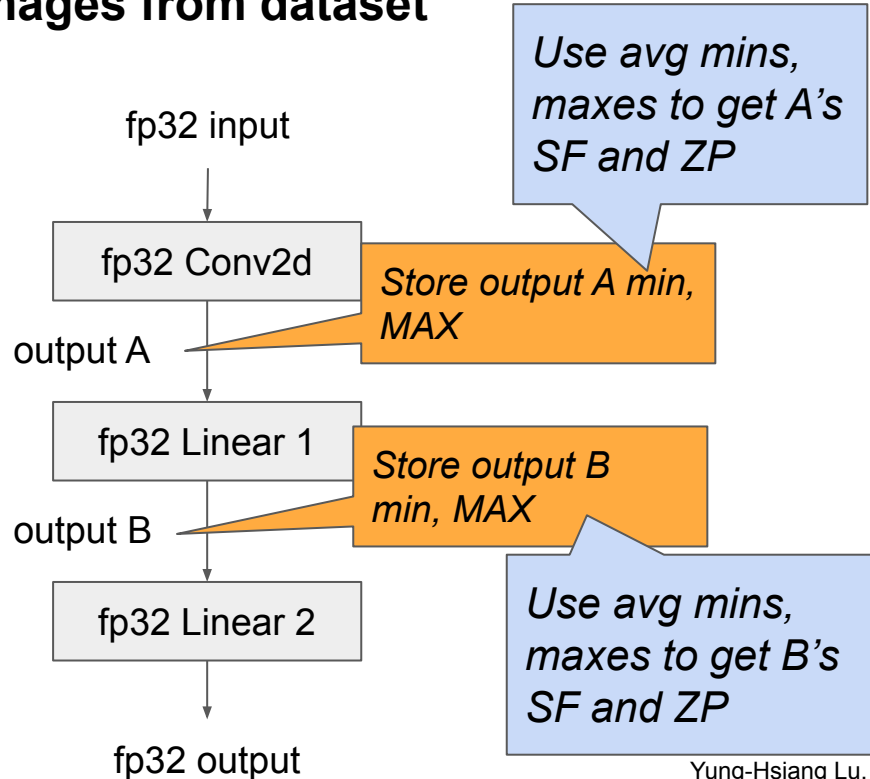
$$\begin{aligned} O_{fp32} &= X_{fp32} \times W_{fp32} \\ \frac{O_{fp32}}{SF_{O,fp32}} (O_{int8} - ZP_{O,int8}) &= SF_{X,fp32} (X_{int8} - ZP_{X,int8}) \times SF_{W,fp32} (W_{int8} - ZP_{W,int8}) \\ O_{int8} &= \left(\frac{1}{SF_{O,fp32}} \right) (SF_{X,fp32} (X_{int8} - ZP_{X,int8}) \times SF_{W,fp32} (W_{int8} - ZP_{W,int8})) + ZP_{O,int8} \end{aligned}$$

All this math is done using only quantized tensors, scale factors, and zero points. No massive FP32 tensors

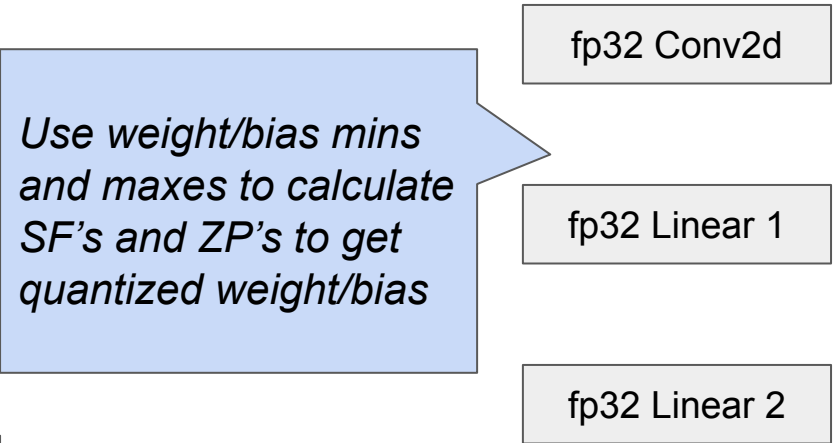
how do we know the scale factor and zero point for **O**?

32 bits to 8 bits: Quantizing an entire NN

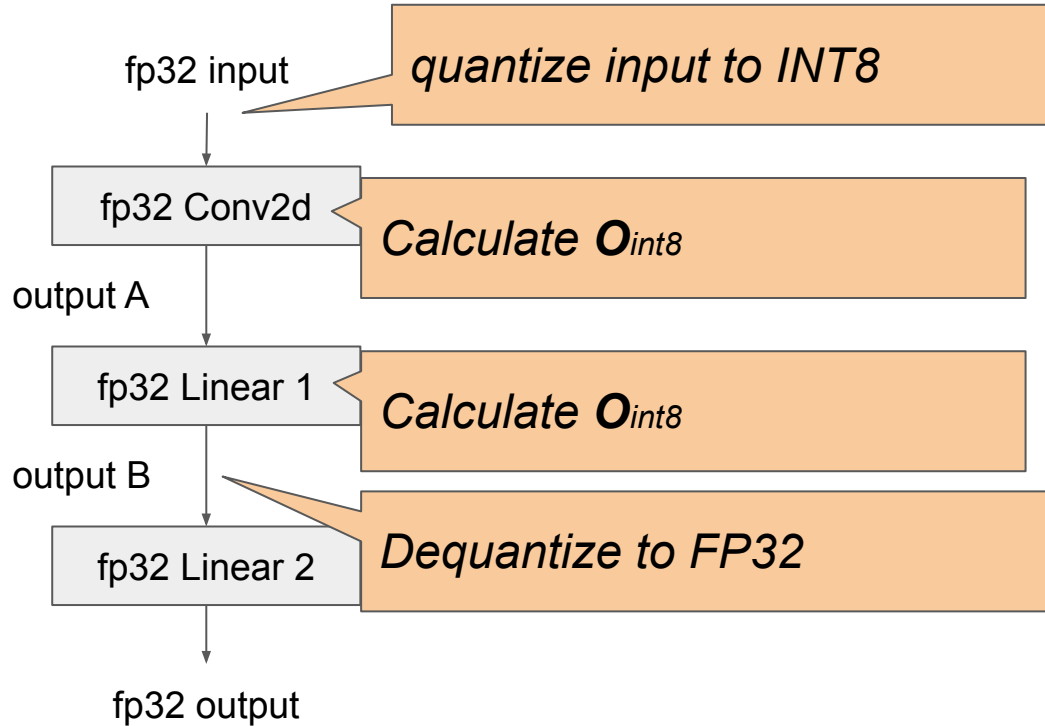
1. Do calibration on multiple images from dataset



2. Quantize weights and biases



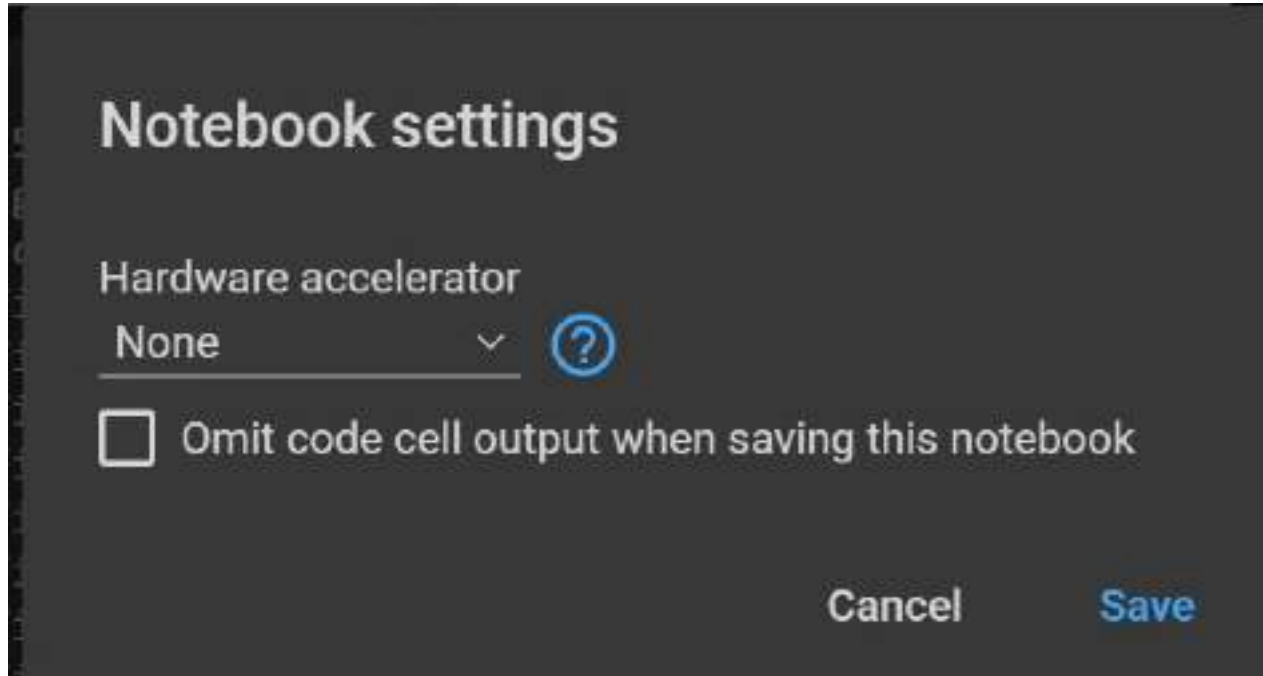
32 bits to 8 bits: Quantized Inference



Post-Training Static Quantization Assignment

Google Colab

Only need to use a CPU runtime in Colab



Get setup

Setup

TASKS:

1. Run these cells to grab the PyPI packages and import the dependencies for the notebook. You can click into the 'Files' explorer on the sidebar to confirm that `./ClassyClassifierParams.pt` was appropriately downloaded.

```
[1] 1 | pip install torchinfo
    2 | pip install gdown
    3 | gdown https://drive.google.com/uc?id=1xTlg4FbV_znqy1KTHt44vbIAqFvwd -O ./ClassyClassifierParams.pt
```

```
Collecting torchinfo
  Downloading torchinfo-1.5.3-py3-none-any.whl (19 kB)
Installing collected packages: torchinfo
Successfully installed torchinfo-1.5.3
Requirement already satisfied: gdown in /usr/local/lib/python3.7/dist-packages (3.6.4)
Requirement already satisfied: six in /usr/local/lib/python3.7/dist-packages (from gdown) (1.15.0)
Requirement already satisfied: requests in /usr/local/lib/python3.7/dist-packages (from gdown) (2.23.0)
Requirement already satisfied: tqdm in /usr/local/lib/python3.7/dist-packages (from gdown) (4.62.3)
Requirement already satisfied: certifi<=2017.4.17 in /usr/local/lib/python3.7/dist-packages (from requests->gdown) (2021.5.30)
Requirement already satisfied: urllib3<=1.25.0, >=1.25.1, <1.26, >=1.21.1 in /usr/local/lib/python3.7/dist-packages (from requests->gdown) (1.24.3)
Requirement already satisfied: chardet<4, >=3.0.2 in /usr/local/lib/python3.7/dist-packages (from requests->gdown) (3.0.4)
Requirement already satisfied: idna<3, >=2.5 in /usr/local/lib/python3.7/dist-packages (from requests->gdown) (2.10)
Downloading...
From: https://drive.google.com/uc?id=1xTlg4FbV\_znqy1KTHt44vbIAqFvwd
To: ./content/ClassyClassifierParams.pt
100% 515k/515k [00:00:00.60, 54.1MB/s]
```

```
[2] 1 | import copy
    2 | import torch.nn.functional as f
    3 | import torch.nn as nn
    4 | import numpy as np
    5 | import os
```

Check out the ClassyClassifier

```
[6] 1 CLASSY_CLASSIFIER_PARAMETERS_FILENAME = "./ClassyClassifierParams.pt"
2
3 class ClassyClassifier(nn.Module):
4     def __init__(self):
5         super().__init__()
6         # INPUT SHAPE: Bx3x32x32 (B is for "batch size," in this case: 5)
7         self.layer1_conv = nn.Conv2d(in_channels=3, out_channels=16, kernel_size=5, stride=1) # output Bx16x28x28
8         self.layer2_pool = nn.MaxPool2d(kernel_size=2, stride=2) # output Bx16x14x14
9         self.layer3_conv = nn.Conv2d(in_channels=16, out_channels=32, kernel_size=5, stride=1) # output Bx32x10x10
10        self.layer4_pool = nn.MaxPool2d(kernel_size=2, stride=2) # output Bx32x5x5
11        self.layer5_flat = nn.Flatten() # output Bx(32x5x5=800)
12        self.layer6_fc = nn.Linear(in_features=800, out_features=128) # output Bx128
13        self.layer7_fc = nn.Linear(in_features=128, out_features=84) # output Bx84
14        self.layer8_fc = nn.Linear(in_features=84, out_features=10) # output Bx10
15
16    def forward(self, x: torch.Tensor):
17        x = F.relu(self.layer1_conv(x))
18        x = self.layer2_pool(x)
19        x = F.relu(self.layer3_conv(x))
20        x = self.layer4_pool(x)
21        x = self.layer5_flat(x)
22        x = F.relu(self.layer6_fc(x))
23        x = F.relu(self.layer7_fc(x))
24        x = self.layer8_fc(x)
25        return x
```

Quantizing by hand...

Quantization formula

To quantize a tensor, you need to **use its min and max values** to calculate the **scale factor** and **zero point**.

$$Value_{fp32} = ScaleFactor_{fp32} \times (Value_{int8} - ZeroPoint_{int8})$$

Generate this

And this



To calculate this

Quantization helper functions

Helper functions to calculate scalefactor/zeropoint, quantize/dequantize

```
1 def calculate_scalefactor_and_zeropoint(min_val: float, max_val: float) -> Tuple[
2     """Calculates the scale factor and zero point to quantize from a [ min_val',
3
4     Follows the quantization formula: FP32 VALUE = scalefactor * (QUANT INT8 VALUE
5
6     Using the formula, float value 'min_val' quantizes to 'INT8_MIN' of 0.
7
8     Args:
9         min_val (float): Minimum value in range of interest.
10        max_val (float): Maximum value in range of interest.
11
12    Returns:
13        tuple(float, int): Scale factor, zero point
14    """
15    INT8_MIN = 0
16    INT8_MAX = 255
17
18    scalefactor = # 1000: fill this in such that you can scale from [min_val, max_val] to [0, 255]
19    zeropoint = # 1000: fill this in
20
21    # This clamps the zero point appropriately into the range [0, 255]
22    if zeropoint < INT8_MIN:
23        zeropoint = INT8_MIN
24    elif zeropoint > INT8_MAX:
25        zeropoint = INT8_MAX
26
27    return scalefactor, int(zeropoint)
```

Read the docstrings

Complete the #TODOs

- Four 1-line #TODOs

Technically, this is *pseudo*-quantization

Yes, we still store as

`FloatTensors`.

But these tensors are still integers, in range `[0, 255]`. The point is to learn quantization arithmetic!

You will use `torch.fx` later
to perform real quantization.

Fill out QuantizedLayer

1. Constructor should call `quantize_tensor()` to get int8 weights/biases with scale factors and zero points
2. `run_quantized_layer` uses the int8 tensors along with their scale factors and zero points to calculate the floating point output. The floating point output is converted to int8 using some precalculated scalefactor and zero point

Again, remember that this is to learn the math. Real hardware implements this symbolically and much quicker!

Tip: refer to the “Example MatMul Layer” slide in the Quantization lecture slide deck to see how the math works

Read the docstrings
Complete the #TODOS

Fill out QuantizerForClassyClassifier

Calibration requires that, for each image batch, you store the min/max of each layer's output (i.e. the preceding layer's input) in lists.

Once this is done, you can average the lists and use the average min/max to calculate scalefactors and zero points (remember, these are needed as parameters for `QuantizedLayer.run_quantized_layer()`)

```
130 # Get calibration data going into layer 1, then use ReLU activations and pooling, store to self.calibration_input
131 # collect the min and max stats for the input of layer 1 (aka the
132 self.calibration_input_stats["layer1_conv"]["mins"].append(x.view(b
133 self.calibration_input_stats["layer1_conv"]["maxes"].append(x.view(b
134 # use ReLU
135 x = f.relu(self.fp32_model.layer1_conv(x))
136 # do layer 1 pooling
137 x = self.fp32_model.layer2_pool(x)
138
139 # Get calibration data going into layer 1, then use ReLU activations and pooling, like above, store to self.calib
140 # TODO: collect the min and max stats for the input of layer 2 (aka the output of layer 2)
```

We've done the first layer for you!

Fill out QuantizerForClassyClassifier

```
188 def run_calibrated_quantized(self, x: torch.Tensor) -> torch.Tensor:
189     """Runs the pseudo-quantized ClassyClassifier on a 32-bit float tensor.
190
191     Quantizes the tensor to [0, 255] prior to forward pass, and then dequantizes back to float before return.
192
193     Args:
194         x (torch.Tensor): Input 32-bit float tensor.
195
196     Returns:
197         torch.Tensor: Output 32-bit float tensor.
198     """
199     x = copy.deepcopy(x)
200
201     # Use the calibration input stats for Layer 1 to quantize the tensor and run the layer
202     x, x_scalefactor, x_zeropoint = quantize_tensor(x, self.calibration_input_stats["layer1_conv"]["avg_min"], self.calibration_input_stats["l
203     x = self.int8_layer1_conv.run_quantized_layer(x, x_scalefactor, x_zeropoint, self.calibration_input_stats["layer3_conv"]["input_scalefacto
204     x = F.relu(x)
205     x = self.fp32_model.layer2_pool(x)
206
207     # TODO: As above, use the calibration input stats for the appropriate layers to complete the rest of the forward pass: layer 3-7 (DON'T DO
```

Accuracy should be similar!

Run the cells.

Your output should look like this image

```
10 outputs = quantizer.from_calibrated_quantized(torch.cat(outputs, 1))
11
12 predicted = torch.max(outputs, 1)
13 print("PREDICTED CLASSES from PSEUDO-QUANTIZED NETWORK should match original neural network outputs (i.e., cat, car, plane, pla
14 print(" " + " ".join([CIFAR10_CLASS_NAMES[idx] for idx in predicted]))
```

SAMPLE IMAGES from CIFAR-10:



CORRESPONDING GROUND TRUTH LABELS:

cat ship ship plane frog

PREDICTED CLASSES from PSEUDO-QUANTIZED NETWORK should match original neural network outputs (i.e., cat, car, plane, plane, frog):

cat car plane plane frog

Trying out a real quantization library

Run the cells.

Your output should look like this image

```
ACCURACY:  
ORIGINAL ClassyClassifier:  
Accuracy: 0.687, 6874 correct/10000 total images  
Model Size (kB): 539.1  
FX-QUANTIZED ClassyClassifier:  
Accuracy: 0.690, 6901 correct/10000 total images  
Model Size (kB): 160.4  
PSEUDO-QUANTIZED ClassyClassifier:  
Accuracy: 0.689, 6890 correct/10000 total images
```

Note that accuracy remains similar, but the FX-Quantized network is much smaller than the original!