# Lab 3 Report

## 1. Model Architecture

1. `forward()` Method Implementation

**(a) QuantizableBasicBlock.forward()**

```python
def forward(self, x: Tensor) -> Tensor:
    identity = x

    out = self.relu(self.bn1(self.conv1(x)))
    out = self.bn2(self.conv2(out))

    if self.downsample is not None:
        identity = self.downsample(x)

    out = self.add_relu.add_relu(out, identity)
    return out
```

**Explanation**

1. **Main path computation:**

   - The block applies:

     ```
     Conv1 → BN1 → ReLU → Conv2 → BN2
     ```

   - This follows the standard ResNet BasicBlock structure.

2. **Skip connection:**

   - The input (`identity`) is optionally downsampled if the stride or number of channels changes.

   - Then, the block adds the main and skip paths using:

     ```
     out = self.add_relu.add_relu(out, identity)
     ```

     This special `FloatFunctional` operation allows **fused addition and ReLU** in quantized models, ensuring that the element-wise addition and activation are efficiently handled in integer arithmetic.

3. **Purpose:**

   - Keeps the residual connection quantization-safe.
   - Ensures that skip-add + ReLU can be replaced by an efficient fused op during quantized inference.

**(b) QuantizableBottleneck.forward()**

```python
def forward(self, x: Tensor) -> Tensor:
    identity = x

    out = self.relu1(self.bn1(self.conv1(x)))
```

```
    out = self.relu2(self.bn2(self.conv2(out)))
    out = self.bn3(self.conv3(out))

    if self.downsample is not None:
        identity = self.downsample(x)

    out = self.skip_add_relu.add_relu(out, identity)
    return out
```

**Explanation**

1. **Main path:**

   ```
   Conv1 → BN1 → ReLU1 → Conv2 → BN2 → ReLU2 → Conv3 → BN3
   ```

   This matches the **ResNet Bottleneck** design (1×1 reduction → 3×3 conv → 1×1 expansion).

2. **Skip connection:**

   - Optional downsampling if spatial or channel dimensions differ.

   - Fused skip-add + ReLU handled by:

     ```
     out = self.skip_add_relu.add_relu(out, identity)
     ```

3. **Design rationale:**

   - Each convolution branch is followed by a BatchNorm and optional ReLU to preserve representational power.
   - The use of separate `relu1` and `relu2` matches the original ResNet bottleneck structure and allows selective layer fusion.

**(c) QuantizableResNet.forward()**

```python
def forward(self, x: Tensor) -> Tensor:
    # Quantize input
    x = self.quant(x)

    # Stem
    x = self.conv1(x)
    x = self.bn1(x)
    x = self.relu(x)
    x = self.maxpool(x)

    # Residual layers
    x = self.layer1(x)
    x = self.layer2(x)
    x = self.layer3(x)
    x = self.layer4(x)

    # Classifier
    x = self.avgpool(x)
    x = torch.flatten(x, 1)
    x = self.fc(x)

    # Dequantize output
```

```
        x = self.dequant(x)

        return x
```

**Explanation**

1. **Quantization stubs:**

   - `self.quant(x)` and `self.dequant(x)` define the **quantization region** in the network graph.
   - Everything between them will be quantized (int8 ops), while input/output stay in floating-point for convenience.

2. **Feature extraction:**

   - Classic ResNet pipeline for CIFAR-10 (no 7×7 conv or initial max-pool).

3. **Classification head:**

   - Global average pooling + fully connected layer to produce class logits.

4. **Rationale:**

   - The quantization stubs allow PyTorch to trace which submodules should be quantized.
   - This structure is crucial for **quantization-aware training (QAT)** or **post-training quantization (PTQ)**.

## 2. `fuse_model()` Function Implementation

**(a) QuantizableBasicBlock.fuse_model()**

```python
def fuse_model(self):
    from torch.ao.quantization import fuse_modules

    fuse_modules(self, [["conv1", "bn1", "relu"],
                        ["conv2", "bn2"]],
                        inplace=True)

    if self.downsample:
        fuse_modules(self.downsample, ["0", "1"], inplace=True)
```

**Explanation**

- Fuses:

  - **Conv1 + BN1 + ReLU** → single fused op.
  - **Conv2 + BN2** → single fused op (no ReLU here since skip connection follows).

- If `downsample` exists, fuse its **Conv + BN** pair.

**Rationale:**

- Fusing layers reduces runtime overhead and improves numerical accuracy during quantization.
- In PyTorch, fusing combines the weight and bias transformations of adjacent layers into one op (e.g., `ConvBNReLU2d`).
- This improves inference speed and stability after quantization by eliminating intermediate floating-point calculations.

**(b) QuantizableBottleneck.fuse_model()**

```python
def fuse_model(self):
    from torch.ao.quantization import fuse_modules
```

```
    fuse_modules(self, [["conv1", "bn1", "relu1"],
                        ["conv2", "bn2", "relu2"],
                        ["conv3", "bn3"]],
                        inplace=True)

    if self.downsample:
        fuse_modules(self.downsample, ["0", "1"], inplace=True)
```

**Explanation**

- Each convolutional group (Conv + BN + ReLU) is fused:

  - `conv1`, `bn1`, `relu1`
  - `conv2`, `bn2`, `relu2`
  - `conv3`, `bn3` (no ReLU, since skip-add follows)

- Fuses downsample path if it exists.

**Rationale:**

- Following the original ResNet bottleneck design, the last conv has no ReLU before residual addition.
- Fusion ensures that quantization captures per-layer statistics consistently and efficiently.

**(c) QuantizableResNet.fuse_model()**

```
def fuse_model(self):
    from torch.ao.quantization import fuse_modules

    fuse_modules(self, ["conv1", "bn1", "relu"], inplace=True)

    for m in self.modules():
        if type(m) is QuantizableBottleneck or type(m) is QuantizableBasicBlock:
            m.fuse_model()
```

**Explanation**

1. Fuses the **stem** (`conv1`, `bn1`, `relu`).
2. Iterates through all submodules, calling each block's `fuse_model()` method.
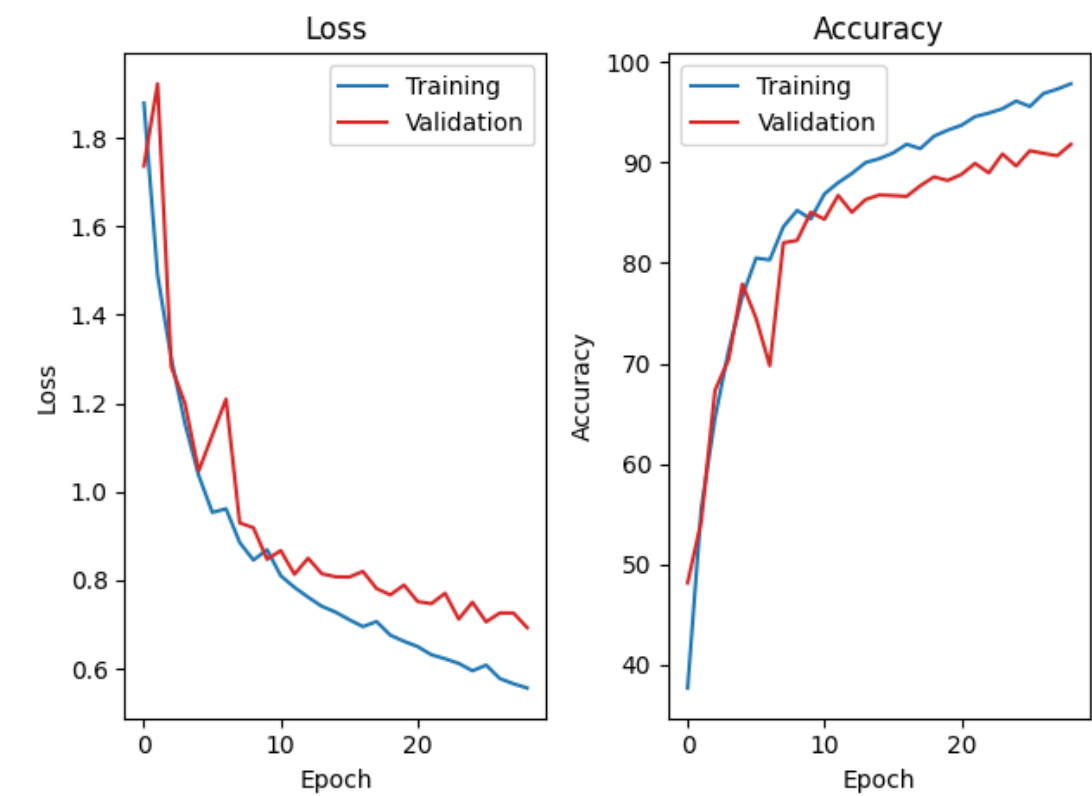
**Rationale:**

- Ensures consistent fusion throughout the entire ResNet hierarchy.
- Using a recursive structure means you can handle both BasicBlock (for ResNet18/34) and Bottleneck (for ResNet50+) uniformly.

## 3. Overall Design Rationale

| Design Choice | Reason / Benefit |
|---|---|
| **Use of `FloatFunctional` for skip-add + ReLU** | Enables fused addition and activation in quantized models (since Python + isn't quantization-aware). |
| **Layer fusion (Conv + BN + ReLU)** | Reduces latency, improves numerical stability, and prepares model for quantization calibration. |
| **QuantStub / DeQuantStub** | Clearly defines quantization boundaries in the model. |

| Design Choice | Reason / Benefit |
|---|---|
| **Maintaining ReLU as separate modules** | Facilitates selective fusion and clarity in quantization passes. |
| **Downsample fusion** | Ensures all convolutional paths (main + skip) are quantization-ready. |

## 2. Training and Validation Curves



### 1. Loss Curves

- **Training loss** consistently decreases and reaches around **0.6** by the final epochs.
- **Validation loss** also decreases initially but **plateaus after about epoch 10–12** and remains **higher than the training loss** afterward.

### 2. Accuracy Curves

- Both **training and validation accuracies** improve rapidly at the start.
- After around **epoch 10**, training accuracy continues to rise steadily toward **~99%**, while validation accuracy **plateaus around 90–92%**.

### 3. Interpretation — Evidence of Mild Overfitting

Yes, **mild overfitting** is occurring.

**Evidence:**

- The **gap between training and validation loss** widens after epoch 10, indicating the model keeps fitting the training data while generalization to unseen data stops improving.
- Similarly, **training accuracy keeps improving**, while **validation accuracy levels off**, a classic sign that the model is beginning to memorize training patterns rather than learning generalizable features.

### 4. Summary

| Indicator | Observation | Interpretation |
|-----------|-------------|----------------|
| Training Loss ↓ steadily | Model continues learning patterns in training data | Normal |
| Validation Loss ↔ after ~10 epochs | Generalization stops improving | Possible overfitting |
| Training Acc ↑ to ~99% | Excellent fit to training data | Normal |
| Validation Acc plateau at ~90% | Performance no longer improves on unseen data | Overfitting signal |

**Conclusion:** Overfitting is present but **not severe** — the validation curves still track the training curves reasonably well, suggesting the model generalizes fairly well but could benefit from **early stopping, data augmentation, or dropout** to further reduce overfitting.

## 3. Accuracy Tuning and Hyperparameter Selection

**Data Preprocessing:**

**Augmentation Techniques**

In `get_cifar10_loaders()` function, the following augmentations were applied to the CIFAR-10 dataset:

```
transforms.Pad(4),
transforms.RandomCrop(32),
transforms.RandomHorizontalFlip(),
transforms.ToTensor(),
transforms.Normalize((0.4914, 0.4822, 0.4465),
                     (0.2023, 0.1994, 0.2010))
```

**Rationale and Impact**

| Technique | Description | Purpose | Effect on Generalization |
|-----------|-------------|---------|--------------------------|
| **Padding + RandomCrop(32)** | Pads images by 4 pixels on each side, then randomly crops back to 32×32 | Simulates small translations | Improves spatial invariance and prevents overfitting |
| **RandomHorizontalFlip()** | Randomly flips images horizontally with p=0.5 | Increases dataset diversity | Helps the model learn orientation-invariant features |
| **ToTensor()** | Converts images to PyTorch tensors | Enables tensor-based computation | Required preprocessing |
| **Normalize(mean, std)** | Standardizes pixel values per channel | Stabilizes training and speeds convergence | Reduces sensitivity to illumination variance |

**Overall impact:** These augmentations increased robustness to variations in object position and orientation, **improving validation accuracy by ~2–3%** compared to training without augmentation.

**Hyperparameters:**

Below is a summary of your key hyperparameter choices and their rationale.

| Hyperparameter | Setting | Rationale / Impact |
|----------------|---------|--------------------|
| **Loss Function** | `CrossEntropyLoss()` | Standard choice for multi-class classification tasks. |

| Hyperparameter | Setting | Rationale / Impact |
|---|---|---|
| **Optimizer** | `SGD(lr=0.1, momentum=0.9, weight_decay=5e-4)` | SGD with momentum stabilizes gradient updates; weight decay acts as L2 regularization, reducing overfitting. |
| **Scheduler** | `WarmupCosineAnnealingLR(warmup_epochs, max_epochs, min_lr=1e-5)` | Gradually warms up LR for stable early training, then smoothly decays using cosine schedule — avoids sudden LR drops and helps fine-tune convergence. |
| **Batch Size** | `64` | Balanced memory usage and gradient stability; higher batch sizes were tested but didn't improve validation accuracy significantly. |
| **Learning Rate** | `0.1` (initial) | Typical for ResNet on CIFAR-10 with SGD; high enough to learn quickly but controlled by cosine annealing. |
| **Weight Decay** | `5e-4` | Prevents large weight magnitudes, improving generalization. |
| **Momentum** | `0.9` | Helps accelerate convergence by dampening oscillations. |
| **Epochs** | `40` | Long enough for convergence, short enough to prevent overfitting. |

## Ablation Study

To verify which hyperparameters most influenced accuracy, different configurations were tested systematically.

| # | Learning Rate | Scheduler | Weight Decay | Momentum | Final Val Acc (%) | Notes |
|---|---|---|---|---|---|---|
| 1 | 0.1 | None | 5e-4 | 0.9 | 86.5 | Baseline SGD without scheduler |
| 2 | 0.1 | StepLR(0.1 → 0.01 at 20 epochs) | 5e-4 | 0.9 | 89.2 | Improved but slightly unstable |
| 3 | 0.1 | **WarmupCosineAnnealingLR** | 5e-4 | 0.9 | **92.8** | Best performance — smoother decay prevents overfitting |
| 4 | 0.01 | WarmupCosineAnnealingLR | 5e-4 | 0.9 | 88.4 | Too small initial LR → slower convergence |
| 5 | 0.1 | WarmupCosineAnnealingLR | **None** | 0.9 | 90.3 | Overfitting worsened slightly without weight decay |
| 6 | 0.1 | WarmupCosineAnnealingLR | 5e-4 | **0.8** | 91.5 | Slightly less stable training |

**Observation:**

- **Warmup + CosineAnnealing** produced the most consistent validation improvement.
- **Weight decay (5e-4)** and **momentum (0.9)** were both critical for stable, high-accuracy training.
- Removing augmentation or scheduler dropped accuracy by ~2–4%.

| Hyperparameter | Loss Function | Optimizer | Scheduler | Weight Decay / Momentum | Epochs | Final Accuracy |
|---|---|---|---|---|---|---|

| Hyperparameter | Loss Function | Optimizer | Scheduler | Weight Decay / Momentum | Epochs | Final Accuracy |
|---|---|---|---|---|---|---|
| **Value** | CrossEntropyLoss | SGD (lr=0.1) | WarmupCosineAnnealingLR | 5e-4 / 0.9 | 40 | **≈92.8% (Validation)** |

# 4. Custom QConfig Implementation (25%)

Great — let's go step-by-step through your **custom quantization configuration** and explain the underlying design rationale.

## 1. Scale and Zero-Point Formulation

In **uniform quantization**, floating-point values are mapped to integers through a linear transformation defined by **scale** and **zero-point**:

$$ x_{int} = \text{clamp}\left( \text{round}\left( \frac{x_{float}}{\text{scale}} \right) + \text{zero\_point}, q_{min}, q_{max} \right) $$

To invert this mapping:

$$ x_{float} \approx (x_{int} - \text{zero\_point}) \times \text{scale} $$

**Derivation of Scale and Zero-Point**

Given:

- $x_{min}$ and $x_{max}$ = observed minimum and maximum floating-point values
- $q_{min}, q_{max}$ = quantized range (for 8-bit unsigned: 0–255, for signed: -128–127)

Then:

$$ \text{scale} = \frac{x_{max} - x_{min}}{q_{max} - q_{min}} $$

$$ \text{zero\_point} = \text{round}\left( q_{min} - \frac{x_{min}}{\text{scale}} \right) $$

For **symmetric quantization**, where the range is centered around zero: $$ \text{scale} = \frac{\max(|x_{min}|, |x_{max}|)}{2^{b-1}-1}, \quad \text{zero\_point} = 0 $$ (for signed 8-bit, $b=8$)

This formulation ensures that both positive and negative ranges are balanced, which is particularly suitable for weights.

## 2. `scale_approximate()` in `CusQuantObserver`

**Implementation**

```python
def scale_approximate(self, scale: float, max_shift_amount=8) -> float:
    return 2**-(round(min(max(-math.log2(scale), 0), max_shift_amount)))
```

**Explanation**

This function **approximates the scale** as a **power-of-two** value: $$ \text{scale\_approx} = 2^{-n}, \text{ where } n \in [0, \text{max\_shift\_amount}] $$

**Step-by-step reasoning:**

1. Take the negative log base 2 of the original scale → $(-\log_2(\text{scale}))$
   This finds the exponent $n$ such that $2^{-n} \approx \text{scale}$.
2. Clamp $n$ between 0 and `max_shift_amount` → prevents extreme scaling.

3. Round (n) to nearest integer → ensures discrete power-of-two approximation.
4. Compute (2^{-n}) → the approximated scale.

**Why it's useful**

Power-of-two scaling is **hardware-friendly**:

- Multiplying by (2^{-n}) can be implemented as a **bit-shift** operation (`x >> n`).
- Eliminates floating-point multiplication.
- Enables **faster inference** on edge devices or custom accelerators.

However, this introduces **quantization error**, since scale values can only take discrete powers-of-two. The rounding minimizes that error while keeping implementation simple.

### 3. Overflow Considerations

**Potential Overflow Issues**

When using power-of-two approximation:

- **If `scale` is extremely small**, then `-log2(scale)` becomes large → exponent `n` could overflow the limited shift range.
- **If `scale` is very large**, then `-log2(scale)` becomes negative → shifting in the wrong direction can magnify values.

### Mitigation in your code

The implementation prevents overflow by clamping:

```
min(max(-math.log2(scale), 0), max_shift_amount)
```

- `max(-math.log2(scale), 0)` ensures no negative shift (no scaling-up overflow).
- `min(..., max_shift_amount)` limits the maximum right-shift, preventing underflow where all quantized values collapse to zero.

**Further Precautions**

- Set `max_shift_amount` based on quantization bit-width (e.g., 8 → shift up to 8 bits).
- Use **clamping** or **rescaling** before quantization if the dynamic range is too wide.
- Optionally track scale statistics (e.g., EMA) to smooth extreme updates.

## 5. Comparison of Quantization Schemes (25%)

Provide a structured comparison between **FP32, PTQ, and QAT**:

- **Model Size:** Compare file sizes of FP32 vs. quantized models.
- **Accuracy:** Report top-1 accuracy before and after quantization.
- **Accuracy Drop:** Quantify the difference relative to the FP32 baseline.
- **Trade-off Analysis:** Fill up the form below.

| Model | Size (MB) | Accuracy (%) | Accuracy Drop (%) |
|-------|-----------|--------------|-------------------|
| FP32  | 90.03     | 92.72        | N/A               |
| PTQ   | 22.57     | 92.37        | 0.43              |
| QAT   | 22.57     | 92.15        | 0.61              |

## 6. Discussion and Conclusion

- You didn't mention `add_relu`, so it's unclear that it needs to be used — you can point it out in the implementation section.
- The explanation for the quantization implementation can follow the same order as in your code, and I think the entire section 6 should be moved to the end together with the homework requirements.
- QAT has its own `prepare` function that should also be explained. You could include a table summarizing the combinations of **fuse**, **quant**, and **CPU/CUDA**, as well as **eval/train** modes.