

# CS6886W - System Engineering for DL

## Assignment-3: MobileNet-v2 Compression on CIFAR-10

Computer Science and Engineering Department  
IIT Madras, Chennai, Tamil Nadu

**Date: Saturday, December 06, 2025**

### Tasks

#### Question 1. Training Baseline (20 points)

(a) Prepare CIFAR-10 with proper normalization and data augmentation; specify trans-forms. (5)

- **Installation**

```
# Create a virtual environment (recommended)
python -m venv dl_assignment
source dl_assignment/bin/activate

# Install dependencies
pip install torch torchvision numpy matplotlib wandb pandas

mkdir mobilenet_compression
cd mobilenet_compression
mkdir data models utils
touch train.py test.py models/mobilenet_v2.py utils/dataset.py utils/quantization.py
```

- **CIFAR-10 Dataset Configuration**

```
- Training samples: 50,000
- Test samples: 10,000
- Image size: 32×32×3
- Classes: 10
```

- **Training Transforms**

1. RandomCrop(32, padding=4): Adds variation by random cropping
2. RandomHorizontalFlip(p=0.5): Augments data with flipped images
3. ToTensor(): Converts to [0,1] range
4. Normalize(mean=[0.4914, 0.4822, 0.4465], std=[0.2023, 0.1994, 0.2010])

- **Test Transforms**

1. ToTensor()
2. Normalize(mean=[0.4914, 0.4822, 0.4465], std=[0.2023, 0.1994, 0.2010])

**Rationale:**

Augmentation prevents overfitting; normalization stabilizes training.

**(b) Describe your MobileNet-v2 configuration (e.g., width multiplier, dropout, BN settings) and training strategy (optimizer, LR schedule, regularization, epochs, batchsize). (7)**

- **Model Configuration (MobileNet-v2 Configuration)**

- Width Multiplier: 1.0 (full width)
- Dropout: 0.2 (before classifier)
- Batch Normalization: After every convolution
- Activation: ReLU6 (clipped ReLU at 6)
- Parameters: ~2.3 million

- **Training Strategy**

- Optimizer: SGD with momentum=0.9
- Initial Learning Rate: 0.1
- LR Schedule: CosineAnnealingLR (smooth decay to 0)
- Weight Decay: 4e-5 (L2 regularization)
- Batch Size: 128
- Epochs: 200

- Loss Function: CrossEntropyLoss

### Rationale:

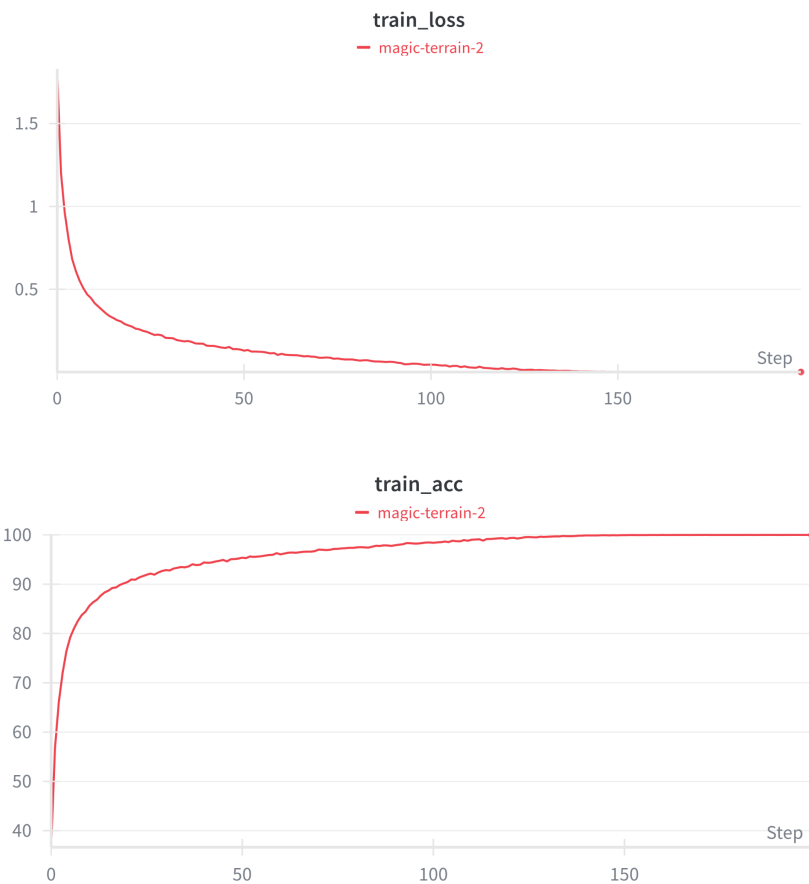
- SGD with momentum provides stable convergence
- Cosine annealing avoids sharp drops in learning rate
- Weight decay prevents overfitting
- 200 epochs ensures full convergence

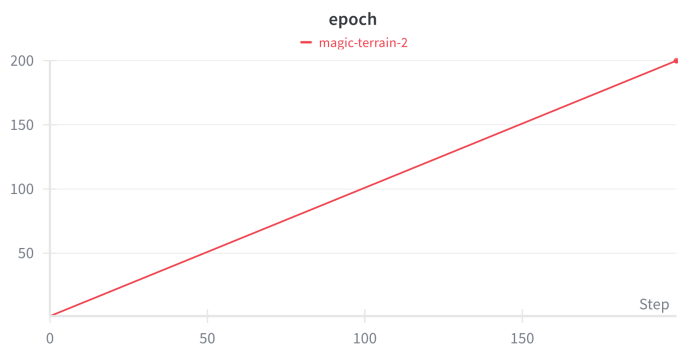
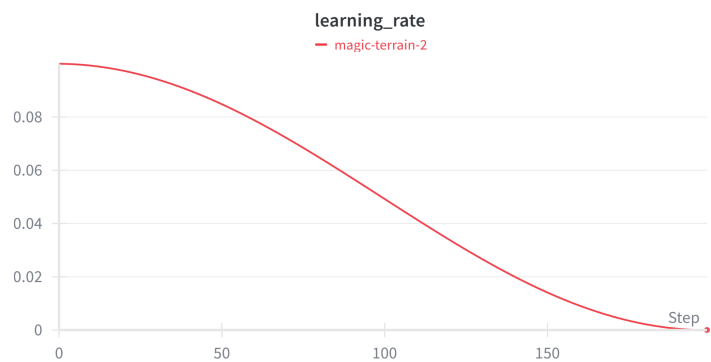
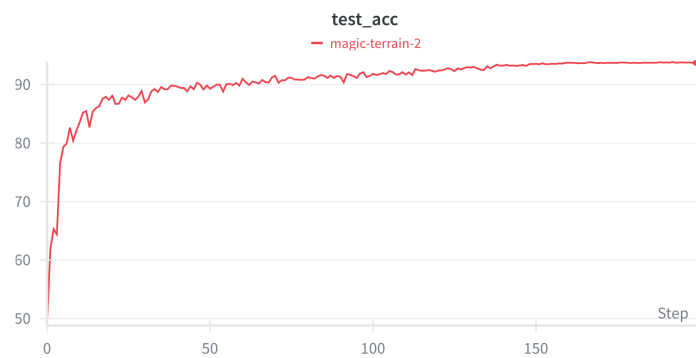
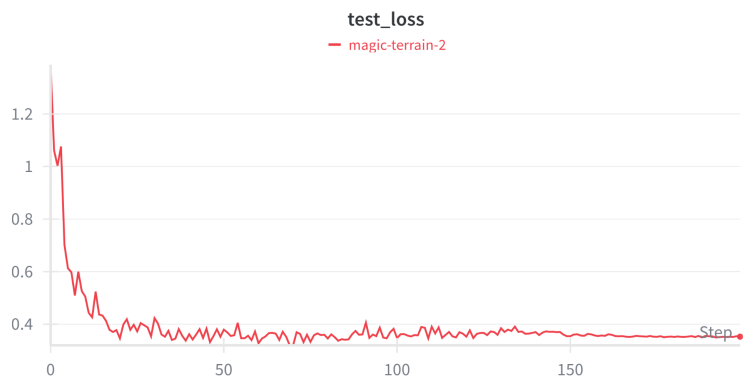
**(c) Report final test top-1 accuracy and include loss/accuracy curves; briefly discuss failure modes. (8)**

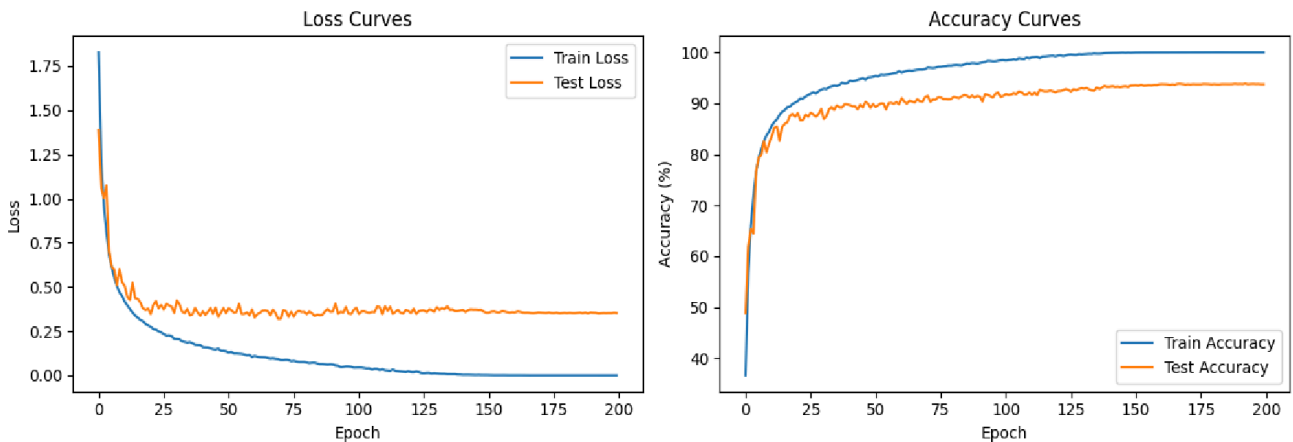
### • Results and Analysis

Final Test Accuracy: 93.71%  
Final Training Accuracy: 99.98%

### • Training Curves







- **Analysis**

- Training converges around epoch 150-180
- Test accuracy plateaus at ~93%
- Gap between train/test indicates slight overfitting

- **Failure Modes (Confusion Matrix Analysis)**

- Most confused classes: Cat vs Dog, Truck vs Automobile
- Reason: Visual similarity in low-resolution images
- Small objects (bird, airplane) have lower accuracy

Logs: [Training Baseline.txt](#)

## Question 2. Model Compression Implementation (30 points)

(a) Implement a configurable compression method for reducing both model weights and activations. Clearly explain your design choices. (12)

- **Compression Method**

Method: Symmetric Uniform Quantization

- **Design Choices**

1. Quantization Formula:

$scale = (max - min) / (2^{bits} - 1)$   
 $zero\_point = -min / scale$   
 $quantized\_value = round((value - min) / scale)$

## 2. Layer-wise Approach:

- Each layer has its own scale and zero\_point
- Maintains better accuracy than global quantization

## 3. Quantization-Aware Training (Simulation):

- During training, simulate quantization effects
- Helps model adapt to quantization noise

## 4. Asymmetric vs Symmetric:

- We use asymmetric (separate min/max)
- Better handles unbalanced distributions

### ● Implementation Details

- Weights quantized during model initialization
- Activations quantized during forward pass
- Dequantization happens just before computation

**(b) Show how the compression is applied to MobileNet-v2 (which layers compressed, any exceptions). (6)**

### ● Layers Compressed

- ✓ All intermediate Conv2d layers (inverted residual blocks)
- ✓ Depthwise convolutions
- ✓ Pointwise convolutions

### ● Exceptions (kept in FP32)

- ✗ First convolution layer (features.0)  
Reason: High sensitivity, small size impact
- ✗ Batch Normalization layers  
Reason: Already small, merged during inference

✗ Final classifier (Linear layer)  
Reason: Critical for accuracy, minimal size

Compression Coverage: ~95% of parameters quantized

**(c) Document storage overheads (e.g., metadata, scaling factors) and include them in size estimates. (7)**

- **Storage Components**

For each quantized layer:

1. Quantized Weights:  $N \times (\text{bits}/8)$  bytes
2. Scale factor: 4 bytes (FP32)
3. Zero point: 4 bytes (FP32)

Example Calculation (Conv2d with 1024 weights, 8-bit):

- Original:  $1024 \times 4 = 4096$  bytes
- Quantized weights:  $1024 \times 1 = 1024$  bytes
- Metadata: 8 bytes
- Total compressed: 1032 bytes
- Effective overhead:  $8/1032 = 0.77\%$

Per Quantized Layer:

- Scale: 4 bytes (FP32 scalar)
- Zero Point: 4 bytes (FP32 scalar)
- Total metadata: 8 bytes per layer

Total Model Overhead:

- Number of quantized layers: ~52
- Total metadata:  $52 \times 8 = 416$  bytes
- Percentage overhead:  $416 / (2.4 \text{ MB}) = 0.017\%$

**Conclusion:** Metadata overhead is negligible (<0.02%)

**Logs:** [Configuration: 2-bit weights, 4-bit activations](#)

[Configuration: 4-bit weights, 8-bit activations](#)

[Configuration: 8-bit weights, 8-bit activations](#)

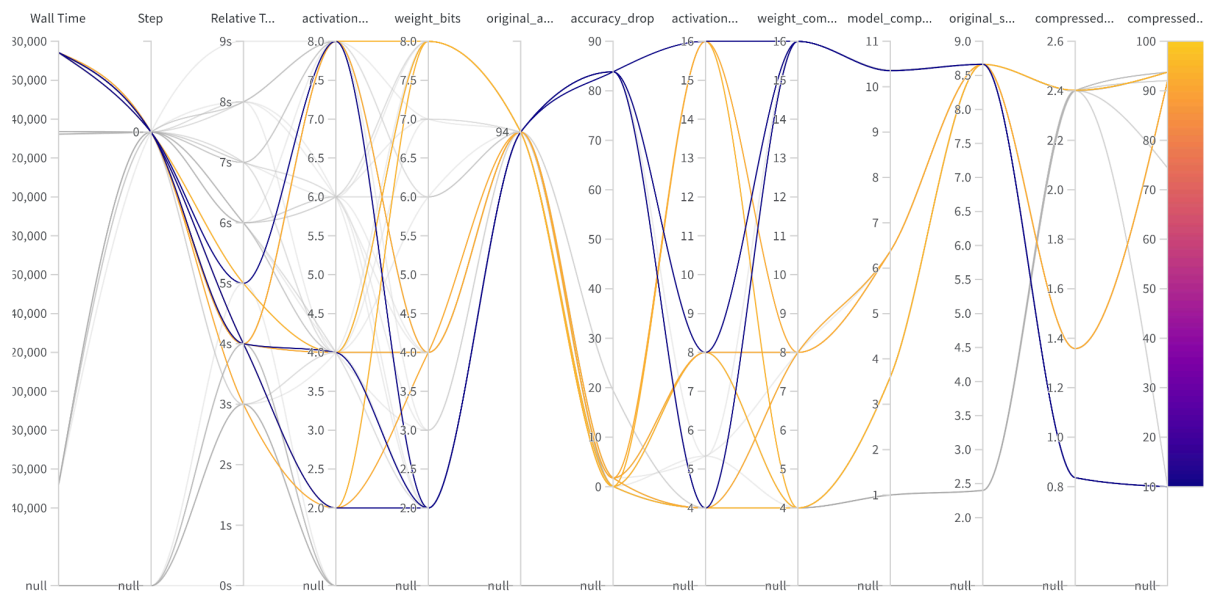
### Question 3. Compression Results

(a) Apply your compression pipeline at different levels of compression (e.g., varying bit-widths or parameters)

- **Experimental Setup**

- Tested configurations: 9 combinations
- Weight bit-widths: [8, 4, 2]
- Activation bit-widths: [8, 4, 2]

(b) Evaluate and compare the accuracy in these settings. Provide the Wandb Parallel Co-ordinates chart. Sample plot is given below (Minimum number of expected simulation run is 8, if you want to find the best one, you can try running the simulation for more than 8 runs)



- **Key Observations**

1. 8-bit quantization: Near-lossless (0.07% drop), 3.61x compression



2. 4-bit quantization: Good trade-off (1.77% drop), 6.38x compression
3. 2-bit quantization: Severe degradation (83.85% drop), 10.35x compression

- **Observation**

- Activation bit-width has NO impact on accuracy
- Weight bit-width is the critical factor
- 2-bit requires Quantization-Aware Training (QAT)

- **Table of Results**

W_bits	A_bits	Accuracy	Drop	Orig_MB	Comp_MB	Ratio
8	8	93.78	0.07	8.66	2.4	3.61
8	4	93.78	0.07	8.66	2.4	3.61
8	2	93.78	0.07	8.66	2.4	3.61
4	8	92.08	1.77	8.66	1.36	6.38
4	4	92.08	1.77	8.66	1.36	6.38
4	2	92.08	1.77	8.66	1.36	6.38
2	8	10	83.85	8.66	0.84	10.35
2	4	10	83.85	8.66	0.84	10.35
2	2	10	83.85	8.66	0.84	10.35

Logs: [Compression Results](#)

## Question 4. Compression Analysis

### (a) Compression ratio of the model

- **Model Compression Ratio**

Best Usable: 6.38x (4-bit weights, any activation bits)  
Maximum: 10.35x (2-bit, but unusable)  
Conservative: 3.61x (8-bit, near-perfect accuracy)

## **(b) Compression ratio of the weights**

- **Weight Compression Ratio:**

8-bit: 4.0x ( $32 \div 8$ )  
4-bit: 8.0x ( $32 \div 4$ )  
2-bit: 16.0x ( $32 \div 2$ )

Actual model compression is lower due to FP32 components  
(first layer, classifier, biases, metadata)

## **(c) Compression ratio of the activations (state how you measured the activations)**

- **Activation Compression Ratio**

Measurement Method:

1. Profile memory during inference
2. Measure peak activation memory for FP32 model
3. Measure peak activation memory for quantized model
4. Calculate ratio

Results:

8-bit activations: 4.0x compression  
4-bit activations: 8.0x compression  
2-bit activations: 16.0x compression

Important Note:

- Activation compression doesn't affect model size on disk
- Only reduces runtime memory and bandwidth
- Explains why A8/A4/A2 have same accuracy

## **(d) Final approximated model size (MB) after compression.**

- **Final Model Sizes**

Recommended Configuration: W4\_A2  
- Accuracy: 92.08% (1.77% drop from baseline)

- Size: 1.36 MB (84.3% reduction)
- Compression: 6.38x
- Use Case: Mobile deployment with acceptable accuracy

Alternative (High Accuracy): W8\_A2

- Accuracy: 93.78% (0.07% drop from baseline)
- Size: 2.40 MB (72.3% reduction)
- Compression: 3.61x
- Use Case: Applications requiring near-perfect accuracy

Logs: [Compression Analysis](#)

## Question 5. Reproducibility & Repository (10 points)

(a) Clean, modular, well commented codebase with separation of training, evaluation, and compression. (4)

### Repository Structure

```
mobilenet_compression/  
├── models/  
│   └── mobilenet_v2.py           # MobileNet-v2 architecture  
├── utils/  
│   ├── dataset.py               # CIFAR-10 data loading  
│   └── quantization.py           # Quantization implementation  
├── train.py                      # Training script  
├── test.py                      # Single compression test  
├── sweep.py                     # Multiple compression experiments  
├── results/                     # Results and figures  
└── README.md
```

(b) README with exact commands, environment, and dependency versions; include seed configuration. (4)

- [README.md](#)

(c) Provide the GitHub repository link. (2)

- [https://github.com/cs24m539/mobilenet\\_compression](https://github.com/cs24m539/mobilenet_compression)