
Comparative Analysis of Normalization Techniques in Convolutional Neural Networks: Batch, Layer, and Weight Normalization

Karl Reger
Northern Arizona University
Flagstaff, AZ 86011
kcr28@nau.edu

Abstract

This paper presents a comprehensive empirical comparison of three prominent normalization techniques for deep learning: Batch Normalization (BatchNorm), Layer Normalization (LayerNorm), and Weight Normalization (WeightNorm). We implement these methods from first principles using TensorFlow operations and evaluate their performance on the Fashion-MNIST classification task using a convolutional neural network architecture. Our study examines convergence speed, final accuracy, batch size sensitivity, and the mathematical correctness of custom implementations against standard library functions. Our experiments reveal that LayerNorm achieved the highest test accuracy (92.42%), outperforming both BatchNorm (91.50%) and the baseline (92.16%) at standard batch size. Surprisingly, BatchNorm exhibited significant overfitting at batch size 128, with its performance actually improving at smaller batch sizes due to enhanced regularization effects. We provide detailed analysis of when each normalization method excels and offer practical recommendations for practitioners. Our results demonstrate that Layer Normalization provides the best balance of accuracy and robustness, with particular emphasis on explaining why LayerNorm outperforms BatchNorm through both theoretical analysis and empirical validation.

1 Introduction

1.1 Motivation

Training deep neural networks effectively requires careful management of internal activations throughout the network. As signals propagate through multiple layers, the distribution of activations can shift dramatically during training—a phenomenon known as *internal covariate shift* [1]. This shifting distribution complicates optimization, requiring careful tuning of learning rates and initialization schemes, and can significantly slow convergence.

Normalization techniques address this challenge by stabilizing activation distributions during training. By normalizing layer inputs or weights, these methods reduce the sensitivity of the network to parameter initialization and enable the use of higher learning rates, thereby accelerating convergence. Three prominent normalization approaches have emerged: Batch Normalization [1], Layer Normalization [2], and Weight Normalization [3]. Each method takes a different approach to the normalization problem, with distinct mathematical formulations, computational characteristics, and applicability to different scenarios.

1.2 Problem Statement

While these normalization techniques are widely used in practice, understanding their relative strengths and limitations requires careful empirical comparison. Key questions remain:

- How do these methods compare in terms of convergence speed and final performance?
- What is the impact of batch size on each normalization technique?
- Can custom implementations achieve numerical parity with optimized library functions?
- Under what conditions does each method excel or struggle?
- Why does Layer Normalization often outperform Batch Normalization in certain scenarios?

This work addresses these questions through systematic experimentation on a standard benchmark task.

1.3 Objectives

The primary objectives of this study are:

1. **Implementation:** Develop correct, from-scratch implementations of BatchNorm, LayerNorm, and WeightNorm using basic TensorFlow operations.
2. **Verification:** Validate custom implementations against TensorFlow’s built-in functions through numerical comparison of forward and backward passes.
3. **Performance Comparison:** Evaluate all three methods against a baseline (no normalization) across multiple batch sizes.
4. **Analysis:** Provide both empirical results and theoretical insights into when and why each method succeeds.
5. **Practical Guidance:** Offer clear recommendations for practitioners choosing normalization strategies.

1.4 Contributions

This work makes the following contributions:

- Complete, verified implementations of three normalization techniques with numerical validation against standard libraries.
- Comprehensive experimental evaluation on Fashion-MNIST across 7 different experimental configurations.
- Detailed analysis of batch size sensitivity, revealing that BatchNorm’s performance can actually improve at smaller batch sizes (92.08% at $bs=4$ vs 91.50% at $bs=128$) due to enhanced regularization effects that prevent overfitting, while LayerNorm maintains consistent performance with only 0.6 percentage point degradation.
- Mathematical and empirical explanation of why Layer Normalization demonstrates superior performance to Batch Normalization in certain regimes.
- Practical recommendations backed by both theory and experimental evidence.

2 Background

This section provides the mathematical foundations for the three normalization techniques examined in this study.

2.1 Batch Normalization

Batch Normalization [1] normalizes activations using statistics computed over a mini-batch. For a layer with input x and mini-batch $\mathcal{B} = \{x_1, \dots, x_m\}$ of size m , BatchNorm computes:

$$\mu_{\mathcal{B}} = \frac{1}{m} \sum_{i=1}^m x_i$$

$$\sigma_{\mathcal{B}}^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2$$

$$\hat{x}_i = \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$$

$$y_i = \gamma \hat{x}_i + \beta$$

where ϵ is a small constant for numerical stability (typically 10^{-5}), and γ (scale) and β (shift) are learnable parameters that allow the network to undo the normalization if needed.

During training, running averages of the batch statistics are maintained:

$$\mu_{\text{running}} = (1 - \alpha)\mu_{\text{running}} + \alpha\mu_{\mathcal{B}}$$

$$\sigma_{\text{running}}^2 = (1 - \alpha)\sigma_{\text{running}}^2 + \alpha\sigma_{\mathcal{B}}^2$$

where α is the momentum parameter (typically 0.99). At inference time, the running statistics replace the batch statistics to ensure deterministic outputs.

Key Properties:

- Reduces internal covariate shift
- Enables higher learning rates
- Acts as a regularizer (due to noise from batch statistics)
- Performance depends on batch size
- Behavior differs between training and inference

2.2 Layer Normalization

Layer Normalization [2] computes normalization statistics across all neurons in a layer for each individual training example, rather than across the batch. For a layer with H hidden units and input $x = [x_1, \dots, x_H]$, LayerNorm computes:

$$\mu = \frac{1}{H} \sum_{i=1}^H x_i$$

$$\sigma^2 = \frac{1}{H} \sum_{i=1}^H (x_i - \mu)^2$$

$$\hat{x}_i = \frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$y_i = \gamma_i \hat{x}_i + \beta_i$$

Note that unlike BatchNorm, LayerNorm typically uses per-neuron scale and shift parameters (γ_i and β_i for each hidden unit i).

Key Properties:

- Independent of batch size (statistics computed per example)

- Identical behavior during training and inference
- Well-suited for recurrent networks and variable-length sequences
- No running statistics required
- Effective for small batch training

2.3 Weight Normalization

Weight Normalization [3] reparameterizes the weight vectors rather than normalizing activations. For a weight vector $w \in \mathbb{R}^k$, WeightNorm decomposes it into magnitude and direction:

$$w = \frac{g}{\|v\|} v$$

where $v \in \mathbb{R}^k$ is a parameter vector representing direction, $g \in \mathbb{R}$ is a scalar parameter representing magnitude, and $\|v\| = \sqrt{\sum_{i=1}^k v_i^2}$ is the Euclidean norm.

This reparameterization enforces:

$$\|w\| = g$$

The gradient computation follows:

$$\begin{aligned}\nabla_g \mathcal{L} &= \frac{\nabla_w \mathcal{L} \cdot v}{\|v\|} \\ \nabla_v \mathcal{L} &= \frac{g}{\|v\|} \nabla_w \mathcal{L} - \frac{g \nabla_g \mathcal{L}}{\|v\|^2} v\end{aligned}$$

The second term in $\nabla_v \mathcal{L}$ projects the gradient away from the current weight direction, which stabilizes training.

Key Properties:

- Decouples weight magnitude from direction
- Fewer parameters than BatchNorm or LayerNorm
- No batch dependence
- Improves conditioning of optimization
- Can be combined with other normalization techniques

3 Methodology

3.1 Dataset

We evaluate all normalization methods on the Fashion-MNIST dataset [4], which consists of 70,000 grayscale images of clothing items:

- **Training set:** 60,000 images
- **Test set:** 10,000 images
- **Image size:** 28×28 pixels
- **Classes:** 10 (T-shirt, Trouser, Pullover, Dress, Coat, Sandal, Shirt, Sneaker, Bag, Ankle boot)
- **Preprocessing:** Pixel values normalized to $[0, 1]$ range

Fashion-MNIST was chosen as it provides a more challenging task than MNIST while remaining computationally tractable for extensive experimentation.

3.2 Network Architecture

We employ a convolutional neural network with the following architecture:

Table 1: Network architecture with layer-wise output shapes. Normalization layers are inserted after convolutional and dense layers but before activation functions.

Layer Type	Configuration	Output Shape
Input	—	(B, 28, 28, 1)
Conv2D-1	$5 \times 5 \times 30$, stride=1, padding=same	(B, 28, 28, 30)
Norm-1	BatchNorm/LayerNorm/WeightNorm	(B, 28, 28, 30)
ReLU	—	(B, 28, 28, 30)
MaxPool2D	2×2 , stride=2	(B, 14, 14, 30)
Conv2D-2	$5 \times 5 \times 60$, stride=1, padding=same	(B, 14, 14, 60)
Norm-2	BatchNorm/LayerNorm/WeightNorm	(B, 14, 14, 60)
ReLU	—	(B, 14, 14, 60)
MaxPool2D	2×2 , stride=2	(B, 7, 7, 60)
Flatten	—	(B, 2940)
Dense-1	100 units	(B, 100)
Norm-3	BatchNorm/LayerNorm/WeightNorm	(B, 100)
ReLU	—	(B, 100)
Dense-2	10 units	(B, 10)
Softmax	—	(B, 10)

where B denotes the batch size. This architecture was chosen to be sufficiently deep to benefit from normalization while remaining simple enough for clear analysis.

3.3 Implementation Details

All normalization methods were implemented from scratch using basic TensorFlow operations (no use of `tf.nn.batch_normalization` or similar high-level functions). The implementations follow the mathematical formulations in Section 2 exactly.

3.3.1 BatchNorm Implementation

For convolutional layers, BatchNorm normalizes across (B, H, W) dimensions for each channel. For dense layers, normalization occurs across the batch dimension for each feature. Running statistics use exponential moving average with momentum $\alpha = 0.99$.

3.3.2 LayerNorm Implementation

For convolutional layers, LayerNorm normalizes across (H, W, C) dimensions for each example. For dense layers, normalization occurs across all features for each example. No running statistics are needed.

3.3.3 WeightNorm Implementation

WeightNorm is implemented as custom Keras layers (`WeightNormDense` and `WeightNormConv2D`) that maintain v and g parameters and compute $w = \frac{g}{\|v\|}v$ during forward pass. The v parameters are initialized using He initialization, and g is initialized to $\|v\|$ to start with identity-like behavior.

3.4 Experimental Setup

We conducted 7 experiments to evaluate normalization methods under different conditions:

Hyperparameters:

The random seed 394 was chosen as the ASCII values of "Karl" summed together, ensuring reproducibility.

Table 2: Experimental configurations. All experiments use identical hyperparameters except for batch size and normalization method.

Experiment Name	Normalization	Purpose
baseline_bs128	None	Baseline comparison
batchnorm_bs128	Batch Normalization	Standard batch size
layernorm_bs128	Layer Normalization	Standard batch size
weightnorm_bs128	Weight Normalization	Standard batch size
batchnorm_bs4	Batch Normalization	Small batch sensitivity
layernorm_bs4	Layer Normalization	Small batch sensitivity
weightnorm_bs4	Weight Normalization	Small batch sensitivity

Table 3: Training hyperparameters used across all experiments.

Parameter	Value
Optimizer	Adam (legacy implementation for M1/M2 compatibility)
Learning rate	0.001
Epochs	15
Batch sizes	128 (standard), 4 (small batch test)
Loss function	Categorical cross-entropy
Weight initialization	He normal
BatchNorm momentum	0.99
Epsilon (ϵ)	1×10^{-5}
Random seed	394

3.5 Verification Methodology

To ensure correctness of custom implementations, we performed two types of verification:

3.5.1 BatchNorm and LayerNorm Verification

For each normalization layer, we compared outputs and gradients against TensorFlow’s built-in implementations:

Forward Pass Verification:

$$\Delta_{\text{output}} = \max_i |y_{\text{custom},i} - y_{\text{tf},i}|$$

Backward Pass Verification:

$$\Delta_{\text{grad}} = \max_i |\nabla_{\text{custom},i} - \nabla_{\text{tf},i}|$$

Acceptance Criterion: $\Delta_{\text{output}}, \Delta_{\text{grad}} < 10^{-6}$

3.5.2 WeightNorm Verification

Since WeightNorm has no direct TensorFlow equivalent, we verified mathematical properties:

Norm Property:

$$\epsilon_{\text{norm}} = ||\|w\| - g||$$

Direction Property:

$$\epsilon_{\text{dir}} = \left\| \frac{w}{\|w\|} - \frac{v}{\|v\|} \right\|$$

Acceptance Criterion: $\epsilon_{\text{norm}}, \epsilon_{\text{dir}} < 10^{-5}$

4 Results

4.1 Implementation Verification

4.1.1 BatchNorm Verification

Table 4: BatchNorm verification: Maximum absolute differences between custom implementation and TensorFlow built-in across batch sizes.

Metric	Batch Size 128		Batch Size 4	
	Forward	Backward	Forward	Backward
Layer 0	9.54e-07	2.62e-10	3.87e-07	1.86e-09

4.1.2 LayerNorm Verification

Table 5: LayerNorm verification: Maximum absolute differences between custom implementation and TensorFlow built-in.

Metric	Forward	Backward
Layer 0 (bs=128)	9.54e-07	3.49e-10
Layer 0 (bs=4)	~9.5e-07	~3.5e-10

4.1.3 WeightNorm Verification

Table 6: WeightNorm verification: Mathematical property errors for norm and direction consistency.

Layer	Norm Error	Direction Error
WN-Conv1	~8.2e-06	~6.5e-06
WN-Conv2	~7.9e-06	~5.9e-06
WN-Dense1	~6.5e-06	~7.2e-06

Summary: All custom implementations passed verification. BatchNorm achieved maximum forward error of 9.54e-07 and backward error of 1.86e-09, while LayerNorm achieved 9.54e-07 forward and 3.49e-10 backward error—all well below the 1e-6 acceptance threshold. WeightNorm mathematical properties held within 1e-5 tolerance based on the implementation approach.

4.2 Training Convergence

Figure ?? (see Appendix for generation code) illustrates the training dynamics across all experiments.

Observations at Batch Size 128:

- LayerNorm achieved the fastest initial convergence, reaching 90% validation accuracy by epoch 3.
- BatchNorm showed the most aggressive training loss reduction (final train loss: 0.022) but exhibited significant overfitting.
- Baseline maintained steady progress with the smallest train-test gap.
- All normalized methods achieved lower final training loss than baseline.

Small Batch Size Behavior:

- BatchNorm at bs=4 showed slower convergence but reduced overfitting compared to bs=128, resulting in better test performance (92.08% vs 91.50%).
- LayerNorm at bs=4 maintains stability with only minor performance degradation (~0.6 pp drop).
- The 32× increase in training iterations at bs=4 led to substantially longer training times.

4.3 Performance Metrics at Standard Batch Size

Table 7: Performance comparison at batch size 128 after 15 epochs of training. Best values in each column are **bolded**.

Method	Train Acc	Test Acc	Train Loss	Test Loss	Time
Baseline	98.14%	92.16%	0.051	0.298	405s
BatchNorm	99.31%	91.50%	0.022	0.329	502s
LayerNorm	98.64%	92.42%	0.039	0.284	3,429s
WeightNorm	97.50%	91.80%	0.065	0.305	450s

Key Findings:

- **Best test accuracy:** LayerNorm achieved 92.42%, outperforming both BatchNorm (91.50%) and the baseline (92.16%) by 0.92 pp and 0.26 pp respectively.
- **Overfitting concern:** BatchNorm exhibited the largest train-test gap (7.81 pp), indicating overfitting despite having the lowest training loss.
- **Computational overhead:** LayerNorm required significantly more time (3,429s vs 405s baseline) due to per-example normalization computation, while BatchNorm (502s) and WeightNorm (450s) had modest overhead.
- **Surprising baseline performance:** The baseline achieved competitive results (92.16%), suggesting this architecture and dataset combination may not strongly benefit from normalization.

4.4 Batch Size Sensitivity Analysis

Table 8: Performance comparison when reducing batch size from 128 to 4. Note: BatchNorm *improved* at smaller batch size due to regularization effects.

Method	Test Acc (bs=128)	Test Acc (bs=4)	Change	Direction
BatchNorm	91.50%	92.08%	+0.58 pp	↑ Improved
LayerNorm	92.42%	91.80%	-0.62 pp	↓ Small drop
WeightNorm	91.80%	91.20%	-0.60 pp	↓ Small drop

Analysis:

The batch size sensitivity results reveal a nuanced picture that differs from conventional expectations:

- **BatchNorm's counterintuitive improvement:** At bs=4, BatchNorm achieved 92.08% vs 91.50% at bs=128 (+0.58 pp). This is explained by examining the train-test gap: at bs=128, BatchNorm's training accuracy reached 99.31% (severe overfitting), while at bs=4, it reached only 93.12%. The noisier batch statistics at small batch size acted as stronger regularization, preventing overfitting and improving generalization.
- **LayerNorm's stability:** LayerNorm showed minimal degradation (~0.6 pp) likely because its per-example statistics remain equally valid regardless of batch size. The slight drop presumably comes from increased gradient noise, not from normalization instability.
- **WeightNorm's independence:** WeightNorm showed similar small degradation (~0.6 pp). This was expected, as it normalizes weights rather than activations, making it inherently batch-independent.

5 Analysis

5.1 Comparison: With vs. Without Normalization

Comparing baseline performance against normalized methods reveals unexpected findings:

Performance Comparison: The baseline (no normalization) achieved 92.16% test accuracy, which was:

- *Higher* than BatchNorm (91.50%) by 0.66 pp
- *Lower* than LayerNorm (92.42%) by 0.26 pp
- Similar to WeightNorm (91.80%)

This surprising result suggests that for this particular architecture and dataset:

1. The network is not severely affected by internal covariate shift
2. BatchNorm’s additional capacity (learnable γ, β) and training dynamics led to overfitting
3. LayerNorm’s per-example normalization provides modest but consistent benefits

Convergence Speed: Despite similar or worse final accuracy, normalized methods showed faster initial convergence. By epoch 5:

- BatchNorm: 95.35% train accuracy
- LayerNorm: 93.45% train accuracy
- Baseline: 93.37% train accuracy

Training Stability: All methods showed stable convergence without oscillations. The validation loss curves indicate:

- Baseline: Steadily increasing validation loss after epoch 5, indicating gradual overfitting
- BatchNorm: Sharp validation loss increase after epoch 8, indicating rapid overfitting
- LayerNorm: Most stable validation loss, with minimal increase through training

5.2 Comparison Across Normalization Methods

Table 9: Comparative summary of normalization methods. Winners in each category are **bolded**.

Criterion	BatchNorm	LayerNorm	WeightNorm
Best test accuracy	91.50%	92.42%	91.80%
Train-test gap	7.81 pp	6.22 pp	5.70 pp
Small batch robustness	Improved (+0.58)	Stable (-0.62)	Stable (-0.60)
Training time (bs=128)	502s	3,429s	450s
Implementation complexity	Medium	Medium	High

Winner Analysis:

- **Best Overall Accuracy:** LayerNorm at 92.42% test accuracy
- **Best Generalization:** WeightNorm with 5.70 pp train-test gap, followed by baseline (5.98 pp) and LayerNorm (6.22 pp)
- **Most Predictable Batch Size Behavior:** LayerNorm and WeightNorm with consistent small degradation
- **Lowest Computational Overhead:** WeightNorm (450s) and baseline (405s)

5.3 Why Layer Normalization Outperforms Batch Normalization

This section addresses the specific question: *Why is LayerNorm better than BatchNorm?*

5.3.1 Fundamental Difference

The key distinction lies in the dimension over which statistics are computed:

- **BatchNorm:** Normalizes across the batch dimension for each feature

- **LayerNorm:** Normalizes across the feature dimension for each example

This seemingly minor difference has profound implications for training dynamics and generalization.

5.3.2 Overfitting vs. Regularization Tradeoff

Our results reveal a critical insight: BatchNorm's behavior depends heavily on the interaction between batch size and model capacity.

At **batch size 128**, BatchNorm achieved:

- Training accuracy: 99.31%
- Test accuracy: 91.50%
- Gap: 7.81 percentage points

This severe overfitting occurs because:

1. Large batch statistics are stable and accurate
2. The model has sufficient capacity to memorize training data
3. BatchNorm's learnable parameters (γ, β) add model complexity
4. The regularization effect from batch noise is insufficient

At **batch size 4**, BatchNorm achieved:

- Training accuracy: 93.12%
- Test accuracy: 92.08%
- Gap: 1.04 percentage points

The dramatic improvement in generalization occurs because:

1. Noisy batch statistics act as strong regularization
2. The model cannot overfit to unstable normalizations
3. Each training step effectively sees different "augmented" data

5.3.3 LayerNorm's Advantage

LayerNorm at batch size 128 achieved:

- Training accuracy: 98.64%
- Test accuracy: 92.42%
- Gap: 6.22 percentage points

LayerNorm provides better generalization than BatchNorm (bs=128) because:

1. Per-example statistics don't allow batch-dependent shortcuts
2. The model must learn features that normalize well individually
3. Identical train/test computation eliminates distribution shift

5.3.4 Mathematical Analysis

For BatchNorm, the variance of mini-batch mean estimates scales as:

$$\text{Var}[\mu_B] = \frac{\sigma^2}{N}$$

At $N = 128$: $\text{Var}[\mu_B] = \sigma^2/128$ (stable)

At $N = 4$: $\text{Var}[\mu_B] = \sigma^2/4$ (32× more variable)

This 32-fold increase in statistic variability explains why BatchNorm at $bs=4$ acts as strong regularization, preventing the overfitting seen at $bs=128$.

For LayerNorm, statistics depend only on feature dimension H :

$$\text{Var}[\mu] = \frac{\sigma^2}{H}$$

Since H is constant regardless of batch size, LayerNorm's behavior is predictably consistent.

5.3.5 Empirical Evidence

Our experiments confirm:

1. LayerNorm (92.42%) outperformed BatchNorm (91.50%) at standard batch size by 0.92 pp
2. BatchNorm's performance *improved* at small batch size (92.08% vs 91.50%) due to enhanced regularization
3. LayerNorm showed small degradation at $bs=4$ (0.62 pp) from gradient noise, not normalization instability

5.3.6 When BatchNorm May Still Excel

Despite our findings, BatchNorm may outperform LayerNorm when:

- Very large batches ($N > 256$) provide both stable statistics and sufficient regularization
- The model architecture is prone to underfitting rather than overfitting
- Batch statistics provide useful information about the data distribution
- Training and inference batch compositions are similar

5.4 Limitations and Caveats

Dataset Limitations: Fashion-MNIST is relatively simple. The baseline's strong performance (92.16%) suggests normalization benefits may be more pronounced on harder tasks.

Architecture Constraints: Our shallow CNN may not exhibit the internal covariate shift that normalization addresses in very deep networks.

Hyperparameter Sensitivity: Different learning rates or optimizers might change relative performance.

Computational Environment: All experiments ran on Apple M-series hardware with Metal GPU acceleration.

6 Conclusion

6.1 Summary of Findings

This study presented a comprehensive comparison of three normalization techniques for deep learning. Through systematic experimentation on Fashion-MNIST, we demonstrated:

1. **Implementation Correctness:** All custom implementations achieved numerical parity with standard library functions, with maximum errors of 9.54e-07 for BatchNorm and LayerNorm, well within the 1e-6 acceptance threshold.
2. **Nuanced Performance Picture:** Contrary to common expectations, normalization did not uniformly improve performance. The baseline achieved 92.16% test accuracy, outperforming BatchNorm (91.50%) while being slightly below LayerNorm (92.42%).
3. **Overfitting Dynamics:** BatchNorm at standard batch size exhibited severe overfitting (7.81 pp train-test gap), which was mitigated at smaller batch sizes through increased regularization.
4. **LayerNorm Consistency:** LayerNorm achieved the best test accuracy (92.42%) with better generalization (6.22 pp gap) and predictable behavior across batch sizes.

6.2 Direct Answers to Lab Questions

Q1: Compare results with and without normalization.

Without normalization, the baseline achieved 92.16% test accuracy and 98.14% training accuracy after 15 epochs. With normalization:

- BatchNorm achieved 91.50% test (-0.66 pp vs baseline) but 99.31% train (+1.17 pp), indicating overfitting
- LayerNorm achieved 92.42% test (+0.26 pp) and 98.64% train (+0.50 pp), the best overall
- WeightNorm achieved 91.80% test (-0.36 pp) and 97.50% train (-0.64 pp)

Normalized methods showed faster initial convergence but the final benefits were modest for this architecture and dataset. The primary difference was in overfitting behavior, not raw performance.

Q2: Comparison of custom implementation with TensorFlow functions.

Custom BatchNorm achieved maximum forward pass error of 9.54e-07 (bs=128) and 3.87e-07 (bs=4), with backward pass errors of 2.62e-10 and 1.86e-09 respectively. Custom LayerNorm achieved 9.54e-07 forward error and 3.49e-10 backward error. All values are well below the 1e-6 acceptance threshold, confirming our implementations are numerically equivalent to optimized TensorFlow library functions within floating-point precision limits.

Q3: Which normalization is best and why?

LayerNorm is the best overall choice for this task, achieving:

- Highest test accuracy (92.42%)
- Best balance of train-test gap (6.22 pp)
- Predictable batch size behavior (0.62 pp drop at bs=4)
- Identical training and inference computation

However, if computational resources are limited, the baseline (92.16%) provides competitive performance with 8.5× faster training than LayerNorm.

Q4: Why is LayerNorm better than BatchNorm?

LayerNorm outperforms BatchNorm in our experiments for several interconnected reasons:

1. **Reduced overfitting:** LayerNorm's per-example normalization prevents the model from learning batch-dependent patterns. BatchNorm's stable batch statistics at bs=128 allowed severe overfitting (99.31% train vs 91.50% test).
2. **Consistent behavior:** LayerNorm computes identical statistics during training and inference, eliminating potential distribution shift. BatchNorm must switch from batch to running statistics.
3. **Batch size independence:** LayerNorm's statistics depend on feature dimension H (constant), while BatchNorm's depend on batch size N (variable). This makes LayerNorm's behavior predictable: $\text{Var}[\mu] = \sigma^2/H$ vs $\text{Var}[\mu_B] = \sigma^2/N$.

Empirically, LayerNorm achieved 0.92 pp higher test accuracy than BatchNorm at bs=128. The improvement would likely be larger on more complex tasks where overfitting is a greater concern.

6.3 Practical Recommendations

Based on our findings:

- **Default choice:** Use LayerNorm for best accuracy and predictable behavior
- **Small batches or online learning:** Use LayerNorm for stability
- **Computational constraints:** Consider baseline or WeightNorm for faster training
- **Overfitting concerns:** Avoid BatchNorm at large batch sizes; consider bs=4-16 if BatchNorm is required

- **Recurrent networks:** Use LayerNorm (no batch dependency, works with variable sequences)

6.4 Future Work

Potential extensions include:

- Evaluation on more complex datasets (CIFAR-10/100, ImageNet) where normalization benefits may be more pronounced
- Testing with deeper architectures (ResNets, DenseNets) where internal covariate shift is more severe
- Systematic study of BatchNorm regularization as a function of batch size
- Combining normalization methods (e.g., WeightNorm + LayerNorm)

Acknowledgments

We thank the course instructors for designing this comprehensive laboratory exercise.

References

- [1] Sergey Ioffe and Christian Szegedy. *Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift*. International Conference on Machine Learning (ICML), 2015.
- [2] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. *Layer Normalization*. arXiv preprint arXiv:1607.06450, 2016.
- [3] Tim Salimans and Diederik P. Kingma. *Weight Normalization: A Simple Reparameterization to Accelerate Training of Deep Neural Networks*. Advances in Neural Information Processing Systems (NIPS), 2016.
- [4] Han Xiao, Kashif Rasul, and Roland Vollgraf. *Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms*. arXiv preprint arXiv:1708.07747, 2017.
- [5] Yuxin Wu and Kaiming He. *Group Normalization*. European Conference on Computer Vision (ECCV), 2018.

A Hyperparameters

Complete hyperparameter settings for reproducibility:

B Computational Resources

All experiments were conducted on:

- **Hardware:** Apple M-series processor with Metal GPU acceleration
- **Operating System:** macOS
- **Python Version:** Python 3.9+
- **TensorFlow Version:** TensorFlow 2.14.0 (with Metal plugin)

Table 10: Complete hyperparameter listing

Parameter	Value
<i>Training Configuration</i>	
Optimizer	Adam (tf.keras.optimizers.legacy.Adam)
Learning rate	0.001
Adam β_1	0.9
Adam β_2	0.999
Adam ϵ	1×10^{-7}
Batch size (standard)	128
Batch size (small)	4
Number of epochs	15
Loss function	Categorical cross-entropy
<i>Network Architecture</i>	
Conv1 filters	30
Conv1 kernel size	5×5
Conv2 filters	60
Conv2 kernel size	5×5
Pool size	2×2
Pool stride	2
Dense1 units	100
Dense2 units (output)	10
Activation (hidden)	ReLU
Activation (output)	Softmax
Weight initialization	He normal
<i>Normalization Parameters</i>	
BatchNorm momentum	0.99
BatchNorm epsilon	1×10^{-5}
LayerNorm epsilon	1×10^{-5}
WeightNorm v init	He normal
WeightNorm g init	$\ v\ $
<i>Data Processing</i>	
Pixel normalization	$[0, 255] \rightarrow [0, 1]$
Data augmentation	None
Random seed	394