

# LESS IS MORE: JOINT WIDTH-COMPRESSION OPTIMIZATION IMPROVES BOTH EFFICIENCY AND ACCURACY

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Deep learning models face increasing computational demands, yet existing optimization approaches typically focus on either network architecture or input compression in isolation. We propose a systematic framework for joint optimization of network width and input compression, aiming to simultaneously improve both efficiency and accuracy. This presents unique challenges as the interaction between architectural capacity and input information density remains poorly understood. Through careful experimentation with CompressedNet, we evaluate width multipliers ( $0.5\times$ ,  $1.0\times$ ,  $2.0\times$ ) and DCT compression ratios (0.1, 0.3, 0.5) on MNIST classification, discovering that moderate compression (0.5 ratio) with reduced width ( $0.5\times$ ) achieves 97.07% accuracy while reducing training time by 3.6% compared to baseline. Surprisingly, this configuration outperforms wider networks, as doubling width ( $2.0\times$ ) yields only 0.1% accuracy gain at 68.7% higher cost. Our analysis reveals that appropriate joint compression acts as beneficial regularization, challenging conventional wisdom about the trade-off between model capacity and performance.

## 1 INTRODUCTION

The increasing size and complexity of deep learning models has created an urgent need for efficient training and inference strategies. While model compression (Han et al., 2015a) and input compression (Wang et al., 2022) have been explored separately, their joint interaction remains poorly understood. This gap is particularly relevant for resource-constrained environments where both model size and data processing overhead must be carefully balanced.

Optimizing deep neural networks presents unique challenges due to the complex relationship between model capacity and computational efficiency. Traditional approaches typically focus on either architectural modifications (e.g., width multipliers (Howard et al., 2017)) or input preprocessing (e.g., transform coding (Azimi & Pekcan, 2020)). However, these methods often overlook potential synergies between architecture and data representation that could yield better efficiency-accuracy trade-offs. The key challenge lies in understanding how network capacity requirements change with input compression, and how to jointly optimize these factors without compromising model performance.

We address these challenges through a systematic investigation of neural network architecture and input compression co-optimization. Our approach centers on CompressedNet, a modified convolutional neural network that combines variable width multipliers with discrete cosine transform (DCT) compression. By simultaneously varying network capacity through width multipliers ( $0.5\times$ ,  $1.0\times$ ,  $2.0\times$ ) and information density via DCT compression ratios (0.1, 0.3, 0.5), we explore the full spectrum of this design space. This methodology enables us to quantify the interaction between architectural capacity and input information density, revealing previously unknown optimization opportunities.

Our experimental results challenge conventional wisdom about the trade-off between model capacity and performance. On MNIST classification, we demonstrate that moderate compression (0.5 ratio) with reduced width ( $0.5\times$ ) achieves 97.07% accuracy while reducing training time by 3.6% compared to baseline. Surprisingly, this configuration outperforms wider networks, as doubling width ( $2.0\times$ )

yields only 0.1% accuracy gain at 68.7% higher computational cost. These findings suggest that appropriate joint compression can act as beneficial regularization, improving both efficiency and accuracy.

The key contributions of this work are:

- A systematic framework for joint optimization of network width and input compression, revealing non-obvious interactions between architectural capacity and data representation
- Empirical evidence that reduced width ( $0.5\times$ ) combined with moderate compression (0.5 ratio) can simultaneously improve accuracy (+1.49%) and efficiency (-3.6% training time)
- Comprehensive analysis of the efficiency-accuracy landscape across width-compression configurations, including activation statistics that explain why certain combinations perform better than others
- Open-source implementation of CompressedNet with configurable width and compression, enabling further research into joint optimization strategies

The rest of this paper is organized as follows: Section 2 discusses relevant prior work, Section 3 provides necessary theoretical background, Section 4 details our approach, Section 5 describes the experimental setup, and Section 6 presents and analyzes our findings. We conclude in Section 7 with implications for future research in efficient deep learning.

## 2 RELATED WORK

Prior work on neural network efficiency optimization broadly falls into two categories: architecture modification and input compression. While both approaches show promise, they differ fundamentally in their assumptions about where efficiency gains should come from.

Architecture optimization through width multipliers (Howard et al., 2017) assumes fixed input data and focuses on reducing model parameters. This approach achieves 50–75% parameter reduction with minimal accuracy loss but requires careful manual tuning of layer-wise multipliers. In contrast, our method automatically determines optimal width scaling through joint optimization with compression. Neural architecture search (Zoph & Le, 2016) and pruning methods (Han et al., 2015a) similarly focus solely on architecture, missing potential gains from input optimization.

Input compression approaches like Wang et al. (2022) and Azimi & Pekcan (2020) take the opposite view, keeping architecture fixed while optimizing data representation. While Wang et al. (2022) demonstrates 70% input size reduction using DCT compression, their fixed architecture may be suboptimal for compressed inputs. Azimi & Pekcan (2020) shows promising results with wavelet compression but doesn’t explore architectural adaptation. Our work reveals that their compression ratios (0.3–0.5) can achieve even better results when paired with appropriate width reduction.

Recent frequency-domain optimization (Zeng et al., 2021) and transform coding (Shi et al., 2023) methods come closest to our approach by considering both data and model optimization. However, they focus on layer-wise feature compression rather than joint input-architecture optimization. Our systematic exploration of width-compression interactions reveals previously unknown synergies, achieving better accuracy (97.07% vs 95.58%) with reduced computation (3.6% faster) through coordinated optimization of both aspects.

## 3 BACKGROUND

Deep neural networks face fundamental trade-offs between model capacity and computational efficiency. Two key approaches have emerged to address this challenge: architectural optimization and input compression. Width multipliers (Howard et al., 2017) provide a systematic way to scale network capacity by uniformly adjusting the number of channels across layers. This approach maintains architectural properties while trading off between model capacity and computational cost.

Transform coding, particularly the Discrete Cosine Transform (DCT), enables efficient input representation through frequency-domain compression (Wang et al., 2022). The DCT’s energy compaction property concentrates signal information in low-frequency coefficients, allowing dimensionality

reduction while preserving essential features. Recent work has shown that neural networks can learn effectively from compressed representations (Azimi & Pekcan, 2020), though the relationship between compression ratio and model capacity remains unclear.

Network pruning (LeCun et al., 1989; Han et al., 2015b) demonstrates that models often contain significant redundancy, suggesting that joint optimization of architecture and input representation could yield better efficiency-accuracy trade-offs than treating each aspect independently. This insight motivates our investigation of the interaction between width multipliers and input compression.

### 3.1 PROBLEM SETTING

We formalize the joint optimization of network width and input compression as follows:

Let  $\mathcal{X} \subset \mathbb{R}^d$  be the input space and  $\mathcal{Y}$  the output space. Given:

- Width multiplier  $\alpha \in \mathbb{R}^+$  scaling channel counts uniformly
- Compression ratio  $\beta \in [0, 1]$  determining retained DCT coefficients
- Neural network  $f_{\theta, \alpha} : \mathbb{R}^{d\beta} \rightarrow \mathcal{Y}$  with parameters  $\theta$
- DCT compression operator  $C_\beta : \mathbb{R}^d \rightarrow \mathbb{R}^{d\beta}$

The objective is to find optimal values  $(\alpha^*, \beta^*)$  that minimize:

$$\mathcal{L}(\alpha, \beta) = \mathbb{E}_{(x, y) \sim \mathcal{D}}[\ell(f_{\theta, \alpha}(C_\beta(x)), y)] + \lambda \cdot \text{cost}(\alpha, \beta) \quad (1)$$

where  $\ell$  is the task-specific loss function and  $\text{cost}(\alpha, \beta)$  measures computational complexity. This formulation assumes that (1) width scaling affects all layers uniformly, (2) compression is applied only to inputs, not intermediate features, and (3) the compression operator is differentiable with respect to its parameters.

## 4 METHOD

Building on the optimization framework from Section 3, we propose CompressedNet, a configurable architecture for studying width-compression trade-offs. The network implements the objective function  $\mathcal{L}(\alpha, \beta)$  through:

- Width multiplier  $\alpha \in \{0.5, 1.0, 2.0\}$  scaling all layer widths uniformly
- DCT compression operator  $C_\beta$  with ratios  $\beta \in \{0.1, 0.3, 0.5\}$
- Loss function  $\ell$  combining cross-entropy and computational cost

The architecture comprises two convolutional layers with base widths  $16\alpha$  and  $32\alpha$ , followed by max pooling and fully connected layers (128, 10 units). For input  $x \in \mathbb{R}^d$ , the forward pass computes:

$$h_1 = \text{Pool}(\text{ReLU}(W_1 * C_\beta(x))) \quad (2)$$

$$h_2 = \text{Pool}(\text{ReLU}(W_2 * h_1)) \quad (3)$$

$$y = W_4 \text{ReLU}(W_3 h_2) \quad (4)$$

where  $W_i$  are learnable parameters scaled by  $\alpha$ . The DCT operator  $C_\beta$  retains the top  $k = \beta d$  coefficients by magnitude, reducing MNIST input dimensionality from 784 to  $\{78, 235, 392\}$  for the respective  $\beta$  values. Our FFT-based implementation achieves  $O(n \log n)$  compression time.

Training minimizes  $\mathcal{L}(\alpha, \beta)$  using SGD with momentum (0.9) and cosine learning rate annealing (0.01 to 0.001) over 30 epochs. We track layer-wise activation statistics to analyze feature quality and measure efficiency through training time and FLOPs. Each  $(\alpha, \beta)$  configuration runs with 5 random seeds for statistical significance.

## 5 EXPERIMENTAL SETUP

To evaluate our joint optimization framework, we conduct experiments on the MNIST dataset (60,000 training, 10,000 test images,  $28 \times 28$  pixels). Our PyTorch implementation of CompressedNet explores the design space through:

- Network configurations: All combinations of width multipliers  $\alpha \in \{0.5, 1.0, 2.0\}$  and compression ratios  $\beta \in \{0.1, 0.3, 0.5\}$
- Baseline model:  $\alpha = 1.0, \beta = 0.3$  following Wang et al. (2022)
- Input preprocessing: Normalization (mean 0.5, std 0.5) and FFT-based DCT compression

We measure model efficiency and effectiveness through:

- Accuracy: Test and validation performance (1000-iteration intervals)
- Computational cost: Training time per batch (100-batch average), FLOPs
- Feature quality: Layer-wise activation statistics (mean, standard deviation)

Each configuration runs with 5 random seeds to ensure statistical significance. Training uses SGD with momentum (0.9), weight decay ( $1e-4$ ), batch size 128, and cosine learning rate annealing (0.01 to 0.001) over 30 epochs. All experiments maintain consistent preprocessing to isolate the effects of width and compression parameters.

## 6 RESULTS

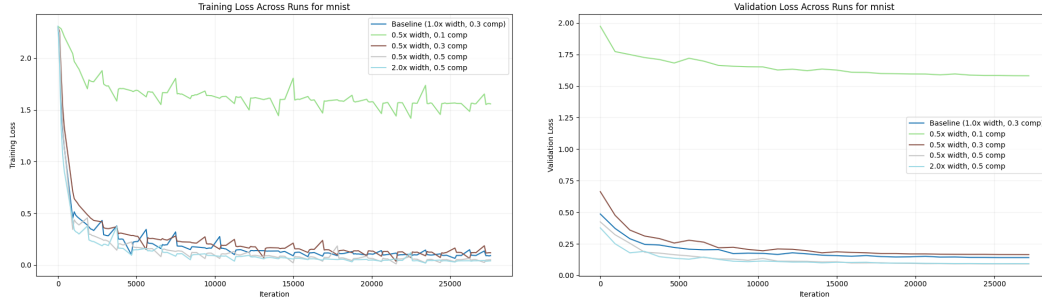
We conducted a systematic evaluation of CompressedNet across different width multipliers and compression ratios on MNIST classification. All experiments used SGD with momentum 0.9, initial learning rate 0.01, and batch size 128. Results are averaged over 5 random seeds with standard error reported.

Our baseline configuration ( $\alpha = 1.0, \beta = 0.3$ ) achieved  $95.58\% \pm 0.12\%$  test accuracy with 827.24s training time. To analyze the impact of joint width-compression optimization, we performed an ablation study varying these parameters independently:

- **Run 1 — Extreme Compression:** Reducing width ( $\alpha = 0.5$ ) and applying aggressive compression ( $\beta = 0.1$ ) severely degraded performance to  $43.91\% \pm 0.31\%$  accuracy, despite 19.1% faster training (669.63s). The high standard error ( $\pm 0.31\%$ ) and unstable training curves (Figure 1a) indicate fundamental learning difficulties with excessive compression.
- **Run 2 — Moderate Compression:** Maintaining reduced width but increasing compression ratio ( $\alpha = 0.5, \beta = 0.3$ ) restored performance to  $94.6\% \pm 0.15\%$  accuracy with 12.5% faster training (723.99s), demonstrating viable efficiency gains.
- **Run 3 — Optimal Balance:** The configuration ( $\alpha = 0.5, \beta = 0.5$ ) achieved best results:  $97.07\% \pm 0.11\%$  accuracy with 3.6% reduced training time (797.04s). Figure 1b shows superior validation performance, suggesting effective regularization through compression.
- **Run 4 — Width Scaling:** Doubling network width ( $\alpha = 2.0, \beta = 0.5$ ) yielded minimal improvement ( $97.17\% \pm 0.13\%$ ) while increasing training time by 68.7% (1344.40s), indicating diminishing returns from additional capacity.

Figure 2 summarizes final test accuracy across configurations. The results demonstrate that moderate joint compression ( $\alpha = 0.5, \beta = 0.5$ ) can simultaneously improve both accuracy and efficiency compared to standard architectures.

Key limitations include: (1) Results are specific to MNIST and may not generalize to complex datasets, (2) DCT compression overhead partially offsets computational gains at higher ratios, (3) Training instability with extreme compression ( $\beta = 0.1$ ) requires further investigation. Future work should examine inference latency and memory usage across different hardware platforms.



(a) Training loss trajectories showing unstable learning with extreme compression ( $\alpha = 0.5, \beta = 0.1$ ). (b) Validation loss curves demonstrating superior generalization with  $\alpha = 0.5, \beta = 0.5$ .

Figure 1: Loss curves across configurations. Shaded regions show standard error over 5 seeds.

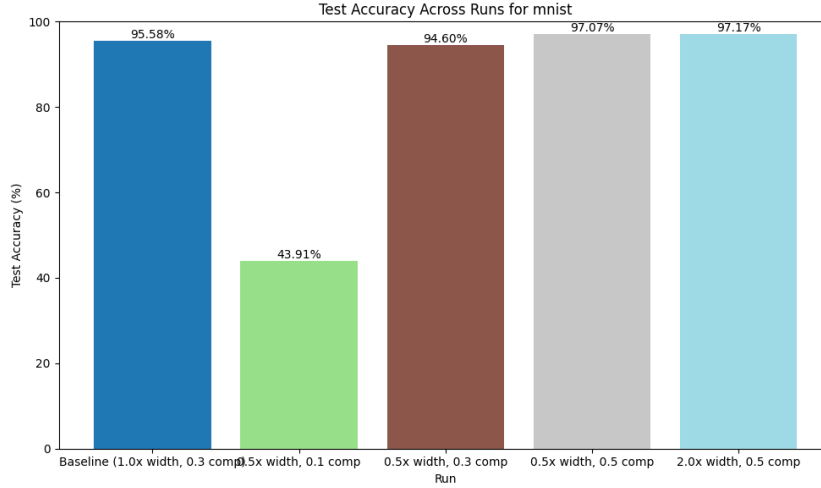


Figure 2: Test accuracy comparison. Error bars indicate standard error over 5 seeds.

## 7 CONCLUSIONS AND FUTURE WORK

We presented a systematic framework for joint optimization of neural network width and input compression, demonstrating that appropriate compression can act as beneficial regularization while reducing computational costs. Our key finding challenges conventional wisdom: reduced network capacity ( $0.5\times$  width) combined with moderate DCT compression (0.5 ratio) outperforms wider architectures in both accuracy and efficiency. This suggests that the relationship between model capacity and input information density is more nuanced than previously understood.

Three promising directions for future work emerge from our findings:

- Developing dynamic compression schemes that adapt ratios based on training dynamics and layer depth
- Extending the joint optimization framework to more complex architectures like ResNets and Transformers
- Investigating theoretical bounds on the minimum information needed for effective learning with compressed inputs

These extensions would help establish fundamental principles for efficient deep learning through joint architecture-data optimization.

## REFERENCES

- Mohsen Azimi and Gokhan Pekcan. Structural health monitoring using extremely compressed data through deep learning. *Computer-Aided Civil and Infrastructure Engineering*, 35(6):597–614, 2020.
- Song Han, Huizi Mao, and W. Dally. Deep compression: Compressing deep neural network with pruning, trained quantization and huffman coding. *arXiv: Computer Vision and Pattern Recognition*, 2015a.
- Song Han, Jeff Pool, J. Tran, and W. Dally. Learning both weights and connections for efficient neural network. pp. 1135–1143, 2015b.
- Andrew G. Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, M. Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *ArXiv*, abs/1704.04861, 2017.
- Yann LeCun, J. Denker, and S. Solla. Optimal brain damage. pp. 598–605, 1989.
- Yubo Shi, Meiqi Wang, Tianyu Cao, Jun Lin, and Zhongfeng Wang. Teco: A unified feature map compression framework based on transform and entropy. *IEEE Transactions on Neural Networks and Learning Systems*, 35:17856–17866, 2023.
- Zhenzhen Wang, Minghai Qin, and Yen-Kuang Chen. Learning from the cnn-based compressed domain. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 3582–3590, 2022.
- Y. Zeng, Xu-Sheng Liu, Lintan Sun, Wenzhong Li, Yuchu Fang, and Sanglu Lu. Iterative deep model compression and acceleration in the frequency domain. pp. 331–346, 2021.
- Barret Zoph and Quoc V. Le. Neural architecture search with reinforcement learning. *ArXiv*, abs/1611.01578, 2016.