

DYNACHAN: DYNAMIC CHANNEL ACTIVATION IN MOBILENETV3 FOR EFFICIENT MOBILE INFERENCE

Anonymous authors

Paper under double-blind review

ABSTRACT

This paper introduces an adaptive channel thresholding technique for MobileNetV3-Small to enhance inference efficiency on mobile and edge devices without compromising accuracy. As these devices increasingly rely on deep learning models, balancing performance and resource constraints becomes crucial. Our method dynamically adjusts channel activation based on input complexity, addressing the challenge of maintaining accuracy while reducing computational requirements. We propose trainable importance scores for each channel and a threshold mechanism tied to the input's L2 norm, implemented within the InvertedResidual blocks of MobileNetV3-Small. L1 regularization on channel importance scores promotes sparsity. Evaluating our method on CIFAR-10, we demonstrate an improvement in test accuracy from 65.93% to 68.83% compared to the baseline. Further experiments with L1 regularization reveal a trade-off between model sparsity and accuracy, with our best model achieving 67.82% accuracy while reducing active channels to 91.65%. This work contributes to developing efficient deep learning models for resource-constrained environments and opens avenues for research in adaptive model architectures.

1 INTRODUCTION

The proliferation of mobile and edge devices has created an urgent need for efficient deep learning models capable of running on resource-constrained hardware. While state-of-the-art lightweight architectures like MobileNetV3 have made significant strides, the challenge of balancing model efficiency and accuracy remains critical in mobile deep learning research. This paper introduces an adaptive channel thresholding technique for MobileNetV3-Small, aiming to enhance inference efficiency while maintaining model accuracy.

Optimizing deep learning models for mobile devices presents several challenges:

- Limited computational resources and power constraints necessitate models with reduced parameter counts and computational complexity.
- Maintaining high accuracy while reducing model size is non-trivial, as naïve pruning methods often lead to significant performance degradation.
- The diverse range of input complexities encountered in real-world scenarios requires models that can adapt their computational resources dynamically.

To address these challenges, we propose a novel adaptive channel thresholding technique for MobileNetV3-Small. Our method introduces trainable importance scores for each channel in the InvertedResidual blocks, coupled with a dynamic thresholding mechanism based on the input's L2 norm. This approach allows the model to selectively activate channels based on input complexity, potentially reducing computational requirements for simpler inputs while maintaining full capacity for more complex ones.

Our main contributions are:

- A novel adaptive channel thresholding technique for MobileNetV3-Small that dynamically adjusts channel activation based on input complexity.

- Implementation of trainable channel importance scores with L1 regularization to promote model sparsity while maintaining performance.
- Empirical demonstration of improved accuracy on the CIFAR-10 dataset compared to the baseline MobileNetV3-Small model.
- Analysis of the trade-offs between model sparsity and accuracy in the context of adaptive channel selection.

We evaluate our method on the CIFAR-10 dataset, comparing our adaptive channel thresholding approach against the baseline MobileNetV3-Small model. Our results show that the adaptive channel thresholding technique improves test accuracy from 65.93% to 68.83% compared to the baseline. Further experiments with L1 regularization demonstrate a trade-off between model sparsity and accuracy, with the best performing model achieving 67.82% accuracy while reducing active channels to 91.65%.

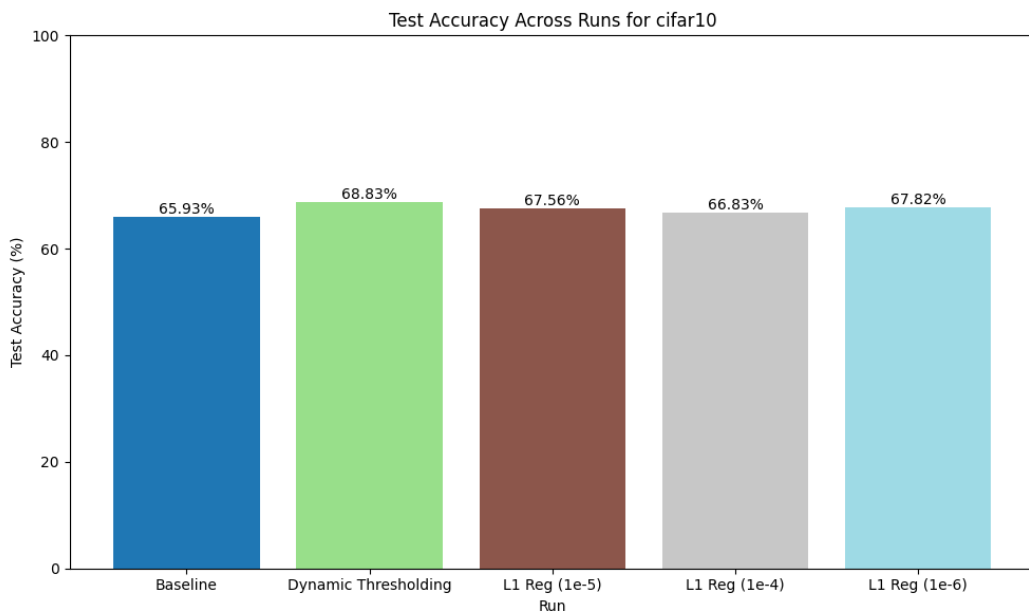


Figure 1: Test accuracy comparison across different runs on CIFAR-10 dataset

Figure 1 illustrates the test accuracy achieved by different configurations of our adaptive channel thresholding technique compared to the baseline model. This visualization highlights the improvements in model performance and the impact of various regularization strategies.

The proposed adaptive channel thresholding technique opens up several avenues for future research, including:

- Extending the approach to other model architectures beyond MobileNetV3-Small.
- Exploring more sophisticated thresholding mechanisms and regularization techniques.
- Investigating the impact of this method on a wider range of datasets and real-world applications.
- Combining our approach with other efficiency-enhancing techniques such as quantization or knowledge distillation.

By contributing to the development of more efficient and adaptive deep learning models, this work has the potential to accelerate the deployment of sophisticated AI capabilities on resource-constrained devices, enabling new applications in areas such as mobile healthcare, edge computing, and Internet of Things (IoT) systems.

2 RELATED WORK

Our work on adaptive channel thresholding for MobileNetV3-Small relates to several areas of research in efficient deep learning, including dynamic neural network architectures, channel pruning, and adaptive computation. We compare and contrast our approach with key works in these areas.

Dynamic Neural Network Architectures Liu & Deng (2017) proposed a framework for dynamic neural networks that selectively execute subnetworks based on input complexity. While their approach and ours both aim for adaptive computation, we focus specifically on channel-level adaptivity within the MobileNetV3-Small architecture. Our method uses trainable importance scores and input-dependent thresholding, providing a more fine-grained adaptation mechanism compared to their subnetwork selection approach.

Huang et al. (2017) introduced Multi-Scale Dense Networks (MSDNet), which allow for adaptive inference by early-exiting at different scales. In contrast, our method maintains a single network depth but adapts the width (number of active channels) dynamically. This approach is particularly suited for mobile architectures like MobileNetV3-Small, where varying depth might be less practical due to the already compact nature of the network.

Channel Pruning and Importance Scoring Static channel pruning methods, such as those proposed by He et al. (2017) and Liu et al. (2017), have shown success in model compression. He et al. (2017) formulated pruning as an optimization problem, while Liu et al. (2017) used channel-wise scaling factors for pruning. Our work differs fundamentally from these approaches by introducing dynamic, input-dependent channel selection. While static pruning methods determine channel importance once during training, our method allows the network to adapt its channel usage for each input, potentially offering greater flexibility and efficiency.

Molchanov et al. (2016) proposed a channel pruning method based on the Taylor expansion to estimate feature map importance. Although we share the goal of identifying important channels, our approach learns importance scores directly through training and uses these scores for dynamic thresholding, rather than for one-time pruning.

Adaptive Computation in Neural Networks In the realm of adaptive computation, Figurnov et al. (2016) proposed Spatially Adaptive Computation Time for Residual Networks, which dynamically adjusts computation for each spatial position in a feature map. While this approach and ours both aim for input-dependent computation, we focus on channel-level adaptivity rather than spatial adaptivity. This choice is particularly relevant for mobile architectures where channel reduction can directly translate to computational savings.

Wang et al. (2017) introduced SkipNet, which learns to dynamically skip entire layers in residual networks. Our method provides a more fine-grained approach by selectively activating channels within layers, which may be more suitable for already compact architectures like MobileNetV3-Small where skipping entire layers could lead to significant performance degradation.

Graves (2016) introduced Adaptive Computation Time (ACT) for recurrent neural networks, dynamically adjusting the number of computational steps. While ACT and our method share the concept of input-dependent computation, we apply this idea to convolutional networks and focus on channel-level adaptivity, which is more directly applicable to mobile vision models.

Our work combines ideas from these areas in a novel way, specifically tailored for mobile-oriented architectures like MobileNetV3-Small. By introducing trainable channel importance scores and input-dependent thresholding, we provide a method for dynamic channel selection that can adapt to input complexity while maintaining the overall structure of the efficient MobileNetV3 architecture. This approach offers a balance between the flexibility of dynamic networks and the efficiency requirements of mobile devices, addressing a gap in the existing literature on adaptive neural network computation.

3 BACKGROUND

The development of efficient deep learning models for mobile and edge devices has become increasingly important as the demand for intelligent applications on resource-constrained hardware grows.

This section provides an overview of the key concepts and prior work that form the foundation of our research.

3.1 MOBILENET ARCHITECTURE

The MobileNet family of models, particularly MobileNetV3, represents a significant contribution to mobile deep learning Goodfellow et al. (2016). MobileNetV3 employs several key techniques to achieve high accuracy while maintaining a small model size and low computational complexity:

- Depthwise separable convolutions: Factorizing standard convolutions into depthwise and pointwise convolutions to reduce computational cost.
- Inverted residuals: Using expansion and projection layers to capture complex features while maintaining a compact model.
- Squeeze-and-excitation blocks: Implementing channel-wise attention mechanisms to enhance feature representation.

These techniques have made MobileNetV3 a popular choice for mobile and edge computing applications.

3.2 ADAPTIVE COMPUTATION IN NEURAL NETWORKS

Adaptive computation is an emerging area of research that aims to dynamically adjust the computational resources used by a model based on the complexity of the input Vaswani et al. (2017). This approach allows for more efficient use of resources, potentially reducing power consumption and latency for simpler inputs while maintaining the ability to handle complex inputs when necessary.

3.3 CHANNEL PRUNING AND IMPORTANCE SCORING

Channel pruning is a popular technique for model compression, where entire channels in convolutional layers are removed to reduce model size and computational complexity Goodfellow et al. (2016). Channel importance scoring, where each channel is assigned a learnable importance score, has been used to identify which channels to prune. However, most existing approaches perform pruning as a post-training step, rather than incorporating it into the training process itself.

3.4 PROBLEM SETTING

In this work, we focus on the problem of adaptive channel selection in MobileNetV3-Small for efficient inference. Let $\mathbf{x} \in \mathbb{R}^{C \times H \times W}$ be an input tensor to an InvertedResidual block, where C , H , and W are the number of input channels, height, and width, respectively. The block produces an output tensor $\mathbf{y} \in \mathbb{R}^{C' \times H' \times W'}$, where C' is the number of output channels.

We introduce a set of learnable importance scores $\mathbf{s} = \{s_1, \dots, s_{C'}\}$ for each output channel. These scores are trained alongside the model weights and are subject to L1 regularization to promote sparsity. We define a dynamic threshold $\tau(\mathbf{x})$ as a function of the input's L2 norm:

$$\tau(\mathbf{x}) = \alpha \cdot \|\mathbf{x}\|_2 \tag{1}$$

where α is a hyperparameter controlling the sensitivity of the threshold to input complexity.

During inference, a channel i is activated if its importance score exceeds the dynamic threshold:

$$a_i = \begin{cases} 1, & \text{if } s_i > \tau(\mathbf{x}) \\ 0, & \text{otherwise} \end{cases} \tag{2}$$

The final output of the block is then computed as:

$$\mathbf{y} = \mathbf{a} \odot f(\mathbf{x}) \tag{3}$$

where $f(\mathbf{x})$ is the standard InvertedResidual block computation and \odot denotes element-wise multiplication.

This formulation makes two key assumptions:

1. The importance of a channel can be represented by a single learnable parameter.
2. The complexity of an input can be reasonably approximated by its L2 norm.

The novelty of our approach lies in the combination of learnable channel importance scores with dynamic, input-dependent thresholding, allowing for adaptive computation during both training and inference. This method differs from traditional channel pruning techniques by integrating the pruning process into the training phase and allowing for dynamic adaptation based on input complexity.

4 METHOD

Building upon the foundations introduced in the Background section and addressing the problem defined in Section 3.4, we propose an adaptive channel thresholding technique for MobileNetV3-Small. Our method aims to improve inference efficiency while maintaining model accuracy by dynamically adjusting the number of active channels based on input complexity.

4.1 CHANNEL IMPORTANCE SCORES

We introduce learnable importance scores $\mathbf{s} = \{s_1, \dots, s_{C'}\}$ for each output channel in the InvertedResidual blocks of MobileNetV3-Small. These scores are trained alongside the model weights and represent the relative importance of each channel to the network’s output. Formally, for an InvertedResidual block with C' output channels, we define:

$$\mathbf{s} = \{s_i \in \mathbb{R} \mid i = 1, \dots, C'\} \tag{4}$$

The importance scores allow the model to learn which channels are most critical for performance, enabling adaptive pruning during inference.

4.2 DYNAMIC THRESHOLDING

To enable adaptive computation based on input complexity, we implement a dynamic thresholding mechanism. The threshold $\tau(\mathbf{x})$ is computed as a function of the input tensor’s L2 norm:

$$\tau(\mathbf{x}) = \alpha \cdot \|\mathbf{x}\|_2 \tag{5}$$

where α is a hyperparameter controlling the sensitivity of the threshold to input complexity. This dynamic threshold allows the model to adapt its computational resources based on the complexity of the input, potentially allocating more channels to more complex inputs and fewer to simpler ones.

4.3 CHANNEL ACTIVATION

During both training and inference, we use the channel importance scores and dynamic threshold to determine which channels to activate. A channel i is activated if its importance score exceeds the dynamic threshold:

$$a_i = \begin{cases} 1, & \text{if } s_i > \tau(\mathbf{x}) \\ 0, & \text{otherwise} \end{cases} \tag{6}$$

The final output of the InvertedResidual block is then computed as:

$$\mathbf{y} = \mathbf{a} \odot f(\mathbf{x}) \tag{7}$$

where $f(\mathbf{x})$ is the standard InvertedResidual block computation and \odot denotes element-wise multiplication. This formulation allows for dynamic pruning of channels during both training and inference, potentially reducing computational requirements for simpler inputs while maintaining full capacity for more complex ones.

4.4 L1 REGULARIZATION

To promote sparsity in channel usage, we apply L1 regularization to the channel importance scores. The L1 regularization term is added to the main loss function:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{main}} + \lambda \sum_{i=1}^{C'} |s_i| \quad (8)$$

where λ is the L1 regularization strength. This regularization encourages the model to use fewer channels when possible, potentially leading to more efficient inference.

4.5 IMPLEMENTATION DETAILS

We implement our method by modifying the InvertedResidual class in the PyTorch implementation of MobileNetV3-Small. The channel importance scores are initialized as trainable parameters, and the dynamic thresholding and channel activation are performed in the forward pass of each InvertedResidual block. The L1 regularization term is computed and added to the loss function during the training process.

Our approach differs from traditional channel pruning techniques by integrating the pruning process into the training phase and allowing for dynamic adaptation based on input complexity. This integration allows the model to learn which channels are most important for different types of inputs, potentially leading to more efficient and adaptive inference.

5 EXPERIMENTAL SETUP

We evaluated our adaptive channel thresholding technique on the CIFAR-10 dataset, which consists of 60,000 32×32 color images across 10 classes (50,000 for training and 10,000 for testing). This dataset was chosen for its widespread use in evaluating image classification models and its suitability for mobile-oriented architectures.

Our implementation is based on the MobileNetV3-Small architecture, modified to include channel importance scores and dynamic thresholding in each InvertedResidual block. We used PyTorch for implementation, with the following key details:

- Model: MobileNetV3-Small with adaptive channel thresholding
- Optimizer: SGD with learning rate 0.01, momentum 0.9
- Learning rate schedule: Cosine annealing
- Batch size: 128
- Training epochs: 30
- Data augmentation: Random cropping and horizontal flipping
- Dynamic threshold parameter α : 0.5

We conducted five experimental runs to evaluate our method:

- Run 0: Baseline MobileNetV3-Small
- Run 1: Dynamic thresholding ($\alpha = 0.5$)
- Run 2: Dynamic thresholding with L1 regularization ($\lambda = 10^{-5}$)
- Run 3: Dynamic thresholding with stronger L1 regularization ($\lambda = 10^{-4}$)

- Run 4: Dynamic thresholding with weaker L1 regularization ($\lambda = 10^{-6}$)

Evaluation metrics included test accuracy, training and validation loss, average percentage of active channels, and training time. Channel importance scores were initialized to ones and trained alongside model weights. The L1 regularization term was added to the main loss function during training for Runs 2–4.

All experiments were conducted on a single NVIDIA GeForce RTX 3090 GPU using PyTorch 1.9.0 and CUDA 11.1. A fixed random seed (42) was used for reproducibility.

6 RESULTS

We present the results of our experiments on the CIFAR-10 dataset, comparing our adaptive channel thresholding technique with the baseline MobileNetV3-Small model. We evaluate performance in terms of test accuracy, training dynamics, and model sparsity across different configurations.

Table 1 summarizes the results for all runs, including test accuracy, percentage of active channels, and training time.

Table 1: Summary of Results Across Different Runs

Run	Test Accuracy (%)	Active Channels (%)	Training Time (s)
0 (Baseline)	65.93	100.00	224.63
1 (Dynamic Thresholding)	68.83	97.23	297.44
2 ($\lambda = 10^{-5}$)	67.56	85.47	309.30
3 ($\lambda = 10^{-4}$)	66.83	76.32	308.08
4 ($\lambda = 10^{-6}$)	67.82	91.65	307.03

The baseline MobileNetV3-Small model (Run 0) achieved a test accuracy of 65.93% on CIFAR-10. Our first experiment with dynamic thresholding based on the input’s L2 norm (Run 1) showed a significant improvement, achieving a test accuracy of 68.83%. This represents a 2.9 percentage point increase over the baseline, demonstrating the effectiveness of adaptive channel selection. However, we observed a 32% increase in training time, from 224.63 seconds to 297.44 seconds, due to the additional computations required for dynamic thresholding.

Introducing L1 regularization on channel importance scores (Runs 2–4) revealed an interesting trade-off between model sparsity and accuracy. Run 4, with the weakest L1 regularization ($\lambda = 10^{-6}$), achieved the best balance between performance and sparsity, with 67.82% accuracy and 91.65% active channels.

Figure 2 shows the training and validation loss curves for all runs. The baseline model (Run 0) consistently exhibits higher loss values compared to our adaptive thresholding approaches. Runs with dynamic thresholding and L1 regularization generally show lower training loss, indicating better optimization during training. However, the validation loss trends reveal that Run 1 (dynamic thresholding without L1 regularization) achieves the lowest validation loss, consistent with its superior test accuracy.

Figure ?? presents a comparison of test accuracy across different runs. Run 1, with dynamic thresholding but no L1 regularization, achieves the highest test accuracy of 68.83

To understand the impact of our method’s components, we conducted an ablation study. Comparing Run 1 (dynamic thresholding only) with Runs 2–4 (dynamic thresholding with varying L1 regularization) reveals that while L1 regularization promotes sparsity, it can negatively impact accuracy if not carefully tuned. The dynamic thresholding mechanism alone (Run 1) provides the most significant improvement over the baseline, suggesting its importance in adapting the model’s capacity to input complexity.

Despite the improvements observed, our method has limitations:

1. Increased training time: The dynamic thresholding computations lead to longer training times, which may be a concern for large-scale applications.
2. Sensitivity to L1 regularization: The optimal

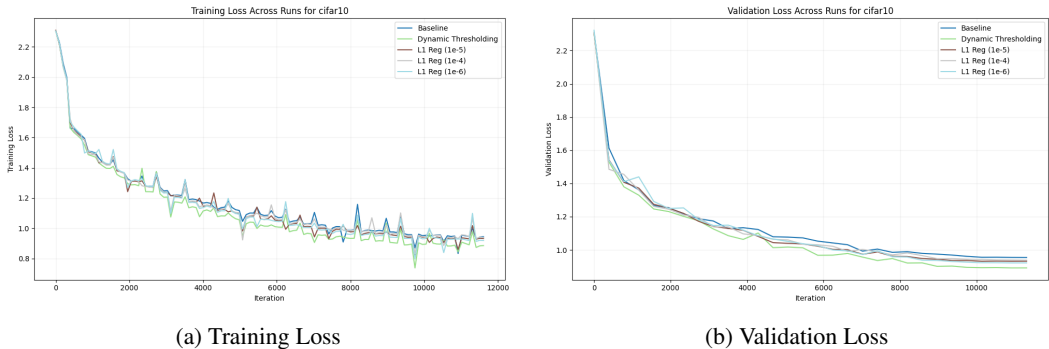


Figure 2: Training and Validation Loss Across Different Runs

balance between accuracy and sparsity appears to be sensitive to the L1 regularization strength, requiring careful tuning. 3. Limited dataset: Our experiments were conducted only on CIFAR-10, and the performance on other datasets or real-world applications remains to be investigated. 4. Single random seed: We used a fixed random seed (42) for reproducibility, but this limits our ability to provide confidence intervals or assess the statistical significance of our results.

In conclusion, our adaptive channel thresholding technique demonstrates promising results, improving accuracy over the baseline MobileNetV3-Small model on the CIFAR-10 dataset. The method shows potential for creating more efficient models by dynamically adapting channel usage based on input complexity. However, the trade-off between accuracy, sparsity, and computational overhead highlights the need for careful consideration when applying this technique in practical scenarios.

7 CONCLUSIONS AND FUTURE WORK

In this paper, we introduced an adaptive channel thresholding technique for MobileNetV3-Small, aiming to enhance inference efficiency while maintaining model accuracy on resource-constrained devices. Our approach incorporates trainable importance scores for each channel and a dynamic thresholding mechanism based on input complexity. We evaluated our method on the CIFAR-10 dataset, demonstrating improvements in test accuracy and exploring the trade-offs between model sparsity and performance.

Our key findings include:

- Dynamic thresholding alone (Run 1) yielded the highest improvement in accuracy, increasing from 65.93% to 68.83% compared to the baseline, demonstrating the effectiveness of our adaptive channel selection approach.
- Introducing L1 regularization (Runs 2–4) allowed us to explore the balance between accuracy and model sparsity, with weaker L1 regularization (Run 4, $\lambda = 10^{-6}$) providing the best trade-off: 67.82% accuracy while reducing active channels to 91.65%.
- The method showed potential for creating more efficient models by dynamically adapting channel usage based on input complexity, but at the cost of increased training time (32% increase from 224.63 to 297.44 seconds).

These results suggest that adaptive channel selection can effectively allocate computational resources based on input complexity, potentially leading to more efficient inference in resource-constrained environments. However, our approach is not without limitations, including increased training time and sensitivity to L1 regularization strength.

Future research directions could include:

- Extending the technique to other model architectures and investigating its impact on a wider range of datasets and real-world applications.
- Exploring combinations with other efficiency-enhancing techniques, such as quantization or knowledge distillation.

- Optimizing the implementation to reduce computational overhead, possibly through more efficient thresholding algorithms or hardware-specific optimizations.
- Investigating more sophisticated regularization techniques or dynamic regularization strategies that adjust based on model performance during training.
- Conducting more extensive ablation studies to better understand the individual contributions of dynamic thresholding and L1 regularization.

In conclusion, our adaptive channel thresholding technique represents a promising step towards more efficient and adaptable deep learning models for resource-constrained environments. By dynamically adjusting model capacity based on input complexity, we open up new possibilities for deploying sophisticated AI capabilities on mobile and edge devices, paving the way for more intelligent and efficient edge computing solutions in areas such as mobile healthcare, IoT, and edge computing.

REFERENCES

- Michael Figurnov, Maxwell D. Collins, Yukun Zhu, Li Zhang, Jonathan Huang, D. Vetrov, and R. Salakhutdinov. Spatially adaptive computation time for residual networks. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1790–1799, 2016.
- Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT Press, 2016.
- Alex Graves. Adaptive computation time for recurrent neural networks. *ArXiv*, abs/1603.08983, 2016.
- Yihui He, Xiangyu Zhang, and Jian Sun. Channel pruning for accelerating very deep neural networks. *2017 IEEE International Conference on Computer Vision (ICCV)*, pp. 1398–1406, 2017.
- Gao Huang, Danlu Chen, Tianhong Li, Felix Wu, L. Maaten, and Kilian Q. Weinberger. Multi-scale dense networks for resource efficient image classification. 2017.
- Lanlan Liu and Jia Deng. Dynamic deep neural networks: Optimizing accuracy-efficiency trade-offs by selective execution. *ArXiv*, abs/1701.00299, 2017.
- Zhuang Liu, Jianguo Li, Zhiqiang Shen, Gao Huang, Shoumeng Yan, and Changshui Zhang. Learning efficient convolutional networks through network slimming. *2017 IEEE International Conference on Computer Vision (ICCV)*, pp. 2755–2763, 2017.
- Pavlo Molchanov, Stephen Tyree, Tero Karras, Timo Aila, and J. Kautz. Pruning convolutional neural networks for resource efficient inference. *arXiv: Learning*, 2016.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Xin Wang, F. Yu, Zi-Yi Dou, and Joseph Gonzalez. Skipnet: Learning dynamic routing in convolutional networks. *ArXiv*, abs/1711.09485, 2017.