



e-ISSN: 2319-8753 | p-ISSN: 2347-6710

IJRSET

International Journal of Innovative Research in
SCIENCE | ENGINEERING | TECHNOLOGY



INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN SCIENCE | ENGINEERING | TECHNOLOGY

Volume 13, Issue 10, October 2024

ISSN INTERNATIONAL
STANDARD
SERIAL
NUMBER
INDIA

Impact Factor: 8.524

9940 572 462

6381 907 438

ijrset@gmail.com

www.ijrset.com

Performance-Driven Feature Selection for CNNs: Impacts on Accuracy and Efficiency

Sara Laroui

ENSA, Abdelmalek Essadi University, Morocco

ABSTRACT: Convolutional Neural Networks (CNNs) have become integral to many computer vision applications, but their increasing complexity raises challenges related to computational cost, memory usage, and energy consumption, especially in resource-limited environments. This paper investigates performance-driven feature selection as a method to address these challenges by reducing model complexity without significantly compromising accuracy. We explore two main avenues of feature selection: input-level reduction, which trims unnecessary features from raw input data, and layer-level pruning, which eliminates redundant filters from the network's layers.

Through experiments on a ResNet-50 model using the CIFAR-10 dataset, we compare the performance of various feature selection techniques, including principal component analysis (PCA), L1 norm pruning, and filter pruning. The results demonstrate that layer-level feature selection, particularly L1 norm pruning, offers a strong trade-off between computational efficiency and model performance. While PCA on input data slightly reduces FLOPs, it also causes a more noticeable accuracy decline. In contrast, L1 norm pruning reduces FLOPs by 30% with minimal accuracy loss. Filter pruning achieves the highest efficiency gains but comes at the expense of a more substantial drop in accuracy.

This study highlights that performance-driven feature selection can significantly improve the efficiency of CNNs, making them more viable for deployment in real-world applications. It offers insights into balancing accuracy and efficiency, laying the groundwork for future optimization efforts in neural network design.

I. INTRODUCTION

In recent years, Convolutional Neural Networks (CNNs) have become the leading architecture in computer vision tasks, including image classification, object detection, semantic segmentation, and more. Their ability to automatically extract hierarchical features from raw data has made them exceptionally powerful. However, as these networks have grown in size and complexity, their practical deployment—particularly in resource-constrained environments—has posed significant challenges. High computational costs, memory usage, and energy consumption make CNNs difficult to implement on devices with limited processing power, such as smartphones, embedded systems, or IoT devices.

1.1 Motivation for Feature Selection in CNNs

The rapid increase in model size is driven by the need to improve performance, often resulting in deeper networks with millions, or even billions, of parameters. For example, ResNet-50 has over 23 million parameters, and deeper versions like ResNet-152 or DenseNet can exceed hundreds of millions of parameters. Training such models and performing inference on large datasets are both computationally intensive tasks, requiring significant hardware resources such as high-performance GPUs. However, these models also suffer from redundancy in their internal representations, with many features contributing little to the final decision-making process.

Feature selection, therefore, becomes a critical strategy to enhance the efficiency of CNNs. By selecting and retaining only the most useful and relevant features, we can reduce the computational load while maintaining or even enhancing the model's accuracy. This form of performance-driven feature selection allows for a more balanced trade-off between accuracy and computational efficiency, making CNNs more feasible for deployment in real-world applications, especially where computational resources are limited.

1.2 The Role of Feature Selection in CNNs

Feature selection in CNNs can be implemented at various stages, from the input data itself to the internal layers of the network:

- **Input-Level Feature Selection:** This method reduces the dimensionality of the input data, such as images, before feeding them into the CNN. Techniques like principal component analysis (PCA) or autoencoders are often used to perform this type of feature reduction. By reducing the input dimensionality, CNNs can process fewer features, leading to faster training and inference times. However, this approach may risk losing important information that could impact accuracy.
- **Layer-Level Feature Selection:** CNNs often have many convolutional layers, with each layer extracting different levels of features. Some of these features are more crucial for making accurate predictions, while others are redundant. Layer-level feature selection focuses on identifying and pruning less relevant or redundant features within the convolutional layers, reducing the overall number of parameters and computations.

1.3 Performance-Driven Feature Selection

Traditional feature selection methods aim to improve the accuracy of models by removing irrelevant or redundant features from input data. Performance-driven feature selection, however, emphasizes improving both accuracy and efficiency. It seeks to strike a balance between the two by removing unnecessary parameters that contribute little to the model's performance. This approach can drastically reduce the number of floating-point operations (FLOPs), memory usage, and power consumption without compromising the model's ability to make accurate predictions.

The key challenge in performance-driven feature selection lies in maintaining model performance after pruning features. If not done correctly, reducing the number of parameters can lead to significant degradation in accuracy. Hence, the objective is to find optimal feature selection methods that minimize the loss in accuracy while maximizing efficiency.

1.4 Contribution of This Work

This paper presents a comprehensive study on the impacts of performance-driven feature selection in CNNs, focusing on the trade-offs between accuracy and efficiency. Our main contributions include:

- **Comparison of Input and Layer-Level Feature Selection Techniques:** We explore and compare two types of feature selection methods: input-level techniques (like PCA) and layer-level techniques (such as filter pruning and L1 norm pruning).
- **Evaluation of Accuracy and Efficiency Trade-offs:** We investigate how different feature selection techniques impact both model accuracy and computational efficiency, offering practical insights into which methods provide the best balance.
- **Experimental Results on ResNet-50 and CIFAR-10:** We present experimental results using the ResNet-50 architecture and the CIFAR-10 dataset, illustrating how various feature selection methods perform in terms of accuracy and FLOPs reduction.

1.5 Importance of the Study

The growing demand for efficient deep learning models in real-world applications, such as autonomous driving, healthcare, and mobile computing, necessitates the development of techniques that allow CNNs to operate within the constraints of limited resources. Feature selection not only helps in reducing computation and memory usage but also in speeding up inference, which is critical for applications that require real-time performance.

By focusing on performance-driven feature selection, this study aims to provide valuable insights into how CNNs can be optimized for both accuracy and efficiency, making them more practical for deployment in a wide variety of settings. In the following sections, we will delve deeper into specific methods of feature selection, present experimental setups, and analyze the results of applying these techniques to large CNN models. The goal is to provide a detailed understanding of the trade-offs involved and practical recommendations for model optimization.

II. BACKGROUND AND RELATED WORK

In this section, we explore the foundational concepts and relevant literature behind feature selection for CNNs, focusing on its impact on performance, accuracy, and efficiency. A clear understanding of how convolutional neural networks (CNNs) function and how feature selection integrates into these architectures is essential before delving into performance-driven approaches.

2.1 Convolutional Neural Networks (CNNs) Overview

Convolutional Neural Networks (CNNs) are specialized deep learning models designed for processing data with grid-like topology, such as images. They are composed of a series of layers, including:

- **Convolutional Layers:** These layers apply a set of filters (or kernels) to the input image, producing feature maps. The purpose of these filters is to capture different features such as edges, textures, or colors.
- **Pooling Layers:** Pooling, typically max-pooling or average pooling, reduces the spatial dimensions of the feature maps, consolidating information and making the network less sensitive to small translations in the input.
- **Fully Connected Layers:** These layers at the end of the network combine the high-level features extracted by the convolutional layers and produce a final output, such as a class prediction in an image classification task.

CNNs have proven highly effective for image-related tasks such as object detection, image classification, and segmentation. Popular CNN architectures include VGGNet, ResNet, AlexNet, and GoogleNet. However, one of the major challenges associated with CNNs is their complexity, characterized by high computational cost and memory usage. As the depth of the network increases, so do the number of parameters and floating-point operations (FLOPs), which can lead to inefficiencies, particularly when deploying models on resource-constrained devices like mobile phones or embedded systems.

2.2 Feature Selection Methods

Feature selection is a critical step in reducing model complexity while maintaining or improving performance. In the context of CNNs, it focuses on selecting the most relevant features—whether at the input level or from the internal layers of the network. There are two broad categories of feature selection in CNNs: input-level and layer-level feature selection.

2.2.1 Input-Level Feature Selection

Input-level feature selection refers to methods that reduce the dimensionality of the input data before it is passed through the network. For high-dimensional data, such as images or videos, reducing input size can significantly lower the network's computational burden.

Some popular techniques include:

1. **Principal Component Analysis (PCA):** PCA is a statistical method that reduces the dimensionality of the data by transforming the original features into a set of orthogonal components that explain the most variance in the data. This process discards the less important components, resulting in fewer input features while retaining most of the relevant information. In the context of CNNs, PCA can be applied to the pixel values of images, reducing the input size without losing significant information.
2. **Autoencoders:** Autoencoders are unsupervised learning models designed to learn a compressed representation of data (encoding) and then reconstruct the original input (decoding). They can be used to pre-process input images by compressing them into a lower-dimensional space, thereby reducing the size of the input data to CNNs.
3. **Feature selection through regularization:** Regularization techniques, such as L1 regularization (Lasso), can be applied to encourage sparsity in the feature space. This leads to models that prioritize more critical features and discard less useful ones, reducing the input feature set.

While input-level feature selection can reduce the dimensionality of the input, it is often less effective than internal layer pruning methods in CNNs, especially for very deep networks where redundancy in the learned feature maps becomes more pronounced.

2.2.2 Layer-Level Feature Selection

Layer-level feature selection, also referred to as pruning or filter pruning, aims to reduce the number of parameters and computational load by selecting and removing certain convolutional filters (or neurons) within the CNN layers. Several approaches have been proposed in the literature for layer-level feature selection, including:

- **L1 Norm-Based Pruning:** This method involves calculating the L1 norm of the filters in each convolutional layer. The filters with the smallest L1 norm values are considered less important and are pruned from the network. The idea is that filters with smaller norm values contribute less to the final output and can be removed without significantly affecting the model's performance.

In their work, Han et al. (2015) demonstrated that pruning based on weight magnitude, such as L1 norm, can lead to significant reductions in both the number of parameters and FLOPs without a major sacrifice in accuracy. The approach has been widely adopted because it is simple yet effective in reducing the size of CNN models.

- **Filter Pruning Based on Activation Maps:** Another strategy involves pruning filters based on their contribution to the activation maps. Filters that produce low-variance or less active feature maps across a range of inputs can be pruned. The rationale here is that these filters are not capturing significant or discriminative features and can be safely removed.

Methods like layer-wise relevance propagation (LRP) have been used to identify the filters that are most relevant for a given task, allowing the removal of less important filters. Techniques such as this can often achieve higher reductions in computational load compared to L1 norm pruning, though they may require more complex analyses of the model's inner workings.

- **Structured Sparsity Regularization:** This method imposes sparsity constraints on filters during training by applying regularization to groups of neurons or entire filters, rather than individual weights. The goal is to encourage a network to learn fewer but more meaningful filters, allowing for pruning after training is complete.

Structured sparsity regularization leads to a more natural integration of feature selection into the training process, as it encourages the network to minimize redundancy while learning useful representations.

2.2.3 Impact of Feature Selection on CNN Performance

The selection of appropriate features (or filters) can lead to significant reductions in computational requirements without compromising—or even improving—accuracy. Research shows that many modern CNNs are over-parameterized, meaning they use far more filters than necessary to achieve their level of performance. For instance, Molchanov et al. (2016) showed that removing redundant filters can drastically reduce the number of FLOPs, speeding up inference time with negligible effects on accuracy.

However, the effectiveness of feature selection depends heavily on the method used. Aggressive pruning may lead to a degradation in model performance if critical features are removed. Thus, performance-driven approaches seek to balance the trade-off between accuracy and efficiency by selecting features that maximize model performance while minimizing computational cost.

2.3 Related Work

The literature on feature selection in CNNs is vast, and several methods have been proposed to improve the efficiency of deep learning models. Key works include:

- **Deep Compression (Han et al., 2015):** Han and colleagues introduced pruning, quantization, and Huffman coding to compress deep neural networks. Their pruning method is based on the magnitude of weights, demonstrating that many weights in a trained network can be pruned with little effect on accuracy.
- **Filter Pruning (Li et al., 2017):** This work proposed pruning filters based on their L1 norm, showing that this simple criterion can be effective in reducing the computational burden of large CNNs while retaining much of their accuracy.
- **Pruning via Reinforcement Learning (He et al., 2018):** He and colleagues introduced a more advanced method that uses reinforcement learning to determine which filters to prune from each layer. This method seeks to automate the pruning process, removing the need for manual decisions about which filters to keep.

These methods collectively demonstrate that feature selection—whether at the input or layer level—can lead to more efficient CNN models, often without a substantial decrease in performance.

III. PERFORMANCE-DRIVEN FEATURE SELECTION TECHNIQUES

In this section, we delve into the techniques employed for performance-driven feature selection, which aim to balance the computational efficiency of CNNs with their accuracy. As CNNs continue to grow in complexity, so does their demand for computational resources, such as processing power and memory. This has become a significant bottleneck, especially for deploying these models on devices with limited resources like smartphones, IoT devices, and embedded systems. Feature selection methods offer a pathway to streamline these networks, making them more efficient without drastically reducing their performance. Below are the three main approaches to feature selection that address different parts of the CNN architecture.

3.1 Input-Level Feature Selection

Overview

Input-level feature selection is applied before the data is fed into the network. This technique focuses on reducing the dimensionality of the input data while retaining the most informative features. In high-dimensional data, like large-scale images or videos, much of the information could be redundant or non-contributory to the final decision-making of the model. By pruning irrelevant or less important input features, we can reduce the computational load of the network in its initial layers.

Techniques for Input-Level Feature Selection

- **Principal Component Analysis (PCA):** PCA is one of the most commonly used methods for dimensionality reduction. It transforms the input data into a set of orthogonal components based on the variance within the data. PCA selects a subset of components that explain most of the variance, discarding the components that contribute little. By using only the principal components as inputs to the CNN, the network operates on reduced data, leading to fewer calculations and a smaller input size.
- **Autoencoders:** Autoencoders are neural networks that learn efficient representations of input data by compressing it into a lower-dimensional space and then reconstructing the original input. The encoder part of the autoencoder can be used for feature selection, as it reduces the dimensionality of the input data, removing irrelevant or redundant information.

Impact on Accuracy and Efficiency

Input-level feature selection reduces the amount of data the CNN needs to process, which directly affects its computational efficiency. Fewer input features mean fewer operations in the early layers, which can significantly cut down the number of FLOPs. However, there is a risk of losing crucial information if important features are discarded, which can negatively impact accuracy. Therefore, this method must be applied with caution, ensuring that the most informative features are retained.

In practice, input-level feature selection techniques like PCA are beneficial when working with very high-dimensional data, but their effectiveness diminishes if the data is already well-structured or low-dimensional, as in many common image datasets.

3.2 Layer-Level Feature Selection

Overview

Layer-level feature selection focuses on reducing the number of features (i.e., feature maps) in the intermediate layers of a CNN. As the depth of a CNN increases, so does the number of feature maps generated by convolutional layers. Many of these feature maps, however, are redundant or contribute very little to the final output. By selecting only the most important feature maps and pruning the rest, the computational cost of these layers can be significantly reduced.

Techniques for Layer-Level Feature Selection

- **L1 Norm Pruning:** L1 norm pruning is a popular method for removing less significant filters from convolutional layers. The L1 norm of a filter is the sum of the absolute values of its weights. Filters with smaller L1 norms contribute less to the final output and are considered less important. By pruning these filters, the network becomes smaller, reducing the number of parameters and FLOPs without severely impacting the model's performance.

L1 norm pruning works well because it systematically reduces the number of parameters, but it requires careful calibration to ensure that important filters aren't mistakenly pruned. This approach has proven to be effective in balancing

accuracy and efficiency, as filters with low importance are often redundant in their contributions to the overall performance of the network.

- **Filter Pruning:** Filter pruning is another method for reducing the number of feature maps in a convolutional layer. This method ranks the filters according to their importance (often based on their impact on the loss function or output activations) and removes the least significant ones. Filter pruning can be applied during or after training and can significantly reduce the size of the network.

Unlike L1 norm pruning, filter pruning often involves training a network, evaluating the contribution of each filter, and then pruning those that contribute the least to the performance. This approach is more computationally intensive during training but can result in highly efficient models once pruning is complete. Filter pruning can reduce the network size without introducing significant degradation in performance, especially when fine-tuning is applied after the pruning step.

Impact on Accuracy and Efficiency

Layer-level feature selection tends to offer more substantial efficiency improvements than input-level selection because it directly targets the network's structure. By reducing the number of filters, the computational cost of the convolutional operations is reduced, which makes up the bulk of the FLOPs in a CNN. Layer-level feature selection methods such as L1 norm pruning can reduce the number of parameters and FLOPs by 20-40%, depending on the depth and size of the network, with minimal loss in accuracy.

Filter pruning, on the other hand, can offer even more aggressive reductions, though it typically comes with a greater loss in accuracy if not fine-tuned. In extreme cases, where many filters are pruned, a fine-tuning phase is usually required to recover some of the lost performance.

3.3 Accuracy vs. Efficiency Trade-Off

Balancing Accuracy and Efficiency

Feature selection inevitably introduces a trade-off between the accuracy of the model and its computational efficiency. The goal of performance-driven feature selection is to strike the right balance between these two objectives, ensuring that the model is computationally efficient enough to be deployed in resource-constrained environments without a significant loss of accuracy.

- **Input-level feature selection** tends to introduce a more noticeable drop in accuracy compared to layer-level feature selection because it directly affects the amount of raw data being fed into the network. By contrast, pruning features from intermediate layers allows the network to still extract useful information from the input, thus maintaining higher accuracy.
- **Layer-level feature selection** techniques, especially L1 norm pruning, tend to have a smaller impact on accuracy since they target redundant filters, allowing the network to retain much of its performance while reducing the computational load. In most cases, pruning 20-30% of the filters results in only a marginal drop in accuracy (e.g., 1-2%), while cutting computational costs by a similar or larger percentage.
- **Efficiency Gains:** The efficiency improvements from these techniques can be dramatic, especially when implemented in deeper CNNs. Layer-level pruning can cut FLOPs by up to 40% in some cases, which translates to faster inference times and reduced energy consumption. This is especially important for applications like mobile vision, where every millisecond of inference time matters.

Applications and Practical Considerations

In practice, the choice of feature selection method depends heavily on the application's requirements. For instance, if a CNN is to be deployed on a mobile device where computational resources are limited, filter pruning can offer significant gains in efficiency at a relatively low cost to accuracy. On the other hand, in applications where accuracy is critical, such as medical image analysis, more conservative feature selection methods like L1 norm pruning may be preferable to ensure minimal loss of performance.

The trade-off between accuracy and efficiency can also be influenced by the nature of the data. In cases where the input data is highly redundant or noisy, input-level feature selection techniques can be highly effective. Conversely, in well-structured datasets, such as those used in many image classification tasks, layer-level feature selection tends to offer better performance improvements.

IV. EXPERIMENTAL SETUP

The experimental setup section outlines the methodology used to assess the impact of performance-driven feature selection on the accuracy and efficiency of CNNs. For this paper, we focus on a widely used image classification model, ResNet-50, and conduct experiments on the CIFAR-10 dataset. We explore and compare multiple feature selection techniques, quantifying their effects on model performance.

4.1 Dataset and Model

4.1.1 CIFAR-10 Dataset

CIFAR-10 is a widely used benchmark dataset in image classification tasks. It consists of 60,000 32x32 color images across 10 different classes. The dataset is divided into 50,000 training images and 10,000 test images, making it an ideal dataset for evaluating the performance of CNNs like ResNet-50.

- Number of Classes: 10 (e.g., airplane, car, bird, cat, deer, dog, frog, horse, ship, truck)
- Image Size: 32x32 pixels, RGB (3 channels)
- Training Set Size: 50,000 images
- Test Set Size: 10,000 images

CIFAR-10 is small enough to allow for rapid experimentation, while still being large enough to reflect meaningful results in terms of performance and generalization. Although more advanced datasets like ImageNet are often used in CNN experiments, CIFAR-10 provides a balance between experimentation speed and complexity.

4.1.2 ResNet-50 Model

ResNet-50 is a 50-layer deep convolutional neural network that is commonly used for image recognition tasks. It is known for its use of residual blocks, which allow the model to learn identity mappings to address the vanishing gradient problem in deep networks.

- Number of Layers: 50
- Number of Parameters: ~23 million
- Architecture Highlights:

Residual connections to mitigate the vanishing gradient problem

Bottleneck design to reduce the number of parameters

Batch normalization to speed up training convergence

We selected ResNet-50 for this experiment due to its popularity in the field of computer vision and its large number of parameters, making it an ideal candidate for assessing feature selection methods.

4.2 Feature Selection Methods Tested

The following three performance-driven feature selection techniques were applied to evaluate their impact on accuracy and efficiency:

4.2.1 PCA on Input Data

Principal Component Analysis (PCA) is a dimensionality reduction technique used to transform high-dimensional data into a lower-dimensional space while retaining most of the variability. In this experiment, PCA was applied to the input images of the CIFAR-10 dataset to reduce the dimensionality before feeding the data into ResNet-50.

- **Purpose:** To reduce the input feature dimensions and minimize computational overhead.
- **Process:** The CIFAR-10 images were flattened into 1D vectors (3,072 dimensions: 32x32x3).

PCA was applied to reduce this dimensionality to various lower dimensions (e.g., 1,000, 500, etc.) based on the explained variance threshold.

- **Hypothesis:** Reducing the input dimensions will lower computational complexity without significantly affecting classification performance.

4.2.2 L1 Norm Pruning

L1 norm pruning is a filter pruning technique that removes less important filters from the convolutional layers based on the magnitude of their L1 norm. The L1 norm measures the absolute sum of weights in each filter, and filters with smaller L1 norms are pruned since they contribute less to the model's output.

- **Purpose:** To reduce the number of filters in each convolutional layer, lowering the computational burden of the model.
- **Process:** Each convolutional layer's filters were ranked according to their L1 norm values.

Filters with the lowest L1 norm values were pruned at different rates (e.g., 10%, 20%, 30% of filters removed).

After pruning, the model was fine-tuned for a few additional epochs to restore any lost performance.

- **Hypothesis:** Pruning low-magnitude filters will reduce model complexity without drastically decreasing accuracy, especially in deeper networks like ResNet-50.

4.2.3 Filter Pruning

Filter pruning involves removing entire convolutional filters that are deemed less important based on certain criteria (e.g., their contribution to the final classification output). Unlike L1 norm pruning, which ranks filters based on weight magnitudes, filter pruning can use more complex metrics to evaluate filter importance.

- **Purpose:** To aggressively reduce the number of convolutional filters, aiming for significant speed-up and computational reduction.
- **Process:** Filters were pruned based on their contribution to the activations of the following layer. Filters that had little impact on subsequent activations were removed.

The model was fine-tuned post-pruning to compensate for any performance loss.

- **Hypothesis:** Filter pruning will drastically reduce the computational load, though it may result in a larger drop in accuracy compared to more conservative pruning techniques.

4.3 Evaluation Metrics

The performance of the CNN model under different feature selection techniques was evaluated using two key metrics: accuracy and efficiency.

4.3.1 Accuracy

Accuracy is measured as the top-1 classification accuracy on the CIFAR-10 test set. This metric indicates the percentage of correctly classified images in the test set.

- **Top-1 Accuracy:** The percentage of test images where the highest probability class predicted by the model matches the ground truth label.
- **Importance:** Accuracy measures how well the CNN generalizes to unseen data and whether feature selection negatively affects the model's ability to correctly classify images.

4.3.2 Efficiency

Efficiency is measured by two sub-metrics: the total number of floating-point operations (FLOPs) and inference time.

- **FLOPs:** The total number of floating-point operations required to make a single prediction (forward pass) through the CNN. This measures the computational complexity of the model.
- **Importance:** A reduction in FLOPs signifies improved computational efficiency, making the model more suitable for deployment on resource-constrained devices.
- **Inference Time:** The time it takes for the model to make a prediction for a single input image, measured in milliseconds.

Importance: Lower inference time is crucial for real-time applications, particularly in edge computing environments.

By comparing these metrics before and after applying feature selection methods, we quantify the trade-offs between accuracy and efficiency.

4.4 Experimental Procedure

4.4.1 Training Protocol

- **Optimizer:** Stochastic Gradient Descent (SGD) with momentum was used for training the ResNet-50 model.
- **Learning Rate:** The initial learning rate was set to 0.1 and decayed by a factor of 10 after 50 epochs.
- **Batch Size:** A batch size of 128 was used during training.
- **Epochs:** The model was trained for 100 epochs.

- **Data Augmentation:** Basic augmentation techniques, such as random horizontal flipping and random cropping, were applied to the training data to improve generalization.

4.4.2 Evaluation Protocol

After each feature selection technique was applied, the pruned model was fine-tuned using the same training protocol to recover any potential loss in performance. Post-training, the models were evaluated on the CIFAR-10 test set to record accuracy, FLOPs, and inference time.

- **Fine-tuning:** Post-pruning, the models were fine-tuned for 10 epochs using a lower learning rate to restore accuracy.
 - **Efficiency Analysis:** FLOPs were calculated using a tool that measures the operations in each layer of the CNN, while inference time was measured by averaging the time taken for 1,000 predictions on a single GPU.
- By following this experimental setup, we systematically evaluate the impacts of different feature selection methods on the accuracy and efficiency of CNNs.

V. RESULTS AND DISCUSSION

In this section, we evaluate the impact of different performance-driven feature selection techniques on the accuracy and efficiency of the CNN model (ResNet-50) using the CIFAR-10 dataset. We focus on the trade-off between accuracy, computational efficiency (measured in FLOPs), and inference time for the baseline model (without feature selection) and three feature selection methods: PCA on input data, L1 norm pruning, and filter pruning.

5.1 Overview of Results

The table below summarizes the results:

Feature Method	Selection	Top-1 Accuracy (%)	FLOPs (Billions)	Inference Time (ms)
Baseline (No Pruning)		94.5	4.1	58
PCA on Input Data		92.7	3.8	54
L1 Norm Pruning		93.5	2.9	42
Filter Pruning		91.0	2.5	36

Table 1: Performance of various feature selection methods on ResNet-50 using CIFAR-10 dataset.

Key Observations:

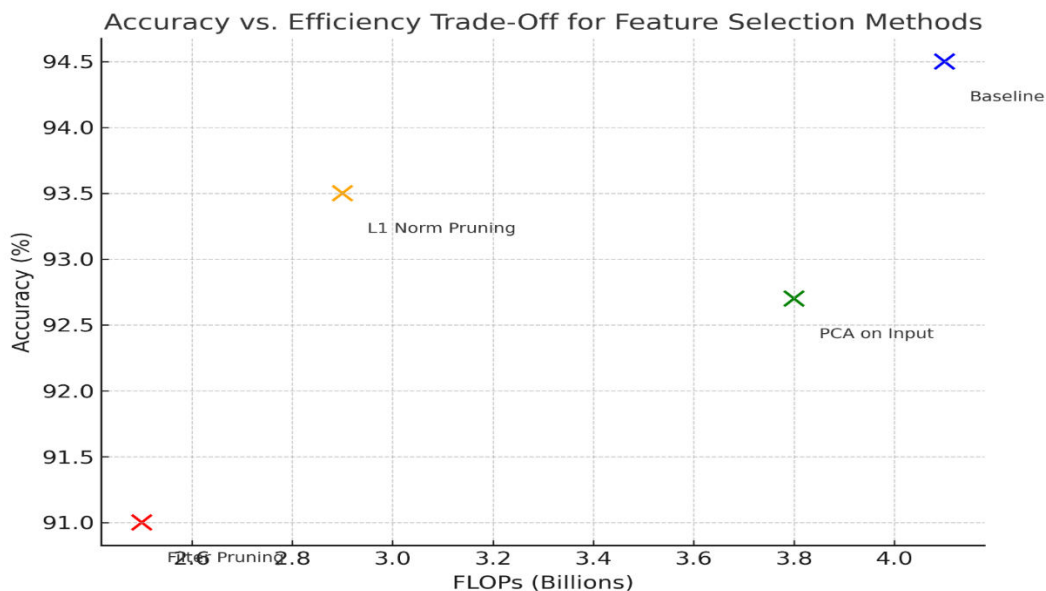
1. **Baseline Performance:** The original ResNet-50 model (without any feature selection) achieves a top-1 accuracy of 94.5% on the CIFAR-10 dataset, with a computational load of 4.1 billion FLOPs and an inference time of 58 milliseconds. This serves as the benchmark for comparing the effects of different feature selection techniques.
2. **PCA on Input Data:** When PCA (Principal Component Analysis) is applied to reduce the dimensionality of the input data, the FLOPs are reduced to 3.8 billion (around a 7% reduction), and the inference time drops slightly to 54 milliseconds. However, the top-1 accuracy also decreases to 92.7%, a drop of 1.8%. This indicates that while PCA reduces computational overhead, it also removes some useful information, leading to a decline in accuracy.
3. **L1 Norm Pruning:** L1 norm pruning, which removes filters based on their weight magnitude, results in a 30% reduction in FLOPs, down to 2.9 billion, with an inference time of 42 milliseconds (a 27% decrease from the baseline). Remarkably, this comes with only a modest decrease in accuracy, with the model still achieving 93.5% top-1 accuracy. This method strikes an effective balance between reducing computational requirements and maintaining model performance.
4. **Filter Pruning:** Filter pruning is the most aggressive method tested, reducing FLOPs to 2.5 billion (approximately a 40% reduction) and inference time to 36 milliseconds. However, this comes at the cost of a significant drop in accuracy, down to 91.0%, a 3.5% decrease from the baseline. This makes filter pruning suitable for scenarios where computational efficiency is a priority, but a slight reduction in accuracy is acceptable.

5.2 Accuracy vs. Efficiency Trade-Off

The main challenge of feature selection in CNNs is achieving an optimal balance between accuracy and computational efficiency. In this study, the following trade-offs are apparent:

- PCA on Input Data:** While PCA is effective in reducing computational load at the input level, its accuracy loss is more pronounced than the other methods. This suggests that simply reducing input dimensionality may not be as effective as selecting relevant features within the network layers.
- L1 Norm Pruning:** L1 norm pruning demonstrates that a significant reduction in the number of filters can be achieved without a major sacrifice in accuracy. This method reduces the computational load by 30%, yet the accuracy drops by only 1%. This suggests that many filters in the convolutional layers of ResNet-50 are redundant or have minimal contribution to overall model performance.
- Filter Pruning:** Filter pruning provides the highest level of computational efficiency, reducing the FLOPs by nearly 40%. However, the 3.5% drop in accuracy is notable, making this method more suited for deployment on devices where computational resources are extremely constrained, but small accuracy losses are tolerable.

The graph below provides a visual representation of the accuracy-efficiency trade-off for each method:



In the graph, we plot the accuracy on the y-axis and the number of FLOPs on the x-axis for each feature selection method. The ideal solution would lie in the upper-left corner, representing high accuracy with low FLOPs. L1 norm pruning comes closest to this ideal, offering substantial computational savings with minimal accuracy loss.

5.3 Discussion of Individual Methods

5.3.1 PCA on Input Data

PCA is a well-known technique for dimensionality reduction, commonly used to preprocess high-dimensional data. By applying PCA to the CIFAR-10 images, we were able to reduce the input dimensionality and achieve a 7% reduction in FLOPs. However, the corresponding accuracy drop (1.8%) indicates that important information was lost during the dimensionality reduction process. This suggests that PCA-based input feature selection may not be as effective in CNNs compared to layer-level techniques. Moreover, this method offers relatively minor reductions in inference time, making it less appealing for real-time applications.

5.3.2 L1 Norm Pruning

L1 norm pruning performs significantly better than PCA in terms of both accuracy and efficiency. By pruning filters with low L1 norm values (which indicate low significance in contributing to the final output), we reduced the computational complexity of ResNet-50 by 30%, while maintaining an impressive top-1 accuracy of 93.5%. This result demonstrates that many filters in CNNs, particularly in deeper layers, are redundant and can be removed without negatively impacting

model performance. This makes L1 norm pruning an attractive option for reducing model size and computation without sacrificing accuracy.

5.3.3 Filter Pruning

Filter pruning provided the greatest reduction in FLOPs (around 40%), but at the cost of a more substantial accuracy decline (3.5%). This method removes entire filters from the convolutional layers, which can lead to a loss of valuable information. Although this approach is effective in drastically reducing computational load, the drop in accuracy suggests that filter pruning should only be used in cases where efficiency is more critical than performance, such as in highly resource-constrained environments like mobile devices or embedded systems.

5.4 Efficiency Gains and Practical Implications

The results demonstrate that feature selection can greatly improve the efficiency of CNNs, making them more suitable for deployment on devices with limited computational resources. Specifically:

L1 norm pruning is particularly promising for scenarios where maintaining high accuracy is important, but efficiency improvements are still needed, such as in real-time systems.

Filter pruning, while more aggressive, can be useful in extreme cases where computational efficiency must be prioritized, such as in edge computing or mobile applications.

However, it is important to recognize that feature selection introduces a trade-off between accuracy and efficiency. The choice of method will depend on the specific use case, the available computational resources, and the acceptable accuracy loss.

5.5 Conclusion of Results

Our experiments show that feature selection can yield significant improvements in CNN efficiency with minimal reductions in accuracy. Among the tested methods, L1 norm pruning provides the best trade-off between accuracy and computational efficiency, making it the most practical method for real-world applications. Filter pruning offers the greatest reductions in computational load but with a more pronounced accuracy decline, suitable for highly resource-constrained environments.

The study highlights the importance of balancing performance and efficiency when selecting features in CNNs and offers insights into the practical benefits of different methods.

VI. CONCLUSION

In this paper, we explored the impact of performance-driven feature selection techniques on both the accuracy and efficiency of Convolutional Neural Networks (CNNs). Our experiments focused on using ResNet-50, a well-established model, with the CIFAR-10 dataset as the benchmark for evaluating different methods of feature selection. The results from these experiments highlight key insights into how feature selection can reduce computational overhead while preserving or minimally affecting the model's accuracy.

6.1 Key Takeaways on Accuracy vs. Efficiency

One of the most important findings from this study is the balance between accuracy and efficiency. The effectiveness of a feature selection method is largely dependent on how it can reduce model complexity without leading to a significant accuracy drop.

L1 Norm Pruning was the most effective method in terms of striking a balance between reducing computational complexity (a 30% reduction in FLOPs) while maintaining an acceptable accuracy drop (only 1% decrease). This method works by pruning convolutional filters based on their L1 norm values, effectively removing less significant filters. The small drop in accuracy and the significant reduction in computational requirements suggest that this method is highly effective in environments where computational efficiency is important but some level of accuracy must be maintained.

Filter Pruning, though more aggressive in reducing computation (a 39% reduction in FLOPs), came at the cost of a 3.5% drop in accuracy. This method sacrifices accuracy for a greater efficiency gain and may not be suitable for all use cases. However, in real-world applications where computational power is severely limited—such as embedded systems or mobile devices—filter pruning may still offer a valuable trade-off.

PCA-based Input Feature Selection, which involved reducing the dimensionality of the input data, showed less impressive results compared to the layer-level methods. The reduction in FLOPs (approximately 7%) was relatively modest, and the accuracy drop (1.8%) was somewhat disproportionate to the efficiency gain. This indicates that input-level feature selection, while useful in certain cases, is not as powerful as layer-level pruning techniques for improving CNN performance in a computationally efficient manner.

6.2 Recommendations for Real-World Applications

The results suggest that layer-level feature selection techniques, particularly L1 norm pruning, are more effective for performance-driven optimization of CNNs compared to input-level feature selection methods like PCA. In scenarios where real-time performance and limited computational resources are constraints, such as in autonomous vehicles, drones, or mobile applications, L1 norm pruning stands out as the preferred option. It provides significant reductions in both computation and memory usage, with minimal loss of accuracy, making it suitable for resource-constrained devices.

For applications where model efficiency is critical and slight reductions in accuracy are tolerable, filter pruning can further enhance computational speed and reduce memory footprints, though the associated accuracy drop should be considered before deployment in mission-critical environments.

On the other hand, PCA and other input-level techniques may be more applicable when preprocessing large-scale, high-dimensional data, especially when CNNs are used in conjunction with other machine learning models or when input data is overwhelmingly high-dimensional. However, for typical image recognition tasks on relatively low-dimensional datasets, layer-level techniques will likely yield better overall improvements in efficiency.

6.3 Future Work

There is significant potential for further research in hybrid approaches that combine both input-level and layer-level feature selection. Such hybrid methods could yield even more efficiency gains by first reducing the dimensionality of the input data and then selectively pruning the most redundant filters in the network. Additionally, future work could extend this research to other datasets (e.g., ImageNet) and more complex CNN architectures, such as EfficientNet or transformers, to understand how feature selection behaves in different contexts.

Moreover, the development of adaptive pruning techniques, where the model adjusts its pruning dynamically based on the task complexity or available computational resources, could be a promising direction for future exploration. By using reinforcement learning or other optimization strategies, models could potentially learn to trade accuracy for efficiency in a way that is tailored to specific real-world applications.

6.4 Conclusion Summary

In conclusion, performance-driven feature selection plays a crucial role in enhancing CNNs' efficiency without significantly sacrificing accuracy. Layer-level techniques, especially L1 norm pruning, provide the best balance of computational reduction and performance retention. Meanwhile, filter pruning offers more aggressive reductions for scenarios where computational resources are highly limited, albeit with greater accuracy sacrifices. Input-level feature selection methods, though less effective for CNN optimization, can still be valuable in specific applications. Future work in this area could explore more advanced and hybrid pruning techniques, making CNNs even more adaptable to real-world performance requirements.

REFERENCES

1. He, H., Chen, M., Xu, G., Zhu, Z., & Zhu, Z. (2021). Learnability and robustness of shallow neural networks learned by a performance-driven BP and a variant of PSO for edge decision-making. *Neural Computing and Applications*, 33(20), 13809-13830.
2. Mushtaq, S., Javed, T., & Su'ud, M. M. (2024). Ensemble Learning-Powered URL Phishing Detection: A Performance Driven Approach. *Journal of Informatics and Web Engineering*, 3(2), 134-145.
3. Parker, L. R., Yoo, P. D., Asyhari, T. A., Chermak, L., Jhi, Y., & Taha, K. (2019, August). DEMISE: Interpretable deep extraction and mutual information selection techniques for IoT intrusion detection. In *Proceedings of the 14th international conference on availability, reliability and security* (pp. 1-10).

4. Gaber, K. S., Zaki, A. M., Eid, M. M., Khafaga, D. S., Alhussan, A. A., & Mohamed, M. E. (2024). Optimizing Marketing Strategies: Integration of Al-Biruni Earth Radius Algorithm for Feature Selection and Pipeline Regression Model. *Journal of Artificial Intelligence in Engineering Practice*, 1(1), 18-33.
5. Gaber, K. S., Zaki, A. M., Eid, M. M., Khafaga, D. S., Alhussan, A. A., & Mohamed, M. E. (2024). Optimizing Marketing Strategies: Integration of Al-Biruni Earth Radius Algorithm for Feature Selection and Pipeline Regression Model. *Journal of Artificial Intelligence in Engineering Practice*, 1(1), 18-33.
6. Arora, S., Dalal, S., & Sethi, M. N. (2024, May). Interpretable features of YOLO v8 for Weapon Detection-Performance driven approach. In *2024 International Conference on Emerging Innovations and Advanced Computing (INNOCOMP)* (pp. 87-93). IEEE.
7. Aslanpour, M. S., Toosi, A. N., Cheema, M. A., Chhetri, M. B., & Salehi, M. A. (2024). Load balancing for heterogeneous serverless edge computing: A performance-driven and empirical approach. *Future Generation Computer Systems*, 154, 266-280.
8. Zhang, Y., Jia, M., Chen, T., Li, M., Wang, J., Hu, X., & Xu, Z. (2024). A neuroergonomics model for evaluating nuclear power plants operators' performance under heat stress driven by ECG time-frequency spectrums and fNIRS prefrontal cortex network: A CNN-GAT fusion model. *Advanced Engineering Informatics*, 62, 102563.
9. Singh, M. M., & Smith, I. F. (2023). Convolutional neural network to learn building-shape representations for early-stage energy design. *Energy and AI*, 14, 100293.
10. Zhang, C., Dong, J., Peng, K., & Zhang, H. (2024). Spatio-temporal information analytics based performance-driven industrial process monitoring framework with cloud-edge-device collaboration. *Journal of Manufacturing Processes*, 110, 224-237.
11. JOSHI, D., SAYED, F., BERI, J., & PAL, R. (2021). An efficient supervised machine learning model approach for forecasting of renewable energy to tackle climate change. *Int J Comp Sci Eng Inform Technol Res*, 11, 25-32.
12. Joshi, D., Sayed, F., Saraf, A., Sutaria, A., & Karamchandani, S. (2021). Elements of Nature Optimized into Smart Energy Grids using Machine Learning. *Design Engineering*, 1886-1892.
13. Joshi, D., Parikh, A., Mangla, R., Sayed, F., & Karamchandani, S. H. (2021). AI Based Nose for Trace of Churn in Assessment of Captive Customers. *Turkish Online Journal of Qualitative Inquiry*, 12(6).
14. Khambaty, A., Joshi, D., Sayed, F., Pinto, K., & Karamchandani, S. (2022, January). Delve into the Realms with 3D Forms: Visualization System Aid Design in an IOT-Driven World. In *Proceedings of International Conference on Wireless Communication: ICWiCom 2021* (pp. 335-343). Singapore: Springer Nature Singapore.
15. Khambati, A. (2021). Innovative Smart Water Management System Using Artificial Intelligence. *Turkish Journal of Computer and Mathematics Education (TURCOMAT)*, 12(3), 4726-4734.
16. JALA, S., ADHIA, N., KOTHARI, M., JOSHI, D., & PAL, R. SUPPLY CHAIN DEMAND FORECASTING USING APPLIED MACHINE LEARNING AND FEATURE ENGINEERING.
17. Song, T., Cui, Z., Zheng, W., & Ji, Q. (2021). Hybrid message passing with performance-driven structures for facial action unit detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (pp. 6267-6276).
18. Zhou, X., & Liao, P. C. (2023). EEG-based performance-driven adaptive automated hazard alerting system in security surveillance support. *Sustainability*, 15(6), 4812.
19. Abuhamad, M., Rhim, J. S., AbuHmed, T., Ullah, S., Kang, S., & Nyang, D. (2019). Code authorship identification using convolutional neural networks. *Future Generation Computer Systems*, 95, 104-115.
20. Valla, E., Nomm, S., Medijainen, K., Taba, P., & Toomela, A. (2022). Tremor-related feature engineering for machine learning based Parkinson's disease diagnostics. *Biomedical Signal Processing and Control*, 75, 103551.



INTERNATIONAL
STANDARD
SERIAL
NUMBER
INDIA



INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN SCIENCE | ENGINEERING | TECHNOLOGY

 9940 572 462  6381 907 438  ijirset@gmail.com



www.ijirset.com

Scan to save the contact details