

# Enhancing Image Classification with Attention-Driven Convolutional Neural Networks on CIFAR Datasets

Gerasimos Vonitsanos  
*Computer Engineering and  
Informatics Department  
University of Patras, Patras, Greece*  
mvonitsanos@ceid.upatras.gr

Emmanouela-Electra Economopoulou  
*Computer Engineering and  
Informatics Department  
University of Patras, Patras, Greece*  
std1057466@ceid.upatras.gr

Spyros Sioutas  
*Computer Engineering and  
Informatics Department  
University of Patras, Patras, Greece*  
sioutas@ceid.upatras.gr

Andreas Kanavos  
*Department of Informatics  
Ionian University, Corfu, Greece*  
akanavos@ionio.gr

Phivos Mylonas  
*Department of Informatics and Computer Engineering  
University of West Attica, Athens, Greece*  
mylonasf@uniwa.gr

**Abstract**—Image classification is a core task in computer vision with wide-ranging applications, from autonomous vehicles to medical diagnostics. While Convolutional Neural Networks (CNNs) have demonstrated strong performance by learning spatial hierarchies of features, they often struggle to capture complex interdependencies among spatial and channel-wise representations. This work proposes an attention-augmented CNN architecture that integrates both Squeeze-and-Excitation (SE) and Convolutional Block Attention Module (CBAM) mechanisms to enhance feature selection and improve classification performance. The proposed model is evaluated on two benchmark datasets, CIFAR-10 and CIFAR-100, and compared against baseline CNN and Artificial Neural Network (ANN) architectures. Experimental results indicate that the attention-enhanced CNN achieves superior classification accuracy and generalization, with notable gains in distinguishing visually similar classes. These findings highlight the effectiveness of combining channel and spatial attention modules to improve the robustness and adaptability of deep learning-based visual recognition systems.

**Index Terms**—Image Classification, Convolutional Neural Networks (CNN), Attention Mechanisms, Channel Attention, Spatial Attention, Deep Learning, CIFAR-10, CIFAR-100, Neural Networks, Visual Recognition

## I. INTRODUCTION

In recent years, the rapid advancement of Artificial Intelligence (AI), particularly in deep learning, has driven transformative progress across various scientific and technological domains. Machine Learning (ML) techniques have revolutionized key areas such as image classification, autonomous driving, natural language processing, and environmental data analytics [21]. Among these, image classification stands out as a fundamental task, supporting numerous real-world applications, including medical image analysis and intelligent surveillance systems. Accurate and efficient object recognition

within images is critical to the functionality and reliability of such technologies.

Traditional image classification methods relied on feature extraction techniques that required manual intervention to identify and process relevant image characteristics. Approaches such as Haar Cascades and Histogram of Oriented Gradients (HOG) depend on prior knowledge of the features to be detected, limiting their adaptability and scalability [4]. With the advent of Artificial Intelligence and the emergence of Convolutional Neural Networks (CNNs), models gained the capability to automatically learn hierarchical features directly from image data, eliminating the need for predefined feature engineering [12].

Despite their success, conventional CNNs have shown limitations in capturing the complex interdependencies among spatial and channel-wise features within an image. To address these shortcomings, more advanced techniques have been introduced, most notably, attention mechanisms. These modules allow models to dynamically focus on the most informative regions of an image, enhancing the precision and relevance of feature extraction. Initially developed for natural language processing tasks in architectures such as the Transformer, attention mechanisms have since been adapted for visual recognition to improve the discriminative power of CNNs [20]. By prioritizing critical image regions and suppressing less relevant information, attention-based CNNs offer improved performance in object classification tasks.

Our study introduces a novel integration of channel and spatial attention mechanisms within a unified CNN framework by combining Squeeze-and-Excitation (SE) and Convolutional Block Attention Module (CBAM) blocks. This design leverages the complementary strengths of both mechanisms to enhance feature selection at multiple representation levels. In addition to evaluating the proposed model against standard CNN

and ANN architectures, our approach systematically examines performance on datasets of varying complexity—CIFAR-10 and CIFAR-100—offering insights into both accuracy gains and computational trade-offs. By focusing on low-dimensional image data, the work addresses a less-explored area of attention-based vision research. At the same time, class-specific error analysis demonstrates the model’s ability to resolve ambiguities between visually similar categories. Furthermore, the inclusion of model complexity and training cost analysis ensures practical relevance for real-world deployments.

The remainder of the paper is organized as follows: Section II reviews related work, highlighting advancements in deep learning and attention mechanisms for image classification. Section III presents the neural network architectures explored in this study, including a traditional Artificial Neural Network (ANN), a CNN, and a CNN enhanced with attention modules. Section IV outlines the technical setup, including data preprocessing, model training configuration, and development tools. Section V delivers a comparative evaluation of the models on the CIFAR-10 and CIFAR-100 datasets using established performance metrics. Finally, Section VII summarizes the key findings and discusses future directions for enhancing accuracy and computational efficiency.

## II. RELATED WORK

The integration of attention mechanisms into CNNs has proven to be highly effective for improving image classification performance. Hu et al. introduced Squeeze-and-Excitation (SE) blocks, which adaptively recalibrate channel-wise feature responses via a squeeze-and-excitation operation, significantly enhancing accuracy in popular architectures like ResNet and Inception [7]. Building on this concept, Wang et al. proposed ECA-Net, an efficient channel attention method that removes fully connected layers to reduce complexity without compromising performance [23]. Similarly, Li et al. developed Selective Kernel Networks, which dynamically choose the optimal convolutional kernel size for each category, improving adaptability [14].

Beyond channel attention, spatial attention mechanisms have been integrated to enhance feature learning further. Woo et al. introduced Convolutional Block Attention Module (CBAM), a lightweight module that sequentially applies channel and spatial attention to improve focus on informative features while maintaining computational efficiency [25]. Roy et al. extended SE blocks with Concurrent Spatial and Channel Squeeze and Excitation (SCSE), which concurrently applies spatial and channel attention, originally for segmentation tasks but also effective for image classification [19].

More advanced approaches include Non-local Neural Networks, which capture long-range dependencies by directly correlating distant regions of an image [24], and Attention Augmented Convolutional Networks, which integrate self-attention mechanisms with convolutional layers for richer feature representations [1]. Zhao et al. explored self-attention as a standalone mechanism for image classification [27].

Residual and bottleneck-based attention mechanisms have also emerged. Wang et al. combined residual learning with attention modules in Residual Attention Networks, enabling the network to suppress less relevant features across layers [22]. Park et al. proposed the Bottleneck Attention Module (BAM) to enhance intermediate feature representations at bottlenecks [18]. Jetley et al. presented Learn to Pay Attention, an end-to-end approach where attention maps are jointly learned with CNN parameters [9].

Other studies have refined the balance between global and local feature information. Zhang et al. introduced SA-Net (Shuffle Attention), which redistributes information across channels and spatial locations [26], while Lyu et al. proposed a coarse-to-fine global/local attention framework for multi-label classification [15].

While numerous studies have investigated attention mechanisms for visual recognition, most have applied either channel attention, such as the SE, or spatial attention, such as the CBAM, independently. Even in cases where both mechanisms are used, they are often evaluated in isolation or sequentially without a systematic comparison against non-attention baselines across datasets of different complexity. In this work, we integrate SE and CBAM within a unified CNN architecture to leverage their complementary strengths in enhancing feature selection. We further benchmark this design against standard CNN and ANN models on CIFAR-10 and CIFAR-100, providing a detailed analysis of classification performance, computational cost, and class-specific improvements — aspects often overlooked in prior attention-based studies.

## III. METHODOLOGY FOUNDATION

This section presents the methodological framework of the study, outlining the neural network architectures employed and the key techniques integrated to enhance model performance.

### A. Artificial Neural Networks

Artificial Neural Networks (ANNs) are computational models inspired by the structure of the human brain. They consist of interconnected layers of nodes (neurons), where each connection has an associated weight that is adjusted during training. ANNs are particularly useful for learning complex nonlinear relationships in data [13]. In this study, ANNs are used as a baseline model to evaluate the performance gains introduced by deeper and attention-augmented architectures.

### B. Convolutional Neural Networks

CNNs are a specialized class of neural networks designed for processing grid-like data such as images. They use convolutional layers to automatically learn spatial hierarchies of features through filters, significantly reducing the need for manual feature extraction. CNNs have become the standard for image classification tasks due to their accuracy, efficiency, and scalability [12]. The convolutional approach has proven effective across a variety of domains, including medical imaging [10] and object recognition in natural images.

### C. Activation Functions

Activation functions play a critical role in the learning capability of neural networks, enabling them to model complex, nonlinear relationships between inputs and outputs. By introducing non-linearity into the network, activation functions allow each neuron to learn and represent different features across layers. Common activation functions include the Sigmoid, Tanh, and ReLU (Rectified Linear Unit).

The Sigmoid function maps input values into the (0,1) interval and was historically used in early neural networks. However, it suffers from vanishing gradient problems, making deep models difficult to train. The Tanh function improves upon this by mapping inputs to  $(-1,1)$ , but still faces similar limitations in deeper architectures.

The most widely used function in modern deep learning is ReLU, which outputs zero for negative inputs and a linear response for positive values. Its simplicity and computational efficiency make it suitable for deep architectures, and it mitigates the vanishing gradient problem, accelerating convergence. Variants such as Leaky ReLU and Parametric ReLU attempt to address ReLU's issue of "dead neurons" by allowing small gradients when the input is negative [5].

In classification tasks, particularly at the output layer, the Softmax function is commonly used to normalize raw output scores into probability distributions over multiple classes, facilitating decision-making and cross-entropy loss optimization. The choice and placement of activation functions have a direct impact on a model's ability to learn meaningful patterns and generalize to unseen data [16].

### D. Attention Mechanisms

Attention mechanisms represent a significant advancement in deep learning, enabling neural networks to focus selectively on the most informative parts of their input rather than processing all input data uniformly. Inspired by human visual attention, these mechanisms allow models to prioritize spatial regions or channels that are most relevant to the task at hand, such as identifying key objects in complex scenes. In CNNs, attention can significantly enhance the model's ability to capture critical features, especially in cluttered or detailed images [7]. Typically integrated after the feature extraction stages of a CNN, attention modules assign dynamic weights to different spatial or channel dimensions, emphasizing relevant patterns while suppressing irrelevant noise [25]. Although initially introduced in Natural Language Processing tasks, such as in Transformer architectures, attention has since been successfully adopted in visual recognition, resulting in improved accuracy and robustness in image classification [1].

## IV. IMPLEMENTATION

This section describes the overall implementation process, including the datasets used, the preprocessing applied to the data, the computational tools employed for model development, and the metrics adopted to evaluate performance.

### A. Datasets

In this study, two well-known benchmark datasets were used for image classification: CIFAR-10 and CIFAR-100. Both datasets were introduced by the Canadian Institute for Advanced Research (CIFAR) and are widely used to evaluate the performance of deep learning models in image recognition tasks [11].

CIFAR-10 consists of 60,000 color images of size  $32 \times 32$  pixels, divided into 10 distinct classes, such as airplanes, automobiles, birds, cats, and more. Each class contains 6,000 images, with 5,000 used for training and 1,000 for testing. The dataset is balanced, ensuring an equal distribution across classes.

CIFAR-100 contains the same number and size of images but introduces greater complexity, as it includes 100 classes grouped into 20 superclasses. Each class has 600 images—500 for training and 100 for testing—resulting in increased granularity and inter-class similarity, which makes it more challenging for classification tasks.

Both datasets are preprocessed and labeled, making them suitable for supervised learning. Their relatively low resolution allows for efficient model training and experimentation, which makes them particularly well-suited for evaluating the effect of attention mechanisms in both convolutional and feedforward neural networks.

### B. Data Preparation and Preprocessing

Before model training, all input images from the CIFAR-10 and CIFAR-100 datasets were preprocessed to ensure consistency and improve learning efficiency. The pixel values of the RGB images, initially in the integer range  $[0, 255]$ , were normalized to the continuous interval  $[0.0, 1.0]$  by dividing each value by 255 [13]. This normalization step enhances numerical stability and accelerates convergence during training.

The class labels, initially provided as integer values, were converted into one-hot encoded vectors to be compatible with the multi-class classification setting [6]. This encoding allows the model to calculate loss using categorical cross-entropy, which compares the softmax probabilities at the output layer to the one-hot target vectors.

To further improve generalization and reduce overfitting, data augmentation techniques were applied using the ImageDataGenerator utility from Keras [3]. During training, this tool performs real-time transformations such as random rotations, shifts, and horizontal flips, generating diverse image variations without increasing the dataset size. This dynamic augmentation enriches the model's exposure to various representations of the training data.

### C. Technology Stack

The experiments and model implementations in this study were conducted using the Python programming language, due to its simplicity, versatility, and widespread use in machine learning applications. The deep learning models were built and trained using the Keras high-level API with a TensorFlow backend [3]. Keras provides an intuitive interface for designing

and training neural networks, while TensorFlow ensures high-performance computations and GPU acceleration.

For data handling, preprocessing, and numerical operations, libraries such as NumPy and Pandas were utilized. Data visualization and evaluation of results were performed using Matplotlib and Seaborn. Additionally, Google Colaboratory was used as the execution environment [2], offering cloud-based access to GPUs and a seamless interface for interactive development and experimentation. The use of this stack enabled efficient training of multiple models and streamlined the experimentation process.

#### D. Evaluation Metrics

To comprehensively assess the performance of the classification models, several evaluation metrics were employed, including accuracy, precision, recall, and F1-score. While accuracy measures the overall correctness of predictions, it may not always reflect performance across all classes, especially in multi-class problems.

Therefore, macro-averaged versions of precision, recall, and F1-score were also calculated to provide a more balanced view of the model's behavior across all categories. These macro metrics compute the arithmetic mean of the corresponding scores for each class, regardless of class imbalance, and highlight how consistently the model performs across different categories [8].

This approach offers a deeper understanding of model robustness, identifying whether the classifier is biased toward specific classes or capable of generalizing across the whole dataset.

### V. EXPERIMENTAL EVALUATION

This section reports the experimental setup and results, detailing the evaluation procedure, performance comparisons, and analysis of the proposed and baseline models.

#### A. Experimental Setup

All experiments were conducted on Google Colaboratory, utilizing its GPU-enabled environment for efficient model training. The implementation was performed in Python using the TensorFlow-Keras deep learning framework. The datasets used were CIFAR-10 and CIFAR-100, and each model was trained for 50 epochs.

A batch size of 32 was used during training, along with the Adam optimizer and categorical cross-entropy as the loss function, which is suitable for multi-class classification. Three architectures were implemented and evaluated: a fully connected ANN, a standard CNN, and an Attention-enhanced CNN incorporating Squeeze-and-Excitation and CBAM blocks.

### VI. RESULTS AND COMPARISON

1) *Artificial Neural Network (ANN)*: The ANN model was evaluated as a baseline architecture, offering a non-convolutional approach to image classification. While relatively lightweight and fast to train, the ANN lacked spatial awareness due to its fully connected structure. This limitation hindered its ability to capture localized patterns and

visual structures, especially in more complex datasets such as CIFAR-100.

On CIFAR-10, the ANN achieved moderate performance, with an overall accuracy of 45.58% and an F1-score of 0.4369. However, on CIFAR-100, its performance dropped significantly to 43.15% accuracy and an F1-score of 0.54, highlighting its inability to generalize effectively across a larger number of classes.

The model also exhibited frequent misclassifications between visually similar categories. For instance, the confusion matrix revealed high error rates in pairs such as “cat” and “dog”, or “automobile” and “truck”, due to the absence of feature hierarchies and convolutional filters.

Despite its limitations, the ANN provided a valuable reference point for evaluating the impact of spatial feature learning and attention mechanisms introduced in the more advanced models.

2) *Convolutional Neural Network (CNN)*: The standard CNN architecture provided a significant improvement over the ANN by introducing convolutional layers capable of learning spatial hierarchies and localized features. This structural advantage enabled the model to achieve better generalization and classification accuracy across both datasets.

On CIFAR-10, the CNN achieved an accuracy of 73.17% and an F1-score of 0.7288, outperforming the ANN by approximately 28%. On CIFAR-100, the performance improved to 47.23% accuracy and an F1-score of 0.60, which confirmed the model's capacity to handle more complex visual patterns.

Despite these gains, the CNN still exhibited weaknesses in differentiating between visually similar or semantically close classes. The confusion matrix revealed persistent misclassifications between classes such as “cat” and “dog”, or “truck” and “automobile”, due to limitations in the model's ability to focus on the most informative regions of the input.

Moreover, while the training process showed stable convergence, the validation loss occasionally indicated minor overfitting, particularly in the later epochs. These observations motivated the incorporation of attention mechanisms in the enhanced CNN model, aiming to provide targeted focus and improved feature representation.

3) *Attention-enhanced CNN*: The Attention-enhanced CNN builds upon the standard CNN architecture by incorporating channel and spatial attention mechanisms, such as SE and CBAM blocks. These additions enable the model to dynamically highlight informative features while suppressing irrelevant ones, thus improving both local and global context understanding during training.

This architecture consistently delivered the best results across all metrics and datasets. On CIFAR-10, it reached 79.98% accuracy and an F1-score of 0.7957. On CIFAR-100, the improvements were even more pronounced, achieving 49.08% accuracy and an F1-score of 0.7957. The confusion matrices confirmed that the attention-enhanced model corrected many of the misclassifications seen in the plain CNN, especially in categories with high visual similarity, such as “cat” and “dog”, or “train” and “truck”.

The attention mechanism also contributed to smoother convergence curves, with less overfitting and lower validation loss compared to the baseline CNN. Although training required slightly more time (approximately 120 seconds more per average run), the increase in computational cost was offset by substantial improvements in classification robustness and generalization.

Overall, the integration of attention modules proved highly effective, offering state-of-the-art performance while preserving efficiency. This confirms their value as a scalable and practical enhancement for convolutional architectures in visual recognition tasks.

#### A. Model Complexity and Training Cost

While the Attention-enhanced CNN introduces additional layers and computations, the training time increased only moderately, from 654.23s (CNN) to 774.82s. This 18% overhead is justified by the observed performance gains, especially in complex classification scenarios like CIFAR-100. In practice, the added complexity remains within feasible limits for real-world applications using mid-range GPUs.

#### B. CIFAR-100 Evaluation

The results for the CIFAR-100 dataset are shown in Table I. The Attention-enhanced CNN significantly outperformed both the CNN and the ANN models in all metrics, achieving the highest accuracy (79.98%) and F1-score (0.7957), with a manageable increase in training time. These results demonstrate the effectiveness of incorporating attention mechanisms, especially in complex classification scenarios with a large number of classes.

TABLE I  
PERFORMANCE COMPARISON OF MODELS ON CIFAR-100

| Model | Accuracy      | Precision     | Recall        | F1-Score      | Training Time  |
|-------|---------------|---------------|---------------|---------------|----------------|
| ANN   | 43.15%        | 0.4618        | 0.4327        | 0.5400        | 377.24s        |
| CNN   | 47.23%        | 0.4900        | 0.4733        | 0.6000        | 654.23s        |
| CNN+A | <b>49.08%</b> | <b>0.5282</b> | <b>0.5098</b> | <b>0.7957</b> | <b>790.82s</b> |

As shown in Figure 1, the Attention-enhanced CNN significantly reduces misclassifications between visually similar classes, such as "deer" and "horse", or "cat" and "dog", which exhibited higher confusion in the CNN and ANN models. Furthermore, the ROC curve in Figure 2 demonstrates the robust discriminative capability of the model across multiple classes, with macro-average AUC values approaching 0.90.

#### C. CIFAR-10 Evaluation

A similar evaluation was conducted on the CIFAR-10 dataset. As shown in Table II, all models demonstrated higher performance compared to CIFAR-100, which can be attributed to the smaller number of classes and reduced classification complexity. The Attention-based CNN once again achieved the best results, with an accuracy of 80.29% and an F1-score of 0.8023, indicating consistent generalization benefits from attention mechanisms.

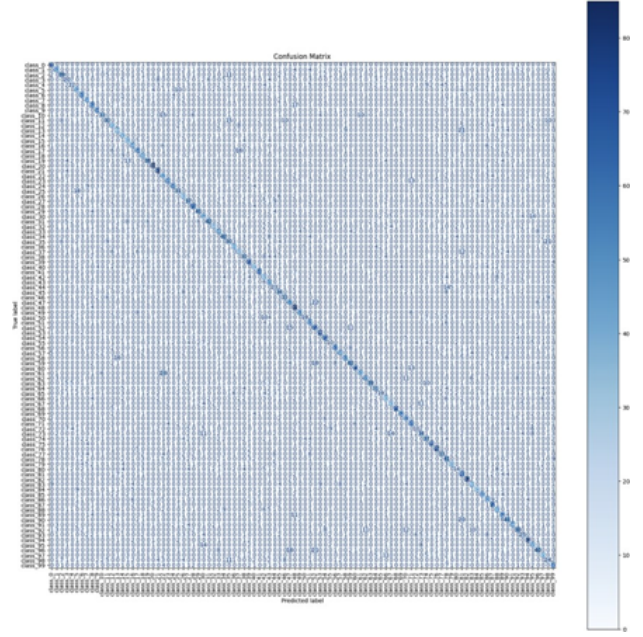


Fig. 1. Confusion Matrix – CNN with Attention on CIFAR-100

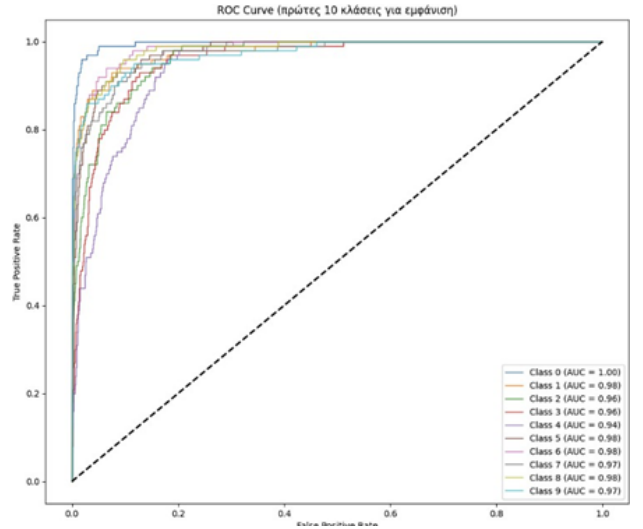


Fig. 2. ROC Curve – CNN with Attention on CIFAR-100

TABLE II  
PERFORMANCE COMPARISON OF MODELS ON CIFAR-10

| Model | Accuracy      | Precision     | Recall        | F1-Score      | Training Time  |
|-------|---------------|---------------|---------------|---------------|----------------|
| ANN   | 45.58%        | 0.4618        | 0.4427        | 0.4369        | 88.10s         |
| CNN   | 73.17%        | 0.7312        | 0.7317        | 0.7288        | 313.28s        |
| CNN+A | <b>79.98%</b> | <b>0.7982</b> | <b>0.7998</b> | <b>0.7957</b> | <b>774.82s</b> |

While the CIFAR-10 dataset posed a less challenging classification task compared to CIFAR-100 due to its reduced number of classes and lower inter-class similarity, the inclusion of attention mechanisms still yielded measurable performance gains. The Attention-enhanced CNN reduced common misclassifications between visually similar categories, most notably the "cat"–"dog" pair, where baseline CNN and ANN



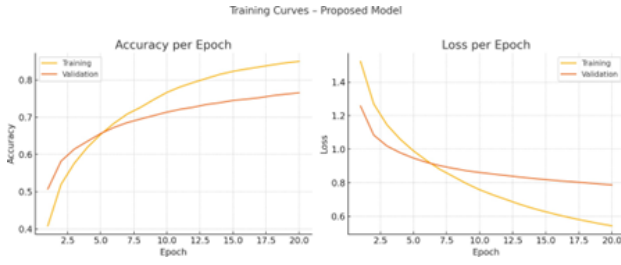


Fig. 3. Training and Validation Loss – CNN with Attention on CIFAR-10

models frequently confused features such as fur texture or body outline. Similar improvements were observed in differentiating “truck” from “automobile,” where attention modules helped focus on distinctive shape cues and background contexts.

## VII. CONCLUSIONS AND FUTURE WORK

This study demonstrated that the integration of attention mechanisms into convolutional neural networks improves performance in image classification tasks. Through comparative experiments on two benchmark datasets—CIFAR-10 and CIFAR-100—it was shown that attention-enhanced CNNs consistently outperform standard CNN and ANN architectures across all evaluation metrics, including accuracy, F1-score, and generalization capability.

Particularly in the more challenging CIFAR-100 dataset, the attention model achieved a notable performance gain, highlighting its ability to capture essential features better and manage inter-class similarities. These improvements stem from the ability of attention modules to emphasize salient image regions while suppressing irrelevant information, thereby enabling more focused and discriminative learning. Despite a moderate increase in training time, the performance benefits justify the added complexity.

Overall, attention-enhanced CNNs offer a scalable and effective solution to image classification problems, improving both predictive accuracy and robustness, especially in scenarios with high visual complexity.

Future work may explore the use of more advanced attention strategies such as self-attention or Transformer-based modules, as well as the evaluation of these models on higher-resolution and domain-specific datasets. Optimizations for real-time deployment, including lightweight attention blocks and training efficiency improvements, also represent promising directions. Future work could also explore hybrid models, inspired by successful approaches in NLP [17].

## REFERENCES

- [1] I. Bello, B. Zoph, A. Vaswani, J. Shlens, and Q. V. Le. Attention augmented convolutional networks. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3286–3295, 2019.
- [2] E. Bisong. Google colab. In *Building machine learning and deep learning models on google cloud platform: a comprehensive guide for beginners*, pages 59–64. Springer, 2019.
- [3] F. Chollet. *Deep learning with Python*. simon and schuster, 2021.
- [4] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *2005 IEEE computer society conference on computer vision and pattern recognition (CVPR’05)*, volume 1, pages 886–893, 2005.
- [5] X. Glorot, A. Bordes, and Y. Bengio. Deep sparse rectifier neural networks. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 315–323. JMLR Workshop and Conference Proceedings, 2011.
- [6] G. E. Hinton and R. R. Salakhutdinov. Reducing the dimensionality of data with neural networks. *science*, 313(5786):504–507, 2006.
- [7] J. Hu, L. Shen, and G. Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7132–7141, 2018.
- [8] G. James, D. Witten, T. Hastie, and R. Tibshirani. *An introduction to statistical learning: with applications in R*, volume 103. 2013.
- [9] S. Jetley, N. A. Lord, N. Lee, and P. H. Torr. Learn to pay attention. *arXiv preprint arXiv:1804.02391*, 2018.
- [10] A. Kanavos, O. Papadimitriou, G. Vonitsanos, M. Maragoudakis, and P. Mylonas. Enhanced brain tumor classification with convolutional neural networks.
- [11] A. Krizhevsky, G. Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [12] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012.
- [13] Y. LeCun, Y. Bengio, and G. Hinton. Deep learning. *nature*, 521(7553):436–444, 2015.
- [14] X. Li, W. Wang, X. Hu, and J. Yang. Selective kernel networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 510–519, 2019.
- [15] F. Lyu, F. Hu, V. S. Sheng, Z. Wu, Q. Fu, and B. Fu. Coarse to fine: Multi-label image classification with global/local attention. In *2018 IEEE International Smart Cities Conference (ISC2)*, pages 1–7, 2018.
- [16] C. Nwankpa. Activation functions: Comparison of trends in practice and research for deep learning. *arXiv preprint arXiv:1811.03378*, 2018.
- [17] O. Papadimitriou, A. Kanavos, G. Vonitsanos, M. Maragoudakis, and P. Mylonas. Advancing sentiment analysis of imdb movie reviews with a hybrid multinomial naive bayes and lstm approach. In *Novel & Intelligent Digital Systems Conferences*, pages 276–285, 2024.
- [18] J. Park, S. Woo, J.-Y. Lee, and I. S. Kweon. Bam: Bottleneck attention module. *arXiv preprint arXiv:1807.06514*, 2018.
- [19] A. G. Roy, N. Navab, and C. Wachinger. Concurrent spatial and channel ‘squeeze & excitation’ in fully convolutional networks. In *International conference on medical image computing and computer-assisted intervention*, pages 421–429, 2018.
- [20] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [21] G. Vonitsanos, T. Panagiotakopoulos, and A. Kameas. Comparative analysis of time series and machine learning models for air quality prediction utilizing iot data. In *IFIP International Conference on Artificial Intelligence Applications and Innovations*, pages 221–235, 2024.
- [22] F. Wang, M. Jiang, C. Qian, S. Yang, C. Li, H. Zhang, X. Wang, and X. Tang. Residual attention network for image classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3156–3164, 2017.
- [23] Q. Wang, B. Wu, P. Zhu, P. Li, W. Zuo, and Q. Hu. Eca-net: Efficient channel attention for deep convolutional neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11534–11542, 2020.
- [24] X. Wang, R. Girshick, A. Gupta, and K. He. Non-local neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7794–7803, 2018.
- [25] S. Woo, J. Park, J.-Y. Lee, and I. S. Kweon. Cbam: Convolutional block attention module. In *Proceedings of the European conference on computer vision (ECCV)*, pages 3–19, 2018.
- [26] Q.-L. Zhang and Y.-B. Yang. Sa-net: Shuffle attention for deep convolutional neural networks. In *ICASSP 2021-2021 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 2235–2239, 2021.
- [27] H. Zhao, J. Jia, and V. Koltun. Exploring self-attention for image recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10076–10085, 2020.