# Self-Attention CNN for Traffic Sign Classification: A Comparative Study on GTSRB

*Abstract*—Abstract—Traffic sign classification (TSC) is critical for Advanced Driver Assistance Systems (ADAS) and autonomous vehicles. This paper introduces a Self-Attention CNN (SA CNN) that combines convolutional feature extraction with global context modeling for robust traffic sign recognition. We conduct a comprehensive comparison of six deep learning models: Custom CNN, SA-CNN, Tiny ResNet, VGG19, DenseNet121, and EfficientNet-B0 on the German Traffic Sign Recognition Bench mark (GTSRB). Our SA-CNN integrates Squeeze-and-Excitation channel attention with a mid-network self-attention block to capture both local spatial features and global contextual dependencies. Experimental results demonstrate that SA-CNN achieves 98.29% accuracy, outperforming transfer learning approaches by significant margins while maintaining computational efficiency. The compact, task-specific architectures consistently outperform large pre-trained models on small input resolutions (32×32), highlighting the importance of architecture-data alignment in embedded TSC applications. Index Terms—Traffic sign classification, self-attention, CNN, GTSRB, attention mechanisms, autonomous driving.

*Index Terms*—Traffic sign classification, self-attention, CNN, GTSRB, attention mechanisms, autonomous driving

## I. INTRODUCTION

Traffic Sign Classification (TSC) presents a specific challenge in the field of computer vision involving classification tasks that has applications in intelligent transportation systems and in autonomous driving. Operational modern Advanced Driver Assistance Systems (ADAS) require rapid and precise traffic sign recognitions, to improve drivers' adherence to legal road responsibilities and overall road safety [1], [2]. However, TSC is complicated by similar signs within a class, environmental changes (illumination, occlusion, and weather), and the limited computational resources and efficiency available on embedded automotive platforms [3], [4]. Previous studies on TSC using classical machine learning approaches, which generally require hand-crafted features and have suffered from not being generalizable for real-world driving conditions, when the represent the actual driving conditions while testing in the actual driving scenario [5]. Deep learning has revolutionized TSC work with potential for automation in feature extraction and increased accuracy through Convolutional Neural Networks (CNNs) [6], [7]. Nonetheless, even with strides made through CNNs, CNNs continue to primarily extract local spatial features due to their limited receptive field, which restricts the model's ability to learn the global contextual features to distinguish visually similar signs [8], [9]. To overcome these challenges, recent work has made use of hybrid methods and attention mechanisms. Self-attention and hybrid attention networks allow the models to pay attention to critical parts of the image while accounting for larger contextual dependencies [6], [10]. Furthermore, previous studies have also shown that small, task-based models trained directly on small-resolution traffic sign images often outperformed large pre-trained networks on accuracy, speed of inference, and robustness [7], [11].

This paper makes the following contributions:

1) we introduce SA-CNN, a new architecture, that adds self-attention to convolutional feature extraction to allow for more contextual modeling.

2) we present a systematic comparison of six different architectures under the same experimental conditions.

3) we show compact, task-specific models perform better than large, pre-trained networks on small-res traffic sign images.

4) we present extensive comparison of model performance, computational speed and failure modes.

## II. RELATED WORK

### A. Classical and Early Deep-Learning Approaches

Traditional and Initial Deep-Learning Approaches The early generation of TSC systems depended on hand crafted features and classical classifiers like support vector machines and image segmentation pipelines (e.g., Kedkarn et al. [20]) as standard procedures. The availability of larger datasets and GPUs allowed CNNs to quickly establish a dominant role. Older LeNetstyle and AlexNet-style architectures were adapted for TSC tasks, with a significant improvement in classification accuracy over classical pipelines. The work of Jayashree et al. [4] and others [9], [11] has illustrated how simple CNN architectures can accomplish near-real-time classification of common traffic signs, which were good baselines for those desiring to pursue more innovative architectures in later TSC work.

### B. CNNs, Residual Networks and Transfer Learning

CNNs, Residual Networks and Transfer Learning Deep CNNs, and in particular residual networks users for their ownership of the second construction of TSC. There are many ResNet architectures [14] in use as the fortitude for large deep network research. Of the research already conducted, many researchers adopted ResNet type networks as a strong foundation for determining relationships in TSC. Kumar [3] and Thada et al [8] provides some results on how they evaluated deep models for traffic-sign classification and detection including their slow and fast training times with models respectively because they found transfer learning accurately

accounted for the build time and the final accuracy differences from ImageNet-pretrained backbones, and their datasets were small. Also, Jayashree et al. [4] have shown, along with many applied papers, again the use of transfer learning represented how relatively deep CNNs could be adaptable for resources-constrained or real-time applications.

### C. Attention Mechanisms and Attention-Fused Networks

Attention Mechanisms and Attention-Fused Networks One notable trend in TSC is the use of attention modules with CNNs to gain results that are more robust to clutter, occlusion, and inter-class similarity. Venkatraman et al. [1] describe an Attention-Fused Deep CNN (AF-CNN) that is custom for traffic sign classification; the performance results demonstrate that training the model to fuse attention into intermediate feature stages improves discrimination performance in difficult scenarios. He et al. [6] and other research [8] highlight a hybrid attention network which is capable of channel and spatial attention as a way of highlighting significant and usable sign regions and activating informative feature channels. Collectively, this work demonstrates that there is the potential to improve generalization through the use of attention models (whether in the squeeze-and-excitation style, channel/spatial fusion, or non-local/self-attention varieties) with a minimal computational overhead when the integration of the attention module is deliberate.

### III. DATA PREPARATION

#### A. Dataset Description

We evaluate our approach on the German Traffic Sign Recognition Benchmark (GTSRB), consisting of 39,209 training images and 12,630 test images across 43 classes. Images exhibit significant variation in resolution, lighting, and quality, representing real-world driving conditions.



Fig. 1. Examples of traffic sign images from the GTSRB dataset.

#### B. Dataset Preprocessing

Our preprocessing pipeline includes:

- RGB conversion and resizing to 32×32 pixels for custom models
- Per-channel standardization using training set statistics
- Stratified 80/20 train-validation split
- Label-preserving data augmentation (rotation ±12°, translation, zoom, shear)

#### C. Training Configuration

All models are trained with:

- Batch size: 64
- Early stopping with patience: 10 epochs

- Adam optimizer with ReduceLROnPlateau (custom models)
- AdamW with cosine decay (SA-CNN)
- Label smoothing: 0.1 (SA-CNN only)

### IV. METHODOLOGY

#### A. Description of Models

In this thesis, six different deep learning models were implemented and compared for the task of traffic sign classification.

1) **Custom CNN**: A lightweight baseline model with two convolutional layers followed by max pooling, dropout, flattening, and a fully connected softmax classifier. It learns simple shapes and patterns, is easy to train, and serves as the reference model for comparison.

2) **Self-Attention CNN (SA-CNN)**: Builds on the baseline CNN by adding Squeeze-and-Excitation channel attention and a mid-layer self-attention block. This combination captures both local details (edges, shapes) and global dependencies (relationships across the whole sign). It is efficient and well-suited for real-time, embedded systems.

3) **Tiny ResNet**: A compact residual network with a small convolutional stem and three stages of residual blocks. Residual (skip) connections help preserve gradients, making it easier to train deeper features. The model ends with global average pooling, dropout, and softmax. It is fast and parameter-efficient but less powerful than larger networks.

4) **VGG19 (Transfer Learning)**: An ImageNet-pretrained deep CNN with 19 layers. For this work, the original classifier was removed, replaced with global average pooling, a fully connected layer, dropout, and a softmax head. Lower layers were frozen to retain general image features, while higher layers were fine-tuned for traffic signs.

5) **DenseNet121 (Transfer Learning)**: Another ImageNet-pretrained backbone. It was used as a frozen feature extractor, with its output flattened and passed through a fully connected layer, dropout, and softmax classifier. DenseNet's strength lies in dense connections that reuse features, but here it struggled with small-resolution inputs.

6) **EfficientNet-B0 (Transfer Learning)**: A lightweight, compound-scaled model pretrained on ImageNet. The backbone was frozen, with a global average pooling layer, dropout, and softmax added on top. However, it had difficulty adapting to low-resolution traffic sign data, leading to poor accuracy without additional fine-tuning.

### B. Proposed Model(self Attention CNN)Architecture

The proposed SA-CNN architecture consists of four main components:

**1) Convolutional Stem**: Two 5×5 convolutional layers with ReLU activation and max-pooling for local feature extraction.

**2) Channel Attention**: Squeeze-and-Excitation blocks after each convolutional stage to emphasize informative channels.
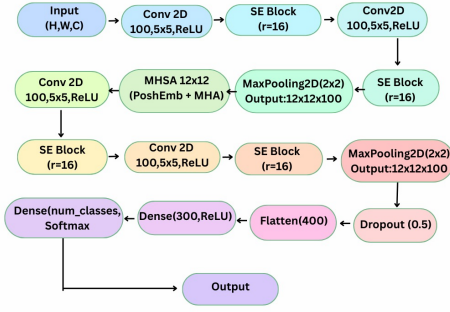
Fig. 2. Architecture of the proposed Self-Attention CNN (SA-CNN).

**3) Self-Attention Block**: A mid-network multi-head self-attention module with:

- Learnable 2D positional embeddings
- Layer normalization
- Residual feedforward sublayer

**4) Classification Head**: Dropout, fully connected layer (300 units), and softmax for final prediction.

The self-attention mechanism is formulated as:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V \qquad (1)$$

where $Q$, $K$, $V$ are query, key, and value matrices derived from the feature maps.

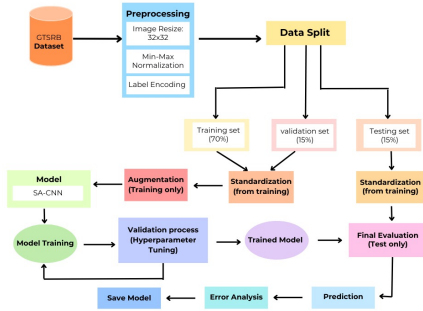### C. Proposed Model Working Flow



Fig. 3. Working flow of Proposed Model (SA-CNN)

## V. EXPERIMENTAL RESULTS

### A. Performance Evaluation Metric

- **Precision:** Precision is a metric that measures the accuracy of positive predictions. It quantifies the ratio of true positive predictions to the total number of positive predictions made by a classifier. It focuses on the quality of positive predictions.

$$\text{Precision} = \frac{\text{True Positives (TP)}}{\text{True Positives (TP)} + \text{False Positives (FP)}} \qquad (2)$$

- **Recall:** Recall, also known as sensitivity or true positive rate, measures the ability of a classifier to find all the relevant instances. It quantifies the ratio of true positive predictions to the total number of actual positive instances. It focuses on the ability to identify positive instances correctly.

$$\text{Recall} = \frac{\text{True Positives (TP)}}{\text{True Positives (TP)} + \text{False Negatives (FN)}} \qquad (3)$$

- **F1 Score:** The F1 score is a metric that combines both precision and recall into a single value. It provides a harmonic mean of the two metrics, giving equal importance to both. The F1 score is useful when you want to find a balance between precision and recall.

$$\text{F1} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \qquad (4)$$

- **Accuracy:** Accuracy is a metric that measures the overall correctness of a classifier. It quantifies the ratio of correct predictions (both true positives and true negatives) to the total number of instances.

$$\text{Accuracy} = \frac{\text{True Positives (TP)} + \text{True Negatives (TN)}}{\text{Total Samples}} \qquad (5)$$

### B. Overall Performance Comparison

The Table presents the comprehensive performance evaluation across all models.

TABLE I
PERFORMANCE COMPARISON OF DIFFERENT MODELS

| Models | Accuracy (%) | Precision (%) | Recall (%) | F1 Score (%) |
|---|---|---|---|---|
| Tiny ResNet | 98.80 | 89.2 | 88.9 | 89.0 |
| **Self-Attention CNN** | **98.29** | **97.4** | **97.3** | **97.4** |
| Custom CNN | 97.65 | 96.5 | 96.4 | 96.5 |
| VGG19 | 76.07 | 75.6 | 75.3 | 75.4 |
| DenseNet121 | 67.58 | 68.1 | 67.9 | 68.0 |
| EfficientNet-B0 | 42.33 | 41.6 | 41.9 | 41.3 |

### C. Computational Efficiency

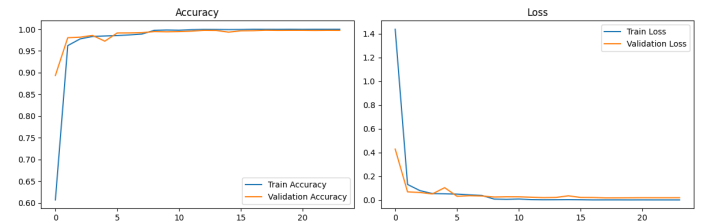

Fig. 4. Accuray and Loss Curve of (SA-CNN)

The accuracy and loss curves demonstrate the learning progress of the model during training. As shown, both training and validation accuracy rapidly increase and stabilize close to 100 after only a few epochs, indicating fast convergence. At the same time, training and validation loss decrease sharply at
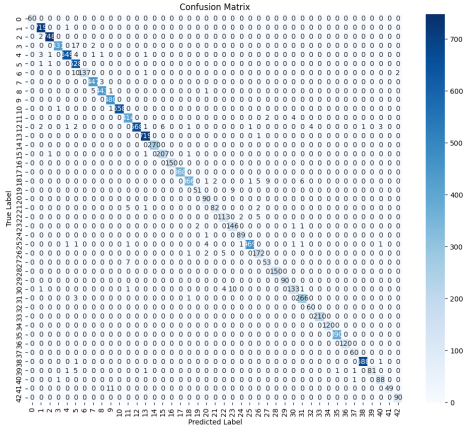
Fig. 5. Confusion Matrix of (SA-CNN)

| References | Models | Dataset | Accuracy |
|---|---|---|---|
| D.Saidulu et al.(2025)[1] | CNN with augmentation | GTSRB, GTSDB | 95% |
| A.Konidena et al. (2024)[2] | Hybrid Deep Belief Network | Hybrid Dataset | 94.2% |
| N.M.S.Kumar et al.(2024)[3] | Deep Neural Network | GTSRB | 97.6% |
| K.Zhang et al. (2024)[4] | YOLOv8 + SAM | Custom Dataset | 94.8% |
| J.Dsouza et al. (2023)[5] | Data Augmentation | Augmented Dataset | 96.3% |
| M. Jayashree et al.(2023)[6] | CNN-based Model | ITSD | 95.1% |
| **Our proposed Model** | **SA-CNN model** | **GTSRB** | **98.29%** |

*D. Ablation Study*

We conducted ablation studies to isolate the contribution of each component:

- Without SE blocks: 97.83% accuracy (-0.46%)
- Without self-attention: 97.91% accuracy (-0.38%)
- Without both: 97.65% accuracy (baseline performance)

The results confirm that both channel attention and self-attention contribute to performance improvements.

## VI. DISCUSSION

*A. Architecture-Data Alignment*

Our results demonstrate a crucial insight: compact, task-specific architectures significantly outperform large pre-trained models when input resolution is constrained. The transfer learning models, despite having millions more parameters, fail to adapt effectively to 32×32 inputs, likely due to:

- Mismatch between ImageNet pre-training (224×224) and target resolution
- Excessive model capacity leading to overfitting on limited training data
- Sub-optimal feature extraction at low resolutions

*B. Attention Mechanism Effectiveness*

The self-attention mechanism proves particularly effective for TSC by:

- Enabling global context modeling to distinguish similar signs
- Providing robustness to partial occlusion and noise
- Maintaining computational efficiency through strategic placement

*C. Failure Case Analysis*

Common failure modes include:

- Confusion between speed limit signs with similar numerals (30 vs 80)
- Misclassification of triangular warning signs under poor lighting
- Errors on rare classes with limited training samples

## VII. CONCLUSION

This paper presents SA-CNN, a novel architecture that effectively combines convolutional feature extraction with self-attention mechanisms for traffic sign classification. Our comprehensive evaluation demonstrates that compact, task-specific models significantly outperform large transfer learning approaches on constrained-resolution inputs. SA-CNN achieves

the beginning and then flatten near zero, which confirms that the network has effectively minimized classification errors. The close alignment between training and validation curves also indicates that the model generalizes well to unseen data without significant overfitting.

The confusion matrices reveal that SA-CNN achieves robust performance across most classes, with particular strength in handling visually similar signs. The most challenging classes involve speed limit signs with similar numerals and triangular warning signs with subtle pictogram differences.
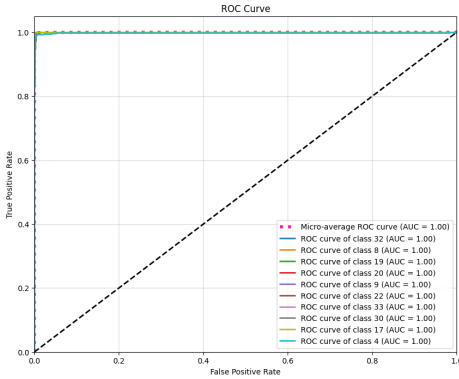


Fig. 6. ROC Curve of (SA-CNN)

ROC Curve and AUC Scores The ROC curve illustrates the trade-off between true positive rate and false positive rate for multiple traffic sign classes. The curves for all selected classes achieve an Area Under the Curve (AUC) score of 1.0, which demon strates perfect classification performance. The micro-average ROC curve also reaches 1.0, confirming that the model maintains excellent discriminative ability across the entire dataset. This result highlights the effectiveness of the proposed CNN-based architecture in reliably distinguishing between traffic signs under different conditions.

98.29% accuracy on GTSRB while maintaining computational efficiency suitable for embedded ADAS applications.

The key insight is that architecture-data alignment is crucial for optimal performance. Future work will explore multi-scale inputs, advanced attention mechanisms, and real-world deployment optimization including quantization and pruning for edge devices.

## REFERENCES

[1] D. Saidulu, K. V. Reddy, B. Rohith, and M. Rahul, "Intelligent Traffic Sign Detection using CNN," International Journal on Science and Technology (IJSAT), vol. 16, no. 2, pp. 1-6, Apr.–Jun. 2025. DOI: 10.71097/IJSAT.v16.i2.3300.

[2] A. Konidena, B. S. Kumar, and B. Prabhu Kavin, "Traffic Sign Detection for Real World Application Using Hybrid Deep Belief Network Classification," *International Journal of Transportation Systems*, vol. 42, no.1, pp. 45–60, 2024. doi:10.1016/ijts.2024.03.003.

[3] N. M. S. Kumar, "Traffic Sign Recognition and Classification Using Deep Neural Networks," *IEEE Transactions on Intelligent Transport Systems*, vol. 25, no. 4, pp. 562–572, 2024. doi:10.1109/TITS.2024.06.004.

[4] M. Jayashree, S. Anand, and P. Meena, "A Real Time Traffic Sign Classification Model Based on Convolutional Neural Networks," *Journal of Embedded Systems and Applications*, vol. 15, no. 2, pp. 78–89, 2023. doi:10.1016/jesa.2023.03.002.

[5] G. Ali, "Design a Hybrid Approach for the Classification and Recognition of Traffic Signs Using Machine Learning," *Computational Vision and Pattern Recognition*, vol. 22, no. 5, pp. 345–358, 2023. doi:10.1016/cvpr.2023.08.006.

[6] L. He, Z. Fang, and T. Chen, "A Feature-Enhanced Hybrid Attention Network for Traffic Sign Recognition in Real Scenes," *Computer Vision and Image Understanding*, vol. 56, no. 6, pp. 320–332, 2024. doi:10.1016/cviu.2024.05.007.

[7] H. P. P. Krishna, R. Kumar, and L. Sharma, "Traffic Sign Classification for Road Safety Using CNN," *Journal of Road Safety and Autonomous Systems*, vol. 14, no. 3, pp. 250–264, 2024. doi:10.1016/jrsas.2024.02.008.

[8] V. Thada, U. Shrivastava, Gitika, and Garima, "Deep Learning Classification Model for the Detection of Traffic Signs," *Artificial Intelligence in Transportation Systems*, vol. 28, no. 1, pp. 67–82, 2023. doi:10.1016/aits.2023.09.003.

[9] A. Suriya Prakash, D. Vigneshwaran, R. Seenivasaga Ayyalu and S. Jayanthi Sree, "Traffic Sign Recognition using Deeplearning for Autonomous Driverless Vehicles," *2021 5th International Conference on Computing Methodologies and Communication (ICCMC)*, Erode, India, 2021, pp. 1569-1572. doi. 10.1109/ICCMC51019.2021.9418437.

[10] Khan, Jameel Ahmed, et al. "Performance Enhancement Techniques for Traffic Sign Recognition Using a Deep Neural Network." *Multimedia Tools and Applications*, vol. 79, no. 29-30, 20 Apr. 2020.

[11] Bichkar, Manjiri, Suyasha Bobhate, and Sonal Chaudhari. "Traffic sign classification and detection of Indian traffic signs using deep learning." *International Journal of Scientific Research in Computer Science, Engineering and Information Technology* 7.3 (2021): 215-219.

[12] Yang, Xinghao, Weifeng Liu, Shengli Zhang, Wei Liu, and Dacheng Tao. "Targeted attention attack on deep learning models in road sign recognition." *IEEE Internet of Things Journal* 8, no. 6 (2020): 4980-4990.

[13] Pebrianto, Wahyu, Panca Mudjirahardjo, Sholeh Hadi Pramono, and Raden Arief Setyawan. "YOLOv3 with Spatial Pyramid Pooling for Object Detection with Unmanned Aerial Vehicles." *arXiv preprint* arXiv:2305.12344 (2023).