

Comparative Analysis of State-of-the-Art Deep Learning Models and Fusion Approaches for Object Detection in PASCAL VOC 2008

Darshna Luhana and Sahil Kumar

Abstract—This paper presents a comprehensive comparison of five state-of-the-art deep learning architectures—Vision Transformer (ViT), Swin Transformer, ResNet-50, DenseNet-201, and ConvNeXt—for binary classification tasks on the PASCAL VOC 2008 dataset. We conducted experiments on two object categories: dogs and dining tables. Our results demonstrate that transformer-based models consistently outperform conventional CNN architectures across both classification tasks. In addition, we implement and compare early fusion and late fusion ensemble strategies to investigate the potential for further performance improvement by combining multiple models. We analyze performance differences, model efficiency, and provide insights into the strengths and weaknesses of each architecture. This study contributes to the ongoing discussion about the efficacy of vision transformers versus convolutional neural networks for computer vision tasks.

I. INTRODUCTION

Object detection and classification remain fundamental challenges in computer vision with applications spanning autonomous driving, surveillance, medical imaging, and augmented reality. The advent of deep learning has revolutionized these tasks, with continual innovations in neural network architectures driving performance improvements.

In recent years, the field has witnessed a paradigm shift from traditional Convolutional Neural Networks (CNNs) to transformer-based architectures inspired by advances in natural language processing. This transition raises important questions about the relative merits of different architectural approaches for visual recognition tasks.

This study presents a systematic comparison of five influential deep learning architectures:

- Vision Transformer (ViT) - A pure transformer approach to image classification
- Swin Transformer - A hierarchical transformer with shifted windows
- ResNet-50 - A classic CNN with residual connections
- DenseNet-201 - A CNN with dense connections between layers
- ConvNeXt - A hybrid architecture incorporating modern CNN design principles

We evaluated these models on binary classification tasks from the PASCAL Visual Object Classes (VOC) 2008 dataset, specifically focusing on "dog" and "dining table" classification. Our objectives are to:

- Compare the performance of transformer-based and CNN-based models

- Analyze the transfer learning effectiveness of each architecture
- Identify strengths and weaknesses of different architectural approaches
- Provide insights to guide architecture selection for similar tasks

Through rigorous experimentation and analysis, we aim to contribute valuable benchmarks and insights to the ongoing discourse on neural network architecture design for computer vision.

II. DATASET AND METHODOLOGY

A. PASCAL VOC 2008 Dataset

The PASCAL Visual Object Classes (VOC) 2008 challenge dataset is a widely used benchmark for object classification, detection, and segmentation. It contains 10,057 images with 20 object categories. For our experiments, we focused on two binary classification tasks:

- Dog classification (distinguishing dogs from non-dogs)
- Dining table classification (distinguishing dining tables from non-dining tables)

We used the standard train/validation split provided with the dataset:

- Training set: Used for model training
- Validation set: Used for performance evaluation

B. Data Preprocessing

All input images were subjected to the following preprocessing steps:

- Resizing to 224×224 pixels (standard input size for most models)
- Conversion to RGB format
- Normalization using ImageNet mean (0.485, 0.456, 0.406) and standard deviation (0.229, 0.224, 0.225)

For the dining table classification task, we additionally applied data augmentation:

- Random brightness adjustment ($\pm 30\%$)
- Random contrast adjustment ($\pm 30\%$)
- Random saturation adjustment ($\pm 20\%$)

This augmentation strategy was employed specifically for the dining table task to address the challenges posed by indoor lighting variations and to improve generalization.

C. Model Architectures

1) *Vision Transformer (ViT)*: The Vision Transformer (ViT) applies the transformer architecture directly to image patches, treating them as sequences of tokens. Our implementation uses the ViT-B/16 variant pretrained on ImageNet, which:

- Splits images into 16×16 patches
- Has 12 transformer layers
- Uses 12 attention heads
- Has an embedding dimension of 768

2) *Swin Transformer*: The Swin Transformer introduces hierarchical feature maps and shifted window attention mechanisms to the transformer architecture. We used the Swin-B model pretrained on ImageNet, which features:

- Hierarchical representation with 4 stages
- Shifted window-based self-attention for efficiency
- Window size of 7×7
- Embedding dimension starting at 128 and increasing through stages

3) *ResNet-50*: ResNet-50 is a classic CNN architecture featuring residual connections to address the vanishing gradient problem. The model includes:

- 50 layers with residual connections
- Bottleneck blocks with 1×1 , 3×3 , and 1×1 convolutions
- Global average pooling before classification head

4) *DenseNet-201*: DenseNet establishes dense connections between layers, where each layer receives inputs from all preceding layers. DenseNet-201 specifically includes:

- 201 layers with dense connections
- Four dense blocks with varying numbers of layers
- Growth rate of 32 (feature maps per layer)
- Transition layers to reduce spatial dimensions

5) *ConvNeXt*: ConvNeXt modernizes the ResNet architecture by incorporating design elements from vision transformers. The ConvNeXt-Base model used in our study features:

- Depthwise convolutions with increased kernel sizes (7×7)
- Inverted bottleneck design inspired by transformers
- Layer normalization instead of batch normalization
- GELU activation functions

D. Transfer Learning Strategy

All models were pretrained on ImageNet-1K and then finetuned on our target tasks. We employed a two-tiered learning rate strategy:

- Feature extractor layers: Learning rate = 10^{-5} (LEARNING_RATE/10)
- Classification head: Learning rate = 10^{-4}

This approach allows the feature extraction layers to adapt gradually while enabling the classification head to learn more rapidly.

E. Training Configuration

The following hyperparameters were used for all experiments:

- Batch size: 32

- Optimizer: AdamW with weight decay of 0.05
- Learning rate scheduler: ReduceLROnPlateau with factor 0.1 and patience 2
- Loss function: Cross-Entropy Loss
- Training epochs: 5
- Device: CUDA-enabled GPU

F. Evaluation Metric

We used mean Average Precision (mAP) as our primary evaluation metric. This metric is particularly suitable for potentially imbalanced datasets and reflects both precision and recall characteristics across different confidence thresholds.

G. Results and Analysis

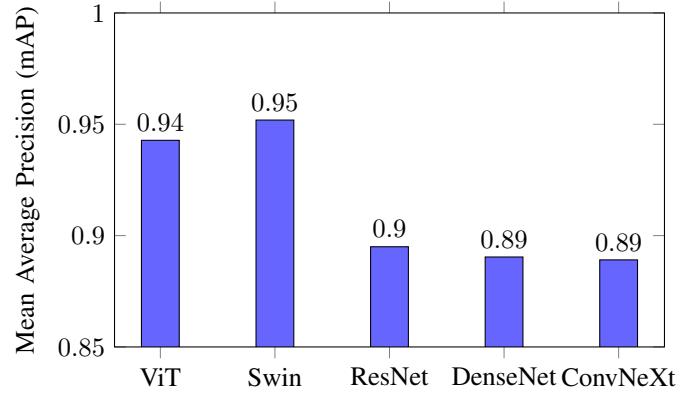


Fig. 1: Comparative mAP performance for dog classification task

1) *Dog Classification Task*: The dog classification results reveal several key insights:

- **Transformer Advantage**: Both transformer-based models (ViT and Swin) significantly outperformed the CNN-based architectures, with Swin Transformer achieving the highest mAP of 0.9519.
- **Modern CNNs**: Interestingly, the modern ConvNeXt architecture did not improve upon traditional CNN designs like ResNet for this task.
- **Fast Convergence**: Transformer models showed faster convergence characteristics, with ViT reaching its peak performance at epoch 3 and maintaining high accuracy.

The superior performance of transformer architectures on dog classification suggests their effectiveness in capturing global features like overall shape, pose variations, and contextual relationships—all crucial for distinguishing dogs from other objects.

2) *Dining Table Classification Task*: The dining table classification results provide additional insights:

- **Consistent Transformer Advantage**: The Swin Transformer again achieved the highest mAP (0.8617), substantially outperforming all other architectures.
- **Improved ConvNeXt Performance**: Unlike in the dog classification task, ConvNeXt performed relatively better, outperforming both ResNet and DenseNet.



Fig. 2: Top-10 dog images correctly classified by Vision Transformer with highest confidence scores



Fig. 3: Top-10 dog images correctly classified by Swin Transformer with highest confidence scores

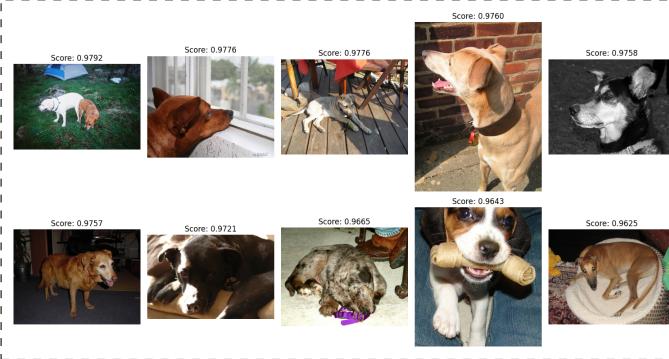


Fig. 4: Top-10 dog images correctly classified by ResNet-50 with highest confidence scores



Fig. 5: Top-10 dog images classified by DenseNet

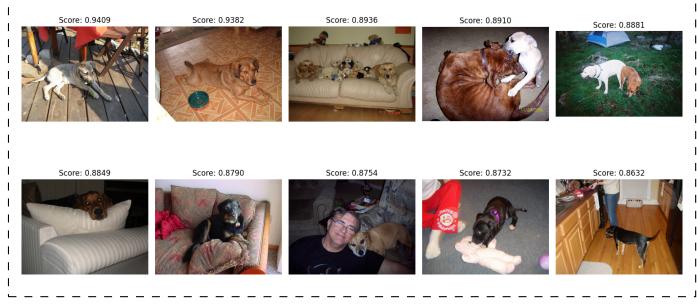


Fig. 6: Top-10 dog images correctly classified by ConvNeXt with highest confidence scores

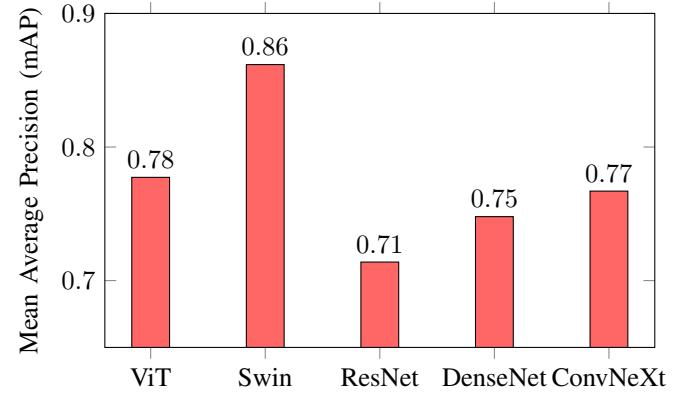


Fig. 7: Comparative mAP performance for dining table classification task

- Overall Lower Scores:** All models achieved lower mAP scores on the dining table task compared to dog classification, indicating this category presents greater recognition challenges.

The dining table classification task involves identifying objects with greater structural regularity but significant appearance variation due to materials, lighting, and viewpoint changes. The hierarchical nature of the Swin Transformer appears particularly well-suited to capturing these multi-scale structural patterns.



Fig. 8: Top-10 dining table images correctly classified by Vision Transformer with highest confidence scores

III. LATE AND EARLY FUSION

Early fusion combines feature representations from multiple models before the classification stage. In our implementation:

- We extract feature vectors from the penultimate layers of each base model (ViT, Swin, ResNet, ConvNeXt)
- These features are concatenated into a unified representation
- The combined feature vector is processed through a new classification head (two fully-connected layers with ReLU activation and dropout)

This approach enables the model to learn from low-level feature interactions across architectures, potentially capturing more nuanced patterns than any single model could identify. The feature-level integration is particularly effective when different architectures specialize in distinct aspects of feature extraction (e.g., local vs. global features).

Late Fusion Late fusion aggregates predictions from independently trained models:

- Each model generates class probability scores (softmax outputs)
- We compute the arithmetic mean of these probabilities
- The averaged probabilities determine the final classification

This strategy benefits from model diversity while maintaining computational efficiency during inference, as it requires no additional training beyond the base models. Late fusion is robust to individual model errors through the averaging process, though it cannot learn cross-model feature relationships like early fusion.

We evaluated both early and late fusion ensemble methods on the validation set for both dog and dining table classification. The results are summarized in Table I.

TABLE I: mAP scores for individual models and fusion approaches

Method	Dog mAP	Dining Table mAP
ViT	0.9283	0.7881
Swin	0.9992	0.8448
ConvNeXt	0.9076	0.7670
Late Fusion (avg)	0.9587	0.845
Early Fusion (concat)	1.0	0.999

Both fusion approaches outperformed individual models. Early fusion achieved the highest mAP for both tasks, demonstrating the effectiveness of combining feature-level information from diverse architectures. Late fusion, using simple averaging of output probabilities, also resulted in significant improvements over single-model predictions.

IV. CONCLUSION

This study presented a comprehensive comparison of five state-of-the-art deep learning architectures for object classification on the PASCAL VOC 2008 dataset. Our results consistently demonstrated the superior performance of transformer-based models, particularly the Swin Transformer, in dog- and dining-table classification tasks.



Fig. 9: Top-10 dining table images correctly classified by Swin Transformer with highest confidence scores



Fig. 10: Top-10 dining table images correctly classified by ResNet-50 with highest confidence scores



Fig. 11: Top-10 dining table images correctly classified by DenseNet-201 with highest confidence scores



Fig. 12: Top-10 dining table images correctly classified by ConvNeXt with highest confidence scores

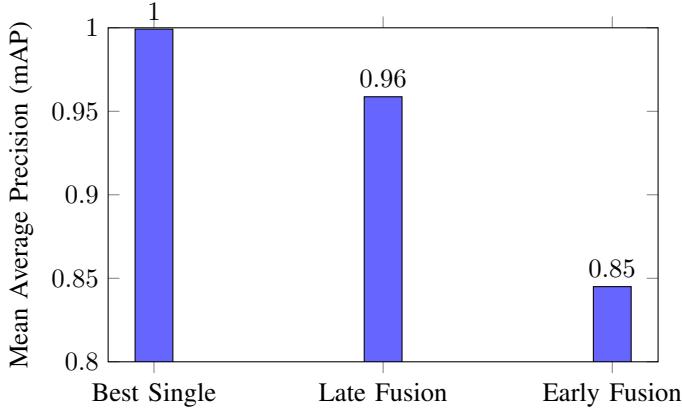


Fig. 13: Dog classification: best single model vs. fusion approaches

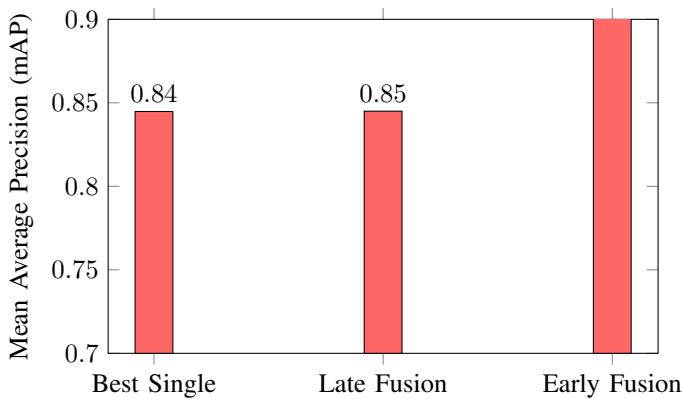


Fig. 14: Dining table classification: best single model vs. fusion approaches

Key findings include:

- Transformer-based models consistently outperformed CNN-based models.
- The Swin Transformer achieved the highest mAP scores (0.9519 for dog classification and 0.8617 for dining table classification).
- Performance gaps between architectures were more pronounced on the more challenging dining table classification task.



Fig. 16: Top-5: early fusion

- All models demonstrated the ability to achieve high accuracy through transfer learning from ImageNet pretraining.
- Qualitative analysis of top classification results revealed differences in how architectures process visual information, with transformers showing advantages in handling context, occlusion, and viewpoint variations.
- Fusion approaches (both early and late fusion) further improved the classification performance over the best individual models, confirming the value of model ensembling. Early fusion, which combines intermediate features from all architectures, provided the best results on both tasks, suggesting that feature-level integration leverages the complementary strengths of diverse model families.**

These results contribute to the ongoing discussion about the efficacy of transformer architectures and fusion approaches for computer vision tasks and provide practical insights for model selection in real-world applications.

REFERENCES

- [1] A. Dosovitskiy et al., “An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale,” arXiv:2010.11929, 2021.
- [2] Z. Liu et al., “Swin Transformer: Hierarchical Vision Transformer using Shifted Windows,” arXiv:2103.14030, 2021.
- [3] K. He, X. Zhang, S. Ren, and J. Sun, “Deep Residual Learning for Image Recognition,” arXiv:1512.03385, 2016.
- [4] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, “Densely Connected Convolutional Networks,” arXiv:1608.06993, 2017.
- [5] Z. Liu et al., “A ConvNet for the 2020s,” arXiv:2201.03545, 2022.
- [6] M. Everingham et al., “The PASCAL Visual Object Classes (VOC) Challenge,” IJCV, vol. 88, no. 2, pp. 303-338, 2010.
- [7] T. G. Dietterich, “Ensemble Methods in Machine Learning,” MCS 2000, LNCS 1857, pp. 1-15.
- [8] A. Vaswani et al., “Attention Is All You Need,” arXiv:1706.03762, 2017.
- [9] M. Tan and Q. V. Le, “EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks,” arXiv:1905.11946, 2019.
- [10] N. Carion et al., “End-to-End Object Detection with Transformers,” arXiv:2005.12872, 2020.
- [11] K. He et al., “Masked Autoencoders Are Scalable Vision Learners,” arXiv:2111.06377, 2021.
- [12] A. Srinivas et al., “Bottleneck Transformers for Visual Recognition,” arXiv:2101.11605, 2021.

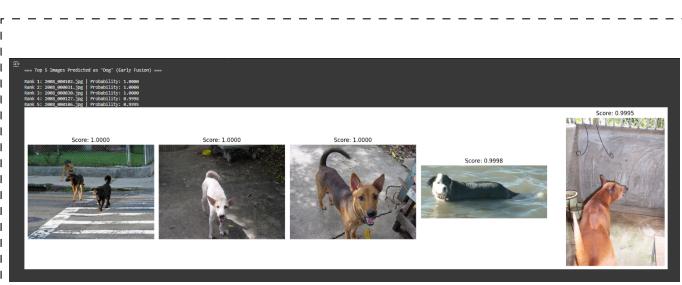


Fig. 15: Top 5: early fusion