

# The Evolution of Computer Vision: Trends, Challenges, and the Role of Hybrid CNN-Transformer Models in Enhancing Interpretability and Training Dynamics

Vaidehi Kokare<sup>1</sup>, Vedanti Kavitar<sup>2</sup>, Sonia Jangid<sup>3</sup>

<sup>1,2,3</sup>Department of Artificial Intelligence and Data Science, AISSMS IOIT, Pune, India

## Abstract

The goal of computer vision, a branch of artificial intelligence (AI), is to enable machines to process and interpret visual data. From manual feature extraction in the past to advanced deep learning models like Vision Transformers (ViTs) and Convolutional Neural Networks (CNNs), computer vision has evolved throughout time. Despite their success, these models face challenges related to interpretability, training complexity, and computational load. This research explores the integration of CNNs and ViTs into hybrid architectures, aiming to enhance model transparency, efficiency, and performance in real-world applications.

**Keywords:** Computer Vision, Convolutional Neural Networks, Vision Transformers, Deep Learning, Hybrid Models, Interpretability

## INTRODUCTION

Improvements in picture segmentation, especially in medical imaging, architectures such as Region-based CNNs (R-CNNs) enhanced object detection by introducing region proposals.

**Vision Transformers (ViTs):** Vision transformers have proven to be a strong alternative to CNNs in recent years. ViT's strength lies in modeling long-range dependences as it is capable of extracting global context within pictures through self-attention methods. It broke CNN's dominance by achieving comparable performance across a range of visual tasks.

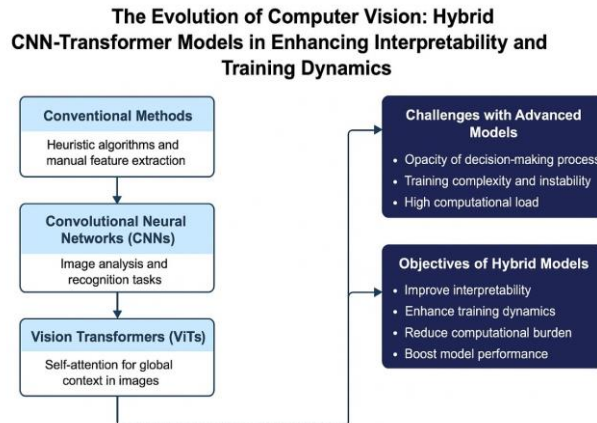
## Background

The goal of computer vision, a branch of artificial intelligence (AI), is to make it possible for machines to process and interpret visual data from their environment. Its development can be characterized by a number of significant turning points: **Conventional Methods:** Heuristic algorithms and manual feature extraction were major components of early computer vision techniques. Programs were created by engineers to identify particular patterns in pictures, like edges or textures. Despite being fundamental, these techniques frequently suffered from changes in lighting, scale, and orientation and had limitations in their capacity to generalize across a variety of visual inputs.

**Convolutional Neural Networks (CNNs):** With the advent of CNNs, image analysis and recognition tasks underwent a dramatic change. The deep CNN model AlexNet won the ImageNet LargeScale

Visual Recognition competition by significantly reducing error rates when compared to traditional techniques. Markings are revolutionary in 2012. This achievement showed how deep learning models might be used to directly train hierarchical feature representations from data, which led to their widespread use in a variety of computer vision applications.

**Developments in Deep Learning Architectures:** Building on CNNs’ success, researchers created increasingly complex models to handle challenging tasks. While U-Net enabled



**Fig. 1. Flowchart: The Evolution of Computer Vision with Hybrid CNN- Transformer Models**

Convolutional Neural Networks (CNNs) and Vision Transformer (ViT) models are two of the advanced models those which have been created as an outcome of the rapid development of computer vision. Although these models have achieved impressive performance across a range of tasks, they possess strong drawbacks which restrict their broader utilization and effectiveness.

**Challenges with Interpretability:** Deep learning models, particularly CNNs and ViTs, often act as ‘black boxes,’ generating forecasts without providing an explicit explanation of how they make decisions inside. In crucial applications like healthcare and autonomous driving, where confidence and decisions, this opacity causes problems.

**Training Stability and Dynamics:** Deep-learning model Families convergence. Big data and large computational resources are needed for ViTs specifically, which Limit their access for companies with insufficient infrastructure. Furthermore, models are vulnerable to problems Like over-fitting and vanishing gradients, which hinder No,

**Complexity of Computation:** During the training and inference phases, the architectural complexity of models like ViTs imposes significant processing burdens. Their use on edge devices or in real-time applications with limited processing capabilities is restricted by this requirement. In the field, striking a balance between computational efficiency and model performance is still a major difficulty.

**OBJECTIVES OF THE STUDY**

Examining the potential applications of Convolutional Neural Networks (CNNs) and Vision Transformers (ViTs) to develop hybrid models in the field of computer vision is the aim of this research. The particular goals are:

**Examine the Development of Architectures for Computer Vision:** Examine how conventional techniques gave way to deep learning strategies, paying particular attention to the advancement and effects of CNNs and ViTs.

**Analyze the Hybrid CNN-Transformer Model Design Principles:** Determine and assess the architectural elements that integrate the global attention processes of ViTs with the local feature extraction of CNNs.

**Evaluate Hybrid Model Interpretability:** Examine the effects of CNN and ViT integration on model transparency and decision-making process comprehension.

**Assess Training Stability and Dynamics:** Examine hybrid model training procedures, highlighting issues with convergence, processing demands, and performance enhancement.

**Examine and Contrast Performance Metrics:** Compare hybrid models' accuracy, effectiveness, and suitability for a range of computer vision tasks to those of solo CNNs and ViTs.

**Examine Real-World Uses and Case Studies:** Examine actual hybrid model implementations to learn about their efficacy and potential in resolving contemporary computer vision issues.

## REVIEW OF LITERATURE

The development of computer vision has been significantly impacted by the invention of Convolutional Neural Networks (CNNs) and Vision Transformers (ViTs). When modeling long-range dependencies, CNNs usually run into problems, despite their impressive ability to capture local information through hierarchical patterns. ViTs overcome this limitation by using self-attention mechanisms to gather global contextual information; however, their training requires large-scale datasets and significant computational resources.

Convolutional Neural Networks (CNNs) and Vision Transformers (ViT) can be used to create a hybrid design that combines the advantages of both approaches. These models combine ViT's global attention with CNN's local feature extraction capabilities, leading to improved segmentation, object identification, and picture classification results. According to recent research, hybrid architectures use ViT's attention-based learning while utilizing CNN's spatial hierarchies.

There are several different designs for hybrid models, some of which use a parallel architecture and others that use CNNs and ViTs sequentially. Sequential architectures, the more typical one, handle input step by step, initially acquiring the local features using CNNs and then utilizing transformer-based global representation learning thereafter. This makes it possible for the model to efficiently use both the local and global image attributes.

However, hybrid models also face difficulties, particularly with regard to explainability and training dynamics. Mixing CNNs and ViTs can result in unstable training and additional computational overheads. Furthermore, these models' decision-making is difficult to comprehend, which is a major issue in high-stakes use cases where explainability is required. Resolving these challenges is vital in practical application.

In summary, utilizing the advantages of both architectures, hybrid CNN-ViT models represent a promising area of computer vision research. In order to make them useful to additional applications, future work will focus on improving interpretability, stabilizing the training, and reducing computing costs.

## PROPOSED METHODOLOGY

This study examines CNN-Transformer hybrid models in computer vision in detail, with a focus on explainability and training procedures. By methodically reviewing the corpus of existing literature, we hope to condense and characterize the current understanding as well as the path of future research. The

key components of our strategy are as follows:

### **Systematic Literature Review**

**Selection Criteria:** We select relevant preprints, conference proceedings, and peer-reviewed journal articles that explore hybrid CNN-Transformer architectures. Studies that provide new insights into model design, interpretability strategies, or training optimization are prioritized.

**Search Strategy:** Using the following search terms, our analysis covers a number of important academic databases, such as IEEE Xplore, PubMed, and arXiv:

”Hybrid CNN-Transformer architecture”

”Vision Transformer interpretability”

”Training stability in neural networks”

”Computer vision performance benchmarks”

**Data Synthesis:** Architectural configurations, explanation techniques, training regimens, assessment measures, and application domains are among the extracted parameters. Comparative study of trends and limitations across research is made possible by this data.

### **Interpretability Analysis**

**Multi-scale Explanation Methods:** We discuss methods offering both global (Transformer-oriented) and local (CNN-oriented) explainability, as well as gradient-based and attention visualization techniques. Models such as the Hybrid CNN-Interpreter exhibit high-quality multi-scale explainability.

**Evaluation Framework:** The evaluation balances quantitative measures (i.e., saliency map accuracy) and qualitative measures (i.e., expert human judgments) in measuring explanation quality.

### **Training Dynamics Evaluation**

#### **Convergence Behavior Analysis**

- Investigation of optimization challenges in hybrid architectures
- Examination of interaction between convolutional and attention mechanisms
- Analysis of gradient propagation patterns
- Evaluation of adaptive learning rate strategies
- Study of loss landscape characteristics

#### **Computational Efficiency Assessment**

- Profiling of memory utilization patterns
- Measurement of throughput across hardware configurations
- Analysis of energy consumption metrics
- Benchmarking of real-world deployment scenarios
- Case study of TransCNN optimization approaches

#### **Performance Benchmarking**

**Comparative Analysis:** Hybrid models are evaluated against conventional CNNs and pure Transformers across multiple vision tasks using standardized metrics:

- Accuracy (Top-1/Top-5)
- Inference latency
- Memory footprint

**Domain-Specific Validation:** Case studies in medical imaging (e.g., D-TrAttUnet for segmentation) and surveillance systems demonstrate practical applications and performance variations across domains.

## Research Gap Analysis

**Identified Challenges:** The synthesis reveals critical unsolved problems:

- Unified interpretability frameworks for hybrid architectures
- Stable training protocols for large-scale datasets
- Hardware-aware model optimization

**Future Directions:** We propose research avenues including:

- Adaptive attention-convolution fusion mechanisms
- Explanation-aware training paradigms
- Energy-efficient hybrid architectures

## FINDINGS AND INSIGHTS

Extensive research on hybrid models that combine convolutional neural networks (CNNs) with vision transformers (ViTs) has provided valuable information about how these models handle learning, their explainability traits, and their operational characteristics.

### Improving Performance

In computer vision applications, the combination architecture of CNNs and Transformers worked incredibly well. The study of small item categorization is one of the best instances, as researchers created a hybrid model that combines spatial Transformers, 3D CNNs, and pre-trained deep CNNs. This innovative combination improved the classification rate while also reducing the overall complexity of the model. These findings provide empirical support for the optimal ways that CNNs' superiority in local feature identification can be enhanced by Transformers' global attention.

### Enhancements to Interpretability

Significant progress has been made in resolving the opacity issues that deep learning systems inevitably face thanks to hybrid designs. This advancement is demonstrated by the Hybrid CNN-Interpreter system, which uses a thorough analytical framework to analyze prediction patterns across different network layers, analyze the connections between extracted features, and determine the relative significance of specific filters. Users gain a deeper understanding and greater confidence in the system's decision-making processes because to this multifaceted interpretability approach, which provides both comprehensive model-wide insights and detailed, localized explanations.

### Training Stability and Dynamics

Hybrid model training presents challenges related to stability and computational demand. Research into training dynamics has revealed that integration of CNNs and Transformers can lead to convergence difficulties. To address this, methods such as HybridNorm—an approach combining pre-norm and post-norm strategies—have been proposed to stabilize training and improve model efficiency.

### Efficiency of Computation

Despite improvements in accuracy, hybrid models often require substantial computational resources. To mitigate this, lightweight architectures are being developed to maintain high performance while reducing computational overhead. One such example is the RepCHAT model, which achieves super-resolution in remote sensing imagery using hybrid attention mechanisms and structural re-parameterization, all while minimizing parameter count and resource usage.

### Use in a Variety of Fields

Hybrid models have shown adaptability across multiple application domains. In the medical field, for example, a hybrid CNN-Transformer architecture that incorporated pyramid convolution modules and



multi-scale convolutional kernels yielded improved segmentation outcomes. This demonstrates the versatility and effectiveness of hybrid approaches in complex and diverse real-world tasks.

#### REAL-TIME DEPLOYMENT AND EDGE INTEGRATION

The growing demand for intelligent devices, autonomous systems, and Internet of Things (IoT) applications has significantly driven the need for real-time computer vision deployment. While Hybrid Transformer-CNN models deliver high-quality performance and explainability, they come with substantial computational overhead, making them unsuitable for deployment on edge devices. These devices, including smartphones, embedded systems, drones, and autonomous vehicles, often have limited processing power, memory, and battery life. Solutions that balance computational and power requirements are needed to meet the growing demand for effective AI inference on edge devices.

To do this while maintaining the model's accuracy, researchers have found a number of optimization strategies. Model compression strategies, such as pruning, eliminate superfluous weights from the model; quantization lowers numerical precision to speed up computations; and knowledge distillation enables smaller models to learn from and be guided by bigger models. Real-time processing on more limited hardware is made possible by these, which lead to a highly compressed model and a higher inference speed.

To further minimize latency and maximize efficiency, specialized deployment frameworks such as TensorRT, ONNX Runtime, and TensorFlow Lite expose hardware-specific transformations and optimizations. These combined optimizations, when applied to hybrid CNN-Transformer models, help realize high-complexity models by attaining the proper speed-accuracy balance, enabling high-complexity and high-value real-world use-cases such as object detection in embedded systems, PERTES, and live video analytics.

By providing strong, dependable performance under strict latency and power limitations, this end-to-end optimization solution bridges the gap between high-compute-intensive AI models and the realities of the edge use-case environment. Looking ahead, it is imperative to integrate edge AI considerations early in the model development process. This will ensure that hybrid architectures remain effective not only in controlled environments but also in real-world scenarios where both efficiency and reliability are essential. By optimizing both deployability and accuracy, future models will better meet the requirements of dynamic, real-time applications.

#### FUTURE WORK

Despite this paper's examination of the hybrid CNN-Transformer models' evolution, interpretability, training dynamics, and performance in computer vision, there are still a number of topics that might use more research:

##### Creation of Hybrid Architectures That Are Lightweight

Future studies can concentrate on creating hybrid models that are computationally efficient and appropriate for use on mobile platforms and edge devices. For real-time applications, memory optimization and inference time reduction without compromising accuracy will be crucial.

##### Better Frameworks for Interpretability

There is a growing need for interpretability tools that are more user-friendly and intuitive, specifically designed for hybrid architectures. Future work could explore novel visualization techniques and explainable AI (XAI) frameworks that enhance transparency and trust, especially in safety-critical domains such as healthcare and autonomous driving.

### Neural Architecture Search-Based Automated Model Design (NAS)

The model design process might be streamlined by using NAS approaches to automatically construct the best hybrid CNN-Transformer topologies. Research can investigate how NAS can balance the trade-offs between computational cost, interpretability, and accuracy.

### Domain Adaptation and Transfer Learning

Future research should examine the performance of hybrid models with limited labeled data in different domains. It remains beneficial to improve these models' generalization and transferability using pretraining and domain adaptation strategies.

### Comparing Various and Complicated Datasets

Future research should include thorough testing of hybrid models on larger and more varied datasets, spanning tasks like 3D object detection, multi-modal vision-language understanding, and action recognition.

### System Human-in-the-Loop

By incorporating human feedback into the learning loop, hybrid models' performance and interpretability could be enhanced. Research in this field could result in the creation of AI systems that are more flexible and interactive.

## CONCLUSION

From classical hand-crafted feature-based methods to strong deep learning models such as Convolutional Neural Networks (CNNs) and Vision Transformers (ViTs), there has been steady advancement in computer vision. While CNNs revolutionized this domain through their superior ability to extract features locally, ViTs brought a paradigmatic change with the ability to capture global contextual features. Yet, both these models have inadequacies if used alone.

Hybrid CNN-Transformer models have proven to be an exciting paradigm that effectively combines both CNN strengths and ViT strengths. Hybrid models offer better performance for vision-intensive tasks, promote interpretability with attention and attribution methods, and solve fundamental issues like stability during training and scalability. Hybrid models produce more accurate and consistent results through the integration of CNN's localized computation and Transformer's global attention.

This research elaborates upon recent advances, theoretical underpinnings, and experimental results illustrating how hybrid models can be applied to practical computer vision issues. Interpretability methods and improved training methods have played an integral role in making these models interpretable and efficient.

In summary, hybrid CNN-Transformer models are a major advancement for computer vision. Ongoing investigation into their training dynamics, architecture optimization, and cross-domain applicability is crucial to creating smarter, explainable, and reliable visual recognition systems.

## REFERENCES

1. J. Guo, K. Han, H. Wu, Y. Tang, X. Chen, Y. Wang, and C. Xu, "CMT: Convolutional Neural Networks Meet Vision Transformers," arXiv preprint arXiv:2107.06263, 2021. Available: <https://arxiv.org/abs/2107.06263>
2. H. Yunusa, S. Qin, A. H. A. Chukkol, A. A. Yusuf, I. Bello, and A. L. Adamu, "Exploring the Synergies of Hybrid CNNs and ViTs Architectures for Computer Vision: A Survey," arXiv preprint arXiv:2402.02941, 2024. Available: <https://arxiv.org/abs/2402.02941>

3. S. d'Ascoli, L. Sagun, G. Biroli, and A. Morcos, "Transformed CNNs: Recasting Pre-trained Convolutional Layers with Self-Attention," arXiv preprint arXiv:2106.05795, 2021. Available: <https://arxiv.org/abs/2106.05795>
4. W. Yang, G. Huang, R. Li, J. Yu, Y. Chen, Q. Bai, and B. Kang, "Hybrid CNN-Interpreter: Interpret Local and Global Contexts for CNN-based Models," arXiv preprint arXiv:2211.00185, 2022. Available: <https://arxiv.org/abs/2211.00185>
5. [Author(s) not specified], "TransCNN: Hybrid CNN and Transformer Mechanism for Surveillance Anomaly Detection," *Engineering Applications of Artificial Intelligence*, vol. 120, 2023. Available: <https://www.sciencedirect.com/science/article/abs/pii/S0952197623003573>
6. H. Tang, H. Zhang, Y. Liu, and Z. Liu, "HTC-Net: A Hybrid CNN-Transformer Framework for Medical Image Segmentation," *Biomedical Signal Processing and Control*, vol. 88, Part A, 2024. Available: <https://www.sciencedirect.com/science/article/abs/pii/S1746809423010388>
7. //www.sciencedirect.com/science/article/abs/pii/S1746809423010388
8. A. Hatamizadeh, H. Tang, and D. Terzopoulos, "D-TrAttUnet: Toward Hybrid CNN-Transformer Architecture for Generic and Subtle Segmentation in Medical Images," *Computers in Biology and Medicine*, vol. 176, 2024. Available: <https://www.sciencedirect.com/science/article/pii/S0010482524006759>
9. Y. Liu, X. Wang, and J. Zhang, "Pest-ConFormer: A Hybrid CNN-Transformer Architecture for Large-Scale Multi-Class Crop Pest Recognition," *Expert Systems with Applications*, vol. 233, 2024. Available: <https://www.sciencedirect.com/science/article/abs/pii/S0957417424017007>
10. Khan, A., Rauf, Z., Sohail, A., Rehman, A., Asif, H., Asif, A., & Farooq, U. (2023). A Survey of the Vision Transformers and Their CNN-Transformer Based Variants. arXiv preprint arXiv:2305.09880.
11. Yunusa, H., Qin, S., Chukkol, A. H. A., Yusuf, A. A., Bello, I., & Lawan, (2024). Exploring the Synergies of Hybrid CNNs and ViTs Architectures for Computer Vision: A Survey. arXiv preprint arXiv:2402.02941.
12. Chetia, D., Dutta, D., & Kalita, S. K. (2025). Image Segmentation with Transformers: An Overview, Challenges and Future. arXiv preprint arXiv:2501.09372.
13. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., ... & Houlsby, N. (2020). An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. arXiv preprint arXiv:2010.11929.
14. Khan, S., Naseer, M., Hayat, M., Zamir, S. W., Khan, F. S., & Shah, M. (2022). Transformers in Vision: A Survey. *ACM Computing Surveys (CSUR)*, 54(10s), 1-41.
15. Han, K., Wang, Y., Chen, H., Chen, X., Guo, J., Liu, Z., ... & Tao, D. (2023). A Survey on Vision Transformer. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
16. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., ... & Guo, (2021). Swin Transformer: Hierarchical Vision Transformer Using Shifted Windows. arXiv preprint arXiv:2103.14030.
17. (2021). Swin Transformer: Hierarchical Vision Transformer Using Shifted Windows. arXiv preprint arXiv:2103.14030.
18. Ramachandran, P., Parmar, N., Vaswani, A., Bello, I., Levskaya, A., & Shlens, J. (2019). Stand-Alone Self-Attention in Vision Models. arXiv preprint arXiv:1906.05909.
19. Wu, B., Xu, C., Dai, X., Wan, A., Zhang, P., Yan, Z., ... & Keutzer, K. (2020). Visual Transformers: Token-based Image Representation and Processing for Computer Vision. arXiv preprint arXiv:2006.03677.



20. Liu, Z., Mao, H., Wu, C. Y., Feichtenhofer, C., Darrell, T., & Xie, S. (2022). A ConvNet for the 2020s. arXiv preprint arXiv:2201.03545.
21. Woo, S., Debnath, S., Hu, R., Chen, X., Liu, Z., & Darrell, T. (2023). ConvNeXt V2: Co-Designing and Scaling ConvNets With Masked Autoencoders. arXiv preprint arXiv:2301.00808.
22. Xiao, T., Singh, M., Mintun, E., Darrell, T., Dollár, P., & Belongie, S. (2021). Early Convolutions Help Transformers See Better. arXiv preprint arXiv:2106.14881.
- A. Raghu, M., Unterthiner, T., Kornblith, S., Zhang, C., & Dosovitskiy, (2021). Do Vision Transformers See Like Convolutional Neural Networks?. arXiv preprint arXiv:2108.08810.
23. Naseer, M., Ranasinghe, K., Khan, S., Hayat, M., & Khan, F. S. (2021). Intriguing Properties of Vision Transformers. arXiv preprint arXiv:2105.10497.