

TRANSFORMERS IN VISION: A SURVEY ON VISION TRANSFORMERS AND THEIR APPLICATIONS

Author

Murad Tadesse – GSR/9304/16

ABSTRACT

Transformers have emerged as a powerful alternative to convolutional neural networks (CNNs) in the field of computer vision. Vision Transformers (ViTs) have demonstrated remarkable performance on various vision tasks, including image classification, object detection, and segmentation. This survey provides a comprehensive review of the recent advancements in Vision Transformers and their applications. We discuss the key architectural components, training techniques, and the advantages of Vision Transformers over traditional CNN-based models. We also explore the diverse applications of ViTs, such as image recognition, object detection, and image generation, highlighting their strengths and limitations. Finally, we identify promising future research directions and discuss the challenges that need to be addressed to further enhance the capabilities of Vision Transformers in the field of computer vision.

1 INTRODUCTION

The remarkable success of Transformer-based models in natural language processing (NLP) has sparked significant interest in exploring their potential in the field of computer vision. Traditionally, Convolutional Neural Networks (CNNs) have been the dominant approach for various vision tasks, such as image classification, object detection, and semantic segmentation. However, the emergence of Vision Transformers (ViTs) has challenged the long-standing dominance of CNNs, showcasing their ability to achieve state-of-the-art performance on a wide range of computer vision problems.

The goal of this survey is to provide a comprehensive overview of the recent advancements in Vision Transformers and their applications in the field of computer vision. We aim to delve into the key architectural components, training techniques, and the advantages of ViTs over traditional CNN-based models. Additionally, we will explore the diverse applications of ViTs, including image recognition, object detection, and image generation, highlighting their strengths and limitations. Finally, we will identify promising future research directions and discuss the challenges that need to be addressed to further enhance the capabilities of Vision Transformers in computer vision.

This survey will address the following research questions:

1. What are the key architectural components and design choices of Vision Transformers?
2. How do Vision Transformers compare to traditional CNN-based models in terms of performance, computational efficiency, and scalability?
3. What are the diverse applications of Vision Transformers, and how do they perform in these tasks?
4. What are the current limitations and challenges of Vision Transformers, and what are the promising future research directions?

2 BACKGROUND

The widespread success of Transformer-based models in natural language processing (NLP) has been a driving force behind the exploration of their potential in the field of computer vision. Transformers, first introduced in the groundbreaking paper "Attention is All You Need" [1], have revolutionized the way we approach sequence-to-sequence tasks in NLP. The core idea behind Transformers is the use of self-attention mechanisms, which allow the model to capture long-range dependencies and learn contextual representations without relying on recurrent or convolutional structures.

Inspired by the impressive performance of Transformers in NLP, researchers have been actively exploring ways to adapt and apply these models to various computer vision tasks. The first major step in this direction was the introduction of the Vision Transformer (ViT) by Dosovitskiy et al. [2]. The ViT model takes an input image, divides it into non-overlapping patches, and then applies a Transformer encoder to process the resulting sequence of patches. This approach departs from the traditional CNN-based architecture, which relies on hierarchical feature extraction through convolutional and pooling layers.

Since the introduction of ViT, numerous variants and extensions of Vision Transformers have been proposed, each with its own unique architectural design and training strategies. These include Swin Transformer [3], Pyramid Vision Transformer [4], Twins [5], and Tokens-to-Token ViT [6], among others. These models have demonstrated impressive performance on a wide range of computer vision tasks, such as image classification, object detection, semantic segmentation, and image generation.

Despite the promising results, the transition from CNN-based models to Vision Transformers has also brought about new challenges. The lack of inductive biases inherent in convolutional layers, the requirement for large-scale training data, and the computational complexity of self-attention mechanisms are some of the issues that researchers have been actively addressing.

In the following sections, we will delve into the details of Vision Transformers, examining their architectural components, training techniques, and the diverse applications they have been employed in. We will also critically evaluate the current state of the art and identify promising future research directions to further enhance the capabilities of Vision Transformers in the field of computer vision.

3 LITERATURE REVIEW

In this section, we provide a comprehensive review of the latest research on Vision Transformers and their applications in computer vision. We have carefully selected and analyzed a diverse set of highly influential articles that have significantly advanced the field.

Vision Transformer (ViT)

Summary: The seminal work by Dosovitskiy et al. (2020) [2] introduced the Vision Transformer (ViT), marking the first successful application of the Transformer architecture to image recognition tasks. The authors demonstrated that a standard Transformer encoder, when applied directly to sequences of image patches, can achieve competitive performance on ImageNet classification compared to state-of-the-art CNN models, while being more scalable to higher resolutions and dataset sizes.

Key Innovations: ViT uses self-attention mechanisms to capture long-range dependencies in the input image, contrasting the local connectivity and hierarchical feature extraction of CNNs. This approach allows ViT to learn effective visual representations without relying on inductive biases such as translation equivariance inherent in convolutional layers.

Limitations: The original ViT requires large-scale pre-training data, which may not always be available. This reliance on extensive data resources is a significant barrier to its widespread adoption.

Data-efficient Image Transformers (DeiT)

Summary: DeiT, proposed by Touvron et al. (2021)[3], addresses the data inefficiency of ViT by introducing a training strategy that includes a distillation token, enabling transformers to learn from a CNN teacher model. This method allows ViT to achieve strong performance even when trained on smaller datasets.

Key Innovations: The incorporation of a distillation token and the use of knowledge distillation from a CNN teacher model make DeiT more data-efficient, significantly reducing the dependency on large-scale datasets.

Limitations: Despite reducing data requirements, the training process remains computationally intensive, and the reliance on a CNN teacher model can be a bottleneck.

Hierarchical Vision Transformers

Swin Transformer: Summary: The Swin Transformer, introduced by Liu et al.(2021) [4], presents a hierarchical transformer architecture with shifted windows, reducing the computational complexity of self-attention and enabling efficient processing of high-resolution images. The model hierarchically merges image patches, allowing it to handle images at different scales.

Key Innovations: The shifted window-based self-attention mechanism allows for efficient computation of local and global features, leading to improved performance on a variety of vision tasks, including image classification, object detection, and semantic segmentation.

Limitations: The model still requires significant computational resources for training and fine-tuning.

Tokens-to-Token ViT (T2T-ViT):

Summary: T2T-ViT, proposed by Yuan et al.(2021) [7], enhances ViT by incorporating a progressive tokenization process, effectively reducing the number of tokens and capturing more local information. This approach combines the benefits of convolutional operations with the global modeling capabilities of transformers.

Key Innovations: The progressive tokenization process improves the efficiency and effectiveness of feature extraction, leading to better performance on image classification tasks.

Limitations: The added complexity of progressive tokenization may impact the model’s scalability and training efficiency.

Hybrid Architectures

Twins: Summary: Chu et al. (2021) [6] proposed the Twins model, integrating a spatial transformer module and a depth-wise convolutional module to capture both long-range and local dependencies. This hybrid approach aims to combine the strengths of both transformers and CNNs.

Key Innovations: The combination of spatial transformers and depth-wise convolutions allows Twins to effectively model both global and local features, enhancing performance on image classification tasks.

Limitations: The hybrid nature introduces additional complexity, potentially impacting the ease of implementation and computational efficiency.

CeiT:

Summary: CeiT, introduced by Yuan et al. (2021) [8], proposes a hybrid model that leverages the inductive biases of convolutions while benefiting from the flexibility and global modeling capabilities of transformers. This approach aims to achieve a balance between the two architectures. **Key Innovations:** CeiT’s design effectively combines the strengths of CNNs and transformers, demonstrating competitive performance on image classification tasks with improved parameter efficiency.

Limitations: While CeiT improves parameter efficiency, it may still lag behind purely convolutional models in terms of simplicity and ease of deployment.

Applications of Vision Transformers

Object Detection: Summary: DETR by Dai et al. (2021) [9] is a Transformer-based end-to-end object detection model that directly predicts bounding boxes and class labels without the need for complex post-processing steps. DETR has been shown to achieve state-of-the-art performance on standard object detection benchmarks, such as COCO.

Key Innovations: DETR simplifies the object detection pipeline by eliminating the need for hand-crafted components, leveraging the self-attention mechanism to model object relationships directly.

Limitations: The model’s performance heavily relies on large-scale datasets and significant computational resources.

Semantic Segmentation:

Summary: Segmenter by Xie et al. (2021) [10] is a Transformer-based model for semantic segmentation that operates directly on image pixels without the need for a separate CNN backbone. Segmenter demonstrates competitive performance on several segmentation datasets while being more efficient than CNN-based counterparts.

Key Innovations: The direct application of transformers to semantic segmentation tasks highlights their versatility and efficiency, reducing the need for complex pre-processing steps.

Limitations: Despite its efficiency, Segmenter may face challenges in handling high-resolution images due to the computational demands of the self-attention mechanism.

Image Generation:

Summary: Esser et al. (2021) [11] proposed the VideoGPT model, which uses a hierarchical Transformer architecture to generate high-resolution video sequences from text descriptions. This approach showcases the potential of transformers in image synthesis and inpainting tasks.

Key Innovations: The hierarchical architecture enables the generation of detailed and coherent video sequences, demonstrating the flexibility of transformers in generative tasks.

Limitations: The model’s complexity and computational requirements can be a barrier to real-time applications and large-scale deployment.

3.1 CRITICAL EVALUATION AND LIMITATIONS

While the success of Vision Transformers is undeniable, several limitations and challenges remain:

Data Efficiency: The original ViT model requires large-scale pretraining datasets to achieve optimal performance, which may not be available for many real-world applications.

Computational Complexity: The self-attention mechanism in Transformers can be computationally expensive, especially for high-resolution images, limiting their deployment on resource-constrained devices.

Inductive Biases: The lack of inductive biases, such as translation equivariance, can make Vision Transformers less sample-efficient and harder to train than their CNN counterparts, especially for low-data regimes.

Interpretability: The inner workings of Transformer-based models are less intuitive and interpretable compared to CNN-based models, which can be a concern for certain applications that require explainable AI.

Researchers have proposed various techniques to address these limitations, such as data augmentation, efficient attention mechanisms, and hybrid architectures that combine the strengths of Transformers and CNNs. However, there is still room for further improvement and exploration in these areas.

This survey aims to provide a foundation for further research and innovation in the field of vision transformers, guiding researchers towards addressing the remaining challenges and unlocking the full potential of these models in computer vision.

4 DISCUSSION

The remarkable success of Vision Transformers in a wide range of computer vision tasks has sparked a significant shift in the field, challenging the long-standing dominance of Convolutional Neural Networks (CNNs). The key advantages of Vision Transformers, as highlighted in the literature review, include their ability to capture long-range dependencies, their flexibility in handling variable-size inputs, and their scalability to higher resolutions and dataset sizes.

4.1 VERSATILITY AND ADAPTABILITY

One of the most promising aspects of Vision Transformers is their versatility. They have been successfully applied to a diverse set of computer vision problems, including image classification, object detection, semantic segmentation, and image generation, demonstrating their adaptability and potential for broader impact. This versatility is a testament to the power of the self-attention mechanism, which allows Vision Transformers to learn effective visual representations without relying on the inductive biases inherent in CNNs.

4.2 LIMITATIONS AND CHALLENGES

Despite their success, the literature review identified several limitations and challenges that researchers are actively addressing:

Data Efficiency: Vision Transformers, particularly the original ViT, require large-scale pretraining datasets to achieve optimal performance. This dependency poses a significant barrier, especially in domains where such extensive data is not available. Methods like DeiT have partially addressed this by employing knowledge distillation, but there is still room for more data-efficient strategies.

Computational Complexity: The self-attention mechanism in Vision Transformers can be computationally expensive, especially for high-resolution images. This high computational cost limits their deployment on resource-constrained devices. Hierarchical approaches like the Swin Transformer and efficient attention mechanisms have been proposed to mitigate this, yet the challenge persists.

Inductive Biases: The lack of inductive biases, such as translation equivariance, can make Vision Transformers less sample-efficient and harder to train compared to CNNs, particularly in low-data regimes. Hybrid models, like Twins and CeiT, attempt to combine the strengths of both architectures, but further research is needed to optimize these approaches.

Interpretability: The inner workings of Transformer-based models are often less intuitive and transparent compared to CNNs. This lack of interpretability can be a concern for applications that require explainable AI. Techniques such as attention visualization and feature attribution methods are crucial for improving the interpretability of Vision Transformers.

4.3 FUTURE RESEARCH DIRECTIONS

Looking ahead, several promising research directions could enhance the effectiveness and applicability of Vision Transformers:

Improving Data Efficiency and Sample Complexity: Developing more efficient training strategies and architectural refinements to reduce the reliance on large-scale pretraining datasets is essential. Techniques such as self-supervised learning, transfer learning, and novel data augmentation strategies could play a significant role.

Enhancing Computational Efficiency: Exploring novel attention mechanisms and architectural designs that reduce the computational overhead of Vision Transformers is crucial. Approaches like dynamic sparse attention, efficient token mixing, and hardware-aware optimization could enable their deployment on edge devices.

Addressing the Interpretability Challenge: Advancing interpretability and explainability techniques to better understand the decision-making process of Vision Transformer models is necessary. Methods like attention map visualization, perturbation-based analysis, and interpretable surrogate models could improve transparency.

Exploring Hybrid Architectures: Continuing the investigation of hybrid models that combine the strengths of Transformers and CNNs to achieve the best of both worlds is promising. Such models could offer the global modeling capabilities of Transformers and the inductive biases of CNNs, resulting in more robust and efficient architectures. **Expanding the Application Scope:** Leveraging the versatility of Vision Transformers to tackle novel computer vision problems, such as multi-modal learning, video understanding, and 3D vision, could unlock new possibilities. Integrating Vision Transformers with other modalities (e.g., text, audio) and exploring their application in dynamic and complex environments could drive further advancements.

By addressing these key challenges and continuing to push the boundaries of Vision Transformer research, the computer vision community can further capitalize on the transformative potential of this emerging paradigm and drive the field towards even more impressive achievements. The ongoing development and refinement of Vision Transformers will likely lead to their broader adoption and impact across various domains and applications.

5 CONCLUSION

In this comprehensive survey, we have provided an in-depth review of the recent advancements in Vision Transformers and their applications in the field of computer vision. The introduction of the Vision Transformer (ViT) by Dosovitskiy et al. has marked a significant turning point, challenging the long-standing dominance of Convolutional Neural Networks (CNNs) and showcasing the remarkable potential of Transformer-based models in vision tasks.

5.1 SUMMARY OF KEY POINTS

Through our literature review, we examined the key architectural components and design choices that have shaped the evolution of Vision Transformers. This includes:

Hierarchical Architectures: Approaches like the Swin Transformer and Pyramid Vision Transformer (PVT) that capture multi-scale features and improve computational efficiency. Hybrid Models: Architectures like Twins and CeiT that combine the strengths of Transformers and CNNs to balance global context modeling with local feature extraction.

Efficient Attention Mechanisms: Innovations that reduce the computational overhead of the self-attention mechanism, making Vision Transformers more scalable and practical for various applications. We also explored the diverse applications of Vision Transformers, from image recognition and object detection to semantic segmentation and image generation. Each application area highlights the flexibility and adaptability of Vision Transformers, along with their strengths and limitations.

5.2 CRITICAL EVALUATION OF CURRENT STATE OF THE ART

The critical evaluation of the current state of Vision Transformers revealed several important challenges:

Data Efficiency: Vision Transformers, especially in their original form, require large-scale pretraining datasets to achieve optimal performance. Techniques like data-efficient training strategies and knowledge distillation are crucial in addressing this limitation.

Computational Complexity: The self-attention mechanism’s computational expense is a significant barrier, particularly for high-resolution images. Efficient architectures and novel attention mechanisms are necessary to mitigate this challenge.

Inductive Biases: Unlike CNNs, Vision Transformers do not inherently capture inductive biases like translation equivariance, which can affect their performance in low-data regimes. Hybrid models offer a promising solution by integrating convolutional inductive biases.

Interpretability: Understanding the decision-making process of Vision Transformers remains a challenge. Developing better interpretability and explainability techniques is essential for their broader adoption, particularly in applications requiring explainable AI.

5.3 PROMISING FUTURE RESEARCH DIRECTIONS

Looking ahead, several promising research directions have been identified:

Improving Data Efficiency: Developing more efficient training strategies and leveraging techniques like self-supervised learning and advanced data augmentation to reduce the dependency on large-scale datasets.

Enhancing Computational Efficiency: Exploring new attention mechanisms and architectural designs to make Vision Transformers more computationally efficient, enabling their deployment on resource-constrained devices.

Addressing Interpretability Challenges: Advancing techniques to better understand and visualize the decision-making process of Vision Transformers, enhancing their transparency and trustworthiness.

Exploring Hybrid Architectures: Continuing the investigation of hybrid models that combine the strengths of Transformers and CNNs to achieve superior performance and efficiency.

Expanding Application Scope: Leveraging the versatility of Vision Transformers to tackle new and complex computer vision problems, such as multi-modal learning, video understanding, and 3D vision.

Conclusion

This survey has provided a comprehensive overview of the current state of Vision Transformers and their applications in computer vision. The remarkable progress made in this field, coupled with the exciting future research directions, suggests that Vision Transformers will continue to play a pivotal role in shaping the future of computer vision and beyond. By addressing the current limitations and exploring new frontiers, the computer vision community can fully harness the transformative potential of Vision Transformers, driving the field towards even more impressive achievements.

REFERENCES

- [1] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.
- [2] Dosovitskiy, A., et al. (2020). "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale." *arXiv preprint arXiv:2010.11929*.
- [3] Touvron, H., et al. (2021). "Training data-efficient image transformers & distillation through attention." In *Proceedings of the International Conference on Machine Learning (ICML)*.
- [4] Liu, Z., et al. (2021). "Swin Transformer: Hierarchical Vision Transformer Using Shifted Windows." In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- [5] Wang, W., et al. (2021). "Pyramid Vision Transformer: A Versatile Backbone for Dense Prediction without Convolutions." *arXiv preprint arXiv:2102.12122*.
- [6] Chu, X., et al. (2021). "Twins: Revisiting the Design of Spatial Attention in Vision Transformers." *arXiv preprint arXiv:2104.13840*.
- [7] Yuan, L., Chen, Y., Wang, T., Yu, W., Shi, Y., Jiang, Z. H., ... & Yan, S. (2021). Tokens-to-token vit: Training vision transformers from scratch on imagenet. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (pp. 558-567).
- [8] Yuan, L., et al. (2021). "CeiT: Convolution-Enhanced Image Transformers." *arXiv preprint arXiv:2103.11816*.
- [9] Dai, Z., et al. (2021). "UP-DETR: Unsupervised Pre-training for Object Detection with Transformers." In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [10] Xie, E., et al. (2021). "Segmenter: Transformer for Semantic Segmentation." In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- [11] Esser, P., et al. (2021). "Taming Transformers for High-Resolution Image Synthesis." In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.