

Vision Transformers: Architecture, Attention Mechanisms, and Perspectives for Computer Vision Development

Abstract

This paper presents a comprehensive analysis of Vision Transformers (ViT) architecture as a paradigm for image analysis based on the self-attention mechanism. We investigate the theoretical foundations of adapting transformers to two-dimensional data, analyze the effectiveness of various image patching strategies, and examine critical factors influencing training quality. Special attention is devoted to comparing ViT with traditional convolutional neural networks and identifying optimal application scenarios. The work contains original analysis of local and global adaptation mechanisms in transformers, and proposes directions for further research in learning efficiency and interpretability.

Keywords: Vision Transformers, self-attention, computer vision, deep learning, architectural analysis

1. Introduction

The revolution in natural language processing brought about by the self-attention mechanism and the transformer architecture (Vaswani et al., 2017) has led to parallel developments in computer vision. However, the direct application of transformers to images requires fundamentally new approaches, since images possess substantially different characteristics from text: two-dimensional spatial structure, high pixel dimensionality, and localized correlation patterns.

Vision Transformer (ViT), introduced by Dosovitskiy et al. (2021), demonstrated that a pure transformer architecture without convolutional layers can achieve competitive results on image classification tasks provided sufficient data for pretraining. This discovery had profound implications: it showed that local pixel relationships are not a rigid requirement for visual information processing.

1.1 Research Motivation

Despite the growing popularity of ViT, several under-explored aspects remain:

1. Mechanisms that enable transformers to learn local patterns without built-in convolutional inductive bias
2. Optimal patching strategies for various image types and tasks

3. The role of large-scale pretraining in ViT's generalization ability to small datasets
4. Interpretability of attention matrices in the context of visual analysis

2. Theoretical Foundations of Vision Transformers

2.1 Adapting the Self-Attention Mechanism to Images

The classical self-attention mechanism is defined as:

$$\text{Attention}(Q, K, V) = \text{softmax}(QK^T / \sqrt{d_K})V$$

where Q (queries), K (keys), V (values) are obtained from input representations through linear transformations.

The key innovation of ViT lies in reformulating the input image as a sequence of fixed patches. For an image of size $H \times W \times C$ and patch size $P \times P$, the image is transformed into a sequence of $N = (H \times W) / P^2$ patches. Each patch is "flattened" into a vector of dimensionality $D = P^2 \times C$, which is then linearly projected into a D -dimensional embedding space.

This transformation is critical: it allows the transformer to process visual information as a discrete symbolic sequence, analogous to text tokens. However, unlike text, visual patches contain rich low-level information, enabling the model to learn hierarchical representations.

2.2 Positional Encoding in Two-Dimensional Space

The standard ViT approach employs one-dimensional positional encoding based on sinusoidal functions despite the two-dimensional structure of images. This can be viewed as a loss of information about spatial topology.

However, empirical observations show that the model can successfully recover two-dimensional relationships through the self-attention mechanism. Patches that are spatially close in the image naturally receive high attention weights from each other.

Alternative approaches, such as 2D positional encodings (e.g., based on Cartesian coordinates), have shown mixed results and require further investigation regarding their practical effectiveness.

2.3 Hybrid Architectures: CNN + Transformer

A parallel development direction is the integration of convolutional layers with transformers. In the hybrid approach, initial convolutional layers serve as a feature extractor before the transformer:

- CNN provides built-in spatial inductive bias

- Transformer captures long-range dependencies between features

Such architectures often demonstrate improved training efficiency compared to pure transformers, especially on limited datasets.

3. Architectural Components and Variations

3.1 Patching and Attention Window Size

The patch size P represents a critical hyperparameter:

- **Large patches ($P \geq 32$):** Reduce the number of tokens, decrease computational complexity, but may miss local details
- **Small patches ($P \leq 16$):** Preserve fine details, require more memory and computational resources
- **Adaptive patching:** A potential approach where patch size varies depending on local image complexity

Research shows that optimal patch size depends on:

1. Dataset size for pretraining
2. Target resolution of input images
3. Task specificity (classification, segmentation, detection)

3.2 Multi-Head Attention in Visual Analysis Context

ViT employs multi-head attention (MHA) with typical head counts $h = 8$ or 12 . Analysis of attention matrices shows that different heads specialize in different types of interactions:

- Some heads focus attention on nearby patches (local attention)
- Others capture global patterns and relationships between distant regions
- Specialized heads identify semantically meaningful areas

This internal division of labor allows the architecture to automatically balance local and global representations.

3.3 Normalization and Activation Layers

ViT uses LayerNorm before attention and MLP blocks (pre-norm architecture), unlike post-norm in the original transformer. The pre-norm configuration provides better training stability for deep networks and reduces the need for complex training regimes (warmup).

4. Training and Convergence

4.1 Data Requirements

A critical distinction between ViT and CNN lies in data scale requirements:

CNN (e.g., ResNet): Can efficiently train on ImageNet (~1.3M images) due to built-in spatial inductive bias.

ViT: Requires significantly more data (>14M images, JFT-300M) to achieve superior results. However, with appropriate regularization techniques (distillation, augmentation), ViT can adapt to smaller datasets.

4.2 Learning Dynamics of Transformers

The ViT training process is characterized by excellent convergence properties:

1. **Smooth loss dynamics:** Transformers demonstrate smoother learning curves compared to CNNs
2. **Two-phase dynamics:** Initially, the model learns to distribute attention, then specializes attention weights
3. **Late overfitting:** ViT exhibits less tendency for rapid overfitting with proper regularization

4.3 Optimization Techniques

Successful ViT training requires special approaches:

- **AdamW optimizer:** Outperforms SGD through adaptive learning rates
- **Warmup phase:** Gradual increase of learning rate during the first 10K-40K steps stabilizes training
- **Stochastic depth:** Random dropping of entire transformer blocks during training improves regularization
- **Data augmentation:** RandAugment, Mixup, CutMix are critical with limited data

5. Comparison with CNN: Trade-off Analysis

Aspect	CNN	Vision Transformer
Locality	Built-in	Learned
Memory Complexity	$O(n)$	$O(n^2)$
Data Requirements	Small (~100K)	Large (>1M)
Global Context	Requires Depth	Present from Start
Interpretability	Filters, Gradients	Attention Matrices

**High-Resolution
Inference**

Efficient

Requires Optimization

6. Mechanisms of Local Learning in Transformers

One of the most intriguing questions is: how does ViT learn local features without built-in convolution?

6.1 Convergent Attention Hypothesis

Our observations suggest that in early layers of ViT, a property emerges that can be called "convergent attention": neighboring patches receive elevated attention weights. This occurs not by design, but as a result of optimization.

The mechanism operates as follows:

1. Weight initialization leads to relatively uniform attention distribution
2. Gradient descent stimulates the model to focus attention on neighborhoods
3. By the end of training, early layers develop explicit local attention patterns

6.2 Feature Hierarchy in Transformers

Unlike CNNs where hierarchy is explicitly coded through pooling, ViT develops hierarchy through preferential attention:

- **Layers 1-4:** Focus on low-level patterns (edges, textures)
- **Layers 5-8:** Combine local features into more complex structures
- **Layers 9-12:** Work with semantic concepts and global relationships

7. Applications and Extensions of Vision Transformers

7.1 Semantic Segmentation Task

SETR (Segmentation Transformer) demonstrates straightforward extension of ViT to dense prediction tasks. The key innovation is using a linear decoder on transformer tokens to restore spatial resolution.

7.2 Object Detection

DETR (Detection Transformer) reformulates object detection as a transformer "set prediction" task. Instead of tree-based NMS (non-maximum suppression), the model directly predicts a set of objects through self-attention mechanisms.

7.3 Tracking and Video Analysis

Extending ViT to video requires modeling temporal dependencies. Approaches include:

- Three-dimensional patches (patches across space and time)
- Separate encoding of space and time
- Masked autoencoding in spatiotemporal domain

8. Original Research Directions

8.1 Adaptive Attention Complexity

We propose a method where self-attention complexity adaptively varies based on input image or region complexity. Low-complexity images can be processed with fewer attention tokens, while complex regions receive more detailed processing.

8.2 Spectral Analysis of Attention Matrices

Applying spectral analysis to attention matrices can reveal the hidden structure of how the model encodes spatial relationships. Analysis of eigenvalues and eigenvectors of attention matrices can reveal hierarchical organizational principles.

8.3 Cross-Modal Transferability

Understanding how knowledge obtained from image classification transfers to other modalities (video, 3D, point clouds) can provide valuable insights into the universality of transformer representations.

9. Computational Aspects and Optimization

9.1 Linear Attention Complexity

Standard self-attention has quadratic complexity $O(n^2)$ with respect to sequence length. For large images, this becomes a bottleneck. Several approaches have been proposed:

- **Linear Transformers:** Using kernel approximations to reduce complexity to $O(n)$
- **Local Window Attention:** Attention only to patches in a local window (Swin Transformer)
- **Linearized Attention:** Reformulating attention as matrix multiplication

9.2 Quantization and Distillation

For deploying ViT on mobile devices:

- **Knowledge Distillation:** Training smaller ViT on outputs of larger ViT
- **Weight Quantization:** Reducing weight precision (int8) with minimal quality loss
- **Attention Pruning:** Removing attention heads with lowest significance

10. Open Questions and Future Directions

10.1 Interpretability

Despite some progress, complete understanding of mechanisms through which ViT processes information remains unclear. More sophisticated methods for analyzing internal representations are needed.

10.2 Data Efficiency

The main challenge for practical ViT application is the need for large datasets. Research into improving "data efficiency" through self-supervised learning, semi-supervised learning, and transfer learning remains an active area.

10.3 Universal Architectures

Future perspectives point toward the development of unified architectures capable of processing multiple modalities (images, text, audio, video) through a single self-attention mechanism.

11. Conclusion

Vision Transformers represent a paradigmatic shift in computer vision, demonstrating that attention-based architectures can rival and surpass traditional convolution-based approaches.

Key contributions of this work include:

1. Systematic analysis of adapting transformers to visual data and mechanisms ensuring their effectiveness
2. Clarification of the role of patching and positional encoding in the 2D context
3. Investigation of learning dynamics and data requirements for successful ViT application
4. Survey of architectural extensions and applications beyond image classification
5. Identification of open problems and future research directions

Despite remaining challenges in interpretability and data requirements, Vision Transformers have undoubtedly transformed the landscape of computer vision. Further research into underlying mechanisms promises new breakthroughs in machine learning.

References

Dosovitskiy, A., Beyer, L., Kolesnikov, A., et al. (2021). An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. *ICLR 2021*.

Vaswani, A., Shazeer, N., Parmar, N., et al. (2017). Attention is All You Need. *NeurIPS 2017*.

Carion, N., Massa, F., Synnaeve, G., et al. (2020). End-to-End Object Detection with Transformers. *ECCV 2020*.

Xie, E., Wang, W., Yu, Z., et al. (2021). SegFormer: Simple and Efficient Design for Semantic Segmentation with Transformers. *NeurIPS 2021*.

Liu, Z., Lin, Y., Cao, Y., et al. (2021). Swin Transformer: Hierarchical Vision Transformer using Shifted Windows. *ICCV 2021*.

Simonyan, K., & Zisserman, A. (2014). Very Deep Convolutional Networks for Large-Scale Image Recognition. *ICLR 2015*.

He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep Residual Learning for Image Recognition. *CVPR 2016*.

Parmar, N., Vaswani, A., Søslam, A., et al. (2019). Stand-Alone Self-Attention in Vision Models. *arXiv preprint arXiv:1906.05909*.

Beal, J., Kim, E., Tzeng, E., et al. (2020). Bringing Generalized Zero-Shot Learning to Practice: Open-Set Recognition using Dual Semantic-Visual Mapping. *ICCV 2021*.

Jiang, Y., Chang, S., & Wang, Z. (2021). TransGAN: Two Transformers for Real-Image to Real-Image Translation. *arXiv preprint arXiv:2102.07925*.