# A Comprehensive Survey on Vision Transformers: Foundations, Advances, Applications, and Future Directions

SurveyForge

**Abstract**— Vision Transformers (ViTs) have redefined computer vision by leveraging self-attention mechanisms from natural language processing to effectively model long-range dependencies in image data, offering notable improvements over traditional convolutional neural networks. This survey meticulously examines ViTs' architectural foundations, spanning core components like patch embeddings and self-attention mechanisms, their integration with other paradigms such as convolutional models for enhanced performance, and their application across varied vision tasks, including image classification, object detection, and multimodal learning. Emphasis is placed on recent advances in hybrid models that enhance task-specific performance, addressing the computational intensity of ViTs through efficient attention mechanisms, quantization, and pruning techniques. The paper identifies prevalent challenges in data efficiency, interpretability, and resource constraints, offering future directions for scalable, robust, and environmentally sustainable deployment. It also highlights the transformative impact of ViTs in cross-domain and emerging applications, underscoring the potential for further innovations in multimodal integration and responsible AI practices. The survey establishes Vision Transformers as a foundational architecture poised to advance the field of computer vision.

**Index Terms**—Vision Transformers Integration, Efficient Attention Innovations, Multimodal Learning Applications

✦

## 1 INTRODUCTION

VISION Transformers have emerged as a groundbreaking advancement in the field of computer vision, fundamentally altering the landscape of visual data processing by leveraging the self-attention mechanism previously successful in natural language processing (NLP) [1]. The genesis of Vision Transformers occurs against a backdrop where convolutional neural networks (CNNs) long dominated due to their proficiency in capturing spatial hierarchies through localized receptive fields. However, despite their efficacy, CNNs face limitations, particularly in modeling long-range dependencies due to the inherently hierarchical and local nature of convolution operations [2], [3].

The revolutionary capabilities of Vision Transformers are underpinned by their ability to effectively capture both local and global contextual information in images, akin to how Transformers in NLP models handle sequential data [4]. Unlike CNNs that convolve feature maps through kernels, Vision Transformers interpret an image as a sequence of patches, engaging self-attention mechanisms to capture relationships among distant patches. This architectural innovation not only replaces the convolution process with a non-local mechanism but also provides the benefit of parallel processing, making Transformers inherently scalable and capable of handling high-dimensional data [2].

A historical outline of Vision Transformers begins with their adaptation from Transformers used in NLP tasks, where their unmatched performance for translation and language understanding tasks inspired the vision community to explore their applicability in the domain of visual recognition. The early significant work, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale,"

demonstrated that, when trained with large data, a vision transformer could indeed surpass traditional convolution-based methods in standard benchmarks like ImageNet [2]. This study opened a new paradigm where visual tasks could be envisaged through a lens of global context awareness, stepping beyond the local inductive biases that place inherent constraints on CNNs.

Beyond addressing single-domain tasks, Vision Transformers possess an intrinsic flexibility to be adapted across multimodal learning frameworks. Studies indicate their symbiotic integration with other data modalities like text and audio, pushing the boundaries further in capturing data complexity for tasks such as visual question answering and cross-modal retrieval [5], [6]. Moreover, the lessened reliance on handcrafted feature extraction significantly lowers the bar for model development and permits Vision Transformers to be employed effectively with data derived from varying sensory inputs, yielding a new horizon of intelligent and unified architectural designs.

The key to the transformative impact of Vision Transformers lies in their capacity to model long-range dependencies without the rigidity of a predefined kernel or architecture. This is facilitated primarily through the self-attention mechanism, which redefined how spatial and sequential data could be interrelated and processed [7]. Compared to CNNs, which often require extensive data to learn local ambiguities through small receptive fields in initial layers, Vision Transformers leverage self-attention to possess a broader field of vision and learn global feature dependencies from the onset [8].

Despite these advantages, the journey of Vision Transformers is accompanied by hurdles related to computational efficiency and robustness against data variation and noise.

The quadratic complexity of the self-attention mechanism can introduce significant overhead, especially for high-resolution images or dense prediction tasks like object detection and semantic segmentation [9]. Research efforts have proposed architectural optimizations like sparse attention mechanisms, dimension reduction, and hybrid models integrating convolutions to mitigate these issues and bolster computational tractability [10], [11].

Moreover, when tasked with robustness and performance consistency across diverse visual datasets, Vision Transformers may wrestle with overfitting, as evidenced in small dataset scenarios where inductive biases and generalization capabilities remain limited [12]. Therein lies the appeal of hybrid architectures that combine the detail-capturing finesse of CNNs with the expansive context understanding of Transformers, offering a balance that capitalizes on the synergies between spatial precision and global context [13].

In synthesis, the essential narrative of Vision Transformers acknowledges their profound potential in advancing computer vision, counterbalanced by present computational constraints and their need for optimization in smaller datasets and varied application domains [14]. Future research directions are poised towards achieving architectural efficiencies, enhancing positional encoding strategies for diverse data types, and exploiting domain-specific adaptations to maximize Vision Transformers' utility across emerging fields like autonomous systems and medical imaging [15].

In essence, Vision Transformers signify not just an incremental improvement but a conceptual shift in the treatment of visual data processing. As ongoing research addresses current limitations and broadens their scope of application through innovations in computational methodologies and hybrid integrations, Vision Transformers are well-positioned to redefine cutting-edge practices in computer vision [16]. This ongoing evolution not only promises enhanced performance and scalability for a variety of vision tasks but also paves the path for more pervasive and robust AI-driven vision systems in the near future.

## 2 ARCHITECTURAL FOUNDATIONS OF VISION TRANSFORMERS

### 2.1 Self-Attention Mechanism and its Evolution

The self-attention mechanism stands as the foundational element of Vision Transformers, marking a paradigm shift in how visual data is processed by deep learning architectures. Originating from its pivotal role in natural language processing, self-attention provides the ability to model global relationships across tokenized data inputs, making it ideal for capturing long-range dependencies in visual information [1]. This mechanism evaluates the inter-dependencies across all parts of a given input, assigning varying importance to different elements, enabling the dynamic identification of salient features irrespective of spatial hierarchy.

Fundamentally, self-attention computes an output based on a set of queries, keys, and values derived from the input data. Given an input $X$, the self-attention mechanism constructs projections into query $Q = XW^Q$, key $K = XW^K$, and value $V = XW^V$ spaces using trainable weight matrices $W^Q, W^K, W^V$. The attention scores are calculated as a weighted sum of values, with weights determined by the similarity of queries and keys, specifically employing a scaled dot-product attention mechanism:

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

where $d_k$ is the dimension of the key vector, serving to mitigate the effect of vanishing gradients in the softmax function by scaling.

A notable advantage of self-attention over traditional convolutional architectures is its ability to capture non-local, global contextual dependencies at scale. While convolutional neural networks inherently excel in modeling localized features through fixed receptive fields, their reliance on multiple stacked layers for global context assimilation is both computationally intensive and suboptimal for large-scale images [1]. Self-attention, conversely, offers a direct and efficient pathway to capture both local and global feature interactions in a single operation.

Emergence of self-attention variants has been pivotal in addressing the mechanism's primary limitation—quadratic complexity with respect to the sequence length. Techniques such as sparse and deformable attention have been introduced to reduce these computational burdens by limiting the attention calculations to a relevant subset of tokens [11]. Sparse attention mechanisms, such as introduced in various transformer adaptations, prune the attention spans, focusing on the most informative tokens, effectively reducing unnecessary computations and enhancing scalability [11].

Moreover, innovations like the Focal Transformers leverage coarse-to-fine granularity by introducing focal attention, where each token attends closely to its immediate neighbors and distant tokens in a hierarchical scheme, allowing the model to efficiently encapsulate both short and long-range dependencies [17]. These approaches contribute to substantial reductions in computational complexity while maintaining or even enhancing performance on high-resolution vision tasks.

The evolution of the self-attention mechanism has also sparked the integration with and enhancement by convolutional designs, resulting in hybrid architectures such as the Convolutional Vision Transformer (CvT) and cross-scale models like CrossFormer [10], [18]. These architectures incorporate the spatial hierarchy potential of convolutions, promoting locality and robust feature extraction alongside the global feature synthesis aided by self-attention.

Emerging trends highlight a growing shift towards modular and flexible attention architectures, capable of adapting computation dynamically to different portions of the input space. For instance, techniques like modular attention, which partition the attention process and focus computational resources dynamically based on the input complexity, show promise in improving trade-offs between efficiency and accuracy [19].

Future developments in self-attention are likely to focus on further reducing complexity, enhancing efficiency, and integrating advanced inductive biases to improve generalization and performance on diverse and data-scarce tasks, aligning with efforts to bridge performance gaps on small

datasets [8]. As research continues, particularly in extending these mechanisms to applications beyond classical image recognition—such as point cloud and multi-modal data—their adaptability and robustness will be vital in driving forward capabilities in spatial understanding and autonomous systems.

Overall, the self-attention mechanism within Vision Transformers not only redefines architectural capabilities in visual recognition but also sets the stage for continued innovation in efficiently managing high-dimensional, complex data inputs across diverse and evolving application landscapes.

## 2.2 Core Design and Components of Vision Transformers

The core design and components of Vision Transformers (ViTs) play an essential role in defining their unique operational paradigm, setting them apart from the traditionally dominant convolutional neural networks (CNNs) in computer vision. Central to the Vision Transformers' approach is a sequence-processing framework that treats an image as a sequence of patches rather than as a grid of pixels. This approach introduces distinctive architectural elements—patch embeddings, positional encoding, and feed-forward networks—that collectively contribute to the Vision Transformer's ability to effectively learn representations across sequences. Each component presents its own mechanisms that equip Vision Transformers with robust and expressive power while also presenting distinct challenges and design considerations.

The process of patch embedding marks the inception of transforming an image into a format suitable for transformer processing. Essentially, this involves dividing an image into small, fixed-size patches, akin to non-overlapping square tiles [20]. Each patch is linearly mapped to a vector that acts as the input token for the transformer, simulating the interpretation of words in language models and enabling the use of transformer architecture [21]. This method allows ViTs to handle non-square images and process them in parallel, treating each patch as an independent element. Despite the advantages, this paradigm introduces challenges such as extensive computation and a potentially large memory requirement due to the detailed information extracted from each patch.

Introducing explicit spatial encoding becomes necessary due to the transformer's lack of inherent order-awareness, as the original transformer architecture was designed for natural language processing, where position is implied by token order. In vision tasks, positional encoding is crucial to capturing spatial relations among patches. Positional encoding strategies such as absolute and relative encodings have evolved to integrate spatial information into patch embeddings [22]. Absolute positional encoding assigns a unique vector to each position, while relative positional encoding emphasizes positional differences between tokens [23]. Although effective, these strategies present computational and representational trade-offs in preserving spatial integrity.

At the architectural core of Vision Transformers are transformer blocks that integrate multi-head self-attention mechanisms with feed-forward networks (FFNs). Conventional FFNs serve as the mortar for the Vision Transformer's capability in learning intricate representations. Typically comprising two linear layers with an intervening non-linearity, FFNs demonstrate an impressive capacity for feature transformation [24]. However, the computational load and parameter count can be substantial, posing scalability challenges when addressing high-resolution images or video sequences [25].

Recent research efforts aim to enhance the efficacy and scalability of FFNs. Incorporating convolutions into the feed-forward stack is one way to counteract the loss of locality inherent in deep networks [26]. Convolutional operations help in embedding spatial hierarchies, providing a natural means to incorporate image-like features that pure linear layers might otherwise overlook.

Despite the promising capabilities demonstrated by Vision Transformers across various tasks in computer vision, the trade-off still leans heavily towards a requirement for extensive computational resources. Emerging trends, such as dynamic sparse attention and hierarchical attention, seek to balance this by optimizing computational resources without compromising the quality of learned representations. Dynamic attention mechanisms, for example, adaptively focus on the most pertinent image areas, reducing unnecessary computations observed in various transformer adaptations [27].

In pursuit of further advancements, integration with other neural models to harness their strengths, such as CNNs' locality-capturing ability, could lead to the development of more efficient and powerful hybrid models [28].

In summary, while Vision Transformers herald a paradigm shift in visual information processing, the ongoing challenge is to strike a balance between computational efficiency and representation richness. Addressing these challenges will likely involve developing novel architectures that seamlessly integrate domain-specific knowledge and resources. Future research directions could include exploring multi-modal transformers that merge vision with other sensory data, leveraging the transformative potential established by the ViT framework, thus laying a foundation that goes beyond conventional vision models [29].

## 2.3 Positional Encoding Strategies

Understanding the fundamental role of positional encoding in Vision Transformers (ViTs) is pivotal for leveraging their full computational potential across diverse tasks. Originally designed for sequence data in natural language processing, positional encoding strategies enable these models to effectively process unordered data by embedding an order or structure into input sequences. This is particularly critical in computer vision, where the inherent spatial relationships between pixels or patches form the bedrock of visual understanding.

Traditional convolutional neural networks (CNNs) inherently encode positional information through their architecture, where spatial locality is embedded within the kernels. In contrast, ViTs treat images as sets of non-sequential patches, necessitating a mechanism to encode spatial relationships. The two predominant strategies for positional

encoding in Vision Transformers are absolute and relative positional encodings, each with distinct characteristics and trade-offs.

Absolute positional encoding assigns a unique position to each token in a sequence, typically using sinusoidal functions or learnable embeddings. These techniques aim to encode positional information independently of the input data, thus eliminating the influence of input content on the positional context. As observed in early ViT applications, this approach largely stems from the Transformer models used in NLP [2].

One of the pioneering methods, sinusoidal embeddings, leverages periodic functions to encode position. This technique offers an elegant, non-learnable method that imbues the model with inductive biases for position, such as translation invariance. Despite its mathematical simplicity, sinusoidal encoding may not fully exploit data-driven learning capacities.

Other iterations, such as learnable absolute positional encodings, optimize positional vectors as part of the training process. This modality can potentially adapt more dynamically to the input data, capturing nuances specific to datasets [30]. However, it also increases the risk of overfitting and decreased generalization if the training set is not representative. The flexibility of learnable encodings exemplifies their use in complex visual tasks, but this can come at the cost of computational efficiency and model interpretability [10].

While absolute encodings define a fixed reference for all tokens, relative positional encoding focuses on the relationship between tokens. This method has proven particularly beneficial for visual tasks that demand a nuanced understanding of local spatial exchanges. By encoding the distance or direction between tokens rather than fixed positions, relative encodings help models maintain spatial hierarchies and contextual relevance [31].

Emergent frameworks, such as the Rotary Position Embedding (RoPE), integrate these principles by allowing the rotational invariance of positional information, thus enhancing the model's capacity to capture geometric transformations [32]. The decomposed relative positional embeddings employed in the improved Multiscale Vision Transformers model further demonstrate this strategy's flexibility and its ability to sustain performance across diverse visual recognition tasks [33].

Despite the advantages of relative positional encoding, challenges persist, particularly in the scalable training of large ViT models. The complexity of modeling pairwise token relationships can lead to increased computational overhead, requiring innovations such as sparse and low-rank attention to control resources [17].

The choice of positional encoding strategy imparts significant impact on the model's adaptation to diverse visual contexts. It directly influences the ability of ViTs to generalize across tasks with varying resolutions and data distributions. As such, researchers are exploring hybrid and adaptable approaches that incorporate elements of both absolute and relative encoding methods. These might leverage conditional encodings that dynamically adjust based on input features or task requirements [1].

The ongoing development of more robust and efficient encodings continues to address the pressing need for scalable ViTs that can deliver consistent results in both low and high-resource settings. Advances might include the integration of convolutional layers within ViT architectures to naturally embed local positional biases while retaining global attention benefits [10], [34].

Future directions in positional encoding for ViTs are likely to focus on enabling finer-grained spatial understanding while minimizing computational costs. Techniques that enhance adaptability and interpretability without sacrificing performance will be pivotal. This includes the design of encoding methods that dynamically accommodate different data modalities and scales.

Expanding the spatial reasoning capabilities of ViTs through innovative encoding schemes will allow these models to transcend current limitations, achieving breakthroughs in complex visual understanding tasks such as medical imaging, autonomous systems, and beyond. As the field progresses, combining principles from both traditional machine learning paradigms and cutting-edge transformer designs promises to offer new vistas of possibility in computer vision research.

## 2.4 Architectural Enhancements for Efficiency and Scalability

The architectural landscape of Vision Transformers (ViTs) has evolved significantly, with researchers actively innovating to enhance efficiency and scalability while maintaining performance. Given their potential for superior outcomes across diverse tasks, addressing ViTs' inherent computational demands is crucial for widespread adoption. This subsection explores architectural modifications aimed at overcoming these challenges, spotlighting techniques for reducing computational complexity and enhancing scalability without sacrificing accuracy.

Traditional ViTs encounter considerable computational overheads due to the intrinsic complexity of self-attention mechanisms, which have a quadratic time complexity relative to the number of input tokens. This makes them especially taxing for high-resolution images. To tackle this issue, researchers have proposed various modifications to optimize the efficiency of self-attention. One such approach is local self-attention mechanisms that restrict attention calculations to local neighborhoods, rather than the entire image. This approach not only reduces computational load but also enhances scalability by efficiently computing attention maps for larger images [35]. These modified architectures have shown notable performance improvements in image classification tasks, maintaining expressiveness while operating efficiently on high-dimensional data.

Beyond attention optimization, patch-wise processing has emerged as a viable strategy to boost ViT efficiency. By processing image patches independently before integrating them within the attention framework, this method reduces the effective resolution processed by transformer layers. Patches are encoded with rich local features using smaller, less resource-intensive networks like CNNs, then combined by the transformer, which focuses on capturing global dependencies. This hierarchical approach is reminiscent of classical image processing techniques and aids in resource management, allowing ViTs to scale gracefully with input size [31].

Significant improvements in ViT efficiency also stem from model compression techniques such as pruning and quantization. Structured pruning involves removing redundant blocks or attention heads, significantly cutting computational costs while maintaining core functionalities. Unstructured pruning, which eliminates individual parameters, provides a finer degree of compression, yielding lightweight architectures for deployment on resource-constrained devices [36]. Fixed-point quantization has similarly been employed to optimize computational requirements by reducing numerical precision, preserving model accuracy during inference.

Emerging trends incorporate sparse and low-rank approximations within attention mechanisms. By introducing sparsity in attention maps or factorizing them into low-rank matrices, these techniques diminish the need for full pairwise computations, significantly reducing overheads [37]. Such approximations balance preserved expressiveness of global attention with minimized computational demands.

These architectural adaptations underscore an evolving understanding of trade-offs inherent in Vision Transformer design. While innovations like local attention and patch-based processing address computational challenges, they may introduce limitations in model generalization and robustness. For instance, local attention centers on smaller regions, which may underperform on tasks requiring comprehension across widely dispersed image regions. Further exploration of these trade-offs may involve dynamic architectures capable of toggling between detailed local processing and comprehensive global analysis based on task requirements.

Prospective directions involve multicore and distributed processing strategies to expand ViT scalability. By scaling both computational capacity and model size through efficient task distribution across GPUs and TPUs, these strategies leverage emerging hardware capabilities effectively. Additionally, training methodologies focused on data efficiency could compound the benefits of ViTs' architectural scalability, as seen in neural network frameworks for multimodal transformers.

In summary, advancements in ViT architectural enhancements emphasize the delicate balance of performance, efficiency, and scalability. Innovations such as local attention, patch-wise processing, and model compression strategically address these challenges, positioning Vision Transformers to extend their impact across computational domains. Ongoing research and convergence of software and hardware optimizations will be vital for realizing ViTs' full potential, facilitating a broader spectrum of practical applications. Collectively, these efforts highlight a promising trajectory toward more efficient, scalable, and versatile Vision Transformers, poised to redefine future computational models in computer vision.

## 2.5 Integration of Vision Transformers with Other Architectures

The intersection of Vision Transformers (ViTs) with other architectural paradigms has given rise to hybrid models that efficiently capture both local and global visual information, leveraging the strengths of distinct methodologies. This strategic integration addresses some inherent limitations of standalone architectures while enriching their capabilities. Here, we explore how Vision Transformers are combined with other architecture types, specifically convolutional neural networks (CNNs) and dynamic/sparse architectures, to achieve improved performance in various visual tasks.

Hybridization with convolutional layers is a prominent strategy, often termed as CNN-Transformer hybrids. These approaches capitalize on the local feature extraction capabilities of CNNs, which, paired with the global context modeling of Transformers, yield powerful representations for visual data. A pivotal work in this regard is Bottleneck Transformers, which integrate self-attention into the latter stages of a ResNet, transforming bottleneck blocks into Transformer-like structures while enhancing accuracy with reduced parameters [38]. This synergy not only improves the spatial precision in tasks such as detection and segmentation but also mitigates the Transformers' heavy computation.

In CNN-Transformer hybrids, the convolutional layers act to preserve locality biases inherent to image data. This integration is realized either by appending Transformer layers over existing CNN stacks or by replacing certain CNN layers with Transformer ones, depending on the target application and computational constraints [39]. When confronting high-resolution images, these hybrids can exploit local features efficiently, thus reducing the computational overhead typical of pixel-based global attention mechanisms.

Dynamic and sparse architecture integrations further enhance the efficiency of Vision Transformers. Dynamic architectures, such as those using dynamic convolutions, adapt their computation based on input data characteristics. This adaptability is key in handling diverse visual inputs without fixed computational budgets. For example, models that use dynamic token sparsification can prune non-essential tokens based on attention scores, thus accelerating processing without significant performance loss in tasks like image classification or object detection [40].

Sparse architectures extend this concept by introducing mechanisms to reduce redundancy during both training and inference phases. Notably, techniques like soft token pruning selectively exclude less important tokens, which minimizes computation while attempting to preserve essential image features [41]. These approaches are particularly beneficial in resource-constrained environments where efficiency is crucial.

Emerging trends in hybrid models focus on leveraging neural architecture search (NAS) to systematically explore the vast design space, seeking optimal configurations of CNN-Transformer hybrids. For instance, GLiT employs a locality module in its search space, favoring the trade-off between global and local data processing [42]. Such techniques significantly automate and optimize the integration process, producing architectures with superior performance and computational efficiency tailored to specific tasks.

The integration of Vision Transformers with other architectures encounters challenges mostly related to balancing computational cost against performance gains. Transformer layers, despite their high representational capacity, introduce considerable complexity. Hence, their strategic

placement within a hybrid architecture must be carefully considered, focusing Transformer use on parts of the model where global context is most beneficial [43].

In conclusion, integrating Vision Transformers with other architectures unlocks significant potential in visual processing tasks, providing a robust framework that benefits from both local and global information synthesis. Future directions point towards even more sophisticated hybrid models incorporating advancements in NAS to fully exploit the complementary strengths of CNNs and ViTs. Expanding these explorations further into multimodal domains, where vision models intersect with textual or auditory data, could define the next frontier for Vision Transformer applications. Engaging with these synergies holds the promise of more versatile, efficient, and powerful vision systems, capable of tackling increasingly complex real-world tasks.

## 3 ARCHITECTURAL VARIANTS AND ENHANCE-MENTS

### 3.1 Hybrid Model Integration

The integration of convolutional neural networks (CNNs) with Vision Transformers (ViTs) has emerged as a formidable strategy in computer vision, marrying the local feature extraction prowess of CNNs with the global contextual capture strength of ViTs. This hybrid approach seeks to exploit the strengths of both architectures while mitigating their individual limitations. CNNs have long been celebrated for their ability to efficiently extract fine-grained local features from images due to their inherent locality biases. In contrast, ViTs excel at modeling long-range dependencies and global context via self-attention mechanisms, a capability that is somewhat restricted in traditional CNNs due to the limited receptive field size. Thus, the amalgamation of these architectures is particularly advantageous, offering a comprehensive representation that can optimize performance across various vision tasks.

One prevalent approach to hybrid model integration is the incorporation of convolutional layers into ViT architectures. The Convolutional Vision Transformer (CvT) is a prime example, which introduces a hierarchical architecture that integrates convolutional token embedding and convolutional projection within transformer blocks [10]. This design enriches the transformer with convolutional CNN properties such as shift-invariance and localized feature extraction while retaining the dynamic attention capabilities of ViTs. The CvT achieves superior performance compared to both pure ViTs and traditional CNNs, demonstrating the potential of this hybrid approach.

Conversely, some hybrid models incorporate transformer blocks into CNN frameworks. The Convolution-enhanced image Transformer (CeiT) exemplifies this methodology by embedding transformer mechanisms into a predominantly convolutional architecture [44]. CeiT employs an Image-to-Tokens module for low-level feature extraction and a Locally-enhanced Feed-Forward network to promote spatial correlations, capitalizing on both architectures' capabilities. This model has shown increased performance on image classification benchmarks without needing extensive datasets, a testament to the synergy obtained through hybridization.

Hybrid models also explore different arrangements of CNN and ViT layers to capitalize on multiscale feature learning. The Multiscale Vision Transformers (MViT) extend this idea by using a pyramidal structure that progressively expands channel dimensions while reducing spatial resolution, allowing different scales of contextual information to be captured at distinct stages [45]. This multiscale approach is particularly effective in tasks like video recognition, where understanding spatial and temporal dimensions simultaneously is crucial.

Despite these promising results, the integration of CNNs and ViTs is not without challenges. One major issue is the increased computational complexity and model size resulting from combining two distinct architectures. Solutions such as dynamic token sparsification in Vision Transformers, which progressively prune tokens to maintain computational efficiency, offer pathways to balance performance and computation [11].

Furthermore, the design and optimization of hybrid architectures necessitate careful consideration of module interaction and sequencing. For instance, ensuring a seamless transition from convolutions to transformers—and vice versa—requires innovative engineering to facilitate effective feature fusion and propagation [13]. This involves selecting appropriate fusion strategies, such as cross-attention mechanisms, that allow different feature sets to be combined without redundancy or loss of detail.

There is also growing interest in exploring how learned representations from CNNs and ViTs can complement each other across varied applications, ranging from medical imaging to autonomous systems. Specifically, Vision Transformers' robustness in modeling long-range dependencies and global interactions complements CNNs' detailed spatial representations. Thus, hybrid models can be tailored to meet the unique requirements of specific tasks, such as precise boundary detection in segmentation tasks or efficient temporal context capture in video analysis.

In terms of future directions, research is increasingly focused on developing smarter, more adaptive fusion techniques that allow hybrid architectures to dynamically adjust their configuration based on input characteristics or computational constraints. This could involve the introduction of adaptive control units that selectively activate different pathways—convolutional or transformer-based—at different stages of processing, effectively merging the flexibility of transformers with the efficiency of CNNs.

In summary, the integration of CNNs with Vision Transformers presents a compelling advancement in enhancing the versatility and capability of vision models. By combining local discrimination with global integration, hybrid architectures promise to set new benchmarks in accuracy and efficiency across a myriad of computer vision tasks. However, the pursuit of optimized architectures demands continued innovation in design principles, computational strategies, and application-specific adaptations, suggesting a productive avenue for ongoing research and development in the field.

### 3.2 Token and Layer Optimization Strategies

In the pursuit of enhancing Vision Transformers' computational efficiency and scalability, token and layer opti-

mization strategies have emerged as critical focal points. These methodologies not only address the computational bottleneck but also ensure that model performance remains uncompromised. This subsection delves into various approaches, such as token pruning, merging, and layer optimization, offering insights into their mechanisms, advantages, and challenges.

A fundamental step in optimizing Vision Transformers is the management of tokens, the basic units of input data that the transformer processes. Token pruning involves selecting the most informative parts of an input sequence, thereby reducing computational load while maintaining essential features. The approach recognizes that not every token contributes equally to the decision-making process of a model. Techniques like attention score-based pruning assess the importance of tokens based on attention weights, discarding those with negligible influence. This method effectively lessens the model's computation without a significant hit to accuracy, as token significance rather than uniformity drives the pruning process [46].

Another strategy is token merging, which consolidates similar tokens into single representations. This process streamlines the data processed by the model, reducing redundancy and enhancing efficiency. The merging is typically guided by criteria such as spatial similarity or similar feature representations among tokens. Models employing this technique can maintain competitive performance levels while operating with a lower memory footprint. The introduction of k-NN attention, which selectively attends to the nearest neighbor tokens, exemplifies this strategy. It naturally introduces a local bias without the explicit use of convolution operations, thus reducing complexity while preserving performance [47].

Layer optimization presents another crucial avenue for improving efficiency. Layer-wise pruning techniques identify and remove less significant components within layers, such as certain attention heads or entire layers, based on their contribution to the model's output. Structured pruning focuses on entire sub-modules, while unstructured pruning targets individual parameters, aiming to preserve model architecture while significantly lowering resource demands [48]. Furthermore, innovations like adaptive layer scaling enable dynamic adjustment of layer operations depending on the input or its complexity, optimizing resource allocation during model inference.

Employing hybrid and hierarchical models is a practical approach to optimization, integrating the strengths of other architectures like CNNs with transformers. These models benefit from the locality of CNNs and the global feature extraction mechanism of transformers. Techniques such as Convolutional Self-Attention (CSA) enrich low-level features, while Recursive Atrous Self-Attention (RASA) enhances high-level feature extraction, demonstrating the impact of hybrid structures on efficiency and scalability [49].

Equally important are emerging trends around reducing the computational burden at the attention mechanism level. Approaches such as sparse and low-rank approximations of the attention matrix cut down the quadratic complexity traditionally associated with self-attention. Sparse attention selects critical parts of input data for focus, reducing the number of operations required, while low-rank approxima-

tions simplify matrix operations by projecting them into lower-dimensional spaces [50]. Recent innovations like Performer leverage Fast Attention Via Orthogonal Random (FAVOR) features, reducing the complexity to linearity, thus supporting the model's deployment in real-time applications without compromising on long-range reliance [51].

The challenges associated with these optimization strategies entail maintaining a delicate balance between computational efficiency and model accuracy. Pruning and merging can overly simplify representations, risking the omission of critical data points, particularly in scenarios involving high-resolution inputs. Moreover, while methods such as sparse attention optimization strive to cut computational costs, they may inadvertently impair the model's capacity to differentiate intricate patterns, illustrating the challenge of aligning simplicity with expressive power [52].

Future directions in token and layer optimization are promising and require navigating these complexities with further refinement of existing methods and innovative design of new architectures. Incorporating learning-based approaches that automatically determine the importance and interdependencies of tokens and layers could lead to more robust implementations. Additionally, real-time adaptation of computational strategies based on dynamic inputs could enhance scalability across different deployment settings without substantial accuracy trade-offs.

In conclusion, token and layer optimization represents a pivotal area in advancing Vision Transformers. Substantial progress has been made through pruning, merging, and strategic layer management, yet ongoing research is necessary to overcome challenges related to accuracy trade-offs and adaptability. As the field advances, these optimizations will be key to realizing Vision Transformers' full potential in practical and resource-constrained environments, paving the way for their broader adoption and utility in diverse application domains.

### 3.3 Efficient Attention Mechanisms

Efficient attention mechanisms in Vision Transformers (ViTs) aim to balance the expansive computational demands of self-attention with the need for sustainable model deployment across a variety of hardware constraints. This subsection delves into innovative strategies targeting computational efficiency, building on the foundational concept of self-attention while reducing overheads through architectural adaptations and approximation methods.

Traditionally, the self-attention mechanism, a cornerstone of the Vision Transformer, computes pairwise interactions among all tokens, leading to quadratic complexity in relation to the input size. This poses significant challenges, particularly for high-resolution images where the number of tokens becomes substantial. To address this, several compelling approaches have emerged that either sparsify the computation or employ strategic approximations to maintain performance while alleviating computational burdens.

One of the prominent methods to enhance efficiency is the incorporation of sparse attention mechanisms. Sparse attention strategically selects which token interactions to compute, effectively reducing the overall complexity. This is achieved by focusing on key interactions while disregarding

less salient ones. Researchers have proposed models such as the SCRAM framework, which leverages spatial coherence and sparsity to achieve an O(n log(n)) complexity from the traditional O(n²) [53]. By utilizing randomized algorithms like PatchMatch for correspondence, it enables a significant reduction in the number of necessary computations, while maintaining the global context effectively.

Another approach involves low-rank approximations. These techniques optimize the attention matrix by identifying the low-dimensional subspaces that encapsulate most of the correlation between tokens. For instance, the Dilate-Former model integrates Multi-Scale Dilated Attention to efficiently model local interactions across sliding windows of features before employing global attention stages [54]. This method capitalizes on the observation that long-range dependencies in early layers are frequently redundant, thus prioritizing computational resources for more critical stages in the network. Such strategies leverage the inherent redundancy in visual data for computational economy without sacrificing performance.

Approximate and localized attention variants represent another significant trend. These methods restrict attention calculations to local neighborhoods or employ various approximations to relax the burden of global attention computation. The idea is to combine the locality bias, advantageous in convolutional networks, with the global context-modeling capacity of transformers. For example, Focal Self-Attention introduces a mechanism where each token attends to its closest surroundings at a higher granularity while considering distant tokens more coarsely [17]. This nuanced attention allows capturing both fine and coarse details efficiently, making it particularly useful for tasks demanding both detailed understanding and broader context.

A comparative examination of these approaches highlights distinct strengths and trade-offs. Sparse attention methods are adept at significantly reducing computation time and resources by leveraging the inherent sparsity in visual data. However, they may require careful parameter tuning to strike a balance between computational gain and performance drop-off. Low-rank approximations, while powerful, introduce additional complexity in terms of matrix factorization and require sophisticated engineering to maintain stability. Conversely, localized attention approaches align more naturally with existing convolutional paradigms and are relatively easier to integrate into hybrid architectures, offering a more accessible entry point for practical deployment.

Emerging trends show a convergence of techniques that blend multiple strategies for optimal efficiency. Hybrid models that merge convolutional operations with attention mechanisms are increasingly adopted to bring down computational requirements further while enhancing model robustness [10]. These hybrids permit efficient modulatory architectures that can easily adjust to varying levels of detail required by different tasks.

Despite these advancements, challenges remain. One of the primary hurdles is achieving seamless integration of efficient attention mechanisms without degrading the quality of learned representations. As models continue to scale, ensuring consistent performance across diverse data scales and domains demands ongoing innovation. More-over, the interpretability and adaptability of such efficient mechanisms to new and unseen data distributions continue to be areas ripe for exploration.

Moving forward, future research could explore deeper integration of data-driven approaches to dynamically adjust attention computation based on input characteristics or task requirements. For instance, employing meta-learning strategies could potentially enable models to learn optimal attention distribution patterns during training, thus reducing computational loads during inference.

In conclusion, efficient attention mechanisms in Vision Transformers reflect a critical push towards making these powerful models accessible and applicable across a broader spectrum of real-world applications. By harnessing the synergies between sparse computation, low-rank approximations, and localized attention, these advancements not only promise to democratize access to cutting-edge vision capabilities but also drive the evolution of more intelligent, responsive, and resource-efficient computational models.

## 3.4 Resolution and Scalability Adaptations

In recent years, Vision Transformers (ViTs) have revolutionized the field of computer vision by capturing global dependencies in visual data. However, their performance can significantly vary across different input resolutions and computational scales, necessitating adaptations that enhance their scalability and maintain their efficacy. This subsection delves into these adaptations, emphasizing techniques that enable ViTs to dynamically adjust to varying input sizes and ensure robustness across resolutions.

ViTs operate on fixed-size image patches, processing inputs at a predefined resolution, which can limit scalability and hinder performance when image sizes fluctuate. Several innovative approaches have emerged to address these limitations. One notable method is the development of resolution-adjustable architectures, allowing models to adapt their processing pipeline based on the input resolution dynamically. The introduction of Naive Dynamic Resolution (NDR) in the Qwen2-VL series exemplifies significant progress in this area. This mechanism supports the processing of variable-resolution images by adjusting the number of visual tokens accordingly, offering a practical strategy for efficient visual representation across resolutions [55].

Another critical aspect of scaling ViTs across different input resolutions is maintaining effective positional encoding. Traditional positional encodings assume a static resolution, which can degrade performance when inputs differ from those seen during training. To address this, Conditional Positional Encodings (CPE) have been introduced, providing dynamically generated encodings that maintain performance across varying resolutions. CPE is conditioned on the local neighborhood of input tokens, ensuring robust positional information against resolution changes [30].

In addition, the application of Rotary Position Embedding (RoPE) extends its effectiveness from language models to vision tasks within ViTs. RoPE maintains precision and enables resolution adjustments while preserving spatial relationships, enhancing the scalability and robustness of ViTs across datasets like ImageNet, COCO, and ADE20k [56].

Lightweight attention mechanisms further contribute to scalability. Hybrid Sparse Attention and local attention mechanisms ensure computational feasibility for high-resolution inputs by prioritizing locality. This is analogous to Vicinity Attention methods where patches are given stronger attention based on proximity, efficiently managing complexity while preserving crucial data relationships. Castling-ViT exemplifies this approach by retaining both global and local context using adjusted linear-angular attention during inference, reducing computational demands without sacrificing accuracy [57], [58].

As these approaches are adopted, several trade-offs emerge. Dynamic resolution processing enhances adaptability but can increase model complexity and training time as networks accommodate broader data variability. Similarly, incorporating flexible positional encodings like CPE and RoPE demands careful architectural tuning to prevent computational overhead. Lightweight attention mechanisms, though reducing demands, might trade off detailed global interactions, affecting tasks relying on entire scene comprehension.

Emerging trends in resolution and scalability adaptations focus on generalization across diverse datasets by integrating positional and dynamic resolution adaptation with efficient computational strategies. Examples include TokenPacker, which reduces redundant tokens while optimizing computational load, facilitating seamless transitions across variable data inputs [59]. Future developments must balance the efficiency of these adaptations with the growing demand for real-time processing in edge devices.

In conclusion, the evolution of resolution and scalability adaptations in ViTs marks a significant progression toward universal visual processing. Key advancements, such as dynamic resolution adjustment and adaptive positional encoding, extend the versatility and robustness of these models across resolutions, enhancing their applicability in varied real-world scenarios. As research converges on integrating these strategies with efficient attention mechanisms, their impact is poised to expand from high-performance computational environments to resource-constrained edge scenarios, ushering in a new era of adaptable and efficient computer vision solutions.

### 3.5 Pruning and Quantization Techniques

In the realm of Vision Transformers (ViTs), pruning and quantization techniques offer pivotal strategies for model compression, enabling their deployment on resource-constrained devices without sacrificing substantial performance. These methodologies aim to tackle the inherent computational burden posed by Vision Transformers, characterized by high parameter counts and intricate self-attention mechanisms. This subsection delves into the various methodologies, empirical results, and theoretical underpinnings associated with pruning and quantization approaches applied to Vision Transformers.

Pruning in the context of Vision Transformers can be categorized into two primary types: structured and unstructured pruning. Structured pruning typically involves the removal of entire attention heads or even complete layers, thus offering a reduction in computation and memory usage that aligns closely with the architectural layout of the model. Unstructured pruning, on the other hand, removes individual weights based on predefined criteria, leading to a potentially more computationally optimal but irregular network structure. For instance, the work of [60] explores an integrated approach that focuses on training sparse subnetworks within Vision Transformers, dynamically identifying unimportant weights during training, which results in a sparse but efficient network.

Structured pruning methods have shown promising results in maintaining model integrity while reducing computational costs [61]. The structured approach is particularly appealing for real-time applications, given it often results in a smaller, albeit more structured, computational graph that hardware accelerators can exploit for increased efficiency. These pruning strategies underscore the importance of maintaining a balance between model size reduction and preserving the core functionality of Vision Transformers, ensuring that the degradation in performance is minimal despite significant reductions in model parameters.

Quantization, another paramount method in model compression strategies, reduces the precision of computations performed within Vision Transformers, hence reducing both memory footprint and computational load. This is often achieved through fixed-point arithmetic, converting high-precision weights and activation values into lower-bit representations. Techniques such as post-training quantization and quantization-aware training have been applied successfully to ViTs, with frameworks such as [62] introducing sophisticated techniques that reframe scale parameters to effectively optimize for low-bit quantization without inducing significant accuracy drops.

The advantages of quantization are further amplified when they are paired with hardware-aware optimizations. The paper [63] demonstrates how the combination of quantization with hardware constraints can yield efficient Vision Transformers that meet specific power and performance metrics relevant for deployment in edge devices. By customizing the quantization strategy to leverage the specific strengths of a target hardware platform, it is possible to develop models that not only perform well in terms of accuracy but also adhere to stringent latency and power consumption requirements.

However, there are challenges that underpin the use of pruning and quantization. A core issue in the utilization of these techniques lies in the trade-off between model accuracy and computational savings—a primary area where novel methodologies continue to focus their efforts. For instance, the complexity introduced by non-uniform pruning and bit-level variations in quantization requires complex retraining or fine-tuning procedures, impacting the ease and speed of deploying these models in practical settings. Techniques that incorporate adaptive learning rates and dynamic quantization levels have shown potential in mitigating these concerns by allowing for continuous adaptation of the network's pruned and quantized states without significant manual intervention.

Emerging trends in this domain are pushing towards the integration of pruning and quantization into end-to-end learning pipelines. This approach aims to unify model reduction techniques with model training, thus optimizing the model's entire lifecycle from initialization to deployment.

Methods that employ co-design strategies [64] consolidate these disparate techniques, ensuring that compression occurs simultaneously with model training, backed by a detailed understanding of the interactions between hardware efficiency and machine learning outcomes.

Looking forward, the focus on evolving these compression methods will likely hone in on developing automated frameworks for pruning and quantization that do not involve manual hyperparameter tuning or extensive validation on new datasets. With sophisticated profiling and optimization tools, such methods promise to automate the delicate balancing act of optimizing performance while constrained by the computational abilities of deployment platforms. Moreover, as Vision Transformers continue to expand into more domains, the ability to easily adapt these models to new tasks while maintaining compact models will be a significant avenue of research and practical exploration.

In summary, the synergetic application of pruning and quantization holds significant promise in enabling Vision Transformers to transcend the current limitations of computational demand and energy use in constrained environments. Through continued innovation in methodologies and integration with advanced deployment strategies, Vision Transformers can achieve widespread utility without compromising their renowned effectiveness in visual tasks.

## 4 TRAINING AND OPTIMIZATION STRATEGIES

### 4.1 Pre-training Strategies

Pre-training strategies represent one of the most pivotal areas of innovation in the development and optimization of Vision Transformers (ViTs). The underlying philosophy of these strategies is to leverage vast amounts of data to instill in the model a foundational set of representations or understandings before fine-tuning for more specific tasks. This subsection delves into the breadth of unsupervised and self-supervised techniques that define the pre-training landscape for Vision Transformers, aiming to enhance model generalization and efficiency.

The early success of Vision Transformers can largely be attributed to robust pre-training methodologies on expansive datasets such as ImageNet or JFT-300M, which provided vivid representations that later fine-tuning stages could refine for specific tasks [8]. This parallels the pre-training techniques applied in natural language processing (NLP), where models benefit from being exposed to a wide corpus before task-specific tuning [14]. Consequently, pre-training facilitates improved sample efficiency and mitigates the risks involved with specialized datasets, such as overfitting or limited generalization [1].

A prominent self-supervised learning paradigm that has gained traction is Masked Image Modeling (MIM). Inspired by Masked Language Modeling in NLP, MIM trains models by hiding parts of the input data and challenging the model to reconstruct the missing pieces. This effectively forces the network to develop a nuanced understanding of context and structure in visual data, which enhances its ability to understand and generate complex patterns post-pre-training [65]. This is substantiated by empirical studies which emphasize the role of pre-training in stabilizing training processes and reducing the need for large labeled datasets [12].

Another significant strategy involves leveraging transfer learning from the domain of NLP. The idea of borrowing methodologies like sequential token training and incorporating textual data into Vision Transformer pre-training has shown promising results [66]. Notably, the transformational capabilities of these linguistic techniques in establishing a semantic understanding transfer smoothly to visual models, resulting in enhanced model adaptability and insight formation [67]. These cross-domain transfers signify a broader trend towards hybrid learning systems that amalgamate insights from diverse model architectures.

Furthermore, domain-specific pre-training is highlighted as a strategy that tailors the foundational knowledge of Vision Transformers to specialized fields like medical imaging or autonomous driving. Domain-specific pre-training entails customizing the training datasets and objectives to cater to the peculiarities inherent in a particular field. The efficacy and applicability of this approach become evident when reflecting on tasks requiring high-resolution image analysis and context-aware decision-making processes [15].

Despite the advantages, these pre-training strategies are not without challenges. The computational costs and time investments associated with pre-training on large datasets remain significant barriers. Efforts are being made to address these through innovative methods such as efficient architectural designs that maintain accuracy with reduced computational requirements [9]. Scalable architectures like the Multi-Scale Vision Transformer hint at new pathways for creating models that economize on resource usage while enhancing multi-resolution feature extraction [45].

Additionally, there are ongoing discussions and research efforts aimed at improving the robustness and stability of pre-training processes. The model's ability to handle distribution shifts and adversarial perturbations are among the chief concerns, promoting investigations into the resilience characteristics embedded during the pre-training stages [68].

Looking forward, the trajectory of pre-training strategies appears to be oriented towards more adaptable and versatile frameworks. There's potential in exploring meta-learning paradigms, where models self-adjust pre-training objectives based on the characteristics of incoming data [69]. Additionally, the integration of multi-modal data sources — such as incorporating auditory, text, and sensory inputs into visual models — offers a fertile area for creating richly interconnected neural networks capable of handling highly intricate and varied information streams [32].

In summary, pre-training strategies for Vision Transformers are integral to establishing powerful, generalizable models that excel with minimum task-specific labeled data. While challenges in computational demands and scalability persist, continuous advances in architectural innovations, cross-modal integration, and domain-specific training promise to broaden the scope and efficiency of these models. As this subsection outlines, the future of pre-training strategies lies in creating synergistic frameworks that not only bridge the gap between varied data streams but also optimize learning processes in resource-constrained environments.

## 4.2 Fine-tuning and Transfer Learning

In the realm of Vision Transformers (ViTs), fine-tuning and transfer learning have emerged as essential strategies for adapting pre-trained models to specialized tasks. This subsection delves into effective methodologies that harness pre-existing models, enhancing the adaptability and performance of Vision Transformers across distinct domains. By building on core principles from pre-training strategies and setting the stage for subsequent optimization techniques, we explore methods for achieving parameter efficiency, domain alignment, and cross-modal interactions, each contributing to the broader applicability of ViTs in diverse environments.

Achieving parameter efficiency is crucial for optimizing ViTs for new tasks without the computational burden of full model retraining. Techniques like adapter layers and prompt tuning are at the forefront of this effort. Adapter layers integrate lightweight, trainable modules into the pre-trained network, allowing for refinement with task-specific data while minimizing alterations to the core model and conserving computational resources. Similarly, prompt tuning involves appending task-specific prompts to input tokens, guiding the model's attention without altering the main architecture. These approaches highlight the transition towards modular learning, enabling rapid adaptation to new contexts while mitigating computational and memory demands.

Domain adaptation plays a pivotal role in transferring Vision Transformers to environments with varying data distributions. Techniques such as feature alignment and adversarial training aim to bridge the gap between source and target domains. Feature alignment fine-tunes the model to minimize discrepancies in feature representations, enhancing generalization. Adversarial training leverages adversarial examples to fortify models against data distribution shifts, bolstering robustness and adaptability. Together, these methods address the challenge of domain discrepancies, promoting robust performance in dynamic environments.

Cross-modal learning expands the versatility of Vision Transformers by integrating diverse data modalities, supporting complex tasks such as visual question answering and multimodal sentiment analysis. By utilizing models pre-trained on textual data, semantic insights are infused into visual tasks, enhancing comprehension. Key strategies involve leveraging multimodal embeddings and designing architectures to effectively fuse visual and textual features [70]. This integration heralds an era where Vision Transformers excel not only in visual recognition but also in bridging semantic gaps across disparate data sources.

While offering substantial benefits in efficiency and domain adaptability, these strategies present certain limitations. Adapter layers and prompt tuning may introduce biases if not adequately calibrated to the target domain. Domain adaptation methods like adversarial training enhance robustness but add training complexity and necessitate comprehensive dataset alignments [23]. Cross-modal frameworks require sophisticated architectures to balance modality-specific nuances, maintaining coherent integration [27].

To succinctly describe parameter-efficient learning adap-

tations, consider augmenting the ViT model $\mathbf{M}$ with an adapter module $A(\cdot)$:

$$\mathbf{M}_{\text{new}}(\mathbf{x}) = \mathbf{M}(\mathbf{x}) + A(\mathbf{x}) \tag{2}$$

Here, $A(\cdot)$ is a parametrized function optimized for the target domain, contributing minimally to the overall parameter count. In contrast, domain adaptation with feature alignment seeks to minimize distance $d$ between source $S$ and target $T$ feature spaces:

$$\min_{\theta_S, \theta_T} d(\mathcal{F}_S(\mathbf{X}_S; \theta_S), \mathcal{F}_T(\mathbf{X}_T; \theta_T)) \tag{3}$$

By treating domain-specific variations as optimization constraints, these methodologies ensure Vision Transformers remain robust yet flexible across diverse applications.

The sophistication of fine-tuning and transfer learning strategies is instrumental in extending Vision Transformers' utility. Emerging trends in zero-shot learning and self-supervised domains demonstrate potential for reducing data dependency, advancing ViTs' scalability. Despite challenges in maintaining model integrity against domain shifts, ongoing innovations promise further refinements. Future explorations may focus on nuanced domain adaptation, leveraging computational paradigms like federated learning to minimize data centralization [25].

In conclusion, fine-tuning and transfer learning significantly enhance Vision Transformers' adaptability, bridging pre-trained efficiencies and specialized requirements. These strategic innovations foster a new era of flexibility in deep learning, seamlessly transitioning to optimization methodologies in subsequent discussions.

## 4.3 Optimization Techniques

As Vision Transformers (ViTs) increasingly become the architecture of choice for many computer vision applications, optimizing their performance, efficiency, and scalability remains a significant focus of research. This subsection elaborates on the optimization techniques employed to refine Vision Transformers, detailing both computational and architectural strategies that push the boundaries of what these models can achieve.

Optimization in Vision Transformers involves a multi-faceted approach, prioritizing both computational efficiency and architectural robustness to ensure the models can handle large-scale data while being deployable in resource-constrained environments. At the heart of this effort is the endeavor to balance the inherent complexity of self-attention mechanisms with practical usability in real-world applications.

Initially, strategies such as sparse computation have been instrumental in making self-attention more scalable. Techniques like SCRAM (Spatially Coherent Randomized Attention Maps) have been pivotal in reducing the quadratic computational complexity traditionally associated with self-attention to a more manageable $O(n \log n)$ complexity [53]. By exploiting spatial coherence and sparser representations, such techniques streamline computation while preserving the attention mechanism's core benefits.

Efficient layer designs form another critical component of optimizing Vision Transformers. Convolutional enhance-

ments, such as those in the CvT model, incorporate convolutions directly within Vision Transformer architectures to ameliorate scale and positional bias challenges [10]. This blending of convolutional efficiency with transformer expressiveness results in models that not only preserve locality but also process global context more cost-effectively. Furthermore, integrating depth-wise convolutions can significantly reduce computational requirements while maintaining the fidelity of feature representations [21].

Pruning techniques, which refine model structures by excising redundant parameters, play a vital role in improving operational efficiency. Methods like structured and unstructured pruning significantly compress ViTs without severely impacting performance. By selectively removing less impactful weights or entire attention heads, these techniques allow a reduction in both memory usage and inference times, paving the way for deployment on edge devices.

Quantization represents another strand of optimization by lowering numerical precision during computations, thus reducing the model's size and energy requirements. Using techniques like the piece-wise linear approximation for GELU within the PEANO-ViT highlights how function-specific shortcuts can drive substantial improvements in power efficiency without compromising accuracy [71]. Similarly, data-free quantization approaches, as demonstrated by PSAQ-ViT, leverage unique properties inherent in Vision Transformers like patch similarity to calibrate quantization parameters effectively [72].

Dynamic adaptation techniques also enhance Vision Transformers' efficiency by adjusting computational strategies based on input variability. For instance, frameworks like CF-ViT employ coarse-to-fine strategies to focus computational resources on the most informative regions of an input image, thus balancing computational load [73]. Additionally, approaches such as token pruning reorganize image patches dynamically during inference, ensuring computation focuses on the most relevant information [74].

Emerging trends indicate a growing interest in transforming the Vision Transformer architecture to leverage hardware constraints effectively. Innovations like LookupViT emphasize token compression to reduce computational overhead while maintaining or even improving model performance across different data domains. Similarly, the development of reversible architectures like Reversible Vision Transformers illustrates a shift towards memory-efficient design paradigms that don't sacrifice model accuracy [75].

From an analytical perspective, optimizing Vision Transformers entails grappling with a triad of key challenges: scalability, practicality, and performance retention. The integration of sparse and efficient computational designs presents trade-offs between precision and resource allocation, needing careful calibration to maintain accuracy across diversified applications. Experimental results consistently suggest the viability of these techniques, as vision tasks on datasets such as ImageNet and COCO consistently exhibit improved performance metrics with reduced computational demands [2].

Looking forward, the optimization landscape for Vision Transformers is poised to advance with increasing sophistication. Research directions hint at hybrid methods that merge principles from varied domains, including state space models and attention alternatives, to create even more scalable solutions [70]. Furthermore, these endeavors align with the broader agenda of sustainable AI, advocating for methods that minimize environmental footprint while maximizing computational throughput.

In conclusion, the optimization techniques for Vision Transformers, as outlined, signify a concerted push towards systems that are as ecologically conscious as they are effective. With ongoing explorations into cross-disciplinary methodologies and hardware-constrained innovations, the development of Vision Transformers continues to evolve, promising refinements that could redefine their role within the AI landscape.

## 4.4 Advanced Training Methods

The advanced training methods for Vision Transformers (ViTs) are pivotal in enhancing their robustness, generalizability, and overall performance. Building on the optimization and fine-tuning methodologies discussed in previous sections, this subsection delves into innovative techniques employed in the training phase of Vision Transformers. It emphasizes data augmentation, regularization strategies, and meta-learning approaches, equipping these models to excel across diverse visual recognition tasks.

Data augmentation is essential for improving the diversity of training data and enhancing model generalization in machine learning models. In the realm of ViTs, emergent methods such as Mixup and CutMix are gaining traction. These techniques involve mixing pixels from two images or segments within images to create composite training samples, effectively enhancing data distribution and potentially reducing overfitting. However, the complexity of optimal mixing parameter settings underscores a key challenge. Moreover, geometric augmentations, like those employed in "Axial Attention in Multidimensional Transformers" [35], utilize image warping and shifting operations to retain spatial context, strengthening the model's capability to learn spatial invariances. These transformations are particularly beneficial in modeling varied viewpoints in multimodal tasks [76].

Regularization is another cornerstone of advanced ViT training. Adversarial training stands out as a robust regularization strategy that improves model resilience against perturbations, thus enhancing robustness to adversarial attacks. By introducing perturbed inputs during training, this approach fosters an extensive exploration of the input space. Nonetheless, balancing clean and adversarial samples while managing computational overhead remains a challenge. Equally noteworthy is high-frequency component augmentation, which involves adding noise or signals to emphasize certain frequency bands in the data, preventing the model from focusing excessively on specific input aspects that may not generalize well.

Meta-learning and few-shot learning contribute significantly to training Vision Transformers, enabling rapid learning of new tasks with minimal data. An innovative approach within this domain involves temporally adaptive modules and synthetic gradient networks, which propose optimization gradients to expedite convergence and enhance adapt-

ability to new data environments. By leveraging a meta-learner that optimizes task-agnostic parameters, ViTs generalize beyond initial training distributions, as evidenced by models attaining state-of-the-art performance on task adaptation benchmarks [77]. As these approaches evolve, managing computational costs and efficiently transitioning from coarse to fine-grained tasks remain crucial research challenges.

The synergy between augmentation, regularization, and meta-learning introduces novel learning landscapes aimed at achieving a synergistic effect. Reinforcement learning paradigms integrate with ViTs, utilizing reward signals based on performance to fine-tune and dynamically adjust learning procedures. This aligns model training with task-specific goals and optimizes decision-making strategies for long-term benefits, differing from static training processes.

Emerging trends include the implementation of hybrid transformer architectures that fuse components from various models to capitalize on their respective strengths. Conditional training modules that switch between transformer layers and convolutional operations based on input characteristics illustrate a promising path toward scalability and enhanced resource efficiency [78]. The challenge in these methodologies lies in managing the complexity of hybridization while ensuring streamlined interactions among disparate components.

Attention-based mechanisms across different ViT layers are crucial for enhancing model performance. Methods such as rotary and cross-attention encodings generalize spatial structure across heterogeneous datasets, offering robustness in tasks like multimodal and spatio-temporal analyses. While attention mechanisms are computationally demanding, ongoing innovations aim to reduce this burden, enhancing scalability and real-world deployment [79].

In conclusion, as Vision Transformers continue to evolve, advanced training methods are driving significant strides in their efficiency and robustness. Future directions should explore adaptive learning principles more extensively across diverse input modalities and environments, potentially leveraging unsupervised relationship embeddings and conditional motifs to streamline pre-training processes. This ongoing work is anticipated to uncover new paradigms for improving ViT's conceptual and semantic understanding, broadening their applicability across complex visual tasks.

## 4.5 Fine-tuning and Robustness Improvement

The continuing evolution of Vision Transformers (ViTs) underscores the necessity for effective fine-tuning strategies that amplify their adaptability to varied and task-specific environments while ensuring robustness against diverse challenges. Fine-tuning pre-trained Vision Transformers is a critical avenue for leveraging existing models' capacities, allowing them to be tailored to new tasks without compromising the overall structural robustness. This subsection delves into the framework and methodologies essential for fine-tuning Vision Transformers, along with strategies aimed at enhancing their resilience against various adversarial and operational challenges.

Fine-tuning involves adjusting a pre-trained Vision Transformer's parameters to adapt to specific tasks, thus necessitating minimal task-specific data. The primary objective remains to balance model optimization for a particular domain while ensuring the preservation of generalized robustness. Methods like parameter-efficient fine-tuning approaches which leverage adapter layers can significantly impact this balance. These methods introduce additional parameters into certain layers to guide task-specific learning without overwriting the original transferable features, ensuring computational and memory efficiency. This approach proves particularly effective when applied to scenarios with limited training data, as seen in studies centered on parameter-efficient transfer learning techniques [80].

Additionally, domain adaptation techniques serve to counter distribution shifts between the pre-training and fine-tuning phases. Domain adaptation techniques, such as adversarial training approaches, are essential for adjusting models to new distributions that may not be encountered during pre-training. These methods improve the stability of Transformer models across different tasks and environments; frameworks like DAE-Former provide a robust mechanism incorporating cross-attention for improved segmentation tasks [81]. Here, domain adaptation not only bridges the gap between different data distributions but also serves to bolster model robustness against potential adversarial challenges.

Rapid task adaptation, another focal point of ViT fine-tuning, can benefit significantly from few-shot learning techniques. Few-shot learning emphasizes swift adaptation to new tasks with minimal data and has demonstrated potential in specific Vision Transformer architectures by using meta-learning frameworks for cross-task learning dynamics [82]. Such frameworks equip models with capabilities to leverage previously learned knowledge to accelerate learning in new environments, thereby enhancing both task adaptability and robustness.

A pivotal concern during fine-tuning is maintaining robustness against adversarial attacks, ubiquitous challenges in deploying Vision Transformers reliably. Various methodologies have emerged to address adversarial robustness, concentrating on making Vision Transformers less susceptible to subtle changes that might deteriorate their predictive accuracy. Techniques such as adversarial training, which involves training models on adversarially perturbed samples, have shown promising results [83]. Adversarial training helps models anticipate and resist changes that might otherwise compromise their functionality, ensuring high performance in pragmatic deployment scenarios.

The interpretability of Vision Transformers also merits attention in the realm of fine-tuning and robustness enhancement. Transparency in decision-making processes not only fosters trust but also aids in isolating and rectifying inefficiencies within the model's operational framework. Mechanisms for attention visualization, such as visual attention saliency maps, have underscored the importance of understanding the decision-making pathways of models during inference [64]. These tools can elucidate which parts of an input image most influence the model's predictions, providing insights that are crucial for further enhancement and refinement of model robustness.

Looking to the future, adaptive fine-tuning techniques are an emerging trend, fostering both robustness and adapt-

ability in Vision Transformers. These methods involve dynamic adjustments to model parameters in real-time, allowing transformations in response to contextual or environmental changes—a capability increasingly desirable as models are deployed in varied, real-time applications. Moreover, continuous learning approaches that integrate lifelong learning principles can significantly enhance a ViT's adaptability and resilience [40]. By enabling models to learn constantly from new inputs without forgetting previously obtained knowledge, they can remain both robust and highly adaptable over time.

In conclusion, the symbiosis between fine-tuning strategies and robustness improvements in Vision Transformers involves a multi-faceted approach, simultaneously ensuring optimal task-specific tuning while mitigating vulnerabilities to adversarial and functional challenges. This subsection highlighted key methodologies and emerging trends, advocating for a comprehensive integration of adaptive, interpretability, and adversarial resilience strategies into fine-tuning processes. As applications of Vision Transformers continue to expand, ongoing research must persist towards developing innovative methodologies that enhance their robustness, ensuring reliable performance across an ever-broadening array of tasks and conditions. Future studies are anticipated to concentrate on hybrid models that synergize aspects of robustness enhancement with fine-tuning, highlighting collaborative strategies for sustaining high adaptability and resilience in complex environments. [60] illustrates the tantalizing possibilities of this direction, suggesting that a robust understanding and implementation of fine-tuning parameters can configure Vision Transformers for even more nuanced and complex applications.

## 5 APPLICATIONS OF VISION TRANSFORMERS IN COMPUTER VISION

### 5.1 Image Classification and Detection

Vision Transformers (ViTs) have emerged as a groundbreaking advancement in computer vision, achieving significant strides in image classification and detection tasks. This subsection explores how ViTs enhance these tasks by leveraging their capability to capture global context and long-range dependencies. The analysis includes a comparative evaluation of various methodologies, presents critical insights into their strengths and limitations, and highlights potential future directions within the domain.

Traditional Convolutional Neural Networks (CNNs) have predominantly shaped the landscape of image classification and detection for decades, owing to their strong inductive biases towards locality and translation invariance. However, these characteristics are both a strength and a limitation, as CNNs inherently struggle to model global context without resorting to resource-intensive architectures. The introduction of Vision Transformers has provided a paradigm shift by capitalizing on self-attention mechanisms to process visual data, offering a holistic perspective that enhances both the accuracy and efficiency of image analysis.

Vision Transformers can model the entire image by dividing it into non-overlapping patches, subsequently treated as tokens akin to words in natural language processing. This conceptual alignment allows the self-attention mechanism within ViTs to consider interactions among all parts of an image simultaneously, thus capturing comprehensive contextual information. This ability is crucial in tasks where holistic understanding supersedes local feature extraction. Dosovitskiy et al. [2] demonstrated the remarkable performance of ViTs over state-of-the-art CNNs in image classification tasks across several benchmarks, such as ImageNet, by aptly leveraging this attribute.

Furthermore, the ViT model's impressive generalization ability from pre-training on massive datasets, followed by fine-tuning on domain-specific tasks, has been documented across multiple studies [1]. The Tokens-to-Token Vision Transformer (T2T-ViT), for example, enhances efficiency by progressively aggregating neighboring tokens into more structured representations, which further optimizes the balance between global contextual awareness and computational demands [66].

In object detection, ViTs exploit their long-range dependency modeling to accurately localize and identify multiple objects within a single image. The self-attention mechanism enables the mapping of complex spatial relationships and hierarchical representation, which are pivotal in distinguishing between closely situated objects or objects with overlapping characteristics. The Multiscale Vision Transformers (MViTs) leverage this by integrating multiscale feature hierarchies that are adept at differentiating between lower-level details and high-level semantic information [45].

These Transformers have also been integrated into hybrid architectures to address remaining gaps in capturing fine local details. Convolutional Vision Transformers (CvT) and Convolutional Neural Networks Meet Vision Transformers (CMT) showcase how incorporating convolutions into the Transformer pipeline can enhance performance on object detection tasks by harnessing the scale and shift invariance properties of CNNs while maintaining the global contextual awareness of ViTs [10], [13]. Such hybridizations systematically exploit the strengths of both paradigms and have set new performance records, particularly in detecting objects from high-resolution images, an area where pure ViTs initially lagged due to their dependency on extensive computational resources.

While the advantages of ViTs are compelling, they are not without challenges. A primary concern is their substantial computational and memory overhead, attributed to the quadratic complexity of self-attention, which escalates with increased input resolution [9]. This aspect necessitates innovations like sparse attention mechanisms and token sparsification to bring down computational cost while retaining accuracy [11].

On a different front, ViTs exhibit a deficiency in localization inductive bias compared to CNNs, which inherently limits their performance on smaller datasets without massive pre-training [8], [12]. Strategies such as locality self-attention and shifted patch tokenization have been developed to imbue ViTs with the requisite local feature learning capabilities, improving their efficacy on smaller datasets [8], [21].

Looking forward, the continued hybridization of Transformers with CNN elements promises to bridge gaps in their limitations. Moreover, leveraging unsupervised and self-supervised learning frameworks can mitigate the data

efficiency issue, enabling ViTs to learn robust representations with fewer annotations [1], [65].

There is also an increasing inclination towards developing dynamic and scalable architectures that adjust according to task-specific requirements, a trend observed in adaptive token pruning and context-aware attentional mechanisms [9], [11]. Further research into enhancing interpretability and optimizing computational pathways will be essential in reducing deployment costs and widening the accessibility of ViTs in real-time applications.

In conclusion, Vision Transformers have substantially enriched image classification and detection paradigms by providing enhanced global context understanding and object interrelation modeling. Their innovative integration into hybrid models and adaptation to various architectural advances continue to propel their capability in addressing complex visual tasks. As the community ventures into refining these models towards efficiency and interpretability, ViTs are poised to play a pivotal role in advancing the frontier of computer vision technologies.

## 5.2 Advanced Vision Tasks

Vision Transformers (ViTs) have significantly advanced the realm of computer vision, extending their influence beyond traditional tasks like image classification and object detection. They now address more complex vision tasks, including semantic segmentation, video analysis, and three-dimensional (3D) vision. These tasks present unique challenges and opportunities, requiring sophisticated approaches to effectively extract and utilize contextual and temporal information. This subsection delves into the application of ViTs in these advanced vision tasks, examining specific advancements, challenges, and future directions.

Semantic segmentation is a critical vision task that involves assigning a class label to every pixel in an image, requiring a comprehensive understanding of spatial hierarchies and contextual relationships. Traditional convolutional neural networks (CNNs) often struggle with capturing global context due to their local receptive fields. In contrast, Vision Transformers leverage self-attention mechanisms to achieve superior context modeling. Techniques such as Convolutional Self-Attention (CSA) enhance low-level features, providing intricate segmentation outputs with fewer parameters [49]. Moreover, the integration of Gabor filter-inspired focal attention mechanisms in models like FViTs demonstrates improved performance in dense prediction tasks by focusing on feature representations across various scales and orientations [34]. However, a significant challenge in applying ViTs to semantic segmentation is the high computational cost associated with the quadratic complexity of self-attention. To address this, efficient attention mechanisms like those in the Focal Transformer balance fine-grained local and coarse-grained global interactions to optimize computational efficiency [17].

In transitioning to video analysis, Vision Transformers are especially well-suited to handling the temporal dynamics intrinsic to video data. The Temporal Patch Shift (TPS) method is a noteworthy innovation, shifting patches in the temporal dimension to convert spatial attention into spatiotemporal attention with minimal computational overhead [84]. This method excels at tasks like action recognition, where capturing temporal sequences is crucial. Similarly, advanced architectures such as the Neighborhood Attention Transformer leverage a scalable sliding-window attention mechanism, ensuring linear complexity and enhancing capability in long-sequence modeling [27]. Despite these advancements, video analysis still faces challenges in efficiently modeling long-range temporal dependencies without incurring prohibitive computational costs. Addressing these concerns, methods like Separable Self-Attention significantly reduce complexity by using element-wise operations, making them ideal for resource-constrained video tasks [85].

In the domain of three-dimensional vision, ViTs have shown promise in tasks like 3D object recognition and reconstruction, which require nuanced integration of spatial and depth information across multiple viewpoints. The deformable self-attention module, as seen in models like Deformable Attention Transformers, selects key positions in a data-dependent manner, efficiently focusing on relevant spatial regions for enhanced feature extraction [20]. This approach effectively combines the benefits of global attention with local flexibility, which is essential for comprehensively understanding 3D spaces.

However, deploying these advanced models across different vision tasks introduces several challenges. Computational demands and memory overhead are major hurdles, particularly for high-resolution and real-time applications. Innovative solutions, such as the Flowformer, address these challenges by leveraging flow networks to linearize attention mechanisms, achieving scalability without sacrificing expressiveness [86]. Furthermore, maintaining interpretability while enhancing performance is an ongoing challenge. Methods to elucidate ViTs' decision-making, such as visualization techniques for cross-task evaluation, are crucial for fostering trust in automated systems [68].

Looking ahead, the potential of Vision Transformers in advanced vision tasks is promising. The evolution of scalable and efficient attention mechanisms, coupled with robust training and fine-tuning strategies, will likely facilitate broader adoption and integration in real-world applications. Future research may focus on developing hybrid architectures that seamlessly integrate the inductive biases of CNNs with the contextual modeling strengths of Transformers, similar to advancements seen in ConViT [26]. Additionally, expanding the application of Vision Transformers in domains such as autonomous systems and environmental monitoring will necessitate continuous improvements in latency reduction, robustness against adversarial attacks, and adaptability to diverse environmental conditions.

In conclusion, Vision Transformers offer a powerful framework for addressing the complexities of semantic segmentation, video analysis, and three-dimensional vision. While they have successfully redefined possibilities in these tasks, ongoing research must address computational and interpretability challenges to realize their full potential in practical applications.

## 5.3 Cross-Domain Adaptations

Vision Transformers (ViTs) have notably transformed the landscape of computer vision by offering a robust frame-

work for tasks across varying domains. The ability of ViTs to excel in diverse applications is primarily facilitated by their capacity for domain adaptation and transfer learning, which enables them to leverage learned representations across different datasets and tasks. This subsection examines the various methodologies utilized by Vision Transformers to perform cross-domain adaptations effectively, highlighting the advancements, challenges, and potential trajectories for future research.

Fundamentally, domain adaptation in Vision Transformers leverages the distinct capability of these models to transfer knowledge across domains, mitigating the challenges posed by domain discrepancies. Unsupervised domain adaptation, in particular, plays a pivotal role where labeled data in the target domain is scarce or unavailable, facilitating ViTs to adapt using source domain knowledge [1]. The approach typically involves pre-training the ViT on a source domain with abundant labeled data, followed by fine-tuning on the target domain using domain adaptation techniques. By utilizing strategies like adversarial domain adaptation and domain-invariant feature learning, ViTs can minimize distribution shifts between source and target domains [87].

Transfer learning further amplifies the cross-domain adaptability of Vision Transformers, enabling models trained on extensive datasets like ImageNet to excel on smaller, domain-specific datasets. This is achieved through fine-tuning, where pre-trained ViTs adjust their learned weights to accommodate the nuances of the target domain [12]. Transfer learning addresses the data scarcity problem by reducing the requirement for large-scale labeled datasets in the target domain, thus enhancing the applicability of ViTs in varied domains such as medical imaging and remote sensing [16].

The adaptability of Vision Transformers is also bolstered by architectural innovations. Hybrid architectures, incorporating convolutional elements into Vision Transformers, introduce locality biases that enhance the model's performance on domain-specific tasks by capturing both local and global contextual features [88]. This integration mitigates challenges associated with pure-transformer architectures, such as the lack of inductive biases, which are critical in domains like object detection and segmentation where precise spatial localization is necessary [44].

Emerging trends in cross-domain adaptations of Vision Transformers emphasize the role of self-supervised learning (SSL) and few-shot learning paradigms. SSL techniques, such as contrastive learning and masked image modeling, allow Vision Transformers to learn robust feature representations without reliance on labeled data [89]. These methods are particularly advantageous for domain adaptation as they enable the model to internalize features that are generalizable across domains, thereby bolstering transferability [32].

Moreover, Vision Transformers leverage few-shot and zero-shot learning to adapt to new domains with minimal data. Few-shot learning frameworks, facilitated by ViTs, incorporate novel methods like prototype learning that guides the model to identify and transfer task-agnostic features efficiently [90]. This approach significantly reduces the dependency on large annotated datasets and enhances the scalability of Vision Transformers for cross-domain adaptations [12].

Despite these advancements, cross-domain adaptations in Vision Transformers come with challenges that warrant further investigation. The computational and memory demands of ViTs pose a significant challenge when deploying them for domain adaptation on resource-constrained devices or scenarios requiring real-time performance [91]. Addressing these bottlenecks requires the development of lightweight and efficient architectures or the application of techniques like pruning and quantization without compromising performance [75].

Furthermore, maintaining interpretability and robustness across domains remains a critical challenge for Vision Transformers. Ensuring model decisions are transparent and reliable is vital, especially in safety-critical domains such as healthcare and autonomous systems. Research into enhancing the interpretability of ViTs using techniques like attention map visualization and causal explanation remains a promising direction [68].

In summary, the cross-domain adaptations of Vision Transformers stand at the frontier of facilitating their application across diverse and challenging environments. As the field progresses, future research will likely focus on enhancing the computational efficiency of ViTs, improving their adaptability through advanced SSL techniques, and ensuring their usability across low-resource scenarios. Additionally, continued exploration into hybrid architectures and their potential in specific domains could significantly impact the deployment of Vision Transformers in real-world applications [32].

### 5.4 Multimodal and Multitask Learning

Multimodal and multitask learning frameworks have witnessed remarkable advancement through the integration of Vision Transformers (ViTs), revolutionizing how these tasks are approached by leveraging their exceptional ability to model complex relationships across various modalities and tasks. This subsection delves into the transformative role of Vision Transformers in these domains, offering a nuanced analysis of their capabilities, challenges, and future potential.

Multimodal learning seeks to harness data from multiple sources or modalities, such as visual, textual, and auditory information, for more comprehensive data interpretation and improved task performance. Vision Transformers are particularly well-suited for this purpose due to their self-attention mechanism, which excels at capturing dependencies across diverse data types. For instance, in visual question answering, ViTs have been employed to simultaneously process images and accompanying text, thereby enhancing the integration of visual and linguistic information [92].

In this context, ViTs facilitate the incorporation of additional modalities, enhancing the performance of visual-linguistic models by offering a unified framework for cross-modal interactions. This approach allows for rich semantic understanding through the attention mechanism, which prioritizes relevant information across modalities, yielding enhanced results in tasks requiring spatial and semantic integration [93].

In multitask learning, Vision Transformers demonstrate the ability to handle multiple related tasks simultaneously by sharing architectural components, such as attention heads and feed-forward layers, across tasks. This shared learning exploits synergies among tasks, leading to improved efficiency and robustness. The multitask capabilities of ViTs are evident in frameworks like SeqTR, where a unified approach to various visual grounding tasks, such as phrase localization and segmentation, reduces architectural complexity while maintaining robust performance [93].

Such architectures optimize learning and inference by employing a single model for multiple tasks, conserving computational resources while facilitating knowledge transfer between tasks. This is especially important in contexts where computational efficiency and memory constraints are critical, such as deployment on edge devices or in real-time applications [94].

Compared to traditional convolutional neural networks (CNNs), ViTs offer distinct advantages in multimodal and multitask learning due to their flexibility and scalability with diverse data sources. CNNs often struggle with integrating long-range dependencies across modalities, whereas ViTs naturally facilitate this integration with their global receptive field. However, challenges remain, such as high computational costs and data efficiency issues, especially when working with high-dimensional data from multiple modalities [36].

Emerging trends involve developing more efficient attention mechanisms that retain critical information while reducing computational overhead. Techniques such as sparsity-inducing modifications and hybrid models that combine the strengths of transformers and CNNs are actively explored [94]. Additionally, integrating ViTs with state space models and leveraging their dynamic capabilities to manage multi-resolution data streams is a growing interest area [95].

Despite their robustness, ViTs face limitations in multimodal and multitask learning. Significant challenges include needing large-scale datasets to fully harness their potential, often a barrier in specialized domains with limited annotated data. Additionally, the interpretability of transformer-based models remains an ongoing concern, as their decision-making processes can be opaque, complicating diagnostics and refinement [96].

There are trade-offs in deploying ViTs, notably between model complexity and inference efficiency. Efforts to simplify architectures may reduce model performance, while more complex models demand significant computational resources [97].

Future research is poised to optimize Vision Transformers' architectural designs for better scalability and enhance data fusion techniques for improved multimodal learning outcomes. This could involve adaptive attention mechanisms that dynamically adjust based on input characteristics, enhancing the model's flexibility and applicability to new tasks [55].

Moreover, evolving capabilities in cross-modal generation and retrieval tasks highlight the importance of developing mechanisms that enhance interaction between modalities at all model levels [98]. Addressing these challenges can extend ViTs' applicability, paving the way for richer and more generalizable multimodal and multitask learning systems.

In conclusion, as Vision Transformers continue to redefine the approach to multimodal and multitask learning, their evolution will hinge on balancing performance improvements with the practical constraints of deployment, achieving a symbiotic relationship between model complexity and computational efficiency.

## 5.5 Medical Imaging Applications

Vision Transformers (ViTs) have emerged as promising tools in the domain of medical imaging, where precise diagnostic and analytical capabilities are crucial. The deployment of ViTs in medical imaging aims to harness their ability to model complex visual patterns, thus enhancing tasks such as diagnosis, segmentation, and image enhancement. This subsection delves into the transformative impact of ViTs in medical imaging applications, assessing various approaches, analyzing their strengths and limitations, and offering predictions on future advancements in this arena.

One of the primary advantages of deploying Vision Transformers in medical imaging is their proficiency in handling high-dimensional data and modeling intricate relationships within medical images. The self-attention mechanism, a cornerstone of ViTs, enables the examination of every pixel's relationships across an image, suiting tasks like segmentation, where understanding both local and global contexts is critical. For instance, the use of ViTs in segmenting anatomical structures often surpasses traditional convolutional neural networks (CNNs) by capturing finer details and underlying patterns that are crucial for clinical analysis. The Recurrent Attention Multi-scale Transformer (RAMS-Trans) exemplifies such capabilities by effectively learning discriminative region attention through multi-scale image patch integration [28].

The diagnostic accuracy afforded by ViTs is another significant contribution to medical imaging. ViTs excel at interpreting complex patterns in medical images such as MRIs or CT scans, facilitating early and more accurate diagnosis of conditions. The inherent capability of Vision Transformers to process high-resolution images supports the detection of subtle visual cues that may go unnoticed by human observation or less advanced models. Vision Transformers' adaptability to a vast range of imaging modalities further underlines their potential in enhancing diagnostic procedures across diverse medical fields.

ViTs have demonstrated particular effectiveness in medical image segmentation tasks, which require precise delineation of anatomical and pathological structures. Techniques like the Dual Attention-guided Efficient Transformer model (DAE-Former) have shown improved segmentation performance by leveraging reformulated self-attention mechanisms to capture spatial and channel relationships while maintaining computational efficiency [81]. Such models highlight the delicate trade-off between improving segmentation precision and maintaining efficiency, crucial for practical clinical applications.

However, the integration of ViTs in medical imaging is not without challenges. High computational and memory demands pose a substantial barrier to the widespread adoption of Vision Transformers in real-time clinical settings. To

tackle these issues, methods such as model pruning and quantization have been employed to compress ViT models without significantly compromising accuracy, allowing for their deployment on resource-constrained devices like edge computing platforms [99]. Additionally, optimization techniques like Neural Architecture Search (NAS) [42] are being explored to develop more adaptive and efficient ViT frameworks tailored to the specific needs of medical imaging tasks.

In considering the practical implications, the deployment of Vision Transformers in medical imaging presents an array of benefits that could revolutionize current diagnostic and analytical processes. Their ability to provide detailed segmentation can significantly aid in surgical planning, disease monitoring, and radiotherapy treatment. Moreover, the reduction in image processing times achieved by optimized ViT models can enhance workflow efficiency in medical facilities, offering faster patient evaluations and diagnoses.

Looking ahead, a notable direction for future research lies in developing hybrid models that combine the strengths of both transformers and CNNs, aiming to exploit local feature extraction capabilities alongside global attention modeling. Such integration could lead to comprehensive models adept at handling the unique challenges of medical imaging, providing robustness and enhanced interpretability. Moreover, exploring transfer learning techniques to adapt pre-trained Vision Transformers to specialized medical domains could facilitate more accurate and data-efficient deployments.

In summary, Vision Transformers stand at the forefront of innovative analysis in medical imaging, offering vast potential for redefining diagnostic paradigms and segmentation precision. While challenges related to computational resources and model efficiency remain, ongoing research and advancements in hybrid architectures, model compression, and adaptive learning strategies promise to further enhance the utility and performance of ViTs in medical settings. Thus, as Vision Transformers continue to evolve, their adoption in medical imaging applications is anticipated to grow substantially, paving the way for more accurate, efficient, and versatile diagnostic tools.

# 6 VISION TRANSFORMERS IN MULTIMODAL LEARNING

## 6.1 Fundamental Concepts in Multimodal Learning with Vision Transformers

The convergence of multimodal learning and Vision Transformers (ViTs) represents a significant evolution in the field of artificial intelligence, blending diverse data sources such as images, text, and audio to develop models capable of holistic understanding. Vision Transformers, with their self-attention mechanisms, are uniquely positioned to enhance multimodal learning due to their capability to capture long-range dependencies and facilitate robust cross-modal integration.

In the realm of multimodal learning, the self-attention mechanism, a cornerstone of Transformers, plays a pivotal role by allowing models to focus on relevant information across various modalities. This selective attention is crucial for effectively handling the intricate relationships that emerge when different data types interact. As highlighted by the foundational work on Transformers, which outlines the efficacy of self-attention in modeling dependencies between sequence elements [7], this mechanism undergirds the potential of ViTs in capturing complex interactions across modalities.

One of the critical challenges in multimodal learning with ViTs is generating comprehensive cross-modal representations that faithfully integrate information from disparate sources while preserving the unique characteristics of each modality. This is where cross-modal fusion techniques come into play, which consolidate distinct streams into a unified representation. Techniques leveraging self-attention to facilitate this process are crucial, as exemplified in studies showcasing ViTs' capability to model long-range interactions efficiently [1]. The inherent ability of ViTs to process sequential information and refine it through layers of attention facilitates the creation of nuanced cross-modal embeddings, which are indispensable for tasks that require an understanding of context across different data types.

Throughout these processes, another essential aspect is multimodal token fusion. This method, detailed in [5], involves dynamically selecting and integrating the most pertinent features from each modality. This approach not only optimizes the information retention across modalities but also mitigates the risk of diluting intra-modal details. By leveraging techniques like dynamic token sparsification [100], the Vision Transformer can effectively manage computational resources while maintaining comprehensive feature integration.

The advent of Vision Transformers has also accelerated advances in cross-modal retrieval and visual question answering (VQA). These tasks necessitate understanding and contextualizing visual data with linguistic inputs, a challenge that ViTs address with their ability to align and process multi-type data seamlessly. For example, when employed in VQA, Vision Transformers excel by mapping and attending to the pertinent visual features in response to textual queries, enabling enhanced context comprehension and facilitating nuanced responses [101].

However, the deployment of Vision Transformers in multimodal learning is not without challenges. A primary limitation is their inherent computational complexity, as the quadratic nature of the self-attention mechanism can be prohibitive when scaling to large datasets or high-dimensional input spaces [102]. Addressing this, the integration of hierarchical and sparse attention mechanisms has been suggested, which significantly reduces resource demands while maintaining model performance [21].

Moreover, the interpretability and transparency of Vision Transformers remain active areas of research. Understanding how these models make cross-modal inferences is crucial for adapting them to tasks where accountability and insight into decision processes are imperative. Progress in visualizing and interpreting attention weights [3] can provide valuable insights into the decision-making pathways of ViTs, potentially leading to more trustworthy and understandable multimodal AI systems.

Looking to the future, there are several promising directions for the use of Vision Transformers in multimodal learning. One such direction is the better integration of emerging

technologies, such as IoT and autonomous systems, where multimodal data is inherently collected and processed. Vision Transformers could serve as the core models for fusion and analysis in such systems, given their flexibility and scalability [9]. Additionally, research could explore extending ViTs to accommodate modalities beyond conventional text and visuals, such as sensor data or biological signals, thus broadening their application domains.

Equally important is the exploration of environmentally sustainable practices in deploying these models. Given the substantial computational resources required for training and inference in Vision Transformers [7], strategies for improving efficiency and reducing energy consumption will be crucial for their sustainable application.

In conclusion, Vision Transformers hold the promise of redefining multimodal learning by bridging multiple data types into a cohesive analytical framework. As the community continues to refine these models and address their limitations, Vision Transformers are poised to unlock new potentials in applications that demand an integrated understanding of the multifaceted world. With ongoing research, they stand to significantly impact a variety of fields, offering more robust, efficient, and intelligible multimodal learning systems.

## 6.2 Vision Transformers in Multimodal Tasks

Vision Transformers (ViTs) have increasingly demonstrated their versatility and effectiveness in transcending traditional unimodal boundaries, emerging as potent tools for various multimodal tasks. This exploration into multimodal scenarios underscores the potential of ViTs to synergistically integrate data from disparate sources, such as textual, auditory, and visual inputs. The subsequent analysis delves into their application across key tasks, drawing comparisons to conventional methods and highlighting emerging trends and future possibilities.

One area where ViTs have notably driven advancement is Visual Question Answering (VQA), which requires a nuanced understanding and integration of visual data with natural language questions. Historically dominated by convolutional neural networks (CNNs) and recurrent neural networks (RNNs), these approaches often struggled to efficiently capture intricate cross-modal dependencies. Vision Transformers, leveraging their self-attention mechanisms, have changed this landscape significantly. They enable explicit modeling of interactions between image regions and language tokens, enhancing VQA accuracy. Innovative strategies such as hierarchical alignment and multimodal fusion are setting new performance benchmarks.

Similarly, cross-modal retrieval tasks, which involve retrieving relevant images based on text queries, benefit from the inherent capabilities of Vision Transformers. The self-attention mechanism allows effective transformation of modality-specific representations into a unified latent space, facilitating improved retrieval accuracy. Unlike CNNs, which are constrained by fixed receptive fields, ViTs permit adaptive focusing of attention across modalities, thereby supporting more precise content alignment and retrieval.

In more complex domains, audio-visual segmentation tasks require concurrent processing of auditory and visual streams to effectively segment and classify multimedia content. Vision Transformers excel in this area by processing audio spectrograms alongside visual data within a unified framework, thus integrating cross-modal cues. This capacity extends beyond what is achievable with traditional convolutional methods, which often necessitate separate pathways for visual and auditory data, leading to more cohesive and accurate segmentation outcomes.

A key strength of Vision Transformers in these multimodal tasks lies in their ability to attend to pertinent features across different modalities, thanks to their robust self-attention mechanisms. ViTs dynamically prioritize relevant features across modalities, offering fine-grained control and synergy—attributes in which CNNs are inherently limited due to their deterministic convolutional layers. The resulting explicit feature interaction boosts task performance significantly.

Nevertheless, Vision Transformers face several challenges in multimodal tasks. Chief among these is the computational complexity associated with processing high-dimensional multimodal data, exacerbated by the quadratic scaling of self-attention operations. To address this, research is focusing on developing sparsification techniques and efficient attention mechanisms that reduce computation without sacrificing performance, as shown in the development of Sparse MLP networks and scalable architectures [103].

Another critical challenge is data scarcity across certain modalities, which hinders large-scale pretraining or fine-tuning. Transfer learning offers a viable solution, where ViTs pretrained on large datasets in specific modalities are adapted to new tasks with limited data availability. This approach mitigates the demand for extensive training datasets and enhances the generalizability of ViTs across tasks, broadening their applicability.

Looking to the future, significant opportunities exist in refining Vision Transformer architectures to further optimize integration capabilities. As domains expand—embracing fields such as augmented reality, autonomous systems, and interactive media—the potential of ViTs to synthesize complex multimodal data streams presents numerous research avenues. Advances in model interpretability and robustness will also be pivotal, ensuring that decisions made by ViTs in multimodal contexts remain transparent and reliable, especially in high-stakes areas like healthcare and autonomous driving.

Ultimately, the journey of Vision Transformers in multimodal tasks not only highlights their transformative impact but also reveals promising prospects and challenges. The evolution of techniques and methodologies will continue to shape this landscape, with Vision Transformers being central figures in the integration of diverse data modalities.

## 6.3 Challenges and Optimization in Vision Transformers for Multimodal Learning

Vision Transformers (ViTs) have established themselves as a formidable architecture in computer vision by leveraging attention mechanisms to model long-range dependencies across visual elements [104]. In the realm of multimodal learning, where the integration of diverse data types such as text, audio, and images is paramount, Vision Transformers

hold significant promise due to their ability to effectively process different modalities through a unified framework. However, several challenges emerge when applying Vision Transformers to multimodal contexts. This subsection delves into these hurdles and examines optimization strategies that can be utilized to address them.

One of the primary impediments to the application of Vision Transformers in multimodal learning is the computational overhead associated with processing large and diverse datasets. Vision Transformers inherently possess a quadratic computational complexity due to their reliance on self-attention mechanisms, which can become a bottleneck when dealing with high-resolution images and extensive sequences, especially in multitasking environments [2].

Various techniques have been proposed to mitigate this issue. Techniques such as sparse attention and low-rank approximations aim to reduce computational load by selectively attending to key elements within the data [17]. Sparse attention, for instance, limits the attention computation to a subset of crucial elements, thus decreasing the number of operations. Similarly, structured self-attention mechanisms leverage inherent patterns within key-query interactions to streamline processing [105].

Moreover, adopting model compression techniques such as pruning and quantization has shown promise in enhancing the scalability of Vision Transformers. Structured pruning strategies systematically remove redundant neurons, thereby reducing the computational burden without compromising performance [72]. Merging tokens through dynamic granularity, as demonstrated in [106], provides further scalability, allowing the model to manage computational resources more effectively.

Another significant challenge in multimodal learning is data scarcity, especially when dealing with modalities that are inherently resource-intensive to annotate, such as audio-visual datasets. Transfer learning, particularly from large-scale pre-trained Vision Transformers, presents an effective remedy for this issue. By leveraging models trained on extensive datasets like ImageNet, the adaptation to new tasks with limited data becomes feasible [12].

Cross-modal transfer learning, where knowledge from one modality aids in interpreting another, presents a sophisticated strategy for boosting performance in data-scarce environments. Techniques from natural language processing, such as text encodings, can enhance visual recognition capabilities [1]. However, challenges remain in effectively preserving and transferring complex semantic relationships across different data modalities.

The interpretability of Vision Transformers in multimodal applications remains a contentious issue. As these models function as black-box predictors, understanding their decision-making processes can be crucial for trust and acceptability, particularly in sensitive domains like healthcare and autonomous driving [107].

Methodologies have been proposed to elucidate the internal workings of Vision Transformers. Saliency maps and attention visualization techniques help demystify where and why models focus on certain parts of the input data, thereby providing insights into their behavior [21]. Moreover, novel interpretability strategies focus on decomposing complex attention mechanisms into more comprehensible

forms [108]. These frameworks not only foster user trust but also facilitate troubleshooting and enhancement of model robustness.

Moving forward, the integration of Vision Transformers into multimodal learning systems holds immense potential. Future research should aim to refine and expand fusion techniques, ensuring that synergy across modalities is maximally capitalized upon. The development of hybrid models that combine convolutional and transformer-based architectures can further harness both local and global features efficiently [10].

Additionally, addressing real-world application constraints, such as latency and edge-processing capabilities, will be critical for the deployment of Vision Transformers in practical scenarios. Effective strategies to tackle these issues include reducing model size and improving inference speed, which have seen promising advancements with emerging techniques like token reorganization and dynamic scalable architectures [109].

In sum, while Vision Transformers present unparalleled opportunities for multimodal learning, overcoming their inherent challenges requires continued innovation and interdisciplinary collaboration. By addressing computational demands, data efficiency, and interpretation transparency, Vision Transformers can revolutionize multimodal systems and significantly advance the field of artificial intelligence.

## 6.4 Future Directions for Vision Transformers in Multimodal Learning

Vision Transformers (ViTs) have significantly advanced the field of multimodal learning, paving the way for systems that understand and process information across various data modalities simultaneously. This subsection delves into the future directions for leveraging ViTs in multimodal learning, concentrating on technological advancements and research opportunities that could enhance their application and effectiveness.

A fundamental trend in the future application of vision transformers is the enhancement of multimodal fusion techniques, as integrating signals from multiple modalities promises deeper understanding and improved performance across tasks. The current state-of-the-art involves utilizing ViTs to combine visual data with textual or auditory information, yielding promising results in fields such as visual question answering and cross-modal retrieval [92]. However, these techniques still face challenges in effectively harnessing and managing diverse types of information to fully exploit multimodal data's potential.

The development of more sophisticated fusion techniques is a promising research avenue. Effective fusion must balance the intricate complexity of each modality while maintaining efficient computational overhead. For example, the Object Relation Transformer utilizes geometric attention to improve correlations between detected objects across modalities, providing a potential framework for comprehensive fusion strategies [110]. Further exploration into hierarchical and distributed fusion methodologies could yield models more adept at understanding nuanced inter-modal relationships, a direction suggested by structured approaches seen in attention-based models that regulate hierarchical layer interactions [35].

Vision Transformers are set to extend their utility by integrating with emerging technologies such as the Internet of Things (IoT) and autonomous systems. In IoT applications, promising research can investigate how ViTs process data from multiple sensors, leading to more intelligent and responsive IoT ecosystems. This integration will involve addressing challenges with latency and real-time processing inherent in edge computing environments, where computational resources are limited, yet the need for rapid and reliable data integration is crucial. Furthermore, implementing transformers in autonomous systems, such as self-driving vehicles, could enhance decision-making processes by providing a richer, multimodal perception of the environment.

Real-world applications necessitate addressing challenges such as computational efficiency, on-edge processing, and robustness to incomplete or noisy data. Future research can focus on refining model architectures through techniques like pruning and quantization to reduce complexity while preserving capability [111]. Moreover, optimizing for latency by inventing novel methods to manage computational demands without sacrificing performance is essential for deploying ViTs in real-time systems.

Data scarcity in some modalities remains a persistent hurdle. Transfer learning remains a vital approach to address this, where pre-trained ViTs on large datasets enable efficient domain adaptation to new environments or tasks, requiring minimal additional labeled data [37]. Domain-specific fine-tuning on Transformer models, employing techniques like cross-modal transfer learning, can aid in enhancing domain adaptation, especially where labeled data is scarce [37].

As these models evolve, another pivotal research area is interpretability and transparency. Understanding the decision-making processes within ViTs is imperative for fostering trust and improving transparency in multimodal settings [112]. Techniques such as attention visualization, causal inference, and other interpretability frameworks can provide insights into model operations, ensuring responsible and ethical deployment in sensitive applications like healthcare.

In summary, the future of Vision Transformers in multimodal learning is brimming with possibilities. The continuous advancement of fusion techniques, integration with new technologies, and addressing real-world challenges such as data scarcity, computational efficiency, and interpretability remain primary focal points. By leveraging these research opportunities, ViTs can substantially contribute to the ongoing development of intelligent systems capable of profound understanding across varied and complex information landscapes. As these models grow more sophisticated, their impact will extend beyond current applications, creating new opportunities to innovate in domains previously seen as challenging or inaccessible.

# 7 CHALLENGES AND LIMITATIONS

## 7.1 Computational Demands and Resource Constraints

Vision Transformers (ViTs) have emerged as pivotal architectures in computer vision, offering the capability to capture long-range dependencies and model comprehensive global contexts more effectively than traditional convolutional neural networks (CNNs). However, these advantages are accompanied by significant computational and memory challenges, threatening their widespread application, especially in contexts with limited resources. This subsection delves into these core challenges, analyzes existing mitigation strategies, and explores future directions that could help bridge these constraints.

One of the primary challenges underlying ViTs is their high computational cost, which stems from the intrinsic nature of the self-attention mechanism. Self-attention operates with a quadratic complexity concerning the input sequence length, leading to sizable computational loads, particularly when processing high-resolution images or large-scale datasets. The scalability issue limits the feasibility of deploying ViTs on resource-constrained platforms such as mobile and edge devices [9]. In response to this, various strategies have been developed to reduce computational loads while retaining model efficacy.

An influential approach in reducing computational demands is the sparse attention mechanism, which selectively limits attention computations to a subset of tokens. Techniques such as the QuadTree Attention focus computation on a dynamically selected subset of tokens, thereby reducing the number of attention operations without significantly degrading performance [102]. Similarly, focal self-attention balances fine-grained local and coarse-grained global interactions, capturing essential dependencies efficiently by adjusting attention granularity [17]. This restructuring of self-attention not only diminishes computational complexity but can also enhance task-specific performance by emphasizing relevant spatial information.

Quantization and pruning stand out as traditional yet effective methods for reducing model size and computational overhead in ViTs. Quantization techniques, including fixed-point quantization, reduce the number of bits used to represent model weights and activations without compromising neural network accuracy critically [14]. Meanwhile, pruning, both structured and unstructured, excises unnecessary components of a network. DynamicViT introduces a dynamic pruning mechanism that harnesses token sparsification, thereby reducing the number of tokens processed progressively without impairing the model's representational capacity [11]. These techniques collaboratively enhance memory efficiency and speed up both training and inference stages, crucial for large-scale model deployment in restricted environments.

Memory constraint issues in ViTs are primarily attributed to the need for storing and manipulating extensive attention weight matrices. To address this, recent efforts have been geared towards more memory-efficient architectures. Implementations like the Tokens-to-Token (T2T) Vision Transformer incorporate architectural ingenuity by progressively aggregating tokens, optimizing both memory utilization and computational operations [66]. Furthermore, integrating convolutional designs within vision transformers, as demonstrated by CvT, can impart shift, scale, and distortion invariance akin to CNNs, which reduces the dependency on large attention matrices [10].

Parallel to architectural advancements, training methodologies also play a crucial role in mitigating resource con-

straints. Multi-scale approaches, like the Multiscale Vision Transformers, efficiently leverage varying resolutions and channel scales to optimize feature extraction across network layers, thereby balancing computational demands against model performance [45]. Distributed training paradigms, facilitated by frameworks such as distributed stochastic gradient descent (SGD), allow large ViT models to be split across multiple computing nodes, harnessing parallel processing to manage datasets and model parameters more effectively [113].

Emerging trends also include the use of hybrid models that combine transformers with CNNs to exploit both local and global feature representation capabilities while managing resource demands. These hybrid architectures, as explored in models like CMT, bring together the strengths of CNNs in capturing local patterns and ViTs' prowess in modeling global contexts, optimizing both performance and computational costs [13].

While considerable advances have been made, challenges persist, particularly concerning maintaining a balance between computational expanse and accuracy. Recent work emphasizes the potential of automated resource optimization through techniques like automated progressive learning, which dynamically adjusts model complexity during training to enhance efficiency without sacrificing model robustness [91]. Additionally, developing more intelligent model architectures that inherently embody efficiency—through reduced parameter counts, sparse matrix computations, or enhanced feature reuse—remains a promising avenue.

Looking ahead, further innovation is required to ensure ViTs can be universally applied in various real-world scenarios, including those with stringent resource limitations. Coupling advances in hardware accelerations, such as more efficient GPUs and TPUs tailored for transformer operations, with architectural and algorithmic improvements offers a promising pathway to achieving this. Additionally, there is an urgent need for sustained research into sustainable and eco-friendly AI practices to minimize the carbon footprint associated with training and deploying massive ViT models. Ultimately, continued cross-disciplinary collaboration will be imperative to unlock the full potential of Vision Transformers, advancing the boundaries of what is computationally feasible and practically beneficial.

## 7.2 Data Efficiency and Training Complexity

The remarkable rise of Vision Transformers (ViTs) in computer vision has concurrently highlighted significant challenges related to data efficiency and training complexity. This subsection elucidates these challenges, focusing on the intricate balance needed between achieving high model performance and the associated demands for data and computational resources, as opposed to traditional models like Convolutional Neural Networks (CNNs).

Vision Transformers inherently differ from traditional convolutional architectures due to their reliance on large datasets to sufficiently capture and model spatial relationships within images. This dependency arises from the self-attention mechanism's need for a comprehensive contextual representation [114]. Empirical evidence suggests that ViTs typically underperform compared to CNNs when trained on limited datasets because the absence of inductive biases—such as locality and translation invariance—impairs effective feature learning [26].

Data efficiency emerges as a pivotal challenge for Vision Transformers. These models are generally expected to deliver exceptional performance only with extensive pre-training on vast datasets. Consequently, their straightforward adoption in data-scarce domains remains problematic. Recent advancements in data augmentation techniques and synthetic data generation have offered marginal improvements in data efficiency. Augmentation strategies that venture beyond basic transformations like rotation and scaling are being progressively tailored to suit the self-attentional architecture [24].

Furthermore, self-supervised learning has surfaced as a potent strategy for addressing data scarcity in ViTs. Techniques such as Masked Image Modeling (MIM) are used to learn representations without labeled data, thereby alleviating the pre-training burden [7]. These approaches exploit self-supervised dynamics by masking portions of input data and challenging the model to predict them, thus fostering a robust understanding of semantic structures. In contrast, transfer learning techniques strive to leverage pretrained models on large datasets, repurposing the learned representations for smaller, domain-specific datasets [7].

The complexity of training Vision Transformers can also be prohibitive. The quadratic computational complexity of the self-attention mechanism with respect to input size is a significant hurdle. This challenge has propelled the development of sparse, hierarchical, and local attention mechanisms to mitigate this complexity [21], [46]. Sparse attention methods selectively concentrate computation on salient data parts, thereby reducing computational overhead while maintaining the ability to model long-range dependencies. Concurrently, hierarchical attention structures endeavor to balance local and global information processing across multiple scales, optimizing both computational demands and performance [48].

The convergence dynamics of Vision Transformer models also pose challenges. Training ViTs requires sophisticated optimization strategies that heavily rely on first-order methods like Adam [115]. Nonetheless, the learning dynamics of visual representations through attention complicate gradient descent paths, often necessitating advanced learning rate schedules and regularization techniques to ensure convergence [52].

Emerging trends focus on hybridizing Vision Transformers with traditional convolutional frameworks to harness global and local feature modeling benefits [21], [26]. Such integrations propose enhancements in data efficiency and exhibit reduced computational burdens during training. The introduction of soft inductive biases allows ViTs to learn more effectively from smaller datasets, bridging the gap toward practical deployment in data-constrained environments.

Future directions in mitigating data efficiency and training complexity challenges may include enhanced integration of domain-specific knowledge into Vision Transformers, akin to how CNNs exploit locality. Innovations in efficient, parallelizable training frameworks are also antic-

ipated to optimize the computational economics of ViTs. Exploring evolutionary paradigms and meta-learning strategies has the potential to revolutionize model adaptation during training, fostering broader applicability across diverse datasets and task domains [116].

In conclusion, while Vision Transformers have ushered in a paradigm shift in computer vision, the inherent challenges of data efficiency and training complexity continue to influence their broader adoption and application scope. Strategic research focusing on judicious augmentation techniques, optimization of computational graphs, and innovative hybrid architectures holds promise for alleviating these challenges, paving the way for more pervasive use of ViTs in the future [68]. The ongoing discourse on these methodologies lays the groundwork for advancing Vision Transformers' capabilities and utility in achieving robust, scalable, and efficient model deployment.

## 7.3 Interpretability and Explainability

Vision Transformers (ViTs) have emerged as formidable models in computer vision, reshaping how tasks are approached and displaying remarkable performance across various domains. However, the interpretability and explainability of these models have become central concerns, given their complex architectures and the black-box nature inherent to deep learning models. Understanding Vision Transformers' decision-making processes is crucial for ensuring their trustworthy deployment, especially in critical applications such as medical imaging and autonomous systems.

The foundational challenge with the interpretability of ViTs lies in the self-attention mechanism, which facilitates long-range dependencies at the cost of increased model complexity and opacity. Self-attention mechanisms aggregate information across the entire input sequence, which often makes it challenging to delineate how specific input features contribute to final predictions. Saliency and attention visualization methods have been developed to address this challenge by exposing which parts of the input data attract more attention during processing [107]. These methods, while insightful, primarily provide a heuristic examination rather than a definitive causal explanation.

Saliency maps serve as a primary tool to visualize attention distributions within Vision Transformers. These maps highlight the relevance of individual input features based on the attention weights assigned during processing. While saliency methods like Gradient-weighted Class Activation Mapping (Grad-CAM) have been traditionally used with CNNs, adapting them to ViTs necessitates recognizing the layered self-attention architecture that encodes hierarchical dependencies through transformer blocks [104]. Furthermore, the inherent global nature of attention maps can lead to diffuse focus, spanning multiple, contextually significant areas of an image, all of which could contribute to the decision-making process, thereby posing a challenge in pinpointing precise causative dimensions of attention.

Causal and post-hoc explanations are another dimension explored in demystifying Vision Transformers. Post-hoc interpretability aims at providing an understanding after the model has made a prediction. Approaches like Integrated Gradients and SHAP (SHapley Additive exPlanations) have been employed to assign importances retrospectively to input features. Such methods, however, remain computationally expensive and may require extensive computational resources to produce fine-grained explanations of model behavior [108].

The incorporation of interactive attention prediction models offers a promising avenue for interpretability by enabling real-time user interaction to modify the attention weights and observe resultant changes. Such techniques reveal the dynamic nature of self-attention mechanisms, allowing stakeholders to understand how perturbations in input influence model predictions [68]. Furthermore, integration with end-to-end trainable methods for vision systems, as seen in works such as Contextual Transformer Networks which leverage context-aware attention, can profoundly enhance interpretability by incorporating human-centric factors directly into model training [22].

The trade-off between model complexity and explainability is pronounced in Vision Transformers compared to their CNN counterparts. Unlike CNNs, which inherently capture local feature hierarchies through their architecture, ViTs' reliance on global sequence processing renders explicit feature hierarchies implicit, thus complicating their interpretability [3]. Research suggests embedding structured attention mechanisms that incorporate domain-specific inductive biases to enhance interpretability without compromising performance. Such structures could leverage hybrid architectures that blend convolutional inductive biases with the transformer attention mechanism to delineate clearer hierarchical processing paths [44].

Despite these advancements, challenges remain in ensuring transparency in more intricate multitask and multimodal systems where ViTs are often deployed. The multimodal integration exacerbates the complexity of the interpretability landscape, necessitating novel approaches that can intuitively explain how models integrate information across different sensory inputs and modalities [16].

Looking forward, there are pivotal future directions for advancing the interpretability and explainability of Vision Transformers. One promising area is the development of explanation-aware frameworks that incorporate interpretability as a fundamental design component of model architectures. Such systems could assess and adjust model components through an interpretability lens during the pre-training and adaptation phases, thus embedding transparency intrinsically rather than retrospectively. Advances in explainability could also be enriched through collaborative efforts with experts in cognitive sciences to bridge gaps between neural network processing and human cognitive understanding, fostering more intuitive and human-aligned model interpretable methods [117].

In conclusion, while substantial strides have been made towards elucidating Vision Transformers' decisions, the path forward requires a multi-faceted approach integrating novel visualization techniques, model architecture innovations, and interdisciplinary insights. Such advances will not only render these models more transparent but also bolster trust and reliability, ensuring these powerful tools realize their potential across diverse domains. Continued focus on integrating interpretability into the design and deployment will be essential as these models become more pervasive in

sensitive areas requiring critical decision-making support.

## 7.4 Generalization and Robustness

Vision Transformers (ViTs) have demonstrated remarkable advances in computer vision tasks, often surpassing traditional convolutional neural networks in terms of performance. However, ensuring that these models generalize well across varied datasets and maintain robustness against adversarial attacks and distributional shifts remains a substantial challenge. This subsection delves into the intricacies of these challenges, examining existing strategies, comparative approaches, technical limitations, and the future directions necessary for advancing ViT robustness and generalization.

A fundamental challenge for Vision Transformers relates to their ability to generalize across different datasets. Unlike CNNs, which incorporate inductive biases such as translation invariance, ViTs operate on a different paradigm with their reliance on self-attention mechanisms. This structural difference can lead to challenges in maintaining consistent performance across diverse datasets. Techniques to enhance generalization often involve architectural innovations and refined training methodologies. One promising approach is domain adaptation, enabling models to better handle shifts in data distributions between training and test datasets. Techniques such as self-supervised learning, as discussed in [15], offer insights by using large volumes of unlabeled data to pre-train models, thus improving their adaptability across varied domains. Additionally, relative positional encoding provides another avenue for improving transformer adaptability. The introduction of directional distance modeling in [111] presents a mechanism that can help ViTs better capture spatial hierarchies, thereby enhancing their domain generalization capabilities.

On the robustness front, Vision Transformers are susceptible to adversarial attacks, which exploit model sensitivities to perturbations in input data. While CNNs have been extensively studied in this context, ViTs present unique vulnerabilities due to their non-local attention mechanisms. Strategies for enhancing adversarial robustness often involve integrating adversarial training and developing more resilient architectures. Enhancing the robustness of ViTs to distributional shifts is crucial for practical applications, requiring models to remain effective when faced with changes in input distributions, which can occur due to variations in environmental conditions or sensor noise. Techniques such as conditional positional encoding, proposed in [30], dynamically adjust to input variations, effectively mitigating performance degradation due to such shifts. Furthermore, methods like the discussed invariant features in [76] emphasize the significance of integrating spatial transformation invariance into the transformer framework to enhance robustness.

Another layer of complexity arises with cross-modal and multimodal inputs, which frequently change in distribution. The innovation described in [118] highlights the potential of cross-modal fusion techniques that strengthen both visual and language feature representations, leading to improved robustness in multimodal settings.

Despite these advancements, several challenges remain. One such issue is the scalability of training methods needed to instill robustness, as transformer models scale with data and computation resources more than their convolutional counterparts. Adopting distributed learning strategies, coupled with memory-efficient training models, as explored in [119], can offer pathways to mitigate these constraints and foster improved performance scaling.

Moreover, the interpretability of ViTs is pivotal in understanding and enhancing their robustness. Efforts to visualize and interpret attention mechanisms are vital for uncovering the model's decision-making processes and subsequently fostering more resilient systems. Additionally, [112] provides methodologies for decomposing and interpreting model representations, which can be instrumental in diagnosing and addressing robustness issues.

In conclusion, improving the generalization and robustness of Vision Transformers is an intricate task that requires multi-faceted approaches, ranging from novel architectural designs to advanced training regimes. While significant progress has been made, continued research must focus on understanding the underlying limitations of these models and developing comprehensive solutions that encompass domain adaptation, adversarial defenses, and robust multimodal integration. The future of ViTs may hinge on innovations that not only enhance performance metrics but also ensure the reliability and trustworthiness of these systems in real-world applications, where they can encounter a myriad of unforeseen challenges. Exploring hybrid models that integrate the strengths of different architectures could also unlock new pathways for achieving robust generalization across diverse visual environments.

## 7.5 Environmental and Hardware Limitations

The deployment of Vision Transformers (ViTs) in resource-constrained environments, such as edge devices and embedded systems, presents significant environmental and hardware challenges that limit their feasibility. These challenges are rooted in the computational intensity and memory demands intrinsic to Transformer architectures, which are exacerbated under constraints of limited processing power, energy availability, and thermal management. This subsection delves into these limitations and explores potential strategies to ameliorate them, supported by recent research and advances in hardware-aware design and sustainability practices.

Vision Transformers are synonymous with significant computational and memory overhead, primarily due to the quadratic complexity of their self-attention mechanisms [120]. This computational burden is further intensified when dealing with high-resolution image inputs, which, while enabling richer feature extraction, also demand more from the hardware [38]. When considering deployment on edge devices, where power availability and thermals are constrained, these demands pose nontrivial obstacles. For instance, typical deployments require balancing throughput and latency with minimal energy consumption, a task that becomes challenging with Vision Transformers due to their extensive resource requirements for both inference and training [121].

Memory constraints become another barricade, where traditional high-precision computations (e.g., 32-bit floating

point) are infeasible due to limited memory bandwidth and the size of memory available on edge devices. The substantial size of ViTs, especially larger models, not only complicates their deployment but also leads to practical issues related to heat dissipation and power draw, which further exacerbate the environmental footprint of these models. Addressing these hardware-related limitations is crucial to advancing the real-world applicability of Vision Transformers.

Several approaches aim to mitigate these limitations through model compression and architectural innovation. Quantization and pruning are prevalent techniques employed to reduce the model size and computational demand by lowering the bit precision of weights and biases, compressing the model, and selectively removing less critical parameters without significantly losing predictive accuracy [99]. Not only do these methods help in reducing the model's size, but they also contribute to faster computations and reduced energy consumption, making them viable for deployment on low-power hardware.

Further, innovations in efficient attention mechanisms streamline computational loads by adopting sparsity techniques and lightweight attention models. For example, the use of structured sparsity and local attention mechanisms allows selectively focusing computational resources on meaningful portions of the input data, thus conserving processing power [43] [27]. These techniques significantly cut down on the resource demand without degrading model performance, offering an eco-friendly alternative for deploying ViTs sustainably at the edge.

Another promising direction is the hardware-aware design and optimization of Vision Transformers, which acknowledges the specific constraints and strengths of different hardware platforms. Hardware-aware neural architecture search (NAS) methods are gaining traction as they optimize model architectures specifically for target hardware platforms, taking into account factors like available memory, parallel processing capabilities, and energy efficiency [63]. These approaches ensure that the resultant models are not only efficient in operations but are also capable of meeting the performance requirements of the edge devices they are meant to run on.

Moreover, sustainable AI practices encourage the adoption of eco-friendly procedures throughout the lifecycle of a Vision Transformer model—from training to deployment. This includes using renewable energy sources to power data centers for model training, optimizing resource scheduling for model inference, and adopting algorithms designed to minimize carbon footprint [122]. By advancing these practices, the AI community not only addresses the environmental concerns associated with Vision Transformers but also paves the way for more conscientious deployment strategies.

In conclusion, while Vision Transformers are emerging as powerful tools in computer vision applications, their deployment in resource-constrained settings presents distinct environmental and hardware challenges that must be addressed to leverage their full potential. The integration of techniques like quantization, pruning, model compression, hardware-aware design, and sustainable AI practices presents a multifaceted approach to overcoming these barriers. Future research opportunities lie in the refinement of these methods and the exploration of new pathways for balance between performance, efficiency, and sustainability. By progressively refining these strategies, Vision Transformers can be harnessed in a manner that not only pushes the boundaries of technology but also adheres to the pressing demands of environmental responsibility and technological adaptability in the face of diverse deployment constraints.

## 8 FUTURE DIRECTIONS AND RESEARCH OPPORTUNITIES

### 8.1 Efficiency and Scalability Innovations

In recent years, Vision Transformers (ViTs) have emerged as a revolutionary architecture in the domain of computer vision, exhibiting state-of-the-art performance across a myriad of tasks. However, the computational intensity and resource demands associated with these models pose significant hurdles for their scalability and efficiency in practical applications. This subsection explores architectural innovations and algorithmic strategies devised to address these challenges, focusing on enhancing computational efficiency while maintaining or even improving model performance.

Architectural redesigns have played a pivotal role in boosting the efficiency of Vision Transformers. A prime example is the integration of convolutional designs within transformer architectures, which leverages the strengths of both convolutional neural networks (CNNs) and transformers. Convolutional Vision Transformers (CvT) [10] exemplify this approach, offering improved performance and efficiency by incorporating a convolutional token embedding and a convolutional projection mechanism. This integration exploits the locality and shift-invariance properties inherent in CNNs while retaining the global contextual understanding facilitated by transformers. Similarly, the Convolution-enhanced image Transformer (CeiT) [44] employs locally-enhanced feed-forward layers, effectively bridging the gap between CNNs and ViTs and addressing limitations related to local feature extraction.

Beyond architectural enhancements, algorithmic innovations have been instrumental in reducing computational load. Dynamic token management represents a significant advancement, as seen in models like DynamicViT [11]. By employing a dynamic token sparsification framework, redundant tokens are pruned progressively based on their informativeness, optimizing computation while minimizing accuracy loss. This strategy highlights the potential of dynamic mechanisms to balance efficiency with performance, offering insights into scalable model deployment.

The concept of sparsity has further guided efficiency advancements. Attentions mechanisms, traditionally hampered by quadratic complexity, have been reimagined through efforts such as Focal Self-attention [17] and QuadTree Attention [102]. Focal Self-attention introduces a mechanism whereby tokens focus on immediate neighbors with fine granularity while attending to distant tokens coarsely. This not only reduces computational overhead but also retains key long-range dependencies crucial for vision tasks. QuadTree Attention, on the other hand, utilizes a hierarchical token pyramid approach, evaluating attention at multiple levels to maintain efficiency and effectiveness.

Another promising avenue of research lies in the optimization of multi-scale processing within transformers. Multiscale Vision Transformers (MViT) [45] and their subsequent iterations have adapted pyramid architectures to process visual data at varying resolutions. This paradigm allows for high-resolution input handling in early layers, gradually abstracting features through successive layers. This strategic scaling enhances model capacity and detail retention without incurring the computational costs associated with full-resolution processing throughout the model's depth.

Model compression techniques such as pruning and quantization have also gained traction as methods for achieving efficiency gains without substantial sacrifices in performance. For instance, structured and unstructured pruning approaches have been adapted to transformer architectures [123], enabling the reduction of model size while maintaining functional integrity. These techniques support the deployment of ViTs in resource-constrained environments by trimming unnecessary parameters and operations.

In tandem with these strategies, innovative training methodologies have equipped Vision Transformers with capabilities to handle limited computational resources more effectively. Automated Progressive Learning [91] exemplifies this by dynamically adjusting model capacity during training, ensuring that computational resources are allocated efficiently throughout the learning process. This approach minimizes training overhead while maintaining robust performance, making ViTs more accessible and scalable.

From this examination of efficiency and scalability advancements in Vision Transformers, several pertinent trends and challenges emerge. There is a growing emphasis on hybrid architectures that amalgamate the best aspects of different neural network paradigms. The fusion of CNN-like locality and ViT's global dependency modeling illustrates a path forward for network design that can achieve superior performance without exorbitant computational costs. Moreover, nearly all efficient designs strive to maintain a delicate equilibrium between resource savings and performance retention, underscoring the importance of adaptive mechanisms that respond dynamically to varying input and task demands.

The pursuit of scalability in Vision Transformers is likely to spearhead further innovations in hierarchical architectures, enabling these models to process information with an unprecedented level of efficiency. Given the evolution of these architectures, it becomes imperative to develop standardized benchmarks for evaluating efficiency advancements. Furthermore, new research into cross-layer interaction, information-sharing mechanisms across hybrid networks, and advanced token management techniques will undoubtedly shape the next generation of scalable Vision Transformers.

As the field progresses, the focus will likely extend to the deployment implications in edge computing environments, where constraints on power, memory, and real-time processing capabilities demand even greater optimization. Vision Transformers hold tremendous potential, and as efficiency and scalability continue to improve, their applicability across diverse sectors—from autonomous driving to augmented reality—expands, paving the way for ViTs to become foundational components in future technological landscapes. This necessitates ongoing investment in innovative strategies that balance efficacy with feasibility, ensuring that Vision Transformers remain at the forefront of advancements in computer vision.

## 8.2 Domain Expansion and Emerging Applications

The growing application of Vision Transformers (ViTs) across various domains marks a significant evolution in computer vision technology, aligning with the architectural and algorithmic innovations that enhance their efficiency and scalability. As these models continue to exhibit exceptional performance in traditional tasks, their expansion into unconventional areas could instigate revolutionary changes across industries, from healthcare to autonomous systems. This subsection explores these emerging applications, emphasizing both potential breakthroughs and the inherent challenges they entail.

In the realm of medical imaging, Vision Transformers offer promising advancements. Their capability to process and analyze complex high-resolution images provides unprecedented advantages in diagnostic accuracy. For instance, ViTs excel in identifying minute pathological patterns in medical imagery, potentially transforming disease diagnostics [81]. The inherent ability of ViTs to model long-range dependencies allows them to enhance feature representation, aiding in more precise segmentation of anatomical structures. However, challenges concerning the need for sensitivity in medical diagnoses persist, as false positives or negatives can have significant implications. Moreover, the reliance on large annotated datasets for training remains an obstacle, necessitating innovative data augmentation and synthesis approaches to traverse this gap [51].

Conversely, autonomous systems and robotics represent another promising application domain for Vision Transformers. Here, the model's ability to integrate and process multiple data streams—visual and otherwise—can augment real-time decision-making capabilities vital for autonomous vehicles and robotics. ViTs may potentially replace existing systems by providing improved environmental perception and navigation accuracy. For instance, the Deformable Attention Transformer [20] has been shown to enable efficient real-time processing by dynamically focusing attention, thereby enhancing the resilience and reliability of autonomous decision systems. Nonetheless, obstacles such as adverse weather conditions, diverse terrains, and varying lighting conditions persist. Further research is necessary to bolster robustness against adversarial inputs and environmental variabilities [29].

In addition to these fields, Vision Transformers are expanding into broader sectors like environmental monitoring and precision agriculture. Within precision agriculture, ViTs can conduct detailed analyses of satellite and aerial imagery to monitor crop health, identify diseases or pest infestations, and optimize resource allocation. Their capacity to model both local and global contexts allows them to outperform traditional methods in deriving meaningful insights from such data [23]. However, addressing computational resource allocation and real-time data processing challenges is

crucial for fully leveraging ViTs in these areas. Furthermore, there is an emerging need for multimodal data integration, combining spectral, thermal, and visual data for comprehensive analysis.

Additional prospective applications of Vision Transformers include their deployment in security and surveillance systems, where their ability to identify and track objects across frames could significantly enhance real-time monitoring capabilities. Moreover, adapting ViTs to interpret low-light or obscured conditions could further integrate them into high-security environments [124].

Yet, despite the promising prospects, deploying Vision Transformers in these novel domains presents technical and ethical challenges. Technically, the computational complexity and high memory demand of ViTs pose significant constraints, particularly within edge computing paradigms where resources are limited [49]. Developing efficient methods to reduce resource demands without sacrificing performance will enhance the feasibility of utilizing ViTs in constrained environments. Furthermore, issues of interpretability and explainability remain crucial in sensitive domains such as healthcare and autonomous systems. Improving the transparency of ViT decision-making processes remains a vital area of research to ensure models provide trustworthy and accountable outputs [117].

From an ethical standpoint, frameworks must be developed to ensure the responsible application of Vision Transformers, mitigating risks related to privacy, bias, and data security [27]. Strategies to address bias, like those explored in Bias-Free ViT, can be pivotal in ensuring fair and equitable AI systems.

In conclusion, the extension of Vision Transformers into new domains possesses transformative potential across numerous sectors. However, challenges such as computational resource allocation, robustness, interpretability, and ethical deployment must be navigated. By focusing on these issues, the research community can optimize Vision Transformers for a range of applications, unlocking their potential while ensuring responsible use. Future research directions should involve creating hybrid models that draw upon various architectures to increase domain adaptability and continue advancing model efficiency and scalability [20]. Such progress is crucial to fully leverage the benefits of Vision Transformers in these burgeoning application domains.

## 8.3 Multimodal and Hybrid Systems

The burgeoning field of Vision Transformers (ViTs) has opened new avenues in the realm of multimodal and hybrid systems, providing versatile platforms for integrating diverse data modalities to enrich computational capabilities and understanding. This subsection delves into the integration of Vision Transformers with other modalities and structural systems, evaluating different approaches, their relative strengths and weaknesses, and potential future directions.

Vision Transformers offer a robust framework characterized by their aptitude for capturing global dependencies, a feature particularly beneficial in multimodal environments where contextual understanding from varied sources is pivotal. ViTs' ability to process sequences allows them to interact seamlessly with other data modalities like textual

and auditory information, facilitating comprehensive integration in tasks such as visual question answering (VQA) and cross-modal retrieval. Such an integration leverages the self-attention mechanism to capture complex interdependencies across modalities, improving model performance and enriching context understanding [1]. In visual question answering systems, for instance, the self-attention mechanism inherent in Vision Transformers enhances the model's understanding of visual context in conjunction with textual queries, leading to more accurate question comprehension and response generation [125].

One prominent approach is the hybridization of Vision Transformers with convolutional neural networks (CNNs), which aims to capitalize on CNNs' proficiency in local feature extraction alongside ViTs' global contextual modeling. This synthesis has inspired architectures like Convolutional ViTs (CvTs), which incorporate convolutional token embedding to introduce locality biases crucial for resolving image context [10]. Such hybrids are shown to outperform their purely convolutional counterparts on tasks requiring intricate feature integration, such as semantic segmentation and object detection [126]. Similarly, efforts like LocalViT embed local attention mechanisms within vision transformer structures, enabling the effective capture of detailed visual cues necessary for finer segmentation and detection precision [21].

Multiscale Vision Transformers further exemplify the utility of integrating traditional neural network elements by implementing multiscale feature hierarchies. MViTs optimize feature extraction across resolutions, enabling the handling of complex visual signals with effective computational resource allocation [45]. The balance achieved through these hybrid structures showcases both architectural depth and computational efficiency, highlighting the potential of such integrations in scaling Vision Transformers to diverse and challenging tasks [33].

However, the integration of Vision Transformers with other modalities is not without challenges. The computational complexity remains a significant hurdle, often impeding scalability and practical deployment in real-time applications. Sparse and local attention mechanisms have been proposed to alleviate computational burdens by refining attention focus to relevant modalities. For instance, the use of sparse attention in Vision Transformers has been shown to reduce computational costs significantly while maintaining performance, addressing one of the core limitations of extensive cross-modal models [53]. Similarly, attention mechanisms optimized for specific modalities, such as the Dual Attention Vision Transformers (DaViT), offer computational efficiency by alternating between spatial and channel tokens to capture both global and local interactions without compromising performance [29].

Emerging challenges include scalability and interpretability in multimodal interactions, where maintaining transparent model processes is critical. Techniques to demystify Vision Transformers' decision-making processes are crucial for their application in sensitive domains like medical imaging and autonomous navigation [75]. Enhancing interpretability and robustness in multimodal models will likely involve developing novel methods for visualizing cross-modal interactions and understanding emergent be-

haviors within these hybrid systems.

Future research directions may explore deeper integration techniques that go beyond hybrid architectures, such as leveraging reinforcement learning to adaptively choose relevant modalities or samples. As Vision Transformers continue to develop, the expansion into unexplored domains such as IoT and edge computing will necessitate innovations that address latency, real-time processing demands, and resource constraints [127].

In summarizing, while Vision Transformers are advancing the capabilities of multimodal and hybrid systems, scientific progress rests on resolving computational inefficiencies and bolstering interpretability. Further exploration into lightweight yet potent architectural redesigns and adaptive algorithmic strategies will be imperative. As researchers continue to harness the transformative potential of Vision Transformers, the collaborative future of diverse multimodal applications becomes not only plausible but promising.

### 8.4 Human-Centric and Responsible AI

In the rapidly advancing landscape of Vision Transformers (ViTs), ensuring human-centric and responsible AI development is increasingly crucial. This subsection explores the imperative advancements necessary to enhance the interpretability, fairness, and ethical implementation of Vision Transformers, fostering a responsible AI paradigm that aligns technology with societal values.

Despite their remarkable strides across various applications, the opacity of Vision Transformers remains a significant barrier to broader trust and adoption. Interpretability is a core aspect of responsible AI, essential for understanding the decision-making processes of ViTs and ensuring stakeholders can trust and validate these systems. Efforts to enhance interpretability often involve demystifying the attention mechanisms intrinsic to transformers. For instance, research on Rotary Position Embedding for Vision Transformers highlights how positional encodings can influence the spatial reasoning capacity of ViTs without sacrificing performance [56]. Such advances not only offer methods to visualize model behavior but also support error analysis and debugging, essential for high-stakes applications like medical imaging and autonomous vehicles.

Beyond technical clarity, achieving fairness in Vision Transformers necessitates addressing biases that emerge during data collection and model training. Bias in AI systems can propagate through skewed datasets or myopic algorithmic designs, potentially leading to discriminatory outcomes. Models like BRAVE, which incorporate diverse inductive biases, demonstrate how leveraging varied visual encoders can mitigate bias by offering more robust and generalized representations [98]. This method reduces overfitting to specific data distributions and provides a foundational step toward equitable AI systems. Additionally, by expanding the visual encoding capacity and diversifying training data, these models can better generalize across cultures, demographics, and environments, promoting fairness in AI applications.

However, merely enhancing model interpretability and mitigating bias is insufficient. Ethical implementation of Vision Transformers involves adhering to legal regulations and societal norms. This requires a multidisciplinary approach where computer scientists, ethicists, and policymakers collaboratively establish guidelines and frameworks to steer AI development responsibly. The Conditional Positional Encodings for Vision Transformers pose an exemplary effort, showing how adaptable encoding schemes can support compliance with ethical standards without losing computational efficiency [30].

Amidst these technical strategies, the role of user-centric design and inclusivity is paramount. Human-centered approaches in AI development ensure systems are accessible and beneficial to end users. Efforts like developing Qwen2-VL models emphasize dynamic processing capabilities, adapting to user needs by generating more efficient visual representations [55]. This highlights the importance of designing interactable AI systems catering to diverse user requirements and contexts, democratizing access to advanced machine learning technologies.

Emerging trends in Vision Transformers highlight increasing attention towards scalable and efficient designs that do not compromise ethical standards. One such trend is the optimization of model architectures to reduce energy consumption, supporting environmentally sustainable AI practices. Techniques illustrated in Vision Mamba showcase efficient visual representation learning while maintaining high performance across computational constraints, reflecting the growing emphasis on resource-conscious AI [94]. This is particularly relevant amidst global initiatives to reduce AI's carbon footprint, aligning technical advancements with broader environmental goals.

Despite these achievements, challenges persist. A primary concern is maintaining the balance between model complexity and interpretability—complex models may offer superior accuracy but at the cost of transparency and explainability. Researchers continue to explore novel architectures like A Transformer-Based Feature Segmentation and Region Alignment Method For UAV-View Geo-Localization, integrating innovative self-attention mechanisms to enhance usability while maintaining architectural simplicity [128]. These approaches represent promising directions in harmonizing performance with clarity.

Furthermore, the trajectory of Vision Transformers lies in integrating contextual knowledge and domain-specific insights into their frameworks. Doing so can significantly mitigate unintended biases and ethical pitfalls. Employing a combination of human-in-the-loop approaches can ensure continuous feedback and validation, aligning the training and deployment of ViTs with ethical principles. Initiatives like Image Translation as Diffusion Visual Programmers underline the potential of such integration, offering frameworks blending algorithmic efficiency with user-centered intuitiveness [129].

In conclusion, while remarkable progress has been made toward richer and more responsible Vision Transformers, enduring vigilance and proactive innovation are needed to address ongoing challenges in interpretability, fairness, and ethics. The future of responsible AI lies in fostering cross-disciplinary collaboration and continued research into scalable, transparent, and ethically aligned Vision Transformer models. By emphasizing human-centric design and

responsible innovation, the field can ensure these sophisticated systems not only advance technologically but also benefit society as a whole, setting a benchmark for ethical AI practices worldwide.

## 8.5 Robustness and Adaptability

The robustness and adaptability of Vision Transformers (ViTs) are critical qualities that enable them to maintain high performance across a variety of challenging conditions. These qualities are especially vital as ViTs are increasingly deployed in real-world scenarios where data can be highly heterogeneous and subject to noise, occlusion, or varying levels of quality. This subsection explores the strategies developed to enhance the robustness and adaptability of ViTs, offering insights into current approaches, trade-offs, and directions for future research.

Vision Transformers have demonstrated remarkable capabilities in visual recognition tasks by leveraging self-attention mechanisms, which allow them to effectively capture long-range dependencies in data. However, their ability to generalize and remain robust under challenging conditions, such as adversarial attacks, abnormal data distributions, and unseen domains, remains an active area of research and development. Improving robustness and adaptability involves enhancing the architectures to tolerate data variability and ensuring that models can be seamlessly transferred across diverse applications without extensive retraining [60], [130].

One promising strategy for enhancing adversarial robustness involves augmenting the self-attention mechanism. The Refiner approach directly refines the self-attention maps, including modifications like attention expansion and convolutional augmentations, aiming to diversify energy across attention heads for better generalization and robustness [131]. Additionally, robust training techniques such as adversarial training have been proposed to fortify ViTs against adversarial inputs, enhancing their stability under attack scenarios.

Ensuring that Vision Transformers generalize well across domains without retraining is a considerable challenge. Techniques such as domain adaptation methods have been explored to adjust the learning process dynamically based on the input data variability. By fine-tuning Vision Transformers on domain-specific datasets using strategies like domain-invariant feature extraction or domain-agnostic learning strategies, models become more adaptable to new domains—both in seen and unseen conditions [64].

Handling visual noise and occlusions requires models that can prioritize salient features and suppress irrelevant data. Vision Transformers can incorporate spatial attention mechanisms that emphasize vital regions while downweighing background noise [74]. These approaches can be complemented by strategic token pruning, which focuses the model's computational resources on the most informative inputs, enhancing performance in cluttered or occluded scenes [41].

Cross-modal adaptability is significantly enhanced through multimodal learning approaches that integrate Vision Transformers with other sensory inputs, such as audio or language, thus facilitating robust perception under multifaceted conditions [132]. These strategies enable Vision Transformers to capture more holistic representations of the environment, thus bolstering their adaptability.

Models like SPViT implement soft token pruning to dynamically prune less informative tokens, significantly reducing computation while preserving model accuracy. This adaptability allows Vision Transformers to maintain efficacy across a range of computational environments, from high-powered servers to resource-constrained edge devices [41].

As Vision Transformers evolve, their integration with hardware accelerators is becoming increasingly important to achieve rapid model adaptation. Approaches such as Hardware-Aware Transformers (HAT) have aimed to optimize ViTs for specific hardware constraints, thereby enabling low-latency inference and execution on limited-resource platforms without sacrificing robustness [63].

Future research directions emphasize the development of sustainable AI models that not only adapt dynamically to diverse conditions but also reduce their environmental footprint. This involves the co-design of algorithms and hardware to reduce energy consumption while boosting adaptability and robustness [121].

Incorporating meta-learning algorithms can endow Vision Transformers with the capability to self-improve continuously by learning from prior experiences across tasks and environments. These approaches could potentially lead to models that require fewer data and less computational resources for adaptation, making them highly desirable for real-time applications.

In conclusion, the robustness and adaptability of Vision Transformers are critical to their application in real-world scenarios where data variability and noise are prevalent. Energetic research is underway to enhance these models by refining self-attention mechanisms, employing domain generalization techniques, and developing cross-modal adaptability. Future innovations in efficient algorithm-hardware integration and sustainable AI practices promise to expand the horizons of Vision Transformers, making them robust, adaptable, and environmentally responsible components of intelligent systems.

## 9 CONCLUSION

Vision Transformers (ViTs) have undoubtedly revolutionized computer vision. This survey has detailed their transformative role, underscoring their ability to model global dependencies and provide a universal framework across different tasks and domains. The strengths of Vision Transformers lie in their flexible architecture, which eschews the rigid structure of convolutional neural networks (CNNs), allowing for a parallelized and more effective modeling of complex visual tasks [1], [2]. This flexibility has broadened the application scope of ViTs, enabling their use in tasks such as image classification, object detection, semantic segmentation, and even multimodal tasks where they integrate with other data types [1], [7].

One of the core insights gleaned from the extensive study of ViTs is their unparalleled ability to capture long-range relationships within image data. Unlike CNNs, where the learning of features is largely local due to convolutional operations, ViTs benefit from the self-attention mechanism enabling them to form global contextual representations [1],

[2]. This attribute is particularly beneficial in high-resolution vision tasks and situations requiring the modeling of complex object relationships within a scene, making Vision Transformers well suited for more nuanced and intricate vision applications [17], [45].

However, Vision Transformers are not without their limitations. Their primary challenge remains the computational demands and memory resources required, a consequence of the self-attention mechanism's quadratic complexity. Recent efforts have seen the introduction of efficient attention mechanisms aimed at reducing these computational loads by means of sparse and locality-focused attention designs [11], [102]. Moreover, hybrid approaches that integrate CNN-like properties into ViTs are developed to better harness local feature extraction while maintaining the benefits of global attention [10], [44].

The emergence of hybrid models suggests a significant research trajectory focused on balancing global and local feature extraction. By incorporating convolutional elements into ViTs, hybrid architectures endeavor to merge the best aspects of CNNs and transformers, particularly addressing the gap in inductive biases that ViTs experience when operating on small datasets without extensive pre-training [13], [21]. Hybridization has been a pivotal strategy in improving model efficiency, reducing computational overhead, and providing robustness across varying scales of data and tasks.

Looking forward, there is ample opportunity to refine these transformers, aiming for greater adaptability and efficiency without compromising their formidable capabilities. Innovative approaches such as automated resource allocation during training, progressive learning methods, and dynamic token pruning are key areas of interest that promise further reductions in computational cost, thereby broadening the applicability and accessibility of Vision Transformers [11], [91].

With the landscape of AI constantly evolving, Vision Transformers are posited to spearhead developments, especially in fields that require the parsing of complex relational data. Their role in multimodal learning environments continues to expand, leveraging their architecture to bridge visual and linguistic modalities seamlessly [5], [133]. The integration of Vision Transformers with advanced multimodal datasets promises further breakthroughs, enabling models to better understand and process diverse streams of information concurrently.

In conclusion, Vision Transformers signify a pivotal advancement in computer vision. By providing a robust, scalable framework, they have already facilitated substantial strides across myriad vision tasks. Challenges remain, particularly in resource optimization and model interpretability [4], [68]. Nonetheless, ongoing research shows promising directions for overcoming these hurdles, ensuring that Vision Transformers continue to lead as a dominant architecture in the AI domain. The future of computer vision, enhanced by Vision Transformers, points toward an era characterized by more intelligent, versatile, and efficient visual processing systems.

## REFERENCES

[1] S. H. Khan, M. Naseer, M. Hayat, S. W. Zamir, F. Khan, and M. Shah, "Transformers in vision: A survey," *ACM Computing Surveys (CSUR)*, vol. 54, pp. 1 – 41, 2021. 1, 2, 4, 10, 14, 15, 16, 18, 20, 27, 29

[2] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," *ArXiv*, vol. abs/2010.11929, 2020. 1, 4, 12, 14, 20, 29

[3] B. Heo, S. Yun, D. Han, S. Chun, J. Choe, and S. J. Oh, "Rethinking spatial dimensions of vision transformers," *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 11 916–11 925, 2021. 1, 18, 23

[4] M. Naseer, K. Ranasinghe, S. H. Khan, M. Hayat, F. Khan, and M.-H. Yang, "Intriguing properties of vision transformers," in *Neural Information Processing Systems*, 2021, pp. 23 296–23 308. 1, 30

[5] Y. Wang, X. Chen, L. Cao, W. bing Huang, F. Sun, and Y. Wang, "Multimodal token fusion for vision transformers," *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 12 176–12 185, 2022. 1, 18, 30

[6] H. Yunusa, S. Qin, A. H. A. Chukkol, A. A. Yusuf, I. Bello, and A. Lawan, "Exploring the synergies of hybrid cnns and vits architectures for computer vision: A survey," *ArXiv*, vol. abs/2402.02941, 2024. 1

[7] Y. Liu, Y. Zhang, Y. Wang, F. Hou, J. Yuan, J. Tian, Y. Zhang, Z. Shi, J. Fan, and Z. He, "A survey of visual transformers," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 35, pp. 7478–7498, 2021. 1, 18, 19, 22, 29

[8] S. H. Lee, S. Lee, and B. Song, "Vision transformer for small-size datasets," *ArXiv*, vol. abs/2112.13492, 2021. 1, 3, 10, 14

[9] X. Zhai, A. Kolesnikov, N. Houlsby, and L. Beyer, "Scaling vision transformers," *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1204–1213, 2021. 2, 10, 14, 15, 19, 21

[10] H. Wu, B. Xiao, N. Codella, M. Liu, X. Dai, L. Yuan, and L. Zhang, "Cvt: Introducing convolutions to vision transformers," *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 22–31, 2021. 2, 4, 6, 8, 12, 14, 20, 21, 25, 27, 30

[11] Y. Rao, W. Zhao, B. Liu, J. Lu, J. Zhou, and C.-J. Hsieh, "Dynamicvit: Efficient vision transformers with dynamic token sparsification," *ArXiv*, vol. abs/2106.02034, 2021. 2, 6, 14, 15, 21, 25, 30

[12] H. Gani, M. Naseer, and M. Yaqub, "How to train vision transformer on small-scale datasets?" *ArXiv*, vol. abs/2210.07240, 2022. 2, 10, 14, 16, 20

[13] J. Guo, K. Han, H. Wu, C. Xu, Y. Tang, C. Xu, and Y. Wang, "Cmt: Convolutional neural networks meet vision transformers," *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 12 165–12 175, 2021. 2, 6, 14, 22, 30

[14] T. Lin, Y. Wang, X. Liu, and X. Qiu, "A survey of transformers," *AI Open*, vol. 3, pp. 111–132, 2021. 2, 10, 21

[15] F. Shamshad, S. H. Khan, S. W. Zamir, M. H. Khan, M. Hayat, F. Khan, and H. Fu, "Transformers in medical imaging: A survey," *Medical image analysis*, vol. 88, p. 102802, 2022. 2, 10, 24

[16] K. Islam, "Recent advances in vision transformer: A survey and outlook of recent work," *ArXiv*, vol. abs/2203.01536, 2022. 2, 16, 23

[17] J. Yang, C. Li, P. Zhang, X. Dai, B. Xiao, L. Yuan, and J. Gao, "Focal self-attention for local-global interactions in vision transformers," *ArXiv*, vol. abs/2107.00641, 2021. 2, 4, 8, 15, 20, 21, 25, 30

[18] W. Wang, L. Yao, L. Chen, B. Lin, D. Cai, X. He, and W. Liu, "Crossformer: A versatile vision transformer hinging on cross-scale attention," in *International Conference on Learning Representations*, 2021. 2

[19] R. Yang, H. Ma, J. Wu, Y. Tang, X. Xiao, M. Zheng, and X. Li, "Scalablevit: Rethinking the context-oriented generalization of vision transformer," *ArXiv*, vol. abs/2203.10790, 2022. 2

[20] Z. Xia, X. Pan, S. Song, L. E. Li, and G. Huang, "Vision transformer with deformable attention," *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4784–4793, 2022. 3, 15, 26, 27

[21] Y. Li, K. Zhang, J. Cao, R. Timofte, and L. Gool, "Localvit: Bringing locality to vision transformers," *ArXiv*, vol. abs/2104.05707, 2021. 3, 12, 14, 18, 20, 22, 27, 30

[22] Y. Li, T. Yao, Y. Pan, and T. Mei, "Contextual transformer networks for visual recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. PP, pp. 1–1, 2021. 3, 23

[23] C.-F. Chen, R. Panda, and Q. Fan, "Regionvit: Regional-to-local attention for vision transformers," *ArXiv*, vol. abs/2106.02689, 2021. 3, 11, 26

[24] L. Melas-Kyriazi, "Do you even need attention? a stack of feed-forward layers does surprisingly well on imagenet," *ArXiv*, vol. abs/2105.02723, 2021. 3, 22

[25] Y. Rao, W. Zhao, Z. Zhu, J. Lu, and J. Zhou, "Global filter networks for image classification," *ArXiv*, vol. abs/2107.00645, 2021. 3, 11

[26] S. d'Ascoli, H. Touvron, M. L. Leavitt, A. S. Morcos, G. Biroli, and L. Sagun, "Convit: improving vision transformers with soft convolutional inductive biases," *Journal of Statistical Mechanics: Theory and Experiment*, vol. 2022, 2021. 3, 15, 22

[27] A. Hassani, S. Walton, J. Li, S. Li, and H. Shi, "Neighborhood attention transformer," *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 6185–6194, 2022. 3, 11, 15, 25, 27

[28] Y. Hu, X. Jin, Y. Zhang, H. Hong, J. Zhang, Y. He, and H. Xue, "Rams-trans: Recurrent attention multi-scale transformer for fine-grained image recognition," *Proceedings of the 29th ACM International Conference on Multimedia*, 2021. 3, 17

[29] M. Ding, B. Xiao, N. Codella, P. Luo, J. Wang, and L. Yuan, "Davit: Dual attention vision transformers," in *European Conference on Computer Vision*, 2022, pp. 74–92. 3, 26, 27

[30] X. Chu, Z. Tian, B. Zhang, X. Wang, and C. Shen, "Conditional positional encodings for vision transformers," in *International Conference on Learning Representations*, 2021. 4, 8, 24, 28

[31] B. Wu, C. Xu, X. Dai, A. Wan, P. Zhang, M. Tomizuka, K. Keutzer, and P. Vajda, "Visual transformers: Token-based image representation and processing for computer vision," *ArXiv*, vol. abs/2006.03677, 2020. 4

[32] A. Khan, Z. Rauf, A. Sohail, A. R. Khan, H. Asif, A. Asif, and U. Farooq, "A survey of the vision transformers and their cnn-transformer based variants," *Artificial Intelligence Review*, vol. 56, pp. 2917 – 2970, 2023. 4, 10, 16

[33] Y. Li, C. Wu, H. Fan, K. Mangalam, B. Xiong, J. Malik, and C. Feichtenhofer, "Mvitv2: Improved multiscale vision transformers for classification and detection," *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4794–4804, 2021. 4, 27

[34] Y. Shi, M. Sun, Y. Wang, R. Wang, H. Sun, and Z. Chen, "Fvit: A focal vision transformer with gabor filter," *ArXiv*, vol. abs/2402.11303, 2024. 4, 15

[35] J. Ho, N. Kalchbrenner, D. Weissenborn, and T. Salimans, "Axial attention in multidimensional transformers," *ArXiv*, vol. abs/1912.12180, 2019. 4, 12, 20

[36] H. Zhang, W. Hu, and X. Wang, "Parc-net: Position aware circular convolution with merits from convnets and transformer," in *European Conference on Computer Vision*, 2022, pp. 613–630. 5, 17

[37] N. Engel, V. Belagiannis, and K. Dietmayer, "Point transformer," *IEEE Access*, vol. 9, pp. 134 826–134 840, 2020. 5, 21

[38] A. Srinivas, T.-Y. Lin, N. Parmar, J. Shlens, P. Abbeel, and A. Vaswani, "Bottleneck transformers for visual recognition," *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 16 514–16 524, 2021. 5, 24

[39] X. Chu, Z. Tian, Y. Wang, B. Zhang, H. Ren, X. Wei, H. Xia, and C. Shen, "Twins: Revisiting the design of spatial attention in vision transformers," *Neural Information Processing Systems*, 2021, pp. 9355–9366. 5

[40] Y. Rao, Z. Liu, W. Zhao, J. Zhou, and J. Lu, "Dynamic spatial sparsification for efficient vision transformers and convolutional neural networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, pp. 10 883–10 897, 2022. 5, 14

[41] Z. Kong, P. Dong, X. Ma, X. Meng, W. Niu, M. Sun, B. Ren, M. Qin, H. Tang, and Y. Wang, "Spvit: Enabling faster vision transformers via soft token pruning," *ArXiv*, vol. abs/2112.13890, 2021. 5, 29

[42] B. Chen, P. Li, C. Li, B. Li, L. Bai, C. Lin, M. Sun, J. Yan, and W. Ouyang, "Glit: Neural architecture search for global and local image transformer," *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 12–21, 2021. 5, 18

[43] Z. Pan, B. Zhuang, H. He, J. Liu, and J. Cai, "Less is more: Pay less attention in vision transformers," *ArXiv*, vol. abs/2105.14217, 2021. 6, 25

[44] K. Yuan, S. Guo, Z. Liu, A. Zhou, F. Yu, and W. Wu, "Incorporating convolution designs into visual transformers," *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 559–568, 2021. 6, 16, 23, 25, 30

[45] H. Fan, B. Xiong, K. Mangalam, Y. Li, Z. Yan, J. Malik, and C. Feichtenhofer, "Multiscale vision transformers," *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 6804–6815, 2021. 6, 10, 14, 22, 26, 27, 30

[46] G. Zhao, J. Lin, Z. Zhang, X. Ren, Q. Su, and X. Sun, "Explicit sparse transformer: Concentrated attention through explicit selection," *ArXiv*, vol. abs/1912.11637, 2019. 7, 22

[47] P. Wang, X. Wang, F. Wang, M. Lin, S. Chang, W. Xie, H. Li, and R. Jin, "Kvt: k-nn attention for boosting vision transformers," *ArXiv*, vol. abs/2106.00515, 2021. 7

[48] A. Hatamizadeh, G. Heinrich, H. Yin, A. Tao, J. Álvarez, J. Kautz, and P. Molchanov, "Fastervit: Fast vision transformers with hierarchical attention," *ArXiv*, vol. abs/2306.06189, 2023. 7, 22

[49] C. Yang, Y. Wang, J. Zhang, H. Zhang, Z. Wei, Z. L. Lin, and A. Yuille, "Lite vision transformer with enhanced self-attention," *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 11 988–11 998, 2021. 7, 15, 27

[50] C. Zhu, W. Ping, C. Xiao, M. Shoeybi, T. Goldstein, A. Anandkumar, and B. Catanzaro, "Long-short transformer: Efficient transformers for language and vision," in *Neural Information Processing Systems*, 2021, pp. 17 723–17 736. 7

[51] K. Choromanski, V. Likhosherstov, D. Dohan, X. Song, J. Davis, T. Sarlós, D. Belanger, L. J. Colwell, and A. Weller, "Masked language modeling for proteins via linearly scalable long-context transformers," *ArXiv*, vol. abs/2006.03555, 2020. 7, 26

[52] H. Li, M. Wang, S. Liu, and P.-Y. Chen, "A theoretical understanding of shallow vision transformers: Learning, generalization, and sample complexity," *ArXiv*, vol. abs/2302.06015, 2023. 7, 22

[53] D. A. Calian, P. Roelants, J. Calì, B. Carr, K. Dubba, J. E. Reid, and D. Zhang, "Scram: Spatially coherent randomized attention maps," *ArXiv*, vol. abs/1905.10308, 2019. 8, 11, 27

[54] J. Jiao, Y. Tang, K.-L. C. Lin, Y. Gao, J. Ma, Y. Wang, and W.-S. Zheng, "Dilateformer: Multi-scale dilated transformer for visual recognition," *IEEE Transactions on Multimedia*, vol. 25, pp. 8906–8919, 2023. 8

[55] P. Wang, S. Bai, S. Tan, S. Wang, Z. Fan, J. Bai, K.-Y. Chen, X. Liu, J. Wang, W. Ge, Y. Fan, K. Dang, M. Du, X. Ren, R. Men, D. Liu, C. Zhou, J. Zhou, and J. Lin, "Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution," *ArXiv*, vol. abs/2409.12191, 2024. 8, 17, 28

[56] B. Heo, S. Park, D. Han, and S. Yun, "Rotary position embedding for vision transformer," *ArXiv*, vol. abs/2403.13298, 2024. 8, 28

[57] H. You, Y. Xiong, X. Dai, B. Wu, P. Zhang, H. Fan, P. Vajda, and Y. Lin, "Castling-vit: Compressing self-attention via switching towards linear-angular attention at vision transformer inference," *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 14 431–14 442, 2022. 9

[58] W. Sun, Z. Qin, H. Deng, J. Wang, Y. Zhang, K. Zhang, N. Barnes, S. Birchfield, L. Kong, and Y. Zhong, "Vicinity vision transformer," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, pp. 12 635–12 649, 2022. 9

[59] W. Li, Y. Yuan, J. Liu, D. Tang, S. Wang, J. Zhu, and L. Zhang, "Tokenpacker: Efficient visual projector for multimodal llm," *ArXiv*, vol. abs/2407.02392, 2024. 9

[60] T. Chen, Y. Cheng, Z. Gan, L. Yuan, L. Zhang, and Z. Wang, "Chasing sparsity in vision transformers: An end-to-end exploration," in *Neural Information Processing Systems*, 2021, pp. 19 974–19 988. 9, 14, 29

[61] A. Chavan, Z. Shen, Z. Liu, Z. Liu, K.-T. Cheng, and E. P. Xing, "Vision transformer slimming: Multi-dimension searching in continuous optimization space," *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4921–4931, 2022. 9

[62] Z. Li, J. Xiao, L. Yang, and Q. Gu, "Repq-vit: Scale reparameterization for post-training quantization of vision transformers," *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 17 181–17 190, 2022. 9

[63] H. Wang, Z. Wu, Z. Liu, H. Cai, L. Zhu, C. Gan, and S. Han, "Hat: Hardware-aware transformers for efficient natural language processing," in *Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 7675–7688. 9, 25, 29

[64] H. You, Z. Sun, H. Shi, Z. Yu, Y. Zhao, Y. Zhang, C. Li, B. Li, and Y. Lin, "Vitcod: Vision transformer acceleration via dedicated algorithm and accelerator co-design," *2023 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*, pp. 273–286, 2022. 10, 13, 29

[65] X. Chen, S. Xie, and K. He, "An empirical study of training self-supervised vision transformers," *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 9620–9629, 2021. 10, 15

[66] L. Yuan, Y. Chen, T. Wang, W. Yu, Y. Shi, F. E. H. Tay, J. Feng, and S. Yan, "Tokens-to-token vit: Training vision transformers from scratch on imagenet," *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 538–547, 2021. 10, 14, 21

[67] M. Walmer, S. Suri, K. Gupta, and A. Shrivastava, "Teaching matters: Investigating the role of supervision in vision transformers," *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 7486–7496, 2022. 10

[68] D. Zhou, Z. Yu, E. Xie, C. Xiao, A. Anandkumar, J. Feng, and J. Álvarez, "Understanding the robustness in vision transformers," *ArXiv*, vol. abs/2204.12451, 2022. 10, 15, 16, 23, 30

[69] S. Chen, C. Ge, Z. Tong, J. Wang, Y. Song, J. Wang, and P. Luo, "Adaptformer: Adapting vision transformers for scalable visual recognition," *ArXiv*, vol. abs/2205.13535, 2022. 10

[70] H. Zhang, Y. Zhu, D. Wang, L. Zhang, T. Chen, and Z. Ye, "A survey on visual mamba," *ArXiv*, vol. abs/2404.15956, 2024. 11, 12

[71] H. Touvron, M. Cord, and H. J'egou, "Deit iii: Revenge of the vit," in *European Conference on Computer Vision*, 2022, pp. 516–533. 12

[72] Z. Li, L. Ma, M. Chen, J. Xiao, and Q. Gu, "Patch similarity aware data-free quantization for vision transformers," in *European Conference on Computer Vision*, 2022, pp. 154–170. 12, 20

[73] W. Zhao, J. Tang, Y. Han, Y. Song, K. Wang, G. Huang, F. Wang, and Y. You, "Dynamic tuning towards parameter and inference efficiency for vit adaptation," *ArXiv*, vol. abs/2403.11808, 2024. 12

[74] Y. Liang, C. Ge, Z. Tong, Y. Song, J. Wang, and P. Xie, "Not all patches are what you need: Expediting vision transformers via token reorganizations," *ArXiv*, vol. abs/2202.07800, 2022. 12, 29

[75] K. Mangalam, H. Fan, Y. Li, C. Wu, B. Xiong, C. Feichtenhofer, and J. Malik, "Reversible vision transformers," *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 10 820–10 830, 2022. 12, 16, 27

[76] J. F. Henriques and A. Vedaldi, "Warped convolutions: Efficient invariance to spatial transformations," in *International Conference on Machine Learning*, 2016, pp. 1461–1469. 12, 24

[77] J. Ding, N. Xue, Y. Long, G.-S. Xia, and Q. Lu, "Learning roi transformer for detecting oriented objects in aerial images," *ArXiv*, vol. abs/1812.00155, 2018. 13

[78] Q. Wan, Z. Huang, J. Lu, G. Yu, and L. Zhang, "Seaformer: Squeeze-enhanced axial transformer for mobile semantic segmentation," *ArXiv*, vol. abs/2301.13156, 2023. 13

[79] A. B. Koyuncu, H. Gao, and E. Steinbach, "Contextformer: A transformer with spatio-channel attention for context modeling in learned image compression," *ArXiv*, vol. abs/2203.02452, 2022. 13

[80] G. Luo, M. Huang, Y. Zhou, X. Sun, G. Jiang, Z. Wang, and R. Ji, "Towards efficient visual adaption via structural reparameterization," *ArXiv*, vol. abs/2302.08106, 2023. 13

[81] R. Azad, R. Arimond, E. K. Aghdam, A. Kazerouni, and D. Merhof, "Dae-former: Dual attention-guided efficient transformer for medical image segmentation," *ArXiv*, vol. abs/2212.13504, 2022. 13, 17, 26

[82] C. Renggli, A. S. Pinto, N. Houlsby, B. Mustafa, J. Puigcerver, and C. Riquelme, "Learning to merge tokens in vision transformers," *ArXiv*, vol. abs/2202.12015, 2022. 13

[83] F. Chen, Z. Luo, L. Zhou, X. Pan, and Y. Jiang, "Comprehensive survey of model compression and speed up for vision transformers," *ArXiv*, vol. abs/2404.10407, 2024. 13

[84] W. Xiang, C. Li, B. Wang, X. Wei, X. Hua, and L. Zhang, "Spatiotemporal self-attention modeling with temporal patch shift for action recognition," in *European Conference on Computer Vision*, 2022, pp. 627–644. 15

[85] S. Mehta and M. Rastegari, "Separable self-attention for mobile vision transformers," *ArXiv*, vol. abs/2206.02680, 2022. 15

[86] H. Wu, J. Wu, J. Xu, J. Wang, and M. Long, "Flowformer: Linearizing transformers with conservation flows," in *International Conference on Machine Learning*, 2022, pp. 24 226–24 242. 15

[87] S. Jamil, M. J. Piran, and O.-J. Kwon, "A comprehensive survey of transformers for computer vision," *ArXiv*, vol. abs/2211.06004, 2022. 16

[88] Q. Hou, C. Lu, M.-M. Cheng, and J. Feng, "Conv2former: A simple transformer-style convnet for visual recognition," *IEEE transactions on pattern analysis and machine intelligence*, vol. PP, 2022. 16

[89] S. Yun, H. Lee, J. Kim, and J. Shin, "Patch-level representation learning for self-supervised vision transformers," *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 8344–8353, 2022. 16

[90] S. Jelassi, M. E. Sander, and Y.-F. Li, "Vision transformers provably learn spatial structure," *ArXiv*, vol. abs/2210.09221, 2022. 16

[91] C. Li, B. Zhuang, G. Wang, X. Liang, X. Chang, and Y. Yang, "Automated progressive learning for efficient training of vision transformers," *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 12 476–12 486, 2022. 16, 22, 26, 30

[92] R. Hu, A. Singh, T. Darrell, and M. Rohrbach, "Iterative answer prediction with pointer-augmented multimodal transformers for textvqa," *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 9989–9999, 2019. 16, 20

[93] C. Zhu, Y. Zhou, Y. Shen, G. Luo, X. Pan, M. Lin, C. Chen, L. Cao, X. Sun, and R. Ji, "Seqtr: A simple yet universal network for visual grounding," in *European Conference on Computer Vision*, 2022, pp. 598–615. 16, 17

[94] L. Zhu, B. Liao, Q. Zhang, X. Wang, W. Liu, and X. Wang, "Vision mamba: Efficient visual representation learning with bidirectional state space model," *ArXiv*, vol. abs/2401.09417, 2024. 17, 28

[95] T. Huang, X. Pei, S. You, F. Wang, C. Qian, and C. Xu, "Localmamba: Visual state space model with windowed selective scan," *ArXiv*, vol. abs/2403.09338, 2024. 17

[96] R. Esfandiarpoor, C. Menghini, and S. H. Bach, "If clip could talk: Understanding vision-language model representations through their preferred concept descriptions," *ArXiv*, vol. abs/2403.16442, 2024. 17

[97] F. Li, A. Zeng, S. Liu, H. Zhang, H. Li, L. Zhang, and L. Ni, "Lite detr : An interleaved multi-scale encoder for efficient detr," *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 18 558–18 567, 2023. 17

[98] O. F. Kar, A. Tonioni, P. Poklukar, A. Kulshrestha, A. Zamir, and F. Tombari, "Brave: Broadening the visual encoding of vision-language models," *ArXiv*, vol. abs/2404.07204, 2024. 17, 28

[99] D. Du, G. Gong, and X. Chu, "Model quantization and hardware acceleration for vision transformers: A comprehensive survey," *ArXiv*, vol. abs/2405.00314, 2024. 18, 25

[100] A. Arnab, M. Dehghani, G. Heigold, C. Sun, M. Lucic, and C. Schmid, "Vivit: A video vision transformer," *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 6816–6826, 2021. 18

[101] A. Ulhaq, N. Akhtar, G. Pogrebna, and A. Mian, "Vision transformers for action recognition: A survey," *ArXiv*, vol. abs/2209.05700, 2022. 18

[102] S. Tang, J. Zhang, S. Zhu, and P. Tan, "Quadtree attention for vision transformers," *ArXiv*, vol. abs/2201.02767, 2022. 18, 21, 25, 30

[103] C. Tang, Y. Zhao, G. Wang, C. Luo, W. Xie, and W. Zeng, "Sparse mlp for image recognition: Is self-attention really necessary?" *ArXiv*, vol. abs/2109.05422, 2021. 19

[104] M. Raghu, T. Unterthiner, S. Kornblith, C. Zhang, and A. Dosovitskiy, "Do vision transformers see like convolutional neural networks?" in *Neural Information Processing Systems*, 2021, pp. 12 116–12 128. 19, 23

[105] M. Kim, P. H. Seo, C. Schmid, and M. Cho, "Learning correlation structures for vision transformers," *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 18 941–18 951, 2024. 20

[106] L. Song, S. Zhang, S. Liu, Z. Li, X. He, H. Sun, J. Sun, and N. Zheng, "Dynamic grained encoder for vision transformers," in *Neural Information Processing Systems*, 2023, pp. 5770–5783. 20

[107] H. Touvron, M. Cord, A. El-Nouby, J. Verbeek, and H. J'egou, "Three things everyone should know about vision transformers," *ArXiv*, vol. abs/2203.09795, 2022. 20, 23

[108] A. Sahiner, T. Ergen, B. M. Ozturkler, J. Pauly, M. Mardani, and M. Pilanci, "Unraveling attention via convex duality: Analysis and interpretations of vision transformers," *ArXiv*, vol. abs/2205.08078, 2022. 20, 23

[109] Z. Chen, B. Li, J. Xu, S. Wu, S. Ding, and W. Zhang, "Towards practical certifiable patch defense with vision transformer," *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 15 127–15 137, 2022. 20

[110] S. Herdade, A. Kappeler, K. Boakye, and J. Soares, "Image captioning: Transforming objects into words," in *Neural Information Processing Systems*, 2019, pp. 11 135–11 145. 20

[111] K. Wu, H. Peng, M. Chen, J. Fu, and H. Chao, "Rethinking and improving relative position encoding for vision transformer," *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 10 013–10 021, 2021. 21, 24

[112] Y. Gandelsman, A. A. Efros, and J. Steinhardt, "Interpreting clip's image representation via text-based decomposition," *ArXiv*, vol. abs/2310.05916, 2023. 21, 24

[113] M. Dehghani, J. Djolonga, B. Mustafa, P. Padlewski, J. Heek, J. Gilmer, A. Steiner, M. Caron, R. Geirhos, I. M. Alabdulmohsin, R. Jenatton, L. Beyer, M. Tschannen, A. Arnab, X. Wang, C. Riquelme, M. Minderer, J. Puigcerver, U. Evci, M. Kumar, S. van Steenkiste, G. F. Elsayed, A. Mahendran, F. Yu, A. Oliver, F. Huot, J. Bastings, M. Collier, A. Gritsenko, V. Birodkar, C. Vasconcelos, Y. Tay, T. Mensink, A. Kolesnikov, F. Paveti'c, D. Tran, T. Kipf, M. Luvci'c, X. Zhai, D. Keysers, J. Harmsen, and N. Houlsby, "Scaling vision transformers to 22 billion parameters," *ArXiv*, vol. abs/2302.05442, 2023. 22

[114] P. Ramachandran, N. Parmar, A. Vaswani, I. Bello, A. Levskaya, and J. Shlens, "Stand-alone self-attention in vision models," *ArXiv*, vol. abs/1906.05909, 2019. 22

[115] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *CoRR*, vol. abs/1412.6980, 2014. 22

[116] J. Zhang, X. Li, Y. Wang, C. Wang, Y. Yang, Y. Liu, and D. Tao, "Eatformer: Improving vision transformer inspired by evolutionary algorithm," *ArXiv*, vol. abs/2206.09325, 2022. 23

[117] K. L. Ong, C. Lee, H. S. Lim, K. Lim, and A. Alqahtani, "Melmvitv2: Enhanced speech emotion recognition with mel spectrogram and improved multiscale vision transformers," *IEEE Access*, vol. 11, pp. 108 571–108 579, 2023. 23, 27

[118] Z. Yang, J. Wang, Y. Tang, K. Chen, H. Zhao, and P. H. S. Torr, "Lavt: Language-aware vision transformer for referring image segmentation," *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 18 134–18 144, 2021. 24

[119] A. Handa, M. Blösch, V. Patraucean, S. Stent, J. McCormac, and A. Davison, "gvnn: Neural network library for geometric computer vision," in *ECCV Workshops*, 2016, pp. 67–82. 24

[120] Y. Tay, M. Dehghani, D. Bahri, and D. Metzler, "Efficient transformers: A survey," *ACM Computing Surveys*, vol. 55, pp. 1 – 28, 2020. 24

[121] R. Pope, S. Douglas, A. Chowdhery, J. Devlin, J. Bradbury, A. Levskaya, J. Heek, K. Xiao, S. Agrawal, and J. Dean, "Efficiently scaling transformer inference," *ArXiv*, vol. abs/2211.05102, 2022. 24, 29

[122] X. Liu, H. Peng, N. Zheng, Y. Yang, H. Hu, and Y. Yuan, "Efficientvit: Memory efficient vision transformer with cascaded group attention," *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 14 420–14 430, 2023. 25

[123] K. Xu, Z. Wang, C. Chen, X. Geng, J. Lin, X. Yang, M. Wu, X. Li, and W. Lin, "Lpvit: Low-power semi-structured pruning for vision transformers," *ArXiv*, vol. abs/2407.02068, 2024. 26

[124] H. Ren, H. Dai, Z. Dai, M. Yang, J. Leskovec, D. Schuurmans, and B. Dai, "Combiner: Full attention transformer with sparse computation cost," in *Neural Information Processing Systems*, 2021, pp. 22 470–22 482. 27

[125] J. He, J. Chen, S. Liu, A. Kortylewski, C. Yang, Y. Bai, C. Wang, and A. Yuille, "Transfg: A transformer architecture for fine-grained recognition," in *AAAI Conference on Artificial Intelligence*, 2021, pp. 852–860. 27

[126] R. Ranftl, A. Bochkovskiy, and V. Koltun, "Vision transformers for dense prediction," *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 12 159–12 168, 2021. 27

[127] M. E. Sadeghi, A. Fayyazi, S. Azizi, and M. Pedram, "Peanovit: Power-efficient approximations of non-linearities in vision transformers," *ArXiv*, vol. abs/2406.14854, 2024. 28

[128] M. Dai, J. Hu, J. Zhuang, and E. Zheng, "A transformer-based feature segmentation and region alignment method for uav-view geo-localization," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 32, pp. 4376–4389, 2022. 28

[129] C. Han, J. Liang, Q. Wang, M. Rabbani, S. Dianat, R. M. Rao, Y. N. Wu, and D. Liu, "Image translation as diffusion visual programmers," *ArXiv*, vol. abs/2401.09742, 2024. 28

[130] Y. Tang, K. Han, Y. Wang, C. Xu, J. Guo, C. Xu, and D. Tao, "Patch slimming for efficient vision transformers," *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 12 155–12 164, 2021. 29

[131] D. Zhou, Y. Shi, B. Kang, W. Yu, Z. Jiang, Y. Li, X. Jin, Q. Hou, and J. Feng, "Refiner: Refining self-attention for vision transformers," *ArXiv*, vol. abs/2106.03714, 2021. 29

[132] X. Ye, Y. Gan, X. Huang, Y. Ge, Y. Shan, and Y. Tang, "Voco-llama: Towards vision compression with large language models," *ArXiv*, vol. abs/2406.12275, 2024. 29

[133] A. Parvaiz, M. A. Khalid, R. Zafar, H. Ameer, M. Ali, and M. Fraz, "Vision transformers in medical computer vision - a contemplative retrospection," *Eng. Appl. Artif. Intell.*, vol. 122, p. 106126, 2023. 30