

iiANET: Inception Inspired Attention Hybrid Network for efficient Long-Range Dependency

Anonymous authors

Paper under double-blind review

Abstract

The recent emergence of hybrid models has introduced a transformative approach to computer vision, gradually moving beyond conventional convolutional neural networks and vision transformers. However, efficiently combining these two approaches to better capture long-range dependencies in complex images remains a challenge. In this paper, we present iiANET (Inception Inspired Attention Network), an efficient hybrid visual backbone designed to improve the modeling of long-range dependencies in complex visual recognition tasks. The core innovation of iiANET is the iiABlock, a unified building block that integrates a modified global r-MHSA (Multi-Head Self-Attention) and convolutional layers in parallel. This design enables iiABlock to simultaneously capture global context and local details, making it effective for extracting rich and diverse features. By efficiently fusing these complementary representations, iiABlock allows iiANET to achieve strong feature interaction while maintaining computational efficiency. Extensive qualitative and quantitative evaluations on some SOTA benchmarks demonstrate improved performance.

1 Introduction

From autonomous drones to urban planning, understanding complex visual scenes is more critical than ever, yet traditional models fall short. Over the last decade, deep Convolutional Neural Network (CNN) architectures have emerged as the de facto standard for solving most computer vision (CV) tasks, including image classification He et al. (2016); Tan & Le (2019), object detection Ren et al. (2015); Redmon & Farhadi (2017) and segmentation Long et al. (2015) with compelling results. The prevalence of CNN architectures is not coincidental, as they excel at capturing spatial features and patterns in images. However, the dominance of CNN architectures is being challenged by the emergence of ViT (Vision in Transformer) Dosovitskiy et al. (2020), presenting a transformative approach to solving CV tasks. Interestingly, this groundbreaking model outperforms SOTA CNN-based models on ImageNet benchmark Dosovitskiy et al. (2020) and emerges as a competitive alternative Han et al. (2022). Practically, ViT works exactly like the text-based Natural Language Processing (NLP) transformers but with patch embedding. It divides the input image into patches, projects them into a high-dimensional feature space through a linear projection layer, adds positional embedding, passes them through a transformer encoder, and finally maps the output to a fixed-length vector for classification tasks.

Significantly, the key component of ViT is the self-attention mechanism Dosovitskiy et al. (2020) within the encoder, which enables the model to capture long-range dependencies by allowing each element in the input sequence to attend to all other elements, considering their relative importance Dosovitskiy et al. (2020). While this capability allows the model to selectively focus on distantly related pixels, facilitating the efficient capture of contextual information across the entire input sequence, it encounters limitations such as increased computational complexity, reduced interpretability, data hungry, and challenges in handling spatial information effectively compared to CNNs Haruna et al. (2025). In contrast, CNN-based models, while effective at capturing local features through parameter sharing and local receptive fields, struggle with capturing long-range dependencies, limiting their ability to integrate distant pixel relationships Haruna et al. (2025). These limitations have led to the development of hybrid models, which combine their strengths to improve performance Haruna et al. (2025).

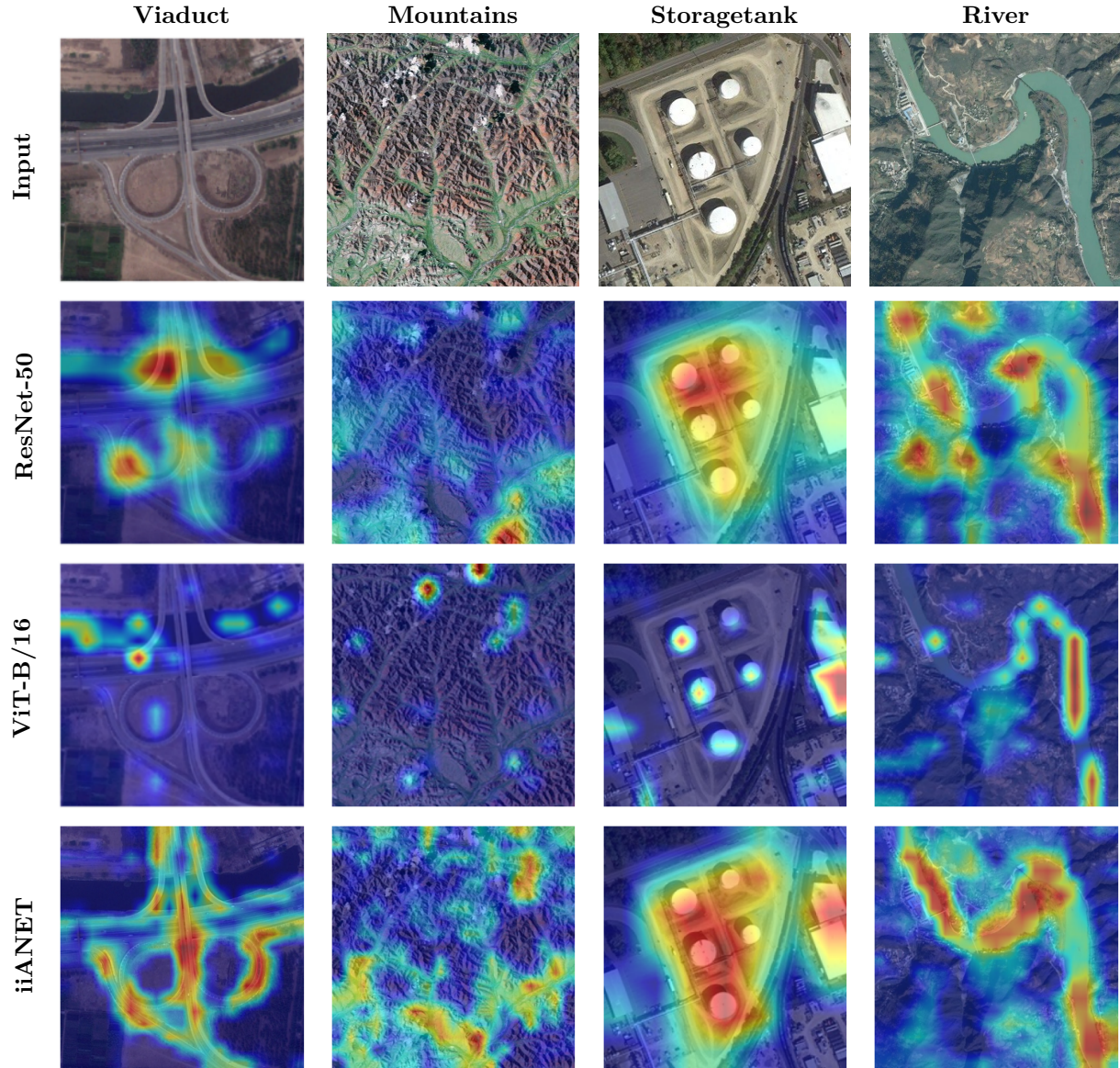


Figure 1: iiANET Grad-CAM comparison and other state-of-the-art models, e.g., (a) shows an aerial image of viaduct, mountain, storage tanks, and river featuring complex infrastructure consisting of multiple spans, roads, and surrounding landscapes. The primary objective is to accurately detect and classify various elements to facilitate efficient maintenance, safety management, and infrastructure planning. Consequently, capturing long-range dependencies in this scenario is crucial for comprehending the spatial layout of different viaducts, mountains, storage tanks, and rivers, their interactions, and potential structural issues. (b) illustrates ResNet-50’s inability to capture long-range dependencies, but local features (c) demonstrate the limitations of ViT-B/16 interpretability, as it primarily highlights tiny spots on the image. (d) It is a hybrid model, depicting notable improvements in capturing long-range dependencies, global context, and improved interpretability.

Specifically, previous hybrid designs have aimed to enhance capturing long-range dependencies for various CV tasks Guo et al. (2022); Srinivas et al. (2021); Dai et al. (2021a). However, the design of hybrid models introduces additional design complexities Haruna et al. (2025), and computational costs compared to monolithic models Khan et al. (2023), while also potentially leading to information loss due to feature fusion of distinct models Haruna et al. (2025). Lastly, more effort is needed to design efficient hybrid models capable of capturing long-range dependencies in complex images, a challenge that remains largely unaddressed.

In this work, we propose a novel architectural block, the iiABlock (Inception-Inspired Attention Block), which serves as the core component of our model, iiANET. The iiABlock is a carefully designed hybrid module that integrates parallel convolutional layers for efficient local feature extraction, and a global 2D Multi-Head Self-Attention (MHSA) mechanism with Registers to effectively model long-range dependencies. The outputs from these branches are then fused via concatenation and feature shuffling, enabling rich interaction between local and global features. By leveraging the complementary strengths of CNNs and transformers in a lightweight design, iiANET offers a simple yet powerful solution for understanding complex visual scenes with long-range dependencies. For example, on the AID (Aerial Image Dataset) Xia et al. (2017), iiANET-B and iiANET-L achieve an accuracy of 80.57% and 83.11% respectively, outperforming ResNet-50 (71.93%), ViT-B/224 (69.93%), and DiNAT-B (79.12%). These results highlight iiANET’s effectiveness in modeling long-range dependencies in challenging datasets. The contributions of this paper are summarized as follows:

- We introduce iiABlock, a novel hybrid module that integrates parallel convolutional branches with global rMHSA, enabling efficient capture of long-range dependencies in complex vision tasks.
- We propose iiANET, a robust and efficient backbone model for visual recognition downstream tasks, such as object detection and semantic segmentation.
- Extensive experimental results on commonly used benchmarks demonstrate that iiANET outperforms some existing SOTA methods.

2 Related Work

CNN-based methods have seen various attempts to enhance their ability to capture long-range dependencies in images. Donahue et al. (2015) introduced the Long-term Recurrent Convolutional Network (LRCN) by fusing CNNs with LSTMs, while Yu et al. (2017) proposed the Dilated Residual Network (DRN) using multiple dilation rates to expand receptive fields, and Yu & Koltun (2015) designed a Dilated Convolution (DC) model to improve global context in semantic segmentation. These approaches advanced CNN capacity for long-range dependencies but face limitations: LRCN increases computational complexity due to recurrent connections, DRN can lose fine-grained spatial details from varying dilation rates, and DC suffers from gridding artifacts. The emergence of ViT Dosovitskiy et al. (2020); Liu et al. (2021a) offered a breakthrough in capturing long-range dependencies via attention mechanisms, achieving SOTA performance. However, their quadratic complexity, high data requirements, and weaker inductive bias compared to CNN demand substantial computational resources. To address these trade-offs, hybrid methods combine CNN feature extraction with ViT global dependency modeling Haruna et al. (2025). Zhang et al. (2022) proposed ELAN, using group-wise multi-scale self-attention for super-resolution; Guo et al. (2022) introduced CMT (CNN Meet ViT) to integrate attention into CNN blocks; and Srinivas et al. (2021) developed BoTNet by replacing the final ResNet block with MHSA. While effective, these methods often apply attention at later stages with smaller spatial dimensions, limiting effectiveness, and face challenges such as memory constraints (CMT-L), structural complexity, and information loss from fusing distinct methods Peng et al. (2021); Dai et al. (2021b); Hassani & Shi (2022); Wu et al. (2021).

Overall, CNN excel at local detail capture but struggle with long-range dependencies Haruna et al. (2025); Khan et al. (2023), RNNs handle such dependencies but lack parallelism and train slowly Banerjee et al. (2019), and ViT capture them efficiently but require more memory Dosovitskiy et al. (2020). However, it remains a challenge to efficiently combine CNN and ViT architectures due to design complexity, higher computational costs, feature fusion losses, and interpretability issues. This paper addresses this gap by proposing a hybrid model that efficiently captures long-range dependencies in complex images.

3 Method

3.1 Our Approach: iiABlock

The iiABlock is the core component of the proposed iiANET, designed to capture both local and global features in complex visual scenes. It combines parallel convolutional layers for efficient extraction of localized features with a global 2D r-MHSA mechanism augmented by Register tokens to model long-range dependencies. The convolutional branch leverages parameter sharing and local receptive fields, while 2D r-MHSA attends to distant pixel relationships across the input, with Register tokens to enhance interpretability. To merge these complementary features, we introduce a lightweight fusion strategy using feature concatenation followed by channel shuffling, enabling rich local-global interactions with minimal computational cost. This balanced design offers a favorable trade-off between speed and accuracy, providing a robust backbone for downstream visual recognition tasks. Figure 2 illustrates the iiABlock architecture.

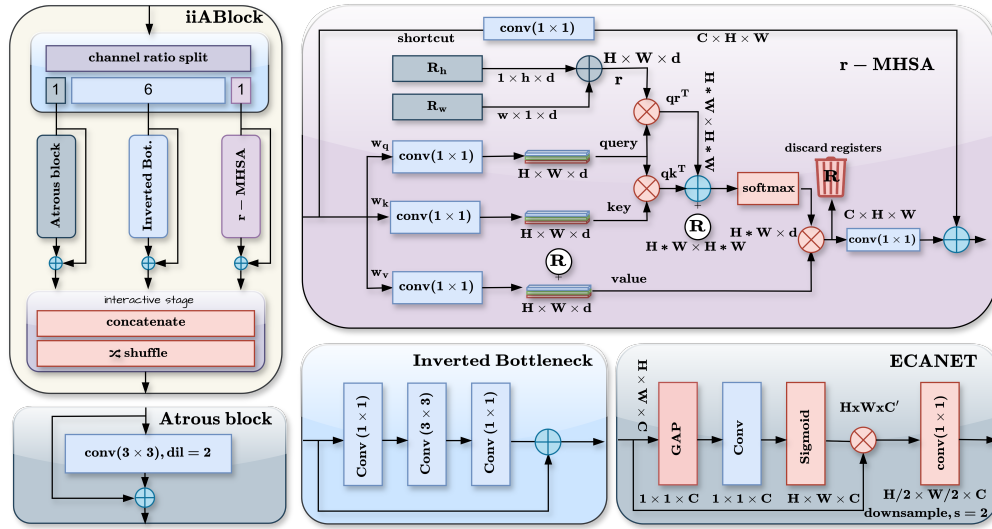


Figure 2: iiABLOCK design showing r-MHSA, inverted bottleneck, ECANET and Atrous block

3.2 Local Details

To extract fine-grained features and spatial patterns from complex images, iiABlock introduces components for modeling local details. This is important for recognizing textures, edges, and region-specific patterns.

Inverted Bottleneck (Efficient Convolutional Block) iiABlock utilizes the inverted bottleneck in parallel to improve computational efficiency and enhance local feature extraction, consisting of 1×1 convolutions for dimensionality reduction, 3×3 depth-wise separable convolution for spatial information extraction, and 1×1 convolution for projection. Notably, this block is limited to capturing local context with a fixed kernel, making it less effective in understanding prevalent global context in complex images Dai et al. (2021a), e.g., road, viaduct, bridge. Given the depth-wise operation in equation 1.

$$y_i = \sum_{j \in \mathcal{L}(i)} w_{i-j} \cdot x_j \quad (1)$$

y_i calculates the output at position i by taking a weighted sum of the input x_i , where $x_i, y_i \in \mathbb{R}^D$. The weights w_{i-j} determine the contribution of each input x_j to the output, and $\mathcal{L}(i)$ represents a local neighborhood, typically a 3×3 grid centered around i . The small size of $\mathcal{L}(i)$ limits the receptive field’s ability to capture intricate details, particularly in complex images with prevalent long-range dependencies. To mitigate the limitation of the inverted bottleneck in capturing long-range dependencies, we also introduce atrous convolution into the iiABlock.

Atrous Convolution (*Expanding Receptive Field*) In iiABlock, a single 3×3 atrous convolution expands the receptive field without increasing parameters. Unlike standard convolution, it applies a dilation rate r to space kernel elements, covering a wider area while preserving resolution and computational efficiency (See equation 2). Given an input feature map x and filter w , the output at location i is:

$$y[i] = \sum_{k=1}^K x[i + r \cdot k] \cdot w[k] \quad (2)$$

This captures mid-range dependencies and enriches contextual understanding, bridging local and global representations. However, it remains insufficient for fully modeling global, long-range dependencies in complex scenes, which are further addressed by integrating a global r-MHSA module into iiABlock.

3.3 Global Details

r-MHSA (*Capturing global context and long-range dependencies*). To capture global context and long-range dependencies in complex images, iiABlock integrates a modified global 2D r-MHSA mechanism. Unlike CNN layers limited to local receptive fields, r-MHSA allows each spatial location to attend to all others, effectively modeling contextual relationships across the entire image. Given a 2D input feature $x \in \mathbb{R}^{H \times W \times C}$ reshaped into $x \in \mathbb{R}^{HW \times d}$ (where d is the feature dimension), linear projections generate queries $Q = xW_Q$, keys $K = xW_K$, and values $V = xW_V$. Attention Z with h heads is computed as shown in Equation 3, enabling each token to attend to all others and capture long-range dependencies.

$$Z_h(Q_i, K, V) = \text{softmax} \left(\frac{Q_i K^\top}{\sqrt{d_k^h}} \right) V \quad (3)$$

Here, Q_i interacts with all keys in K across the entire sequence, unlike the standard MHSA, which attends within a limited context window. This enables the model to consider the relationships between all tokens regardless of their positional distance, capturing the global context and long-range dependencies prevalent in complex images. The **softmax** operation normalizes these scores and produces attention weights, which, when applied to the value matrix V , compute the final attended values. However, this interaction is order-agnostic and doesn't capture positional relationships in the input sequence. Therefore, in image data where spatial information is essential, integrating positional encodings is necessary to effectively complement MHSA.

Relative Position Encoding Srinivas et al. (2021) MHSA is permutation equivariant with no positional encoding. This characteristic limits its representational power, particularly for vision tasks involving highly structured data like images. Notably, it is added to the input image representation before the MHSA is applied, and it is used to guide the attention weights to focus on relevant pixels based on their relative positions in the input image.

$$Z_h(Q_i, K, V) = \text{softmax} \left(\frac{Q_i K^\top + Q_i R}{\sqrt{d_k^h}} \right) V \quad (4)$$

Where R is a trainable matrix. Lastly, reshape $Z(x)_h$ back to its original spatial shape of $x \in \mathbb{R}^{H \times W \times C}$. This addresses MHSA's order-agnostic nature, enhancing its representation power.

Registers (*Improving interpretability*). While the self-attention mechanism significantly improves the network's ability to capture long-range dependencies, it struggles with poor interpretability Darcet et al. (2023). We instead add additional learnable tokens to mitigate prevalent artifacts in the attention mechanism caused by high norms in image areas with low information during inference or training, similar to the implementation by Darcet et al. (2023). In this case, it is a global MHSA where the attention mechanism has a single input image, in contrast to having several patches. We initialize the register tokens for queries and keys as $R_{qk} \in \mathbb{R}^{N \times HW \times HW}$, where N is the number of register tokens and HW is the spatial dimension, then

value register tokens as $R_v \in \mathbb{R}^{N \times \frac{D}{\text{head}} \times HW}$, where D is the dimension. The operations expand both R_{qk} and R_v to $\mathbb{R}^{B \times N \times H \times W}$ and $\mathbb{R}^{N \times \frac{D}{\text{head}} \times HW}$, respectively (equation 5 and 6) effectively creating B copies of the register tokens for each batch, where B is the batch size. This step ensures that each batch has its own set of register tokens, facilitating batch-wise parallel processing in the attention mechanism.

$$R_{qk} = \text{repeat}(R'_{qk}, nhw \rightarrow bnhw', b = B) \quad (5)$$

$$R_v = \text{repeat}(R'_v, nhw \rightarrow bnhw', b = B) \quad (6)$$

Then integrate the register tokens into the attention mechanism as $Q_i K_R^\top = Q_i K^\top + R_{qk}$ and $V_R = V + R_v$ before computing the attention, where Z_R is the final 2D-MHSA output with registers. After computation, the register tokens are discarded (See equation 7).

$$Z_R^h(Q_i, K, V) = \text{softmax} \left(\frac{Q_i K_R^\top + Q_i R}{\sqrt{d_k^h}} \right) V_R \quad (7)$$

3.4 Features Recalibration

ECANET (*Channel-wise Recalibration and Down-Sampling*). While long-range dependencies in computer vision span both spatial and channel dimensions, traditional CNNs and attention mechanisms often emphasize spatial adaptivity, neglecting channel-wise adaptability. To address this, iiABlock integrates ECANET Wang et al. (2020), which efficiently computes channel attention weights by modeling inter-channel relationships. Compared to SENET Hu et al. (2018), ECANET offers improved efficiency, scalability, and accuracy. Given input $x \in \mathbb{R}^{H \times W \times C}$, adaptive average pooling reduces spatial dimensions to $p \in \mathbb{R}^{1 \times 1 \times C}$, followed by a 1×1 convolution and sigmoid activation to produce attention weights $Z \in [0, 1]$. These weights modulate channel significance through element-wise multiplication: $ECA = x \otimes Z$. Feature maps are then down-sampled via a 1×1 convolution with stride 2 after each stage.

3.5 Feature Interaction Fusion

iiABlock enables efficient multi-scale feature learning through a structured feature interaction mechanism that splits input channels into three branches—r-MHSA, inverted bottleneck, and atrous convolution processes them separately, then fuses outputs via concatenation followed by channel shuffling to enhance cross-branch interaction. This design captures short-range, long-range, and channel-specific dependencies in parallel while maintaining simplicity and low computational cost, outperforming more complex fusion strategies like heavy cross-attention Chen et al. (2021) or multi-level stacking Li et al. (2019). Empirically, we found the channel ratio $r = (1:6:1)$ optimal. For input $x \in \mathbb{R}^{H \times W \times C}$, channels are split as x_1, x_2, x_3 with dimensions approximately $\lfloor C/8 \rfloor$, $\lfloor 3C/4 \rfloor$, and $\lfloor C/8 \rfloor$, respectively. Corresponding functions f_1, f_2 , and f_3 process each branch, whose outputs are concatenated (Figure 3). Finally, channel shuffling is applied to strengthen inter-branch interactions, defined as $\hat{x} = \text{shuffle}(f(x)) \in \mathbb{R}^{H \times W \times C}$. The fused output \hat{x} effectively combines diverse receptive fields and attention mechanisms with efficiency and expressiveness.

3.6 iiANET Architectural Overview

iiANET is a hybrid visual recognition backbone architecture designed to enhance capturing long-range dependencies in complex images while maintaining computational efficiency. At each stage, iiANET stacks iiABlock in parallel across four stages. Each iiABlock captures both short-range and long-range dependencies while incorporating global context, enabling the model to process intricate patterns effectively. The architecture is also non-isotropic as it down-samples spatial features at every stage, allowing for a progressive abstraction of features. By stacking iiABlock in parallel at each stage, iiANET efficiently captures both fine-grained local details and global context, making it well-suited for complex vision tasks. Figure 3 shows the iiANET architectural design in detail. **Stem (Initial stage)** Given the higher resolutions of complex images, this component serves to compress the computational costs of iiANET by shrinking the spatial dimensions of the input image to half, trading off spatial details for improved model efficiency and

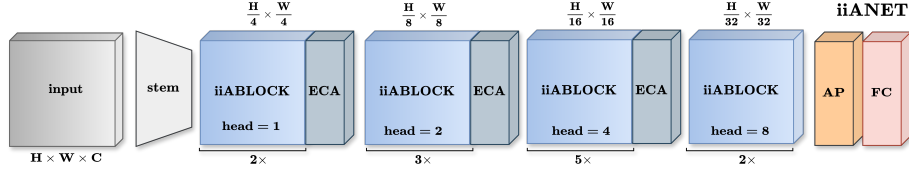


Figure 3: iiANET architectural overview.

basic feature extraction Bello et al. (2021). Let the input image be $x \in \mathbb{R}^{H \times W \times 3}$, we apply two sequential 3×3 convolutional layers, each followed by a batch normalization and ReLU activation function, with the initial layer using a stride of 2.

iiABlock (Main Building Blocks) The iiANET-B and iiANET-L variants stack iiABlocks in non-isotropic configurations of $[2, 3, 5, 2]$ and $[2, 3, 10, 4]$ across four stages to capture multi-scale features while reducing spatial resolution. At each stage, input feature maps are down-sampled, yielding spatial dimensions of $H/4 \times W/4$, $H/8 \times W/8$, $H/16 \times W/16$, and $H/32 \times W/32$. Inspired by Inception designs, multiple iiABlocks run in parallel per stage, enabling simultaneous processing of features with varied receptive fields. These branches are fused via concatenation followed by channel shuffling to promote cross-path interaction. This hierarchical, multi-branch structure allows iiANET to efficiently capture both short- and long-range dependencies while maintaining a compact, computationally efficient architecture suited for complex vision tasks.

Output Layer. After the final stage, Adaptive Pooling reduces the feature map to a fixed size, which is then passed through a Fully Connected layer for classification.

4 Experimental Results and Comparisons

We evaluate iiANET qualitatively and quantitatively on several widely used benchmark datasets, comparing it with state-of-the-art CNN, ViT, and hybrid models. The evaluation covers classification performance and effectiveness as a backbone for object detection and segmentation, focusing on iiANET’s ability to capture long-range dependencies in complex images.

Datasets and Metrics. Experiments are conducted on diverse datasets: AID Xia et al. (2017) (10,000 images, 30 scene classes), Oxford-III Parkhi et al. (2012) (4,978 images, 37 cat and dog breeds), and RLD Sathy et al. (2020) (5,932 images, 4 disease categories). Additionally, COCO-2017 Lin et al. (2014) and ImageNet1K Deng et al. (2009) assess generalization and robustness. Metrics include top-1/top-5 accuracy, Average Precision (AP), FLOPs, and throughput, providing comprehensive performance insights.

Experimental Setup. Training was performed on a Linux system with an Intel Core i7-8700K CPU, 2 NVIDIA Titan XP GPUs (12GB), and 32GB RAM. Models were trained for 90 or 150 epochs with batch size 16 using the AdamW optimizer, an initial learning rate of 0.0001, and a decay rate of 0.05. For fair comparison, baseline models were re-trained on AID, Oxford-III, and RLD using the authors’ default settings.

4.1 Qualitative Evaluation and Comparison: iiANET Visual Inspection

We applied Grad-CAM Selvaraju et al. (2017) on the final layer of iiANET and several state-of-the-art models—ResNet-50/100 He et al. (2016), EfficientNet-B4/B5 Tan & Le (2019), DenseNet-169/201 Huang et al. (2017), ViT-B/L-16 Dosovitskiy et al. (2020), CoatNet-3 Dai et al. (2021a), and BoTNet Srinivas et al. (2021) with visualizations shown in Figure 4. CNN-based models generate localized heatmaps around objects due to their limited receptive fields, while ViT-based models produce scattered attention spots, indicating interpretability challenges. Hybrid models better capture long-range dependencies and improve interpretability. Notably, iiANET excels at precisely outlining complex objects with minimal background noise, suggesting enhanced accuracy and reliability in tasks reliant on long-range dependencies, such as medical imaging, autonomous driving, remote sensing, and security surveillance.

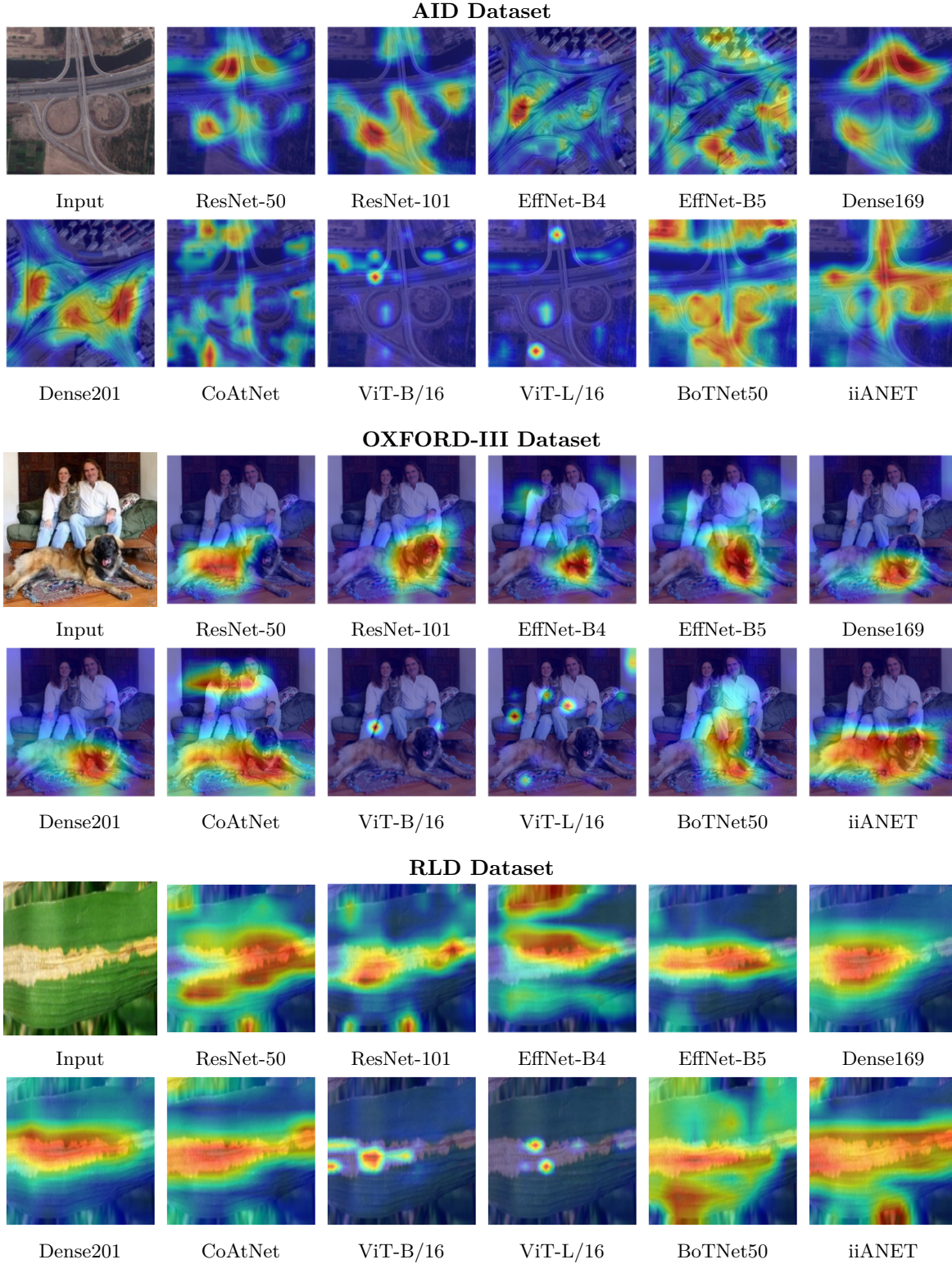


Figure 4: Visual inspection of iiANET compared to some SOTA models using Grad-CAM Selvaraju et al. (2017) highlights the model’s ability to focus on complex object regions. The heatmaps demonstrate iiANET’s improved capacity to capture long-range dependencies and provide more interpretable attention on relevant spatial structures compared to ResNet-50, BoTNet, ViT-B/16, and CoAtNet.

4.2 Quantitative Evaluation and Comparison

Table 1 shows iiANET classification performance across four standard benchmark datasets. The quantitative evaluation involves comparing iiANET-B and iiANET-L with several recent state-of-the-art models in terms of both accuracy and computational cost on ImageNet-1K, AID, Oxford-III, and RLD datasets.

(1) ImageNet-1K: iiANET-L achieves a competitive top-1 accuracy of 84.9% while maintaining computational efficiency with only 17.07 GFLOPs and 52.3M parameters. This demonstrates a decent accuracy-efficiency tradeoff, outperforming larger and more computationally expensive models such as ViT-L/16 (76.5%, 59.69 GFLOPs) and CoAtNet-3 (84.5%, 32.53 GFLOPs). These results show iiANET’s ability to learn robust, discriminative features with minimal complexity.

(2) AID Dataset: Contains complex spatial structures and long-range dependencies, which challenge traditional models. iiANET shows good performance, achieving 80.57% top-1 accuracy using only 8.22 GFLOPs. This surpasses some established convolutional backbones such as ResNet-101 (68.93%) and DenseNet-201 (71.60%), as well as some transformer-based models like ViT-B/16 (69.93%), indicating iiANET’s enhanced ability for modeling spatial complexity in aerial imagery.

(3) Oxford-IIIT Pets: In the fine-grained Oxford-IIIT Pets classification task, which requires distinguishing subtle visual differences between classes, iiANET-L achieves the highest top-1 accuracy of 76.23%, outperforming some strong baselines including CrossViT-B (73.1%) and DeiT-B (71.0%). The smaller variant, iiANET-B, also performs well with 74.04%, also has a higher score than ViT-B/16 (51.23%) and ResNet-101 (59.25%). These results demonstrate iiANET’s strong generalization and fine-grained discriminative power.

This analysis shows iiABlock as a strong visual backbone, enabling both iiANET-B and iiANET-L to capture long-range dependencies while maintaining a favorable trade-off among accuracy, parameter count, and computational cost. These characteristics make iiANET a promising backbone for deployment in resource-constrained environments across domains such as medical imaging, autonomous driving, and remote sensing.

Object detection on COCO val2017 in Table 2, we experiment with iiANET as a backbone, the YOLOv8 on the COCO dataset, our results show that both iiANET-B and iiANET-L demonstrate better performance in comparison to other methods in terms of mAP across all evaluation metrics. Specifically, iiANET-B and iiANET-L achieve mAP of (62.6%, 63.1%) and (63.0%, 64.4%) for AP val2017 and AP test2017, indicating their proficiency in detecting objects with high precision and recall.

Instance segmentation on COCO val2017 in Table 3, we also experiment with the effectiveness of iiANET in capturing long-range dependencies on instance segmentation tasks. Table 3 demonstrates better performance in comparison to other models in instance segmentation tasks. Both iiANET-B and iiANET-L achieve better AP box of (45.3% and 45.8%) and AP mask of (39.5% and 42.1%), indicating their capability to accurately segment instances in images with varying complexities and occlusions.

4.3 Ablation

We performed series of ablation studies to investigate iiANET-B capacity from different aspects, we use image classification tasks.

Analysis. Table 4 highlights the effectiveness of the proposed combination of modules within the **iiABlock**. Variant (c), which integrates MBConv, Dilated Convolution, and MHSA, achieves the best performance with a top-1 accuracy of **80.57%**. This demonstrates that the synergy between local depth-wise convolutions, dilated receptive fields, and global attention significantly improves feature representation. Variants (a), (b), and (e) show progressive improvements, while (d), which lacks MBConv, suffers a performance drop, underscoring the importance of lightweight local feature extraction.

Table 1: Classification result comparison on ImageNet-1K, AID, and Oxford-III datasets.

Backbone	Size	Train	#Params	FLOPs	Throughput	Top-1 Acc.	Top-5 Acc.
IMAGENET-1K Deng et al. (2009)							
ResNet-101	224 ²	90	44.5M	14.58G	-	78.0%	94.0%
EffNet-B5	224 ²	90	30.4M	4.49G	-	83.6%	96.7%
Dense201	224 ²	90	20.0M	7.35G	-	77.42%	93.6%
ViT-L/16	224 ²	150	304.3M	59.69G	-	76.53%	93.2%
MobileViT-S	256 ²	90	6M	2G	-	77.0%	94.6%
Dilate-B	224 ²	120	48M	9.96G	-	84.9%	-
BoT50	256 ²	90	25.6M	3.18G	-	84.4%	-
CoAtNet-3	224 ²	90	168M	32.53G	-	84.5%	-
Vim-S	224 ²	90	26M	5.3G	811	80.3%	-
S4ND-ViT-B	224 ²	90	89M	17.1G	397	80.4%	-
VMamba-T	224 ²	90	30M	4.9G	1686	82.6%	-
Swin-T	224 ²	300	29M	4.5G	755.2	81.3%	-
DeiT-B	224 ²	120	86M	17.5G	292.3	81.8%	-
Cross-ViT-B	224 ²	120	105M	20.1G	1321	82.2%	-
iiANET-B	299 ²	90	25.2M	8.22G	-	79.34%	94.71%
iiANET-L	299 ²	120	52.3M	17.07G	-	84.9%	96.83%
AID Xia et al. (2017)							
ResNet-101	224 ²	90	44.5M	14.58G	-	68.93%	93.37%
EffNet-B5	224 ²	90	30.4M	4.49G	81.85	65.73%	92.07%
Dense201	224 ²	90	20.0M	7.35G	76.38	71.60%	94.37%
ViT-B/16	224 ²	150	86.6M	16.86G	48.01	69.93%	93.27%
MobileViT-S	256 ²	90	6M	2G	1986	66.76%	91.23%
DiNAT-B	224 ²	90	90M	13.7G	764	79.12%	93.27%
BoT50	256 ²	90	25.6M	3.18G	95.20	72.50%	94.27%
CoAtNet-3	224 ²	90	168M	32.53G	27.92	80.17%	94.93%
Vim-S	224 ²	90	26M	5.3G	811	74.2%	-
VMamba-T	224 ²	90	30M	4.9G	1686	78.9%	-
Swin-T	224 ²	300	29M	4.5G	755.2	78.0%	-
DeiT-B	224 ²	120	86M	17.5G	292.3	81.2%	-
Cross-ViT-B	224 ²	120	105M	20.1G	1321	80.6%	-
iiANET-B	299 ²	90	25.2M	8.22G	-	80.57%	95.67%
iiANET-L	299 ²	90	52.3M	17.07G	-	83.11%	96.07%
OXFORD-III Parkhi et al. (2012)							
ResNet-50	224 ²	90	25.6M	7.71G	119.75	59.07%	86.61%
ResNet-101	224 ²	90	44.5M	14.58G	105.80	59.25%	87.38%
EffNet-B5	224 ²	90	30.4M	4.49G	96.84	47.54%	78.92%
Dense201	224 ²	90	20.0M	7.35G	98.82	62.55%	86.84%
ViT-B/16	224 ²	150	86.6M	16.86G	-	51.23%	82.31%
ViT-L/16	224 ²	150	304.3M	59.69G	-	53.19%	83.28%
MobileViT-S	256 ²	90	6M	2G	1986	51.89%	84.52%
DiNAT-B	224 ²	90	90M	13.7G	764	69.37%	92.09%
Dilate-B	224 ²	120	48M	9.96G	-	64.85%	88.18%
BoT50	256 ²	90	25.6M	3.18G	118.97	64.13%	90.00%
CoAtNet-3	224 ²	90	168M	32.53G	43.80	67.57%	91.36%
VMamba-T	224 ²	90	30M	4.9G	1686	67.8%	-
Swin-T	224 ²	300	29M	4.5G	755.2	72.3%	-
DeiT-B	224 ²	120	86M	17.5G	292.3	71.0%	-
Cross-ViT-B	224 ²	120	105M	20.1G	1321	73.1%	-
iiANET-B	299 ²	90	25.2M	8.22G	51.44	74.04%	93.98%
iiANET-L	299 ²	90	52.3M	17.07G	-	76.23%	94.54%

Table 2: Object detection results on the COCO dataset.

Backbone	Object Detector	AP ^b val2017	AP ^b test2017
ResNet-50 Carion et al. (2020)	Faster R-CNN	36.7	37.9
FD-SwinV2-G Wei et al. (2022)	HTC++	-	64.2
Florence-CoSwin-H Yuan et al. (2021)	DyHead	62.0	62.4
Swin-L Liu et al. (2021b)	DINO	63.2	63.3
BEiT-3 Wang et al. (2022)	ViTDet	-	63.7
Swin V2-G Liu et al. (2022)	HTC++	62.5	63.1
iiANET-B	YOLOv8	62.6	63.0
iiANET-L	YOLOv8	63.1	64.4

Table 3: Instance Segmentation on COCO dataset with Mask R-CNN (1x schedule).

Backbone	AP ^{box}	AP ^{box} ₅₀	AP ^{box} ₇₅	AP ^{mask}	AP ^{mask} ₅₀	AP ^{mask} ₇₅
ResNet-50 He et al. (2016)	38.0	58.6	41.4	34.4	55.1	36.7
PVT-M Hassani et al. (2023)	42.0	64.4	45.6	39.0	61.6	42.1
TRT-ViT-C Wang et al. (2021)	44.7	66.9	48.8	40.8	63.9	44.0
Focal-T Xia et al. (2022)	44.8	67.7	49.2	41.0	64.7	44.2
UniFormer-S/h14 Yang et al. (2021)	45.6	68.1	49.7	41.6	64.8	45.0
Swin-T Li et al. (2023)	42.2	64.6	46.2	39.1	61.6	42.0
Dilate-S Jiao et al. (2023)	45.8	68.2	50.1	41.7	65.3	44.7
BoT50 Srinivas et al. (2021)	43.7	-	-	37.9	-	-
iiANET-B	45.3	65.1	49.8	39.5	58.9	58.0
iiANET-L	45.8	68.3	51.7	42.1	65.1	59.5

Table 5: Ablation study on the effect of changing the MHSA head size and iiABlock ratio on AID dataset.

Settings	Model (Components)	Size	Top-1 Accuracy
iiANET	Ratio: 1.6.1 \rightarrow 2.4.2	299	74.21%
iiANET	Head Size: 8 \rightarrow 16	299	76.85%

Table 6: Ablation study on the effect of adding Register tokens to the 2D-MHSA mechanism on AID dataset.

Number of Registers	Top-1 Accuracy
0	80.57%
1	80.61%
2	80.62%
4	80.77%

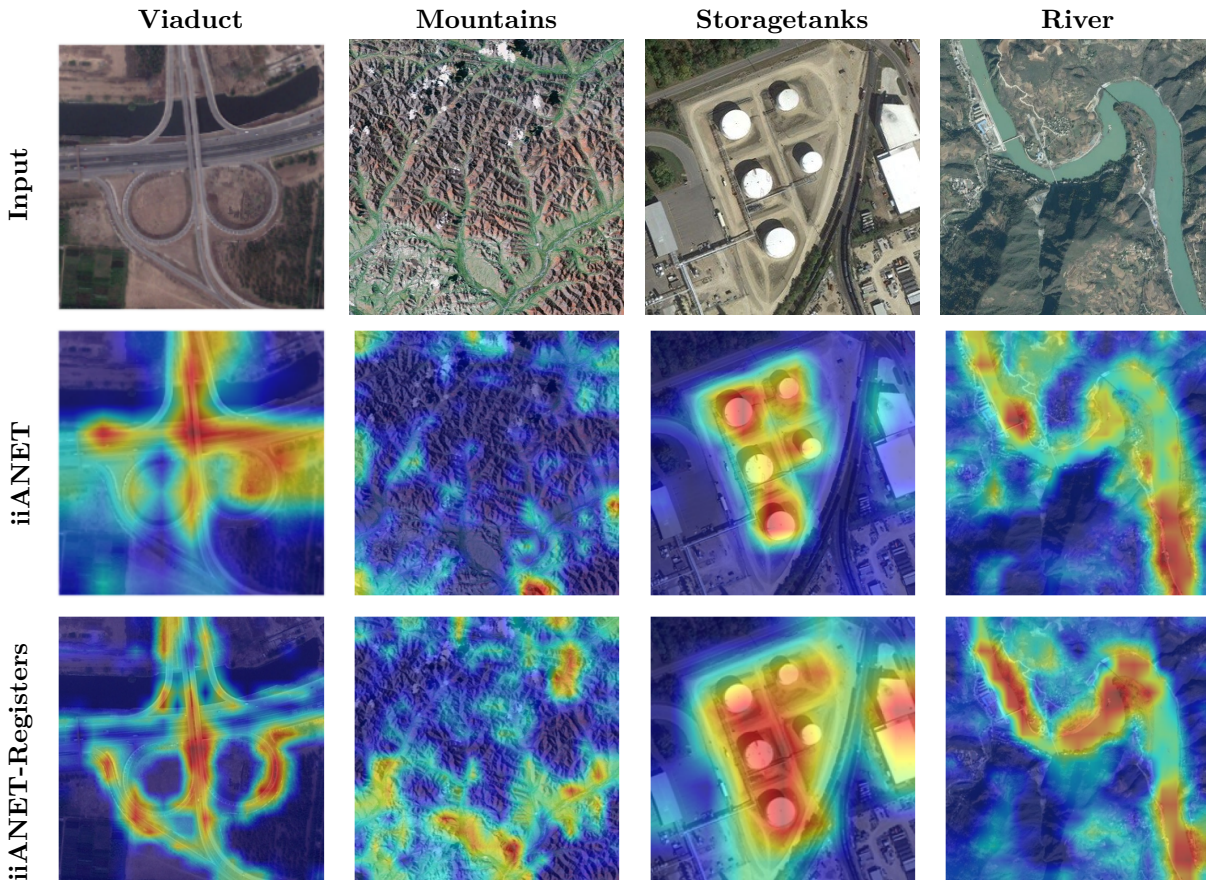
Analysis. Table 6 demonstrates that introducing register tokens into the 2D-MHSA improves performance on the AID dataset. With four registers, the model achieves the highest top-1 accuracy of **80.77%**, compared to **80.57%** without any. The consistent improvement suggests that registers help compensate for MHSA’s order-agnostic nature by enhancing the contextual richness of learned representations.

Analysis. Table 5 shows that increasing the MHSA head size from 8 to 16 significantly improves the model’s performance, boosting top-1 accuracy to **76.85%**. This gain comes with a minor GPU processing increase of approximately 3%–5%. Additionally, adjusting the internal ratio of the iiABlock from 1.6.1 to 2.4.2 also enhances accuracy, though at the cost of increased computational complexity. These results emphasize the importance of careful architectural tuning for achieving optimal performance-efficiency tradeoffs.

Table 4: Ablation studies on various variants of iiABlock using the AID dataset.

Settings	Model (Components)	Size	Top-1 Accuracy
(a)	MBCConv2	299	69.23%
(b)	MBCConv + Dilated Conv	299	74.85%
(c)	MBCConv + Dilated Conv + MHSA	299	80.57%
(d)	Dilated + MHSA	299	67.01%
(e)	MBCConv + MHSA	299	78.72%

Figure 5: Visual effect of registers on iiANET for the AID dataset, showing how adding registers to the r-MHSA module enhances interpretability.



5 Conclusion

This work proposed iiANET, a novel hybrid model designed to efficiently improve long-range dependencies in complex images by integrating CNN layers and the MHSA mechanism with registers in parallel. Comprehensive qualitative and quantitative results show improvements in capturing long-range dependencies compared to some previous SOTA models. Additionally, we validate the performance of our model across diverse datasets and highlight its potential as a backbone in object detection and segmentation models.

Limitations. iiANET performs better on images with long-range dependencies but is less effective for datasets like ImageNet-1K with localized objects. Its multi-branch design may also cause scaling issues, increasing memory and computation on very high-resolution inputs or limited hardware.

References

- Imon Banerjee, Yuan Ling, Matthew C. Chen, Sadid A. Hasan, Curtis P. Langlotz, Nathaniel Moradzadeh, Brian Chapman, and et al. Comparative effectiveness of convolutional neural network (cnn) and recurrent neural network (rnn) architectures for radiology text report classification. *Artificial Intelligence in Medicine*, 97:79–88, 2019.
- Irwan Bello, William Fedus, Xianzhi Du, Ekin Dogus Cubuk, Aravind Srinivas, Tsung-Yi Lin, Jonathon Shlens, and Barret Zoph. Revisiting res-nets: Improved training and scaling strategies. In *Advances in Neural Information Processing Systems*, volume 34, pp. 22614–22627, 2021.
- Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European Conference on Computer Vision*, pp. 213–229, Cham, 2020. Springer International Publishing.
- Chun-Fu Richard Chen, Quanfu Fan, and Rameswar Panda. Crossvit: Cross-attention multi-scale vision transformer for image classification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 357–366, 2021.
- Zihang Dai, Hanxiao Liu, Quoc V. Le, and Mingxing Tan. Coatnet: Marrying convolution and attention for all data sizes. *Advances in Neural Information Processing Systems*, 34:3965–3977, 2021a.
- Zihang Dai, Hanxiao Liu, Quoc V. Le, and Mingxing Tan. Coatnet: Marrying convolution and attention for all data sizes. In *Advances in neural information processing systems*, volume 34, pp. 3965–3977, 2021b.
- Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision transformers need registers. *arXiv preprint arXiv:2309.16588*, 2023.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. IEEE, 2009.
- Jeffrey Donahue, Lisa Anne Hendricks, Sergio Guadarrama, Marcus Rohrbach, Subhashini Venugopalan, Kate Saenko, and Trevor Darrell. Long-term recurrent convolutional networks for visual recognition and description. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2625–2634, 2015.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, and et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- Jianyuan Guo, Kai Han, Han Wu, Yehui Tang, Xinghao Chen, Yunhe Wang, and Chang Xu. Cmt: Convolutional neural networks meet vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12175–12185, 2022.
- Kai Han, Yunhe Wang, Hanting Chen, Xinghao Chen, Jianyuan Guo, Zhenhua Liu, Yehui Tang, and et al. A survey on vision transformer. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(1): 87–110, 2022.
- Y. Haruna, S. Qin, A. H. A. Chukkol, A. A. Yusuf, I. Bello, and A. Lawan. Exploring the synergies of hybrid convolutional neural network and vision transformer architectures for computer vision: A survey. *Engineering Applications of Artificial Intelligence*, 144:110057, 2025.
- Ali Hassani and Humphrey Shi. Dilated neighborhood attention transformer. *arXiv preprint arXiv:2209.15001*, 2022.
- Ali Hassani, Steven Walton, Jiachen Li, Shen Li, and Humphrey Shi. Neighborhood attention transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6185–6194, 2023.

- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778, 2016.
- Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 7132–7141, 2018.
- Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q. Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4700–4708, 2017.
- Jiayu Jiao, Yu-Ming Tang, Kun-Yu Lin, Yipeng Gao, Andy J. Ma, Yaowei Wang, and Wei-Shi Zheng. Dilateformer: Multi-scale dilated transformer for visual recognition. *IEEE Transactions on Multimedia*, 25:8906–8919, 2023.
- Asifullah Khan, Zunaira Rauf, Anabia Sohail, Abdul Rehman Khan, Hifsa Asif, Aqsa Asif, and Umair Farooq. A survey of the vision transformers and their cnn-transformer based variants. *Artificial Intelligence Review*, 56(Suppl 3):2917–2970, 2023.
- Hanchao Li, Pengfei Xiong, Haoqiang Fan, and Jian Sun. Dfanet: Deep feature aggregation for real-time semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9522–9531, 2019.
- Kunchang Li, Yali Wang, Junhao Zhang, Peng Gao, Guanglu Song, Yu Liu, Hongsheng Li, and Yu Qiao. Uniformer: Unifying convolution and self-attention for visual recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(10):12581–12600, 2023.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V*, volume 13, pp. 740–755. Springer International Publishing, 2014.
- Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 10012–10022, 2021a.
- Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 10012–10022, 2021b.
- Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, et al. Swin transformer v2: Scaling up capacity and resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12009–12019, 2022.
- Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3431–3440, 2015.
- Omkar M. Parkhi, Andrea Vedaldi, Andrew Zisserman, and C. V. Jawahar. Cats and dogs. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3696–3703. IEEE, 2012.
- Zhiliang Peng, Wei Huang, Shanzhi Gu, Lingxi Xie, Yaowei Wang, Jianbin Jiao, and Qixiang Ye. Conformer: Local features coupling global representations for visual recognition. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 367–376, 2021.
- Joseph Redmon and Ali Farhadi. Yolo9000: Better, faster, stronger. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7263–7271, 2017.
- Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in Neural Information Processing Systems*, 28, 2015.

- Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pp. 618–626, 2017.
- Prabira Kumar Sethy, Nalini Kanta Barpanda, Amiya Kumar Rath, and Santi Kumari Behera. Deep feature-based rice leaf disease identification using support vector machine. *Computers and Electronics in Agriculture*, 175:105527, 2020.
- Aravind Srinivas, Tsung-Yi Lin, Niki Parmar, Jonathon Shlens, Pieter Abbeel, and Ashish Vaswani. Bottleneck transformers for visual recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16519–16529, 2021.
- Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International Conference on Machine Learning*, pp. 6105–6114. PMLR, 2019.
- Qilong Wang, Banggu Wu, Pengfei Zhu, Peihua Li, Wangmeng Zuo, and Qinghua Hu. Eca-net: Efficient channel attention for deep convolutional neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11534–11542, 2020.
- Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 568–578, 2021.
- Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, Zhiliang Peng, Qiang Liu, Kriti Aggarwal, et al. Image as a foreign language: Beit pretraining for all vision and vision-language tasks. *arXiv preprint arXiv:2208.10442*, 2022.
- Yixuan Wei, Han Hu, Zhenda Xie, Zheng Zhang, Yue Cao, Jianmin Bao, Dong Chen, and Baining Guo. Contrastive learning rivals masked image modeling in fine-tuning via feature distillation. *arXiv preprint arXiv:2205.14141*, 2022.
- Haiping Wu, Bin Xiao, Noel Codella, Mengchen Liu, Xiyang Dai, Lu Yuan, and Lei Zhang. Cvt: Introducing convolutions to vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 22–31, 2021.
- Gui-Song Xia, Jingwen Hu, Fan Hu, Baoguang Shi, Xiang Bai, Yanfei Zhong, Liangpei Zhang, and Xiaoqiang Lu. Aid: A benchmark data set for performance evaluation of aerial scene classification. *IEEE Transactions on Geoscience and Remote Sensing*, 55(7):3965–3981, 2017.
- Xin Xia, Jiashi Li, Jie Wu, Xing Wang, Xuefeng Xiao, Min Zheng, and Rui Wang. Trt-vit: Tensorrt-oriented vision transformer. *arXiv preprint arXiv:2205.09579*, 2022.
- Jianwei Yang, Chunyuan Li, Pengchuan Zhang, Xiyang Dai, Bin Xiao, Lu Yuan, and Jianfeng Gao. Focal self-attention for local-global interactions in vision transformers. *arXiv preprint arXiv:2107.00641*, 2021.
- Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*, 2015.
- Fisher Yu, Vladlen Koltun, and Thomas Funkhouser. Dilated residual networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 472–480, 2017.
- Lu Yuan, Dongdong Chen, Yi-Ling Chen, Noel Codella, Xiyang Dai, Jianfeng Gao, Houdong Hu, et al. Florence: A new foundation model for computer vision. *arXiv preprint arXiv:2111.11432*, 2021.
- Xindong Zhang, Hui Zeng, Shi Guo, and Lei Zhang. Efficient long-range attention network for image super-resolution. In *European Conference on Computer Vision*, pp. 649–667. Springer Nature Switzerland, 2022.