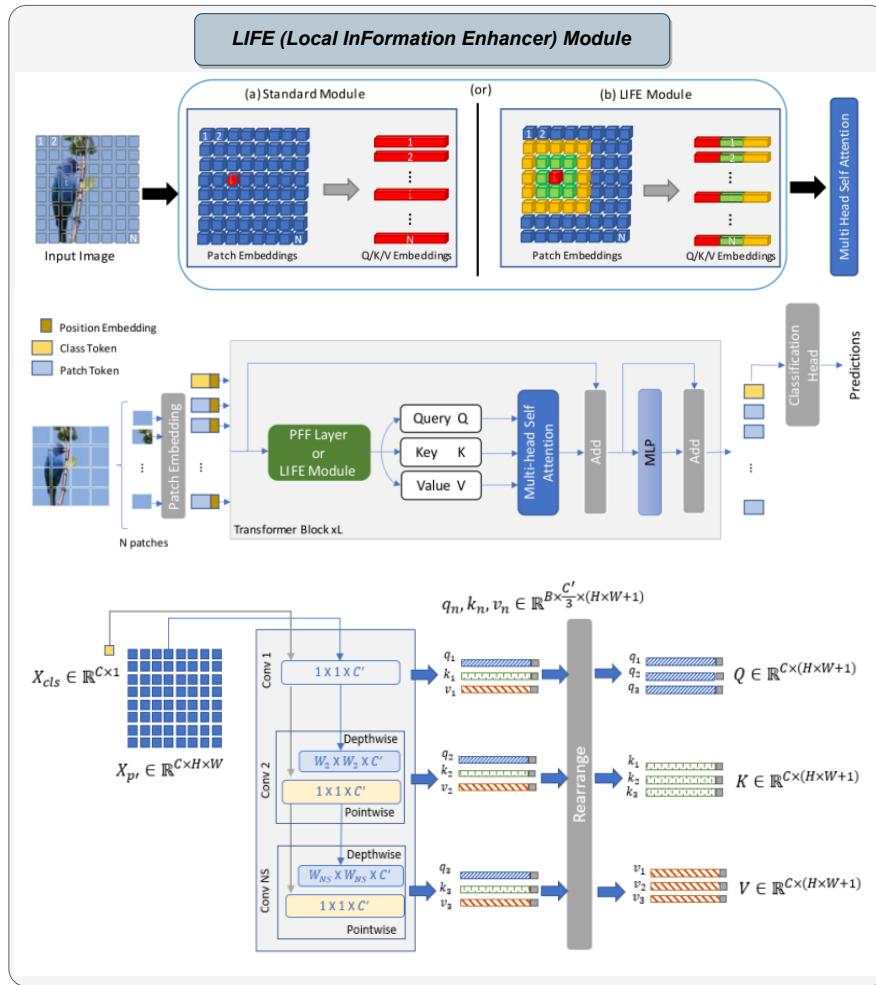


Graphical Abstract

Enhancing Performance of Vision Transformers on Small Datasets through Local Inductive Bias Incorporation

Ibrahim Batuhan Akkaya^{*}, Senthilkumar S. Kathiresan, Elahe Arani[†], Bahram Zonooz[†]



^{*}Corresponding author: bthakkaya@gmail.com. [†]Equal contribution.

Highlights

Enhancing Performance of Vision Transformers on Small Datasets through Local Inductive Bias Incorporation

Ibrahim Batuhan Akkaya^{*}, Senthilkumar S. Kathiresan, Elahe Arani[†], Bahram Zonooz[†]

- Introducing Local InFormation Enhancer (LIFE) module to complement global information with local context in the ViT architecture
- Versatile and Efficient: LIFE module can be easily integrated into different ViT architectures with minimal computational and memory costs, including auxiliary tokens.
- Boost Small-Dataset Results: LIFE module boosts the performance of ViTs on small image classification datasets and dense prediction tasks like object detection and semantic segmentation.
- Introducing a novel visualization method, Dense Attention Roll-Out, to visualize attention for dense prediction tasks.

^{*}Corresponding author: bthakkaya@gmail.com. [†]Equal contribution.

Enhancing Performance of Vision Transformers on Small Datasets through Local Inductive Bias Incorporation

Ibrahim Batuhan Akkaya^{*,a}, Senthilkumar S. Kathiresan^a, Elahe Arani^{†,a,b},
Bahram Zonooz^{†,a,b}

^a*Advanced Research Lab, NavInfo Europe, Eindhoven, 5657 DB, Netherlands*

^b*Department of Mathematics and Computer Science, Eindhoven University of
Technology, Eindhoven, 5612 AZ, Netherlands*

Abstract

Vision transformers (ViTs) achieve remarkable performance on large datasets, but tend to perform worse than convolutional neural networks (CNNs) when trained from scratch on smaller datasets, possibly due to a lack of local inductive bias in the architecture. Recent studies have therefore added locality to the architecture and demonstrated that it can help ViTs achieve performance comparable to CNNs in the small-size dataset regime. Existing methods, however, are architecture-specific or have higher computational and memory costs. Thus, we propose a module called *Local InFormation Enhancer (LIFE)* that extracts patch-level local information and incorporates it into the embeddings used in the self-attention block of ViTs. Our proposed module is memory and computation efficient, as well as flexible enough to process auxiliary tokens such as the classification and distillation tokens. Empirical results show that the addition of the LIFE module improves the performance of ViTs on small image classification datasets. We further demonstrate how the effect can be extended to downstream tasks, such as object detection and semantic segmentation. In addition, we introduce a new visualization method, Dense Attention Roll-Out, specifically designed for dense prediction tasks, allowing the generation of class-specific attention maps utilizing the attention maps of all tokens.*

Keywords: Vision Transformer, Inductive Bias, Locallity, Small Dataset

*Corresponding author: bthakkaya@gmail.com. †Equal contribution.

*The code will be publicly available upon acceptance.

1. Introduction

Transformers, a new kind of encoder-decoder model that uses a self-attention mechanism to process input data [1], were initially proposed for sequence modeling in the natural language processing (NLP) domain. The success of transformers in NLP has led to the development of these architectures for a wide range of vision tasks [2, 3, 4]. Vision transformers (ViTs), when trained or pre-trained on a large dataset, outperform their CNN counterparts. However, for many real-world vision tasks, a large amount of annotated data is either too expensive or not feasible. As a result, the data-hungry nature of ViTs prevents them from being applied to a number of crucial real-world problems for which a limited amount of annotated data are available.

In the majority of ViTs, an image is divided into a sequence of non-overlapping patches, from which the self-attention layer learns the global context. Information in patches that are spatially adjacent to a given patch can be used to create its local context. ViTs do not, however, exploit this information due to a low inductive bias that is only coming from strided convolution in the patch embedding layer [2]. Convolutional layers, on the other hand, enable CNN architectures to utilize local information at the pixel level, thereby enhancing their data efficiency [5, 6]. The local context may therefore be essential in enabling ViTs to learn vision tasks with fewer data samples.

Recent studies have improved the use of the local context in ViTs through architectural modifications [3, 7], token pre-processing [8], or the addition of convolutional layers [9]. Despite the fact that these findings support the importance of utilizing local information, their design is not adaptable enough to be annexed to other ViT architectures and/or has a negative impact on other performance factors such as memory consumption and computational cost. Therefore, it is advantageous to devise a method to incorporate local information effectively and modularly into the architecture of ViTs.

The functionality of a transformer architecture depends on its self-attention layers, which require a sequence of embeddings. These embeddings are typically generated using a point-wise feedforward layer with a receptive field consisting of only one patch from the input image or one token from the previous layer (Figure 1(a)). However, the use of larger receptive fields can result in embeddings that contain local information from spatially adjacent patches, as shown in Figure 1(b). We hypothesize that by enriching the embeddings with this local context, vision transformer models (ViTs) will be

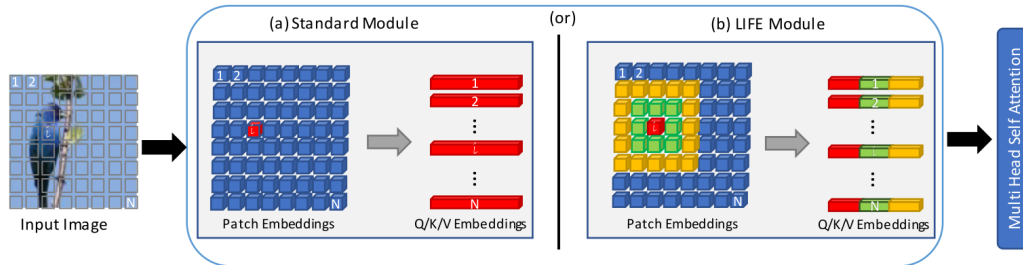


Figure 1: In vision transformers, the input image is divided into N non-overlapping patches, which are then transformed into embeddings. (a) These patch embeddings are then passed through a point-wise feedforward layer to generate query (Q), key (K), and value (V) embeddings. The Q, K, and V embeddings are then inputted to a self-attention layer. (b) Alternatively, the LIFE module integrates local information from adjacent patches into the Q, K, and V embeddings using sequentially larger receptive fields. For example, the standard module uses only the i^{th} patch to generate the i^{th} Q/K/V embedding, while the LIFE module also includes information from adjacent patches (shown in green and orange).

able to learn the global context with fewer data points and perform better on smaller datasets.

We propose the LIFE (*Local InFormation Enhancer*) module to improve the performance of vision transformers on smaller datasets by creating embeddings for self-attention layers with larger receptive fields. The LIFE module reshapes the input tokens from the previous layer into an image format and applies convolution layers with multiple kernel sizes. After the feature maps from the convolutional layers are transformed back into tokens, they are sent to the self-attention layers. To add local context to all patch tokens and any other auxiliary tokens in the ViT architecture, we use depthwise separable convolutional (DSC) layers [10] in the LIFE module. Unlike a standard convolution layer, a DSC layer performs the computations in two steps: a depth-wise convolution layer followed by a point-wise convolution layer, which is more computationally efficient. Note that auxiliary tokens, such as classification and distillation tokens, in ViTs, are processed by the pointwise convolution layer in the DSC.

We evaluate the efficacy and versatility of the LIFE module by integrating it into different ViT architectures with varying capacities. Using DeiT, T2T, and Swin transformers as base architectures, we evaluate classification performance on the ImageNet-100, CIFAR10, CIFAR100, and Tiny-ImageNet

datasets. We also evaluate the addition of LIFE to ViTs on dense prediction tasks using the VOC dataset for object detection and the Cityscapes dataset for semantic segmentation tasks. Our extensive empirical experiments demonstrate that the LIFE module can be easily integrated into various ViT architectures and consistently improves performance, regardless of the task at hand, with negligible memory and computation overhead. In addition, we qualitatively assess the contribution of the LIFE module to each task. To visualize the attention for dense prediction tasks, we propose a dense attention roll-out. Our results further support the notion that LIFE can enhance local context learning by guiding the network to attend to more specific regions. The contributions of our work can thus be summarized as follows;

- Introducing *Local InFormation Enhancer (LIFE)* module, which complements global information by adding local context to the embeddings used in ViT.
- Demonstration of the ability of the LIFE module to be easily integrated into different ViT architectures with minimal memory and computation costs overhead, even in the architecture that contains auxiliary tokens.
- Employing the LIFE module in different ViTs results in performance gains on smaller datasets such as ImageNet-100, Tiny-ImageNet, CIFAR10, and CIFAR100.
- LIFE module is versatile and also results in a boost in performance for dense prediction tasks.
- Proposing a novel method, *Dense Attention Roll-Out*, to visualize attention for dense prediction tasks.
- Qualitative evaluation of the contribution of the LIFE module to each task using our proposed visualization method.

2. Related Work

Vision Transformers (ViT). have demonstrated competitiveness with convolutional neural networks (CNNs) in various vision tasks, such as image classification [2, 11, 3], object detection [12], and semantic segmentation [13, 14]. These models utilize self-attention at the early levels to construct a convolution-free neural network. Following the introduction of the original ViT model [2], numerous studies have focused on improving classification

performance through architectural modifications, knowledge distillation, or advanced data augmentation techniques [11, 15, 8, 16, 17, 3, 18].

Locality. Convolutional neural networks (CNNs) have become a common architecture for visual tasks. They typically contain a stack of convolution layers with small kernel sizes that utilize information from neighboring feature vectors. This architectural design exploits the spatial correlation in the natural images, making CNNs more data efficient. Evidence from CNNs suggests that it is essential to use local information to improve performance on small-scale datasets [5, 6]. In contrast, the self-attention mechanism in the transformer block establishes a global relation between tokens, but ignores the locality. To address this, many approaches have been proposed to introduce locality bias into transformer architectures through architectural modifications.

Several recent studies have proposed hybrid networks to incorporate locality bias from convolution operations into transformer architectures. CeiT [19] uses low-level features from CNNs rather than patches extracted from raw images, and introduces a depth-wise separable convolution into the feed-forward network within the transformer block to improve locality. CvT [20] integrates convolution into token embeddings, allowing a progressively decreasing number of tokens while increasing the dimension of the features. It also uses a depth-wise separable convolution to compute query, key, and value. Crossformer [21] addresses the lack of multiscale information in transformer architectures by processing tokens using short- and long-distance attention and combining multiscale information from neighboring tokens using multiple convolution layers with different kernel sizes within a pyramid-like architecture.

Another line of work proposes to incorporate a pyramid structure derived from CNNs into the transformer architecture in order to induce locality in the network. This is achieved through various techniques, such as the use of a gradual shrinking technique in PVT [16, 17], self-attention within windows in Swin [3, 18], and hierarchical aggregation of transformer blocks in NesT [22]. All of these approaches involve the gradual combination of neighboring tokens, enabling the network to consider the local context in its processing.

Other approaches aim to introduce locality while maintaining the pure transformer architecture. T2T [8] replaces the standard patch embedding layer with progressive tokenization to combine neighboring tokens; while TNT [15] divides the input image into large patches called visual sentences

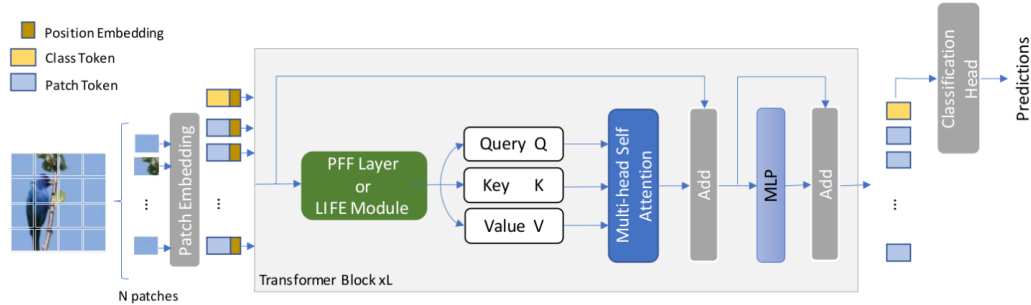


Figure 2: Architecture of a standard ViT. In order to incorporate local context in the embeddings for multi-head self-attention layers, we replace the pointwise feedforward (PPF) layer with our proposed LIFE module (details in Figure 3), where it generates a multiscale query, key, and value.

and small patches called visual words. TNT uses a shared sub-transformer architecture to extract the relation and similarity between the small patches, and passes the information generated by the sub-transformer to the visual sentence for further processing by the standard transformer block.

However, the aforementioned approaches are architecture-specific and focus on designing effective models for large-scale datasets, but may not necessarily perform well on smaller datasets. Liu et al. [23] propose a method that introduces a self-supervised task to predict the geometric distance between pairs of output tokens, while Li et al. [24] propose a distillation-based training mechanism that uses a lightweight pre-training CNN model to distill its features into a ViT. Although these approaches are flexible and can be integrated into different architectures, they may also come with additional memory and computational costs, and may not be as effective when applied to smaller datasets.

In contrast, the LIFE module is a modular and efficient approach that uses locality to generate richer embeddings, which is more effective for learning with less data. It introduces locality in query, key, and value, making it flexible for integration with any transformer architecture. Unlike CvT, the LIFE module benefits from multiscale information and is more efficient compared to T2T and Crossformer, as it uses dense concatenation and depthwise separable convolution. Additionally, unlike previous work, the LIFE module utilizes multi-scale locality in every block and can be used in conjunction with training mechanisms for small-scale datasets.

3. Methodology

In ViTs, an input image is divided into N non-overlapping square patches [11]. These patches are then flattened and embedded in patch tokens X_{patch} of length C using a linear layer. Positional embeddings are added to the token patches and an additional classification token X_{cls} , which is a learnable embedding of the same length, forming a matrix X :

$$X = [X_{patch}|X_{cls}]; \quad X \in \mathbb{R}^{C \times (N+1)}, X_{patch} \in \mathbb{R}^{C \times N}, X_{cls} \in \mathbb{R}^{C \times 1} \quad (1)$$

The resulting matrix is then passed to a series of L transformer blocks (Figure 2), each consisting of a point-wise feedforward layer (PFF), followed by a multi-head self-attention layer (MHSA) and a multi-layer perceptron (MLP):

$$\begin{aligned} Q, K, V &\leftarrow PFF(X), \\ X_{out} &= MLP(MHSA(Q, K, V)), \end{aligned} \quad (2)$$

The final prediction is usually obtained by processing the classification token at the end of the last transformer block or by using the global average pooling of the patch tokens. However, standard ViTs lack local inductive bias; therefore, to introduce local context, we replace the PFF in Eq. 2 with our proposed LIFE module.

3.1. LIFE Module

To incorporate local information, the LIFE module uses multiple hierarchically arranged convolutional layers, as depicted in Figure 3. The first layer in the hierarchy is always a point-wise convolution with the smallest receptive field, while subsequent layers have progressively larger receptive fields due to increasing kernel sizes. The output features of these layers represent local information gleaned from various receptive fields, and are rearranged to obtain the final query Q , key K , and value V . The kernel sizes and paddings are configured to maintain constant spatial resolution throughout all layers, and for similar reasons, the channel size of all convolutional layer output feature maps is fixed at a constant value C' .

Except for the first layer, the LIFE module uses depth-wise separable convolutions [10] for all other layers, which consist of a depth-wise convolution followed by a point-wise convolution. This type of convolution is efficient in terms of memory and computation, allowing the LIFE module to have

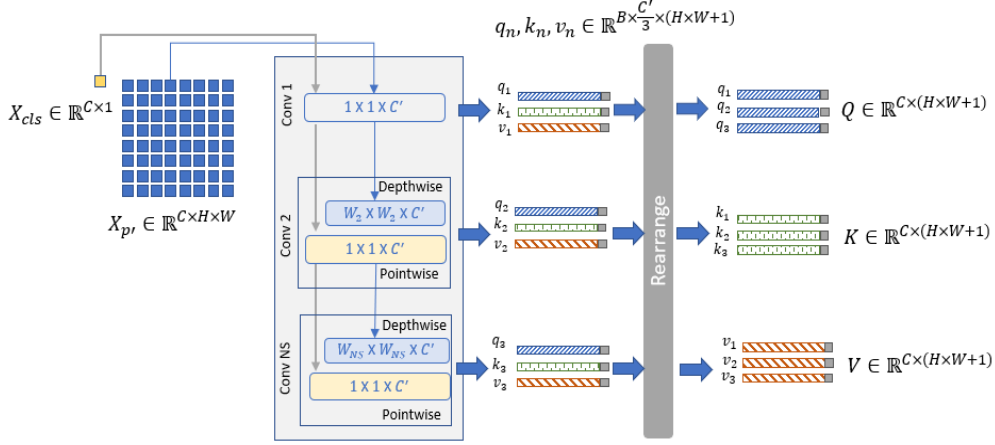


Figure 3: Overview of *LIFE* module. It consists of a hierarchy of convolutional layers, with the first layer being a point-wise convolution and the remaining layers being depth-wise separable convolutions. Both the classification token X_{cls} and patch tokens $X_{p'}$ (arranged in the input image format) are passed through these layers (marked in \downarrow and \downarrow , respectively). Each layer in the hierarchy outputs local information from a relatively larger receptive field. The information is then divided into query q_i , key k_i , and value v_i (including the processed class token represented as a gray box). Then, these hierarchical characteristics are rearranged to form the final query Q , key K , and value V .

minimal overhead compared to the original PFF layer. In addition, point-wise convolution can be used to process any number of auxiliary tokens, such as the classification token and the distillation token. The LIFE module is described in more detail in the following.

Processing Patch Tokens in the LIFE Module. The patch tokens $X_{patch} \in \mathbb{R}^{C \times N}$ are rearranged into an image format of shape $X_{p'} \in \mathbb{R}^{C \times H \times W}$, where $N = H \times W$. The $X_{p'}$ is then passed through a series of convolutional layers with progressively increasing receptive fields. Here, we use three layers:

$$F_{p_1} = Conv_1(X_{p'}); \quad F_{p_2} = Conv_2(F_{p_1}); \quad F_{p_3} = Conv_3(F_{p_2}); \quad (3)$$

where $F_{p_1}, F_{p_2}, F_{p_3} \in \mathbb{R}^{C \times H \times W}$ represent three different scales of local features obtained from the patch tokens. Each of these features is divided into three in the channel dimension and rearranged to form embeddings $Q_p, K_p, V_p \in$

$\mathbb{R}^{C \times H \times W}$;

$$\begin{aligned} Q_p &= [F_{p_1}^{(C/3 \times H \times W)}, F_{p_2}^{(C/3 \times H \times W)}, F_{p_2}^{(C/3 \times H \times W)}]^T \\ K_p &= [F_{p_1}^{(C/3 \times H \times W)}, F_{p_2}^{(C/3 \times H \times W)}, F_{p_2}^{(C/3 \times H \times W)}]^T \\ V_p &= [F_{p_1}^{(C/3 \times H \times W)}, F_{p_2}^{(C/3 \times H \times W)}, F_{p_2}^{(C/3 \times H \times W)}]^T \end{aligned} \quad (4)$$

Processing Auxiliary Tokens in the Proposed Module. The class token $X_{cls} \in \mathbb{R}^{C \times 1}$ and any other auxiliary tokens can be processed using point-wise convolutions in the LIFE module:

$$\begin{aligned} F_{c_1} &= Conv_1(X_{cls}); \\ F_{c_2} &= PointConv_2(F_{c_1}); \\ F_{c_3} &= PointConv_3(F_{c_2}); \end{aligned} \quad (5)$$

where $F_{c_1}, F_{c_2}, F_{c_3} \in \mathbb{R}^{C \times 1}$ represent different linear projections of the classification token for different scales. Next, each of the features is divided into three parts in the channel dimension and rearranged to form embeddings $Q_c, K_c, V_c \in \mathbb{R}^{C \times 1}$:

$$\begin{aligned} Q_c &= [F_{c_1}^{(C/3 \times 1)}, F_{c_2}^{(C/3 \times 1)}, F_{c_2}^{(C/3 \times 1)}]^T \\ K_c &= [F_{c_1}^{(C/3 \times 1)}, F_{c_2}^{(C/3 \times 1)}, F_{c_2}^{(C/3 \times 1)}]^T \\ V_c &= [F_{c_1}^{(C/3 \times 1)}, F_{c_2}^{(C/3 \times 1)}, F_{c_2}^{(C/3 \times 1)}]^T \end{aligned} \quad (6)$$

Subsequently, the features of the patch tokens $\in \mathbb{R}^{C \times H \times W}$ are flattened in spatial dimensions. The embeddings Q, K, V obtained from the patch tokens and the auxiliary tokens are concatenated to form the final embeddings $Q, K, V \in \mathbb{R}^{C \times (N+1)}$, which contain local features from different scales. The global information obtained through self-attention is complemented by the local context information encoded in these embeddings, resulting in improved performance. By combining global and local information, the model can better understand the context and relationships between different parts of the input.

4. Experimental Settings

We analyze the LIFE module by integrating it into different state-of-the-art transformer architectures. The main experiments are conducted on

small-scale image datasets, as our study focuses primarily on the performance of the transformer with limited data. We evaluate our model for the image classification task. However, many real-world applications are based on object detection or semantic segmentation, so we also examine how the module affects downstream dense prediction tasks. We demonstrate the effectiveness of our module for small datasets with quantitative and qualitative results in classification, detection, and segmentation tasks.

Datasets. Small-scale datasets used in our experiments include CIFAR-100 [25], CIFAR-10 [25], TinyImageNet [26], and ImageNet-100 [27]. These datasets contain 50k, 60k, 100k, and 130k training samples, respectively. ImageNet-100 and TinyImageNet are subsets of the ImageNet-1k dataset [27], with 200 and 100 classes, respectively. CIFAR-10 and CIFAR-100 have 10 and 100 classes, respectively. Except for the ImageNet-100 datasets, all other datasets have small image resolutions, either 32×32 or 64×64 . The VOC dataset [28] is used for object detection, and the Cityscapes dataset [29] is used for semantic segmentation tasks. The VOC and Cityscapes datasets contain 1,464 and 5000 samples for training, respectively.

Implementation Details. In our experiments, we used the Tiny and Small variants of the DeiT [11], T2T-ViT-12, and Swin transformers as baseline models for the classification task. We obtained the LIFE variants of these models by replacing the linear projection that generates the query, key, and value with the LIFE module. We employ three scales with kernel sizes of 1, 3, and 5 and zero paddings of sizes of 0, 1, and 2, respectively.

We use the original image size in our experiments. For the Swin transformer, we select a window size of 8, and for T2T-ViT-12, we decrease the token dimension from 255 to 252 in order to process the token in a multiscale manner in the LIFE module for all datasets. Baseline models were designed for an input image size of 224×224 . Therefore, we keep the network configurations the same as the baselines for the ImageNet-100 dataset. Only for the Swin architecture, we resize the input to 256×256 . For other datasets with small image sizes, we used a patch size of 1 for Swin and 4 for DeiT architectures. To adapt T2T-ViT-12, we replaced the first unfold operation with a 3×3 kernel with a stride of 2 and the last unfold operation with a 1×1 kernel.

For detection and segmentation tasks, we use a tiny version of DeiT and Swin as the backbone. In order to highlight the effect of the LIFE module,

we employ simple heads for dense prediction tasks. For the object detection task, we exploit DETR [12] and evaluate different combinations of backbone and head with and without the LIFE module. For segmentation tasks, we use simple upsampling after the linear transformation operation as the segmentation head. The input image size for both tasks is 512×512 .

For the classification, detection, and segmentation tasks, we followed the training details in [11], [30], and [13]. We re-train all models from scratch with random initialization in our framework for a fair comparison. We use a batch size of 512 for all datasets. For the dense prediction tasks, we also present results for ImageNet pre-trained initialization, as it is believed that the initialization can significantly impact the final performance.

Unlike transformer models, we observed that the ResNet [5] architecture benefited more from simple augmentations for small datasets. We, therefore, apply no augmentation other than random crop and random horizontal flip. ResNet models are trained with SGD with a momentum of 0.9 for 200 epochs. The initial learning rate is set to 0.1 and adjusted with a multi-stage scheduler, which multiplied the learning rate by 0.2 at epochs 60, 120, and 160.

5. Quantitative Results

In the following, we present the results of the evaluation of the LIFE module on various tasks and datasets. We first examine the efficacy of the LIFE module in addressing the performance gap between the transformer and the CNN counterpart when trained on smaller datasets in the image classification task. We then evaluate the versatility of the LIFE module by employing it for dense prediction tasks, including object detection and semantic segmentation.

5.1. Image Classification

We examine the efficacy of the LIFE module in addressing the performance gap between the transformer and the CNN counterpart when trained on smaller datasets. The LIFE module aims to address this issue by introducing locality bias into the transformer architecture through the use of convolutional layers. We train models on a small image classification dataset and compare their performance with that of CNNs, such as ResNet, of similar size. As shown in Table 1, the efficiency of the LIFE module is demonstrated by integrating it into multiple models with different sizes and architectures. We

Table 1: Performance comparison of ViTs on small-scale image classification datasets with and without the addition of the LIFE module. The results include Top-1 accuracy, number of parameters, and GMAC for DeiT and T2T architectures with an input size of $224 \times 224 \times 3$, and for the Swin architecture with an input size of $256 \times 256 \times 3$.

	#Params	GMAC	CF-10	CF-100	Tiny-IM	IM-100
CNN						
Resnet-18	11.23	2.37	94.14	75.10	60.86	80.74
Resnet-50	23.71	5.35	94.32	74.25	63.45	82.70
Transformers						
DeiT-T	5.54	1.26	85.65	51.45	54.89	63.56
DeiT-T-LIFE	5.62	1.27	89.28	71.74	60.51	68.32
T2T-ViT-12	6.66	1.74	88.41	52.51	57.89	83.58
T2T-ViT-12-LIFE	6.59	1.71	89.96	54.51	60.52	83.96
DeiT-S	21.70	4.61	87.24	70.39	55.08	80.84
DeiT-S-LIFE	21.85	4.64	90.38	73.97	59.78	81.52
Swinv2-T	27.72	5.95	95.03	76.65	64.78	86.86
Swinv2-T-LIFE	27.88	6.01	95.14	77.48	65.84	86.98

use accuracy as a performance metric for comparison. The results show that the use of the LIFE module in the transformer architecture improves performance on small-scale datasets across different architectures and model sizes. The improvement is more significant for smaller model sizes; for instance, the LIFE module improves the performance of the DeiT-Tiny architecture by approximately 15% and the DeiT-Small architecture by approximately 5% on average of small dataset.

Although the T2T-ViT-12 and Swin transformers already incorporate some degree of locality through their modified architectural designs (i.e., the token-to-token module in T2T-ViT-12 processes overlapping windows; and the Swin has a pyramid architecture that combines neighboring tokens at each stage), the integration of the LIFE module further improves performance by providing multi-scale locality in every block. Overall, the improvement gain is more significant for DeiT, which lacks local information in its architecture, compared to the T2T and Swin Transformers.

We also report the number of parameters and GMACs required to infer an image size of 224×224 for T2T and DeiT, and 256×256 for the Swin transformer. The LIFE module has a negligible impact on the number of parameters and computations. In fact, T2T-ViT-12-LIFE is even more efficient than the baseline, as it has an embedding size of 252 compared to 256 for T2T-ViT-12.

Overall, these results demonstrate that the addition of the LIFE module improves the performance of ViTs on various small datasets and architectures. This effect is more prominent when the model is smaller and the dataset is more complex (with a higher number of classes and fewer samples per class). Additionally, the LIFE module can be easily integrated into different architectures with minimal memory and computational overhead.

5.2. Impact and Importance of Multiscale Embeddings

To demonstrate the effectiveness of multiscale information, we considered a modified version of the LIFE module called LIFE-OneScale, which encodes information using a single scale with a kernel size of three, similar to CvT. We evaluated this configuration by integrating it into DeiT-Tiny on the CIFAR-100 dataset. Using a single scale, we observed over 9% improvement over the baseline. However, encoding multiscale information with three kernels of sizes 1, 3, and 5 resulted in over 20% improvement over the baseline. These results emphasize the advantage of using multi-scale information.

Table 2: Comparison of performance on the CIFAR-100 dataset for models without locality (DeiT-T), with one-scale locality (DeiT-T-OneScale), and with multi-scale locality (DeiT-T-LIFE).

Model	DeiT-T	DeiT-T-OneScale	DeiT-T-LIFE
Accuracy	51.45	60.64	71.74

5.3. Object Detection

For the object detection task, we use the DETR architecture [12], which consists of a feature extractor as the backbone and a transformer encoder-decoder (EncDec). In the encoder, the features extracted by the backbone are flattened and used as query Q , key K , and value V in the self-attention layer of the encoder. In the decoder, the encoder output is used as Q and K . V is defined as a learnable parameter that is later used to predict the final details of the object.

To evaluate the effect of the LIFE module, we integrate it into both the backbone and the encoder-decoder. The DETR encoder-decoder does not include a linear transformation to generate Q , K , and V . To ensure a fair comparison, we first add a linear layer to the encoder-decoder just before self-attention to generate Q , K , and V from the backbone features, which is denoted as EncDec-Linear. Then, we replace this linear layer with the LIFE

module, referred to as EncDec-LIFE. We use the encoder-decoder without any transformation as a baseline, which is referred to as EncDec. We use the DeiT-T and Swinv2-T architectures and their LIFE variants as the backbone.

Table 3: Evaluation of the effectiveness of adding the LIFE module to ViTs in the object detection task using the DETR architecture with DeiT-T as the backbone, trained and tested on the VOC dataset.

Backbone	Head	# params	GMAC	mAP	F1-Score
Random Initialization					
DeiT-T	EncDec	27.20	14.20	27.31	37.52
	EncDec-Linear	28.38	14.50	29.76	40.24
	EncDec-LIFE	28.19	14.45	31.37	41.80
DeiT-T-LIFE	EncDec	27.27	14.28	27.09	37.49
	EncDec-Linear	28.46	14.58	29.22	39.60
	EncDec-LIFE	28.26	14.53	32.21	42.52
ImageNet-1k Initialization					
DeiT-T	EncDec	27.20	14.20	72.12	79.32
	EncDec-Linear	28.38	14.50	73.02	80.12
	EncDec-LIFE	28.19	14.45	73.51	80.42
DeiT-T-LIFE	EncDec	27.27	14.28	72.43	79.68
	EncDec-Linear	28.46	14.58	73.05	80.16
	EncDec-LIFE	28.26	14.53	73.92	80.73
Swinv2-T	EncDec	50.07	27.32	78.21	84.48
	EncDec-Linear	51.26	27.62	78.94	85.01
	EncDec-LIFE	51.06	27.57	80.57	86.13
Swinv2-T-LIFE	EncDec	50.23	27.57	80.24	85.92
	EncDec-Linear	51.41	27.88	81.27	86.73
	EncDec-LIFE	51.22	27.83	81.27	86.68

Table 3 shows the mAP and F1 scores for different combinations of the backbone and encoder-decoder. We observe that adding linear layers to generate Q , K , and V leads to improvement. Replacing the linear layer with the LIFE module leads to additional improvement, indicating that the local information from neighboring patches can aid in more accurate object detection in a scene. The best results are obtained when the LIFE module is used in both the backbone and the encoder-decoder.

We also evaluate the performance of the LIFE module with two different initializations. When models are randomly initialized, the LIFE module improves performance by $\sim+4$ mAP. When ImageNet-1k pre-trained weights are used for initialization, the LIFE module improves performance by $\sim+2$ mAP. These results suggest that the inclusion of locality bias can be beneficial when training a model from scratch with a small dataset.

Table 4: Evaluation of the effectiveness of adding the LIFE module to ViTs on the semantic segmentation task using the DeiT-T architecture as the backbone, trained and tested on the Cityscapes dataset. The results include the mean IoU (mIoU), the number of parameters, and the computation cost (in GMAC).

Model	# params	GMAC	mIoU
Random Initialization			
DeiT-T	10.13	10.50	31.89
DeiT-T-LIFE	10.21	10.58	33.73
ImageNet-1k Initialization			
DeiT-T	10.13	10.50	58.62
DeiT-T-LIFE	10.21	10.58	59.17
Swinv2-T	32.36	23.95	59.37
Swinv2-T-LIFE	32.52	24.20	60.25

5.4. Semantic Segmentation

For the semantic segmentation task, we use a simple segmentation architecture to assess the effectiveness of the LIFE module. We use the DeiT-T and Swinv2-T architectures as the backbone and a simple linear layer followed by upscaling to match the dimension of the features to the number of classes as the segmentation head. DeiT-T-LIFE and Swinv2-T-LIFE refer to models in which the LIFE module is integrated into all the attention layers in the backbone.

Table 3 shows the performance in terms of the number of parameters, GMAC, and mIoU results for both models and two initialization methods. The LIFE module improves performance in both cases. When the model is initialized with random weights, the LIFE module improves performance by $\sim 2\%$. If ImageNet-1k pre-trained weights are used for initialization, the improvement is 0.55%. Similar to the object detection task, we observe a greater improvement when the model is trained from scratch.

6. Qualitative Results

To understand the decision-making process of a transformer model, we present the results of qualitative analyses using attention maps derived from transformer architectures. First, we generate attention visualizations for the classification task employing the existing *Attention Roll-Out* method [31]. Then, we propose the *Dense Roll-Out* method, which generates class-specific attention maps for dense prediction tasks, and demonstrate the effectiveness of the LIFE module utilizing our proposed visualization method.

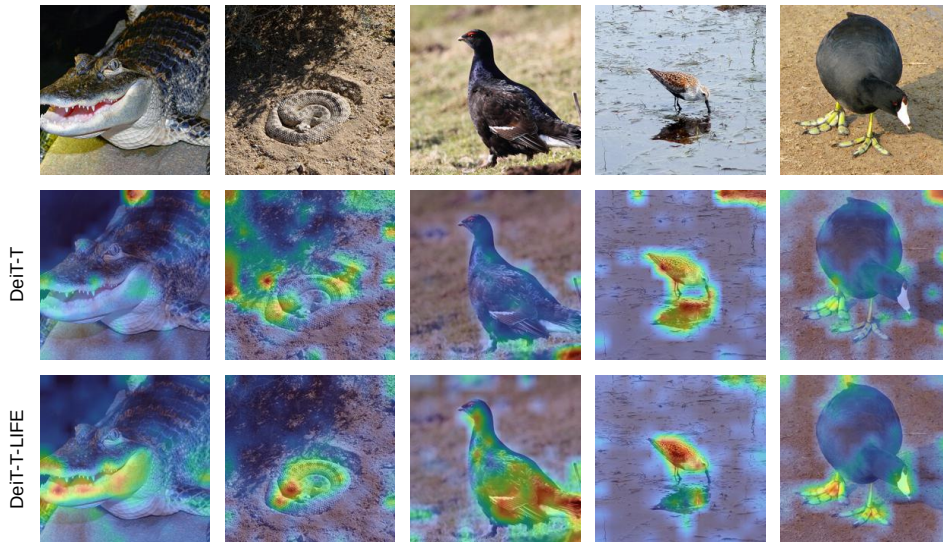


Figure 4: Comparison of DeiT-T attention maps with and without the LIFE module, trained on the ImageNet-100 dataset, using the Attention Roll-Out method [31].

6.1. Image Classification

We present a qualitative analysis of attention maps generated from DeiT-S with and without the LIFE module, trained on the ImageNet-100 dataset. As shown in Figure 4, the attention of a standard DeiT-T is scattered between the background and the foreground. However, when the LIFE module is used, the ViT architecture gains more local information, enabling it to better identify the foreground object and focus more on it.

6.2. Dense Prediction Tasks

Dense Attention Roll-Out. Attention mechanisms in transformer architectures allow the model to consider contextual information about the relationships between input tokens within a transformer block. Attention visualizations, which display the attention weights assigned to different input tokens, can provide insight into how the model makes its decisions and how well it can generalize to unseen data. There have been several methods proposed in the literature for generating attention visualizations in classification tasks, such as Attention Roll-Out [31] and Gradient Attention Roll-Out [32]. These methods have been effective for classification tasks in vision transformers, but there is currently no method proposed for dense prediction tasks, such as object detection and semantic segmentation.

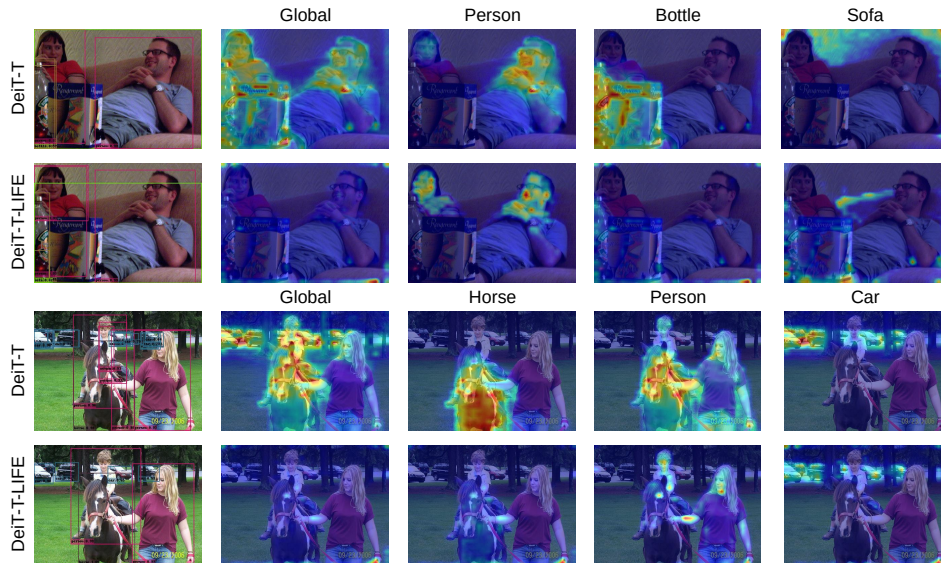


Figure 5: Comparison of attention maps of object detection models with and without the LIFE module using the proposed Dense Roll-Out method.

To address this gap, we propose a new method, *Dense Attention Roll-Out*, to generate attention visualizations in dense prediction tasks. This method is based on the Attention Roll-Out method, which creates a pairwise attention graph by linearly combining attention from all blocks in a transformer architecture. However, we modified the method to specifically target dense prediction tasks. In contrast to the standard method, which only uses the attention map corresponding to the classification token, our method utilizes attention maps of all tokens, since all patch tokens are utilized for dense prediction tasks. Additionally, we use simple heads to generate attention maps using the features of the backbone, where most of the relevant information for prediction is located.

To generate class-specific attention maps, we use network predictions to identify relevant tokens for predicting a particular class by aligning the tokens spatially with the predictions, using either a segmentation map or a bounding box depending on the task. We then take the average of the attention maps for these corresponding tokens, remove the global content, and calculate the final class-specific attention map. The global content is common for all tokens that contain global information from the input image and is obtained by taking the average of all tokens.

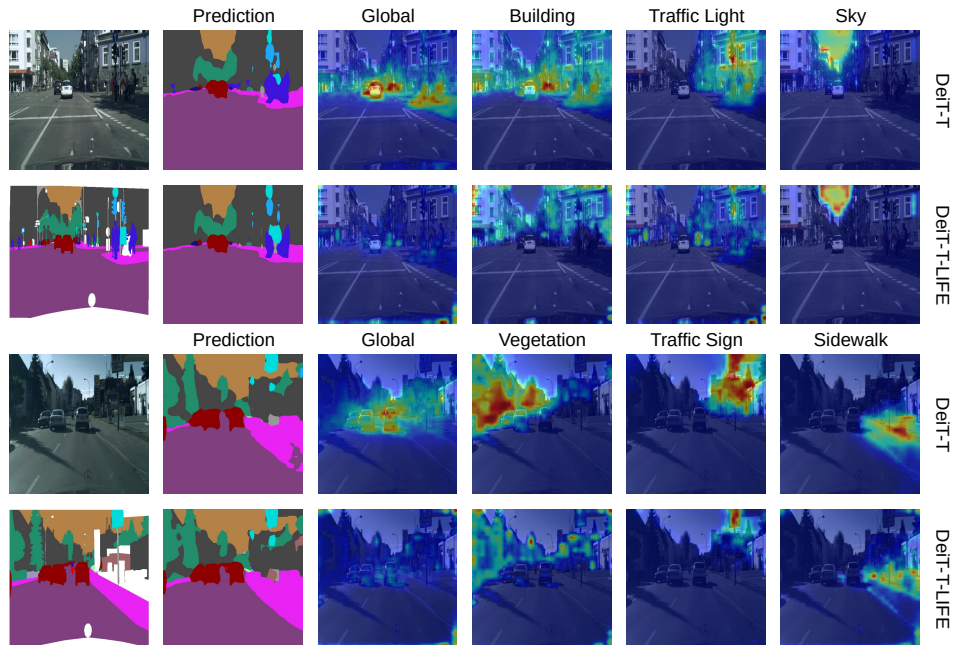


Figure 6: Comparison of attention maps of segmentation models with and without the LIFE module using the proposed Dense Roll-Out method.

Result. In order to demonstrate the effect of the LIFE module, we chose simple heads for both object detection and segmentation tasks. The use of simple heads allows us to generate visualizations using the backbone features, where the information for prediction is primarily located. Figures 5 and 6 show a comparison of models with and without the LIFE module in the backbone for object detection and segmentation tasks, respectively. The first three columns depict the input image, prediction, and global attention. The remaining columns show the class-specific attention maps, with the class names displayed at the top. The attention maps illustrate that the LIFE module helps the model focus more on the relevant regions of the input image, leading to more accurate predictions.

7. Conclusion

We proposed a novel module, *Local InFormation Enhancer (LIFE)*, for vision transformers (ViTs) that effectively leverages local information from images to generate more informative embeddings for the self-attention layers

in each transformer block. We evaluated the impact of the LIFE module on models with different sizes (DeiT-Tiny and DeiT-Small) and architectures (DeiT, T2T, and Swin). The classification results showed that our proposed method consistently improved performance on small datasets. Additionally, we found that the improvement was more significant when the model had a lower capacity, indicating that the model effectively benefits from the locality bias introduced by the LIFE module. The LIFE module was also effective in the hierarchical transformer architecture, which inherently utilizes multi-scale information, suggesting that locality is important even at the window level. Furthermore, the incorporation of the LIFE module into ViTs for object detection and segmentation tasks improved performance, demonstrating the versatility and effectiveness of the LIFE module in various tasks. We also observed that initializing the model with parameters learned from the ImageNet dataset led to greater improvement compared to training the model from scratch, indicating that the local information encoded by the LIFE module helps ViTs achieve improved performance with minimal computational or memory overhead in a range of vision applications, including classification, object detection, and semantic segmentation. Finally, we introduced a new visualization method, Dense Attention Roll-Out, specifically tailored for dense prediction tasks such as object detection and semantic segmentation. This method allows for the generation of class-specific attention maps using attention maps of all tokens, providing insight into the decision-making process and generalization capabilities of a transformer architecture.

Appendix A. Comparison with Related Works

The design of several architectures, such as ConViT, Crossformer, and CvT, incorporates locality. This section aims to clarify the novelty of the LIFE model compared to these existing approaches.

ConViT utilizes a gated positional self-attention mechanism (GPSA) that integrates a positional self-attention component with a “soft” convolutional inductive bias. The self-attention block is enhanced by including positional information to achieve this objective. In contrast, LIFE alters the creation of query, key, and value embeddings, incorporating various convolutional operations to incorporate multiscale information into these embeddings before the self-attention operation.

Crossformer proposes two modules, namely the Cross-scale Embedding Layer (CEL) and the Long-Short Distance Attention (LSDA) modules, to

incorporate locality into their design. The CEL module integrates various patches with different scales into each embedding to provide cross-scale features to the self-attention module. The LSDA module divides the self-attention module into short-distance and long-distance components to reduce computational load while maintaining small- and large-scale features in the embeddings. Our LIFE module shares a similar objective to the CEL module, as both leverage multiscale information. However, the LIFE module offers a more efficient and adaptable design compared to the CEL module. Specifically, the LIFE module utilizes depthwise separable convolution, which provides superior efficiency compared to the CEL module. Additionally, the LIFE module can handle auxiliary tokens, improving its flexibility for integration with diverse architectures. On the contrary, the Crossformer design does not incorporate auxiliary tokens, which prevents the CEL module from having this functionality.

The introduction of locality in CvT has been accomplished through two mechanisms: a convolutional patch embedding layer and a convolutional query, key, and value projection. The convolutional projection layer in CvT is similar to our own work; however, CvT employs only a single scale with a kernel size of three. In contrast, our LIFE module encodes multiscale information utilizing three scales with kernel sizes of 1, 3, and 5.

To perform a comparative analysis of the efficacy of the LIFE module with state-of-the-art approaches, we evaluated the performance of the CvT-21 model, which has 30 million parameters, on the CIFAR-100 dataset. Our results indicate that CvT-21 achieves an accuracy of 76.07%. The Swin Transformer with the LIFE module achieved an accuracy of 77.48%. This finding suggests that the LIFE module, in combination with a strong backbone, can outperform the CvT-21 model even with fewer parameters.

It is important to note that the objective of our study is to introduce locality into the transformer architecture in a modular and efficient manner. To achieve this goal, we proposed a LIFE module that efficiently leverages multi-scale information and adapts to a variety of architectures. We demonstrated its successful integration into a diverse set of architectures, including the standard transformer architecture (DeiT), the pyramid architecture (Swin), and a modified patch embedding of a standard transformer (T2T). Our findings indicate that the integration of the LIFE module effectively incorporates local information and leads to improved performance on small datasets.

References

- [1] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, I. Polosukhin, Attention is all you need, *Advances in neural information processing systems* 30 (2017).
- [2] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al., An image is worth 16x16 words: Transformers for image recognition at scale, *arXiv preprint arXiv:2010.11929* (2020).
- [3] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, B. Guo, Swin transformer: Hierarchical vision transformer using shifted windows, in: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 10012–10022.
- [4] K. Jeeveswaran, S. Kathiresan, A. Varma, O. Magdy, B. Zonooz, E. Arani, A comprehensive study of vision transformers on dense prediction tasks, *arXiv preprint arXiv:2201.08683* (2022).
- [5] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [6] M. Tan, Q. Le, Efficientnet: Rethinking model scaling for convolutional neural networks, in: *International conference on machine learning*, PMLR, 2019, pp. 6105–6114.
- [7] S. d’Ascoli, H. Touvron, M. L. Leavitt, A. S. Morcos, G. Biroli, L. Sagun, Convit: Improving vision transformers with soft convolutional inductive biases, in: *International Conference on Machine Learning*, PMLR, 2021, pp. 2286–2296.
- [8] L. Yuan, Y. Chen, T. Wang, W. Yu, Y. Shi, Z.-H. Jiang, F. E. Tay, J. Feng, S. Yan, Tokens-to-token vit: Training vision transformers from scratch on imagenet, in: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 558–567.
- [9] C.-F. R. Chen, Q. Fan, R. Panda, Crossvit: Cross-attention multi-scale vision transformer for image classification, in: *Proceedings of the*

- IEEE/CVF International Conference on Computer Vision, 2021, pp. 357–366.
- [10] F. Chollet, Xception: Deep learning with depthwise separable convolutions, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 1251–1258.
 - [11] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, H. Jégou, Training data-efficient image transformers & distillation through attention, in: International Conference on Machine Learning, PMLR, 2021, pp. 10347–10357.
 - [12] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, S. Zagoruyko, End-to-end object detection with transformers, in: European conference on computer vision, Springer, 2020, pp. 213–229.
 - [13] S. Zheng, J. Lu, H. Zhao, X. Zhu, Z. Luo, Y. Wang, Y. Fu, J. Feng, T. Xiang, P. H. Torr, et al., Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2021, pp. 6881–6890.
 - [14] B. Cheng, A. Schwing, A. Kirillov, Per-pixel classification is not all you need for semantic segmentation, *Advances in Neural Information Processing Systems* 34 (2021).
 - [15] K. Han, A. Xiao, E. Wu, J. Guo, C. Xu, Y. Wang, Transformer in transformer, *Advances in Neural Information Processing Systems* 34 (2021) 15908–15919.
 - [16] W. Wang, E. Xie, X. Li, D.-P. Fan, K. Song, D. Liang, T. Lu, P. Luo, L. Shao, Pyramid vision transformer: A versatile backbone for dense prediction without convolutions, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2021, pp. 568–578.
 - [17] W. Wang, E. Xie, X. Li, D.-P. Fan, K. Song, D. Liang, T. Lu, P. Luo, L. Shao, Pvt v2: Improved baselines with pyramid vision transformer, *Computational Visual Media* 8 (2022) 415–424.
 - [18] Z. Liu, H. Hu, Y. Lin, Z. Yao, Z. Xie, Y. Wei, J. Ning, Y. Cao, Z. Zhang, L. Dong, et al., Swin transformer v2: Scaling up capacity and resolution,

- in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2022, pp. 12009–12019.
- [19] K. Yuan, S. Guo, Z. Liu, A. Zhou, F. Yu, W. Wu, Incorporating convolution designs into visual transformers, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2021, pp. 579–588.
 - [20] H. Wu, B. Xiao, N. Codella, M. Liu, X. Dai, L. Yuan, L. Zhang, Cvt: Introducing convolutions to vision transformers, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2021, pp. 22–31.
 - [21] W. Wang, L. Yao, L. Chen, D. Cai, X. He, W. Liu, Crossformer: A versatile vision transformer based on cross-scale attention, arXiv e-prints (2021) arXiv–2108.
 - [22] Z. Zhang, H. Zhang, L. Zhao, T. Chen, S. Ö. Arik, T. Pfister, Nested hierarchical transformer: Towards accurate, data-efficient and interpretable visual understanding, in: Proceedings of the AAAI Conference on Artificial Intelligence, volume 36, 2022, pp. 3417–3425.
 - [23] Y. Liu, E. Sangineto, W. Bi, N. Sebe, B. Lepri, M. Nadai, Efficient training of visual transformers with small datasets, Advances in Neural Information Processing Systems 34 (2021) 23818–23830.
 - [24] K. Li, R. Yu, Z. Wang, L. Yuan, G. Song, J. Chen, Locality guidance for improving vision transformers on tiny datasets, arXiv preprint arXiv:2207.10026 (2022).
 - [25] A. Krizhevsky, G. Hinton, et al., Learning multiple layers of features from tiny images (2009).
 - [26] Y. Le, X. Yang, Tiny imagenet visual recognition challenge (????).
 - [27] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, L. Fei-Fei, Imagenet: A large-scale hierarchical image database, in: 2009 IEEE conference on computer vision and pattern recognition, Ieee, 2009, pp. 248–255.
 - [28] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, A. Zisserman, The pascal visual object classes (voc) challenge, International journal of computer vision 88 (2010) 303–338.

- [29] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, B. Schiele, The cityscapes dataset for semantic urban scene understanding, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 3213–3223.
- [30] Anonymous, A comprehensive study of real-time object detection networks across multiple domains: A survey, Submitted to Transactions on Machine Learning Research (2022). URL: <https://openreview.net/forum?id=ywr5sWqQt4>, under review.
- [31] S. Abnar, W. Zuidema, Quantifying attention flow in transformers, arXiv preprint arXiv:2005.00928 (2020).
- [32] H. Chefer, S. Gur, L. Wolf, Transformer interpretability beyond attention visualization, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2021, pp. 782–791.