
Peripheral Vision Transformer

Juhong Min¹ Yucheng Zhao^{2,3} Chong Luo² Minsu Cho¹

¹Pohang University of Science and Technology (POSTECH)

²Microsoft Research Asia (MSRA)

³University of Science and Technology of China (USTC)

<http://cvlab.postech.ac.kr/research/PerViT/>

Abstract

Human vision possesses a special type of visual processing systems called *peripheral vision*. Partitioning the entire visual field into multiple contour regions based on the distance to the center of our gaze, the peripheral vision provides us the ability to perceive various visual features at different regions. In this work, we take a biologically inspired approach and explore to model peripheral vision in deep neural networks for visual recognition. We propose to incorporate peripheral position encoding to the multi-head self-attention layers to let the network learn to partition the visual field into diverse peripheral regions given training data. We evaluate the proposed network, dubbed PerViT, on ImageNet-1K and systematically investigate the inner workings of the model for machine perception, showing that the network learns to perceive visual data similarly to the way that human vision does. The performance improvements in image classification over the baselines across different model sizes demonstrate the efficacy of the proposed method.

1 Introduction

For the last ten years, convolution has been a dominant feature transformation in neural networks for visual recognition due to its superiority in modelling spatial configurations of images [21, 32, 34]. Despite the efficacy in learning visual patterns, the local and stationary nature of convolutional kernels limited the maximum extent of representation ability in flexible processing, *e.g.*, dynamic transformations with global receptive fields. Originally devised for natural language processing (NLP), self-attention [63] shed a light on this direction; equipped with adaptive input processing and the ability to capture long-range interactions, it has emerged as an alternative feature transform for computer vision, being widely adopted as a core building block [18]. The stand-alone self-attention models, *e.g.*, ViT [18], however, demand significantly more training data [57] for competitive performance with its convolutional counterparts [6, 25, 27, 72] since they miss certain desirable property which convolution possesses, *e.g.*, locality. These inherent pros and cons of convolution and self-attention encourage recent researches toward combinations of both so as to enjoy the best of the both worlds but which one suits the best for effective visual processing is yet controversial in literature [8, 9, 10, 12, 35, 37, 38, 40, 41, 49, 50, 59, 60, 62, 67, 69, 71, 73, 77].

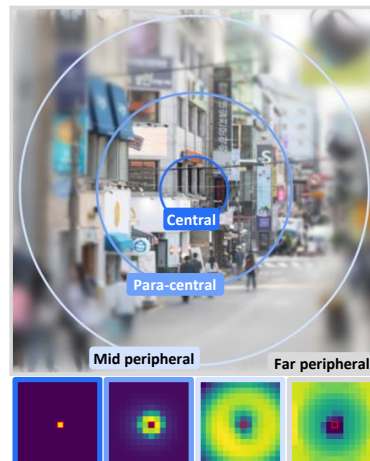


Figure 1: This work explores blending human peripheral vision (top) with attention-based neural networks (bottom) for visual recognition.

Unlike the dominant visual feature transformations in machine vision, human vision possesses a special type of visual processing systems called *peripheral vision* [36]; it partitions the entire visual field into multiple contour regions based on the distances to the center of a gaze where each region identifies different visual aspects. As seen in Fig. 1, we have high-resolution processing near the center of our gaze, *i.e.*, central and para-central regions, to identify highly-detailed visual elements such as geometric shapes, and low-level details. For the regions more distant from the gaze, *i.e.*, mid and far peripheral regions, the resolution decreases to recognize abstract visual features such as motion, and high-level contexts. This systematic strategy enables us to effectively perceive important details within a small fraction (1%) of the visual field while minimizing unnecessary processing of background clutter in the rest (99%), thus facilitating efficient visual processing for human brain.

According to recent study on inner workings of vision transformers [12, 18, 49, 58, 71], their behaviors are related to the aforementioned visual processing strategies in the following respect; the attention maps of early layers are learned to locally capture *fine-grained geometric details* at *central regions* while those of later layers perform global attentions to identify *coarse-grained semantics and contexts* from the whole visual field, covering *peripheral regions*. This finding reveals that imitating biological designs may potentially help in modelling an effective machine vision, and also support recent approaches towards a hybrid [5, 10, 12, 24, 39] of convolution and self-attention beyond stand-alone visual processing to take advantages of the two different perception strategies: fine-grained/local and coarse-grained/global, similarly to the human visual processing as in Fig. 1.

In this work, we take a biologically inspired approach and propose to inject the peripheral inductive biases¹ to deep neural networks for image recognition. We propose to incorporate peripheral attention mechanism to the multi-head self-attention [63] to let the network learn to partition the visual field into diverse peripheral regions given training data where each region captures different visual features. We experimentally show that the proposed network models effective visual periphery for reliable visual recognition. Our main contributions can be summarized as follows:

- This work explores to narrow the gap between human and machine vision by injecting peripheral inductive biases to self-attention layers, and presents a new form of feature transformation named **Multi-head Peripheral Attention (MPA)**.
- Based on the MPA, we introduce **Peripheral Vision Transformer (PerViT)**, and systematically study the inner workings of PerViT by qualitatively and quantitatively analyzing its learned attentions, which reveal that the network learns to perceive visual elements similarly to the way that human vision does without any special supervisions.
- The performance improvements in image classification over columnar Transformer baselines, *e.g.*, DeiT, across different model sizes demonstrate the efficacy of the proposed approach.

2 Related Work

Feature transformations in computer vision. With notable success in NLP, Transformers [14, 63] introduced a paradigm shift in computer vision [3, 4, 7, 18, 30, 50, 55, 56, 58, 59, 62]. Despite their generalization capability, pure Vision Transformers [18] require extensive amount of training data to capture spatial layout of images due to lack of certain desirable property, *e.g.*, locality. This encouraged many recent ViT work to incorporate local inductive biases via distillation of convolutional biases [58], local self-attention [35, 40, 73], a hybridization [10, 12, 37], and augmenting convolutions [8, 60, 67, 69], all of which convey a unified message: “Despite high generalizability of self-attentions, a sufficient amount of convolutional processing must be incorporated to capture the spatial configurations of images for reliable visual processing.”

Position encoding for Transformer. Witnessing the efficacy of position encoding in capturing input structures in NLP [11, 28, 54], recent vision models [18, 50, 70] have begun employing position encodings for images to model spatial structures of images. In particular, relative position encoding (RPE) plays a vital role for the purpose: The work of Cordonnier *et al.* [9] proves that self-attention has close relationship with convolution when equipped with certain form of RPE. Wu *et al.* [70] explore the existing RPE methods used in vision transformers [11, 28, 54, 65] and draws a conclusion that RPEs impose convolutional processing on vision transformers. Dai *et al.* [10] observe that

¹We refer peripheral inductive bias as the prioritization of our hypothesis which any attention-based neural networks can use to mimic human peripheral vision by modelling torus-shaped attentions as illustrated in Fig. 1.

depthwise convolution and self-attention can naturally be unified via RPEs. While offering promising directions, the previous RPE work, however, is limited in the sense that the focus of RPE utilization is restricted to only local attention, *e.g.*, convolution. This work exploits RPEs to devise an original visual feature transformation which naturally generalizes convolution and self-attention layers, thus enjoying the benefits of both via imitation of human visual processing system: *peripheral vision*.

Peripheral vision for machine perception. Along with central vision, peripheral vision plays a vital role in a wide range of visual recognition tasks [33]. The fundamental mechanisms of peripheral visual processing, however, have not been fully disclosed in human vision literature [52] which stimulated many researchers to reveal its inner workings and deep implications [2, 15, 16, 43, 53, 68]. The work of Rosenholtz [52] discusses pervasive myths and current findings about peripheral vision, suggesting that peripheral vision is more crucial for human perception than previously thought to perform diverse important tasks. Inspired by its importance, a number of pioneering work [16, 20, 22, 23, 44, 66] investigate the linkage between peripheral vision and machine vision, *e.g.*, CNNs, while some [31, 64] devise biologically-inspired models for the creation of stronger machine vision. Continuing previous study, this paper explores to blend human peripheral vision with attention-based neural networks, *e.g.*, vision transformer [18, 58], and introduces a new network called Peripheral Vision Transformer.

3 Our Approach

In this section, we introduce the Peripheral Vision Transformer (PerViT) which learns to model peripheral vision for effective image recognition. We first revisit the mathematical formulation of a self-attention layer and then describe how we improve it with peripheral inductive biases.

Background: Multi-Head Self-Attention. The multi-head self-attention (MHSA) [63] with N_h heads performs an attention-based feature transformation by aggregating N_h self-attention outputs:

$$\text{MHSA}(\mathbf{X}) := \text{concat}_{h \in [N_h]} [\text{Self-Attention}^{(h)}(\mathbf{X})] \mathbf{W}_{\text{out}} + \mathbf{b}_{\text{out}}, \quad (1)$$

where $\mathbf{X} \in \mathbb{R}^{HW \times D_{\text{emb}}}$ is a set of input tokens and $\mathbf{W}_{\text{out}} \in \mathbb{R}^{N_h D_h \times D_{\text{emb}}}$ and $\mathbf{b}_{\text{out}} \in \mathbb{R}^{D_{\text{emb}}}$ are the transformation parameters. The N_h outputs of self-attention are designed to extract a diverse set of features from the input representation. Formally, the self-attention at head h is defined as

$$\text{Self-Attention}^{(h)}(\mathbf{X}) := \text{Normalize}[\Phi^{(h)}(\mathbf{X})] \mathbf{V}^{(h)}, \quad (2)$$

where $\text{Normalize}[\cdot]$ denotes a row-wise normalization and $\Phi^{(h)}(\cdot) \in \mathbb{R}^{HW \times HW}$ is a function that provides spatial attentions based on content information to aggregate the values $\mathbf{V}^{(h)} = \mathbf{X} \mathbf{W}_{\text{val}}^{(h)}$:

$$\Phi^{(h)}(\mathbf{X}) := \exp(\tau \mathbf{X} \mathbf{W}_{\text{qry}}^{(h)} (\mathbf{X} \mathbf{W}_{\text{key}}^{(h)})^\top) = \exp(\tau \mathbf{Q}^{(h)}, \mathbf{K}^{(h)\top}), \quad (3)$$

using linear projections of $\mathbf{W}_{\text{qry}}^{(h)}$, $\mathbf{W}_{\text{key}}^{(h)}$, $\mathbf{W}_{\text{val}}^{(h)} \in \mathbb{R}^{D_{\text{emb}} \times D_h}$ for queries, keys, and values respectively where τ is softmax temperature and $\exp(\cdot)$ applies an element-wise exponential to the input matrix.

3.1 Multi-head Peripheral Attention

Based on the formulation of MHSA in Eq. 1, we define **Multi-head Peripheral Attention** (MPA) as

$$\text{MPA}(\mathbf{X}) := \text{concat}_{h \in [N_h]} [\text{Peripheral-Attention}^{(h)}(\mathbf{X}, \mathbf{R})] \mathbf{W}_{\text{out}} + \mathbf{b}_{\text{out}}, \quad (4)$$

where $\mathbf{R} \in \mathbb{R}^{HW \times HW \times D_r}$ is the relative position encoding with D_r channel dimension. The self-attention in MHSA is now replaced with **Peripheral-Attention**(\cdot), consisting of content- and position-based attention functions $\Phi_c^{(h)}(\mathbf{X})$, $\Phi_p^{(h)}(\mathbf{R}) \in \mathbb{R}^{HW \times HW}$, which is formulated as follows:

$$\text{Peripheral-Attention}^{(h)}(\mathbf{X}, \mathbf{R}) := \text{Normalize} \left[\Phi_c^{(h)}(\mathbf{X}) \odot \Phi_p^{(h)}(\mathbf{R}) \right] \mathbf{V}^{(h)}, \quad (5)$$

where \odot is Hadamard product which mixes the given pair of attentions to provide a mixed attention $\Phi_a^{(h)}(\mathbf{X}, \mathbf{R}) := \Phi_c^{(h)}(\mathbf{X}) \odot \Phi_p^{(h)}(\mathbf{R}) \in \mathbb{R}^{HW \times HW}$. For the content-based attention Φ_c , we use exponentiated (scaled) dot-product between queries and keys as in Eq. 3: $\Phi_c^{(h)}(\mathbf{X}) := \exp(\tau \mathbf{Q}^{(h)} \mathbf{K}^{(h)\top})$. For the position-based attention Φ_p , we design a neural network that aims to imitate human visual processing system, *e.g.*, peripheral vision, which we discuss next.

Modelling peripheral vision: a Roadmap. Human visual field can be grouped into several regions based on the Euclidean distances from the center of gaze, each forming ring-shaped region as seen in Fig. 1, where each region captures different visual aspects; the closer to the gaze, the more complex features we process, and further from the gaze, the simpler visual features we perceive. In the context of 2-dimensional attention map $\Phi_*(\cdot)_{\mathbf{q},:} \in \mathbb{R}^{HW}$, we refer the query position $\mathbf{q} \in \mathbb{R}^2$ as the center of gaze, *i.e.*, the position where feature of our interest lies at for the transformation. We refer the local regions around the query \mathbf{q} as central/para-central regions and the rest as mid/far peripheral regions.

Perhaps the simplest approach to divide the visual field into multiple subregions is to perform a single linear projection on the Euclidean distances, *i.e.*, $\Phi_p^{(h)}(\mathbf{R}) = \sigma[\mathbf{R}\mathbf{W}_p^{(h)}]$ where $\mathbf{W}_p^{(h)} \in \mathbb{R}^{D_r}$ and $\sigma[\cdot]$ is non-linearity, similarly to the previous work of Wu *et al.* [70]². For straightforward imitation of peripheral vision, we use Euclidean distance for relative position input \mathbf{R} and weigh the distances in D_r different ways for the network to learn the mapping in multiple scales: $\mathbf{R}_{\mathbf{q},\mathbf{k},:} := \text{concat}_{r \in [D_r]} [w_r \cdot \mathbf{R}_{\mathbf{q},\mathbf{k}}^{\text{euc}}] \in \mathbb{R}^{D_r}$ where $\{w_r\}_{r \in [D_r]}$ is a set of learnable parameters shared across layers and heads, and $\mathbf{R}_{\mathbf{q},\mathbf{k}}^{\text{euc}} = \|\mathbf{q} - \mathbf{k}\|_2$ is the Euclidean distance between query and key positions $\mathbf{q}, \mathbf{k} \in \mathbb{R}^2$. In our experiments, we choose sigmoid function for σ to provide normalized (positive) weights to the content-based attention Φ_c . A main drawback of this single-layer formulation is that Φ_p is only able to provide Gaussian-like attention map as seen in top-left in Fig. 2, thus being unable to represent peripheral regions in diverse shapes. For the encoding function to represent various (torus-shaped) peripheral regions, the distances must be processed by an MLP:

$$\Phi_p^{(h)}(\mathbf{R}) = \sigma \left[\text{Linear}(\text{ReLU}(\text{Linear}(\mathbf{R}; \mathbf{W}_{p1})) \mathbf{W}_{p2}^{(h)}) \right] = \sigma \left[\text{ReLU}(\mathbf{R}\mathbf{W}_{p1}) \mathbf{W}_{p2}^{(h)} \right], \quad (6)$$

where $\mathbf{W}_{p1} \in \mathbb{R}^{D_r \times D_{\text{hid}}}$ and $\mathbf{W}_{p2}^{(h)} \in \mathbb{R}^{D_{\text{hid}}}$ are the linear projection parameters³, and ReLU gives non-linearity to the function. The first projection \mathbf{W}_{p1} is shared across the heads in order to exchange information so each function is able to provide attention that are effective or complementary to other heads' attention. Note that given identical relative distances between a fixed query point $\mathbf{q} \in \mathbb{R}^2$ and key points $\mathbf{k}_i, \mathbf{k}_j \in \mathbb{R}^2$, *i.e.*, $\mathbf{R}_{\mathbf{q},\mathbf{k}_i} = \mathbf{R}_{\mathbf{q},\mathbf{k}_j}$, Eq. 6 provides the same attention scores: $\Phi_p(\mathbf{R})_{\mathbf{q},\mathbf{k}_i} = \Phi_p(\mathbf{R})_{\mathbf{q},\mathbf{k}_j}$ as seen in top-right of Fig. 2. This property, however, is not always desired in practical scenarios because the rotational symmetric property hardly holds for most real-world objects. To break the symmetric property in Eq. 6 while preserving peripheral design to sufficient extent, we introduce **peripheral projections** in which the transformation parameters are given small spatial resolutions, *e.g.*, $K \times K$ window, such that $\mathbf{W}_{p1} \in \mathbb{R}^{K^2 \times D_r \times D_{\text{hid}}}$ and $\mathbf{W}_{p2}^{(h)} \in \mathbb{R}^{K^2 \times D_{\text{hid}}}$ so that they provide similar but different attention scores, $\Phi_p(\mathbf{R})_{\mathbf{q},\mathbf{k}_i} \neq \Phi_p(\mathbf{R})_{\mathbf{q},\mathbf{k}_j}$, given $\mathbf{R}_{\mathbf{q},\mathbf{k}_i} = \mathbf{R}_{\mathbf{q},\mathbf{k}_j}$, by aggregating neighboring relative distances around the key location \mathbf{k} as follows:

$$\Phi_p^{(h)}(\mathbf{R})_{\mathbf{q},\mathbf{k},:} := \sigma \left[\sum_{\mathbf{n} \in \mathcal{N}(\mathbf{k})} \text{ReLU} \left(\sum_{\mathbf{m} \in \mathcal{N}(\mathbf{k})} \mathbf{R}_{\mathbf{q},\mathbf{m},:} \mathbf{W}_{p1} \mathbf{m} - \mathbf{k}, : \right) \mathbf{W}_{p2}^{(h)}_{\mathbf{n} - \mathbf{k},:} \right], \quad (7)$$

where $\mathcal{N}(\mathbf{k}) := [\mathbf{k} - \lfloor \frac{K}{2} \rfloor, \dots, \mathbf{k} + \lfloor \frac{K}{2} \rfloor] \times [\mathbf{k} - \lfloor \frac{K}{2} \rfloor, \dots, \mathbf{k} + \lfloor \frac{K}{2} \rfloor]$ is a set of K^2 neighbors around input position \mathbf{k} . We set $K = 3$ for all layers and heads as $K > 3$ hardly brought improvements. Note that each linear projection in Eq. 7 is equivalent to a 4-dimensional convolution [51], taking 4-dimensional input $\mathbf{R} \in \mathbb{R}^{HW \times HW \times D_r}$ to process in convolutional manner using 4-dimensional kernels in size of $K \times K \times 1 \times 1$, *i.e.*, $\mathbf{W}_{p1} \in \mathbb{R}^{K \times K \times 1 \times 1 \times D_r \times D_{\text{hid}}}$. Precisely, the peripheral projection considers a small subset of 4D local neighbors that pivots the query position \mathbf{q} , similarly to the center-pivot 4D convolution [45, 46]. After each peripheral projection, we add an instance normalization layer [61] for stable optimization:

$$\mathbf{R}' = \text{ReLU}(\text{IN}(\text{PP}(\mathbf{R}; \mathbf{W}_{p1}); \gamma_{p1}, \beta_{p1})), \quad \Phi_p^{(h)}(\mathbf{R}) = \sigma \left(\text{IN}(\text{PP}(\mathbf{R}'; \mathbf{W}_{p2}^{(h)}); \gamma_{p2}^{(h)}, \beta_{p2}^{(h)}) \right), \quad (8)$$

where $\gamma_{p1}, \beta_{p1} \in \mathbb{R}^{D_{\text{hid}}}$ and $\gamma_{p2}^{(h)}, \beta_{p2}^{(h)} \in \mathbb{R}$ are weights/biases of the instance norms and $\text{PP}(\cdot)$ denotes the peripheral projection: $\text{PP}(\mathbf{R}, \mathbf{W})_{\mathbf{q},\mathbf{k},:} := \sum_{\mathbf{n} \in \mathcal{N}(\mathbf{k})} \mathbf{R}_{\mathbf{q},\mathbf{n},:} \mathbf{W}_{\mathbf{n} - \mathbf{k},:}$. The middle row of Fig. 2 depicts learned attentions of Φ_p with peripheral projections, which provides peripheral attention maps in greater diversity compared to single- and multi-layer counterparts without \mathcal{N} .

²Given $\sigma[\cdot] := \exp(\cdot)$, Peripheral-Attention^(h) = Normalize $[\exp(\mathbf{Q}^{(h)} \mathbf{K}^{(h)\top}) \odot \exp(\mathbf{R}\mathbf{W}_p^{(h)})] \mathbf{V}^{(h)}$ = softmax $(\mathbf{Q}^{(h)} \mathbf{K}^{(h)\top} + \mathbf{R}\mathbf{W}_{\text{rpe}}) \mathbf{V}^{(h)}$, which is equivalent to the *bias mode* RPE presented in [70].

³We omit the bias terms in the linear layers for brevity.

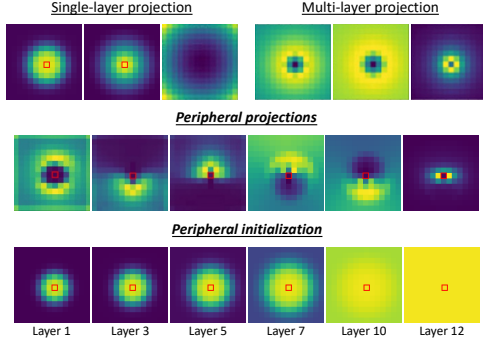


Figure 2: Top: representation ability of Φ_p under varying # layers. Middle: peripheral projections. Bottom: peripheral initialization.

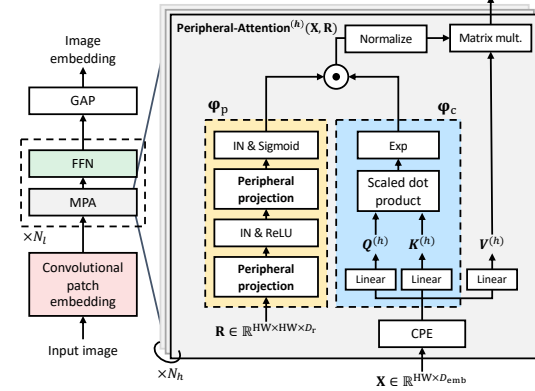


Figure 3: Overall architecture of PerViT which is based on DeiT [12] architecture.

Peripheral initialization. Recent study [49, 58] observe that early layers of trained vision transformer learn to attend locally whereas late layers perform global attentions. To facilitate training of our network, we inject this property in the beginning of the training stage by initializing parameters of Φ_p for the purpose, making attention scores near the queries larger than the distant ones in the early layers while uniformly distributing them in the late layers as seen in bottom row of Fig. 2. We refer this method as *peripheral initialization* for its resemblance to the functioning of peripheral vision [36] which also operates either locally or globally to perceive different visual features [76]. To formally put, given two arbitrarily chosen distances $\mathbf{R}_{q,k_i}^{\text{euc}}, \mathbf{R}_{q,k_j}^{\text{euc}} \in \mathbb{R}$ which satisfy $\mathbf{R}_{q,k_i}^{\text{euc}} < \mathbf{R}_{q,k_j}^{\text{euc}}$, we want $\Phi_p^{(l,h)}(\mathbf{R})_{q,k_i} \gg \Phi_p^{(l,h)}(\mathbf{R})_{q,k_j}$ ⁴ as $l \rightarrow 1$, *i.e.*, local attention in early layers, and $\Phi_p^{(l,h)}(\mathbf{R})_{q,k_i} \approx \Phi_p^{(l,h)}(\mathbf{R})_{q,k_j}$ as $l \rightarrow N_l$, *i.e.*, global attention in late layers, where N_l is the total number of MPA layers. We first initialize the parameters of $\Phi_p^{(l,h)}$ and $\{w_r\}_{r \in [D_r]}$ to particular values. Specifically, for all layers $l \in [N_l]$ and heads $h \in [N_h]$,

$$w_r := -c_1, \forall r \in [D_r] \quad \mathbf{W}_{p1}^{(l)} := c_2 J_{K^2, D_r, D_{\text{hid}}} \quad \mathbf{W}_{p2}^{(l,h)} := c_2 J_{K^2, D_{\text{hid}}} \quad \gamma_{p1}^{(l)} := \mathbf{1} \quad \beta_{p1}^{(l)} := \mathbf{0} \quad (9)$$

where $c_1, c_2 \in \mathbb{R}^+$ are positive reals, and $J_{N,M} \in \mathbb{R}^{N \times M}$ refers to all-one matrix in size $N \times M$. The above initialization provides local attention after the second peripheral projection, *i.e.*, $\text{PP}(\mathbf{R}'; \mathbf{W}_{p2}^{(h)})_{q,k_i} > \text{PP}(\mathbf{R}'; \mathbf{W}_{p2}^{(h)})_{q,k_j}$: given $\mathbf{R}_{q,k_i}^{\text{euc}} < \mathbf{R}_{q,k_j}^{\text{euc}}$. Next, based on our findings that biases $\beta_{p2}^{(l,h)}$ and the weights $\gamma_{p2}^{(l,h)}$ in the second instance norm control the sizes and strengths of local attention respectively, we simulate peripheral initialization by setting their initial values as $\beta_{p2}^{(l,h)} := s_l$ and $\gamma_{p2}^{(l,h)} := v_l$ where respective $\{s_l\}_{l \in [N_l]}$ and $\{v_l\}_{l \in [N_l]}$ are sets of initial values for attention sizes and strengths. We set their values collected from uniform intervals: $s_l \in [-5.0, 4.0]$ and $v_l \in [3.0, 0.01]$ where $s_{l-1} < s_l$ and $v_{l-1} > v_l$ to give stronger local attentions to shallow layers compared to deep ones as seen in bottom row of Fig. 2. We set $c_1, c_2 = 0.02$ in our experiments. We refer the readers to the supplementary for the complete derivation of the peripheral initialization.

3.2 Peripheral Vision Transformer

Based on the proposed peripheral projections and initialization, we develop image classification models, dubbed Peripheral Vision Transformer, which is illustrated in Fig. 3. We follow similar architecture to DeiT [58] with convolutional patch embedding stem; as the original patchify stem [18] exhibits substandard optimizability due to its coarse-grained early visual processing [71], many recent ViT models adopt multi-resolution *pyramidal designs* [40, 67, 69, 73] to mitigate the issue. While the pyramidal models have shown their efficacy in learning reliable image embeddings, we stick with the original single-resolution *columnar design* for PerViT because features in multiple resolution make our study less interpretable, which further requires additional techniques for combining our

⁴We now use the superscript to indicate both layer and head indices for the ease of demonstration.

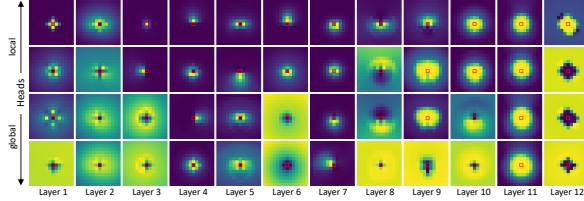


Figure 4: Learned attentions Φ_p of PerViT-T.

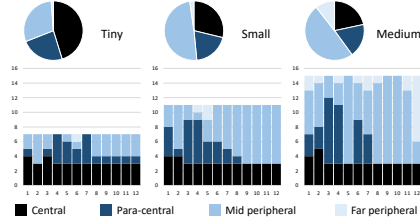


Figure 5: Peripheral region classification.

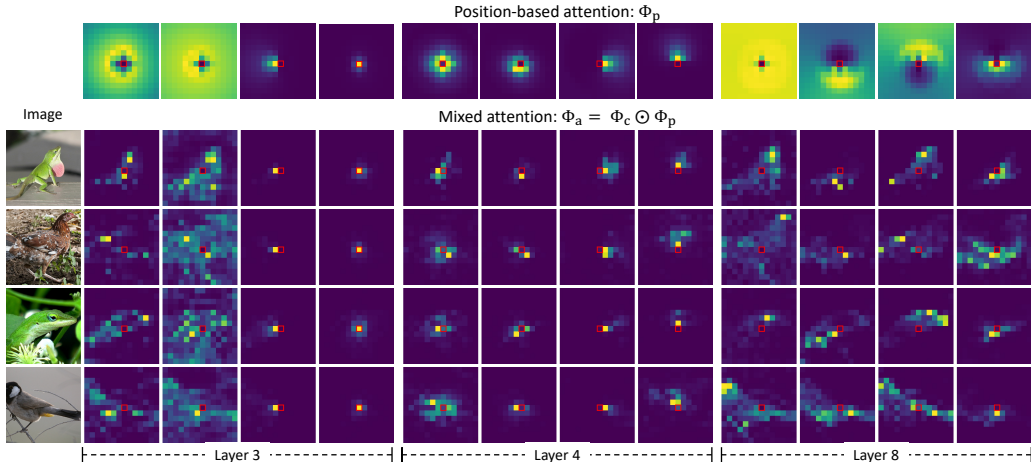


Figure 6: Visualization of learned attentions Φ_p and Φ_a for $l \in \{3, 4, 8\}$. Best viewed in electronics.

peripheral attention Φ_p with the existing cost-effective self-attentions mechanisms such as factorized attention [73], shifted-window [40], and cross-shaped window [17]. To carry out fine-grained early processing while keeping single-resolution features across the layers, we adopt convolutional patch embedding layer [71] with multi-stage layouts for channel dimensions. The convolutional embedding layer consists of four 3×3 and one 1×1 convolutions where the 3×3 convolutions are followed by batch norm [29] and ReLU [47]. For additional details, we refer to the supplementary materials.

Overall pipeline. Given an image, the convolutional patch embedding provides token embeddings $\mathbf{X}^{(1)} \in \mathbb{R}^{HW \times D_{\text{emb}}}$. Similarly to [18, 58], the embeddings are fed to a series of N_l blocks each of which consists of an MPA layer and a feed-forward network with residual pathways:

$$\mathbf{X}^{(l')} = \text{MPA}(\text{LN}(\text{CPE}(\mathbf{X}^{(l)}))) + \mathbf{X}^{(l)}, \quad \mathbf{X}^{(l+1)} = \text{FFN}(\text{LN}(\mathbf{X}^{(l')})) + \mathbf{X}^{(l')}, \quad (10)$$

where LN is layer normalization [1], and FFN is an MLP consisting of two linear transformations with a GELU activation [26]. Following the work of [35, 73], we adopt convolutional position encoding (CPE), *i.e.*, a 3×3 depth-wise convolution, before first layer norm for its efficacy with negligible computational cost. The output $\mathbf{X}^{(N_h)}$ is global-average pooled to form an image embedding.

4 Experiments

In this section, we first investigate the inner workings of PerViT trained on ImageNet-1K classification dataset to examine how it benefits from the proposed peripheral projections and initialization, and then compare the method with previous state of the arts under comparable settings.

Experimental setup. Our experiments focus on image classification on ImageNet-1K [13]. Following training recipes of DeiT [58], we train our model on ImageNet-1K from scratch with batch size of 1024, learning rate of 0.001 using AdamW [42] optimizer, cosine learning rate decay scheduler, and the same data augmentations [14] for 300 epochs, including warm-up epochs. We evaluate our model with three different sizes, *e.g.*, Tiny (T), Small (S), and Medium (M). We use stochastic depths of 0.0, 0.1, and 0.2 for T, S, and M respectively. We refer to the supplementary for additional details.

4.1 The inner workings of PerViT

Learning peripheral vision. We begin by investigating how PerViT models peripheral vision by qualitatively analyzing its learned attention of Φ_p . Figure 4 depicts the learned attention map of $\Phi_p^{(l,h)} \in \mathbb{R}^{HW}$ for all layers and heads where the query position is given at the center, *i.e.*, $\mathbf{q} = [7, 7]^T$ ⁵. We observe that the attentions are learned to be in diverse shapes of peripheral regions. Interestingly, without any special supervisions, the four attended regions ($N_h = 4$) in most layers are learned to complement each other to cover the entire visual field, capturing different visual aspects at each region (head), similarly to human peripheral vision illustrated in Fig. 1. For example, first two heads in Layer 3 attend the central regions while the others cover the rest peripheral regions. The second and third heads in Layer 8 cover top and bottom hemicircles respectively, forming a circular-shaped semi-global receptive field. Moreover, a large number of early attentions is in form of central/para-central regions while those of late layers are learned to cover mid to far peripheral regions. To quantitatively inspect how PerViT models the peripheral visual system, we classify every feature transformation layer in the network into one of the four visual regions $\mathbb{P} \in \{c, p, m, f\}$ where respective elements refer to central, para-central, mid, and far peripheral regions. PerViT-Attention of head h at layer l is classified as peripheral region p if the average of its attention scores which fall in visual region p is the largest among the others:

$$\text{PeripheralRegion}(l, h) := \arg \max_{p \in \mathbb{P}} \left[\frac{1}{|\mathcal{P}|^2} \sum_{(\mathbf{q}, \mathbf{k}) \in \mathcal{P} \times \mathcal{P}} \Phi_p^{(l,h)} \cdot \mathbb{1} [\|\mathbf{q} - \mathbf{k}\|_2 \in \mathcal{I}_p] \right], \quad (11)$$

where \mathcal{P} is a set of spatial positions ($|\mathcal{P}| = HW$) and \mathcal{I}_p is distance range (real-valued interval) of peripheral region p ⁶. The pie charts of Fig. 5 describe the proportions of peripheral regions for Tiny, Small, and Medium models where the bar graphs show them in layer-wise manner⁷. Similarly to the visualized attention maps in Fig. 4, the early layers attends central/para-central regions whereas deeper ones focus on outer region. We observe that, as the model size grows, the number of mid/far peripheral attention increases whereas that of central/para-central attention stays similar, suggesting that the models no longer require local attentions once sufficient amount of processing is done in the central region because, we hypothesize, identifying geometric patterns, *e.g.*, corners and edges, is relatively simpler process than understanding high-level semantics.

Inspecting the impact of attentions (static vs. dynamic). To study how position-based attentions Φ_p contribute to the mixed attentions $\Phi_a = \Phi_c \odot \Phi_p$, we collect sample images and visualize their attention maps of Layers 3, 4 and 8 in Fig. 6. The mixed attentions Φ_a at Layer 4 are formed dynamically (Φ_c) within statically-formed region (Φ_p) while the attentions Φ_a at Layer 8 weakly exploit position information (Φ_p) to form dynamic attentions (Φ_c). The results reveal that Φ_p plays two different roles; it imposes *semi-dynamic attention* if the attended region is focused in a small area whereas it serves as *position bias injection* when the attended region is relatively broad. In the supplementary, we constructively prove that *an MPA layer in extreme case of semi-dynamic attention/position bias injection is in turn convolution/multi-head self-attention*, naturally generalizing the both transformations. To quantitatively examine the contributions of Φ_c and Φ_p to the mixed attention Φ_a , we define a measure of ‘impact’ by taking inverse of difference between two attentions:

$$\Psi_p^{(l,h)} := [\|\Phi_a^{(l,h)} - \Phi_p^{(l,h)}\|_F]^{-1}, \quad (12)$$

where $\|\cdot\|_F$ is Frobenius norm. The higher the measure $\Psi_p^{(l,h)}$, the larger the impact of position-based attention $\Phi_p^{(l,h)}$. Being averaged over all test samples, $\Psi_c^{(l,h)}$ is similarly defined. As seen in Fig. 7, we observe a clear tendency that the impact of position-based attention is significantly higher in early processing, transforming features *semi-dynamically*, while the later layers require less position information, regarding Φ_p as a minor *position bias*. This tendency becomes more visible with larger models as seen in right of Fig 7; Small and Medium models exploit dynamic transformations much

⁵The columnar design of PerViT provides identical spatial resolution for every intermediate feature map in the network: $H, W = 14$, thus facilitating the ease of qualitative/quantitative analyses of the learned attentions.

⁶We use $\mathcal{I}_c = [0, 1.19)$, $\mathcal{I}_p = [1.19, 3.37)$, $\mathcal{I}_m = [3.37, 5.83)$, and $\mathcal{I}_f = [5.83, 7.9)$. We refer the readers to the supplementary materials for the justification on the these interval choices.

⁷We classify the 3×3 depth-wise convolution in CPE and the two linear projections in FFN as central regions as their receptive fields approximately fall in the interval of $\mathcal{I}_c = [0, 1.19)$.

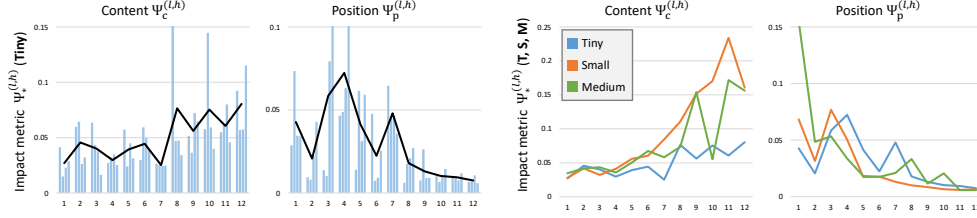


Figure 7: The measure of impact (x-axis: layer index, y-axis: the impact metric Ψ_*). Each bar graph shows the measure of a single head (4 heads at each layer), and the solid lines represent the trendlines which follow the average values of layers. (left: results of PerViT-T. right: results of T, S, and M.)

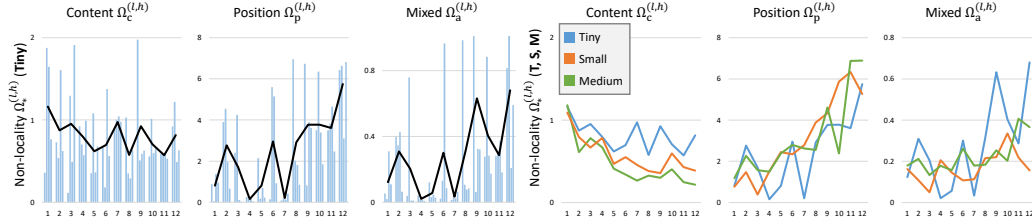


Figure 8: The measure of nonlocality (x-axis: layer index, y-axis: the nonlocality metric Ω_*).

more compared to Tiny model, especially in the later layers. Moreover, we note that the impact measures of four heads (bar graphs) within each layer are unevenly distributed, showing high variance, which imply that the network evenly utilizes both position and content information simultaneously within each MPA layer as seen in Layer 3 in Fig. 6, *performing both static/local and dynamic/global transformation in a single-shot*. These results reveal that feature transformations for effective visual recognition should not be restricted to be in either position-only [25, 41] or content-only [18, 58] design but they should be in the form of a hybrid [10, 12].

Inspecting the locality (local vs. global). We further investigate the inner workings of PerViT by quantifying how locally Φ_* attends. Following the work of [12], we define the measure of nonlocality for $\Phi_p^{(l,h)}$ by summing all pair-wise query-key distances weighted by their attention scores:

$$\Omega_p^{(l,h)} := \frac{1}{|\mathcal{P}|^2} \sum_{(\mathbf{q}, \mathbf{k}) \in \mathcal{P} \times \mathcal{P}} \Phi_p^{(l,h)} \mathbf{R}_{\mathbf{q}, \mathbf{k}}^{\text{euc}}. \quad (13)$$

The metrics Ω_c and Ω_a are similarly defined, being averaged over all test samples. As seen in Fig. 8, we observe a similar trend of locality between Φ_p and Φ_a , which reveals the position information play more dominant role over the content information in forming spatial attentions (Φ_a) for feature transformation. Interestingly, we also observe that content- and position-based attentions behave conversely; Φ_c attends globally in early layers, *i.e.*, large scores are distributed over the whole spatial region, while being relatively local in deeper layers. We hypothesize that the proposed Φ_p in early layers is trained to effectively suppress unnecessary scores of Φ_c at distant positions, thus exploiting only a few relevant ones within its local region of interest. Meanwhile, Φ_p at later layers gives Φ_c higher freedom in forming the spatial attention Φ_a as described in the plots of Fig. 7, which allows the attention scores of Φ_c to be clustered in semantically meaningful parts, *e.g.*, eyes of the animals as seen in attentions of the first head at Layer 8 (Fig. 6), which makes Φ_c relatively local.

4.2 Quantitative evaluation on ImageNet-1K

Ablation study on main components. In Tab. 1, we analyze the impact of each component in PerViT, which is denoted as (a), where C-stem refers to convolutional patch embedding stem⁸. We observe that the proposed attention Φ_p brings consistent gains to models (b, f, g, i) with relative improvements of 1.4~4.2%p for (a, d, e, h) respectively. Among three main components (Φ_p , C-stem, CPE), Φ_p has the most significant impact on PerViT (a), losing 1.5%p Top-1 accuracy without Φ_p , *i.e.*, model (b).

⁸We increase feature dimensions of the models (c) (without Φ_c) and (f, h, i) (without C-stem) accordingly to make FLOPs comparable to the others (a-i) to ensure the accuracy drops are not simply due to lower FLOPs.

Table 1: Study on the effect of each component in PerViT.

Reference	Φ_p	Φ_c	C-stem	CPE	Top-1	Top-5	FLOPs (G)
(a)	✓	✓	✓	✓	78.8	94.3	1.6
(b)	✗	✓	✓	✓	77.3	94.1	1.6
(c)	✓	✗	✓	✓	76.8	93.5	1.6
(d)	✓	✓	✗	✓	77.8	94.0	1.5
(e)	✓	✓	✓	✗	78.1	94.0	1.6
(f)	✗	✓	✗	✓	76.3	93.2	1.5
(g)	✗	✓	✓	✗	76.7	93.3	1.6
(h)	✓	✓	✗	✗	76.5	93.4	1.5
(i)	✗	✓	✗	✗	72.3	93.4	1.5

Table 2: Comparisons between different relative position encodings with DeiT-Tiny [58] as a baseline.

Method	Top-1	FLOPS (G)
DeiT-T [58]	72.2	1.3
+ CPVT [8]	73.4	2.1
+ iRPE [70]	73.7	1.1
+ PPE (ours)	74.4	1.1

Table 3: Ablation on PerViT-T/S/M: the effects of Φ_p , C-Stem, and CPE.

Ref.	Φ_p	C-Stem & CPE	T	S	M
(a)	✓	✓	78.8	82.1	82.9
(b)	✗	✓	77.3	81.1	81.9
(c)	✗	✗	72.2	79.8	81.8

Table 4: Top-1 accuracy comparisons with DeiT-S [58] under different subsampling ratios: {100%, 50%, 25%}.

Subsampling ratio	Top-1		Top-5	
	DeiT-S	PerViT-S	DeiT-S	PerViT-S
100%	79.9	82.1	95.0	95.8
50%	74.6	77.4	91.8	93.1
25%	61.8	67.5	82.6	86.9

Surprisingly, PerViT without content-based attention Φ_c , model (c), achieves decent performance, almost equalling to the accuracy of PerViT without position-based attention Φ_p , model (b) (-0.5%p). The results verify that the proposed peripheral attention, which achieves comparable level of efficacy to the content-based attention, learns to generate reliable spatial attentions for visual recognition. We also implement the proposed position-based attention Φ_p on DeiT [58] baseline and compare the results with recent state-of-the-art RPE methods. As seen in Tab. 2, the large improvements over the previous RPE methods [8, 70] further verify the efficacy of the proposed peripheral position encoding (PPE). To confirm that the impact of Φ_p is consistent with large models, we conduct similar ablations using PerViT-S/M in Tab. 3; without Φ_p , the accuracy consistently drops for all the three models. Comparing (b) with (c), we observe that C-Stem and CPE are less effective for large models, bringing 1.3%p and 0.1%p gains for Small and Medium respectively whereas they improve the Tiny model by 5.1%p. In contrast, the impact of Φ_p is consistent across different model sizes, bringing 1%p gains for all the three models. The better efficacy of Φ_p for larger models, we hypothesize, is due to its flexibility in modeling local/global spatial attentions while C-Stem/CPE are designed only to be local.

Sample-efficiency of PerViT. To investigate the training sample efficiency of our model, we train PerViT-S with ImageNet subsampled by fractions of 50% and 25%⁹ and evaluate it on full-sized test set of ImageNet-1K. Table 4 compare our results with DeiT [58]; our model consistently surpasses the baseline for all subsampled datasets, showing its robustness under limited training data.

Ablation study on Φ_p . The top section of Tab. 5 reports results of PerViT-T with different parameter initialization methods for Φ_p where peripheral denotes the proposed peripheral initialization, conv refers to convolutional initialization such that $s_l = -5.0$ and $v_l = 3.0$ for all $l \in [N_l]$, and rand refers to random initialization for all parameters in Φ_p : w_r , $\mathbf{W}_{p1}^{(l)}$, $\mathbf{W}_{p2}^{(l,h)}$, $\gamma_{p1}^{(l)}$, $\beta_{p1}^{(l)}$, $\gamma_{p2}^{(l,h)}$, and $\beta_{p2}^{(l,h)}$. The results show the efficacy of our peripheral initialization which is also supported by the results in Fig. 4 and 8: Φ_p provides early local and late global attentions, suggesting that peripheral initialization effectively reduces burden in learning such form of attentions. The bottom section of Tab. 5 studies network designs for Φ_p where \mathcal{N} represents the proposed peripheral projection, *i.e.*, projecting input distance representation by referring neighbors \mathcal{N} , ML refers to multi-layer design of Φ_p , and Euc & Lrn indicate the type of embedding \mathbf{R} : Euc is relative Euclidean distances ($\mathbf{R}_{q,k,:} = \text{concat}_{r \in [D_r]} [w_r \cdot \mathbf{R}]$) where Lrn is relative distances between learnable vectors ($\mathbf{R} \in \mathbb{R}^{HW \times HW \times D_r}$). Without \mathcal{N} and ML, we observe consistent accuracy drops for Euc and Lrn by 0.8%p and 0.2%p respectively. A sole multi-layer projection hardly improves accuracy but the model performs the best when \mathcal{N} is jointly used, meaning that both need to complement each other to provide diverse attention shapes as in Fig 4. Furthermore, Euc models consistently surpasses Lrn models, implying the Euclidean distance is more straightforward encoding type than learnable vectors in capturing spatial configurations of images.

⁹For each subsamples, we increase the number of epochs to present models with a fixed number of images.

Table 5: Ablation study on different initialization methods (top section) and network designs (bottom section) for the position-based attention Φ_p .

Init. method for Φ_p	Top-1	Top-5
Peripheral (ours)	78.8	94.3
Conv	78.6	93.8
Rand	78.5	93.6
Network design for Φ_p	Top-1	Top-5
Euc + \mathcal{N} + ML (ours)	78.8	94.3
Euc + ML	77.9	94.0
Euc	78.0	94.0
Lrn + \mathcal{N} + ML	77.8	94.0
Lrn + ML	77.5	93.8
Lrn	77.6	93.8

Table 6: Model performance on ImageNet-1K [13].

	Model	Size (M)	FLOPs (G)	Top-1 (%)
Pyramidal Vision Transformers (<i>multi-resolution</i>)	PVT-T [67]	13	1.9	75.1
	CoaT-Lite-T [73]	5.7	1.6	77.5
	Swin-T [40]	28	4.5	81.3
	CoaT-Lite-S [73]	20	4.0	81.9
	Focal-T [74]	29	4.9	82.2
	Swin-S [40]	50	8.7	83.0
Columnar Vision Transformers (<i>single-resolution</i>)	CoaT-Lite-M [73]	45	9.8	83.6
	Focal-S [74]	51	9.1	83.5
	DeiT-T [58]	5.7	1.3	72.2
	XCiT-T12/16 [19]	7.0	1.2	77.1
	PerViT-T (ours)	7.6	1.6	78.8
	DeiT-S [58]	22	4.6	79.8
Columnar Vision Transformers (<i>single-resolution</i>)	T2T-ViT _t -14 [75]	22	6.1	81.7
	XCiT-S12/16 [19]	26	4.8	82.0
	PerViT-S (ours)	21	4.4	82.1
	DeiT-B [58]	86	18	81.8
	T2T-ViT _t -24 [75]	64	15	82.6
	XCiT-S24/16 [19]	48	9.1	82.6
	PerViT-M (ours)	44	9.0	82.9

Comparison with state of the arts. Table 6 summarizes the results of our method and recent state of the arts. For fair comparison, the baselines used in our comparison are trained using 224×224 input resolution without distillations, and are grouped into either pyramidal or columnar ViT based on the network designs, *i.e.*, multi- or single-resolution feature processing, where the results are partitioned according to model sizes within each group. As shown in the bottom section of Tab 6, the proposed method achieves consistent improvements over the recent columnar ViT methods [12, 19, 58, 75] while showing competitive results to the pyramidal counterparts.

5 Scope and Limitations

Despite the interpretability and effectiveness of PerViT, it still leaves much room for improvements. First, PerViT-Attention (Eq. 5) is based on the original self-attention formulation [63], thus directly inheriting its limitations [18, 58], *e.g.*, quadratic complexity w.r.t. input resolution. The computational efficiency could be further improved by approximating low-rank matrices as in [7, 37, 73]. Second, given the ability to process high-resolution input with feasible complexity, the efficacy of PerViT could be improved by adopting multi-resolution pyramidal design following recent trend of ViT designs [17, 35, 37, 40, 67, 69, 73, 74]. Third, the focus of this paper is model development & exploration for image classification task but we believe the proposed idea is broadly generalizable to other vision applications such as object detection and segmentation. We leave this to future work.

6 Conclusion

This paper explores blending human peripheral vision with machine vision for effective visual recognition, and introduces Peripheral Vision Transformer which learns to provide diverse position-based attentions to model peripheral vision using peripheral projections and initialization. We have systematically investigated the inner workings of the proposed network and observed that the network enjoys the benefits of both convolution and self-attention by learning to decide level of the locality and dynamicity for the feature transformations, by the network itself given training data. The consistent improvements over the baseline models on ImageNet-1K classification across different model sizes and in-depth ablation study confirm the efficacy of the proposed approach.

7 Acknowledgments and Disclosure of Funding

This work was supported by the IITP grants (No.2021-0-01696: High-Potential Individuals Global Training Program (40%), No.2022-0-00290: Visual Intelligence for Space-Time Understanding and Generation based on Multi-layered Visual Common Sense (50%), No.2019-0-01906: AI Graduate School Program - POSTECH (10%)) funded by Ministry of Science and ICT, Korea. This work was done while Juhong Min was working as an intern at Microsoft Research Asia.

References

- [1] J. Ba, J. Kiros, and G. Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 07 2016. 6
- [2] B. Balas, L. Nakano, and R. Rosenholtz. A summary-statistic representation in peripheral vision explains visual crowding. *Journal of vision*, 2009. 3
- [3] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko. End-to-end object detection with transformers. In *Proc. European Conference on Computer Vision (ECCV)*, 2020. 2
- [4] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin. Emerging properties in self-supervised vision transformers. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2021. 2
- [5] C.-F. R. Chen, R. Panda, and Q. Fan. RegionViT: Regional-to-Local Attention for Vision Transformers. In *International Conference on Learning Representations (ICLR)*, 2021. 2
- [6] F. Chollet. Xception: Deep learning with depthwise separable convolutions. *arXiv preprint arXiv:1610.02357*, 2016. 1
- [7] K. Choromanski, V. Likhoshesterov, D. Dohan, X. Song, A. Gane, T. Sarlos, P. Hawkins, J. Davis, A. Mohiuddin, L. Kaiser, D. Belanger, L. Colwell, and A. Weller. Rethinking attention with performers. In *International Conference on Learning Representations (ICLR)*, 2021. 2, 10
- [8] X. Chu, Z. Tian, B. Zhang, X. Wang, X. Wei, H. Xia, and C. Shen. Conditional positional encodings for vision transformers. *arXiv preprint arXiv:2102.10882*, 2021. 1, 2, 9
- [9] J.-B. Cordonnier, A. Loukas, and M. Jaggi. On the relationship between self-attention and convolutional layers. In *International Conference on Learning Representations (ICLR)*, 2020. 1, 2
- [10] Z. Dai, H. Liu, Q. V. Le, and M. Tan. Coatnet: Marrying convolution and attention for all data sizes. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. 1, 2, 8
- [11] Z. Dai, Z. Yang, Y. Yang, J. Carbonell, Q. V. Le, and R. Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. In *Proc. Association for Computational Linguistics (ACL)*, 2019. 2
- [12] S. d’Ascoli, H. Touvron, M. Leavitt, A. Morcos, G. Biroli, and L. Sagun. Convit: Improving vision transformers with soft convolutional inductive biases. In *Proc. International Conference on Machine Learning (ICML)*, 2021. 1, 2, 5, 8, 10
- [13] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009. 6, 10
- [14] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*, 2019. 2, 6
- [15] A. Deza and M. P. Eckstein. Can peripheral representations improve clutter metrics on complex scenes? *arXiv preprint arXiv:1608.04042*, 2016. 3
- [16] A. Deza and T. Konkle. Emergent properties of foveated perceptual systems. *arXiv preprint arXiv:2006.07991*, 2021. 3
- [17] X. Dong, J. Bao, D. Chen, W. Zhang, N. Yu, L. Yuan, D. Chen, and B. Guo. Cswin transformer: A general vision transformer backbone with cross-shaped windows. *arXiv preprint arXiv:2107.00652*, 2021. 6, 10
- [18] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations (ICLR)*, 2021. 1, 2, 3, 5, 6, 8, 10, 15
- [19] A. El-Nouby, H. Touvron, M. Caron, P. Bojanowski, M. Douze, A. Joulin, I. Laptev, N. Neverova, G. Synnaeve, J. Verbeek, et al. Xcit: Cross-covariance image transformers. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. 10, 15
- [20] L. Fridman, B. Jenik, S. Keshvari, B. Reimer, C. Zetsche, and R. Rosenholtz. Sideeye: A generative neural network based simulator of human peripheral vision. *arXiv preprint arXiv:1706.04568*, 2017. 3
- [21] K. Fukushima and S. Miyake. Neocognitron: A self-organizing neural network model for a mechanism of visual pattern recognition. *Competition and cooperation in neural nets*, 1982. 1

- [22] S. Gould, J. Arfvidsson, A. Kaehler, B. Sapp, M. Messner, G. Bradski, P. Baumstarck, S. Chung, and A. Y. Ng. Peripheral-foveal vision for real-time object recognition and tracking in video. In *Proc. International Joint Conferences on Artificial Intelligence Organization (IJCAI)*, 2007. 3
- [23] A. Harrington and A. Deza. Finding biological plausibility for adversarially robust features via metameric tasks. In *International Conference on Learning Representations (ICLR)*, 2022. 3
- [24] H. He, J. Liu, Z. Pan, J. Cai, J. Zhang, D. Tao, and B. Zhuang. Pruning self-attentions into convolutional layers in single path. *arXiv preprint arXiv:2111.11802*, 2022. 2
- [25] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 1, 8
- [26] D. Hendrycks and K. Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2020. 6
- [27] J. Hu, L. Shen, and G. Sun. Squeeze-and-excitation networks. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 1
- [28] Z. Huang, D. Liang, P. Xu, and B. Xiang. Improve transformer models with better relative position embeddings. In *Annual Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2020. 2
- [29] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015. 6
- [30] A. Jaegle, F. Gimeno, A. Brock, A. Zisserman, O. Vinyals, and J. Carreira. Perceiver: General perception with iterative attention. In *Proc. International Conference on Machine Learning (ICML)*, 2021. 2
- [31] A. Jonnalagadda, W. Y. Wang, B. S. Manjunath, and M. P. Eckstein. Foveater: Foveated transformer for image classification. *arXiv preprint arXiv:2105.14173*, 2022. 3
- [32] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2012. 1
- [33] A. M. Larson and L. C. Loschky. The contributions of central versus peripheral vision to scene gist recognition. *Journal of vision*, 2009. 3
- [34] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 1998. 1
- [35] Y. Lee, J. Kim, J. Willette, and S. J. Hwang. Mpvit: Multi-path vision transformer for dense prediction. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 1, 2, 6, 10
- [36] J. Lettvin. On seeing sidelong. *The Sciences*, 16, 07 1976. 2, 5
- [37] K. Li, Y. Wang, J. Zhang, P. Gao, G. Song, Y. Liu, H. Li, and Y. Qiao. Uniformer: Unifying convolution and self-attention for visual recognition. *arXiv preprint arXiv:2201.04676*, 2022. 1, 2, 10
- [38] H. Liu, X. Jiang, X. Li, Z. Bao, D. Jiang, and B. Ren. Nommer: Nominate synergistic context in vision transformer for visual recognition. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 1
- [39] J. Liu, H. Li, G. Song, X. Huang, and Y. Liu. Uninet: Unified architecture search with convolution, transformer, and mlp. *arXiv preprint arXiv:2110.04035*, 2021. 2
- [40] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2021. 1, 2, 5, 6, 10, 15
- [41] Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T. Darrell, and S. Xie. A convnet for the 2020s. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 1, 8
- [42] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations (ICLR)*, 2019. 6
- [43] C. I. Lou, D. Migotina, J. P. Rodrigues, J. Semedo, F. Wan, P. U. Mak, P. I. Mak, M. I. Vai, F. Melicio, J. G. Pereira, and A. Rosa. Object recognition test in peripheral vision: A study on the influence of object color, pattern and shape. In *Brain Informatics*, 2012. 3

- [44] H. Lukanov, P. König, and G. Pipa. Biologically inspired deep learning model for efficient foveal-peripheral vision. *Frontiers in Computational Neuroscience*, 15, 2021. 3
- [45] J. Min and M. Cho. Convolutional hough matching networks. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2940–2950, June 2021. 4
- [46] J. Min, D. Kang, and M. Cho. Hypercorrelation squeeze for few-shot segmentation. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2021. 4
- [47] V. Nair and G. E. Hinton. Rectified linear units improve restricted boltzmann machines. In *Proc. International Conference on Machine Learning (ICML)*, 2010. 6
- [48] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. 15
- [49] M. Raghu, T. Unterthiner, S. Kornblith, C. Zhang, and A. Dosovitskiy. Do vision transformers see like convolutional neural networks? In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. 1, 2, 5
- [50] P. Ramachandran, N. Parmar, A. Vaswani, I. Bello, A. Levskaya, and J. Shlens. Stand-alone self-attention in vision models. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. 1, 2
- [51] I. Rocco, M. Cimpoi, R. Arandjelović, A. Torii, T. Pajdla, and J. Sivic. Neighbourhood consensus networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018. 4
- [52] R. Rosenholtz. Capabilities and limitations of peripheral vision. *Annual review of vision science*, 2016. 3
- [53] R. Rosenholtz. Demystifying visual awareness: Peripheral encoding plus limited decision complexity resolve the paradox of rich visual experience and curious perceptual failures. *Attention, Perception, & Psychophysics*, 2020. 3
- [54] P. Shaw, J. Uszkoreit, and A. Vaswani. Self-attention with relative position representations. In *Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*, 2018. 2
- [55] A. Srinivas, T.-Y. Lin, N. Parmar, J. Shlens, P. Abbeel, and A. Vaswani. Bottleneck transformers for visual recognition. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [56] R. Strudel, R. Garcia, I. Laptev, and C. Schmid. Segmenter: Transformer for semantic segmentation. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2021. 2
- [57] C. Sun, A. Shrivastava, S. Singh, and A. Gupta. Revisiting unreasonable effectiveness of data in deep learning era. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2017. 1
- [58] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, and H. Jegou. Training data-efficient image transformers & distillation through attention. In *Proc. International Conference on Machine Learning (ICML)*, 2021. 2, 3, 5, 6, 8, 9, 10, 15
- [59] H. Touvron, M. Cord, A. El-Nouby, P. Bojanowski, A. Joulin, G. Synnaeve, J. Verbeek, and H. Jegou. Augmenting convolutional networks with attention-based aggregation. *arXiv preprint arXiv:2112.13692*, 2021. 1, 2
- [60] Z. Tu, H. Talebi, H. Zhang, F. Yang, P. Milanfar, A. Bovik, and Y. Li. MaxViT: Multi-Axis Vision Transformer. In *Proc. European Conference on Computer Vision (ECCV)*, 2022. 1, 2
- [61] D. Ulyanov, A. Vedaldi, and V. Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*, 2017. 4
- [62] A. Vaswani, P. Ramachandran, A. Srinivas, N. Parmar, B. Hechtman, and J. Shlens. Scaling local self-attention for parameter efficient visual backbones. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 1, 2
- [63] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017. 1, 2, 3, 10
- [64] M. R. Vuyyuru, A. Banburski, N. Pant, and T. Poggio. Biologically inspired mechanisms for adversarial robustness. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 3

- [65] H. Wang, Y. Zhu, B. Green, H. Adam, A. Yuille, and L.-C. Chen. Axial-deeplab: Stand-alone axial-attention for panoptic segmentation. In *Proc. European Conference on Computer Vision (ECCV)*, 2020. [2](#)
- [66] P. Wang and G. W. Cottrell. Central and peripheral vision for scene recognition: A neurocomputational modeling exploration. *arXiv preprint arXiv:1705.00816*, 2017. [3](#)
- [67] W. Wang, E. Xie, X. Li, D.-P. Fan, K. Song, D. Liang, T. Lu, P. Luo, and L. Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2021. [1](#), [2](#), [5](#), [10](#), [15](#)
- [68] M. W. Wijnjtjes and R. Rosenholtz. Context mitigates crowding: Peripheral object recognition in real-world images. *Cognition*, 2018. [3](#)
- [69] H. Wu, B. Xiao, N. Codella, M. Liu, X. Dai, L. Yuan, and L. Zhang. Cvt: Introducing convolutions to vision transformers. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2021. [1](#), [2](#), [5](#), [10](#)
- [70] K. Wu, H. Peng, M. Chen, J. Fu, and H. Chao. Rethinking and improving relative position encoding for vision transformer. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2021. [2](#), [4](#), [9](#)
- [71] T. Xiao, P. Dollar, M. Singh, E. Mintun, T. Darrell, and R. Girshick. Early convolutions help transformers see better. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. [1](#), [2](#), [5](#), [6](#)
- [72] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He. Aggregated residual transformations for deep neural networks. *arXiv preprint arXiv:1611.05431*, 2016. [1](#)
- [73] W. Xu, Y. Xu, T. Chang, and Z. Tu. Co-scale conv-attentional image transformers. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2021. [1](#), [2](#), [5](#), [6](#), [10](#), [15](#)
- [74] J. Yang, C. Li, P. Zhang, X. Dai, B. Xiao, L. Yuan, and J. Gao. Focal self-attention for local-global interactions in vision transformers. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. [10](#), [15](#)
- [75] L. Yuan, Y. Chen, T. Wang, W. Yu, Y. Shi, Z.-H. Jiang, F. E. Tay, J. Feng, and S. Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2021. [10](#)
- [76] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *Proc. European Conference on Computer Vision (ECCV)*, 2014. [5](#)
- [77] Y. Zhao, G. Wang, C. Tang, C. Luo, W. Zeng, and Z.-J. Zha. A battle of network structures: An empirical study of cnn, transformer, and mlp. *arXiv preprint arXiv:2108.13002*, 2021. [1](#)

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [\[Yes\]](#)
 - (b) Did you describe the limitations of your work? [\[Yes\]](#) See Section 5.
 - (c) Did you discuss any potential negative societal impacts of your work? [\[Yes\]](#) We describe potential negative societal impacts in the supplementary.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [\[Yes\]](#)
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [\[Yes\]](#) We describe the full set of assumptions of theoretical results in the supplementary.
 - (b) Did you include complete proofs of all theoretical results? [\[Yes\]](#) We include complete proofs of our theoretical results in the supplementary.
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [\[Yes\]](#) We include the code and instructions for reproduction in the supplementary. The training and validation data, e.g., ImageNet (ILSVRC2012), is available online.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [\[Yes\]](#) We specify all the training details in the supplementary.
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [\[No\]](#) Training large-scale image classification models is very expensive so we only performed one run. This is common practice in the field as seen in cited references [\[18, 19, 40, 58, 67, 73, 74\]](#).
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [\[Yes\]](#) We include the amount and type of resources used for training in the supplementary.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [\[Yes\]](#) As our baseline code, we use DeiT [\[58\]](#) which is implemented with PyTorch [\[48\]](#) framework, all of which are open-sourced. We cite all the papers that are relevant to the code, models, and datasets.
 - (b) Did you mention the license of the assets? [\[Yes\]](#) Our code refers to the licenses of the assets it relies on.
 - (c) Did you include any new assets either in the supplemental material or as a URL? [\[Yes\]](#) We include trained models in the supplementary for reproduction of our main results.
 - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [\[No\]](#) We are using only publicly available, benchmark datasets.
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [\[No\]](#) We are using only publicly available, benchmark datasets.
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [\[N/A\]](#) We did not crowdsource any information for this research.
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [\[N/A\]](#) We did not crowdsource any information for this research.
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [\[N/A\]](#) We did not crowdsource any information for this research.