## LOCALITY-ATTENDING VISION TRANSFORMER

**Anonymous authors** 

000

001 002 003

004

006

008 009

010

011

012

013

014

016

017

018

019

021

024

025

026

028 029

031

033

034

037

040

041

042

043

044

046 047

048

051

052

Paper under double-blind review

#### **ABSTRACT**

Vision transformers have demonstrated remarkable success in classification by leveraging global self-attention to capture long-range dependencies. However, this same mechanism can obscure fine-grained spatial details crucial for tasks such as segmentation. In this work, we seek to enhance the segmentation performance of vision transformers after being trained using the usual image-level classification objective. More specifically, we present a simple yet effective addon for vision transformers that improve their performance on segmentation tasks while retaining their image-level recognition capabilities. In our approach, we modulate the self-attention with a learnable Gaussian kernel that biases the attention toward neighboring patches. We further refine the patch representations to learn better embeddings at patch positions. These modifications ensure meaningful representations at spatial positions and encourage tokens to focus on local surroundings, while still preserving the model's ability to incorporate global information. Experiments demonstrate the effectiveness of our modifications, evidenced by substantial segmentation gains on three benchmarks (e.g., over 6% and 4% on ADE20K for ViT Tiny and Base), without changing the training regime or sacrificing classification performance. The code is available at https://anonymous.4open.science/r/LocAtViTRepo/.

## 1 Introduction

Vision transformers (ViT, Dosovitskiy et al., 2021) have emerged as powerful visual backbones by modeling images as sequences of patch tokens, processed with self-attention. Unlike convolutional neural networks (CNN, LeCun et al., 2015), which aggregate local information in a restricted receptive field, ViTs can capture long-range dependencies at any layer. This global attention mechanism has proven highly effective for image classification, enabling ViT models to surpass CNN performance when sufficient data is available (Touvron et al., 2021). A key factor behind this success is the ability to integrate global context that leads to more uniform and holistic representations across layers, which enhances the recognition of high-level image semantics (Raghu et al., 2021).

The same global focus that makes ViTs excel in classification, however, poses challenges for dense prediction tasks such as semantic segmentation. These tasks require precise localization and fine-grained spatial detail, properties that convolutional inductive biases naturally encourage but vanilla ViTs lack (Hassani et al., 2023). As a result, the design of spatial attention and feature hierarchy has been found critical for adapting transformers to dense tasks (Wang et al., 2021; Liu et al., 2021). Still, a tension remains between capturing global context and preserving local detail. Global attention can dilute local cues or incur high computation, whereas purely local schemes may miss long-range dependencies needed for holistic understanding. Besides, the classification training objective used by models often neglects the necessities of dense prediction, motivating a need for a segmentation-in-mind pretraining.

More recently, foundation models trained at large-scale (Radford et al., 2021; Oquab et al., 2023), which learn versatile visual representations, have seen broad adoption in a breadth of visual tasks. Despite the availability of more intricate designs, these models still mostly adopt vanilla ViT due to its simplicity and ease of integration. This widespread reliance underscores the practical value of enhancing ViT's capabilities rather than pursuing arguably more complex new designs. A prominent example is CLIP (Radford et al., 2021), which couples a ViT-based image encoder with a text encoder to align image-text representations, enabling zero-shot classification and open-vocabulary recognition. Such representations can be repurposed for dense predictions, for instance, by com-

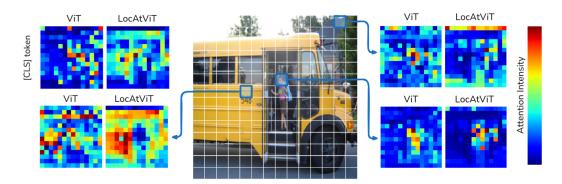


Figure 1: **Qualitative evaluation on the attention maps.** The final attention map of ViT and LocAtViT for the [CLS] token and three patches are illustrated for an image with label *school bus*.

paring local features to text prompts, but this adaptation is non-trivial. Furthermore, recent studies try to harness CLIP's knowledge for segmentation without any task-specific training (Zhou et al., 2022; Wang et al., 2024; Hajimiri et al., 2025). However, as CLIP and similar models are not trained for quality local representations, their features often lack the spatial granularity needed for precise dense prediction.

**Contributions.** In this paper, we propose a modular Locality-Attending (LocAt) add-on, which incorporates two ideas: (i) We modulate the attention logits with a learnable Gaussian kernel centered on each query token's location, ensuring that patches closer to the query receive higher attention. This acts as an explicit inductive bias encouraging each token to attend to its local neighborhood while still allowing global interactions. We denote the resulting self-attention module as the Gaussian-Augmented (GAug) attention (Sec. 4.1). (ii) We enhance patch representations for segmentation by introducing minor changes prior to the classification head, preserving the meaningfulness of spatial tokens, that are most important for dense prediction. We term this procedure as *Patch Represen*tation Refinement (PRR) that addresses the gradient flow issue in ViTs for segmentation, which is overlooked in the literature (see Sec. 4.2). Hence, LocAt refers to the combination of GAug and PRR. Figure 2 demonstrates that it improves different baselines, yielding significant segmentation performance gains (arrows pointing upward), while preserving or improv-

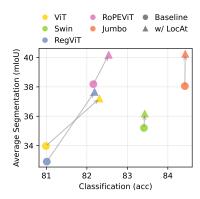


Figure 2: **LocAt considerably enhances different baselines** in segmentation, while preserving or even improving classification.

ing classification accuracy (no arrow pointing to the left). The proposed add-on also enhances the quality of attention maps, as illustrated in Fig. 1. LocAt is a lightweight and objective-agnostic add-on, also compatible with self-supervised pretraining. Importantly, the minimal architectural changes required to integrate LocAt make it readily applicable to any ViT with marginal changes, facilitating its usage in foundation models. To the best of our knowledge, we are the first to offer this perspective on ViT pretraining: designing pretraining with downstream dense prediction in mind, while being faithful to vanilla ViT's training regime and architecture.

#### 2 Related Work

Hierarchical ViT backbones for dense prediction. While the original ViT targets image classification and produces low-resolution features with weak locality priors (Dosovitskiy et al., 2021), dense prediction has motivated backbones that retain or recover spatial detail across stages. Some works use pyramid and token-merging designs to introduce multi-scale features and lightweight decoders for segmentation (Wang et al., 2021; Xie et al., 2021), while others build parallel branches for local and global processing (Chu et al., 2021). These works show that topology substantially

helps dense tasks. However, they typically require non-trivial architectural changes (new stages or merging blocks) and may rely on local window attention that limits full-image interaction.

Convolution-based hybrids. Another line injects convolutional priors either inside attention or in the feed-forward network to encourage local bias while keeping global modeling. Works use convolutional projections (Wu et al., 2021a), add gated positional self-attention to softly bias toward convolutional behavior (d'Ascoli et al., 2021), couple local convolutional features with global representations (Peng et al., 2021), or add convolutions in the feed-forward network (Li et al., 2021). These hybrid models add extra modules that require tuning, and they can reduce plug-and-play compatibility with off-the-shelf ViTs, as they often introduce branches or replace core components. Besides, convolution offers a spatially-shared kernel which is independent of patch information.

Locality mechanisms inside attention. Orthogonal to backbone design, many papers modify the attention pattern itself to introduce locality. Many of the works use fixed or structured windows (Liu et al., 2021; Dong et al., 2021; Yang et al., 2021). Other ideas include utilizing sliding or dilated neighborhoods to expand receptive fields efficiently (Hassani et al., 2023; Hassani & Shi, 2023), sampling content-relevant keys (Xia et al., 2023), selecting regions using dynamic sparse routing (Zhu et al., 2023), or using explicit global-local mixers to balance context with locality (Ding et al., 2022; Tu et al., 2022; Chen et al., 2022; Hatamizadeh et al., 2023). Most of these approaches restrict or mask interactions (using windows or patterns) or add mixing subsystems that complicate design, impeding their widespread adoption.

**Positional encodings that strengthen locality.** Beyond absolute embeddings, relative positional encoding (RPE), and rotary positional encodings (RoPE) improve spatial awareness in ViTs (Shaw et al., 2018; Liu et al., 2021; Wu et al., 2021b; Su et al., 2021; Heo et al., 2024). These approaches are orthogonal to attention locality, and we briefly mentioned them to emphasize that they encode locality as well. Our work complements rather than replaces them, as we show in the experiments.

**Improving patch tokens representation.** Recent work on *register tokens* augments ViTs with dedicated auxiliary tokens that absorb non-informative computation and yield smoother feature maps helpful for dense prediction (Darcet et al., 2024). Unlike this approach, we do not require auxiliary tokens, and we also address the issue of gradient flow to spatial patch outputs, overlooked in the prior work.

Foundation models for dense prediction. Large pre-trained foundation models, such as CLIP (Radford et al., 2021), demonstrate impressive zero-shot generalization on image-level recognition by leveraging ViT backbones. The preference for the standard ViT backbone can be attributed to its strong global attention, predictable scaling behavior with data and model size, and a uniform architecture that avoids the need for complex stage-wise tuning as the model grows (Zhai et al., 2022; Alabdulmohsin et al., 2023). However, despite excelling on image-level benchmarks, such models remain less effective for dense prediction because their representations are predominantly global and task-agnostic (Shao et al., 2024). As a result, additional adaptation or decoding layers are usually required to repurpose them for segmentation or detection (Li et al., 2022; Xu et al., 2023; Luo et al., 2023). While these adaptations yield improvements, they do not fully address the core issue: foundation-model ViTs—trained with classification objectives—tend to emphasize global semantics over local detail (Liang et al., 2023).

A ViT backbone that natively preserves both local detail and global context could enable foundation models to excel at dense prediction without extra adaptation layers or specialized fine-tuning. In this work, we take a step in that direction by refining the ViT backbone itself. Our approach aims to potentially bridge the gap between the powerful image-level understanding and the requirements of pixel-level prediction tasks.

## 3 PRELIMINARIES

Each ViT layer l takes a sequence of tokens  $\mathbf{x}^{(l-1)} \in \mathbb{R}^{(1+hw)\times C}$  as input, containing a [CLS] token and hw spatial patch tokens. Each token is a C-dimensional vector, and h and w denote the number of patches in each column and row.  $\mathbf{x}^{(0)}$  is the partitioned and flattened input after adding

the positional embeddings. At each layer l, the following operations are applied, where LN, attn, and MLP denote layer normalization, self-attention, and feed-forward network, respectively:

$$\mathbf{x}' = \mathbf{x}^{(l-1)} + \operatorname{attn}\left(\operatorname{LN}(\mathbf{x}^{(l-1)})\right),$$
 (1)

$$\mathbf{x}^{(l)} = \mathbf{x}' + \text{MLP}(\text{LN}(\mathbf{x}')). \tag{2}$$

Each self-attention module (attn) consists of two sets of weight matrices:  $\mathbf{W}^{qkv} \in \mathbb{R}^{C \times d \times 3}$  to compute d-dimensional query, key, and value matrices (i.e.,  $\mathbf{q}, \mathbf{k}, \mathbf{v} \in \mathbb{R}^{(1+hw) \times d}$ ) based on the input, and  $\mathbf{W}^o \in \mathbb{R}^{d \times C}$  for the final projection. After obtaining  $\mathbf{q}, \mathbf{k}$ , and  $\mathbf{v}$ , we calculate:

$$\mathbf{A} = \operatorname{softmax} \left( \mathbf{q} \mathbf{k}^{\top} / \sqrt{d} \right) \mathbf{v}. \tag{3}$$

Matrix  $\mathbf{A} \in \mathbb{R}^{(1+hw)\times d}$  is then transformed by  $\mathbf{W}^o$  to form the output of the layer. The *attention logits* of a patch p are represented by the  $p^{\text{th}}$  row of  $\mathbf{q}\mathbf{k}^{\top}/\sqrt{d}$ . Note that for simplicity, we present the formulation of a single-head self-attention.

## 4 METHOD

We now present **LocAtViT**, which enhances ViT with two modular components, GAug attention (Sec. 4.1) and PRR (Sec. 4.2), and is trained with the same classification objective as ViT.

#### 4.1 Gaussian-Augmented attention

We aim to introduce explicit local attention into layers of vision transformer (for all tokens but <code>[CLS]</code>) via adding a patch-specific Gaussian kernel to attention logits. We first discuss the altered self-attention formulation, followed by details on computation of the kernel, and then the final form of the attention.

**Modified self-attention.** At every layer's self-attention, we add a *supplement* matrix S to the attention logits, aiming to emphasize the attention of each patch to its surrounding. With this addition, the self-attention formulation of Eq. (3) is modified as follows, which is also depicted in Fig. 3a:

$$\mathbf{A} = \operatorname{softmax} \left( \frac{\mathbf{q} \mathbf{k}^{\top}}{\sqrt{d}} + \mathbf{S} \right) \mathbf{v}. \tag{4}$$

We construct S so that a patch p attends more to its immediate surroundings, with increment gradually decreasing with distance from p. A natural choice is an unnormalized Gaussian centered at p. Gaussian kernels are infinitely differentiable and they smoothly decay as spatial distance increases, aligning with human visual perception mechanisms, where sensitivity decreases smoothly from local to distant regions. We parameterize variance of the Gaussian kernel for each patch by a 2D vector, stored in the  $p^{th}$  row of  $\Sigma \in \mathbb{R}^{hw \times 2}_+$ , controlling the attention span along both axes for each patch. Since patches might have different needs in how far they should attend to their neighbors, we compute the variances based on the query matrix derived from the input, using a learnable weight matrix  $\mathbf{W}^{\sigma} \in \mathbb{R}^{d \times 2}$  (with f being a scaled sigmoid function ensuring positive, bounded values):

$$\Sigma = f(\mathbf{q}\mathbf{W}^{\sigma}),\tag{5}$$

**Gaussian kernel.** For a patch grid of size  $h \times w$ , we denote the set of coordinate vectors:

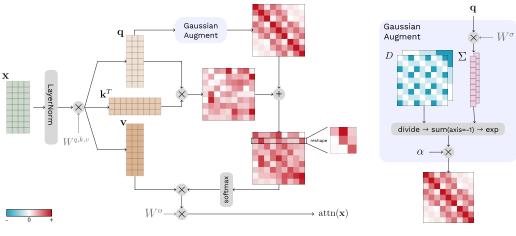
$$\mathbf{P} = \begin{bmatrix} i & j \end{bmatrix}_{i \in \{1, 2, \dots, h\}, \ j \in \{1, 2, \dots, w\}}, \tag{6}$$

in an  $hw \times 2$  matrix. The  $hw \times hw \times 2$  pairwise squared difference **D** is computed as:

$$\mathbf{D}_{ptm} = \left(\mathbf{P}_{pm} - \mathbf{P}_{tm}\right)^2, \quad \text{for } m \in \{1, 2\},\tag{7}$$

where p and t denote indices of source and target patch, and m indexes the coordinate dimensions. Given  $\Sigma$ , the elements in the Gaussian kernel matrix  $\mathbf{G} \in \mathbb{R}_+^{(hw+1)\times (hw+1)}$  are calculated as:

$$\mathbf{G}_{pt} = \exp\left(-\frac{1}{2}\sum_{m=1}^{2} \frac{\mathbf{D}_{ptm}}{\mathbf{\Sigma}_{pm}}\right),\tag{8}$$



(a) Modified self-attention.

(b) Supplement matrix S.

Figure 3: **Illustration of the Gaussian-Augmented attention** for a  $3 \times 3$  grid. For simplicity, the <code>[CLS]</code> token is not shown. (a) The Gaussian addition, i.e., **S** in Eq. (4), is obtained based on **q** and is added to the attention logits. The p-th row in the attention logits matrix presents the attention of patch p to all patch tokens. The reshaped matrix illustrates that with the GAug add-on, both local and global attentions are integrated. (b) The supplement matrix **S** encourages attending to the locality and is computed using the pairwise squared difference tensor **D** from Eq. (7).

which determines the addition to attention logits from patch p to t. To keep the global attention intact, we set the first row/column of G to zero, as they correspond to attention from/to the [CLS] token. By pre-computing D, *i.e.*, the numerator, we can efficiently compute G during training.

**Supplement matrix.** Based on Eq. (8), each entry in G lies in [0,1]. Directly setting S = G in Eq. (4) causes a scale mismatch between S and the attention logits. To mitigate this discrepancy, we assume a learnable weight matrix  $W^{\alpha} \in \mathbb{R}^{d \times 1}$  that computes the desired scaling for each patch, based on its  $\mathbf{q}$  vector. Entries in  $\boldsymbol{\alpha}$  scale rows of the Gaussian kernel, more specifically:

$$\alpha = \text{softplus}(\mathbf{q}\mathbf{W}^{\alpha}) \in \mathbb{R}^{hw}_{+},$$
 (9)

$$\mathbf{S} = \operatorname{diag}(\boldsymbol{\alpha}) \, \mathbf{G},\tag{10}$$

in which softplus ensures positive coefficients. We refer to our modified self-attention as *Gaussian-Augmented* (GAug) attention. Figure 3b illustrates the generation process of the supplement matrix.

## 4.2 PATCH REPRESENTATION REFINEMENT

**Problem statement.** In a classification task using ViT, only the [CLS] token's output of the model is used for computing the loss. While effective for classification, this approach has fundamental limitations for dense prediction from a gradient flow perspective. More concretely, the patch positions' outputs receive no *direct* supervision, *i.e.*, it is not important to the model what ViT's final outputs are at those positions. However, these output representations are crucial for further dense prediction. This is problematic because the fine-grained spatial information carried by individual patch tokens is not effectively learned at the final layer.

Some subsequent methods, such as Swin (Liu et al., 2021), remove the <code>[CLS]</code> token and use global average pooling (GAP) before the classification head. However, this forces an undesirable behavior from a dense prediction standpoint, *i.e.*, a *uniform gradient flow* across all positions. For example, in an image of a bird with other objects in the background, GAP compels the model to match all patch representations—including background regions—with the classifier's prototype of bird. The uniform gradient flow means that all patch tokens receive equal importance, regardless of their relevance, potentially leading to representations particularly suboptimal for tasks like segmentation. Moreover, GAP has been shown to reduce localization in higher layers (Raghu et al., 2021).

**Proposed solution.** To encourage meaningful patch representations at the final layer's output,  $\mathbf{x}^{(l)}$ , we propose the following operation before the classification head:

$$\mathbf{x}^{+} = \operatorname{softmax}\left(\frac{\mathbf{x}^{(l)}\mathbf{x}^{(l)\top}}{\sqrt{d}}\right)\mathbf{x}^{(l)},\tag{11}$$

which acts like a *parameter-free* self-attention. This operation, which introduces no new parameter, aggregates information from all patch positions in a non-uniform manner, thereby preserving their unique contributions and ensuring diverse gradient flow across patch locations. The resulting representation at the [CLS] token,  $\mathbf{x}_0^+$ , is then passed to the classification head. We refer to this strategy as *Patch Representation Refinement* (PRR). PRR can be seen as an alternative to GAP, suitable for segmentation-in-mind pretraining.

Our components share a common objective: enhancing representation quality for dense prediction. GAug contributes by injecting a localized bias into attention, prompting the model to capture richer local information. Complementarily, PRR contributes by ensuring that positional patch outputs—crucial for segmentation—receive gradients, enabling the model to effectively learn and refine those representations.

## 5 EXPERIMENTS

## 5.1 EXPERIMENTAL SETUP

**Datasets.** For the main experiments, where we assess both classification and segmentation performance, we first train models on ImageNet-1K (Deng et al., 2009; Russakovsky et al., 2015), which contains 1.28M training images from 1,000 classes. Then, we further utilize these models for training on segmentation datasets: ADE20K (Zhou et al., 2019), PASCAL Context (Mottaghi et al., 2014), and COCO Stuff (Caesar et al., 2018; Lin et al., 2014), which contain 150, 59, and 171 semantic categories, respectively. ADE20K and COCO Stuff images are resized to  $512 \times 512$  and PASCAL Context images to  $480 \times 480$ . Furthermore, we also assess classification performance on smaller scale datasets: CIFAR-100 (Krizhevsky & Hinton, 2009) and mini-ImageNet (Vinyals et al., 2016), a subset of ImageNet-1K, consisting of 100 classes with 500 training and 100 validation examples each. In all classification experiments, images are resized to  $224 \times 224$ .

Implementation details. Our method is implemented using the PyTorch Image Models (timm) (Wightman, 2019) library. We train models on ImageNet-1K for 300 epochs, with initial learning rate (LR) 0.001, and on CIFAR-100 and mini-ImageNet for 600 epochs, with LR 0.0005. Global batch size is set to 1024, linear warm-up to 20 epochs, and we use AdamW (Kingma & Ba, 2014; Loshchilov & Hutter, 2019) optimizer with a weight decay of 0.05. As in Ding et al. (2022), a simple triangular learning rate scheduler (Smith & Topin, 2018) is applied, and the stochastic depth drop rates (Huang et al., 2016) for the Tiny, Small, and Base backbones are set to 0.1, 0.2, and 0.4, respectively. We follow Liu et al. (2021) for data augmentation and use RandAugment (Cubuk et al., 2020), Mixup (Zhang et al., 2018), Cutmix (Yun et al., 2019), and random erasing (Zhong et al., 2020). The sigmoid function f in Eq. (5) is scaled to have a maximum of  $\max(h, w)$ , and shifted to satisfy f(0) = 1.

For semantic segmentation, we utilize the MMSegmentation toolbox (OpenMMLab, 2020) and employ a simple 1-layer MLP on top of the frozen classification-trained models. This configuration ensures that segmentation performance mainly reflects the discriminative power of the classification-trained backbone in dense prediction. This setup aligns with our goal of isolating and assessing patch representation quality under a low-tuning regime. Training on segmentation datasets is performed over 20K iterations with a batch size of 32.

#### 5.2 Main results

**Segmentation performance.** The LocAt add-on can be applied on several ViT-based models, and Tab. 1 evaluates its effect, in terms of classification performance on ImageNet-1K, as well as segmentation performance on three benchmarks, when applied to five models: ViT (Dosovitskiy et al., 2021), Swin Transformer (Liu et al., 2021), ViTs with registers (denoted as RegViT, we use 4 registers, Darcet et al., 2024), Rotary Position Embedding for ViTs (denoted as RoPEViT, Heo et al.,

Table 1: **Segmentation performance** of models and their counterparts with our LocAt extension (in gray), along with their **classification performance** on ImageNet-1K, which the models are initially trained on. Results demonstrate that (i) LocAt substantially boosts segmentation performance (our primary focus), while preserving or even improving the classification performance, and (ii) this effect holds for a variety of methods, for different backbone sizes. Furthermore, (iii) the segmentation gains appear not only in weaker baselines, but also in strong, high-performing models, where classification improvements are harder to achieve.

	Method	Segm ADE	entation mIol P-Context	U (%) C-Stuff	Top-1 (%) ImageNet	#Params (M)	FLOPs (G)
	ViT	17.30	33.71	20.29	72.39	6	1.26
	+ LocAt	$23.47_{+6.17}$	$38.57_{\pm 4.86}$	$26.15_{+5.86}$	$73.94_{+1.55}$	6	1.27
	Swin	25.58	36.78	28.34	81.18	28	4.50
	+ LocAt	$26.52_{\pm 0.94}$	$37.65_{\pm 0.87}$	$29.09_{\pm 0.75}$	$81.43_{\pm 0.25}$	28	4.51
Tiny	RegViT	15.98	33.45	19.58	72.90	6	1.29
Τ	+ LocAt	$24.39_{+8.41}$	$39.90_{+6.45}$	$27.38_{+7.80}$	$74.08_{\pm 1.18}$	6	1.30
	RoPEViT	19.17	38.16	22.75	73.60	6	1.26
	+ LocAt	$24.48_{+5.31}$	$40.79_{+2.63}$	$27.98_{+5.23}$	$74.34_{\pm 0.74}$	6	1.27
	Jumbo	20.33	36.36	22.13	78.71	9	1.40
	+ LocAt	$21.62_{+1.29}$	$37.22_{\pm 0.86}$	$23.87_{\pm 1.74}$	$78.78_{\pm0.07}$	9	1.42
	ViT	28.40	43.10	30.43	80.99	86	17.58
	+ LocAt	$32.64_{+4.24}$	$45.35_{+2.25}$	$33.62_{+3.19}$	$82.31_{+1.32}$	86	17.64
	Swin	31.90	40.11	33.60	83.41	88	15.46
	+ LocAt	$32.89_{\pm 0.99}$	$41.44_{+1.33}$	$34.20_{\pm 0.60}$	$83.43_{\pm 0.02}$	88	15.47
Base	RegViT	27.93	41.81	28.99	81.01	86	17.95
	+ LocAt	$32.71_{\pm 4.78}$	$46.14_{+4.33}$	$34.12_{+5.13}$	$82.19_{\pm 1.18}$	86	18.02
	RoPEViT	31.38	48.83	34.35	82.16	86	17.58
	+ LocAt	$34.94_{+3.56}$	$49.24_{\pm0.41}$	$36.37_{+2.02}$	$82.54_{\pm0.38}$	86	17.64
	Jumbo	32.20	47.31	34.65	84.42	130	19.74
	+ LocAt	$35.69_{+3.49}$	$49.20_{+1.89}$	$35.84_{+1.19}$	$84.43_{\pm0.01}$	130	19.81

2024), and the recent Jumbo (Fuller et al., 2025). Comparing each baseline with its enhanced counterpart (gray row below), proves LocAt's addition is useful in improving the segmentation performance of all. For instance, LocAtViT Tiny achieves a substantial improvement of +6.17%, +4.86%, and +5.86%, over ViT on ADE20K, PASCAL Context, and COCO Stuff, respectively. Importantly, LocAt-enhanced models' superior segmentation performance is achieved without compromising classification performance; in fact, they deliver comparable or even improved accuracy across different models (e.g., LocAtViT outperforms ViT by +1.55% in Tiny backbone).

LocAt improves baselines that are architecturally close to ViT significantly, e.g., RoPEViT, and interestingly, it brings improvements over Swin as well. We believe this is not trivial as the add-on was designed for ViT's architecture, in which there exists a <code>[CLS]</code> token and the attention width is not limited, while in Swin the windowed attention mechanism severely affects the extent to which LocAt can play a role. Furthermore, our add-on incurs a negligible increase in computational efficiency in terms of number of FLOPs over the corresponding counterparts (measured at  $224 \times 224$  using Sovrasov, 2018-2024). Additional experiments are presented in Appendix B.

Classification performance. In addition to the ImageNet-1K classification results in Tab. 1, Tab. 2 investigates LocAt's classification effectiveness on small-scale datasets: mini-ImageNet (Vinyals et al., 2016) and CIFAR-100 (Krizhevsky & Hinton, 2009). Although designed to enhance segmentation, these results demonstrate LocAt's classification effectiveness even when trained on small-scale datasets. LocAt improves ViT's performance by 3-6% on mini-ImageNet and 4-7% on CIFAR-100, while only introducing 2, 340 new parameters (0.003% increase for Base). Please note that segmentation results are not included for models trained on these datasets since, due to their scale and number of classes, representations are not expected to generalize well to segmentation benchmarks.

**Foundation models.** In the previous sections, we described our interest in improving ViT's segmentation capabilities without changing their training scheme. Our experiments support that our

Table 2: Classification top-1 accuracy of ViT and LocAtViT for different backbone sizes, on mini-ImageNet and CIFAR-100, showcasing LocAt's effectiveness on small-scale datasets.

C:	mini	-ImageNet	CIFAR-100		
Size	ViT	LocAtViT	ViT	LocAtViT	
Tiny	74.94	78.47 <sub>+3.53</sub>	73.84	80.43+6.59	
Small	78.98	84.30+5.32	76.33	$81.13_{\pm 4.80}$	
Base	79.91	$84.86_{+4.95}$	76.90	$82.20_{+5.30}$	

Table 3: **Self-supervised performance of Lo-cAtViT used in DINO**, showcasing LocAt's effectiveness in the self-supervised regime.

Expe	riment	ViT-S/16	LocAtViT-S/16
Linear classification		65.52	$67.65_{+2.13}$
Nearest neighbor	10-NN 20-NN 100-NN 200-NN	61.69 61.53 59.30 57.90	63.96 <sub>+2.27</sub> 63.74 <sub>+2.21</sub> 61.19 <sub>+1.89</sub> 59.78 <sub>+1.88</sub>

minor modifications lead to better dense prediction performance, while performing on par or superior to the vanilla models in classification. One reason for our interest in the mentioned problem is that ViTs have been widely used across computer vision foundation models and are the go-to choice for many of the recent methods (Radford et al., 2021; Kirillov et al., 2023; Caron et al., 2021; Oquab et al., 2023). One of the popular models that yields versatile image representations and transfers well to different computer vision tasks is DINO (Caron et al., 2021), which is trained in a self-supervised manner and can serve as a general-purpose backbone. Two of the main evaluation protocols used by Caron et al. (2021) are learning a linear classifier on top of the frozen backbone and nearest neighbor classification (*k*-NN) on top of the features.

We train DINO ViT-S/16 and DINO LocAtViT-S/16 on ImageNet-1K for 50 epochs using the setting provided in the official repository, and evaluate them on the mentioned tasks. Table 3 demonstrates that replacing ViT with LocAtViT in DINO actually improves its performance on both linear and k-NN classification. We report the k-NN performance on  $k \in \{10, 20, 100, 200\}$  as advised by Caron et al. (2021). These findings reveal our objective-agnostic modifications' effectiveness in the self-supervised regime and the potential of our method on backbones that learn general-purpose representations. While interesting, further investigation on larger foundation models is beyond our computational reach and lies outside the scope of this work.

## 5.3 QUALITATIVE ANALYSIS

An interesting implication of our proposed modifications is the refinement of ViT's patch outputs (through PRR), which makes it more suitable for use cases on dense prediction tasks. Figure 1 offers a visual comparison of attention maps from a vanilla ViT and our LocAtViT, both trained for classification, for an image labeled as *school bus*. From the [CLS] token's attention, we observe that ViT's focus is broadly dispersed, whereas LocAtViT shows more concentrated and coherent activation on key features of the bus. Furthermore, we present the attention maps of three patch tokens to other patches. For instance, a patch on the bus side attends to nearly the entire bus in LocAtViT, whereas ViT's map is harder to interpret. A patch covering the child's face generates meaningful attention in both models, but ViT seems to highlight unrelated regions more. Interestingly, for a patch near the top-right corner, LocAtViT not only focuses on some tree patches, but also extends attention to the sky and road, all corresponding to the image background. Despite being trained solely for classification, LocAtViT exhibits an improved ability to detect some scene structures, suggesting that our proposed local interactions can enrich the model's contextual understanding without sacrificing global attention. Further qualitative examples are presented in Appendix C.

#### 5.4 ABLATION STUDY

In this section, we provide an ablation study on the architectural choices we made. We also provide ablation study on the self-attention module's design in the Appendix D.

**Effect of GAug and PRR.** Part **1** of Tab. 4 ablates on GAug and PRR defined in Secs. 4.1 and 4.2. Results demonstrate that both GAug and PRR indeed enhance the performance of the model in both classification and segmentation, and their combination pushes the performance even further.

**Effect of positional embeddings.** Part ② of Tab. 4 evaluates the impact of the default absolute positional embeddings (PE) on our proposed LocAt add-on. For both backbone sizes, LocAtViT

Table 4: **Ablations on model's architecture.** We report segmentation performance (mIoU %) over three benchmarks and classification accuracy (top-1 %) on ImageNet-1K. PE and GAP stand for positional embeddings and global average pooling.

		Tiny				Base			
	Method		P-Context	C-Stuff	ImageNet	ADE	P-Context	C-Stuff	ImageNet
	ViT	17.30	33.71	20.29	72.39	28.40	43.10	30.43	80.99
	ViT + GAug	18.98	34.97	21.51	73.16	30.26	44.36	32.21	82.00
0	ViT + PRR	21.60	37.93	25.85	73.71	29.89	44.03	32.16	82.19
	LocAtViT	23.47	38.57	26.15	73.94	32.64	45.35	33.62	82.31
2	ViT - PE	15.13	31.94	19.35	69.36	24.59	40.18	28.79	79.39
•	LocAtViT - PE	22.69	38.15	26.05	73.10	29.73	44.69	32.17	82.17
	ViT	17.30	33.71	20.29	72.39	28.40	43.10	30.43	80.99
3	ViT + GAP	19.65	34.94	22.86	72.50	27.99	41.97	29.88	81.84
	ViT + PRR	21.60	37.93	25.85	73.71	29.89	44.03	32.16	82.19

without PE not only outperforms ViT without PE, but also surpasses ViT with PE. This indicates that LocAt captures the spatial information embedded into PE and more, with much fewer learnable parameters. It is worth noting that our approach is not an alternative to positional encoding and we did not intend to propose a new PE method. Therefore, these results are included just to demonstrate empirically that LocAt indeed captures the spatial information that the default PE captures, which is the agent for capturing locality in vanilla ViT. We have shown in Tab. 1 that LocAt is applicable alongside other, newer positional encoding approaches, such as RoPE, as well.

Comparison between PRR and GAP. As discussed in Sec. 4.2, PRR addresses patch locations' gradient flow issues while overcoming GAP's limitations in segmentation. Part 3 of Tab. 4 compares how vanilla ViT performs when equipped with PRR versus GAP. Results show PRR's superior segmentation performance and interestingly, it improves classification accuracy more than GAP. Moreover, although GAP helps ViT in classification, it hurts the segmentation performance in the Base backbone, which is in line with the discussions in Sec. 4.2 about GAP's problems in segmentation.

## 6 Conclusion

**Summary.** We present the *Locality-Attending Vision Transformer*, a modular framework designed to enhance vision transformers for dense prediction, while preserving their image-level capabilities. Our work introduces a new perspective on segmentation-in-mind pretraining of ViTs. By introducing the *GAug* attention, our method biases self-attention toward local regions, enabling ViTs to capture fine-grained spatial details, while *PRR* ensures meaningful gradient flow to patch tokens, thereby strengthening representation quality and improving their suitability for dense prediction. Extensive experiments show that LocAt delivers superior segmentation performance without compromising classification accuracy. Its modular design integrates seamlessly into existing ViTs, proving effective across different baselines. Importantly, our objective in this paper has not been to surpass state-of-the-art performance, but rather to improve classification-trained ViT backbones for segmentation, due to the trend of them being widely used, *e.g.*, by foundation models, with a method orthogonal to many prior advancements. Consistent with Heo et al. (2024), we therefore emphasized comparisons between baselines and their LocAt-enhanced counterparts.

While more advanced, specially designed architectures with great performance exist, we intentionally focused on improving ViT's architecture, due to their wide use cases and to facilitate further usage in foundation models. We hope that with these lightweight modifications, our work will be adopted in ViT-based foundation models.

**Limitations.** We evaluated our method on multiple classification and segmentation datasets. However, these datasets all only contain natural images, and we have left evaluation on other domains such as medical imaging or remote sensing as future work. Furthermore, while we have demonstrated the effectiveness of LocAtViT used in a small foundation model, evaluation on large foundation models, such as CLIP, has been out of our computational reach.

## REFERENCES

- Ibrahim M Alabdulmohsin, Xiaohua Zhai, Alexander Kolesnikov, and Lucas Beyer. Getting vit in shape: Scaling laws for compute-optimal model design. *Advances in Neural Information Processing Systems*, 36:16406–16425, 2023.
- Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Coco-stuff: Thing and stuff classes in context. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018.
- Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2021.
- Chun-Fu (Richard) Chen, Rameswar Panda, and Quanfu Fan. RegionViT: Regional-to-Local Attention for Vision Transformers. In *International Conference on Learning Representations*, 2022.
- Xiangxiang Chu, Zhi Tian, Yuqing Wang, Bo Zhang, Haibing Ren, Xiaolin Wei, Huaxia Xia, and Chunhua Shen. Twins: Revisiting the design of spatial attention in vision transformers. *Advances in neural information processing systems*, 34:9355–9366, 2021.
- Xiangxiang Chu, Zhi Tian, Bo Zhang, Xinlong Wang, and Chunhua Shen. Conditional positional encodings for vision transformers. In *International Conference on Learning Representations*, 2023.
- Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pp. 702–703, 2020.
- Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision transformers need registers. In *International Conference on Learning Representations*, 2024.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Mingyu Ding, Bin Xiao, Noel Codella, Ping Luo, Jingdong Wang, and Lu Yuan. Davit: Dual attention vision transformers. In *European conference on computer vision*, pp. 74–92. Springer, 2022.
- Xiaoyi Dong, Jianmin Bao, Dongdong Chen, Weiming Zhang, Nenghai Yu, Lu Yuan, Dong Chen, and Baining Guo. Cswin transformer: A general vision transformer backbone with cross-shaped windows, 2021.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021.
- Stéphane d'Ascoli, Hugo Touvron, Matthew L Leavitt, Ari S Morcos, Giulio Biroli, and Levent Sagun. Convit: Improving vision transformers with soft convolutional inductive biases. In *International conference on machine learning*, pp. 2286–2296. PMLR, 2021.
- Anthony Fuller, Yousef Yassin, Daniel G Kyrollos, Evan Shelhamer, and James R Green. Simpler fast vision transformers with a jumbo cls token. 2025.
- Sina Hajimiri, Ismail Ben Ayed, and Jose Dolz. Pay attention to your neighbours: Training-free open-vocabulary semantic segmentation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2025.
- Ali Hassani and Humphrey Shi. Dilated neighborhood attention transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023.
- Ali Hassani, Steven Walton, Jiachen Li, Shen Li, and Humphrey Shi. Neighborhood attention transformer. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 6185–6194, 2023.

- Ali Hatamizadeh, Hongxu Yin, Greg Heinrich, Jan Kautz, and Pavlo Molchanov. Global context vision transformers. In *Proceedings of the 40th International Conference on Machine Learning*, pp. 12633–12646, 2023.
  - Byeongho Heo, Song Park, Dongyoon Han, and Sangdoo Yun. Rotary position embedding for vision transformer. In *European Conference on Computer Vision*, pp. 289–305. Springer, 2024.
  - Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q Weinberger. Deep networks with stochastic depth. In *European conference on computer vision*, pp. 646–661. Springer, 2016.
  - Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint* arXiv:1412.6980, 2014.
  - Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything. *arXiv:2304.02643*, 2023.
  - Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. 2009.
  - Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444, 2015.
  - Boyi Li, Kilian Q Weinberger, Serge Belongie, Vladlen Koltun, and René Ranftl. Language-driven semantic segmentation. *arXiv preprint arXiv:2201.03546*, 2022.
  - Yawei Li, Kai Zhang, Jiezhang Cao, Radu Timofte, and Luc Van Gool. Localvit: Bringing locality to vision transformers. *arXiv preprint arXiv:2104.05707*, 2021.
  - Feng Liang, Bichen Wu, Xiaoliang Dai, Kunpeng Li, Yinan Zhao, Hang Zhang, Peizhao Zhang, Peter Vajda, and Diana Marculescu. Open-vocabulary semantic segmentation with mask-adapted clip. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7061–7070, 2023.
  - Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, pp. 740–755. Springer, 2014.
  - Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 10012–10022, 2021.
  - Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019.
  - Huaishao Luo, Junwei Bao, Youzheng Wu, Xiaodong He, and Tianrui Li. SegCLIP: Patch aggregation with learnable centers for open-vocabulary semantic segmentation. *ICML*, 2023.
  - Roozbeh Mottaghi, Xianjie Chen, Xiaobai Liu, Nam-Gyu Cho, Seong-Whan Lee, Sanja Fidler, Raquel Urtasun, and Alan Yuille. The role of context for object detection and semantic segmentation in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2014.
  - OpenMMLab. MMSegmentation: Openmmlab semantic segmentation toolbox and benchmark. https://github.com/open-mmlab/mmsegmentation, 2020.
  - Maxime Oquab, Timothée Darcet, Theo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Russell Howes, Po-Yao Huang, Hu Xu, Vasu Sharma, Shang-Wen Li, Wojciech Galuba, Mike Rabbat, Mido Assran, Nicolas Ballas, Gabriel Synnaeve, Ishan Misra, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision, 2023.
  - Zhiliang Peng, Wei Huang, Shanzhi Gu, Lingxi Xie, Yaowei Wang, Jianbin Jiao, and Qixiang Ye. Conformer: Local features coupling global representations for visual recognition. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 367–376, 2021.

- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pp. 8748–8763, 2021.
  - Maithra Raghu, Thomas Unterthiner, Simon Kornblith, Chiyuan Zhang, and Alexey Dosovitskiy. Do vision transformers see like convolutional neural networks? *Advances in neural information processing systems*, 34:12116–12128, 2021.
  - Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y.
  - Tong Shao, Zhuotao Tian, Hang Zhao, and Jingyong Su. Explore the potential of clip for training-free open vocabulary semantic segmentation. In *European Conference on Computer Vision*. Springer, 2024.
  - Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position representations. In Marilyn Walker, Heng Ji, and Amanda Stent (eds.), *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pp. 464–468, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-2074.
  - Leslie N. Smith and Nicholay Topin. Super-convergence: Very fast training of residual networks using large learning rates, 2018.
  - Vladislav Sovrasov. ptflops: a flops counting tool for neural networks in pytorch framework. https://github.com/sovrasov/flops-counter.pytorch, 2018-2024.
  - Jianlin Su, Yu Lu, Shengfeng Pan, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding, 2021.
  - Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In International conference on machine learning, pp. 10347–10357. PMLR, 2021.
  - Zhengzhong Tu, Hossein Talebi, Han Zhang, Feng Yang, Peyman Milanfar, Alan Bovik, and Yinxiao Li. Maxvit: Multi-axis vision transformer. In *European conference on computer vision*, pp. 459–479. Springer, 2022.
  - Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. *Advances in neural information processing systems*, 29, 2016.
  - Feng Wang, Jieru Mei, and Alan Yuille. Sclip: Rethinking self-attention for dense vision-language inference. *European Conference on Computer Vision*, 2024.
  - Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 568–578, 2021.
  - Ross Wightman. Pytorch image models. https://github.com/rwightman/pytorch-image-models, 2019.
  - Haiping Wu, Bin Xiao, Noel Codella, Mengchen Liu, Xiyang Dai, Lu Yuan, and Lei Zhang. Cvt: Introducing convolutions to vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 22–31, 2021a.
  - Kan Wu, Houwen Peng, Minghao Chen, Jianlong Fu, and Hongyang Chao. Rethinking and improving relative position encoding for vision transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 10033–10041, 2021b.

- Zhuofan Xia, Xuran Pan, Shiji Song, Li Erran Li, and Gao Huang. Dat++: Spatially dynamic vision transformer with deformable attention. *arXiv preprint arXiv:2309.01430*, 2023.
  - Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in neural information processing systems*, 34:12077–12090, 2021.
  - Mengde Xu, Zheng Zhang, Fangyun Wei, Han Hu, and Xiang Bai. San: Side adapter network for open-vocabulary semantic segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
  - Jianwei Yang, Chunyuan Li, Pengchuan Zhang, Xiyang Dai, Bin Xiao, Lu Yuan, and Jianfeng Gao. Focal self-attention for local-global interactions in vision transformers, 2021.
  - Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 6023–6032, 2019.
  - Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 12104–12113, 2022.
  - Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *International Conference on Learning Representations*, 2018.
  - Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random erasing data augmentation. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pp. 13001–13008, 2020.
  - Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Semantic understanding of scenes through the ade20k dataset. *International Journal of Computer Vision*, 2019.
  - Chong Zhou, Chen Change Loy, and Bo Dai. Extract free dense labels from clip. In *European Conference on Computer Vision*, 2022.
  - Lei Zhu, Xinjiang Wang, Zhanghan Ke, Wayne Zhang, and Rynson WH Lau. Biformer: Vision transformer with bi-level routing attention. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10323–10333, 2023.

# LOCALITY-ATTENDING VISION TRANSFORMER APPENDIX

## A TECHNICAL DETAILS

## A.1 CODE

Our code is anonymously and publicly available at https://anonymous.4open.science/r/LocAtViTRepo/. The README.md file provides guidelines on how to set up the environment, train the models, and perform different evaluations. For ViT (Dosovitskiy et al., 2021), Swin Transformer (Liu et al., 2021), RegViT (Darcet et al., 2024), and RoPEViT (Heo et al., 2024), we used the implementation provided by Wightman (2019), and for Jumbo (Fuller et al., 2025) we used their official repository. All of these models are reproduced. Jumbo is a new work the repository is incomplete, hence, we used the available code and implemented some of the components based on the paper.

#### A.2 COMPUTE RESOURCES

Our experiments were conducted using NVIDIA RTX A6000 48GB, V100 32GB, and A100 40GB GPUs. The Tiny, Small, and Base backbones of LocAtViT require 15GB, 29GB, and 29GB of GPU memory with a local batch size of 512, 512, and 256, respectively.

#### A.3 LLM USAGE

We used LLMs to aid or polish writing. Adhering to ICLR's author guideline, we include additional information here. We used LLMs to generate codes for plotting figures, tables, and other code or LaTeX related issues. We also used LLMs to improve the writing, polish, or shorten the paragraphs, while double checking the output.

## B LOCATVIT COMPARISON WITH RELATED WORK

In Tab. 1, we included five baseline methods and implemented LocAt for each. Table 5 compares LocAtViT to multiple related works from Sec. 2: CvT-21 (Wu et al., 2021a), Conformer (Peng et al., 2021), ConViT (d'Ascoli et al., 2021), Twins (Chu et al., 2023; 2021), DaViT (Ding et al., 2022), and GCViT (Hatamizadeh et al., 2023). We utilized the publicly available code and checkpoints, and evaluated the models on our segmentation pipeline, as described in Sec. 5. Although LocAtViT does not achieve the best classification performance, LocAt helps ViT outperform methods like Twins across all three segmentation benchmarks.

Table 5: **Segmentation and classification performance** of the Base backbone of related works and the proposed LocAtViT.

Method	Segn ADE	nentation mIo P-Context	Top-1 (%) ImageNet	
CvT-21	21.40	40.91	29.29	82.50
Conformer	22.11	40.03	26.37	83.83
ConViT	23.08	44.82	25.20	82.30
Twins	30.47	44.55	32.27	82.71
DaViT	30.68	44.87	32.38	84.64
GCViT	30.91	44.71	32.77	84.47
LocAtViT	32.64	45.35	33.62	82.31

## C ADDITIONAL QUALITATIVE EXPERIMENTS

Figure 4 provides three additional images from the mini-ImageNet dataset, alongside the attention maps of the <code>[CLS]</code> token and several patches for ViT and LocAtViT.

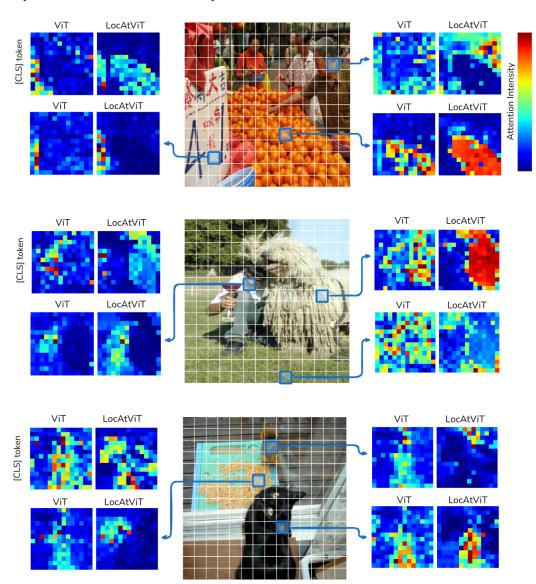


Figure 4: **Qualitative evaluation on the attention maps.** The final attention map of ViT and LocAtViT for the <code>[CLS]</code> token and three different patches are illustrated for three different images from mini-ImageNet with labels: *orange*, *Komondor*, and *corn*.

## D ABLATION STUDY ON SELF-ATTENTION

In this section, we perform ablations on the design choices inside the GAug self-attention module.

#### D.1 GAUSSIAN BASED ON INPUT

In the original ViT, a query vector intuitively determines the information a patch should be looking for. Since the Gaussian variance controls how far a patch attends to its surroundings, we compute  $\Sigma$  based on the query matrix in Eq. (5). Table 6 compares this approach to computing  $\Sigma$  based on x,

Table 6: **Ablations on GAug attention components**.  $\Delta \#Params$  shows the difference in the number of the parameters of each model compared to LocAtViT (first row). Experiments are conducted on mini-ImageNet.

	Tiny	Base	$\Delta$ #Params
LocAtViT (Sec. 4) Gaussian from x	78.47 79.10	84.86 85.18	+18,504,+329,868
Isotropic Gaussian Fixed width $\sigma = 1$ $\sigma = 5$ $\sigma = 10$	78.71 75.20 76.41 75.53	84.66 82.81 82.65 82.42	-780 $-2,340$ $-2,340$ $-2,340$
No scaling Auto $\alpha$	76.26 78.48	83.07 84.54	$-780 \\ -780$

the self-attention input. While the latter improves performance, it significantly increases the number of parameters.

#### D.2 VARIANCE MATRIX

To comply with a more general setting, we assigned separate variances for each image axis. An alternative is to use a single variance per patch, forming an isotropic Gaussian kernel. This simplifies Eq. (8) to:

$$\mathbf{G}_{pt} = \exp\left(-\frac{\sum_{m=1}^{2} \mathbf{D}_{ptm}}{2\sigma_{p}^{2}}\right). \tag{12}$$

The result of this modification is referred to as *Isotropic Gaussian* in Tab. 6. This table also compares this approach with another experiment where the Gaussian kernel width is fixed different constant values, instead of being patch-specific and query-based. These results indicate that an isotropic Gaussian kernel performs comparably, but a fixed kernel width substantially diminishes performance, demonstrating the importance of our dynamic input-dependent kernel width.

#### D.3 NO SUPPLEMENT MATRIX SCALING

In Sec. 4.1, we discussed the procedure to obtain  $\alpha$  using learnable parameters. We motivated this decision by the fact that the scale of attention logits and the supplement matrix might not match. Table 6 we show results for an experiment ( $No \alpha$ ) in which the supplement matrix in Eq. (10) is not scaled, i.e., S = G, i.e.,  $\alpha = 1$ . This setting reduces the accuracy, confirming the need to scale the supplement matrix properly.

## D.4 AUTOMATIC SCALING OF THE SUPPLEMENT MATRIX

As mentioned, we motivated the need for scaling the supplement matrix before adding it to the attention logits in Sec. 4.1. We now propose a parameter-free, input-dependent scheme,  $Auto \alpha$ , that automatically matches the scale of S to that of the original attention logits. Concretely, let N = 1 + hw,  $\mathbf{q}, \mathbf{k} \in \mathbb{R}^{N \times d}$ , and define the row-wise  $\ell_2$ -norm vectors:

$$\mathbf{r} = \left[ \|\mathbf{q}_1\|_2, \dots, \|\mathbf{q}_n\|_2 \right]^\top, \tag{13}$$

$$\mathbf{u} = \left[ \|\mathbf{k}_1\|_2, \dots, \|\mathbf{k}_n\|_2 \right]^{\top}. \tag{14}$$

Then the standard attention logits satisfy:

$$\frac{\mathbf{q}\mathbf{k}^{\top}}{\sqrt{d}} = \left(\frac{\mathbf{r}\mathbf{u}^{\top}}{\sqrt{d}}\right) \circ \cos(\mathbf{q}, \mathbf{k}),\tag{15}$$

where  $\circ$  denotes the Hadamard product, and  $\cos(\mathbf{q}, \mathbf{k}) \in \mathbb{R}^{N \times N}$  has entries  $\cos(\mathbf{q}_i, \mathbf{k}_j)$ . Hence, if we set

$$\alpha = \frac{\mathbf{r}\mathbf{u}^{\top}}{\sqrt{d}} \in \mathbb{R}^{N \times N},\tag{16}$$

then the modified logits in Eq. (4) can be rewritten as

$$\frac{\mathbf{q}\,\mathbf{k}^{\top}}{\sqrt{d}} + \mathbf{S} = \boldsymbol{\alpha} \circ (\cos(\mathbf{q}, \mathbf{k}) + \mathbf{G}), \tag{17}$$

where both terms inside the parentheses are bounded (in [-1,1] and [0,1], respectively), ensuring that S scales comparably to the original logits.

However, using  $\alpha \circ G$  would independently scale each entry of G, destroying the Gaussian kernel structure (each row of G is a kernel centered at one patch). To preserve each kernel's shape, we average  $\alpha$  across columns:

$$\bar{\alpha}_i = \frac{1}{N} \sum_{j=1}^N \alpha_{ij}, \quad \bar{\alpha} = [\bar{\alpha}_1, \dots, \bar{\alpha}_n]^\top \in \mathbb{R}^N,$$
 (18)

and then form:

$$\mathbf{S} = \operatorname{diag}(\bar{\boldsymbol{\alpha}}) \mathbf{G}, \tag{19}$$

similar to Eq. (10). This row-wise scaling applies a single factor to each Gaussian kernel, preserving its shape while matching its magnitude to the attention logits.

Auto  $\alpha$  performs close to learnable  $\alpha$  in the original LocAtViT, with slightly fewer parameters. However, we decided to use the learnable setting to make the formulation simpler.