REVISITING [CLS] AND PATCH TOKEN INTERACTION IN VISION TRANSFORMERS

Anonymous authorsPaper under double-blind review

ABSTRACT

Vision Transformers have emerged as powerful, scalable and versatile representation learners. To capture both global and local features, a learnable [CLS] class token is typically prepended to the input sequence of patch tokens. Despite their distinct nature, both token types are processed identically throughout the model. In this work, we investigate the friction between global and local feature learning under different pre-training strategies by analyzing the interactions between class and patch tokens. Our analysis reveals that standard normalization layers introduce an implicit differentiation between these token types. Building on this insight, we propose specialized processing paths that selectively disentangle the computational flow of class and patch tokens, particularly within normalization layers and early query-key-value projections. This targeted specialization leads to significantly improved patch representation quality for dense prediction tasks. Our experiments demonstrate segmentation performance gains of over 2 mIoU points on standard benchmarks, while maintaining strong classification accuracy. The proposed modifications introduce only an 8% increase in parameters, with no additional computational overhead. Through comprehensive ablations, we provide insights into which architectural components benefit most from specialization and how our approach generalizes across model scales and learning frameworks.

1 Introduction

In recent years, significant progress has been made in developing vision foundation models capable of generating rich and highly generalizable visual representations for images. Notably, latest state-of-the-art results have been achieved using Vision Transformer (ViT) models (Dosovitskiy et al., 2021) trained under various paradigms, including fully-supervised (Touvron et al., 2022), weakly supervised (Radford et al., 2021; Bolya et al., 2025), and self-supervised learning (Zhou et al., 2021; Oquab et al., 2023; Siméoni et al., 2025). These models capture a wide spectrum of visual semantics, enabling robust performance across a diverse range of downstream tasks and data domains.

The ViT architecture (Dosovitskiy et al., 2021) processes images by dividing them into fixed-size patches, which are then embedded and fed to a sequence of transformer blocks. Typically, a trainable

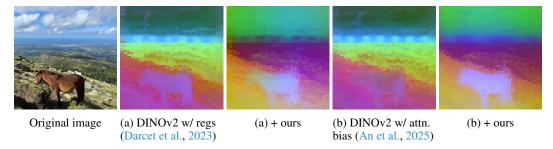


Figure 1: **Visualization of the impact** of our proposed layer specialization for [CLS] and patch tokens on the patch features obtained with DINOv2 when using two strategies to mitigate artifacts, namely registers ('regs') (Darcet et al., 2023) and attention bias ('attn. bias') (An et al., 2025). We display the first PCA components of model outputs in RGB.

class token [CLS] is prepended to the sequence of patch embeddings and is designed to aggregate information from all patches. Patches and [CLS] tokens are trained with different objectives, if any. For instance, most pre-training methods apply a loss function solely on the [CLS] token (Chen et al., 2020; Grill et al., 2020; Caron et al., 2021; Radford et al., 2021; Touvron et al., 2022). Some employs an objective on patch tokens only (He et al., 2022), while others train the [CLS] and patch tokens with separate losses (Zhou et al., 2021; Oquab et al., 2023; Siméoni et al., 2025). Regardless of the specific training paradigm, recent works (Darcet et al., 2023; An et al., 2025; Siméoni et al., 2025) show that there is a persistent imbalance between the [CLS] and patch tokens. Proposed solutions to this issue include introducing additional storage tokens into the input sequence (Darcet et al., 2023), modifying the attention mechanism (An et al., 2025), or incorporating additional loss terms to explicitly constrain patch locality (Siméoni et al., 2025). In contrast, we hypothesize that the imbalance arises because models process the [CLS] and patch tokens through identical computational pipelines, despite their fundamentally different roles and nature, and propose disentangling their treatment to overcome the imbalance.

In this work, we analyze the model statistics in order to better understand the internal mechanisms that govern the interaction between the <code>[CLS]</code> and patch tokens. Our analysis reveals a surprising finding: normalization layers are already implicitly learning to distinguish between the <code>[CLS]</code> and patch tokens before the attention mechanism. Building on this insight, we introduce a simple yet effective architectural modification that explicitly separates the processing of the <code>[CLS]</code> and patch tokens, as illustrated in Fig. 5. With just a minimal set of specialized layers, our approach leads to noticeably richer dense features (see Fig. 1) and delivers substantial gains on dense prediction tasks. For instance, we improve the average mIoU scores on segmentation benchmarks by as much as 2.2 points with a ViT-L. This work sheds light on hidden dynamics within transformer models and also demonstrates how targeted architectural changes can translate into significant real-world performance improvements. We make the following contributions:

- We analyze the interactions between [CLS] and patch tokens within Vision Transformers, and show that models implicitly attempt to distinguish them through normalization layers.
- We propose an architectural modification that specializes their computations to reduce the friction between them, while keeping the number of operations constant. We study different specialization strategies for transformer block components.
- We demonstrate the generalizability of our approach across model scales and learning frameworks, showing significant improvements in dense prediction tasks without compromising classification performance.

2 RELATED WORK

Vision Transformers Inspired by Vaswani et al. (2017) and first introduced by Dosovitskiy et al. (2021), Vision Transformer has become an architecture of choice when building vision models. A typical ViT model consists of a patch embedder and a stack of transformer blocks. Given an image, the patch embedder divides it into equally-sized square patches and transforms them into patch tokens that represent local information in the image. Optionally, a learnable [CLS] token is added to the set of patch tokens in order to capture global information. All tokens are then passed through the transformer blocks which process them with various transformations, most notably the multi-head self-attention operator (Vaswani et al., 2017) that allows tokens to attend to each other. Built on the original architecture, subsequent works have introduced additional components to improve various aspects of ViTs such as data efficiency (Touvron et al., 2020; Yuan et al., 2021), computational cost (Liu et al., 2021; Bolya et al., 2023) and normalization (Touvron et al., 2021). ViT architecture has enabled state-of-the-art performance in various tasks (Carion et al., 2020; Strudel et al., 2021), simplified multi-modal learning (Radford et al., 2021; Fini et al., 2024), and led to excellent local and global representation in foundation models (Oquab et al., 2014; Tschannen et al., 2025; Siméoni et al., 2025). In most ViTs, [CLS] and patch tokens are functionally interchangeable in transformer blocks – they are processed in identical manner using the same operators – despite their distinctive nature. We show in our analysis that the identical treatment of these tokens is suboptimal and that disentangling them leads to better local features for dense tasks.

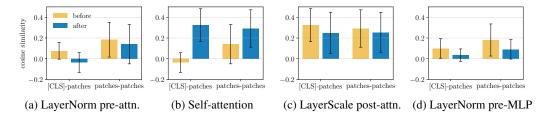


Figure 2: **[CLS]** -patches separation effect within transformer blocks in vanilla DINOv2 ViT-L model. We show mean and standard deviation of cosine similarity between [CLS] and all patches, and all-to-all patches, before and after each transformer layers. 'attn.' stands for attention.

Improving dense feature learning Visual representation learning approaches have mostly focus on optimizing the global representation by primarily training the [CLS] token to summarize the image content either in supervised (Touvron et al., 2022), weakly supervised (Radford et al., 2021; Tschannen et al., 2025) or self-supervised settings (Caron et al., 2021; Liu et al., 2021). As a byproduct, they also produce local representation that perform well on tasks that require fine-grained features such as object detection, semantic segmentation or depth estimation. Most notably, the selfsupervised method DINO (Caron et al., 2021) produces excellent patch features that supercharge research on unsupervised object detection and segmentation. iBoT (Zhou et al., 2021) augments DINO with masked image modeling (He et al., 2022) to optimize both global and local representation. DINOv2 (Oquab et al., 2023) introduces new technical components such as Sinkhorn-Knopp centering and untying heads to successfully scale DINO to large datasets and model sizes, achieving excellent performance on dense tasks. Learning meaningful dense features with Vision Transformers is not without challenges. Darcet et al. (2023) discusses the noisy attention maps produced by models trained at scale during longer training periods. This issue, which degrades dense prediction performance, is caused by some patch tokens losing their local context after being repurposed by the model to store global information. They propose an architectural solution with registers to mitigate these issues. Other successful attempts to enhance the quality of local features include regularizing similarity to neighbor patches post-training (Pariza et al., 2024) or recovering patch similarity with Gram anchoring mechanism (Siméoni et al., 2025). Similar to these works, we improve dense features quality during training by specializing [CLS] and patch tokens treatment within the Transformer blocks of ViTs and thus reducing the friction between them.

3 FRICTION BETWEEN [CLS] AND PATCHES

Vision Transformers are typically trained with a trainable [CLS] token which encodes global information about the image prepended to the sequence of patch tokens. Despite the distinct nature of the [CLS] and patch tokens, current models treat them equivalently, applying the exact same operations to both. However, (Darcet et al., 2023) has highlighted potential communication issues between these two types of tokens, leading to a severe loss of locality of patch tokens and the appearance of undesirable outliers in the attention maps. While registers help to mitigate the appearance of artifacts, we argue that more could be done. Our observations indicate that a degree of friction persists between [CLS] and patch tokens, as discussed below.

ViTs differentiate [CLS] and patch tokens for the attention We analyze the interplay between [CLS] and patch tokens by computing their similarity at different points within the model, before and after principal layers in each transformer block. In Fig. 2, we visualize the mean and standard deviation of these similarities. Our results are averaged over patches of 1000 images and across all model blocks. Additionally, we present the same statistics between patches. While certain operations—such as the LayerScale applied post-attention—have little effect on the similarity between [CLS] and patch tokens, the self-attention layer markedly increases their similarity. This increase is expected, as self-attention realigns the different token types. However, our analysis uncovers a surprising phenomenon: the representations of [CLS] and patch tokens naturally diverge at specific stages of the computational pipeline, particularly just before attention operations. Indeed, the LayerNorm applied before attention drastically reduces the similarity between [CLS] and patch tokens, bringing it close to zero. This implicit differentiation indicates that the model attempts to adapt to

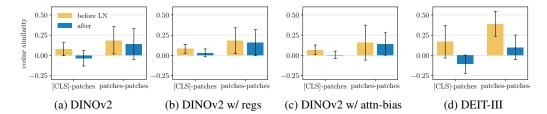


Figure 3: **Impact of LayerNorm before attention layer** for different pre-trained models. We show mean and standard deviation of cosine similarity between [CLS] and all patches, and between all patches. Statistics visualized **before** and **after** LayerNorm (LN).

the distinct functional roles of these token types before the attention mechanism, despite their shared parameterization. We plot the statistics of more layers in Appendix A.2.

Role of the pre-attention LayerNorm In Fig. 3, we focus on the impact of the pre-attention LayerNorm to the similarity between the [CLS] and patch tokens in different pre-trained models, including DINOv2 (Oquab et al., 2023) and its variants with registers (Darcet et al., 2023), noted 'regs', and attention bias (An et al., 2025), noted 'attn. bias', and supervised DEIT-III (Touvron et al., 2022). It can be observed that in all cases, prior to the attention mechanism, the Layer-Norm disentangles the [CLS] and patch tokens, enabling them to serve distinct functions within the attention process. This phenomenon appears in all pre-trained models with different extent. For instance, the LayerNorm strongly enforces negative correlation between [CLS] and patches in DI-NOv2 and DEIT-III while keeping the correlation close to zeros in the variants of DINOv2. In contrast, the similarity among patch tokens remains positive and largely stable, with only a slight decrease observed—a phenomenon we interpret as a regularization effect. This effect likely prevents rank collapse and promotes a more uniform distribution of tokens on the unit sphere, consistent with observations reported in Wu et al. (2024).

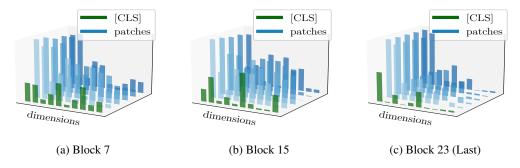


Figure 4: **Dimensions with biggest magnitudes** early (a), in the middle (b), at the end (c) of the model for [CLS] and patches. Tokens taken at the output of blocks. The considered model is a DINOv2 ViT-L with attention bias

Dimension separation The LayerNorm layer performs a point-wise normalization and a dimension-wise affine transformation, therefore a separation effect can appear when inputs have very different magnitudes in each dimension. In Fig. 4, we plot the dimensions with biggest absolute magnitudes—averaged over patches and [CLS]—at the output of different blocks. We observe that some specific dimensions are leveraged only by a certain token type and that the deeper we are in the model, the fewer token types share dimensions. This enables normalization layers to perform distinctive operations. More than just regularizing, they specialize and separate the tokens.

All the observations above indicate that treating [CLS] and patch tokens identically compels the model to allocate resources towards implicitly separating them, which could be used to learn more meaningful features. We argue that disentangling their treatments would facilitate the model in learning better representation, as discussed in next section.

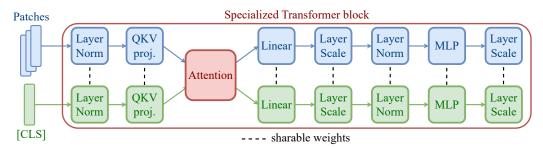


Figure 5: **Architecture specialization.** We investigate how [CLS] and patch tokens can be processed through specialized layers, while preserving their interactions within the attention mechanism.

4 [CLS] - PATCHES SPECIALIZATION: ANALYSIS

In this section, we first define our proposed layer specialization in Sec. 4.1 and set the experimental setting in Sec. 4.2. In Sec. 4.3, we discuss the benefits of splitting normalizations for [CLS] and patch tokens. We also investigate which part of the model needs specialization in Sec. 4.4, and more specifically which layers in Sec. 4.5.

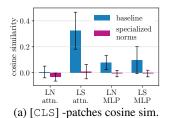
4.1 OUR PROPOSAL: LAYER SPECIALIZATION

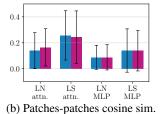
Based on observations made in the previous section, we explore disentangling the computation of global and local representations in ViTs. Taking inspiration from the success of double-stream architectures to handle different modalities (Esser et al., 2024), we explore a similar approach for the [CLS] and patch tokens. More specifically, inside a classic transformer block, [CLS] and patch tokens go through several layers: projections, some normalizations and a MLP. We propose to decouple the [CLS] and patch tokens by processing them with different weights for certain layers. Indeed, instead of using a single layer to process both token types, we introduce two distinct layers—each with its own set of weights-specialized for either [CLS] or patch tokens. This allows each layer to better capture the unique characteristics of its respective token type. However, the tokens continue to interact through the attention mechanism as usual, ensuring information flow between [CLS] and patch tokens is preserved. An illustration of this specialized architecture is provided in Fig. 5. While this approach introduces some additional memory overhead, our experiments show that the increase in model size remains small—approximately 8%—when layer specialization is applied selectively to achieve optimal performance. More importantly, layer specialization does not increase inference FLOPs, as the model continues to perform the same computational operations during inference. This ensures that the efficiency of the model is maintained, even as we enhance its representational capacity through targeted specialization.

4.2 EXPERIMENTAL SETTING: TRAINING AND EVALUATION

Training We investigate layer specialization with different pre-training paradigms including the popular self-supervised strategy DINOv2 (Oquab et al., 2023) and the fully-supervised DeiT-III (Touvron et al., 2022). We also investigate different model sizes (ViT-B, L, H). Unless specified otherwise, we produce results with a ViT-L model trained following DINOv2 recipe. Following An et al. (2025), we integrate the attention bias strategy, which mitigates high-norm anomalies (Darcet et al., 2023) without introducing additional tokens, in all models and attention operations. More discussion can be found in Appendix A.1. For DINOv2, we train our models on ImageNet-22K (Ridnik et al., 2021) dataset for 600k training steps. For DeiT-III, we train our models on ImageNet-1K (Deng et al., 2009) for respectively 400 epochs on ViT-B and 800 epochs on ViT-L. For both training paradigms, we pre-train using the first pre-training phase, and drop the high-resolution fine-tuning step. We report more details in Appendix A.3.

Evaluation Following Oquab et al. (2014), we assess model representations via linear probing on global, with ImageNet-1k (Deng et al., 2009), and dense prediction tasks. For semantic segmentation, we use ADE20K (Zhou et al., 2017), Cityscapes (Cordts et al., 2016) and PASCAL VOC (Everingham et al., 2010), reporting mIoU. For depth estimation, we use KITTI (Geiger et al.,





Spec.	Linear Acc.	Avg Seg.	Avg Depth↓
_	85.4	64.5	1.232
norms	85.1	65.6	1.178

(c) Norm specialization

Figure 6: **Normalization specialization.** Mean and standard deviation of the cosine similarity computed between (a) [CLS] and all patches and (b) all to all patches. We compare post-normalization statistics between the baseline architecture (DINOv2 ViT-L with attn. bias) and when specializing the normalization layers for [CLS] and patch token (specialized norms). (c) Quantitative results with specialized normalization layers. 'LN' stands for LayerNorm and 'LS' for LayerScale.

2013), NYU Depth v2 (Nathan Silberman & Fergus, 2012) and SUN RGB-D (Song et al., 2015), reporting RMSE. Some tables show average segmentation and depth scores across corresponding benchmarks. More details in Appendix A.3.

4.3 Specializing Normalization Layers

As discussed in Sec. 3, ViTs attempt to separate the [CLS] and patch tokens with the LayerNorm applied prior to the attention operation. Building on this observation, our initial experiment focuses on specializing the normalization layers (LayerNorms and Layer Scales) within the model, with the aim of further supporting the model's inherent tendency to separate these feature types.

We specialize the normalization layers in all blocks of the model. This lightweight modification introduces only 0.05% additional parameters, yet significantly alters the feature distributions. In Fig. 6a, we report the mean and standard deviation of the cosine similarity between the [CLS] and patch tokens, computed after each normalization layer. We compare our variant with specialized normalization weights to the baseline. Conversely, Fig. 6b shows the corresponding statistics when using all patches instead of the [CLS]. We observe that specializing the normalization layers further amplifies the disentanglement of the [CLS] and patch tokens, resulting in a more distinct separation of their embeddings after each normalization step.

The impact of these specialized normalizations is quantified in Fig. 6c. The specialization leads to significant improvements on dense prediction tasks, yielding an average increase of +1.1 mIoU points on segmentation benchmarks and an improvement of -0.054m on depth estimation. These results highlight that a better specialization of the token types benefits the patches representations. On the other side, global results are slightly degrading. We however show in the next section that this loss can be mitigated. Unless otherwise specified, in the remainder of the paper, we apply specialized normalization layers to all transformer blocks.

4.4 BLOCK-LEVEL TARGETED SPECIALIZATION

While normalization layers in ViTs show in overall a [CLS] -patch separation effect, we have observed that the extent of their impact is not uniform across all blocks. It can be seen from Fig. 7, which depicts [CLS] -patches cosine similarity before and after the first LayerNorm in each block, that the separation effect of the normalization varies depending on its position within the model. Indeed, blocks at the beginning and near the end see the most impact. We hypothesize that the importance of separation in the early blocks stems from their proximity to the different inputs. Although the [CLS] token is trained to summarize information from the patches, it is initialized as a learned parameter and thus has

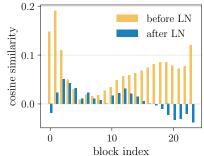
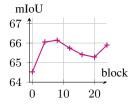
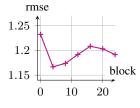


Figure 7: Mean cosine similarity between [CLS] and all patches before and after the first LayerNorm (LN) of each block.

an input distribution very distinct from that of the patch tokens. Later in the model, separation be-

Specialized	Linear	Avg	Avg
blocks	Acc.	Seg.	Depth \downarrow
Ø	85.4	64.5	1.232
1st half	85.1	65.7	1.191
2nd half	85.2	64.5	1.236
All	85.0	65.9	1.191





(a) Which model's part to specialize

(b) Avg. segmentation perf.

(c) Avg. depth perf. ↓

Figure 8: **Block specialization.** (a) Performance metrics when specializing specific parts of the model (b) Average segmentation scores and (c) average depth rmse (\downarrow) vs number of specialized blocks at the beginning of the model. Normalization layers are specialized in all blocks. The base model is a DINOv2 ViT-L with attention bias.

comes important again as tokens are closer to the final representations and the training objectives. The observations above suggest that we can benefit from more targeted specialization within the model. We study next which blocks should be specialized to optimize the model's performance.

We first quantitatively compare the impact of specializing different sections of the model, as shown in Fig. 8a. To this end, we train DINOv2 while specializing either the first half, the second or all of the transformer blocks, on top of specializing all normalization layers. Within a block, all layers are specialized. Our findings indicate that the best performance is achieved when specializing the early layers, which are closest to the input. Specifically, specializing the first half of the layers improves the segmentation results by an average of 1.2 mIoU points, with only a negligible decrease in linear accuracy. In contrast, specializing the late layers yields no improvement compared to the baseline. We attribute this to the fact that [CLS] and patch tokens share the representation space in the first part; once this interaction is established, further specialization has limited effect. Finally, while specializing all layers produces the highest segmentation performance, it comes with a higher memory cost and a larger drop in linear accuracy.

We further analyze how the number of specialized blocks, starting from the first, affects performance (Fig. 8b and 8c). We vary the number of blocks specialized from 0 to 24 (total number of blocks in ViT-L) in steps of 4, and observe that specializing the first third of the model yields the best results, while specializing later layers degrades performance. Notably, the optimal point at one third of the model's depth coincides with a marked shift in the statistics of similarity scores shown in Fig. 7, which might explain the effectiveness of specializing the early layers.

4.5 TARGETED SPECIALIZATION WITHIN TRANSFORMERS BLOCKS

The previous section has shown that careful selection of transformer blocks for specialization is important for optimizing the performance. We now explore whether further improvements can be achieved with a targeted selection of *specific layers* to specialize within the transformer blocks. In the following experiments, we specialize different layers of blocks in the first third of the model while also applying specialization to the normalization layers in all blocks.

Table 1 shows model performance on global and dense prediction tasks when specializing different layers (QKV projection, Linear and MLP, see Fig. 5). We observe that the performance on global task remain largely stable independently of the selected layers. In contrast, results on dense segmentation tasks get further improvements beyond what is achieved with normalization specialization alone. Interestingly, the gains do not increase monotonically with the number of specialized layers or additional parameters, as might be expected from typical scaling laws (Touvron et al., 2021). Specializing either or both QKV and post-attention projections consistently yields improvements. In particular, the greatest performance gains are achieved by specializing the QKV projection, which introduces only 8% additional parameters while delivering an average increase of +1 mIoU point over normalizations alone. In contrast, specializing the post-attention projection does not offer further benefits, and specialization of the MLP layer either has no effect or negatively impacts performance. Note that we could ease this 8% memory cost overhead with Low Rank Adaptation used when specializing the QKV projection. We produce encouraging preliminary results with different ranks in Appendix A.5 and leave further investigation as future work.

Table 1: **Layer specialization ablation.** Performance and increase in parameter count of models trained with different layer specialization strategies, applied to the first third of the transformer blocks. In all cases, the normalization specialization described in Sec. 4.3 is applied in all blocks. The base model is a DINOv2 ViT-L with attention bias.

QKV proj.	Linear	MLP	Parameter Increase (%)	Linear Accuracy	Avg. Seg.	Avg. Depth ↓
			0.05	85.1	65.6	1.178
	\checkmark		3	85.3	66.1	1.191
\checkmark			8	85.2	66.6	1.165
\checkmark	\checkmark		11	85.1	66.1	1.180
		\checkmark	22	85.1	65.2	1.189
	\checkmark	\checkmark	25	85.1	65.9	1.163
\checkmark		\checkmark	30	85.3	65.6	1.185
\checkmark	\checkmark	\checkmark	33	85.3	<u>66.1</u>	1.174

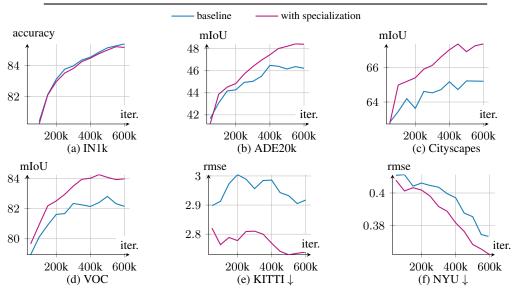


Figure 9: **Performances and training dynamics.** Performance vs. training iterations on global task ImageNet classification (IN1k) and dense tasks—segmentation (ADE20k, Cityscapes, VOC) and depth (KITTI, NYU)—with linear probing. We compare baseline DINOv2 ViT-L with attn. bias and when specializing QKV projection in the first 1/3 of the model and all normalizations.

Our overall results show that increasing the disentanglement between [CLS] and patch tokens before the attention mechanism (with separated normalizations and projection) contributes to improved dense prediction performance. We hypothesize that encouraging the [CLS] and patch tokens to assume more distinct roles in the attention mechanism enhances their interactions, ultimately improving overall model effectiveness. We also report results in Appendix A.6 for the setting where only the QKV projections are specialized, while the normalization layers remain shared. In this configuration, performance is comparable to the baseline, indicating that the specialization of normalization layers is critical to achieve improvements, as shown in Sec. 4.3.

In Fig. 9, we compare performance dynamics of normalization and QKV projection specialization against baseline DINOv2 with attention bias. Across all dense benchmarks—of both segmentation and depth estimation—specialization consistently enhances results. These improvements are evident from early stages of training and continue to increase over time. This trend suggests that employing specialization not only boosts performance but also contributes to more stable training dynamics.

4.6 GENERALIZATION RESULTS

We investigate the generalizability of our specialization approach across different variants of DI-NOv2, as presented in the upper part of Table 2. Specifically, we train models using the DINOv2 recipe with two high-norm handler strategies—registers (Darcet et al., 2023) ("4 registers") and at-

Table 2: **Generalizability** of the specialization on (a) DINOv2 when using different high-norm handling strategies (4 registers (Darcet et al., 2023), attention bias (An et al., 2025) ('attn. bias') or none (\emptyset)), (b) different ViT sizes and (c) a supervised framework: DeiT-III. Relative difference between baseline and our specialization ('+ours') is shown in green if improvement and red otherwise.

Method Size		Classif.	Segmentation		Depth \downarrow			
1,1041104	Size	ImNet	ADE	City	VOC	KITTI	NYU	SUN
DINOv2 - With high-norm handling strategies								
Ø	L	85.3	45.7	64.2	82.1	2.868	0.389	0.410
+ours		85.3+0.0%	47.3+3.5%	66.6+3.7%	83.7+1.9%	2.787 -2.8%	0.369-5.1%	0.390-4.9%
4 registers	L	85.3	45.6	64.9	82.2	2.893	0.372	0.411
+ours		85.3+0.0%	47.5+4.2%	65.9+1.5%	83.6+1.7%	2.906+0.4%	0.367-1.3%	0.395-3.9%
Attn. bias	L	85.4	46.2	65.2	82.2	2.917	0.373	0.406
+ours		85.2-0.2%	48.4+4.8%	67.4+3.4%	84.0+2.2%	2.739-6.1%	0.362-2.9%	0.393 -3.2%
DINOv2 - C	DINOv2 - Other sizes							
Attn. bias	В	80.4	38.3	58.4	76.6	3.250	0.462	0.464
+ours		80.6+0.2%	38.5+0.5%	60.3+3.3%	76.5 -0.1%	3.236-0.4%	0.448-3.0%	0.470+1.3%
Attn. bias	Н	86.2	48.1	67.0	83.1	2.717	0.359	0.387
+ours		86.1 -0.1%	49.2+2.3%	67.1+0.1%	83.5+0.5%	2.752+1.3%	0.344-4.2%	0.386-0.3%
DeiT-III								
Attn. bias	В	81.8	25.4	61.7	48.9	5.040	0.747	0.823
+ours		81.7-0.1%	26.3+3.5%	62.7+1.6%	50.7+3.7%	4.900-2.8%	0.732-2.0%	0.809-1.7%

tention bias (An et al., 2025) ("attn. bias")—as well as without any handler. In all cases, we observe that specialization consistently boosts dense prediction results by up to 4.8% on ADE20k, while having a negligible effect on classification performance (no decrease greater than 0.2%). This shows that better separating the treatment of [CLS] and patch tokens is complementary to both high-norm handling strategies to improve dense features. We also investigate the specialization on DINOv2 ViT-B and ViT-H models and present results in the middle section of Table 2. It can be seen that our proposed specialization leads to improvements on most benchmarks, confirming its generalizability across different ViT model sizes.

We further explore the fully-supervised training setting by applying specialization to a ViT-B trained with DEIT-III (Touvron et al., 2022) strategy. We observe consistent improvements in dense prediction tasks, with gains reaching up to 3.7% on VOC. For ViT-L, specialization does not yield benefits, likely due to the training dynamics related to the absence of a local loss to guide dense feature learning which causes the dense performance to degrade over time (we provide more details in Appendix A.7). These results suggest that the effectiveness of the specialization may depend on training objectives, highlighting promising directions for future research.

Finally, we visualize the learned patch representations using PCA in Fig. 1 for DINOv2 models trained with either registers or attention bias. In both settings, incorporating our specialization strategy produces cleaner and more semantically meaningful patch representations. Specifically, this approach reduces artifacts in textures and uniform regions, resulting in more accurate object segmentation. More visualizations can be found in A.8.

5 CONCLUSION

In this work, we investigate the disentanglement of [CLS] and patches computations in Vision Transformers, focusing on their distinct roles and interactions. Through a comprehensive analysis, we demonstrate that disentangling their processing pathways and selectively specializing architectural layers leads to significant improvements in dense prediction tasks, including segmentation and depth estimation, while maintaining strong global performance. Our approach achieves these gains without increasing computational overhead, with minimal additional parameter cost, and generalizes across multiple ViT architectures and frameworks. These findings highlight the importance of tailored architectural designs and suggest promising directions for future research, including further exploration of efficient specialization strategies and applications to broader modalities and tasks.

REFERENCES

- Yongqi An, Xu Zhao, Tao Yu, Ming Tang, and Jinqiao Wang. Systematic outliers in large language models. *ICLR*, 2025.
- Daniel Bolya, Cheng-Yang Fu, Xiaoliang Dai, Peizhao Zhang, Christoph Feichtenhofer, and Judy Hoffman. Token merging: Your vit but faster. In *ICLR*, 2023.
 - Daniel Bolya, Po-Yao Huang, Peize Sun, Jang Hyun Cho, Andrea Madotto, Chen Wei, Tengyu Ma, Jiale Zhi, Jathushan Rajasegaran, Hanoona Rasheed, et al. Perception encoder: The best visual embeddings are not at the output of the network. *arXiv preprint arXiv:2504.13181*, 2025.
 - Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020.
 - 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 IEEE/CVF international conference on computer vision*, pp. 9650–9660, 2021.
 - Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pp. 1597–1607. PMLR, 2020.
 - Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3213–3223, 2016.
 - Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision transformers need registers. *arXiv preprint arXiv:2309.16588*, 2023.
 - Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
 - Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021.
 - Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first international conference on machine learning*, 2024.
 - Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2): 303–338, 2010.
 - Enrico Fini, Mustafa Shukor, Xiujun Li, Philipp Dufter, Michal Klein, David Haldimann, Sai Aitharaju, Victor Guilherme Turrisi da Costa, Louis Béthune, Zhe Gan, Alexander T Toshev, Marcin Eichner, Moin Nabi, Yinfei Yang, Joshua M. Susskind, and Alaaeldin El-Nouby. Multimodal autoregressive pre-training of large vision encoders, 2024.
 - Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *The international journal of robotics research*, 32(11):1231–1237, 2013.
 - Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Álvaro Bernardo, Bilal Pires, Zhaohan Guo, et al. Bootstrap your own latent-a new approach to self-supervised learning. In *Advances in Neural Information Processing Systems*, volume 33, pp. 21271–21284, 2020.
 - Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16000–16009, 2022.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022.

- 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 (ICCV)*, pp. 10012–10022, October 2021.
- Pushmeet Kohli Nathan Silberman, Derek Hoiem and Rob Fergus. Indoor segmentation and support inference from rgbd images. In *ECCV*, 2012.
- Maxime Oquab, Leon Bottou, Ivan Laptev, and Josef Sivic. Learning and transferring mid-level image representations using convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1717–1724, 2014.
- Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. In *International Conference on Machine Learning*, pp. 27408–27438. PMLR, 2023.
- Valentinos Pariza, Mohammadreza Salehi, Gertjan Burghouts, Francesco Locatello, and Yuki M Asano. Near, far: Patch-ordering enhances vision foundation models' scene understanding. *arXiv* preprint arXiv:2408.11054, 2024.
- 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. PmLR, 2021.
- Tal Ridnik, Emanuel Ben-Baruch, Asaf Noy, and Lihi Zelnik-Manor. Imagenet-21k pretraining for the masses. *arXiv preprint arXiv:2104.10972*, 2021.
- Oriane Siméoni, Huy V Vo, Maximilian Seitzer, Federico Baldassarre, Maxime Oquab, Cijo Jose, Vasil Khalidov, Marc Szafraniec, Seungeun Yi, Michaël Ramamonjisoa, et al. Dinov3. *arXiv* preprint arXiv:2508.10104, 2025.
- Shuran Song, Samuel P Lichtenberg, and Jianxiong Xiao. Sun rgb-d: A rgb-d scene understanding benchmark suite. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 567–576, 2015.
- Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Segmenter: Transformer for semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 7262–7272, October 2021.
- Mingjie Sun, Xinlei Chen, J Zico Kolter, and Zhuang Liu. Massive activations in large language models. *arXiv preprint arXiv:2402.17762*, 2024.
- Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. *arXiv* preprint arXiv:2012.12877, 2020.
- Hugo Touvron, Matthieu Cord, Alexandre Sablayrolles, Gabriel Synnaeve, and Hervé Jégou. Going deeper with image transformers. In *ICCV*, pp. 32–42, 2021.
- Hugo Touvron, Matthieu Cord, and Hervé Jégou. Deit iii: Revenge of the vit. In *ECCV*, pp. 516–533. Springer, 2022.
- Michael Tschannen, Alexey Gritsenko, Xiao Wang, Muhammad Ferjad Naeem, Ibrahim Alabdulmohsin, Nikhil Parthasarathy, Talfan Evans, Lucas Beyer, Ye Xia, Basil Mustafa, Olivier Hénaff, Jeremiah Harmsen, Andreas Steiner, and Xiaohua Zhai. Siglip 2: Multilingual vision-language encoders with improved semantic understanding, localization, and dense features. *arXiv preprint arXiv:2502.14786*, 2025.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

Xinyi Wu, Amir Ajorlou, Yifei Wang, Stefanie Jegelka, and Ali Jadbabaie. On the role of attention masks and layer norm in transformers. *Advances in Neural Information Processing Systems*, 37: 14774–14809, 2024.

Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Zi-Hang Jiang, Francis E.H. Tay, Jiashi Feng, and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.

Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 633–641, 2017.

Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Yuille, and Tao Kong. ibot: Image bert pre-training with online tokenizer. *arXiv preprint arXiv:2111.07832*, 2021.

A APPENDIX

A.1 ADDRESSING TOKEN INTERACTION ANOMALIES

In this work, we examine how the distinct roles of [CLS] and patch tokens affect their interactions within the model. Previous works (Darcet et al., 2023; Sun et al., 2024) show that despite sharing computational pathways, these different token types develop inter-dependencies that can lead to token anomalies, manifested as high-norm outliers in the patch features space. These anomalies suggest an underlying tension in how information flows between global and local representations.

Registers. In order to mitigate such artifacts, observed when using different pre-training strategies (Oquab et al., 2023; Touvron et al., 2022; Radford et al., 2021), Darcet et al. (2023) propose to add learnable register tokens to the input sequence, whose roles are to replace the high-norm patches in the internal communication between patches and the [CLS] token. Doing so mitigates the appearance of such artifacts and boost overall results.

Attention bias. The recent study on artifacts in Large Language Models by An et al. (2025) investigates the systematic appearance of outliers which they link to the attention mechanism. They propose a solution consisting in adding learnable biases to the keys and values in each attention head. They analyze the equivalence of their solution compared to registers.

Table 3: The impact of norm handling strategies on DINOv2 results.

Norm. method	IN	ADE	City.	NYU↓
Ø	85.3	45.7	64.2	0.389
4 registers	85.3	45.6	64.9	0.372
attn. bias	85.4	46.2	65.2	0.373

In our experiments, we observe that both strategies have a similar impact on high-norm artifacts and as seen in Table 3, best overall performance is achieved when using the attention bias ('attn. bias') strategy, with a significant improvement on segmentation benchmarks (e.g. ADE20k and Cityscapes). To minimize confounding factors that could affect the interaction between the [CLS] and patch tokens, we adopt the attention bias strategy, which mitigates high-norm anomalies without introducing additional tokens.

A.2 EFFECT OF OTHER LAYERS ON [CLS] -PATCHES SIMILARITIES

We report in Fig. 10, the effect on [CLS] -patches similarity of the MLP and post-MLP LayerScale layers within transformer blocks. The MLP layer, similar to the self-attention layer, increases the

similarity between [CLS] and patches as it aligns the features. The post-MLP LayerScale, similar to other normalization layers, shows a stronger disentangling effect.

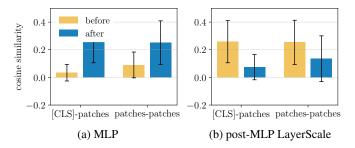


Figure 10: Effect of the MLP and post-MLP LayerScale layers on [CLS] -patches similarity in vanilla DINOv2 pre-trained model. We show mean and standard deviation of the cosine similarity between [CLS] and all patches ('CLS-patches'), and between patches ('patches-patches'), before and after the considered layers.

A.3 TRAINING AND EVALUATION DETAILS

Throughout this work, we follow the experimental protocol of Oquab et al. (2023) and evaluate the performance of the trained models on a set of global and dense task benchmarks.

Classification For the global task, we perform linear probing on ImageNet classification (Deng et al., 2009). We train a linear layer with SGD for 12500 iterations, using random-resized-crop data augmentation, and the [CLS] token as input for the linear layer. We also perform the following grid search on learning rate : $\{1.0e^{-5}, 2.0e^{-5}, 5.0e^{-5}, 0.0001, 0.0002, 0.0005, 0.001, 0.002, 0.005, 0.01, 0.02, 0.05, 0.1\}$

We then report the highest accuracy value obtained on the validation set as is common practice.

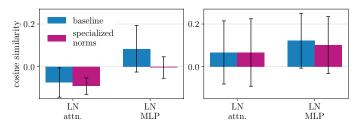
Segmentation For semantic segmentation, we use ADE20K (Zhou et al., 2017), Cityscapes (Cordts et al., 2016), and VOC (Everingham et al., 2010), and report the mean Intersection over Union (mIoU) scores for each. When we report the average performance of segmentation tasks, we average the scores across these 3 datasets. We train a linear classifier on the training set of each benchmark for iterations with a learning rate of $1e^{-3}$. This linear layer is applied on top of the patch output features (after the last layer normalization) of the frozen backbone, with the features further normalized using a trained batch normalization layer.

Depth estimation For depth estimation, we evaluate on KITTI (Geiger et al., 2013), NYU Depth v2 (Nathan Silberman & Fergus, 2012), and SUN RGB-D (Song et al., 2015), reporting the average Root Mean Squared Error (rmse) scores. When we report the average performance of dense tasks, we average the scores across these 3 datasets. We train a linear classifier on the training set of each benchmark for iterations with a learning rate of $1e^{-3}$. For the input of this linear layer, we take patch and [CLS] output features from four evenly spaced layers of the backbone, not applying the last layer normalization.

When training with DINOv2 and DeiT-III models, we follow the default configurations provided in the official repository, modified to add biases in the attention and to specialize layers or blocks.

A.4 SPECIALIZATION OF NORMALIZATION LAYERS IN DEIT-III

We report in Fig. 11 the impact of the specialization of the normalization layers when using DeiT-III pre-training strategy. Similar to the case of DINOv2, the average [CLS] -patch cosine similarity significantly reduces when employing the specialized normalization, showing the disentanglement effect.

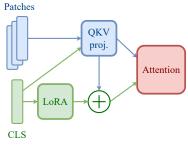


(a) [CLS] -patch cosine similarity (b) all-all patch cosine similarity

Figure 11: **Specialization of normalization layers.** Mean and standard deviation of the cosine similarity computed between (a) the [CLS] and all patches and (b) all patch to all patches. The average is computed over 1k images and all model blocks. We compare post-normalization statistics between the standard architecture ('Baseline') and the model after normalization specialization ('Specialized norms'). 'LN' stands for LayerNorm and 'LS' for LayerScale.

A.5 LORA APPROXIMATION

As the parameters increase can be a bottleneck for training efficiency, we explore the use of Low-Rank Adaptation (LoRA) (Hu et al., 2022) techniques to reduce the number of trainable parameters while maintaining performance. Additionally, we hypothetise that [CLS] and patches representations share common features. Hence we consider [CLS] stream as a specialization of patches stream instead of a complete different stream. Then, for a layer f that we choose to specialize, we compute the operation on the class token x_{cls} as the sum of the patches layer f_{patch} and a low-rank adaptation (LoRA) decomposition $f_{cls}^{(r)}$ of rank r.



Specialization	Param.	Linear	Avg	Avg
	Incr. (%)	Acc.	Seg.	Depth ↓
Ø	-	85.4	64.5	1.232
norms	0.05	85.1	65.6	1.178
+QKV	8.3	85.2	66.6	1.165
+LoRA QKV r=16	0.2	85.3	65.8	1.188
+LoRA QKV r=128	1.4	85.2	65.9	1.193

(a) LoRA design

(b) Results with LoRA

Figure 12: **LoRA impact.** (a) Visualization of LoRA design: [CLS] as an approximation of patches. (b) Performance metrics and parameter increase for different LoRA configurations (rank 16 and 128) during first third of the model. In all cases, the normalization specialization described in Sec. 4.3 is applied, corresponding to 'norms' row.

We conduct experiments in which we specialize normalization layers and the QKV projections with LoRA approximations of ranks 16 and 128 (over an embedding dimension of 1024). The results presented in Fig. 12b shows improvements (+0.2 and +0.3 in segmentation tasks) over specializing only the normalization layers, while adding a limited number of parameters (+0.15% and +1.35% respectively). We leave further investigations as future work.

A.6 NORMALIZATIONS ARE NEEDED

Additionally to the specialization experiments we produced in Sec. 4.5, we also conduct an experiment specializing QKV projection during the first third of the model, *but not the normalization layers*. We plot the results of this experiment in Table 4 compared to the baseline and to our best model specializing normalization layers and QKV projection during third of the model. We observe that specializing only QKV projections brings little improvement over the baseline, e.g. +0.2 mIoU pt in average on segmentation tasks. This shows that specializing the normalization layers is crucial for best performance.

Table 4: **Importance of specializing norms.** Performance of for different layer specialization (Spec.) strategies applied on the first third of the transformer blocks). Normalization layers are specialized in all blocks. Baseline is a ViT-L DINOv2 with attention bias.

Model	Linear Acc.	Avg. Seg.	Avg. Depth ↓
Baseline	85.4	64.5	1.232
Specialized norms	85.1	65.6	<u>1.178</u>
Specialized norms & QKV proj.	85.2	66.6	1.165
Specialized QKV proj.	85.4	64.7	1.211

A.7 ADDITIONAL RESULTS ON DEIT-III

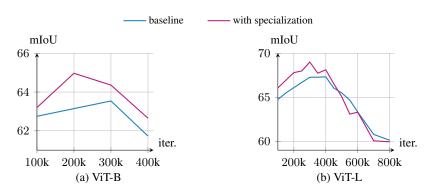


Figure 13: **DeiT-III training evolution.** We visualize VOC segmentation performance (mIoU) throughout training for (a) ViT-B and (b) ViT-L pre-trained with DeiT-III ('baseline') and when adding our layer specialization.

We report in Fig. 13 the performance curves on VOC segmentation task during the training of ViT-B and ViT-L models when following DeiT-III. We observe that the performance reaches its peak in the middle of the pre-training, then drops significantly towards the end. We attribute this behavior to the lack of a local loss to drive dense performances. We observe a significant gain with our specialization in the first half of the training, but the gains are then diluted in the drop, particularly in the case of ViT-L.

A.8 OTHER QUALITATIVE RESULTS

We produce in Fig. 14, 15 and 16 more qualitative results when pre-training the model following DINOv2 with the vanilla architecture, four registers or attention bias and when integrating our specialization. Each figure shows the first three components, computed with the patch features, and mapped to RGB. In all cases, we observe that the specialization helps to produce more precise patch features with less artifact. For instance, we invite the reader to pay attention to the back of the dog (first row), where the artifacts visible in the original pre-training are notably reduce with our specialization.

A.9 WRITING DETAILS

We have used Large Language Models (LLMs) to help write and proofread this paper. More specifically, they have helped to rephrase some parts of the text, propose synonyms, and check the grammar. We have carefully checked all the outputs of the LLMs to ensure that they are accurate and faithful to our work.

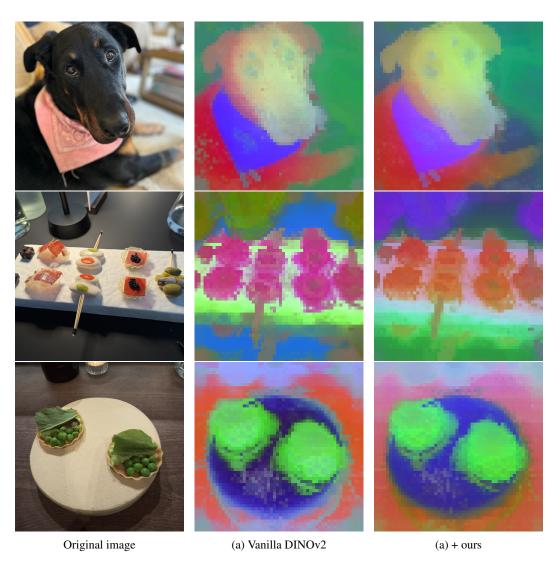


Figure 14: First PCA components of model outputs in RGB. Specialization of normalizations and QKV projections is made during 1/3 of the model. ViT-L with vanilla DINOv2.

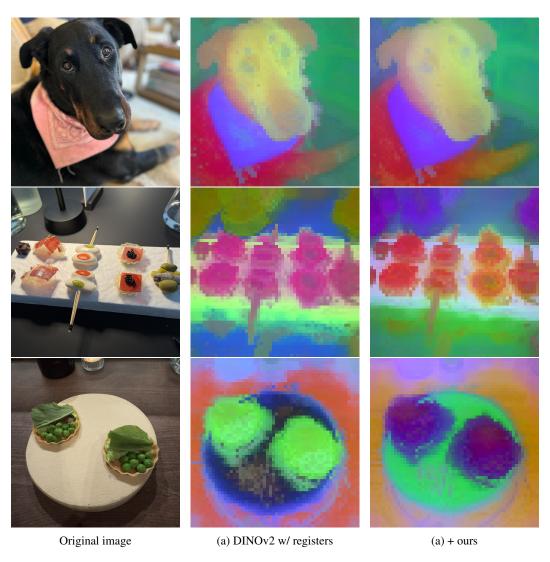


Figure 15: First PCA components of model outputs in RGB. Specialization of normalizations and QKV projections is made during 1/3 of the model. ViT-L DINOv2 with four registers.

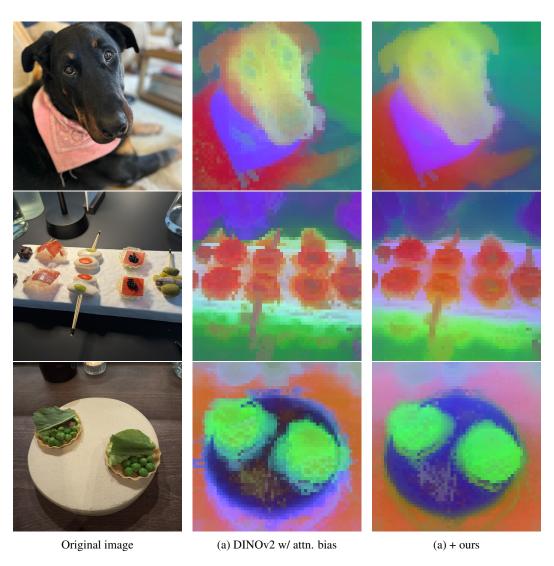


Figure 16: First PCA components of model outputs in RGB. Specialization of normalizations and QKV projections is made during 1/3 of the model. ViT-L DINOv2 with attention bias.