

What Makes for Good Tokenizers in Vision Transformer?

Shengju Qian, Yi Zhu, Wenbo Li, Mu Li, Jiaya Jia, *Fellow, IEEE*

Abstract—The architecture of transformers, which recently witness booming applications in vision tasks, has pivoted against the widespread convolutional paradigm. Relying on the tokenization process that splits inputs into multiple tokens, transformers are capable of extracting their pairwise relationships using self-attention. While being the stemming building block of transformers, what makes for a good tokenizer has not been well understood in computer vision. In this work, we investigate this uncharted problem from an information trade-off perspective. In addition to unifying and understanding existing structural modifications, our derivation leads to better design strategies for vision tokenizers. The proposed Modulation across Tokens (MoTo) incorporates inter-token modeling capability through normalization. Furthermore, a regularization objective TokenProp is embraced in the standard training regime. Through extensive experiments on various transformer architectures, we observe both improved performance and intriguing properties of these two plug-and-play designs with negligible computational overhead. These observations further indicate the importance of the commonly-omitted designs of tokenizers in vision transformer.

Index Terms—Vision Transformer, Tokenization, Representation Learning.

1 INTRODUCTION

Serving as a prevalent model in natural language processing (NLP), advances in transformers have driven progress in computer vision. With a great leap made by Vision Transformer (ViT) [1], we have witnessed an increasing effort on adapting this dominating language paradigm to vision. Meanwhile, performance on many downstream tasks such as video recognition [2], [3], object detection [4], [5], and semantic segmentation [6], [7], [8] have upgraded considerably, suggesting the potential of transformer as a primary backbone for vision applications.

While convolution excels at capturing local interactions with its inherent inductive bias such as translation invariance and local connectivity, self-attention in transformers introduces appealing advantages of long-range context modeling and parameter efficiency. In contrast to the feature extraction pipeline in convolutional neural networks (CNNs), we identify the forward process of transformers as two separate stages: the tokenization head that splits inputs into multiple tokens, and the follow-up transformer body that models the pair-wise correlations among the obtained tokens. This reliance on specific tokenization method possibly introduces a bottleneck into both language and vision transformers that limits their capabilities, as discussed in [9], [10], [11].

Being a fundamental prompt in NLP, tokenization strategies have evolved rapidly from the standard rigid tokenization [12]. To cope with variation in language, probabilistic segmentation algorithms like subword regularization [13] have been proposed. Character-level tokenizations lately emerge as a powerful solution in

dealing with languages without whitespace separation [9], [14]. Compared with language sequences, natural images are more complex and have no concise grammar. However, the naive patchify strategy adopted in ViT splits images into non-overlapping patches during tokenization, which are then fed into transformer blocks. Compared to the more continuous and natural distribution of pixels, tokenization and self-attention performed on the discontinuous token embeddings are coarse-grained and may hinder the modeling power of transformers. Recent works alternatively exploited a convolutional stem [11], [15] or an overlapping patch embedding [16], [17] to replace the naive tokenization in ViT [1]. While existing works commonly offer system-level development, the influence from different tokenizer designs still lacks principled discussion and analysis. Meanwhile, the training instability and data-inefficiency still exist as major curse for vision transformer. As discussed in recent studies [11], [18], the naive tokenization stem possibly accounts for these drawbacks.

In this work, we provide an alternative perspective on designing good vision tokenizers. We first conjecture that a good vision tokenizer serves as the information bridge for follow-up transformer blocks, while connects the representations from input images to split tokens. Then we demonstrate that existing structural modifications fit in this formulation. In addition to offering and defending this holistic understanding that bridges different designs of tokenizers, we further exploit the aforementioned assumption, and conduct in-depth analysis of a good vision tokenizer towards its better normalization and optimization objectives. To summarize, our contributions are:

- S. Qian, W. Li, and J. Jia are with the Department of Computer Science and Engineering, The Chinese University of Hong Kong. E-mail: sjqian@cse.cuhk.edu.hk
- Y. Zhu and M. Li are with Amazon AI.

Manuscript received April 19, 2005; revised August 26, 2015.

- We investigate the vision tokenizer design from a information trade-off perspective, and analyze existing design choices from this viewpoint. We also demonstrate that unlike language, naively increasing

mutual diversity across tokens doesn't guarantee better performance.

- We tailor a novel regional normalization strategy for tokenizer, which models better inter-token and intra-token interactions in input images.
- Inspired by the findings, we further incorporate an additional objective that regularizes the original optimization process of tokenizers, which leads to better performance and adaptability to simpler training recipes.
- We evaluate our strategy on various transformer families. Experiments show that our simple yet effective design boosts sophisticated architectures with negligible overhead. The comparisons also demonstrate that vision tokenizer, while lacks in-depth analysis, should be featured prominently in the transformer pipeline.

2 RELATED WORKS

2.1 Vision Transformers.

The transformer architecture was introduced in machine translation [19], and gradually became a primary model for various NLP tasks. The ViT [1] structure closes the gap with CNN models on image recognition, and further demonstrates data efficiency in DeiT [20] with a distillation token and stronger augmentations [21], [22]. In order to obtain suitable architectures that generalize better to vision, researchers have made progress that introduces more hand-crafted designs [16], [17], [23], [24]. Recent structural modifications include improving the positional encoding as in [25], [26], [27], reducing the computational burden via either a local self-attention [28], [29] or a refined global self-attention [15], [30], [31], [32]. Meanwhile, the feature redundancy in vision transformers is observed, which leads to diversified strategies like adaptive length [33] and redundancy reduction [34], [35].

2.2 Tokenization in Language and Vision.

Tokenization serves as a stemming stage in transformers, for both language and vision. Being a long-standing linguistics problem, recent tokenizers such as FRAGE [36] and CharacterBERT [9], [10] have demonstrated superiority in deep transformers when modeling languages over the widespread BPE [12] and WordPiece [37]. The heavy reliance on proper tokenizers for vision transformers has also been observed in the training instability of MoCov3 [18] and ViT stem [11]. Further evidence appears in the redundancy among visual tokens [34], [35], [38]. Therefore, a convolutional stem [11] and overlapped token embeddings [17], [24], [39] are explored to replace the naive patchify tokenization. Meanwhile, PnP-DETR [40], TokenLearner [41], PS-ViT [42], and Token Labeling [43] also explore sampling and learning strategies. While these system-level designs are achieving better downstream performance and suggesting the importance of tokens, there lacks insightful discussions beyond structures for vision tokenizers.

3 PRELIMINARY

3.1 Mutual Information between Random Variables

Mutual information initially measures dependencies between random variables. Given \mathbf{A} and \mathbf{B} , $I(\mathbf{A}; \mathbf{B})$ estimates the "amount of information" learned from \mathbf{B} about the other variable \mathbf{A} and vice versa. The mutual information can be formally defined as:

$$I(\mathbf{A}; \mathbf{B}) = H(\mathbf{A}) - H(\mathbf{A}|\mathbf{B}) = H(\mathbf{B}) - H(\mathbf{B}|\mathbf{A}) \quad (1)$$

Let $\mathbf{A} \in \mathcal{R}^{h \times w \times 3}$, $\mathbf{B} \in \mathcal{R}^{n_{\text{token}} \times l_{\text{token}}}$ denote an input image and its partitioned tokens, where n_{token} and l_{token} respectively refer to the number and dimension of tokens.

$$\begin{aligned} I(\mathbf{A}; \mathbf{B}) &= H(\mathbf{A}) - H(\mathbf{A}|\mathbf{B}) \\ &\geq H(\mathbf{A}) - \mathbb{R}(\mathbf{A}|\mathbf{B}) \end{aligned} \quad (2)$$

where $\mathbb{R}(\mathbf{A}|\mathbf{B})$ denotes the expected error for reconstructing \mathbf{A} from \mathbf{B} . $H(\mathbf{A})$ refers to \mathbf{A} 's marginal entropy, which is treated as a constant in this formulation [44], [45], [46].

In this work, we strive to understand whether vision tokenization fits in this information bottleneck [47], [48], [49] principle, which disentangles vision transformers into two distinct phases of representation learning.

3.2 Estimating Potential Bottleneck in Tokenizers with Conditional Entropy

For vision transformer, the tokenizer divides input image to multiple patch tokens and the cascaded transformer blocks perform global context modeling across the tokens. Compare to the "continuous" representation of image, self-attention performed can be regarded as subsampled operations on its discontinuous token embeddings. Analogous to linguistic analysis [50], we attempt to testify whether vision tokenizer poses an information bottleneck.

Different from the mathematical definition of MI, the bottleneck in vision tokenizers denotes the post-tokenizer information accessibility. Empirically, we characterize this bottleneck as the information lost in tokenization. Therefore, we estimate $I(\mathbf{A}; \mathbf{B})$ by a parametric decoder that obtains the minimal $\mathbb{R}(\mathbf{A}|\mathbf{B})$ between images and their respective tokens. In practice, we adopt a simple three-layer decoder and optimize it using L_2 loss towards the inputs following [46], [53]. The final reconstruction error reflects the empirical conditional entropy and information accessibility between the input image and its split tokens.

4 STRUCTURAL DESIGNS FOR VISION TOKENIZER

Due to the quadratic complexity of self-attention operation, naive extension of "per-word" tokenization in NLP doesn't scale to image pixels. A widespread practice is to process each input image into a sequence of non-overlapping patches as in ViT [1]. Despite alleviating computational costs, the "patchify" operation causes other problems such as training instability [11], [18] and aliasing [54]. To mitigate the drawbacks, architectural modifications have been made towards more suitable vision tokenizers [11], [17], [52]. To facilitate understanding, we propose to summarize existing structure-level improvements to three major aspects:

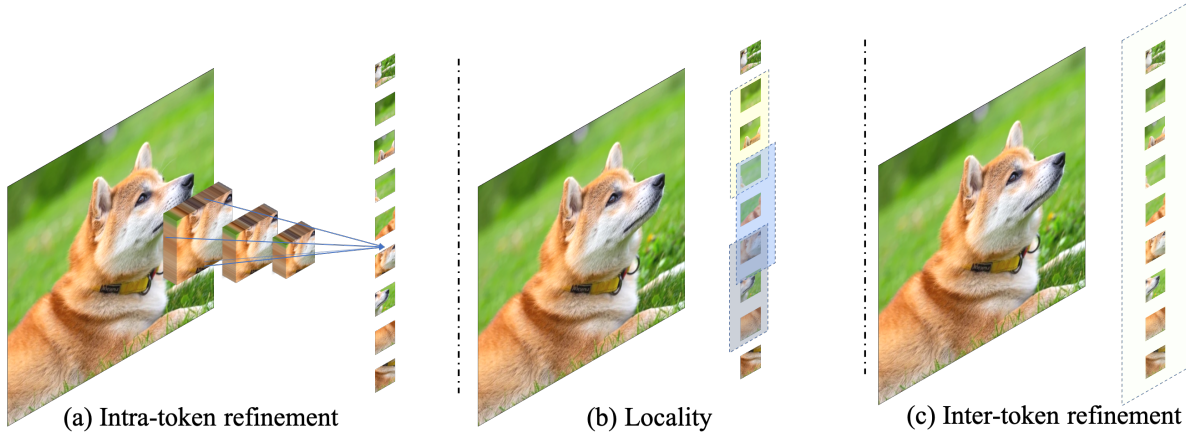


Fig. 1. **Illustration on three major aspects of structural improvements for vision tokenizers.** Note that we make these modifications to exemplify our perspective, which is not designated for superior performance or novelty. For intra-token refinement, we adopt the multi-scale feature extraction as in [51], [52]. We utilize the overlapping embedding as in [17], [39] and self-attention to represent locality and intra-token refinements. The tokens in dotted regions with the same color perform interactions. Locality could be viewed as a local case of Inter-token refinement, where Inter-token refinement emphasize more on global modeling.

- 1) **Intra-token refinement:** As implemented by a stride- p , $p \times p$ convolution, naive “patchifying” fails to capture rich context inside the tokens. It’s straightforward to assign the tokenization encoder with stronger capability of feature extraction. Recent progress that exploited convolutional stem [11], [15], multi-scale embeddings [24], [51], [52], or an additional transformer that models sub-token [55] can be classified in this category.
- 2) **Locality:** The partition strategy splits the continuous image content into relatively divided portions. As pixel statistics of nearby tokens indicate better positions and connectivity, boosting patch tokens with locality also helps. Overlapping [17], [39] and uneven [56] token embeddings have been tailored.
- 3) **Inter-token refinement:** Due to the semantic complexity, similar instances in different images might require distinctive features. Modeling inter-token relationship intuitively helps tokenizers encode better features. This refinement suggests performing global context modeling inside the tokenizers, which has demonstrated its effectiveness in learning [17], [41], [43] or sampling [38], [40], [42] strategies.

While diverse structural refinements have been explored in recent works, MoCov3 [18] empirically shows that naive frozen tokenization enhances contrastive training stability. As such randomly-initialized projection possibly maximize the information accessibility, it’s crucial to validate the existence of potential bottleneck in vision tasks.

4.1 Comparison across Potential Structural Designs

In order to study the corresponding influence, we adopt four modifications from existing transformer frameworks to the original ViT tokenizer, which includes the three categories illustrated in Figure 1.

Experimental Settings. In contrast to the single scale tokenization, the first enhancement exploits a cross-scale

embedding layer which composes of four different kernel sizes used in [51]. The second variant uses the similar overlapping patch embedding with PVTv2 [39]. The third modification further incorporates self-attention that allows for inter-token modeling, as presented in T2T ViT [17]. We also compare with the frozen randomly-initialized variant as in MoCov3 [18]. We compare all variants on three tasks, including supervised classification on ImageNet, linear probing with unsupervised pre-training as in [18], and semantic segmentation on ADE20K [57]. More details including specific designs and experimental settings are provided in Appendix A.2.

4.2 Takeaways from the Empirical Analysis

To illustrate how the information accessibility between image and tokens changes after tokenization, we follow the practice used in [44], [46] and reflect it using reconstruction error. When training completes, we construct an encoder-decoder framework using the pre-trained tokenizer and a random-initialized lightweight decoder. Given the frozen tokenizer, the decoder is then trained to reconstruct the input images in 64×64 using L_2 loss for 10 epochs on ImageNet training set. Note that we find the decoder, albeit with a simple three-layer structure, is easy to optimize with good convergence. In addition to reconstruction, we also measure the token similarity by averaging the cosine similarity between obtained tokens, where lower similarity represents more diversity in the token embedding.

From Table 1 and Figure 2, we perform in-depth analysis and provide several important takeaways for designing better vision tokenizers:

- 1) **Empirical correlations exist between increasing token diversity and better results.** The results in Table 1 verify that three refinement strategies consistently benefit the transformer’s capability. By connecting it with Figure 2, we observe that more diversity and interactions across tokens gradually reduce reconstruction error.

TABLE 1

Performance on various vision tasks with different tokenizers. “Intra”, “Local”, and “Inter” respectively refer to applying the intra-token, locality, and inter-token refinement strategies. “Frozen” refers to the frozen randomly-initialized tokenization in [18].

Architecture	Tokenization				Params	GFLOPs	Classification		Segmentation
	Intra	Local	Inter	Frozen			Linear [18]	Supervised	
DeiT-S	-	-	-	-	22.1	4.6	68.1	79.8	44.0
DeiT-B	-	-	-	-	86.6	17.6	69.6	81.8	45.2
DeiT-S				✓	22.1	4.6	69.0	79.4	42.9
DeiT-S	✓				22.3	4.7	68.4	80.5	44.3
DeiT-S	✓	✓			22.3	5.1	68.5	80.9	44.6
DeiT-S	✓	✓	✓		25.7	6.9	68.7	82.0	45.0

- 2) **Maximizing conditional entropy doesn't guarantee better performance.** While the randomly-initialized frozen tokenizer or raw pixel encoding maximizes conditional entropy and token diversity, it gives inferior performance on vision tasks. Considering the modality of image, certain properties such as multi-scale are absent with these encodings, which is crucial for vision tasks, especially segmentation.
- 3) **Tokenizers might influence optimization across tasks.** Similar to the findings in [18], we also observe better linear performance with frozen tokenizations. However, such trick doesn't apply to supervised classification and segmentation. As has been analyzed in [11], convolutional stem also provides more stable training process.
- 4) **The goal of tokenizer is to maintain a trade-off between feature expressiveness and information accessibility.** Unlike the semantically-structured language representations, images consist of sensory pixels and have lower and imbalanced information density. Such characteristics require the tokenizer to perform feature extraction and “patchifying” simultaneously, instead of eagerly minimizing the information lost as in language.

5 NORMALIZATION IN VISION TOKENIZER

Normalization is vital in network design, which empowers complex architectures and accelerates convergence. Although the importance of normalization has been acknowledged in transformer for both language [58], [59] and vision [60], it has rarely been explored inside the tokenizer. As analyzed in Section 4 that a good tokenizer is critical for transformer, proper normalization is tailored in this section.

Considering the semantic variations between images, patch tokens tend to encounter semantic variation as shown in Figure 3. This large semantic gap across tokens manifests the difficulties in optimizing the tokenizer, which may explain the witnessed training instability of a naive patch embedding [11], [18]. However, the addition of batch normalization (BN) to the patchify stem [11] worsens its accuracy. According to previous analysis [61], [62], normalization layers tend to “wash away” texture and semantic information. Such “filtered” effect inevitably reduces token diversity and possibly cuts off necessary diversity in transformer. Meanwhile, based on our findings in Section 4, naively maximizing token accessibility as in MoCov3 [18] might be suboptimal for different vision tasks.

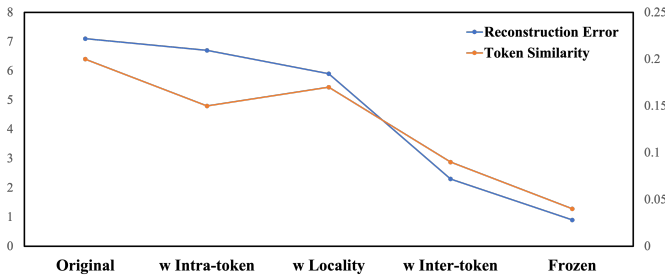


Fig. 2. Reconstruction error and token similarity from different tokenizer designs. Reconstruction error measures the MSE error from decoding the tokens back to images and reflects their conditional entropy. Token similarity averages the cosine similarity between tokens and reflect the diversity of obtained tokens. Both metrics are evaluated on ImageNet validation set.

Motivated by these findings, we expect our tokenizer to simultaneously maintain the information accessibility and feature quality. Therefore, we further study versatile and unexplored design strategies for vision tokenizers, through normalization and optimization.

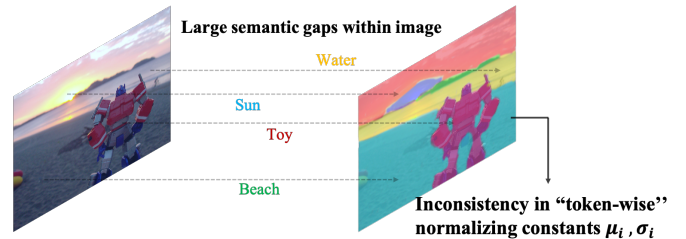


Fig. 3. Example about the large semantic gap across regions and tokens. This difference leads to distinctive statistics. Using fixed normalizing constants would inevitably “filter” the original information and reduces the token diversity.

5.1 Modulation Across Tokens (MoTo)

Therefore, we propose to normalize the input content in a spatial-aware manner called **MoTo**, short for **Modulation across Tokens**. The core idea is to modulate the diverse semantics in input using regional statistics. This strategy not only formulates feature distributions, but also provides the tokenizer with more plausible semantic content. Unlike the spatial normalization used in conditional image generation [62], [63], there exists no given semantic layout

in our pipeline. As shown in Figure 4, MoTo consists of the soft semantic partition and the spatial-aware modulation.

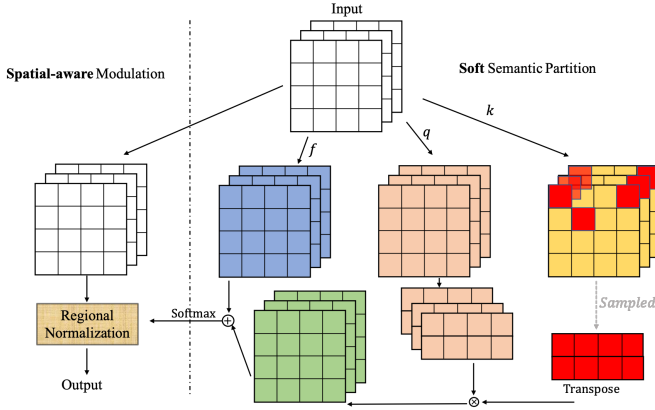


Fig. 4. **Forward illustration of MoTo.** The input first proceeds soft semantic partition to obtain the layout. Then spatial-aware modulation is performed based on the soft layout. The soft semantic partition is shown in the right pattern and the spatial-aware modulation is on the left.

5.1.1 Soft Semantic Partition.

Assume each input image consists of n semantic entities, this module predicts the soft semantic layout $\mathbf{L} \in \mathcal{R}^{h \times w \times n}$ from the input $\mathbf{X} \in \mathcal{R}^{h \times w \times c}$. A convolution layer f with n output channels first extracts the semantic feature $f(\mathbf{X}) \in \mathcal{R}^{h \times w \times n}$. Similarly, two feature extractors k and q obtain $k(\mathbf{X})$ and $q(\mathbf{X})$. Then n feature points $k_n(\mathbf{X})$ are uniformly sampled from $k(\mathbf{X})$ following the practice in [64] to compute the feature correlation. The semantic activation map $\mathbf{Z} \in \mathcal{R}^{h \times w \times n}$ is calculated using matrix multiplication and accumulated with semantic feature

$$\mathbf{Z} = \mathbf{u} \cdot k_n(\mathbf{X})^T q(\mathbf{X}) + f(\mathbf{X}), \quad (3)$$

where $\mathbf{u} \in \mathcal{R}^{1 \times 1 \times n}$ denotes a randomly-initialized learnable dictionary that integrates the correlation with features.

Then \mathbf{Z} is normalized using softmax to generate the soft layout at k -th semantic entity

$$\mathbf{L}_k = \frac{\exp(\tau \mathbf{Z}_k)}{\sum_{i=1}^n \exp(\tau \mathbf{Z}_i)}, \quad (4)$$

where $k \in [0, n - 1]$ and τ is the temperature coefficient set to 0.1. Each $\mathbf{L}_k \in \mathcal{R}^{h \times w}$ indicates the probability of the spatial pixels belonging to entity k .

5.1.2 Spatial-aware Modulation.

As the soft layout provides distributions of potential semantic entities, this component modulates the input with regional normalization. In contrast to instance normalization [61], [65], it models the spatial correlations and treats semantic entities as ‘‘instances’’. Formulating similar semantics with shared means and variances, our method modulates the interactions in input better, as

$$\text{MoTo}(\mathbf{X}) = \sum_{i=1}^n \left(\frac{\mathbf{X} - \mu(\mathbf{X} \odot \mathbf{L}_i)}{\sigma(\mathbf{X} \odot \mathbf{L}_i) + \epsilon} \times \beta_i + \alpha_i \right) \odot \mathbf{L}_i, \quad (5)$$

where $\mu(\cdot)$ and $\sigma(\cdot)$ respectively denote computing the mean and standard deviation from the selected instance. β_i and α_i in Eq. 5 are the learnable parameters for affine transformation in normalization layers. We use a default $n = 8$ in most experiments.

5.2 Experimental Analysis with MoTo

5.2.1 Improvements on various architectures.

To further validate that MoTo is versatile with different tokenizers and transformer architectures, we integrate the proposed module to several state-of-the-art methods, including DeiT [20], T2T-ViT [17], PVT [15], and Swin Transformer [28]. Since these frameworks contain distinctive training pipelines, we utilize their publicly-available scripts [66], [67], [68], [69], [70] and keep the training parameters consistent with the provided ones. When applied with our module, the variants maintain the consistent configurations including augmentation and optimizer as the baselines. We train all models on ImageNet [71] training set using 8 Tesla V100 GPUs for 300 epochs, except for T2T-ViT which originally scheduled 310 epochs. From Table 2, we see that the consistent improvements are made across different transformer architectures. For the wide-adopted baseline DeiT [20], our strategy improves its naive tokenization process. The improvements can be observed on both PVT [15] and PVTv2 [39], where PVTv2 introduces overlapping patch embedding. While T2T-ViT exploits a transformer in tokenizer to perform re-structurization, MoTo is still effective.

Complexity Analysis. MoTo incorporates global context modeling to vision tokenizer using normalization. The channel, height, weight of the feature maps and number of semantic entities are respectively denoted as C, H, W, N . Implemented with convolution layers and sampling, soft semantic partition and spatial-aware modulation cost a time complexity of $\mathcal{O}(NCHW)$. Comparing to the $\mathcal{O}(C(HW)^2)$ complexity of self-attention module, MoTo provides a computational-friendly choice to perform inter-token modulation. As shown in Table 3, the actual computation of injecting pixel-wise regional information with self-attention scales quadratically, which is unfeasible for processing high-resolution images.

5.2.2 Visualization of semantic layout.

To understand the effectiveness of the semantic partition, we visualize the obtained semantic entities using a n -color palette. Given its soft semantic vector $\mathbf{L}_{i,j} \in \mathcal{R}^{1 \times 1 \times k}$, we colorize each location (i, j) using the index with the highest probability. As shown in Figure 5, the semantic entities highlight the difference between foreground, background, and instances. Notably, the entities represent some details on cat face in Figure 5.(a), which switch into instances such as human and mountains in Figure 5.(b). The results demonstrate the effectiveness of our module in capturing semantic variations in images. Normalizing these feature points with semantic grouping improves the feature quality of tokens.

TABLE 2

Performance of image recognition on ImageNet validation set with MoTo. The column of “Params” denotes the number of parameter, “Acc” reports Top-1 accuracy. All experiments and GFLOPs computations use the input size of 224×224 .

Transformer Architecture	Model	Params	GFLOPs	Accuracy	Δ
ViT [1], [20]	DeiT-S	22.1M	4.6	79.8	-
	w MoTo	22.5M	4.8	81.6	+1.8
	DeiT-B	86.6M	17.6	81.8	-
	w MoTo	86.9M	17.9	82.9	+1.1
T2T-ViT [17]	T2T-ViT-14	21.5M	5.2	81.5	-
	w MoTo	21.8M	5.4	82.3	+0.8
PVT [15], [39]	PVT-Small	24.5M	3.8	79.8	-
	w MoTo	24.7M	4.0	81.0	+1.2
	PVT-Medium	44.2M	6.7	81.2	-
	w MoTo	44.5M	6.9	82.1	+0.9
	PVTv2-B2	25.4M	4.0	82.0	-
	w MoTo	25.6M	4.2	82.7	+0.7
Swin [28]	Swin-T	28.3M	4.5	81.2	-
	w MoTo	28.6M	4.7	82.2	+1.0
	Swin-S	49.6M	8.7	83.0	-
	w MoTo	49.9M	8.9	83.7	+0.7

TABLE 3

Inference wall time (ms) with different input scales. Features are fed into the same GPU with a batch size of 1 and channel number of 16. OOM denotes out-of-memory. Note that the definition of self-attention here pixel refers to the pixel-wise self-attention in tokenizer.

Component	224×224	384×384	512×512	1024×1024
Self-attention	4.12	23.58	OOM	OOM
MoTo	0.52	1.67	4.23	18.42

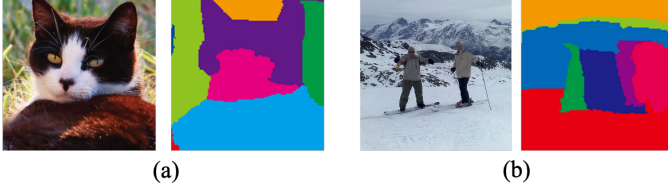


Fig. 5. **Visualizations of the self-learned semantic layout.** Both image (a) and (b) are testing images. Each color shown in the layout denotes a semantic entity in the soft layout L.

5.2.3 Ablations about MoTo.

Number of semantic entities. The hyper-parameter n number of the semantic entities determines our semantic

TABLE 4

Ablation study about entity numbers and partition strategy. The baseline architecture uses DeiT-S. Each row represents a model trained with different number of semantic entities or partition strategy. The gray rows refer to models with hard partition.

Semantic Entities	Partition Strategy	Top-1 Accuracy	
		Val	Δ
-	-	79.8	-
$n = 2$	Soft	80.8	+1.0
$n = 4$	Soft	81.1	+1.3
$n = 8$	Soft	81.6	+1.8
$n = 16$	Soft	81.7	+1.9
$n = 16$	Hard	80.5	+0.5
$n = 32$	Hard	80.8	+0.8

partition process. We gradually increase n to study the influence of varying quantities of semantic entities. When the number is small, Table 4 shows a trend of improved accuracy with more entities. However, growing grouping numbers bring relatively marginal improvements when $n > 8$. As the quantity controls the fineness of semantic modeling, increasing entities contribute to stronger context modeling while inevitably bring redundancy and difficulties in optimization, unable to guarantee better results. We thus choose $n = 8$ for most experiments. considering the trade-off between effectiveness and complexity.

Soft and hard partition strategy. To understand our partition strategy, we also compare the difference between soft and hard partition. In contrast to the soft probability layout in Eq. 4, we implement $\arg \max$ on the activation maps spatially and further obtain a class map to perform hard partition, which resembles the visualization process in Figure 5. As shown in Table 4, hard partition performs much worse than soft partition even with more semantic entities, exhibiting the superiority of soft partition.

Absorbing MoTo into transformer blocks. While we are mainly discussing replacing the original patch embedding layer with better tokenization, it remains unexplored to absorb such tokenization into the transformer blocks. We further conduct a pilot study to ensemble MoTo at the end of each transformer layer in Table 5. Despite with clear improvements with MoTo as tokenization and increased complexity, fusion into transformer blocks gives limited elevations. We can also see the importance of tokenizer-level normalization, apart from its follow-up transformer blocks. The results not only suggest the prominence of vision tokenization, but also elaborate more towards the information flow of transformer.

5.2.4 Discussions on normalization.

Image synthesis [61], [62], [72] and domain adaptation [73], [74], [75], where large domain gap and variations exist, recently find normalizations benefit the networks in encoding rich content. Motivated by these findings, our

TABLE 5

Ablation study on absorbing MoTo into transformer blocks. The baseline architecture uses DeiT-S. We adopt a MoTo layer with 8 semantic entities and ensemble it into transformer blocks. The placement denotes the layer number of transformer block that adopts MoTo.

Tokenizer	Placement			Top-1 Accuracy	
	1-4	4-8	8-12	Val	Δ
-	-	-	-	79.8	-
✓				81.6	+1.8
✓	✓			81.8	+2.0
✓	✓	✓		81.9	+2.1
✓	✓	✓	✓	81.4	+1.6

method investigates normalization inside vision tokenizer and supplements tokenization process with capability of context modeling. The clear boost over multiple structures brought by this lightweight design further demonstrates the importance of a good tokenizer.

In contrast to normalization layers inside transformer blocks, our work investigates the normalization inside vision tokenizers. We emphasize that these two directions are inherently different and both crucial, while the modulation in tokenizers has been rarely explored.

Normalization in Transformer. Previous works have demonstrated the importance of normalizations in transformers [19], [58], where Layer Normalization [76] plays a key role in their success. With further developments in normalization such as PowerNorm [59], transformers on language modeling have received a performance boost. As the paradigm shift becomes success in computer vision, analysis and attempts on normalizations in vision transformers have been made in [60], [77].

Normalization in Tokenizer. Unlike the normalizations conducted inside transformers which are capable of maintaining the attention magnitudes, normalizations in tokenizers perform a different role of token feature extraction. In language pipelines, normalizations in tokenizers refer to the operations that make a raw string cleaner, including traditional Unicode normalization and BERT normalizer [78]. Comparing to highly structured languages, images contain more complex semantics and diversified pixel variations. Thus tokenization on images is more challenging and might require additional regularization. However, normalizations in vision tokenizers have been rarely explored. Our proposed MoTo targets at fixing this omitted ingredient in vision tokenizer.

Based on the different roles between normalization in tokenizers and transformers, their design strategies should be made different accordingly. Considering the semantic variations across images and tokens, MoTo normalizes the input content in a spatial-aware manner that modulates the regional variations while not “wash away” the diverse texture and semantic content of inputs. Based on the analysis in our main paper, MoTo is capable of maintaining the information accessibility and diversity between input images and extracted tokens. With minimal computational costs through modulations, MoTo also incorporates global information into the tokenization process.

Comparison with different standard normalizations. To better illustrate the effectiveness of MoTo, we present more ablative analysis by comparing MoTo to other choices of normalizations in tokenizers. We use DeiT-S architecture as our baseline and follow the experiment settings in Section 5.2. We choose Layer Normalization (LN), Batch Normalization (BN) and instance normalization (IN) as the counterparts of MoTo.

TABLE 6

Comparison of different normalizations in tokenizer. The baseline architecture uses DeiT-S. Each row denotes result trained with different normalization strategy. Top-1 Accuracy denotes the validation accuracy on ImageNet.

Model	Normalization in Tokenizer	Top-1 Accuracy
DeiT-S	-	79.8
DeiT-S	Layer Norm	79.6
DeiT-S	Batch Norm	79.5
DeiT-S	Instance Norm	80.1
DeiT-S	MoTo	81.6

As shown in Table 6, both BN and LN slightly harm the model’s performance. As analyzed, such normalization methods bring “filtered” effect on tokens’ features. In comparison with IN, MoTo performs soft semantic partition and incorporates spatial-aware property into the instance-wise modulation. The results further indicate that different from normalization in transformer blocks, normalization in tokenizers need to take both “intra-token” and “inter-token” modeling into considerations.

6 OBJECTIVES FOR VISION TOKENIZER

From the information trade-off perspective, both structural modifications in Section 4 and MoTo in Section 5 can be considered as refinements that assign the tokenizer with stronger capability in preserving the conditional entropy between image and token representations. In addition to network-level analysis, we explore optimization-level refinement for vision tokenizers in this section.

6.1 The “greedy” training paradigm.

Recent studies [79], [80] have observed a “greedy” trend in the models trained with supervised signals. The network is optimized to encode task-relevant features under the guidance of a loss function, e.g. cross-entropy loss for classification. As the latent features become more discriminative, the conditional entropy between inputs and layer-wise features gradually decreases during training according to the estimation in [79]. This “greedy” characteristic in locally-supervised learning pipeline often collapses task-relevant feature in earlier layers, which further leads to inferior performance. Although not pushed by additional supervision, the tokenizer inevitably reduces information accessibility between the input and tokens during training. Considering the difficulty of vision tokenization and witnessed instability [18], this “greedy” phenomenon possibly exists in vision tokenizers, washing out some valuable information. With similar observations in language understanding [50], [81], a potential solution

is to modify the objectives to regularize the tokenizer from being “short-sighted”.

The objective of Masked language modeling (MLM) in BERT [78] has been demonstrated with its property in guiding the evolution of representations of individual tokens proceed in two stages, including context modeling and token reconstruction [50]. InfoWord [81] further analyzes how MLM differs from InfoNCE [82] and incorporates negative sampling in measuring more concise mutual information, hereby learning better language representations.

Therefore, we attempt to design a less “greedy” training regime, where the *head*, vision tokenizer is able to preserve more context from inputs that can potentially leverage by later *body*, transformer blocks.

6.2 TokenProp Objective

While performing different roles in the information flow of transformers, the vision tokenizer and transformer blocks have been treated equally during training. The core idea of TokenProp is to provide optimization objectives for vision tokenizer, which lead to better token representations by maintaining the trade-off between feature fineness and information accessibility.

Following the notations in Section 4 where **A** and **B** refer to the input image and tokens, we show that the conditional entropy between them could be maximized by optimizing $\mathbb{R}(\mathbf{A}|\mathbf{B})$

$$\begin{aligned} \arg \max I(\mathbf{A}; \mathbf{B}) &= \arg \max -\mathbb{R}(\mathbf{A}|\mathbf{B}) \\ &= \arg \max \mathbb{E}_{q(\mathbf{A}, \mathbf{B})} [\log q(\mathbf{A}|\mathbf{B})] \end{aligned} \quad (6)$$

where we estimate the parameters of $q(\mathbf{A}|\mathbf{B})$ using a decoder network G_θ . Ideally, zero information loss through the tokenization process retains all information including the useful one.

However, solely optimizing the conditional entropy collapses into another extreme case, where the tokenizer merely performs spatial partition and doesn’t capture any task-specific feature. Therefore, we utilize a locally-supervised paradigm that jointly optimizes the target loss and the conditional entropy during tokenization. The surrogate optimization objective could be defined as

$$\underset{\phi, \omega, \theta}{\text{minimize}} \mathcal{L}(F_\omega(F_\phi(\mathbf{x})); \mathbf{y}) + \lambda \mathcal{L}_{rec}(G_\theta(F_\phi(\mathbf{x})); \mathbf{x}) \quad (7)$$

where \mathcal{L}_{rec} and \mathcal{L} represent the reconstruction loss and the standard task loss, e.g. cross entropy loss for classification in our case. F_ϕ and F_ω denote the tokenizer and the follow-up transformer architecture. $\mathbf{x}, \mathbf{y}, \phi, \omega, \theta$ refer to the input, label, and parameters of respective modules.

6.3 Experiments with TokenProp

Implementation Details. We adopt a lightweight decoder structure to reconstruct the input from tokens. We find that a simple three-layer decoder works well in most cases. The decoder first transforms the split tokens into connected spatial feature maps, which are further refined using convolutional layers and upsampled

to 64×64 . This implementation introduces negligible computational overhead to the transformer pipeline. The default hyper-parameter λ , which combines standard loss and reconstruction loss, is set 0.001 in our experiments. All the experiments are trained for 300 epochs on 8 GPUs if not additionally specified.

Choices of \mathcal{L}_{rec} . The reconstruction loss \mathcal{L}_{rec} evaluates the conditional entropy between image and token representations. As there exists multiple choices to compute the pixel-wise distance, we perform a comparison in Table 7. Apart from the commonly used L_1 and L_2 , we also adopt the perceptual loss [83] and contextual loss [84]. Interestingly, the L_2 distance shows a competitive result over other alternatives. While perceptual loss excels at capturing style information and better visual quality [83], [85], it doesn’t perform well as expected. Contextual loss, which is effective at maintaining pixel statistics, boosts slightly more than L_2 . Note that both perceptual loss and contextual loss exploit a pre-trained VGG-19, which might bring information leakage into training. Considering the complexity of contextual loss, we use L_2 distance in our experiments. Meanwhile, the results indicate that better reconstructed quality doesn’t necessarily mean higher accuracy.

TABLE 7
Ablations about \mathcal{L}_{rec} . The baseline architecture uses DeiT-S. Each row denotes result trained with different \mathcal{L}_{rec} .

Model	Loss Type	Top-1 Accuracy	
		Val	Δ
Baseline	-	79.8	-
-	L_1	80.2	+0.4
-	L_2	80.7	+0.9
-	Perceptual Loss [83]	80.4	+0.6
-	Contextual Loss [84]	80.8	+1.0

Reconstruction loss weight of λ . As the target features for reconstruction and classification are completely different, the reconstruction loss weight λ servers as another crucial hyper-parameters. If the model focuses too much on reconstruction, the final classification performance will inevitably be affected. We perform a study on how a different λ influences our model. We adopt different loss weights of 0.001, 0.01, 0.1 and 1.0 and report the respective performance of these four variants. The training recipes are kept the same as the initial one. From Table 8, we can find that our method produces similar improvements when given different reconstruction weights at a reasonable range. If the weight is set too large, instability will be observed in training. The results also demonstrate the robustness of TokenProp and indicates the potential linkage between generative and discriminative signals.

Decoder Structure. We also study whether a complex decoder is favoured in our setting. By enlarging the convolutional channel for n times, we denote the decoder variant as $\times n$. We also compute reconstruction loss using outputs with higher resolutions, which are obtained by stacking more upsampling layers to the decoder. The training overhead is measured on the same 8 GPUs by comparing the final 300 epochs’ training time with the baseline trained without TokenProp. While more

TABLE 8

Ablations about the reconstruction loss weights λ . The baseline architecture uses DeiT-S. Each row denotes result trained with different λ . NaN means the model faces NaN in training loss.

Model	Loss Weight λ	Top-1 Accuracy	
		Val	Δ
Baseline	-	79.8	-
-	0.001	80.7	+0.9
-	0.01	80.4	+0.6
-	0.1	80.6	+0.8
-	1.0	NaN	-

complex decoders introduce extra overhead in Table 9, we see no clear improvement with larger decoders. Similar with sophisticated designs, reconstructions with higher resolutions tend to preserve more details. However, a performance drop is observed when we use the output size of 256×256 , which emphasizes the importance of a proper objective. These comparisons, as well as findings from Table 7, also indicate that a simple decoder is sufficient for learning our objective, which is not designated for generating high quality outputs.

TABLE 9

Analysis about the decoder structures. The baseline architecture uses DeiT-S. The sign of \downarrow denotes the accuracy is lower than the baseline trained without TokenProp.

Decoder Channel	Output Scale	Accuracy	Training Overhead
$\times 1$	64×64	80.7	4.1%
$\times 1$	128×128	80.4	7.8%
$\times 1$	256×256	77.3 (\downarrow)	15.1%
$\times 2$	64×64	80.4	6.9%
$\times 4$	64×64	80.5	9.1%

6.3.1 Improvements on Various Architectures.

In order to validate the generalization of TokenProp, we apply it to various transformer architectures. We maintain the original training recipes of these state-of-the-arts as mentioned in Section 5.2. Except for the training-only decoder, our method incurs no extra parameter or computation to the original framework. Therefore, these models receive negligible overhead to compute our TokenProp objective with a light-weight decoder during training. As shown in Table 10, different models in four families of transformer architectures receive a steady improvement on accuracy, with roughly 5% training overhead. From another perspective, TokenProp serves as a regularization term that provides the encoded features from tokenizers with more diversity.

6.3.2 Compatibility with Optimizers.

Instability is another major curse for vision transformers. AdamW [86] optimizer has been dominant in mitigating the training difficulty of transformers. Recent findings ascribe transformers’ vulnerability to their tokenizers, and try to provide a fix, such as the convolutional stem in [11] and the random patch projection [18] for self-supervised learning. While these modifications make transformers less sensitive towards training recipes, the AdamW counterparts still exhibit considerable superiority over the ones trained with simpler optimizer such as SGD [11], [20]. To testify the

influence of TokenProp on optimizations, we propose to validate our method using SGD in Table 11. We also re-implement a DeiT-S variant with the convolutional stem in [11], as well as its combination with TokenProp. Since the frozen randomly-initialized embedding [18] has been adopted with better stability, we also include such variant with replaced optimizers. We tune the optimal learning rate and weight decay for each variant through multiple runs, while other training recipes such as augmentation are kept consistent. As shown in Table 11, TokenProp enables the adaptability of the model towards SGD, which significantly reduces the performance drop. Notably, by incorporating TokenProp with the convolutional stem in [11], the variant only loses 0.2% validation accuracy. These results further suggest that an optimized objective provides more reliable supervisory signals, which leads to both better performance and stability. Comparing to the additional memory costs from Adam-based methods, TokenProp only incurs limited computation and memory overhead during training.

Comparison with frozen embedding [18]. Although the practice in [18] observes with better stability in self-supervised learning, we observe inferior performance on supervised classification and downstream tasks in Section 4, as well as heavy reliance on optimizers. We believe there exists a connection between the frozen embeddings and our TokenProp objective, as both methods target to maintain the stability in the tokenization process. Nevertheless, feature fineness is neglected in such raw projection, which accounts for its poor performance on downstream tasks. In contrast, TokenProp serves as a surrogate objective that balances optimization stability and representation quality.

7 DISCUSSIONS.

To summarize, we propose two plug-and-play designs in vision tokenizers: Modulation across Tokens (MoTo) that incorporates inter-token modeling capability through normalization, and a surrogate optimization objective TokenProp. In this section, we provide experimental discussions to further demonstrate the effectiveness of our method.

7.1 Ensembles of MoTo and TokenProp.

Since we’ve demonstrated that both MoTo and TokenProp are versatile with different transformer architectures, we also validate whether they benefit from each other as well. In Table 12, we ensemble both strategies into different frameworks and follow the aforementioned training and evaluation recipes. We can see that the models could be further boosted, suggesting the compatibility of structural and objective designs.

7.2 Improvements on Data Efficiency

The data inefficiency has been a major problem of vision transformers. Since the original training regime of ViT [1] contains hundreds of millions of training images, efforts have been made to improve its data efficiency that allows training feasibility. Therefore, we investigate whether our proposed strategies benefit the model from better utilization of training data. We incorporate both our strategies to

TABLE 10

Performance of image recognition on ImageNet validation set with TokenProp. The column of “ Δ ” reports the improvements from the models trained with TokenProp.

Transformer Architecture	Model	Accuracy	Δ	Training Overhead
ViT [1], [20]	DeiT-S	80.7	+0.9	4.1%
	DeiT-B	82.5	+0.7	3.5%
T2T-ViT [17]	T2T-ViT-14	82.2	+0.7	4.9%
	T2T-ViT-19	82.8	+0.9	4.4%
PVT [15], [39]	PVT-Small	80.9	+1.1	5.1%
	PVT-Medium	81.8	+0.6	4.6%
Swin [28]	Swin-T	82.3	+1.1	6.7%

TABLE 11

Analysis about the compatibility with optimizers. The sign of \downarrow shows how much the accuracy is lower than the baseline trained using AdamW. We re-implement DeiT with [11] and [18], as denoted by DeiT_C* and w Frozen. We highlight top-2 results in bold font.

Model Variants	Optimizer	Top-1 Accuracy	
		Val	Δ
DeiT-S	AdamW	79.8	-
DeiT-S	SGD	76.7	\downarrow 3.1
DeiT-S w Frozen [18]	AdamW	79.4	-
DeiT-S w Frozen [18]	SGD	76.0	\downarrow 3.4
DeiT _C -S* [11]	SGD	78.2	\downarrow 1.6
DeiT-S w TokenProp	SGD	78.9	\downarrow 0.9
DeiT_C-S* w TokenProp	SGD	79.6	\downarrow 0.2

TABLE 12

Combinations of MoTo and TokenProp.

Model Architecture	w MoTO	w TokenProp	Accuracy
DeiT-S	-	-	79.8
Ours	\checkmark	\checkmark	82.6
T2T-ViT-14	-	-	81.5
Ours	\checkmark	\checkmark	82.8
Swin-T	-	-	81.2
Ours	\checkmark	\checkmark	82.9

another dominating transformer of Swin transformer [28], which has already demonstrated strong capability of data exploitation.

We follow the original training hyper-parameters in [28], while only leverage a limited portion of ImageNet-1k training set and 100 epochs for training. We train the baseline model under a specific amount of training data, as well as the boosted version with our modules. The

TABLE 13

Performance when the model is trained under the resource-limited setting. The baseline structure uses Swin-T [28] architecture. Each row represents the results trained under a certain portion (percentage) of the original ImageNet-1k training data, for both the baseline and our refined counterpart. Top-1 Accuracy denotes the validation accuracy on original ImageNet validation set.

Training Dataset	Percentage (%)	Val Top-1 Accuracy		
		Baseline	w Ours	Δ
ImageNet-1k	50	73.7	75.1	+1.4%
	40	71.6	73.2	+1.6%
	20	61.2	63.9	+2.7%
	10	43.5	46.5	+2.9%

results under this resource-limited setting are reported in Table 13, where a consistent boost is observed under each training portion. The results show that our refinement strategies on vision tokenizer are able to improve the data efficiency, which are not specially designed for this purpose. Meanwhile, better preservation of information accessibility in the extracted tokens builds up their connection and helps the transformer capture more valuable content.

7.3 Improvements on Downstream Tasks

To further validate the effectiveness of the proposed strategies, we perform evaluations with our modules on downstream tasks including semantic segmentation and object detection. For these tasks, we utilize DeiT-S [20], PVT [15], and Swin-T [28] as the baseline architecture, where our variants include the proposed MoTo and TokenProp.

7.3.1 Semantic Segmentation

We evaluate our strategies on semantic segmentation using the ADE20K [57] dataset. Details regarding the dataset and hyper-parameters can be found in Appendix B.3. Similar to the experiments on classification, we validate separately with their respective performance, as well as their ensembles. As shown in Table 14, our method easily improves the model with Swin-T backbone that is already very strong comparing to previous methods. Similar with TokenProp, the randomly-initialized embedding in MoCov3 [18] demonstrates better stability in self-supervised learning. To better understand this with our findings in Section 4, we also implement it in semantic segmentation and compare with TokenProp. The frozen encoding lacks feature expressiveness, which is more crucial for downstream tasks like segmentation. The results further suggest the importance of proper tokenizer structure, as well as the optimization trade-off between feature quality and accessibility.

7.3.2 Object Detection and Instance Segmentation

We also perform evaluations on object detection and instance segmentation using the COCO 2017 dataset. Based on the implementation in MMDetection [91], we testify the performance on different detection pipelines, including RetinaNet [88], Mask R-CNN [89], and Cascade Mask R-CNN [90]. We also compare the performance with different backbones, where the hyper-parameters are provided in Appendix B.4.

TABLE 14

Downstream performance of semantic segmentation on ADE20K dataset. We modify the tokenizer of DeiT-S and Swin-T with our proposed modules and denote it in bold font. The reported mIoU exploits multi-scale and flip testing.

Method	Backbone	Module		Pre-trained	Crop Size	LR Schedule	mIoU
		MoTo	TokenProp				
OCRNet [87]	HRNet-w48			ImageNet-1k	512 × 512	150K	45.7
UperNet	DeiT-S			ImageNet-1k	512 × 512	160K	44.0
	DeiT-S w Frozen [18]			ImageNet-1k	512 × 512	160K	42.9
	DeiT-S	✓		ImageNet-1k	512 × 512	160K	44.5
	DeiT-S		✓	ImageNet-1k	512 × 512	160K	44.3
	DeiT-S	✓	✓	ImageNet-1k	512 × 512	160K	44.7
UperNet	Swin-T			ImageNet-1k	512 × 512	160K	45.8
	Swin-T	✓		ImageNet-1k	512 × 512	160K	46.3
	Swin-T		✓	ImageNet-1k	512 × 512	160K	46.4
	Swin-T	✓	✓	ImageNet-1k	512 × 512	160K	46.7
UperNet	Swin-S			ImageNet-1k	512 × 512	160K	49.1
	Swin-S	✓		ImageNet-1k	512 × 512	160K	49.4
	Swin-S		✓	ImageNet-1k	512 × 512	160K	49.4
	Swin-S	✓	✓	ImageNet-1k	512 × 512	160K	49.6

TABLE 15

Downstream performance of object detection and instance segmentation on COCO 2017. We modify the tokenizer of PVT and Swin-T with our proposed modules and denote it in bold font. The backbones are pre-trained on ImageNet-1k dataset.

Framework	Backbone	Pre-trained	LR Schedule	Box mAP	Mask mAP
RetineNet [88]	PVTv2-b1	ImageNet-1k	1x	41.2	-
	PVTv2-b1 w MoTo	ImageNet-1k	1x	41.8	-
	PVTv2-b1 w TokenProp	ImageNet-1k	1x	41.5	-
	PVTv2-b1 w Both	ImageNet-1k	1x	42.0	-
Mask R-CNN [89]	PVTv2-b1	ImageNet-1k	1x	41.8	38.8
	PVTv2-b1 w MoTo	ImageNet-1k	1x	42.4	39.1
	PVTv2-b1 w TokenProp	ImageNet-1k	1x	42.3	39.4
	PVTv2-b1 w Both	ImageNet-1k	1x	42.9	39.4
Mask R-CNN [89]	Swin-T	ImageNet-1k	1x	43.7	39.8
	Swin-T w MoTo	ImageNet-1k	1x	44.1	40.1
	Swin-T w TokenProp	ImageNet-1k	1x	44.2	40.0
	Swin-T w Both	ImageNet-1k	1x	44.6	40.4
Cascade Mask R-CNN [90]	Swin-T	ImageNet-1k	1x	48.1	41.7
	Swin-T w MoTo	ImageNet-1k	1x	48.7	42.1
	Swin-T w TokenProp	ImageNet-1k	1x	48.6	41.9
	Swin-T w Both	ImageNet-1k	1x	49.1	42.4

From Table 15, we can see that both our strategies are beneficial for object detection and instance segmentation. We observe a consistent improvement over different backbones and detection pipelines.

8 CONCLUSION AND FUTURE WORK

In this work, we demonstrate the important role tokenizer plays in vision transformers. Based on the “trade-off” perspective that formulates prominent structural designs, we propose to incorporate better normalization and objective for tokenizers. Extensive experimental results manifest their effectiveness across different transformer structures.

Our findings further indicate that proper generative supervisory signals help improve discriminative performance. Notably, concurrent works on self-supervised learning with MAE [92] and BEiT [93] validates the potential of reconstruction objectives. Generalizing such generative signals to different tasks remains an open problem.

REFERENCES

- [1] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, “An image is worth 16x16 words: Transformers for image recognition at scale,” in *ICLR*, 2021.
- [2] Y. Zhang, X. Li, C. Liu, B. Shuai, Y. Zhu, B. Brattoli, H. Chen, I. Marsic, and J. Tighe, “Vidtr: Video transformer without convolutions,” in *ICCV*, 2021.
- [3] H. Shao, S. Qian, and Y. Liu, “Temporal interlacing network,” *AAAI*, 2020.
- [4] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, “End-to-end object detection with transformers,” in *ECCV*, 2020.
- [5] X. Zhu, W. Su, L. Lu, B. Li, X. Wang, and J. Dai, “Deformable detr: Deformable transformers for end-to-end object detection,” in *ICLR*, 2021.
- [6] S. Zheng, J. Lu, H. Zhao, X. Zhu, Z. Luo, Y. Wang, Y. Fu, J. Feng, T. Xiang, P. H. Torr, and L. Zhang, “Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers,” in *CVPR*, 2021.
- [7] E. Xie, W. Wang, Z. Yu, A. Anandkumar, J. M. Alvarez, and P. Luo, “Segformer: Simple and efficient design for semantic segmentation with transformers,” *arXiv preprint arXiv:2105.15203*, 2021.
- [8] W. Li, X. Lu, S. Qian, J. Lu, X. Zhang, and J. Jia, “On efficient

- transformer-based image pre-training for low-level vision," *arXiv preprint arXiv:2112.10175*, 2021.
- [9] Y. Tay, V. Q. Tran, S. Ruder, J. Gupta, H. W. Chung, D. Bahri, Z. Qin, S. Baumgartner, C. Yu, and D. Metzler, "Charformer: Fast character transformers via gradient-based subword tokenization," *arXiv preprint arXiv:2106.12672*, 2021.
- [10] H. El Boukkouri, O. Ferret, T. Lavergne, H. Noji, P. Zweigenbaum, and J. Tsujii, "CharacterBERT: Reconciling ELMo and BERT for word-level open-vocabulary representations from characters," in *COLING*, 2020.
- [11] T. Xiao, M. Singh, E. Mintun, T. Darrell, P. Dollár, and R. Girshick, "Early convolutions help transformers see better," in *NeurIPS*, 2021.
- [12] R. Sennrich, B. Haddow, and A. Birch, "Neural machine translation of rare words with subword units," in *ACL*, 2016.
- [13] T. Kudo, "Subword regularization: Improving neural network translation models with multiple subword candidates," in *ACL*, 2018.
- [14] J. H. Clark, D. Garrette, I. Turc, and J. Wieting, "Canine: Pre-training an efficient tokenization-free encoder for language representation," *arXiv preprint arXiv:2103.06874*, 2021.
- [15] W. Wang, E. Xie, X. Li, D.-P. Fan, K. Song, D. Liang, T. Lu, P. Luo, and L. Shao, "Pyramid vision transformer: A versatile backbone for dense prediction without convolutions," in *ICCV*, 2021.
- [16] H. Wu, B. Xiao, N. Codella, M. Liu, X. Dai, L. Yuan, and L. Zhang, "Cvt: Introducing convolutions to vision transformers," *arXiv preprint arXiv:2103.15808*, 2021.
- [17] L. Yuan, Y. Chen, T. Wang, W. Yu, Y. Shi, Z. Jiang, F. E. Tay, J. Feng, and S. Yan, "Tokens-to-token vit: Training vision transformers from scratch on imagenet," in *ICCV*, 2021.
- [18] X. Chen, S. Xie, and K. He, "An empirical study of training self-supervised visual transformers," in *ICCV*, 2021.
- [19] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *NIPS*, 2017.
- [20] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, and H. Jégou, "Training data-efficient image transformers & distillation through attention," *arXiv preprint arXiv:2012.12877*, 2020.
- [21] E. D. Cubuk, B. Zoph, J. Shlens, and Q. V. Le, "Randaugment: Practical automated data augmentation with a reduced search space," in *CVPR Workshops*, 2020.
- [22] E. D. Cubuk, B. Zoph, D. Mane, V. Vasudevan, and Q. V. Le, "Autoaugment: Learning augmentation policies from data," *arXiv preprint arXiv:1805.09501*, 2018.
- [23] Z. Chen, L. Xie, J. Niu, X. Liu, L. Wei, and Q. Tian, "Visformer: The vision-friendly transformer," in *ICCV*, 2021.
- [24] C.-F. Chen, Q. Fan, and R. Panda, "Crossvit: Cross-attention multi-scale vision transformer for image classification," *arXiv preprint arXiv:2103.14899*, 2021.
- [25] X. Chu, Z. Tian, B. Zhang, X. Wang, X. Wei, H. Xia, and C. Shen, "Conditional positional encodings for vision transformers," *arXiv preprint arXiv:2102.10882*, 2021.
- [26] W. Xu, Y. Xu, T. Chang, and Z. Tu, "Co-scale conv-attentional image transformers," in *ICCV*, 2021.
- [27] Z. Dai, H. Liu, Q. V. Le, and M. Tan, "Coatnet: Marrying convolution and attention for all data sizes," 2021.
- [28] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," in *ICCV*, 2021.
- [29] A. Vaswani, P. Ramachandran, A. Srinivas, N. Parmar, B. Hechtman, and J. Shlens, "Scaling local self-attention for parameter efficient visual backbones," in *CVPR*, 2021.
- [30] J. Yang, C. Li, P. Zhang, X. Dai, B. Xiao, L. Yuan, and J. Gao, "Focal self-attention for local-global interactions in vision transformers," in *NeurIPS*, 2021.
- [31] J. Lu, J. Yao, J. Zhang, X. Zhu, H. Xu, W. Gao, C. Xu, T. Xiang, and L. Zhang, "Soft: Softmax-free transformer with linear complexity," in *NeurIPS*, 2021.
- [32] X. Chu, Z. Tian, Y. Wang, B. Zhang, H. Ren, X. Wei, H. Xia, and C. Shen, "Twins: Revisiting the design of spatial attention in vision transformers," in *NeurIPS 2021*, 2021.
- [33] Y. Wang, R. Huang, S. Song, Z. Huang, and G. Huang, "Not all images are worth 16x16 words: Dynamic vision transformers with adaptive sequence length," in *NeurIPS*, 2021.
- [34] Y. Rao, W. Zhao, B. Liu, J. Lu, J. Zhou, and C.-J. Hsieh, "Dynamicvit: Efficient vision transformers with dynamic token sparsification," in *NeurIPS*, 2021.
- [35] B. Pan, R. Panda, Y. Jiang, Z. Wang, R. Feris, and A. Oliva, "Ia-red2: Interpretability-aware redundancy reduction for vision transformers," in *NeurIPS*, 2021.
- [36] C. Gong, D. He, X. Tan, T. Qin, L. Wang, and T.-Y. Liu, "Frage: Frequency-agnostic word representation," in *NeurIPS*, 2018.
- [37] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey, J. Klingner, A. Shah, M. Johnson, X. Liu, Łukasz Kaiser, S. Gouws, Y. Kato, T. Kudo, H. Kazawa, K. Stevens, G. Kurian, N. Patil, W. Wang, C. Young, J. Smith, J. Riesa, A. Rudnick, O. Vinyals, G. Corrado, M. Hughes, and J. Dean, "Google's neural machine translation system: Bridging the gap between human and machine translation," 2016.
- [38] D. Marin, J.-H. R. Chang, A. Ranjan, A. Prabhu, M. Rastegari, and O. Tuzel, "Token pooling in vision transformers," *arXiv preprint arXiv:2110.03860*, 2021.
- [39] W. Wang, E. Xie, X. Li, D.-P. Fan, K. Song, D. Liang, T. Lu, P. Luo, and L. Shao, "Pvtv2: Improved baselines with pyramid vision transformer," *arXiv preprint arXiv:2106.13797*, 2021.
- [40] T. Wang, L. Yuan, Y. Chen, J. Feng, and S. Yan, "Pnp-detr: Towards efficient visual analysis with transformers," in *ICCV*, 2021.
- [41] M. S. Ryoo, A. Piergiovanni, A. Arnab, M. Dehghani, and A. Angelova, "Tokenlearner: What can 8 learned tokens do for images and videos?" in *NeurIPS*, 2021.
- [42] X. Yue, S. Sun, Z. Kuang, M. Wei, P. Torr, W. Zhang, and D. Lin, "Vision transformer with progressive sampling," in *ICCV*, 2021.
- [43] Z. Jiang, Q. Hou, L. Yuan, D. Zhou, Y. Shi, X. Jin, A. Wang, and J. Feng, "All tokens matter: Token labeling for training better vision transformers," in *NeurIPS*, 2021.
- [44] P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, P.-A. Manzagol, and L. Bottou, "Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion." *JMLR*, 2010.
- [45] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.
- [46] R. D. Hjelm, A. Fedorov, S. Lavoie-Marchildon, K. Grewal, P. Bachman, A. Trischler, and Y. Bengio, "Learning deep representations by mutual information estimation and maximization," in *ICLR*, 2019.
- [47] A. M. Saxe, Y. Bansal, J. Dapello, M. Advani, A. Kolchinsky, B. D. Tracey, and D. D. Cox, "On the information bottleneck theory of deep learning," in *ICLR*, 2018.
- [48] S. Qian, K.-Y. Lin, W. Wu, Y. Liu, Q. Wang, F. Shen, C. Qian, and R. He, "Make a face: Towards arbitrary high fidelity face manipulation," in *ICCV*, 2019.
- [49] S. Qian, K. Sun, W. Wu, C. Qian, and J. Jia, "Aggregation via separation: Boosting facial landmark detector with semi-supervised style translation," in *ICCV*, 2019.
- [50] E. Voita, R. Sennrich, and I. Titov, "The bottom-up evolution of representations in the transformer: A study with machine translation and language modeling objectives," in *EMNLP*, Nov. 2019.
- [51] W. Wang, L. Yao, L. Chen, B. Lin, D. Cai, X. He, and W. Liu, "Crossformer: A versatile vision transformer hinging on cross-scale attention," *arXiv preprint arXiv:2108.00154*, 2021.
- [52] H. Fan, B. Xiong, K. Mangalam, Y. Li, Z. Yan, J. Malik, and C. Feichtenhofer, "Multiscale vision transformers," in *ICCV*, 2021.
- [53] M. I. Belghazi, A. Baratin, S. Rajeshwar, S. Ozair, Y. Bengio, A. Courville, and D. Hjelm, "Mutual information neural estimation," in *ICML*, 2018.
- [54] S. Qian, H. Shao, Y. Zhu, M. Li, and J. Jia, "Blending anti-aliasing into vision transformer," in *NeurIPS*, 2021.
- [55] K. Han, A. Xiao, E. Wu, J. Guo, C. Xu, and Y. Wang, "Transformer in transformer," in *NeurIPS*, 2021.
- [56] Z. Chen, Y. Zhu, C. Zhao, G. Hu, W. Zeng, J. Wang, and M. Tang, "Dpt: Deformable patch-based transformer for visual recognition," in *ACMMM*, 2021.
- [57] B. Zhou, H. Zhao, X. Puig, T. Xiao, S. Fidler, A. Barriuso, and A. Torralba, "Semantic understanding of scenes through the ade20k dataset," *IJCV*, 2019.
- [58] R. Xiong, Y. Yang, D. He, K. Zheng, S. Zheng, C. Xing, H. Zhang, Y. Lan, L. Wang, and T. Liu, "On layer normalization in the transformer architecture," in *ICML*, 2020.

- [59] S. Shen, Z. Yao, A. Gholami, M. Mahoney, and K. Keutzer, "Powernorm: Rethinking batch normalization in transformers," in *International Conference on Machine Learning*, 2020.
- [60] H. Touvron, M. Cord, A. Sablayrolles, G. Synnaeve, and H. Jégou, "Going deeper with image transformers," *arXiv preprint arXiv:2103.17239*, 2021.
- [61] X. Huang and S. Belongie, "Arbitrary style transfer in real-time with adaptive instance normalization," in *ICCV*, 2017.
- [62] T. Park, M.-Y. Liu, T.-C. Wang, and J.-Y. Zhu, "Semantic image synthesis with spatially-adaptive normalization," in *CVPR*, 2019.
- [63] T. Miyato and M. Koyama, "cgans with projection discriminator," *arXiv preprint arXiv:1802.05637*, 2018.
- [64] M. Caron, P. Bojanowski, A. Joulin, and M. Douze, "Deep clustering for unsupervised learning of visual features," in *ECCV*, 2018.
- [65] D. Ulyanov, A. Vedaldi, and V. Lempitsky, "Instance normalization: The missing ingredient for fast stylization," *arXiv preprint arXiv:1607.08022*, 2016.
- [66] R. Wightman, "Pytorch image models," <https://github.com/rwightman/pytorch-image-models>, 2019.
- [67] F. Research, "Deit," <https://github.com/facebookresearch/deit>, 2021.
- [68] Microsoft, "Swin transformer," <https://github.com/microsoft/Swin-Transformer>, 2021.
- [69] YITU, "T2t vit," <https://github.com/yitu-opensource/T2T-ViT>, 2021.
- [70] P. V. Transformer, "Pvt," <https://github.com/whai362/PVT>, 2021.
- [71] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *CVPR*, 2009.
- [72] Y. Wang, Y.-C. Chen, X. Zhang, J. Sun, and J. Jia, "Attentive normalization for conditional image generation," in *CVPR*, 2020.
- [73] X. Pan, P. Luo, J. Shi, and X. Tang, "Two at once: Enhancing learning and generalization capacities via ibn-net," in *ECCV*, 2018.
- [74] Y. Li, N. Wang, J. Shi, J. Liu, and X. Hou, "Revisiting batch normalization for practical domain adaptation," *arXiv preprint arXiv:1603.04779*, 2016.
- [75] Z. Wu, X. Han, Y.-L. Lin, M. G. Uzunbas, T. Goldstein, S. N. Lim, and L. S. Davis, "Dcan: Dual channel-wise alignment networks for unsupervised scene adaptation," in *ECCV*, 2018.
- [76] J. L. Ba, J. R. Kiros, and G. E. Hinton, "Layer normalization," *arXiv preprint arXiv:1607.06450*, 2016.
- [77] Z. Yao, Y. Cao, Y. Lin, Z. Liu, Z. Zhang, and H. Hu, "Leveraging batch normalization for vision transformers," in *ICCV Workshop*, 2021.
- [78] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.
- [79] Y. Wang, Z. Ni, S. Song, L. Yang, and G. Huang, "Revisiting locally supervised learning: an alternative to end-to-end training," in *ICLR*, 2021.
- [80] E. Belilovsky, M. Eickenberg, and E. Oyallon, "Greedy layerwise learning can scale to imagenet," in *ICML*, 2019.
- [81] L. Kong, C. de Masson d'Autume, L. Yu, W. Ling, Z. Dai, and D. Yogatama, "A mutual information maximization perspective of language representation learning," in *ICLR*, 2019.
- [82] M. U. Gutmann and A. Hyvärinen, "Noise-contrastive estimation of unnormalized statistical models, with applications to natural image statistics," *JMLR*, 2012.
- [83] J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual losses for real-time style transfer and super-resolution," in *ECCV*, 2016.
- [84] R. Mechrez, I. Talmi, and L. Zelnik-Manor, "The contextual loss for image transformation with non-aligned data," *arXiv preprint arXiv:1803.02077*, 2018.
- [85] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *ICCV*, 2017.
- [86] I. Loshchilov and F. Hutter, "Fixing weight decay regularization in adam," 2018.
- [87] Y. Yuan, X. Chen, and J. Wang, "Object-contextual representations for semantic segmentation," in *ECCV*, 2020.
- [88] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *ICCV*, 2017.
- [89] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *ICCV*, 2017.
- [90] Z. Cai and N. Vasconcelos, "Cascade r-cnn: Delving into high quality object detection," in *CVPR*, 2018.
- [91] K. Chen, J. Wang, J. Pang, Y. Cao, Y. Xiong, X. Li, S. Sun, W. Feng, Z. Liu, J. Xu *et al.*, "Mmdetection: Open mmlab detection toolbox and benchmark," *arXiv preprint arXiv:1906.07155*, 2019.
- [92] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick, "Masked autoencoders are scalable vision learners," *arXiv preprint arXiv:2111.06377*, 2021.
- [93] H. Bao, L. Dong, and F. Wei, "Beit: Bert pre-training of image transformers," *arXiv preprint arXiv:2106.08254*, 2021.
- [94] T. Xiao, Y. Liu, B. Zhou, Y. Jiang, and J. Sun, "Unified perceptual parsing for scene understanding," in *ECCV*, 2018.
- [95] F. Research, "Mocov3," <https://github.com/facebookresearch/moco-v3>, 2019.
- [96] M. Contributors, "MMSegmentation: Openmmlab semantic segmentation toolbox and benchmark," <https://github.com/open-mmlab/mms Segmentation>, 2020.

APPENDIX A ADDITIONAL DETAILS

A.1 Details on Estimating Conditional Entropy

In Section 3 and Section 4, we estimate the empirical mutual information between input images and split tokens with conditional entropy. The metric is computed by training an additional decoder to estimate the reconstruction error. We train the decoder using the Adam optimizer, with the learning rate of 0.001, beta1 0.5, beta2 0.999 on V100 GPUs. As shown in Table 16, we use a simple 3-layer architecture for the decoder. Note that we also adopt this architecture as decoder in TokenProp for simplicity.

TABLE 16

Details parameters of decoders. The convolutional layer denotes residual blocks. We also adopt this architecture for TokenProp.

Layer	Kernel Size	Output Channel	Output Size
Input w Concat	-	196	16
Convolution	3	256	16
Convolution	3	256	16
Pixel Shuffle	-	64	32
Convolution	3	64	32
Pixel Shuffle	-	16	64
Convolution	3	3	64

A.2 Details on Comparison across structural designs.

Specifically, the intra-token refinement layer consists of kernels with sizes 4×4 , 8×8 , and 16×16 , further concatenating the extracted multi-scale feature and projecting it to the original dimension. Following [39], the locality strategy enlarges the token window and performs overlapping tokenization with a half-patch step size. Inter-token refinement further adds a self-attention layer before the output.

For the experiments, we adopt DeiT-S as our baseline architecture and use the consistent training recipes [20] for 300 epochs on ImageNet [71] training set for supervised classification. Strictly following MoCov3 [18], we also conduct unsupervised pre-training and use the linear probing accuracy to test the performance. Semantic segmentation is evaluated consistently using UperNet [94] backbone and the hyper-parameters in Appendix B.3.

APPENDIX B IMPLEMENTATIONS AND HYPER-PARAMETERS

B.1 Linear Classification with Self-supervised Training

In Section 4, we study the performance of linear classification with self-supervised pre-training as in MoCov3 [18]. Here we provide additional details on running these experiments.

Self-supervised Pre-training. Following the practice in their official implementation [95], we faithfully implement our variants with DeiT-S backbone. By default we use AdamW [86], a batch size of 4096, and a warmup schedule for 10k steps. The learning rate and weight decay is also swept for multiple runs. We also follow the consistent architectures and implementations of the MLP heads and contrastive loss in MoCov3 [18] for fair comparisons.

Linear Probing. As a common practice in estimating the representation quality, we evaluate the trained models with linear probing. After the self-supervised pre-training completes, we remove the MLP layers and train an additional linear classification based on the obtained frozen features. We adopt the SGD optimizer with a batch size of 4096. Similarly, the learning rate is swept by multiple runs for better performance.

B.2 Supervised Classification on ImageNet

As we’ve mentioned in Section 5.2, we integrate our proposed module to various state-of-the-art models, including DeiT [20], Token-2-Token ViT [17], PVT [15], and Swin Transformer [28]. As we don’t intend to propose a new training pipeline in this work, we adopt the common practice in training vision transformers. Since these methods both contain distinctive designs and different pipelines for training, we adopt their original training recipes and keep the training parameters consistent with the provided ones, including data augmentations, batch sizes, and optimizers. As DeiT [20] and Swin serve as two major backbones in this work, we provide the ingredients and hyper-parameters for training in Table 17.

TABLE 17

Ingredients and hyper-parameters for training DeiT and Swin. For training Swin-S, we use a smaller batch size and a scaled learning rate. EMA represents Exponential Moving Average.

Methods	DeiT-S	Swin-T/S
Training Epochs	300	300
Batch Size	1024	512
Optimizer	AdamW	AdamW
Learning rate	0.001	0.0005
Weight decay	0.05	0.05
Warmup epochs	5	20
EMA	✓	
Stochastic depth	✓	✓
Repeated augmentation	✓	
Gradient Clip		✓
Rand Augment	✓	✓
Mixup Prob	0.8	0.8
Cutmix Prob	1.0	1.0
Erasing Prob	0.25	0.25

B.3 Semantic Segmentation on ADE20K

ADE20K dataset consists of 150 different semantic categories, where there are 20210 training images, 2000 validation images, and 3000 testing images. We follow the consistent training recipes in [28] and use UperNet [94] in MMsegmentation [96]. The models are trained for 160K iterations with the AdamW optimizer, where the initial learning rate and weight decay are set to 6×10^{-5} and 0.001. During inference, we also follow the multi-scale testing strategy in [28], where $[0.5, 0.75, 1.0, 1.25, 1.5, 1.75] \times$ of the training resolution is employed.

B.4 Object Detection and Instance Segmentation

The dataset consists of 118K training images, 5K validation images, and 20K test-dev images. Here we exploit three widely-used detection pipelines, RetinaNet [88], Mask R-CNN [89], and Cascade Mask R-CNN [90]. Following the settings and hyper-parameters in [28], we use AdamW optimizer with an initial learning rate of 10^{-4} and a weight decay of 0.05. The training process includes multi-scale training, a batch size of 16, and 1x(12 epochs) training schedule. We re-implement both the baseline models and our variants using their open-source script in MMDetection [91].