Expediting Vision Transformer for Dense Prediction without Fine-tuning

Anonymous Author(s) Affiliation Address email

Abstract

1 Vision transformers have recently achieved competitive results across various 2 vision tasks but still suffer from heavy computation costs when processing a large number of tokens. Many advanced approaches have been developed to reduce 3 the total number of tokens in the large-scale vision transformers, especially for 4 image classification tasks. Typically, they select a small group of essential tokens 5 according to their relevance with the [class] token, then fine-tune the weights of 6 the vision transformer. Such fine-tuning is less practical for dense prediction due 7 to the much heavier computation and GPU memory cost than image classification. 8 In this paper, we focus on a more challenging problem, *i.e.*, accelerating large-scale 9 vision transformers for dense prediction without any additional re-training or fine-10 tuning. In response to the fact that high-resolution representations are necessary 11 for dense prediction, we present two non-parametric operators, a *token clustering* 12 layer to decrease the number of tokens and a token reconstruction layer to increase 13 the number of tokens. The following steps are performed to achieve this: (i) we 14 use the token clustering layer to cluster the neighboring tokens together, resulting 15 in low-resolution representations that maintain the spatial structures; (ii) we apply 16 the following transformer layers only to these low-resolution representations or 17 18 clustered tokens; and (iii) we use the token reconstruction layer to re-create the high-resolution representations from the refined low-resolution representations. 19 The results obtained by our method are promising on five dense prediction tasks 20 21 including object detection, semantic segmentation, panoptic segmentation, instance segmentation, and depth estimation. Accordingly, our method accelerates 40% \uparrow 22 FPS and saves $30\% \downarrow$ GFLOPs of "Segmenter+ViT-L/16" while maintaining 99.5%23 of the performance on ADE20K without fine-tuning the official weights. 24

25 1 Introduction

Transformer [64] has made significant progress across various challenging vision tasks since pi-26 oneering efforts such as DETR [5], Vision Transformer (ViT) [16], and Swin Transformer [44]. 27 By removing the local inductive bias [17] from convolutional neural networks [25, 61, 57], vi-28 sion transformers armed with global self-attention show superiority in scalability for large-scale 29 models and billion-scale dataset [16, 78, 58], self-supervised learning [24, 72, 1], connecting vi-30 sion and language [50, 31], etc. We can find from recent developments of SOTA approaches that 31 vision transformers have dominated various leader-boards, including but not limited to image clas-32 33 sification [70, 14, 15, 78], object detection [80, 43, 35], semantic segmentation [29, 12, 3], pose estimation [73], image generation [79], and depth estimation [39]. 34

Although vision transformers have achieved more accurate predictions in many vision tasks, largescale vision transformers are still burdened with heavy computational overhead, particularly when

Submitted to 36th Conference on Neural Information Processing Systems (NeurIPS 2022). Do not distribute.



Figure 1: Illustrating the improvements of our approach: we report the results of applying our approach to Segmenter [60] for semantic segmentation and DPT [51] for depth estimation on the 1-st and 2-ed row respectively. Without any fine-tuning, our proposed method reduces the GFLOPs and accelerates the FPS significantly with a slight performance drop on both dense prediction tasks. \uparrow and \downarrow represent higher is better and lower is better respectively. Refer to Section 4 for more details.

processing high-resolution inputs [18, 43], thus limiting their broader application on more resource-37 constrained applications and attracting efforts on re-designing light-weight vision transformer ar-38 chitectures [11, 48, 81]. In addition to this, several recent efforts have investigated how to decrease 39 the model complexity and accelerate vision transformers, especially for image classification, and 40 introduced various advanced approaches to accelerate vision transformers. Dynamic ViT [52] and 41 EVIT [41], for example, propose two different dynamic token sparsification frameworks to reduce 42 the redundant tokens progressively and select the most informative tokens according to the scores 43 predicted with an extra trained prediction module or their relevance with the [class] token. To-44 kenLearner [55] learns to spatially attend over a subset of tokens and generates a set of clustered 45 tokens adaptive to the input for video understanding tasks. Most of these token reduction approaches 46 are carefully designed for image classification tasks and require fine-tuning or retraining. These 47 approaches might not be suitable to tackle more challenging dense prediction tasks that need to 48 process high-resolution input images, e.g., 1024×1024 , thus, resulting heavy computation and 49 GPU memory cost brought. We also demonstrate in the supplemental material the superiority of our 50 method over several representative methods on dense prediction tasks. 51

Rather than proposing a new lightweight architecture for dense prediction or token reduction scheme 52 for only image classification, we focus on how to expedite well-trained large-scale vision transformers 53 and use them for various dense prediction tasks without fine-tuning or re-training. Motivated by 54 these two key observations including (i) the intermediate token representations of a well-trained 55 vision transformer carry a heavy amount of local spatial redundancy and (ii) dense prediction tasks 56 require high-resolution representations, we propose a simple yet effective scheme to convert the 57 "high-resolution" path of the vision transformer to a "high-to-low-to-high resolution" path via two 58 non-parametric layers including a token clustering layer and a token reconstruction layer. Our method 59 can produce a wide range of more efficient models without requiring further fine-tuning or re-training. 60 We apply our approach to expedite two main-stream vision transformer architectures, e.g., ViTs and 61 Swin Transformers, for five challenging dense prediction tasks, including object detection, semantic 62 segmentation, panoptic segmentation, instance segmentation, and depth estimation. We have achieved 63 encouraging results across several evaluated benchmarks and Figure 1 illustrates some representative 64 results on both semantic segmentation and depth estimation tasks. 65

66 2 Related work

Pruning Convolutional Neural Networks. Convolutional neural network pruning [2, 27, 67] is
 a task that involves removing the redundant parameters to reduce the model complexity without
 a significant performance drop. Pruning methods typically entail three steps: (i) training a large,
 over-parameterized model to convergence, (ii) pruning the trained large model according to a certain

criterion, and (iii) fine-tuning the pruned model to regain the lost performance [46]. The key idea is 71 to design an importance score function that is capable of pruning the less informative parameters. 72 We follow [8] to categorize the existing methods into two main paths: (i) unstructured pruning 73 (also named weight pruning) and (ii) structured pruning. Unstructured pruning methods explore the 74 absolute value of each weight or the product of each weight and its gradient to estimate the importance 75 scores. Structured pruning methods, such as layer-level pruning [68], filter-level pruning [45, 75], and 76 image-level pruning [23, 62], removes the model sub-structures. Recent studies [6, 7, 28] further 77 extend these pruning methods to vision transformer. Unlike the previous pruning methods, we explore 78 how to expedite vision transformers for dense prediction tasks by carefully reducing & increasing the 79 number of tokens without removing or modifying the parameters. 80 Efficient Vision Transformer. The success of vision transformers has incentivised many recent 81

efforts [55, 59, 47, 26, 69, 52, 34, 22, 65, 10, 32, 41, 54] to exploit the spatial redundancies of 82 intermediate token representations. For example, TokenLearner [55] learns to attend over a subset 83 of tokens and generates a set of clustered tokens adaptive to the input. They empirically show that 84 very few clustered tokens are sufficient for video understanding tasks. Token Pooling [47] exploits 85 a nonuniform data-aware down-sampling operator based on K-Means or K-medoids to cluster 86 the similar tokens together to reduce the number of tokens while minimizing the reconstruction 87 error. Dynamic ViT [52] observes that the accurate image recognition with vision transformers 88 89 mainly depends on a subset of the most informative tokens, and hence it develops a dynamic token sparsification framework for pruning the redundant tokens dynamically based on the input. EViT 90 (expediting vision transformers) [41] proposes to calculate the attentiveness of the [class] token with 91 respect to each token and identify the top-k attentive tokens according to the attentiveness score. Patch 92 Merger [53] uses a learnable attention matrix to merge and combine together the redundant tokens, 93 therefore creating a much more practical and cheaper model with only a slight performance drop. 94 Refer to [63] for more details on efficient transformer architecture designs, such as Performer [13] 95 and Reformer [33]. In contrast to these methods that require either retraining or fine-tuning the 96 modified transformer architectures from scratch or the pre-trained weights, our approach can reuse 97 98 the once trained weights for free and produce light-weight models with a modest performance drop.

Vision Transformer for Dense Prediction. In the wake of early success of the representative 99 pyramid vision transformers [44, 66] for object detection and semantic segmentation, more and more 100 101 efforts have explored many different advanced vision transformer architecture designs [36, 9, 37, 38, 19, 77, 74, 82, 21, 40, 71, 76] suitable for various dense prediction tasks. For example, MViT [37] 102 focuses more on multi-scale representation learning, while HRFormer [77] examines the benefits of 103 combining multi-scale representation learning and high-resolution representation learning. Instead of 104 designing a novel vision transformer architecture for dense prediction, we focus on how to accelerate 105 a well-trained vision transformer while maintaining the prediction performance as much as possible. 106

107 **3** Our Approach

Preliminary. The conventional Vision Transformer [16] first reshapes the input image $\mathbf{X} \in \mathbb{R}^{H \times W \times 3}$ into a sequence of flatten patches $\mathbf{X}_p \in \mathbb{R}^{N \times (P^23)}$, where (P, P) represents the resolution of each patch, (H, W) represents the resolution of the input image, $N = (H \times W)/P^2$ represents the number of resulting patches or tokens, *i.e.*, the input sequence length. The Vision Transformer consists of alternating layers of multi-head self-attention (MHSA) and feed-forward network (FFN) accompanied with layer norm (LN) and residual connections:

$$\mathbf{Z}_{l} = \text{MHSA}\left(\text{LN}\left(\mathbf{Z}_{l-1}\right)\right) + \mathbf{Z}_{l-1},
\mathbf{Z}_{l} = \text{FFN}\left(\text{LN}\left(\mathbf{Z}_{l}^{'}\right)\right) + \mathbf{Z}_{l}^{'},$$
(1)

where $l \in \{1, ..., L\}$ represents the layer index, $\mathbf{Z}_l \in \mathbb{R}^{N \times C}$, and \mathbf{Z}_0 is based on \mathbf{X}_p . The computation cost $\mathcal{O}(\text{LNC}(N + C))$ mainly depends on the number of layers L, the number of tokens N, and the channel dimension C.

Despite the great success of transformer, its computation cost increases significantly when handling high-resolution representations, which are critical for dense prediction tasks. This paper attempts to resolve this issue by reducing the computation complexity during the inference stage, and presents a



Figure 2: (a) Plain high-resolution vision transformer with L layers. (b) U-shape High-to-lowto-high-resolution vision transformer with α , β , and γ layers respectively (L = $\alpha + \beta + \gamma$). (c) Illustrating the details of using our approach to plain ViTs: we insert a token clustering layer and a token reconstruction layer into a trained vision transformer in order to decrease and then increase the spatial resolution, respectively. The weights of modules marked with \square are trained once based on the configuration of (a). The token clustering layer and token reconstruction layer are marked with \square are non-parametric, thus do not require any fine-tuning and can be included directly during evaluation.

very simple solution for generating a large number of efficient vision transformer models directly
 from a single trained vision transformer, requiring no further training or fine-tuning.

We demonstrate how our approach could be applied to the existing standard Vision Transformer in 122 Figure 2. The original Vision Transformer is modified using two non-parametric operations, namely 123 a token clustering layer and a token reconstruction layer. The proposed token clustering layer is 124 utilized to convert the high-resolution representations to low-resolution representations by clustering 125 the locally semantically similar tokens. Then, we apply the following transformer layers on the 126 low-resolution representations, which greatly accelerates the inference speed and saves computation 127 resources. Last, a token reconstruction layer is proposed to reconstruct the feature representations 128 back to high-resolution. 129

Token Clustering Layer. We construct the token clustering layer following the improved SLIC
 scheme [30], which performs local k-means clustering as follows:

¹³² *-Initial superpixel center*: We apply adaptive average pooling (AAP) over the high-resolution ¹³³ representations from the α -th layer to compute the h × w initial cluster center representations:

$$\mathbf{S}_{\alpha} = \mathrm{AAP}(\mathbf{Z}_{\alpha}, (\mathsf{h} \times \mathsf{w})), \tag{2}$$

134 where $\mathbf{S}_{\alpha} \in \mathbb{R}^{\mathsf{hw} \times \mathsf{C}}$, $\mathbf{Z}_{\alpha} \in \mathbb{R}^{\mathsf{N} \times \mathsf{C}}$, and $\mathsf{hw} \ll \mathsf{N}$.

¹³⁵ *-Iterative local clustering*: (i) Expectation step: compute the normalized similarity between each pixel ¹³⁶ p and the surrounding superpixel i (we only consider the neighboring λ positions), (ii) Maximization ¹³⁷ step: compute the new superpixel centers:

$$\mathbf{Q}_{p,i} = \frac{\exp\left(||\mathbf{Z}_{\alpha,p} - \mathbf{S}_{\alpha,i}||^2/\tau\right)}{\sum_{j=1}^{\lambda} \exp\left(||\mathbf{Z}_{\alpha,p} - \mathbf{S}_{\alpha,j}||^2/\tau\right)}, \quad \mathbf{S}_{\alpha,i} = \sum_{p=1}^{\mathsf{N}} \mathbf{Q}_{p,i} \mathbf{Z}_{\alpha,p}, \tag{3}$$

where we iterate the above Expectation step and Maximization step for κ times, τ is a temperature hyper-parameter, and $i \in \{1, 2, \dots, \lambda\}$. We apply the following β transformer layers on \mathbf{S}_{α} instead of \mathbf{Z}_{α} , thus results in $\mathbf{S}_{\alpha+\beta}$ and decreases the computation cost significantly.

Token Reconstruction Layer. We implement the token reconstruction layer by exploiting the relations between the high-resolution representations and the low-resolution clustered representations:



Figure 3: Illustrating the details of using our approach for Swin Transformer [43, 44]: we apply four groups of token clustering layer and token reconstruction layer within the four non-overlapped windows marked with \blacksquare , \blacksquare , and \blacksquare respectively, in the example referred to as *window token clustering layer* and *window token reconstruction layer*. We apply the intermediate swin transformer layers equipped with window size k × k) on the clustered window tokens. Window sizes before and after token clustering layer are K × K and k × k, and vice versa, for the token reconstruction layer.

143

$$\mathbf{Z}_{\alpha+\beta,p} = \sum_{\mathbf{S}_{\alpha,i} \in k\text{-NN}(\mathbf{Z}_{\alpha,p})} \frac{\exp(||\mathbf{Z}_{\alpha,p} - \mathbf{S}_{\alpha,i}||^2/\tau)}{\sum_{\mathbf{S}_{\alpha,j} \in k\text{-NN}(\mathbf{Z}_{\alpha,p})} \exp(||\mathbf{Z}_{\alpha,p} - \mathbf{S}_{\alpha,j}||^2/\tau)} \mathbf{S}_{\alpha+\beta,i},\tag{4}$$

where τ is the same temperature hyper-parameter as in Equation 3. k-NN($\mathbf{Z}_{\alpha,p}$) represents a set of the k nearest, a.k.a, most similar, superpixel representations for $\mathbf{Z}_{\alpha,i}$. We empirically find that choosing the same neighboring positions as in Equation 3 achieves close performance as the k-NN scheme while being more easy to implementation.

In summary, we estimate their semantic relations based on the representations before refinement with the following β transformer layers and then reconstruct the high-resolution representations from the refined low-resolution clustered representations accordingly.

Finally, we apply the remained γ transformer layers to the reconstructed high-resolution features and the task-specific head on the refined high-resolution features to predict the target results such as semantic segmentation maps or monocular depth maps.

Extension to Swin Transformer. We further introduce the *window token clustering layer* and *window* 154 token reconstruction layer, which are suitable for Swin Transformer [43, 44]. Figure 3 illustrates 155 an example usage of the proposed window token clustering layer and window token reconstruction 156 layer. We first cluster the $K \times K$ window tokens into $k \times k$ window tokens and then reconstruct $K \times K$ 157 window tokens according to the refined $k \times k$ window tokens. We apply the swin transformer layer 158 equipped with smaller window size $k \times k$ on the clustered representations, where we need to bi-linear 159 interpolate the pre-trained weights of relative position embedding table from $(2K - 1) \times (2K - 1)$ to 160 $(2k-1) \times (2k-1)$ when processing the clustered representations. In summary, we can improve the 161 efficiency of Swin Transformer by injecting the window token clustering layer and the window token 162 reconstruction layer into the backbones seamlessly without fine-tuning the model weights. 163

164 4 Experiment

We verify the effectiveness of our method across five challenging dense prediction tasks, including object detection, semantic segmentation, instance segmentation, panoptic segmentation, and monocular depth estimation. We carefully choose the advanced SOTA methods that build the framework based on either the plain ViTs [16] or the Swin Transformers [43, 44]. We can integrate the proposed token clustering layer and token reconstruction layer seamlessly with the provided official trained checkpoints with no further fine-tuning required. More experimental details are illustrated as follows.

171 4.1 Datasets

COCO [42]. This dataset consists of 123K images with 896K annotated bounding boxes belonging to 80 thing classes and 53 stuff classes, where the train set contains 118K images and the val set contains 5K images. We report the object detection performance of SwinV2 + HTC++ [43] and the instance/panoptic segmentation performance of Mask2Former [12] on the val set.

ADE20K [83]. This dataset contains challenging scenes with fine-grained labels and is one of the most challenging semantic segmentation datasets. The train set contains 20, 210 images with 150

semantic classes. The val set contains 2,000 images. We report the segmentation results with
Segmenter [60] on the val set.

PASCAL-Context [49]. This dataset consists of 59 semantic classes plus a background class, where
 the train set contains 4,996 images with and the val set contains 5,104 images. We report the
 segmentation results with Segmenter [60] on the val set.

Cityscapes [4]. This dataset is an urban scene understanding dataset with 30 classes while only 19 classes are used for parsing evaluation. The train set and val set contains 2, 975 and 500 images respectively. We report the segmentation results with Segmenter [60] on the val set.

KITTI [20]. This dataset provides stereo, optical flow, visual odometry (SLAM), and 3D object detection of outdoor scenes captured by equipment mounted on a moving vehicle. We choose DPT [51] as our baseline to conduct experiments on the monocular depth prediction tasks, which consists of around 26K images for train set and 698 images for val set, where only 653 images have the ground-truth depth maps and the image resolution is of $1, 241 \times 376$.

NYUv2 [56]. This dataset consists of 1, 449 RGBD images with resolution of 640×480 , which captures 464 diverse indoor scenes and contains rich detailed dense annotations such as surface normals, segmentation maps, depth, 3D planes, and so on. We report the depth prediction results of DPT [51] evaluated on 655 val images.

195 4.2 Evaluation Metrics

We report the numbers of AP (average precision), mask AP (mask average precision), PQ (panoptic quality), mIoU (mean intersection-over-union), and RMSE (root mean squared error) across object detection, instance segmentation, panoptic segmentation, semantic segmentation, and depth estimation tasks respectively. Since the dense prediction tasks care less about the throughput used in image classification tasks [41], we report FPS to measure the latency and the number of GFLOPs to measure the model complexity during evaluation. FPS is tested on a single V100 GPU with Pytorch 1.10 and CUDA 10.2 by default. More details are provided in the supplementary material.

203 4.3 Ablation Study Experiments

We conduct the following ablation experiments on ADE20K semantic segmentation benchmark with the official checkpoints of Segmenter+ViT-L/16 [60]¹ by default if not specified.

Hyper-parameters of token clustering/reconstruction layer. We first study the influence of the 206 hyper-parameters associated with the token clustering layer, i.e., the number of neighboring pixels λ 207 used in Equation 3, the number of EM iterations κ , and the choice of the temperature τ in Table 1. 208 According to the results, we can see that our method is relatively less sensitive to the choice of 209 210 both λ and κ compared to τ . In summary, we choose λ as 5×5 , κ as 5, and τ as 50 considering both performance and efficiency. Next, we also study the influence of the hyper-parameters within 211 the token clustering layer, i.e., the number of nearest neighbors k within k-NN. We do not observe 212 obvious differences and thus set k as 20. More details are provided in the supplementary material. 213

Influence of cluster size choices. We study the influence of different cluster size $h \times w$ choices 214 based on input feature map of size $\frac{H}{P} \times \frac{W}{P} = 40 \times 40$ (N = 1,600)² in Table 2. According to the 215 results, we can see that choosing too small cluster sizes significantly harms the dense prediction 216 performance, and setting $h \times w$ as 28×28 achieves the better trade-off between performance drop 217 and model complexity. Therefore, we choose 28×28 on Segmenter+ViT-L/16 by default. We also 218 empirically find that selecting the cluster size h \times w around N/4 \sim N/2 performs better on most of 219 the other experiments. We conduct the following ablation experiments under two typical settings, 220 including 20×20 (~ N/4) and 28×28 (~ N/2). 221

Comparison with adaptive average pooling and bi-linear upsample. We report the comparison results between our proposed token clustering/reconstruction scheme and adaptive average pooling/bilinear upsample scheme in Table 3 and Table 4 under two cluster size settings respectively. We choose to compare with adaptive average pooling and bi-linear upsample instead of strided convolution or deconvolution as the previous ones are non-parametric and the later ones require re-training or

¹https://github.com/rstrudel/segmenter#ade20k, MIT License

²We choose $H \times W = 640 \times 640$ and P = 16, thus, $\frac{H}{P} \times \frac{W}{P} = 40 \times 40$ or N = 1,600, on ADE20K.

Table	1٠	Influence	of the	hu	nor	noromate	are of	tokan	chue	toring	reconst	truct	ion	101/0	
rabic	1.	mnuence	or the	iny	DCI-	paramen	15 01	token	cius	tering/	recons	uci	1011	Iayc.	л.

Parameter		λ			κ				1	Г	k			
1 arameter	3×3	5×5	7×7	5	10	15	20	10	25	50	75	10	20	50
GFLOPs	438.6	438.9	439.6	438.9	439.6	440.2	440.8	438.9				438.9	438.9	439.0
mIoU	51.46	51.56	51.55	51.56	51.59	51.62	51.62	51.18	51.35	51.56	51.07	51.34	51.56	51.56

Table 2: Influence of the cluster size $h \times w$.							
Cluster size	8×8	12×12	16×16	20×20	24×24	28×28	32×32 (baseline)
GFLOPs	274.0	290.9	315.1	347.2	388.2	438.9	659.0
mIoU	32.13	44.01	48.21	50.17	51.32	51.56	51.82

Table 3: Comparison with adap-Table 4: Comparison with bi-Table 5: Combination with lighter vision tive average pooling(AAP). linear upsample. transformer backbone.

Cluster size	20×20	28×28	Cluster size	20×20	28×28	Cluster size	mIoU	GFLOPs
AAP	46.45	46.54	Bi-linear	44.68	44.74	Segmenter+ViT-B/16	48.48	124.7
Ours	50.17	51.56	Ours	50.17	51.56	Segmenter+ViT-B/16+Ours	48.40	91.9



Figure 4: Influence of the inserted position α of token Figure 5: Influence of the inserted position $\alpha + \beta$ of clustering layer. token reconstruction layer.

fine-tuning, which are not the focus of this work. Specifically, we keep the inserted position choices 227 the same and only replace the token cluster or token reconstruction layer with adaptive average 228 pooling or bi-linear upsampling under the same cluster size choices. According to the results, we 229 can see that our proposed token clustering and token reconstruction consistently outperform adaptive 230 average pooling and bi-linear upsampling under different cluster size choices. 231

Influence of inserted position of token clustering/reconstruction layer. We investigate the 232 influence of the inserted position of both token clustering layer and token reconstruction layer and 233 summarize the detailed results in Figure 4 and Figure 5 under two different cluster size choices. 234 According to the results shown in Figure 4, our method achieves better performance when choosing 235 α larger than 10, therefore, we choose $\alpha = 10$ as it achieves a better trade-off between model 236 complexity and segmentation accuracy. Then we study the influence of the inserted positions of 237 the token reconstruction layer by fixing $\alpha = 10$. According to Figure 5, we can see that our 238 method achieves the best performance when setting $\alpha + \beta = 24$, in other words, we insert the token 239 reconstruction layer after the last transformer layer of ViT-L/16. We choose $\alpha = 10, \alpha + \beta = 24$, 240 and $\gamma = 0$ for all ablation experiments on ADE20K by default if not specified. 241

Combination with lighter vision transformer architecture. We report the results of applying 242 our method to lighter vision transformer backbones such as ViT-B/16, in Table 5. Our approach 243 consistently improves the efficiency of Segmenter+ViT-B/16 at the cost of a slight performance drop 244 without fine-tuning. Specifically speaking, our approach saves more than $26\% \downarrow \text{GFLOPs}$ of a trained 245 "Segmenter+ViT-B/16" with only a slight performance drop from 48.48% to 48.40%, which verifies 246 our method also generalizes to lighter vision transformer architectures. 247

Comparison with uniform downsampling. We compare our method with the simple uniform 248 downsampling scheme, which directly downsamples the input image into a lower resolution. Figure 6 249 summarizes the detailed comparison results. For example, on ADE20K, we downsample the input 250 resolution from 640×640 to smaller resolutions (e.g., 592×592 , 576×576 , 560×560 , and 251 544×544) and report their performance and GFLOPs in Figure 6. We also plot the results with our 252 method and we can see that our method consistently outperforms uniform sampling on both ADE20K 253 and PASCAL-Context under multiple different GFLOPs budgets. 254



sampling on ADE20K and PASCAL-Context with Segmenter+ViT-L/16.

Figure 6: Comparison to uniform down- Figure 7: Illustrating the improvements of our approach on object detection task with SwinV2-L+ HTC++. \uparrow and \downarrow represent higher is better and lower is better respectively.



Figure 8: Illustrating the improvements of our approach on panoptic segmentation and instance segmentation tasks with Mask2Former+Swin-L.



Figure 9: Illustrating the improvements of our approach and the comparisons with EViT [41] on ImageNet classification tasks with SWAG+ViT-H/14 and SWAG+ViT-L/16.

4.4 Object Detection 255

We use the recent SOTA object detection framework SwinV2-L + HTC++ [43] as our baseline. We 256 summarize the results of combining our method with SwinV2-L + HTC++ on COCO object detection 257 and instance segmentation tasks in Figure 7. 258

Implementation details. The original SwinV2-L consists of {2,2,18,2} shifted window transformer 259 blocks across the four stages. We only apply our method to the 3-rd stage with 18 blocks considering 260 it dominates the computation overhead, which we also follow in the "Mask2Former + Swin-L" 261 experiments. We insert the window token clustering/reconstruction layer after the 8-th/18-th block 262 within the 3-rd stage, which are based on the official checkpoints ³ of SwinV2-L + HTC++. In other 263 words, we set $\alpha = 12$ and $\alpha + \beta = 22$ for SwinV2-L. The default window size is K × K=32 × 32 and 264 we set the clustered window size as $k \times k=23 \times 23$. We choose the values of other hyperparameters 265 following the ablation experiments. According to the results summarized in Figure 7, compared to 266 SwinV2-L + HTC++, our method improves the FPS by 21% \uparrow and saves the GFLOPs by nearly 267 $20\% \downarrow$ while maintaining around 98% of object detection & instance segmentation performance. 268

4.5 Semantic/Instance/Panoptic Segmentation 269

270 We first apply our method to a plain ViT-based segmentation framework Segmenter [60] and illustrate 271 the semantic segmentation results across three benchmarks including ADE20K, PASCAL-Context, and Cityscapes on the first row of Figure 1. Then, we apply our method to a very recent framework 272 Mask2Former [12] that is based on Swin Transformer and summarize the instance and panoptic 273 segmentation results on COCO in Figure 8. 274

Implementation details. The original ViT-L first splits an image into a sequence of image patches of 275 size 16×16 and applies a patch embedding layer to increase the channel dimensions to 1024, then 276

³https://github.com/microsoft/Swin-Transformer, MIT License

Table 6: Depth estimation results based on DPT [51] with ResNet-50+ViT-B/16.

Dataset	Method	GFLOPs	FPS	$\delta \!\!>\!\! 1.25$	$\delta > 1.25^2$	$\delta > 1.25^{3}$	AbsRel	SqRel	RMSE	RMSElog	SILog	log10
KITTI	DPT	810	11.38	0.959	0.995	0.999	0.062	0.222	2.573	0.092	8.282	0.027
	DPT+Ours	627	14.75	0.958	0.995	0.999	0.062	0.226	2.597	0.093	8.341	0.027
NYUv2	DPT	560	17.58	0.904	0.988	0.998	0.110	0.054	0.357	0.129	9.522	0.045
	DPT+Ours	404	24.03	0.900	0.987	0.998	0.113	0.056	0.363	0.132	9.532	0.046

applies 24 consecutive transformer encoder layers for representation learning. To apply our method
to the ViT-L backbone of Segmenter, we use the official checkpoints ⁴ of "Segmenter + ViT-L/16"
and insert the token clustering layers and token reconstruction layer into the ViT-L/16 backbone
without fine-tuning.

For the Mask2Former built on Swin-L with window size as 12×12 , we use the official checkpoints ⁵ of "Mask2Former + Swin-L" and insert the window token clustering layer and the window token reconstruction layer into the empirically chosen positions, which first cluster 12×12 tokens into 8×8 tokens and then reconstruct 12×12 tokens within each window. Figure 8 summarizes the detailed comparison results. Accordingly, we can see that our method significantly improves the FPS by more than 35% \uparrow with a slight performance drop on COCO panoptic segmentation task.

287 4.6 Monocular Depth Estimation

To verify the generalization of our approach, we apply our method to depth estimation tasks that measure the distance of each pixel relative to the camera. We choose the DPT (Dense Prediction Transformer) [51] that builds on the hybrid vision transformer, i.e., R50+ViT-B/16, following [16].

Implementation details. The original R50+ViT-B/16 ⁶ consists of a ResNet50 followed by a ViT-B/16, where the ViT-B/16 consists of 12 transformer encoder layers that process $16 \times$ downsampled representations. We insert the token clustering layer & token reconstruction layer into ViT-B/16 and summarize the results on both KITTI and NYUv2 on the second row of Figure 1. We also report their detailed depth estimation results in Table 6, where we can see that our method accelerates DPT by nearly 30%/37% \uparrow on KITTI/NYUv2, respectively.

297 4.7 ImageNet-1K Classification

Finally, we apply our method to the ImageNet-1K classification task and compare our method with a very recent SOTA method EViT [41]. The key idea of EViT is to identify and only keep the top-ktokens according to their attention scores relative to the [class] token. We empirically find that applying EViT for dense prediction tasks directly suffers from significant performance drops. More details are illustrated in the supplementary material.

Implementation details. We choose the recent SWAG [58] as our baseline, which exploits 3.6 billion 303 weakly labeled images associated with around 27K categories (or hashtags) to pre-train the large-scale 304 vision transformer models, i.e., ViT-L/16 and ViT-H/14. According to their official implementations ⁷ 305 SWAG + ViT-H/14 and SWAG + ViT-L/16 achieve 88.55% and 88.07% top-1 accuracy on ImaegNet-306 1K respectively. We apply our approach and EViT to both baselines and summarize the comparison 307 results in Figure 9. According to the results, our method achieves comparable results as EViT while 308 being more efficient, which further verifies that our method also generalizes to the image classification 309 tasks without fine-tuning. 310

311 5 Conclusion

In this paper, we present a simple and effective mechanism to improve the efficiency of large-scale vision transformer models for dense prediction tasks. In light of the relatively high costs associated with re-training or fine-tuning large vision transformer models on various dense prediction tasks, our study provides a very lightweight method for expediting the inference process while requiring no additional fine-tuning. We hope our work could inspire further research efforts into exploring how to accelerate large-scale vision transformers for dense prediction tasks without fine-tuning.

⁴https://github.com/rstrudel/segmenter#model-zoo, MIT License

⁵https://github.com/facebookresearch/Mask2Former/blob/main/MODEL_ZOO.md, CC-BY-NC
4.0

⁶https://github.com/isl-org/DPT, MIT License

⁷https://github.com/facebookresearch/SWAG, CC-BY-NC 4.0

318 **References**

- [1] Hangbo Bao, Li Dong, and Furu Wei. Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*, 2021.
- [2] Davis Blalock, Jose Javier Gonzalez Ortiz, Jonathan Frankle, and John Guttag. What is the state of neural network pruning? *Proceedings of machine learning and systems*, 2:129–146, 2020.
- [3] Walid Bousselham, Guillaume Thibault, Lucas Pagano, Archana Machireddy, Joe Gray,
 Young Hwan Chang, and Xubo Song. Efficient self-ensemble framework for semantic segmen tation. arXiv preprint arXiv:2111.13280, 2021.
- [4] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. COCO-Stuff: Thing and stuff classes in context. In *CVPR*. IEEE, 2018.
- [5] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and
 Sergey Zagoruyko. End-to-end object detection with Transformers. In *ECCV*, 2020.
- [6] Arnav Chavan, Zhiqiang Shen, Zhuang Liu, Zechun Liu, Kwang-Ting Cheng, and Eric Xing.
 Vision transformer slimming: Multi-dimension searching in continuous optimization space.
 arXiv preprint arXiv:2201.00814, 2022.
- [7] Minghao Chen, Houwen Peng, Jianlong Fu, and Haibin Ling. Autoformer: Searching trans formers for visual recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12270–12280, 2021.
- [8] Tianlong Chen, Yu Cheng, Zhe Gan, Lu Yuan, Lei Zhang, and Zhangyang Wang. Chasing
 sparsity in vision transformers: An end-to-end exploration. *Advances in Neural Information Processing Systems*, 34, 2021.
- [9] Wuyang Chen, Xianzhi Du, Fan Yang, Lucas Beyer, Xiaohua Zhai, Tsung-Yi Lin, Huizhong
 Chen, Jing Li, Xiaodan Song, Zhangyang Wang, et al. A simple single-scale vision transformer
 for object localization and instance segmentation. *arXiv preprint arXiv:2112.09747*, 2021.
- [10] Wuyang Chen, Wei Huang, Xianzhi Du, Xiaodan Song, Zhangyang Wang, and Denny Zhou.
 Auto-scaling vision transformers without training. In *International Conference on Learning Representations*, 2021.
- [11] Yinpeng Chen, Xiyang Dai, Dongdong Chen, Mengchen Liu, Xiaoyi Dong, Lu Yuan,
 and Zicheng Liu. Mobile-former: Bridging mobilenet and transformer. *arXiv preprint arXiv:2108.05895*, 2021.
- Bowen Cheng, Ishan Misra, Alexander G. Schwing, Alexander Kirillov, and Rohit Girdhar.
 Masked-attention mask transformer for universal image segmentation. 2022.
- [13] Krzysztof Choromanski, Valerii Likhosherstov, David Dohan, Xingyou Song, Andreea Gane,
 Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, et al. Rethinking
 attention with performers. *arXiv preprint arXiv:2009.14794*, 2020.
- [14] Zihang Dai, Hanxiao Liu, Quoc V Le, and Mingxing Tan. Coatnet: Marrying convolution and
 attention for all data sizes. *Advances in Neural Information Processing Systems*, 34:3965–3977,
 2021.
- [15] Mingyu Ding, Bin Xiao, Noel Codella, Ping Luo, Jingdong Wang, and Lu Yuan. Davit: Dual
 attention vision transformer. *arXiv preprint arXiv:2204.03645*, 2022.
- [16] 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. *arXiv preprint arXiv:2010.11929*, 2020.
- [17] Stéphane d'Ascoli, Hugo Touvron, Matthew L Leavitt, Ari S Morcos, Giulio Biroli, and Levent
 Sagun. Convit: Improving vision transformers with soft convolutional inductive biases. In
 International Conference on Machine Learning, pages 2286–2296. PMLR, 2021.

- [18] Alaaeldin El-Nouby, Hugo Touvron, Mathilde Caron, Piotr Bojanowski, Matthijs Douze,
 Armand Joulin, Ivan Laptev, Natalia Neverova, Gabriel Synnaeve, Jakob Verbeek, et al. Xcit:
 Cross-covariance image transformers. *arXiv preprint arXiv:2106.09681*, 2021.
- [19] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and
 Christoph Feichtenhofer. Multiscale vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6824–6835, 2021.
- [20] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving?
 the kitti vision benchmark suite. In 2012 IEEE conference on computer vision and pattern
 recognition, pages 3354–3361. IEEE, 2012.
- Jiaqi Gu, Hyoukjun Kwon, Dilin Wang, Wei Ye, Meng Li, Yu-Hsin Chen, Liangzhen Lai, Vikas
 Chandra, and David Z Pan. Hrvit: Multi-scale high-resolution vision transformer. *arXiv preprint arXiv:2111.01236*, 2021.
- John Guibas, Morteza Mardani, Zongyi Li, Andrew Tao, Anima Anandkumar, and Bryan
 Catanzaro. Efficient token mixing for transformers via adaptive fourier neural operators. In
 International Conference on Learning Representations, 2021.
- [23] Kai Han, Yunhe Wang, Qiulin Zhang, Wei Zhang, Chunjing Xu, and Tong Zhang. Model
 rubik's cube: Twisting resolution, depth and width for tinynets. *Advances in Neural Information Processing Systems*, 33:19353–19364, 2020.
- [24] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked
 autoencoders are scalable vision learners. *arXiv preprint arXiv:2111.06377*, 2021.
- [25] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image
 Recognition. In *CVPR*, 2016.
- [26] Byeongho Heo, Sangdoo Yun, Dongyoon Han, Sanghyuk Chun, Junsuk Choe, and Seong Joon
 Oh. Rethinking spatial dimensions of vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11936–11945, 2021.
- [27] Torsten Hoefler, Dan Alistarh, Tal Ben-Nun, Nikoli Dryden, and Alexandra Peste. Sparsity
 in deep learning: Pruning and growth for efficient inference and training in neural networks.
 Journal of Machine Learning Research, 22(241):1–124, 2021.
- [28] Zejiang Hou and Sun-Yuan Kung. Multi-dimensional model compression of vision transformer.
 arXiv preprint arXiv:2201.00043, 2021.
- Jitesh Jain, Anukriti Singh, Nikita Orlov, Zilong Huang, Jiachen Li, Steven Walton, and
 Humphrey Shi. Semask: Semantically masked transformers for semantic segmentation. *arXiv preprint arXiv:2112.12782*, 2021.
- [30] Varun Jampani, Deqing Sun, Ming-Yu Liu, Ming-Hsuan Yang, and Jan Kautz. Superpixel
 sampling networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*,
 pages 352–368, 2018.
- [31] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation
 learning with noisy text supervision. In *International Conference on Machine Learning*, pages
 404 4904–4916. PMLR, 2021.
- [32] Kumara Kahatapitiya and Michael S Ryoo. Swat: Spatial structure within and among tokens.
 arXiv preprint arXiv:2111.13677, 2021.
- [33] Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer.
 arXiv preprint arXiv:2001.04451, 2020.
- [34] Zhenglun Kong, Peiyan Dong, Xiaolong Ma, Xin Meng, Wei Niu, Mengshu Sun, Bin Ren,
 Minghai Qin, Hao Tang, and Yanzhi Wang. Spvit: Enabling faster vision transformers via soft
 token pruning. *arXiv preprint arXiv:2112.13890*, 2021.

- [35] Liunian Harold Li*, Pengchuan Zhang*, Haotian Zhang*, Jianwei Yang, Chunyuan Li, Yiwu
 Zhong, Lijuan Wang, Lu Yuan, Lei Zhang, Jenq-Neng Hwang, Kai-Wei Chang, and Jianfeng
 Gao. Grounded language-image pre-training. In *CVPR*, 2022.
- [36] Yanghao Li, Hanzi Mao, Ross Girshick, and Kaiming He. Exploring plain vision transformer
 backbones for object detection. *arXiv preprint arXiv:2203.16527*, 2022.
- [37] Yanghao Li, Chao-Yuan Wu, Haoqi Fan, Karttikeya Mangalam, Bo Xiong, Jitendra Malik,
 and Christoph Feichtenhofer. Improved multiscale vision transformers for classification and
 detection. *arXiv preprint arXiv:2112.01526*, 2021.
- [38] Yanghao Li, Saining Xie, Xinlei Chen, Piotr Dollar, Kaiming He, and Ross Girshick. Bench marking detection transfer learning with vision transformers. *arXiv preprint arXiv:2111.11429*, 2021.
- [39] Zhenyu Li, Zehui Chen, Xianming Liu, and Junjun Jiang. Depthformer: Depthformer: Exploiting long-range correlation and local information for accurate monocular depth estimation.
 arXiv preprint arXiv:2203.14211, 2022.
- [40] Jingyun Liang, Jiezhang Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. Swinir:
 Image restoration using swin transformer. *arXiv preprint arXiv:2108.10257*, 2021.
- [41] Youwei Liang, GE Chongjian, Zhan Tong, Yibing Song, Jue Wang, and Pengtao Xie. Evit:
 Expediting vision transformers via token reorganizations. In *International Conference on Learning Representations*, 2021.
- [42] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr
 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.
- [43] Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng
 Zhang, Li Dong, et al. Swin transformer v2: Scaling up capacity and resolution. *arXiv preprint arXiv:2111.09883*, 2021.
- [44] 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, pages 10012–10022, 2021.
- [45] Zhuang Liu, Jianguo Li, Zhiqiang Shen, Gao Huang, Shoumeng Yan, and Changshui Zhang.
 Learning efficient convolutional networks through network slimming. In *Proceedings of the IEEE international conference on computer vision*, pages 2736–2744, 2017.
- [46] Zhuang Liu, Mingjie Sun, Tinghui Zhou, Gao Huang, and Trevor Darrell. Rethinking the value of network pruning. *arXiv preprint arXiv:1810.05270*, 2018.
- [47] Dmitrii Marin, Jen-Hao Rick Chang, Anurag Ranjan, Anish Prabhu, Mohammad Rastegari, and
 Oncel Tuzel. Token pooling in vision transformers. *arXiv preprint arXiv:2110.03860*, 2021.
- ⁴⁴⁷ [48] Sachin Mehta and Mohammad Rastegari. Mobilevit: light-weight, general-purpose, and mobile-⁴⁴⁸ friendly vision transformer. *arXiv preprint arXiv:2110.02178*, 2021.
- [49] Roozbeh Mottaghi, Xianjie Chen, Xiaobai Liu, Nam-Gyu Cho, Seong-Whan Lee, Sanja Fidler,
 Raquel Urtasun, and Alan L. Yuille. The role of context for object detection and semantic
 segmentation in the wild. In *CVPR*, 2014.
- [50] 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*,
 pages 8748–8763. PMLR, 2021.
- [51] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction.
 In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12179–12188, 2021.

- 459 [52] Yongming Rao, Wenliang Zhao, Benlin Liu, Jiwen Lu, Jie Zhou, and Cho-Jui Hsieh. Dy 460 namicvit: Efficient vision transformers with dynamic token sparsification. *arXiv preprint* 461 *arXiv:2106.02034*, 2021.
- [53] Cedric Renggli, André Susano Pinto, Neil Houlsby, Basil Mustafa, Joan Puigcerver, and Carlos
 Riquelme. Learning to merge tokens in vision transformers. *arXiv preprint arXiv:2202.12015*, 2022.
- [54] Aurko Roy, Mohammad Saffar, Ashish Vaswani, and David Grangier. Efficient content-based
 sparse attention with routing transformers. *Transactions of the Association for Computational Linguistics*, 9:53–68, 2021.
- [55] Michael S Ryoo, AJ Piergiovanni, Anurag Arnab, Mostafa Dehghani, and Anelia Angelova.
 Tokenlearner: What can 8 learned tokens do for images and videos? *arXiv preprint arXiv:2106.11297*, 2021.
- [56] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation
 and support inference from rgbd images. In *European conference on computer vision*, pages
 746–760. Springer, 2012.
- ⁴⁷⁴ [57] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale ⁴⁷⁵ image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [58] Mannat Singh, Laura Gustafson, Aaron Adcock, Vinicius de Freitas Reis, Bugra Gedik, Raj Prateek Kosaraju, Dhruv Mahajan, Ross Girshick, Piotr Dollár, and Laurens van der Maaten.
 Revisiting weakly supervised pre-training of visual perception models, 2022.
- [59] Lin Song, Songyang Zhang, Songtao Liu, Zeming Li, Xuming He, Hongbin Sun, Jian Sun,
 and Nanning Zheng. Dynamic grained encoder for vision transformers. *Advances in Neural Information Processing Systems*, 34, 2021.
- [60] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Segmenter: Transformer for
 semantic segmentation. *arXiv preprint arXiv:2105.05633*, 2021.
- [61] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov,
 Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions.
 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9,
 2015.
- ⁴⁸⁸ [62] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural ⁴⁸⁹ networks. In *International conference on machine learning*, pages 6105–6114. PMLR, 2019.
- [63] Yi Tay, Mostafa Dehghani, Dara Bahri, and Donald Metzler. Efficient transformers: A survey.
 arXiv preprint arXiv:2009.06732, 2020.
- [64] 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.
- [65] Tao Wang, Li Yuan, Yunpeng Chen, Jiashi Feng, and Shuicheng Yan. Pnp-detr: towards efficient
 visual analysis with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4661–4670, 2021.
- [66] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping
 Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction
 without convolutions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 568–578, 2021.
- [67] Wenxiao Wang, Minghao Chen, Shuai Zhao, Long Chen, Jinming Hu, Haifeng Liu, Deng Cai,
 Xiaofei He, and Wei Liu. Accelerate cnns from three dimensions: A comprehensive pruning
 framework. In *International Conference on Machine Learning*, pages 10717–10726. PMLR,
 2021.

- [68] Wenxiao Wang, Shuai Zhao, Minghao Chen, Jinming Hu, Deng Cai, and Haifeng Liu.
 Dbp: discrimination based block-level pruning for deep model acceleration. *arXiv preprint arXiv:1912.10178*, 2019.
- [69] Yulin Wang, Rui Huang, Shiji Song, Zeyi Huang, and Gao Huang. Not all images are worth
 16x16 words: Dynamic transformers for efficient image recognition. *Advances in Neural Information Processing Systems*, 34, 2021.
- [70] Mitchell Wortsman, Gabriel Ilharco, Samir Yitzhak Gadre, Rebecca Roelofs, Raphael Gontijo Lopes, Ari S Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, et al.
 Model soups: averaging weights of multiple fine-tuned models improves accuracy without
 increasing inference time. *arXiv preprint arXiv:2203.05482*, 2022.
- [71] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo.
 Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in Neural Information Processing Systems*, 34, 2021.
- [72] Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai,
 and Han Hu. Simmim: A simple framework for masked image modeling. *arXiv preprint arXiv:2111.09886*, 2021.
- [73] Yufei Xu, Jing Zhang, Qiming Zhang, and Dacheng Tao. Vitpose: Simple vision transformer
 baselines for human pose estimation. *arXiv preprint arXiv:2204.12484*, 2022.
- Jianwei Yang, Chunyuan Li, Pengchuan Zhang, Xiyang Dai, Bin Xiao, Lu Yuan, and Jianfeng
 Gao. Focal self-attention for local-global interactions in vision transformers. *arXiv preprint arXiv:2107.00641*, 2021.
- ⁵²⁷ [75] Jiahui Yu, Linjie Yang, Ning Xu, Jianchao Yang, and Thomas Huang. Slimmable neural networks. *arXiv preprint arXiv:1812.08928*, 2018.
- ⁵²⁹ [76] Yuhui Yuan, Xilin Chen, and Jingdong Wang. Object-contextual representations for semantic ⁵³⁰ segmentation. In *European conference on computer vision*, pages 173–190. Springer, 2020.
- [77] Yuhui Yuan, Rao Fu, Lang Huang, Weihong Lin, Chao Zhang, Xilin Chen, and Jingdong Wang.
 Hrformer: High-resolution transformer for dense prediction. *arXiv preprint arXiv:2110.09408*, 2021.
- [78] Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transform ers. *CVPR*, 2022.
- [79] Bowen Zhang, Shuyang Gu, Bo Zhang, Jianmin Bao, Dong Chen, Fang Wen, Yong Wang, and
 Baining Guo. Styleswin: Transformer-based gan for high-resolution image generation. *arXiv preprint arXiv:2112.10762*, 2021.
- [80] Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun Zhu, Lionel M Ni, and Heung-Yeung Shum. Dino: Detr with improved denoising anchor boxes for end-to-end object detection.
 arXiv preprint arXiv:2203.03605, 2022.
- [81] Haokui Zhang, Wenze Hu, and Xiaoyu Wang. Edgeformer: Improving light-weight convnets by
 learning from vision transformers. *arXiv preprint arXiv:2203.03952*, 2022.
- [82] Pengchuan Zhang, Xiyang Dai, Jianwei Yang, Bin Xiao, Lu Yuan, Lei Zhang, and Jianfeng Gao.
 Multi-scale vision longformer: A new vision transformer for high-resolution image encoding. In
 Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 2998–3008,
 2021.
- [83] Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio
 Torralba. Semantic understanding of scenes through the ADE20K dataset. *IJCV*, 2019.

550 Checklist

551	1. For all authors
552 553	(a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
554 555	(b) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
556 557	(c) Did you discuss any potential negative societal impacts of your work? [N/A](d) Did you describe the limitations of your work? [N/A]
558	2. If you are including theoretical results
559	(a) Did you state the full set of assumptions of all theoretical results? [N/A]
560	(b) Did you include complete proofs of all theoretical results? [N/A]
561	3. If you ran experiments
562 563 564	(a) Did you include the code, data, and instructions needed to reproduce the main exper- imental results (either in the supplemental material or as a URL)? [TODO]We will release the code soon.
565 566	(b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes]
567 568	(c) Did you report error bars (e.g., with respect to the random seed after running experi- ments multiple times)? [TODO]
569 570	(d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes]
571	4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets
572 573	(a) If your work uses existing assets, did you cite the creators? [Yes] We use the official checkpoints provided by:
574	 https://github.com/rstrudel/segmenter#model-zoo, MIT License
575 576	 https://github.com/facebookresearch/Mask2Former/blob/main/MODEL_ZOO. md. CC-BY-NC 4.0
577	 https://github.com/microsoft/Swin-Transformer, MIT License
578	• https://github.com/facebookresearch/SWAG, CC-BY-NC 4.0
579	 https://github.com/isl-org/DPT, MIT License.
580 581	(b) Did you mention the license of the assets? [Yes] We mark the license of these assets as above.
582 583	(c) Did you include any new assets either in the supplemental material or as a URL? [Yes] We add the URL of these assets in the footnote.
584 585	(d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
586 587	(e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
588	5. If you used crowdsourcing or conducted research with human subjects
589 590	(a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
591 592	(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
593 594	(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]