

001

SimPLR: A Simple and Plain Transformer for 002 Object Detection and Segmentation

003

003 Anonymous ECCV 2024 Submission

004 Paper ID #4126

005 **Abstract.** The ability to detect objects in images at varying scales has played a
006 pivotal role in the design of modern object detectors. Despite considerable progress
007 in removing hand-crafted components and simplifying the architecture with trans-
008 formers, multi-scale feature maps and/or pyramid design remain a key factor for
009 their empirical success. In this paper, we show that this reliance on either feature
010 pyramids or an hierarchical backbone is unnecessary and a transformer-based
011 detector with scale-aware attention enables the plain detector ‘SimPLR’ whose
012 backbone *and* detection head are both non-hierarchical and operate on single-scale
013 features. The plain architecture allows SimPLR to effectively take advantages of
014 self-supervised learning and scaling approaches with ViTs, yielding competitive
015 performance compared to hierarchical and multi-scale counterparts. We demon-
016 strate through our experiments that when scaling to larger ViT backbones, SimPLR
017 indicates better performance than end-to-end segmentation models (Mask2Former)
018 and plain-backbone detectors (ViTDet), while consistently being faster. The code
019 will be released.

020

1 Introduction

021 After its astonishing achievements in natural language processing, the transformer [39]
022 has quickly become the neural network architecture of choice in computer vision, as
023 evidenced by recent success in image classification [13, 31], object detection [3, 33, 47]
024 and segmentation [7, 40, 46]. Unlike natural language processing, where the same pre-
025 trained network can be deployed for a wide range of downstream tasks with only minor
026 modifications [2, 11], computer vision tasks such as object detection and segmentation
027 require a different set of domain-specific knowledge to be incorporated into the network.
028 Consequently, it is commonly accepted that a modern object detector contains two
029 main components: a pre-trained backbone as the *general* feature extractor, and a *task-*
030 *specific* head that conducts detection and segmentation tasks using domain knowledge.
031 For transformer-based vision architectures, the question remains whether to add more
032 inductive biases or to learn them from data.

033 The spatial nature of image data lies at the core of computer vision. Besides learning
034 long range feature dependencies, the ability of capturing local structure of neighboring
035 pixels is critical for representing and understanding the image content. Building upon
036 the successes of convolutional neural networks, a line of research biases the transformer
037 architecture to be *multi-scale* and *hierarchical* when dealing with the image input, *i.e.*,
038 Swin Transformer [31] and others [14, 20, 41]. The hierarchical design makes it easy to
039 create multi-scale features for dense vision tasks and allows pre-trained transformers to

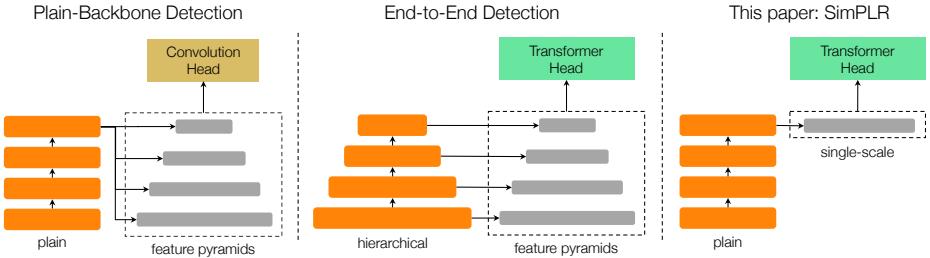


Fig. 1: Object detection architectures. **Left:** The plain-backbone detector from [26] whose input (denoted in the dashed region) are multi-scale features. **Middle:** State-of-the-art end-to-end detectors [7, 33] utilize a hierarchical backbone (*i.e.*, Swin [31]) to create multi-scale inputs. **Right:** Our simple single-scale detector following the end-to-end framework. Where existing detectors require feature pyramids to be effective, we propose a plain detector, SimPLR, whose backbone and detection head are non-hierarchical and operate on a single-scale feature map. The plain detector, SimPLR, achieves on par or even better performance compared to hierarchical and/or multi-scale counterparts while being more efficient.

040 be seamlessly integrated into a convolution-based detection head with a feature pyramid
 041 network [27], yielding impressive results in object detection and segmentation. However,
 042 the inductive biases in the architectural design make it benefit less from self-supervised
 043 learning and the scaling of model size [26].

044 An alternative direction pursues the idea of a simple transformer with “less inductive
 045 biases” and emphasizes learning vision-specific knowledge directly from image data.
 046 Specifically, the Vision Transformer (ViT) [13] stands out as a *plain* architecture with a
 047 constant feature resolution, and acts as the feature extractor in plain-backbone detection.
 048 This is motivated by the success of ViTs scaling behaviours in visual recognition [1, 9, 17].
 049 In addition, the end-to-end detection framework proposed by Carion *et al.* [3] with a
 050 transformer-based detection head further removes many hand-designed components, like
 051 non-maximum suppression or intersection-over-union computation, that encodes the
 052 prior knowledge for object detection.

053 The plain design of ViTs, however, casts doubts about its ability to capture information
 054 of objects across multiple scales. While recent studies [13, 26] suggest that
 055 ViTs with global self-attention could potentially learn translation-equivariant and scale-
 056 equivariant properties during training, leading object detectors still require multi-scale
 057 feature maps and/or an hierarchical backbone. This observation holds true for both
 058 convolutional [18, 26, 35] and transformer-based detectors [7, 33, 47] (see ??). Unlike
 059 hierarchical backbones, the creation of feature pyramids conflicts with the original
 060 design philosophy of ViTs. Therefore, our goal is to pursue a plain detector whose
 061 backbone *and* detection head are both *single-scale* and *non-hierarchical*. This further
 062 simplifies the architecture for detection and segmentation in the pursuit of learning object
 063 representations from data.

064 In this paper, we introduce SimPLR, a plain detector for *both* end-to-end detection
 065 and segmentation frameworks [3, 7, 33, 47]. In particular, our detector extracts the
 066 single-scale feature map from a ViT backbone, which is then fed into the transformer

encoder-decoder via a simple projection to make the prediction. To deal with objects of various sizes, we propose to incorporate scale information into the attention mechanism, resulting in an *adaptive-scale* attention. The proposed attention mechanism learns to capture adaptive scale distribution from training data. This eliminates the need for multi-scale feature maps from the ViT backbone, yielding a simple and efficient detector.

The proposed detector, SimPLR, turns out to be an effective solution for the plain detection and segmentation. We find that multi-scale feature maps are not necessary and the scale-aware attention mechanism adequately captures objects at various sizes from output features of a ViT backbone. Despite the plain architecture, our detector (SimPLR) shows competitive performance compared to the strong hierarchical-backbone or multi-scale detectors (*e.g.*, ViTDet [26] and Mask2Former [7]), while being consistently faster. Moreover, the effectiveness of our detector is observed not only in object detection but also in instance and panoptic segmentation. Interestingly, the efficient design allows SimPLR to take advantages of the significant progress in self-supervised learning and scaling with ViTs (*e.g.*, with MAE [17] and BEiT [34]), indicating plain detectors to be a promising direction for dense prediction tasks.

2 Related Work

Backbones for object detection. Inspired by R-CNN [16], modern object detection methods utilize a task-specific head on top of a pre-trained backbone. Initially, object detectors [19, 21, 36, 42] were dominated by convolutional neural network (CNN) backbones [24] pre-trained on ImageNet [10]. With the success of the transformer in learning from large-scale text data [2, 11], many studies have explored the transformer for computer vision [4, 13, 31]. Most recently, the Vision Transformer (ViT) [13] with a simple design demonstrated the capability of learning meaningful representation for visual recognition. By removing the need for labels, methods with self-supervised learning have emerged as an even more powerful solution for pre-training general vision representations [5, 17]. We show through experiments that SimPLR can take advantages of the significant progress in representation learning and scaling of ViTs.

End-to-end detection and segmentation. The end-to-end framework for object detection proposed in DETR [3] aims to remove the need for many hand-crafted modules. This is made possible by adopting Transformer as the detection head to directly give the prediction. Follow-up works [12, 33, 40] extended the transformer-based head in end-to-end frameworks for instance segmentation, and panoptic segmentation. This inspired MaskFormer [8] and K-Net [46] to unify segmentation tasks with a class-agnostic mask prediction. Pointing out that MaskFormer and K-Net lag behind specialized architectures, Cheng *et al.* [7] introduce Mask2Former, reaching strong performance on segmentation tasks. Yu *et al.* [43] replace self-attention with k -means clustering, boosting the effectiveness of the network. Another direction is to improve the object query in the decoder via a denoising process [25, 45]. While simplifying the detection and segmentation framework, these architectures still require an hierarchical backbone along with feature pyramids. The use of feature pyramids increases the sequence length of the input to the transformer-based detection head, making the detector less efficient. In this work, we enable a plain detector by removing hierarchical and multi-scale constraints.

Plain detectors. Following the goal of less inductive biases in the architecture, recent studies focus on the non-hierarchical and single-scale detector. Motivated by the success of ViT, a line of research considers plain-backbone detectors that replace the hierarchical backbone with a ViT. Initially, Chen *et al.* [6] present UViT as a plain detector that contains a ViT backbone and a single-scale convolutional detection head. Since the backbone architecture is modified *during pre-training* to adopt the progressive window attention, UViT is unable to take the advantages of existing pre-training approaches with ViTs. ViTDet [26] tackles this problem with simple adaptations of the ViT backbone *during fine-tuning*. These simple modifications allow ViTDet to benefit directly from recent self-supervised learning with ViTs (*i.e.*, MAE [17]), resulting in strong results when scaling to larger models. Despite enabling plain-backbone detectors, feature pyramids are still an important factor in ViTDet to detect object at various scales.

Most recently, Lin *et al.* [29] introduce the transformer-based PlainDETR detector, which also removes the multi-scale input. However, it still relies on multi-scale features to generate the object proposals for its decoder. In the decoder, PlainDETR also uses hybrid matching [22] to strengthen its prediction, while our decoder preserves a simple design as in [33, 47]. We believe to be the first to remove the hierarchical and multi-scale constraints which appear in the backbone *and* the input of the transformer encoder for *both* detection and segmentation tasks. Our proposed scale-aware attention can further plug into current end-to-end frameworks without significant architectural changes.

3 SimPLR: A Simple and Plain Detector

Multi-scale feature maps in a hierarchical backbone can be easily extracted from the pyramid structure [27, 30, 47]. When moving to a ViT backbone with a constant feature resolution, the creation of multi-scale feature maps requires complex backbone adaptations. Moreover, the benefits of multi-scale features in object detection frameworks using ViTs remain unclear. Recent studies on plain-backbone detection [6, 26] conjecture the high-dimensional ViT with self-attention and positional embeddings [39] is able to preserve important information for localizing objects¹. From this conjecture, we hypothesize that a proper design of the transformer-based head will enable a plain detector.

Our proposed detector, SimPLR, is conceptually simple: a pre-trained ViT backbone to extract plain features from an image, which are then fed into a single-scale encoder-decoder to make the final prediction (See Fig. 1). Thus, SimPLR is a natural idea as it eliminates the non-trivial creation of feature pyramids from the ViT backbone. But the single-scale encoder-decoder requires an effective design to deal with objects at different scales. First, we review box-attention in [33] as our baseline in end-to-end detection and segmentation using feature pyramids. Then, we introduce the key elements of our plain detector, SimPLR, including its *scale-aware* attention that is the main factor for learning of adaptive object scales.

¹ ViT-B and larger ($\text{dim} \geq 768$) can maintain information with a patch size of $16 \times 16 \times 3$ in the input image.

149 **3.1 Background**

149

150 Our goal is to further simplify the detection and segmentation pipeline from [26, 33,
 151 47], and to prove the effectiveness of the plain detector in *both* object detection and
 152 segmentation tasks. As a result, we utilize the sparse attention mechanism, box-attention
 153 in [33], as strong baselines due to its effectiveness in learning discriminative object
 154 representations while being lightweight in computation.

150

155 In box-attention, each query vector $q \in \mathbb{R}^d$ in the input feature map is assigned
 156 a reference window $r = [x, y, w, h]$, where x, y indicate the query coordinate and w, h
 157 are the size of the reference window. The box-attention refines the reference window
 158 into a region of interest. During the attention head computation, a 2×2 feature grid is
 159 sampled from the corresponding region of interest, resulting in a set of value features
 160 $v_i \in \mathbb{R}^{2 \times 2 \times d_h}$. The 2×2 attention scores are efficiently generated by computing a dot-
 161 product between $q \in \mathbb{R}^d$ and relative position embeddings ($k_i \in \mathbb{R}^{2 \times 2 \times d}$) followed by a
 162 softmax function. The attended feature head $_i \in \mathbb{R}^{d_h}$ is a weighted sum of 2×2 value
 163 features in v_i with the corresponding attention weights: To capture objects at different
 164 scales, the box-attention [33] takes t multi-scale feature maps, $\{e^j\}_{j=1}^t$, as its inputs
 165 in order to produce head $_i$. The multi-scale box-attention shows strong performance in
 166 end-to-end object detection and instance segmentation.

155

167 The sparse attention like box-attention lies at the core of recent end-to-end detection
 168 and segmentation models due to its ability of capturing object information with lower
 169 complexity. The effectiveness and efficiency of these attention mechanisms bring up
 170 the question: *Is multi-scale object information learnable within the detector which is*
 171 *non-hierarchical and single-scale?*

167

172 **3.2 Scale-aware attention**

172

173 The output features of the encoder should capture objects at different scales. Therefore,
 174 unlike the feature pyramids where each set of features encode a specific scale, predicting
 175 objects from a plain feature map requires its feature vectors to reason about dynamic
 176 scale information based on the image content. This can be addressed effectively by a
 177 multi-head attention mechanism that capture different scale information in each of its
 178 attention heads. In that case, global self-attention is a potential candidate because of its
 179 large receptive field and powerful representation. However, its computational complexity
 180 is quadratic w.r.t. the sequence length of the input, making the operation computationally
 181 expensive when dealing with high-resolution images. The self-attention also leads to
 182 worse performance and slow convergence in end-to-end detection [47]. This motivated
 183 us to develop a multi-head *scale-aware* attention mechanism for single-scale input.

173

184 **Scale-aware attention.** In the sparse attention mechanism such as deformable atten-
 185 tion [47] or box-attention [33], each feature vector is assigned to a single scale in feature
 186 pyramids. As a result, feature vectors learn to adapt to that specific scale assigned to
 187 them. While this behaviour may not impact the multi-scale deformable attention or
 188 box-attention – which utilizes feature pyramids for detecting objects – it poses a big
 189 challenge in learning scale-equivariant features on a single-scale input.

174

190 To address this limitation, we propose two variants of multi-head *scale-aware* at-
 191 tention (*i.e.*, *fixed-scale* and *adaptive-scale*) that integrate different scales into each

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

attention head, allowing query vectors to choose the suitable scale information during training. Our proposed attention mechanism is simple: we assign reference windows of m different scales to attention heads of each query. We use m reference windows with size $w=h \in \{s \cdot 2^j\}_{j=0}^{m-1}$, where s is the size of the smallest window, and m is the number of scales. Surprisingly, our experiments show that the results are *not* sensitive to the size of the window, as long as *enough number of scales* are used.

- i) *Fixed-Scale Attention.* Given reference windows of m scales, we distribute them to n attention heads in a round-robin manner. Thus, in multi-head fixed-scale attention, $\frac{n}{m}$ attention heads are allocated for each of the window scales. This uniform distribution of different scales enables fixed-scale attention to learn diverse information from local to more global context. The aggregation of n heads results in scale-aware features, that is suitable for predicting objects of different sizes.
- ii) *Adaptive-Scale Attention.* Instead of uniformly assigning m scales to n attention heads, the adaptive-scale attention learns to allocate a scale distribution based on the context of the query vector. This comes from the motivation that the query vector belonging to a small object should use more attention heads for capturing fine-grained details rather than global context, and vice versa.

Given the query vector $q \in \mathbb{R}^d$ in the input feature map and m reference windows of different scales, $\{r^j\}_{j=0}^{m-1}$, the adaptive-scale attention predicts offsets of all reference windows, $\{\Delta_{x_j}, \Delta_{y_j}, \Delta_{w_j}, \Delta_{h_j}\}_{j=0}^{m-1}$, in each attention head. Besides, we apply a scale temperature to each set of offsets before the transformations:

$$F_{\text{scale}}(r_j, q) = [x, y, w + \Delta_w \cdot \frac{2^j}{\lambda}, h + \Delta_h \cdot \frac{2^j}{\lambda}], \quad (1)$$

$$F_{\text{translate}}(r_j, q) = [x + \Delta_x \cdot \frac{2^j}{\lambda}, y + \Delta_y \cdot \frac{2^j}{\lambda}, w, h], \quad (2)$$

where $\frac{2^j}{\lambda}$ is the scale temperature corresponding to r_j . The scale temperature allows the transformation functions to capture regions of interest corresponding to the scale of reference windows. It then samples feature grids from m regions of interest and generates attention scores for these feature grids followed by softmax normalization. This makes our attention mechanism to focus on feature grids of suitable scale. The adaptive-scale attention provides efficiency due to sparse sampling and strong flexibility to control scale distribution via its attention computation.

3.3 Object Representation for Panoptic Segmentation

Panoptic segmentation proposed by [23] requires the network to segment both “thing” and “stuff”. To enable the plain detector on panoptic segmentation, we make an adaptation in the mask prediction of SimPLR. Following [7], we predict segmentation masks of both types by computing the dot-product between object queries and a high-resolution feature map (*i.e.*, $\frac{1}{4}$ feature scale).

As the ViT and SimPLR encoder features are of lower resolution, we simply interpolate the last encoder layer to $\frac{1}{4}$ scale and add a single scale-aware attention layer on top. This simple modification produces a high resolution feature map that is beneficial for learning fine-grained details.

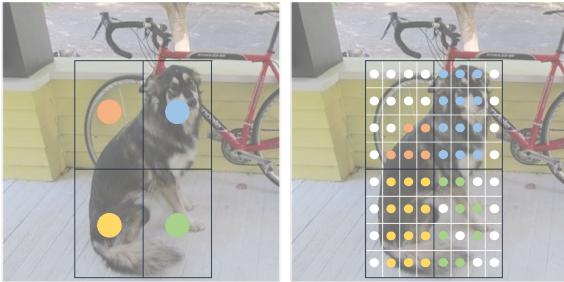


Fig. 2: Masked Instance-Attention. **Left:** The box-attention [33] which samples 2×2 grid features in the region of interest. **Right:** Our masked instance-attention for dense grid sampling. The 2×2 attention scores are denoted in four colours and the masked attention scores are in white.

Masked instance-attention. In order to learn better object representation in the decoder, we propose masked instance-attention that is also a sparse attention mechanism for efficiency. The masked instance-attention follows the grid sampling strategy of the box-attention in [33], but improves the attention computation to better capture objects of different shapes.

To be specific, the region of interest r'_i is divided into 4 bins of 2×2 grid, each of which contains a $\frac{m}{2} \times \frac{m}{2}$ grid features sampled using bilinear interpolation. Instead of assigning an attention weight to each feature vector, a linear projection ($\mathbb{R}^d \rightarrow \mathbb{R}^{2 \times 2}$) is adopted to generate the 2×2 attention scores for 4 bins. The $\frac{m}{2} \times \frac{m}{2}$ feature vectors within the same bin share the same attention weight. Inspired by [7], we utilize the mask prediction of the previous decoder layer as the attention mask to guide the attention scores to better capture the boundary of objects (see Fig. 2).

$$\text{head}_i = \sum_{k=0}^{2 \times 2} \sum_{j=0}^{\frac{m}{2} \times \frac{m}{2}} \frac{\alpha_k}{\frac{m}{2} \cdot \frac{m}{2}} v_{i_{k,j}}, \quad (3)$$

where α_k is the attention weight corresponding to k -th bin, $v_{i_{k,j}}$ is the j -th feature vector inside k -th bin.

3.4 Network Architecture

SimPLR follows the end-to-end detection and segmentation framework in [33] with the two-stage design. Specifically, we use a plain ViT as the backbone with 14×14 windowed attention and four equally-distributed global attention blocks as in [26]. In the detection head, the SimPLR encoder receives input features via a projection of the last feature map from the ViT backbone. The object proposals are then generated using single-scale features from the encoder and top-scored features are initialized as object queries for the SimPLR decoder to predict bounding boxes and masks.

Formally, we apply a projection f to the last feature map of the pre-trained ViT backbone, resulting in the input feature map $e \in \mathbb{R}^{H_e \times W_e \times d}$ where H_e, W_e are the size of the feature map, and d is the hidden dimension of the detection head. In SimPLR, the projection f is simply a single convolution projection, that provides us a great flexibility to control the resolution and dimension of the input features e to the encoder. The projection allows SimPLR to decouple the feature scale and dimension between its

262 backbone and detection head to further improve the efficiency. This practice is different
 263 from the creation of SimpleFPN in [26] where a different stack of multiple convolution
 264 layers is used for each feature scale (more details shown in Fig. A in the supplementary
 265 material). We show by experiments that this formulation is key for plain detection and
 266 segmentation while keeping our network efficient.

267 **Plain backbone.** SimPLR deploys ViT as its plain backbone for feature extraction. We
 268 show that SimPLR can take advantages of recent progress in self-supervised learning with
 269 ViTs. To be specific, SimPLR generalizes to ViT backbones initialized by MAE [17] and
 270 BEiT v2 [34]. The efficient design of SimPLR allows us to effectively scale to larger ViT
 271 backbones which recently show to be even more powerful in learning representations [9,
 272 17, 44]. We provide more details in the supplementary material.

273 4 Experiments

274 **Experimental setup.** In this study, we evaluate our method on COCO [28], a commonly
 275 used dataset for object detection, instance segmentation, and panoptic segmentation tasks.
 276 By default, we use plain features with adaptive-scale attention as described in Sec. 3
 277 due to its strong performance and ability to perform both detection and segmentation;
 278 and initialize the ViT backbone from MAE [17] pre-trained on ImageNet-1K without
 279 any labels. In both fixed-scale and adaptive-scale attention, we set the number of scale
 280 $m = 4$ and the window size $s = 32$. Unless specified, the hyper-parameters are the same
 281 as in [33].

282 For all experiments, our optimizer is AdamW [32] with a learning rate of 0.0001.
 283 The learning rate is linearly warmed up for the first 250 iterations and decayed at 0.9 and
 284 0.95 fractions of the total number of training steps by a factor 10. ViT-B [13] is set as
 285 the backbone. The input image size is 1024×1024 with large-scale jitter [15] between
 286 a scale range of $[0.1, 2.0]$. We employ query denoising [45] when compared with other
 287 methods. Due to the limit of our computational resources, we report the ablation study
 288 using the standard $5 \times$ schedule setting with a batch size of 16 as in [33]. In the main
 289 experiments, we use the finetuning recipe from [26].

290 **SimPLR is an effective single-scale detector.** In Tab. 1, we show the comparison
 291 between SimPLR and recent object detectors using the plain backbone ViT. We also
 292 implement a strong baseline of DeformableDETR [47] with SimpleFPN under our
 293 end-to-end framework for better comparison. The plain detector, SimPLR, with both
 294 scale-aware box-attention (SAB) and scale-aware deformable attention (SAD) removes
 295 the need for multi-scale adaptation of the ViT.

296 While SimPLR with scale-aware deformable attention lags behind its multi-scale
 297 counterpart from our implementation, we observe a much smaller gap compared to
 298 standard deformable attention on single-scale input (*e.g.*, 0.3 vs. 2 AP point). When
 299 equipped with scale-aware box-attention, SimPLR reaches similar performance as BoxeR
 300 and outperforms other multi-scale detectors in both detection and segmentation. In
 301 addition, the plain detector is more efficient than multi-scale counterparts. When moving
 302 to larger models with higher dimension, we find that multi-scale detectors like BoxeR
 303 require significant memory optimization and becomes challenging for our computational
 304 resources. We also compare with PlainDETR [29] which is a recent plain detection

	FPS	AP ^b	AP ^m
Feature pyramids			
DETR reported in [29]	-	46.5	n/a
DeformableDETR reported in [29]	12	52.1	n/a
DeformableDETR (our impl.)	12	54.6	n/a
BoxeR [33]	12	55.4	47.7
ViTDet w/ Cascade head [26]	11	54.0	46.7
Plain detector			
PlainDETR [†] [29]	12	53.8	n/a
SimPLR	17	55.7	47.9

†: Multi-scale features are used to generate object proposals.

Table 1: SimPLR is an effective plain detector. All methods use ViT-B as backbone. Methods that take feature pyramids as input employ SimpleFPN with ViT from [26]. Our plain detector, SimPLR, shows competitive performance compared to multi-scale alternatives, while being faster during inference.

method. As discussed in Sec. 2, PlainDETR and our work approach the plain detector in different ways. PlainDETR aims at designing a strong decoder that compensates for single-scale input, while our goal is to learn scale equivariant features in backbone and encoder. Despite the different approaches, both PlainDETR and our work indicate that plain detection holds a great potential.

Ablation of scale-aware attention. Here, our baseline is the standard box-attention from [33] with single-scale feature input directly taken from the last feature of the ViT backbone (denoted as “base”).

From Tab. 2a, we first conclude that *both* scale-aware attention strategies are substantially better than the naïve baseline, increasing AP by up to 1.8 points. We note that while fixed-scale attention distributes 25% of its attention heads into each of the window scales, adaptive-scale attention decides the scale distribution based on the query content. By choosing feature grids from different window scales adaptively, the adaptive-scale attention is able to learn a suitable scale distribution through training data, yielding better performance compared to fixed-scale attention. This is also verified in Fig. 2e where queries corresponding to *small* objects tend to pick reference windows of small sizes for its attention heads. Interestingly, queries corresponding to *medium* and *large* objects pick not only reference windows of their sizes, but also ones of smaller sizes. One of reasons may come from the fact that performing instance segmentation of larger objects still requires the network to faithfully preserve the per-pixel spatial details.

In Tab. 2b, we compare the performance of SimPLR across several sizes (s) of the reference window. They all improve over the baseline, while the choice of a specific base size makes only marginal differences. Our ablation reveals that the number of scales rather than the window size plays an important role to make our network more *scale-aware*. Indeed, in Tab. 2c, the use of 4 or more window scales shows improvement up to 0.8 AP over 2 window scales; and clearly outperforms the naïve baseline. Last,

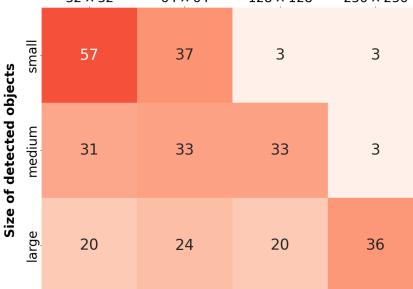
attention	AP ^b	AP ^m
base	53.6	46.1
i) fixed-scale	55.0	47.2
ii) adaptive-scale	55.4	47.6

(a) Scale-aware attention.

s	AP ^b	AP ^m
base	53.6	46.1
16	55.1	47.4
32	55.4	47.6
64	55.1	47.4

(b) Window size.

Adaptive-scale attention (with 4-scale)



(e) Visualization of scale distribution learnt in multi-head adaptive-scale attention of object proposals. Objects are classified into small, medium, and large based on their area.

scale	AP ^b	AP ^m
base	53.6	46.1
1/4	55.4	47.7
1/8	55.4	47.6
1/16	54.3	46.7

(d) Scales of input features.

Table 2: Ablation of scale-aware attention in SimPLR using a plain ViT backbone on COCO val. **Table (a-d):** Compared to the naïve baseline, which employs BoxER and box-attention [33] with *single-scale* features, our plain detector, SimPLR, with scale-aware attention improves the performance consistently for all settings, default setting highlighted. **Figure e:** our adaptive-scale attention captures different scale distribution in its attention heads based on the context of query vectors. Specifically, queries of *small* objects tends to focus on reference windows of small scales (*i.e.*, mainly 32×32), while query vectors of *medium* and *large* objects distribute more attention computation into larger reference windows. All experiments are with $5 \times$ schedule.

we show in Tab. 2d that the *decouple* between feature scale and dimension of the ViT backbone and the detection head features helps to boost the performance of our plain detector by ~ 1 AP point, while keeping its efficiency (*i.e.*, both in terms of FLOPs and FPS). This practice makes scaling of SimPLR to larger ViT backbones more practical.

Ablation on more pre-training data and strategies. Tab. 3 compares the ViT backbone when pre-trained using different strategies with different sizes of pre-training data. SimPLR with the ViT backbone benefits from better pre-training methods even with supervised approaches. Among supervised pre-training methods, DEiTv3 [38] shows better results than DEiT [37], and the pre-training on more data (*i.e.*, ImageNet-21K) further improves the performance of DEiTv3.

However, self-supervised methods like MAE [17] provides strong pre-trained backbones when only pre-trained on ImageNet-1K. The best performance is reported with self-supervised method, BEiTv2 [34] with more pre-training data of ImageNet-21K. This further confirms that our plain detector, SimPLR, enjoys the significant progress of self-supervised learning and scaling ViTs with data. A similar observation is also

pre-train	Object Detection				Instance Segmentation			
	AP ^b	AP _S ^b	AP _M ^b	AP _L ^b	AP ^m	AP _S ^m	AP _M ^m	AP _L ^m
IN-1K, DEiT	53.6	33.7	58.1	71.5	46.1	24.5	50.4	67.2
IN-1K, DEiT _{v3}	54.0	34.3	58.8	70.5	46.4	24.8	51.1	66.7
IN-21K, DEiT _{v3}	54.8	35.4	59.0	72.4	47.1	25.8	51.2	68.5
IN-1K, MAE	55.4	36.1	59.1	70.9	47.6	26.8	51.4	67.1
IN-21K, BEiT _{v2}	55.7	36.5	60.2	72.4	48.1	26.7	52.7	68.9

Table 3: Ablation on scaling with more pre-training data and strategies of SimPLR with the plain ViT-B backbone evaluated on COCO object detection and instance segmentation. We compare the plain backbone pre-trained using supervised methods (*top row*) vs. self-supervised methods (*bottom row*) with different sizes of pre-training dataset (ImageNet-1K vs. ImageNet-21K). Here, we use the $5 \times$ schedule as in [33]. It can be seen that SimPLR with the plain ViT backbone benefits with more pre-training data (*e.g.*, ImageNet-1K vs. ImageNet-21K) and better pre-training approaches (supervised learning vs. self-supervised learning).

346 pointed out in ViTDet [26] where the plain ViT backbone initialized with MAE shows 346
 347 better improvement over hierarchical backbones. 347
 348 **State-of-the-art comparison and scaling behavior.** We show in Tabs. 4 and 5 that 348
 349 SimPLR indicates strong performance on object detection, instance segmentation and 349
 350 panoptic segmentation. To be specific, our plain detector combined with a ViT, pre- 350
 351 trained using MAE [17] or BEiT_{v2} [34], presents good scaling behavior. When moving 351
 352 to large and huge models, our method outperforms multi-scale counterparts including 352
 353 the recent end-to-end Mask2Former segmentation model [7]. Despite involving more 353
 354 advanced attention blocks designs, *i.e.*, shifted window attention in Swin [31] and 354
 355 pooling attention in MViT [14], detectors with hierarchical backbones benefit less from 355
 356 larger backbones. SimPLR is better than the plain-backbone detector, ViTDet, across all 356
 357 backbones in terms of both accuracy and inference speed. 357
 358 **Limitations.** Our final goal is to simplify the detection pipeline and to achieve com- 358
 359 petitive results at the same time. In Sec. 4, we find that the *adaptive-scale* attention 359
 360 mechanism that adaptively learns scale-aware information in its computation plays a 360
 361 key role for a plain detector. However, our adaptive-scale attention still encodes the 361
 362 knowledge of different scales. In the future, we hope that with the large-scale training 362
 363 data, a simpler design of the attention mechanism could also learn the scale equivariant 363
 364 property. Furthermore, SimPLR faces difficulties in detecting and segmenting large 364



Fig. 3: Qualitative results for object detection and panoptic segmentation on the COCO [28] 2017 val set generated by SimPLR. Note that SimPLR gives good predictions on *small* objects.

method	backbone	pre-train	Object Detection				Instance Segmentation				FPS
			AP ^b	AP _S ^b	AP _M ^b	AP _L ^b	AP ^m	AP _S ^m	AP _M ^m	AP _L ^m	
Feature pyramids											
Swin [31]	Swin-L	sup-21K	55.0	38.3	59.4	71.6	47.2	28.7	50.5	66.0	10
Mask2Former [7]	Swin-L	sup-21K	n/a				50.1	29.9	53.9	72.1	4
MViT [14]	MViT-L	sup-21K	55.7	40.3	59.6	71.4	48.3	31.1	51.2	66.3	6
ViTDet [26]	ViT-L	MAE	57.6	40.5	61.6	72.6	49.9	30.5	53.3	68.0	7
MViT [26]	MViT-H	sup-21K	55.9	40.8	59.8	70.8	48.3	30.1	51.1	66.6	6
ViTDet [26]	ViT-H	MAE	58.7	41.9	63.0	73.9	50.9	32.0	54.3	68.9	5
Plain features											
SimPLR	ViT-L	MAE	58.5	42.2	62.5	73.4	50.6	32.1	54.2	69.8	9
SimPLR	ViT-L	BEiTv2	58.7	40.4	63.2	74.8	50.9	30.4	55.1	70.9	9
SimPLR	ViT-H	MAE	59.8	42.2	63.8	74.9	51.9	32.2	55.7	71.0	7

Table 4: State-of-the-art comparison and scaling behavior for object detection and instance segmentation. We compare methods using feature pyramids vs. plain features on COCO val (n/a: entry is not available). Backbones with MAE pre-trained on ImageNet-1K while others pre-trained on ImageNet-21K. With only single-scale features, SimPLR shows strong performance compared to multi-scale detectors including transformer-based detectors like Mask2Former, while being $\sim 2\times$ faster.

365 objects in the image. To overcome this limitation, we think that a design of attention 365
 366 computation which effectively combines both global and local information is necessary. 366

367 5 Conclusion

368 We presented SimPLR, a simple and plain object detector that eliminates the requirement 368
 369 for handcrafting multi-scale feature maps. Through our experiments, we demonstrated 369
 370 that a transformer-based detector, equipped with a scale-aware attention mechanism, 370
 371 can effectively learn scale-equivariant features through data. The efficient design of 371
 372 SimPLR allows it to take advantages of significant progress in scaling ViTs, reaching 372
 373 highly competitive performance on three tasks on COCO: object detection, instance 373
 374 segmentation, and panoptic segmentation. This finding suggests that many handcrafted 374
 375 designs for convolution neural network in computer vision could be removed when 375
 376 moving to transformer-based architecture. We hope this study could encourage future 376
 377 exploration in simplifying neural network architectures especially for dense vision tasks. 377

378 References

- 379 1. Bao, H., Dong, L., Piao, S., Wei, F.: Beit: Bert pre-training of image transformers. In: ICLR 379
 380 (2022) 2
- 381 2. Brown, T.B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., 381
 382 Shyam, P., Sastry, G., Amanda Askell, S.A., Herbert-Voss, A., Krueger, G., Henighan, T., 382

method	backbone	pre-train	Panoptic Segmentation			FPS
			PQ	PQ th	PQ st	
Feature pyramids						
MaskFormer [8]	Swin-B	sup-21K	51.8	56.9	44.1	-
Mask2Former [7]	Swin-B	sup-21K	56.4	62.4	47.3	-
Mask2Former [7]	Swin-L	sup-21K	57.8	64.2	48.1	4
Plain features						
SimPLR	ViT-B	BEiT ^{v2}	56.5	62.6	47.3	13
SimPLR	ViT-L	BEiT ^{v2}	58.5	65.1	48.6	8

Table 5: State-of-the-art comparison and scaling behavior for panoptic segmentation. We compare between methods using feature pyramids (*top row*) vs. single-scale (*bottom row*) on COCO val. SimPLR with single-scale input shows better results when scaling to larger backbones, while being $\sim 2\times$ faster compared to Mask2Former.

- 383 Child, R., Ramesh, A., Ziegler, D.M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., 383
 384 Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, 384
 385 I., Amodei, D.: Language models are few-shot learners. In: NeurIPS (2020) 1, 3 385
 386 3. Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., Zagoruyko, S.: End-to-end 386
 387 object detection with transformers. In: ECCV (2020) 1, 2, 3 387
 388 4. Chen, M., Radford, A., Child, R., Wu, J., Jun, H., Dhariwal, P., Luan, D., Sutskever, I.: 388
 389 Generative pretraining from pixels. In: ICML (2020) 3 389
 390 5. Chen, T., Kornblith, S., Norouzi, M., Hinton, G.: A simple framework for contrastive learning 390
 391 of visual representations. In: ICML (2020) 3 391
 392 6. Chen, W., Du, X., Yang, F., Beyer, L., Zhai, X., Lin, T.Y., Chen, H., Li, J., Song, X., Wang, 392
 393 Z., Zhou, D.: A simple single-scale vision transformer for object localization and instance 393
 394 segmentation. In: ECCV (2022) 4 394
 395 7. Cheng, B., Misra, I., Schwing, A.G., Kirillov, A., Girdhar, R.: Mask2former: Masked-attention 395
 396 mask transformer for universal image segmentation. In: CVPR (2022) 1, 2, 3, 6, 7, 11, 12, 13 396
 397 8. Cheng, B., Schwing, A.G., Kirillov, A.: Per-pixel classification is not all you need for semantic 397
 398 segmentation. In: NeurIPS (2021) 3, 13 398
 399 9. Dehghani, M., Djolonga, J., Mustafa, B., et al.: Scaling vision transformers to 22 billion 399
 400 parameters. In: ICML (2023) 2, 8 400
 401 10. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical 401
 402 image database. In: CVPR (2009) 3 402
 403 11. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional 403
 404 transformers for language understanding. In: ACL (2019) 1, 3 404
 405 12. Dong, B., Zeng, F., Wang, T., Zhang, X., Wei, Y.: SOLQ: Segmenting objects by learning 405
 406 queries. In: NeurIPS (2021) 3 406
 407 13. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., De- 407
 408 hghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N.: An image is 408
 409 worth 16x16 words: Transformers for image recognition at scale. In: ICLR (2021) 1, 2, 3, 8 409
 410 14. Fan, H., Xiong, B., Mangalam, K., Li, Y., Yan, Z., Malik, J., Feichtenhofer, C.: Multiscale 410
 411 vision transformers. In: ICCV (2021) 1, 11, 12 411
 412 15. Ghiasi, G., Cui, Y., Srinivas, A., Qian, R., Lin, T., Cubuk, E.D., Le, Q.V., Zoph, B.: Simple 412
 413 copy-paste is a strong data augmentation method for instance segmentation. In: CVPR (2021) 413
 414 8 414

- 415 16. Girshick, R., Donahue, J., Darrell, T., Malik, J.: Rich feature hierarchies for accurate object
416 detection and semantic segmentation. In: CVPR (2014) 3 415
- 417 17. He, K., Chen, X., Xie, S., Li, Y., Dollár, P., Girshick, R.: Masked autoencoders are scalable
418 vision learners. In: CVPR (2022) 2, 3, 4, 8, 10, 11 416
- 419 18. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask R-CNN. In: ICCV (2017) 2 417
- 420 19. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR
421 (2016) 3 418
- 422 20. Heo, B., Yun, S., Han, D., Chun, S., Choe, J., Oh, S.J.: Rethinking spatial dimensions of
423 vision transformers. In: ICCV (2021) 1 419
- 424 21. Huang, G., Liu, Z., van der Maaten, L., Weinberger, K.Q.: Densely connected convolutional
425 networks. In: CVPR (2017) 3 420
- 426 22. Jia, D., Yuan, Y., He, H., Wu, X., Yu, H., Lin, W., Sun, L., Zhang, C., Hu, H.: Detrs with
427 hybrid matching. In: CVPR (2023) 4 421
- 428 23. Kirillov, A., He, K., Girshick, R., Rother, C., Dollar, P.: Panoptic segmentation. In: CVPR
429 (2019) 6 422
- 430 24. LeCun, Y., Bengio, Y.: Convolutional Networks for Images, Speech and Time Series, pp.
431 255–258. MIT Press (1995) 3 430
- 432 25. Li, F., Zhang, H., xu, H., Liu, S., Zhang, L., Ni, L.M., Shum, H.Y.: Mask dino: Towards a
433 unified transformer-based framework for object detection and segmentation. In: CVPR (2023)
434 3 431
- 435 26. Li, Y., Mao, H., Girshick, R., He, K.: Exploring plain vision transformer backbones for object
436 detection. In: ECCV (2022) 2, 3, 4, 5, 7, 8, 9, 11, 12 435
- 437 27. Lin, T.Y., Dollár, P., Girshick, R., He, K., Hariharan, B., Belongie, S.: Feature pyramid
438 networks for object detection. In: CVPR (2017) 2, 4 436
- 439 28. Lin, T., Maire, M., Belongie, S.J., Bourdev, L.D., Girshick, R.B., Hays, J., Perona, P., Ramanan,
440 D., Dollár, P., Zitnick, C.L.: Microsoft COCO: common objects in context. In: ECCV (2014)
441 8, 11 437
- 442 29. Lin, Y., Yuan, Y., Zhang, Z., Li, C., Zheng, N., Hu, H.: DETR does not need multi-scale or
443 locality design. In: ICCV (2023) 4, 8, 9 438
- 444 30. Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C.Y., Berg, A.C.: SSD: Single
445 shot multibox detector. In: ECCV (2016) 4 439
- 446 31. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer:
447 Hierarchical vision transformer using shifted windows. In: ICCV (2021) 1, 2, 3, 11, 12 440
- 448 32. Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. In: ICLR (2019) 8 441
- 449 33. Nguyen, D., Ju, J., Booij, O., Oswald, M.R., Snoek, C.G.M.: Boxer: Box-attention for 2d and
450 3d transformers. In: CVPR (2022) 1, 2, 3, 4, 5, 7, 8, 9, 10, 11 442
- 451 34. Peng, Z., Dong, L., Bao, H., Ye, Q., Wei, F.: BEiT v2: Masked image modeling with vector-
452 quantized visual tokenizers. arXiv preprint arXiv:2208.06366 (2022) 3, 8, 10, 11 443
- 453 35. Ren, S., He, K., Girshick, R., Sun, J.: Faster R-CNN: Towards real-time object detection with
454 region proposal networks. In: NeurIPS (2015) 2 444
- 455 36. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recog-
456 nition. In: ICLR (2015) 3 445
- 457 37. Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., Jegou, H.: Training data-
458 efficient image transformers and distillation through attention. In: International Conference
459 on Machine Learning (2021) 10 446
- 460 38. Touvron, H., Cord, M., Jegou, H.: Deit iii: Revenge of the vit. arXiv preprint arXiv:2204.07118
461 (2022) 10 447
- 462 39. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L.,
463 Polosukhin, I.: Attention is all you need. In: NeurIPS (2017) 1, 4 448
- 464 40. Wang, H., Zhu, Y., Adam, H., Yuille, A., Chen, L.C.: Max-deeplab: End-to-end panoptic
465 segmentation with mask transformers. In: CVPR (2021) 1, 3 449

- 466 41. Wang, W., Xie, E., Li, X., Fan, D.P., Song, K., Liang, D., Lu, T., Luo, P., Shao, L.: Pyramid
467 vision transformer: A versatile backbone for dense prediction without convolutions. In: ICCV
468 (2021) 1 466
- 469 42. Xie, S., Girshick, R.B., Dollár, P., Tu, Z., He, K.: Aggregated residual transformations for
470 deep neural networks. In: CVPR (2017) 3 467
- 471 43. Yu, Q., Wang, H., Qiao, S., Collins, M., Zhu, Y., Adam, H., Yuille, A., Chen, L.C.: k-means
472 mask transformer. In: ECCV (2022) 3 468
- 473 44. Zhai, X., Kolesnikov, A., Houlsby, N., Beyer, L.: Scaling vision transformers. In: CVPR
474 (2022) 8 469
- 475 45. Zhang, H., Li, F., Liu, S., Zhang, L., Su, H., Zhu, J., Ni, L., Shum, H.Y.: Dino: Detr with
476 improved denoising anchor boxes for end-to-end object detection. In: ICLR (2023) 3, 8 470
- 477 46. Zhang, W., Pang, J., Chen, K., Loy, C.C.: K-net: Towards unified image segmentation. In:
478 NeurIPS (2021) 1, 3 471
- 479 47. Zhu, X., Su, W., Lu, L., Li, B., Wang, X., Dai, J.: Deformable DETR: Deformable transformers
480 for end-to-end object detection. In: ICLR (2021) 1, 2, 4, 5, 8 472