

# Dual-Decoder Network for End-to-End Detection and Segmentation

Anonymous CVPR submission

Paper ID \*\*\*\*\*

## Abstract

## 1. Introduction

After its astonishing success in natural language processing, the transformer [?] has quickly become the neural network architecture of choice in computer vision, as evidenced by recent success in image classification [?, ?], object detection [?, ?, ?] and segmentation [?, ?]. Different from natural language processing, where the same pre-trained network can be deployed for a wide range of downstream tasks with only minor modifications [?, ?], computer vision tasks such as object detection and segmentation require a different set of domain-specific knowledge to be incorporated into the network. Consequently, it is commonly accepted that a modern object detector contains two main components: a pre-trained backbone as the *general* feature extractor, and a *task-specific* head that conducts object detection and instance segmentation tasks using domain knowledge. For transformer-based vision architectures, the question remains whether to add more inductive biases or to learn them from data.

The spatial nature of image data lies at the core of computer vision. Besides learning long range feature dependencies, the ability of capturing structure and relationships between neighboring pixels or features is critical for representing and understanding the image content. Building upon the successes of convolutional neural networks in computer vision, a line of research biases the transformer architecture to be *multi-scale* and *hierarchical* when dealing with the image input, (*i.e.*, Swin Transformer [?] and others [?, ?, ?]). They follow the design principles of convolutional neural networks by progressively increasing the receptive fields in each stage while spatially reducing the feature maps. This conventional design has proven effective in supporting multi-scale features for dense vision tasks and allows pre-trained transformers to be seamlessly integrated into convolution-based architectures with feature pyramid network [?], yielding impressive results in object detection and segmentation. Since the necessity of these complex adaptations in the transformer remains uncertain, an alternative research direction

pursues the idea of a simple transformer with “less inductive biases” and to emphasize learning vision-specific knowledge directly from image data. Starting with the end-to-end detection framework proposed in [?], a transformer-based prediction head removes many hand-designed components, like non-maximum suppression or intersection-over-union (IoU) computation, that encodes the prior knowledge. Recently, the Vision Transformer (ViT) [?], stands out as a *plain* and *non-hierarchical* architecture with a constant feature resolution, and acts as the feature extractor in plain-backbone detection. This is motivated by the success of ViTs in visual recognition [?, ?, ?, ?]. However, the plain design of ViTs casts doubts about its ability to capture information of objects across multiple scales. While recent studies [?, ?] suggest that ViTs with non-local self-attention computation could potentially learn translation-equivariant and scale-equivariant property during training, leading object detectors still require multi-scale feature maps. This observation holds true for both convolutional detectors [?, ?] and transformer-based detectors [?, ?, ?]. We postulate that an attention mechanism for scale-equivariant learning would enable the transformer-based detector to predict multi-scale objects using only the final feature map from a ViT backbone, thus eliminating the need for explicit multi-scale feature input.

With the goal of reducing inductive biases and emphasizing simplicity, we introduce BoxeR, an end-to-end object detector where both backbone and prediction head operate on plain features (See ??). Similar to [?], we focus on the *fine-tuning stage* and keep our approach independent of the pre-training phase. Our network architecture follows the recently introduced end-to-end object detection framework with transformer prediction head. In particular, BoxeR extracts the single-scale feature map from a pre-trained ViT which is then fed into the transformer encoder-decoder via a simple projection to make the final prediction. In order to capture objects of various sizes, we propose to incorporate scale information into the attention computation of the encoder, encouraging the network to learn scale-equivariant features. This eliminates the need for multi-scale feature maps and further simplifies the detection framework. For panoptic segmentation, our plain detector constructs a high-

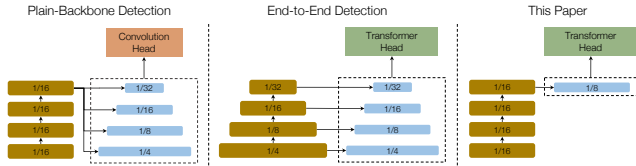


Figure 1. **Object detection architectures.** **Left:** The plain-backbone detector from [?] whose input (denoted in dashed region) are multi-scale features. **Middle:** state-of-the-art end-to-end detectors [?, ?] utilize a hierarchical backbone (*i.e.*, SwinTransformer [?]) to create multi-scale inputs. **Right:** Our simple single-scale detector following the end-to-end framework. Where existing detectors require multi-scale feature maps to be effective, we propose a plain detector whose backbone and detection head require only a single-scale feature.

resolution feature map to capture fine-grained details and incorporates masked instance-attention in the decoder to improve object boundary delineation. We show the effectiveness of our architecture by experimental results on COCO object detection, instance segmentation, and panoptic segmentation. Notably, with a small computational budget, our plain detector shows competitive performance with leading detectors that employ hierarchical backbone and/or multi-scale feature map as input. We hope that our study of a simple and plain detector will open up more possibilities and flexibility for future research on dense vision and multi-modal vision models, where the alignment between visual regions and multi-modal tokens can be easily learnt via a plain feature map.

## 2. Related Work

**Backbones for object detection.** Inspired by R-CNN [?], modern object detection methods utilize a task-specific head on top of a pre-trained backbone for their prediction. Initially, object detectors were dominated by convolutional neural network (CNN) backbones [?] pre-trained on ImageNet [?], with consistent accuracy and efficiency improvements over subsequent generations of architectures, *e.g.* [?, ?, ?, ?]. By removing the need of labels, methods with self-supervised learning have emerged as an even more powerful solution for pre-training general vision representations [?, ?]. With the success of the transformer in learning from large-scale text data [?, ?], many studies have explored the transformer for computer vision [?, ?]. Recently, the Vision Transformer (ViT) [?] demonstrated the capability of learning meaningful representation for visual recognition. We show through experiments that the pre-trained ViT when combined with a transformer-based prediction head is able to learn scale-equivariant features throughout the architecture.

**Multi-scale object detection.** A key challenge in object detection is to detect objects within an image across multiple scales. Early works tackled this problem by applying a CNN

detector with a sliding window strategy on an image pyramid [?] or with a generated region proposal [?] to extract each scale-normalized region from the input image [?]. To reduce computational costs, Faster R-CNN [?] generates object proposals on the feature map using a region proposal network. As deeper features of CNNs tends to capture more high-level information at the expense of fine-grained details for small objects, SSD [?] predicts objects from multiple layers of the feature hierarchy. The Feature Pyramid Network [?] further improves the creation of multi-scale feature maps with top-down fusion and lateral connections. With the recent evidence that non-hierarchical transformer architectures (*i.e.*, ViTs) are able to learn convolution-like behaviour through image data (*i.e.*, translation-equivariance), the necessity of multi-scale feature maps becomes questionable. Recent studies on plain-backbone detection [?, ?] show that the last feature of ViTs captures the necessary information to detect objects at multiple scales. However, these detectors still rely on multi-scale features to achieve better performance [?]. We remove this requirement in the prediction head, allowing us to rely on a simple and plain transformer for both the backbone and the task-specific head.

**End-to-end detection and segmentation.** A universal architecture solves multiple vision tasks without many architectural changes. This is made possible by the introduction of transformers for object detection [?]. Follow-up works [?, ?] extended the end-to-end detection framework to object detection and instance segmentation. This inspired MaskFormer [?] and K-Net [?] to unify segmentation tasks with a class-agnostic mask prediction. Pointing out that MaskFormer and K-Net still lag behind specialized architectures, Cheng *et al.* introduce Mask2Former [?], which outperforms specialized architectures on instance, semantic, and panoptic segmentation. Most recently, Yu *et al.* [?] replace self-attention with *k*-means clustering, further boosting the effectiveness of the network. Similar to convolution-based detection heads, these approaches utilize multi-scale feature maps from a hierarchical backbone. In this work, we remove the multi-scale feature map constraint and enable a universal architecture using a *single-scale* feature map from a plain vision transformer.

## 3. Method

Our goal is to further simplify the object detection pipeline from [?, ?], and to arrive at a universal object detection and segmentation model using plain features. To do so, we take advantage of recent progress in end-to-end object detection. Specifically, we adopt the box-attention mechanism [?] due to its effectiveness in learning discriminative object representations while being light-weight in computation. We investigate the design of a plain detector<sup>1</sup>

<sup>1</sup>In this paper, “backbone” refers to the components that we inherit from the pre-training stage, “detection head” refers to the components that are

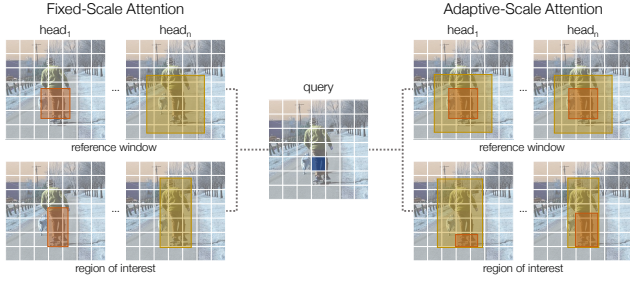


Figure 2. **Scale-aware attention mechanism.** **Left:** The fixed-scale attention with attention computation of a single scale per head. **Right:** The adaptive-scale attention with reference windows of multiple scales in each attention head. The 2-scale reference windows (denoted in yellow and orange) are used in this case. Unlike fixed-scale attention that transforms a single-scale reference window, adaptive-scale attention generates multi-scale regions of interest which are then selected using attention weights.

using *only a single-scale feature* from a plain transformer backbone during the *fine-tuning stage*.

In this section, we first review the box-attention mechanism. Then, the single-scale detector is presented for end-to-end object detection and segmentation. Finally, we propose several adaptations to make our model compatible with panoptic segmentation using a plain backbone.

### 3.1. Background

We first revisit box-attention proposed by Nguyen *et al.* [?]. Given the input feature map from the backbone, the encoder layers with box-attention will output contextual representations. The contextual representations are utilized to predict object proposals and to initialize object queries for the decoder. We denote the input feature map of an encoder layer as  $e \in \mathbb{R}^{H_e \times W_e \times d}$  and the query vector  $q \in \mathbb{R}^d$ , with  $H_e, W_e, d$  denoting height, width, and dimension of the input features respectively. Each query vector  $q \in \mathbb{R}^d$  in the input feature map is assigned a reference window  $r = [x, y, w, h]$ , where  $x, y$  indicate the query coordinate and  $w, h$  are the size of the reference window both being normalized by the image size. The box-attention refines the reference window into a region of interest,  $r'$ , as:

$$r' = F_{\text{scale}}(F_{\text{translate}}(r, q), q), \quad (1)$$

$$F_{\text{scale}}(r, q) = [x, y, w + \Delta_w, h + \Delta_h], \quad (2)$$

$$F_{\text{translate}}(r, q) = [x + \Delta_x, y + \Delta_y, w, h], \quad (3)$$

where  $F_{\text{scale}}$  and  $F_{\text{translate}}$  are the scaling and translation transformations,  $\Delta_x, \Delta_y, \Delta_w$  and  $\Delta_h$  are the offsets regarding to the reference window  $r$ . A linear projection ( $\mathbb{R}^d \rightarrow \mathbb{R}^4$ ) is

initialized from scratch, and “plain” refers to the single-scale property. The “plain detector” is the detector whose backbone and detection head both operate on single-scale features.

applied on  $q$  to predict offset parameters (*i.e.*,  $\Delta_x, \Delta_y, \Delta_w$  and  $\Delta_h$ ) w.r.t. the window size.

Similar to self-attention [?], box-attention aggregates  $n$  multi-head features from regions of interest:

$$\text{MultiHeadAttention} = \text{Concat}(\text{head}_1, \dots, \text{head}_n) W^O, \quad (4)$$

During the  $i$ -th attention head computation, a  $2 \times 2$  feature grid is sampled from the corresponding region of interest  $r'_i$ , resulting in a set of value features  $v_i \in \mathbb{R}^{2 \times 2 \times d_h}$ . The  $2 \times 2$  attention scores are efficiently generated by computing a dot-product between  $q \in \mathbb{R}^d$  and relative position embeddings ( $k_i \in \mathbb{R}^{2 \times 2 \times d}$ ) followed by a softmax function. The attended feature  $\text{head}_i \in \mathbb{R}^{d_h}$  is a weighted average of the  $2 \times 2$  value features in  $v_i$  with the corresponding attention weights:

$$\alpha = \text{softmax}(q^\top k_i), \quad (5)$$

$$\text{head}_i = \text{BoxAttention}(q, k_i, v_i) = \sum_{j=0}^{2 \times 2} \alpha_j v_{i_j}, \quad (6)$$

where  $q \in \mathbb{R}^d$ ,  $k_i \in \mathbb{R}^{2 \times 2 \times d}$ ,  $v_i \in \mathbb{R}^{2 \times 2 \times d_h}$  are query, key and value vectors of box-attention,  $\alpha_j$  is the  $j$ -th attention weight, and  $v_{i_j}$  is the  $j$ -th feature vector in the feature grid  $v_i$ . To better capture objects at different scales, the box-attention in [?] takes  $s$  multi-scale feature maps,  $\{e^j\}_{j=1}^s$  ( $s = 4$ ), as its inputs. In the  $i$ -th attention,  $s$  feature grids are sampled from each of multi-scale feature maps in order to produce  $\text{head}_i$ .

By computing the attended feature within regions of interest in each attention head, box-attention shows strong performance in end-to-end object detection with a small computational budget. The transformation functions (*i.e.*, translation and scaling) allow box-attention to capture long-range dependencies. The effectiveness and efficiency of box-attention leads us to a question whether we can learn multi-scale object information within a single-scale feature map.

### 3.2. BoxeR: A Simple and Plain Single-Scale Detector

**Single-scale detector.** Multi-scale feature maps have proven to be an effective solution to detect objects at multiple scales, from convolution-based object detectors [?, ?, ?] to transformer-based object detectors [?, ?, ?]. Multi-scale feature maps in a hierarchical backbone can be easily extracted from the pyramid structure [?, ?, ?]. 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 framework using ViTs remain unclear. Recent studies on plain-backbone detection [?, ?] conjecture the high-dimensional ViT with self-attention and positional



embeddings [?] is able to preserve important information for localizing objects<sup>2</sup>. From this conjecture, we hypothesize that a properly designed transformer-based prediction head to encourage scale-equivariant learning is capable of capturing the necessary information for detecting multi-scale objects from a single-scale feature map.

We denote the *last* feature map from the plain backbone as  $p \in \mathbb{R}^{H_p \times W_p \times d_p}$ , where  $H_p = \frac{H}{16}$ ,  $W_p = \frac{W}{16}$  are the size of feature map,  $H, W$  are the size of the input image, and  $d_p$  is the hidden dimension of the backbone. The input feature,  $e \in \mathbb{R}^{H_e \times W_e \times d}$ , for the prediction head is simply created by a projection  $f$ , where  $H_e, W_e$  are the size of the input feature map and  $d$  is the hidden dimension of the prediction head. We consider the single-scale input feature map with different scales in  $\{\frac{1}{32}, \frac{1}{16}, \frac{1}{8}, \frac{1}{4}\}$ . Different from [?], we empirically find that  $H_e = \frac{H}{8}$ ,  $W_e = \frac{W}{8}$  provides strong information for our detector to detect objects and helps to reduce the computational cost compared to higher resolution feature map. In the default setting ( $H_e = \frac{H}{8}$ ,  $W_e = \frac{W}{8}$ ), the projection  $f$  is a simple deconvolution layer.

In order to capture objects at different scales, global self-attention is a potential candidate because of its large receptive field and powerful representation. However, its computational complexity is quadratic w.r.t. the sequence length of the input, making the operation computationally expensive when dealing with high-resolution images. The self-attention mechanism also leads to slow convergence in end-to-end detection [?]. Therefore, we prefer to adopt the box-attention in the encoder in order to learn multi-scale information in the input features.

Given the input feature map,  $e \in \mathbb{R}^{H_e \times W_e \times d}$ , and the query vector  $q \in \mathbb{R}^d$ , we assign multi-scale reference windows to  $q$ , denoted by  $\hat{r}_t = [x, y, \hat{w}_t, \hat{h}_t]$ , where the multi-scale reference windows corresponding to  $q$  are centered at the query spatial position. By default, we use 4-scale reference windows ( $\hat{r}_1, \dots, \hat{r}_4$ ) and set the size of  $t$ -th reference window to  $\hat{w}_t = \hat{h}_t = 32 \cdot 2^{t-1}$  pixels. The use of multi-scale reference windows per query vector  $q$  is related to the multi-scale anchors in [?]. However, the multi-scale reference windows in our scenario are incorporated into multi-head attention mechanism in the encoder to learn scale-equivariant features, unlike [?], which utilizes them as the regression references in the final prediction layer. To incorporate scale information into multi-head attention, we study two scale-aware attention mechanisms (see ??):

[leftmargin=\*]*Fixed-Scale Attention.* We assign a reference window of one scale to each attention head. For example, with 8 attention heads and 4-scale reference windows, we assign each of attention heads with a reference window in a round robin manner,  $r = \{\hat{r}_1, \hat{r}_2, \hat{r}_3, \hat{r}_4, \hat{r}_1, \hat{r}_2, \hat{r}_3, \hat{r}_4\}$ . In the  $i$ -th at-

tention head, we predict offset parameters w.r.t. to  $r_i = [x, y, w_i, h_i]$ :

$$\Delta_{x_i} = (qW_{x_i}^\top + b_{x_i}) * w_i, \quad \Delta_{y_i} = (qW_{y_i}^\top + b_{y_i}) * h_i, \quad (7)$$

$$\Delta_{w_i} = \max(qW_{w_i}^\top + b_{w_i}, 0) * w_i, \quad \Delta_{h_i} = \max(qW_{h_i}^\top + b_{h_i}, 0) * h_i, \quad (8)$$

where  $W_{x_i}, W_{y_i}, W_{w_i}, W_{h_i}$  are learnable weights,  $b_{x_i}, b_{y_i}, b_{w_i}, b_{h_i}$  are learnable biases in the  $i$ -th attention head. This is equivalent to a single attention operation per scale, which allows each attended feature, head <sub>$i$</sub> , to capture a specific scale information w.r.t. the reference window  $r_i$ .

- ii) *Adaptive-Scale Attention.* Here, each attention head needs to process reference windows of all scales ( $\hat{r}_1, \dots, \hat{r}_4$ ). In the  $i$ -th attention head, we predict 4 sets of offset parameters regarding to 4-scale reference windows. The corresponding 4 grid features,  $v_i \in \mathbb{R}^{4 \times 2 \times 2 \times d_h}$ , are then sampled from the input feature map. The attention scores,  $\alpha \in \mathbb{R}^{4 \times 2 \times 2}$ , are generated for each feature in  $v_i$  in order to perform the attention operation. By computing the attention map on reference windows of multiple scales in each attention head, the attention mechanism is able to dynamically select necessary scale information for object prediction. We found that adaptive-scale attention leads to better performance.

**Overall transformer architecture.** Our network follows the design of a two-stage detector that performs end-to-end detection. Following [?], we simply use a plain ViT as the backbone with  $14 \times 14$  windowed attention and four equally-distributed global attention blocks. Our prediction head consists of a transformer-based encoder and decoder for end-to-end detection as proposed in [?, ?]. We use output features from the encoder with scale-aware attention to produce object proposals using the last features of ViT backbone. The top-scored feature vectors will be initialized as object queries and the corresponding bounding box proposals will serve as its reference window in the decoder. The decoder utilizes instance-attention [?], an extension of box-attention for dense grid sampling, to decode object queries into bounding boxes for object detection and instance masks for instance segmentation. More details are provided in the supplementary document.

### 3.3. SimPLR for Universal Detection and Segmentation

To enable the plain detector on panoptic segmentation, we make an adaptation in the mask prediction which allows Boxer to segment both “thing” and “stuff” categories effectively. To be specific, we predict segmentation masks of both

<sup>2</sup>ViT-B and larger ( $d \geq 768$ ) can maintain the information with a patch size of  $16 \times 16 \times 3$  in the input image.

types by computing the dot-product between object queries and a high-resolution feature map (*i.e.*,  $\frac{1}{4}$  feature scale) following [?]. Our goal is to have *minimal* adaptation which allows us to use the same training recipe for all tasks, including object detection, instance segmentation, and panoptic segmentation. As the output feature of our encoder is in  $\frac{1}{8}$  scale of the original image, we simply use a transposed convolution layer to generate the high-resolution feature for mask prediction. We found that it is beneficial to stack  $K$  encoder layers on top of the high-resolution feature to learn more contextual information. In addition, we incorporate the decoder with masked instance-attention to obtain discriminative representation for mask prediction. Different from masked cross-attention in [?], our masked instance-attention extends box-attention in [?] which yields a small complexity when dealing with feature map of  $\frac{1}{8}$  scale from the encoder. More details are provided in the supplementary document.

**Masked instance-attention.** The masked instance-attention follows the grid sampling strategy of the box-attention in [?], but differs in the computation of attention scores to better capture objects of different shapes. To be specific, the region of interest  $r'_i$  is divided into 4 bins of  $2 \times 2$  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  $2 \times 2$  attention scores for 4 bins. The  $\frac{m}{2} \times \frac{m}{2}$  feature vectors within the same bin share the same attention weight. This is equivalent to the *average* aggregation of feature values covered by each bin, which shows to reduce misalignments in RoIAlign [?]:

$$\text{head}_i = \sum_{k=0}^{2 \times 2} \sum_{j=0}^{\frac{m}{2} \times \frac{m}{2}} \frac{\alpha_k}{\frac{m}{2} \cdot \frac{m}{2}} v_{i_{k,j}}, \quad (9)$$

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

Inspired by [?], we utilize the mask prediction before the sigmoid of the previous decoder layer  $\mathcal{M}_q \in \mathbb{R}^{H_m \times W_m}$  corresponding to the object query  $q$ . Given the coordinates of grid features within the region of interest  $r'_i$ , we sample the corresponding mask scores using bilinear interpolation. The sampled mask scores are binarized with the 0.5 threshold and applied into the attention computation.

$$\text{head}_i = \sum_{k=0}^{2 \times 2} \sum_{j=0}^{\frac{m}{2} \times \frac{m}{2}} \left( \frac{\alpha_k}{\frac{m}{2} \cdot \frac{m}{2}} + m_{i_{k,j}} \right) v_{i_{k,j}}, \quad (10)$$

$$m_{i_{k,j}} = \begin{cases} 0 & \text{if } \text{sigmoid}(\mathcal{M}_q(v_{i_{k,j}})) \geq 0.5 \\ -\infty & \text{otherwise} \end{cases}, \quad (11)$$

where  $\mathcal{M}_q(v_{i_{k,j}})$  is the mask score sampled at the location of the feature  $v_{i_{k,j}}$ . The masked instance-attention along with high-resolution layers allow Boxer to capture more

[]table

Table 1. **Ablation on single-scale vs. multi-scale feature map** on COCO [?] object detection and instance segmentation using ViT-B as backbone. Compared to ViTDet, Boxer performs better for both tasks with less FLOPs. Interestingly, Boxer performs on par with Boxer under the same FLOPs despite using only single-scale input ( $\ddagger$ : Boxer with SimpleFPN [?]).

plain-backbone detector	AP <sup>b</sup> ↑	AP <sup>m</sup> ↑	FLOPs ↓
<b>Multi-scale head</b>			
ViTDet [?]	54.0	46.7	1.1T
Boxer <sup>‡</sup> [?]	55.4	47.7	0.5T
<b>Single-scale head</b>			
Boxer	55.4	47.6	0.5T

fine-grained details, which is beneficial for panoptic segmentation.

## 4. Experiments

attention computation	AP <sup>b</sup> ↑	AP <sup>m</sup> ↑
base	54.5	46.8
i) fixed-scale attention	55.0	47.2
ii) adaptive-scale attention	55.4	47.6

(a) **Strategy for learning scale equivariance.**

feature scale	AP <sup>b</sup> ↑	AP <sup>m</sup> ↑	window size	AP <sup>b</sup> ↑	AP <sup>m</sup> ↑
base	54.5	46.8	base	54.5	46.8
1/4	55.4	47.7	16	55.1	47.4
1/8	55.4	47.6	32	55.4	47.6
1/16	54.2	46.6	64	55.1	47.4

(b) **Scales of input feature map.**

(c) **Reference window size.**

Table 2. **Ablation on the design of attention mechanism** using a plain ViT backbone on COCO [?] object detection and instance segmentation. Boxer receives the single-scale input from ViT and makes predictions. Compared to the naïve baseline which employs Boxer and box-attention [?] with single-scale features, Boxer improves performance of the plain detector.

**Dataset, tasks & evaluation.** In this study, we evaluate our method on COCO [?], a common detection and segmentation dataset, catering for object detection, instance segmentation, and panoptic segmentation tasks. The COCO dataset contains 118,000 training images and 5,000 validation images of 80 “thing” and 53 “stuff” categories. We report the average precision metric AP<sup>b</sup> for object detection and AP<sup>m</sup> for instance segmentation, and the panoptic quality metric PQ for panoptic segmentation. In panoptic segmentation, the union of “thing” and “stuff” categories are considered whereas in object detection and instance segmentation, the “thing” categories are the main focus. We train our network on the train split and evaluate on the val split.

**Implementation details.** By default, we use a single-scale feature map with adaptive-scale attention as described in

attention	PQ↑	PQ <sup>Th</sup> ↑	PQ <sup>St</sup> ↑
instance-attention [?]	54.9	61.2	45.3
masked instance-attention	55.3	61.6	45.8

(a) Instance-attention vs. Masked instance-attention.

K	PQ↑	PQ <sup>Th</sup> ↑	PQ <sup>St</sup> ↑
0	54.6	60.9	45.1
1	55.3	61.6	45.8

(b) Number of high-resolution layers.

Table 3. **Ablation on mask prediction adaptation** in a plain detector on COCO [?] panoptic segmentation. The masked instance-attention along with high-resolution layers improve the performance of BoxeR on panoptic segmentation, showing its effectiveness in capturing fine-grained information.

?? We initialize the backbone from MAE [?] pre-trained on ImageNet-1K without any labels. Unless otherwise specified, the hyper-parameters are the same as in [?]. For all experiments, our optimizer is AdamW [?] with a learning rate of 0.0001 and a weight decay of 0.05. The learning rate is linearly warmed up for the first 250 iterations and decayed at 0.9 and 0.95 fractions of the total number of training steps by a factor 10. ViT-B [?] is set as the backbone. The layer-wise learning rate decay of 0.7 and drop path rate of 0.3 is applied to the ViT backbone. The input image size is  $1024 \times 1024$  with large-scale jitter [?] between a scale range of  $[0.1, 2.0]$ . Due to the limit of our computational resources, we report the ablation study using the standard  $5 \times$  schedule setting [?] with batch size of 16. Following [?], we also finetune our networks within 100 epochs and a batch size of 64 in the main experiments. As the main goal is to keep the architecture and training process simple, these settings are applied to all tasks (*i.e.*, object detection, instance segmentation, and panoptic segmentation). The code will be released.

**Ablation 1: A single-scale feature map is enough.** In ??, we first compare plain-backbone detectors with single-scale vs. multi-scale features. The baseline of our study is the ViTDet: which employs the plain ViT as the backbone and Cascade R-CNN [?] with *multi-scale* feature maps from a simple feature pyramid (SimpleFPN) [?] as its detection head. We also train BoxeR with the plain backbone ViT and SimpleFPN using a  $5 \times$  schedule for a better comparison.

From our observation, the single-scale feature map is sufficient and our adaptive-scale attention mechanism enables BoxeR to perform detection using plain features. As a result, BoxeR matches the performance of BoxeR under the same amount of FLOPs. This observation is different from BoxeR and recent end-to-end detection approaches [?, ?] where a multi-scale feature map is needed to achieve competitive performance. We note that BoxeR follows a different perspective compared to convolution-based detection in [?] when using a single-scale feature map. The convolution-based detection operates region-of-interest (RoI) pooling on object

proposals to encode multi-scale information. Instead, our transformer-based encoder performs attention computation per feature vector to learn scale equivariant features. The significance of multiple feature maps diminishes when the encoder is equipped with a powerful attention mechanism.

**Ablation 2: Scale-aware attention mechanism is an important factor.** ?? ablates the importance of the proposed attention mechanism. Here, we study the baseline which is the original box-attention from [?] with a single-scale feature map of  $\frac{1}{8}$  scale and window size of 32 pixels (denoted in “base”) and the design choices of our scale-aware attention mechanism for detection.

It can be seen from ?? that both attention strategies of encoding multi-scale information in attention computation demonstrate strong detection performance. By allowing the feature vector to choose suitable scale information in its attended feature, our adaptive-scale attention obtains an improvement of 0.9 AP point over the baseline.

In ??, we verify the performance of BoxeR across multiple feature scales. The  $\frac{1}{16}$  feature is created by applying a  $1 \times 1$  convolution layer to the last ViT feature. For  $\frac{1}{8}$  and  $\frac{1}{4}$  feature, we use a transposed convolution layer. By default, we use BoxeR with adaptive-scale attention. Interestingly, the proposed method works well on  $\frac{1}{16}$  features and performs worse than the baseline with  $\frac{1}{8}$  features by a small gap of 0.3 AP point. The best trade-off between accuracy and efficiency comes from the features of  $\frac{1}{8}$  scale. This results may align with DETR [?] where DETR-DC5 with dilated convolution in the C5 block shows better performance. In our case, we show that a simple transposed convolution layer is sufficient with plain ViT backbone.

?? compares the performance of BoxeR across several base sizes of the reference window. Our attention mechanism has only marginal differences in term of the base size. The ablation reveals that the number of scales in attention computation rather than the reference window size plays an important role to make our network scale-aware. Indeed, in ??, the use of multiple window scales shows improvement over the baseline. Increasing from 2 to 4 scales improves detection performance. This suggests that it is not trivial to be scale equivariant with a plain feature map and our adaptive-scale attention mechanism helps to strengthen features for detecting multi-scale objects.

**Ablation 3: A single-scale architecture for panoptic seg-**



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

method	backbone	AP <sup>b</sup>	AP <sup>b</sup> <sub>S</sub>	AP <sup>b</sup> <sub>M</sub>	AP <sup>b</sup> <sub>L</sub>	AP <sup>m</sup>	AP <sup>m</sup> <sub>S</sub>	AP <sup>m</sup> <sub>M</sub>	AP <sup>m</sup> <sub>L</sub>	PQ	PQ <sup>Th</sup>	PQ <sup>St</sup>	FLOPs
<b>Hierarchical, Multi-scale</b>													
Swin <sup>†</sup> [?]	Swin-B	54.0	-	-	-	46.5	-	-	-	n/a	n/a	0.9T	0.9T
MViT <sup>†</sup> [?]	MViT-B	55.6	-	-	-	48.1	-	-	-	n/a	n/a	0.8T	0.8T
MaskFormer(5×) [?]	Swin-B			n/a				n/a		51.1	56.3	43.2	0.4T
Mask2Former(5×) [?]	Swin-B			n/a		46.7	26.1	50.5	68.8	55.1	61.9	46.1	0.5T
<b>Plain, Multi-scale</b>													
ViTDet [?]	ViT-B	54.0	-	-	-	46.7	-	-	-	n/a	n/a	1.1T	1.1T
BoxeR <sup>‡</sup> (5×) [?]	ViT-B	55.4	-	-	-	47.7	-	-	-	n/a	n/a	0.5T	0.5T
<b>Simple, Plain</b>													
UViT [?]	UViT-B	53.9	-	-	-	46.1	-	-	-	n/a	n/a	1.1T	1.1T
BoxeR (5×)	ViT-B	55.4	36.1	59.1	70.9	47.6	26.8	51.4	67.1	55.3	61.6	45.8	0.3T
BoxeR	ViT-B	<b>55.6</b>	<b>37.1</b>	<b>59.2</b>	<b>71.6</b>	48.0	<b>27.5</b>	<b>51.5</b>	67.8	55.6	62.1	45.8	0.3T

Table 4. **Universal detection and segmentation comparison** for object detection, instance segmentation, and panoptic segmentation on the COCO [?] val set. By default, all backbones are pre-trained using ImageNet1K. (†: backbones pre-trained on ImageNet-22K; ‡: BoxeR with SimpleFPN from [?]; 5×: methods fine-tuned using 5× schedule; n/a: method is not applicable to the task). Methods in gray color are with a convolution-based detection head. BoxeR as a plain detector demonstrates competitive performance compared to many specialized or end-to-end models with a small amount of computation and is the only one suited for all three tasks.

**mentation.** We evaluate the effectiveness of our adaptations on panoptic segmentation in ???. The masked instance-attention outperforms the instance attention, an extension of box-attention for dense grid in [?], when keeping the same efficiency. Moreover, adding encoder layer on top of high-resolution feature map delivers a large gain in panoptic segmentation. This evidence confirms the capability of the encoder layer in learning both fine-grained and semantically meaningful features.

**Universal detection and segmentation.** ?? lists the performance of BoxeR vs. previous methods in object detection, instance segmentation and panoptic segmentation. These approaches are divided into three parts: detectors with hierarchical backbone and multi-scale head, plain-backbone detectors using multi-scale features, and plain detector with both plain backbone and head. For plain-backbone methods, we report the performance of ViTDet [?] and BoxeR w/ SimpleFPN [?] fine-tuned on input image size of 1024 × 1024. Our method accounts for less FLOPs than convolution-based detection head. Across different detectors, BoxeR achieves strong performance on COCO object detection and instance segmentation. Interestingly, BoxeR outperforms Mask2Former [?] in detecting small and medium objects but lags behind in large category. This makes BoxeR to be a competitive plain detector. In panoptic segmentation, we observe the same behaviour in which BoxeR improves over “thing” categories but lags behind in “stuff” categories. Overall, BoxeR, as a plain detector, can achieve highly accurate results on all three tasks and can compete with hierarchical

Visualizations of BoxeR predictions are shown in ??.

**Limitations.** Our final goal is to simplify the detection pipeline and to achieve competitive results at the same time. In [?] and [?], we find that the attention mechanism which can effectively learn scale-aware information in its computation plays a key role for a plain detector. However, our adaptive-scale attention mechanism still encodes the knowledge of different scales. In the future, we hope that with the help of large-scale training data, a simpler design of the attention mechanism could also learn the scale equivariant property. Furthermore, BoxeR faces difficulties in detecting and segmenting large objects in the image. To overcome this limitation, we think that a design of attention computation which effectively combines both global and local information is necessary.

## 5. Conclusion and Limitations