

# You Only Look at One Sequence: Rethinking Transformer in Vision through Object Detection

**Yuxin Fang**<sup>1\*</sup> **Bencheng Liao**<sup>1\*</sup> **Xinggang Wang**<sup>1</sup>  $\boxtimes$  **Jiemin Fang**<sup>2,1</sup>  
**Jiayang Qi**<sup>1</sup> **Rui Wu**<sup>3</sup> **Jianwei Niu**<sup>3</sup> **Wenyu Liu**<sup>1</sup>

<sup>1</sup> School of EIC, Huazhong University of Science & Technology

<sup>2</sup> Institute of Artificial Intelligence, Huazhong University of Science & Technology

<sup>3</sup> Horizon Robotics

{yxf, bcliao, xgwang}@hust.edu.cn

## Abstract

Can Transformer perform 2D object-level recognition from a pure sequence-to-sequence perspective with minimal knowledge about the 2D spatial structure? To answer this question, we present You Only Look at One Sequence (YOLOS), a series of object detection models based on the naïve Vision Transformer with the fewest possible modifications as well as inductive biases. We find that YOLOS pre-trained on the mid-sized ImageNet-1k dataset only can already achieve competitive object detection performance on COCO, *e.g.*, YOLOS-Base directly adopted from BERT-Base can achieve 42.0 box AP. We also discuss the impacts as well as limitations of current pre-train schemes and model scaling strategies for Transformer in vision through object detection. Code and model weights are available at <https://github.com/hustvl/YOLOS>.

## 1 Introduction

Transformer [56] is born to transfer. In natural language processing (NLP), the dominant approach is to first pre-train Transformer on large, generic corpora for general language representation learning, and then fine-tune the model on specific target tasks [17]. Recently, Vision Transformer (ViT)<sup>1</sup> [20] demonstrates that typical Transformer encoder architecture directly inherited from NLP can perform surprisingly well on image recognition at scale using modern vision transfer learning recipe [32]. Taking sequences of image patch embeddings as inputs, ViT can successfully transfer pre-trained general visual representations from sufficient scale to more specific image classification tasks with fewer data points from a pure sequence-to-sequence perspective.

Since a pre-trained Transformer can be successfully fine-tuned on sentence-level tasks [7, 18] in NLP, as well as token-level tasks [44, 49], where models are required to produce fine-grained output at the token level [17]. A natural question is: Can ViT transfer to more complex target tasks in computer vision such as object detection other than image-level recognition?

ViT-FRCNN [5] is the first to use a pre-trained ViT as the backbone for an R-CNN [22] object detector. However, this design cannot get rid of the reliance on convolutional neural networks (CNNs) and

\*Yuxin Fang and Bencheng Liao contributed equally.  $\boxtimes$  Xinggang Wang is the corresponding author. This work was done when Yuxin Fang was interning at Horizon Robotics mentored by Rui Wu.

<sup>1</sup>Recently, there are various sophisticated or hybrid architectures termed as “Vision Transformer”. For disambiguation, in this paper, “Vision Transformer” and “ViT” refer to the naïve or vanilla Vision Transformer architecture proposed by Dosovitskiy et al. [20] unless specified.

strong 2D inductive biases, as ViT-FRCNN re-interprets the output sequences of ViT to 2D spatial feature maps and depends on region-wise pooling operations (*i.e.*, RoIPool [21, 24] or RoIAvg [26]) as well as region-based CNN architectures [47] to decode ViT features for object-level perceptions. Inspired by modern CNN design, some recent works [37, 57, 60, 63] introduce the pyramidal feature hierarchy and locality to Vision Transformer design, which largely boost the performance in dense prediction tasks including object detection. However, these architectures are performance-oriented and cannot reflect the properties of the naïve or vanilla Vision Transformer [20] that directly inherited from Vaswani et al. [56]. Another series of work, the DETR families [9, 70], use a random initialized Transformer to encode & decode CNN features, which does not reveal the transferability of a pre-trained Transformer in object detection.

Intuitively, ViT is designed to model long-range dependencies and global contextual information instead of local and region-level relations. Moreover, ViT lacks hierarchical architecture as modern CNNs [25, 34, 50] to handle the large variations in the scale of visual entities [1, 36]. Based on the available evidence, it is still unclear whether a pure ViT can transfer pre-trained general visual representations from image-level recognition to the much more complicated 2D object detection task.

To answer this question, we present You Only Look at One Sequence (YOLOS)<sup>2</sup>, a series of object detection models based on the canonical ViT architecture with the fewest possible modifications as well as inductive biases injected. The change from a ViT to a YOLOS detector is simple: (1) YOLOS drops the [CLS] token in ViT and appends one hundred learnable [DET] tokens to the input sequence for object detection. (2) YOLOS replaces the image classification loss in ViT with the bipartite matching loss to perform object detection in a set prediction manner following Carion et al. [9], which can avoid re-interpreting the output sequences of ViT to 2D feature maps as well as prevent manually injecting heuristics and prior knowledge of object 2D spatial structure during label assignment [69].

Directly inherited from ViT [20], YOLOS is not designed to be yet another high-performance object detector, but to unveil the versatility and transferability of Transformer from image recognition to object detection. Concretely, our main contributions are summarized as follows:

- We use the mid-sized ImageNet-1k [48] as the sole pre-training dataset, and show that a naïve ViT [20] can be successfully transferred to perform the challenging object detection task and produce competitive COCO [35] results with the fewest possible modifications, *i.e.*, by only looking at one sequence (YOLOS).
- For the first time, we demonstrate that 2D object detection can be accomplished in a pure sequence-to-sequence manner by taking a sequence of fixed-sized non-overlapping image patches as input. Among existing object detectors, YOLOS utilizes minimal 2D inductive biases.
- For ViT, we find the object detection results are quite sensitive to the pre-train scheme and the detection performance is far from saturating. Therefore the proposed YOLOS can be used as a challenging benchmark task to evaluate different pre-training strategies for ViT.
- We also discuss the impacts of prevalent pre-train schemes and model scaling strategies for Transformer in vision through transferring to object detection.

## 2 You Only Look at One Sequence

In model design, YOLOS closely follows the original ViT architecture [20], and is optimized for object detection in the same vein as Carion et al. [9]. YOLOS can be easily adapted to various Transformers available in NLP as well as in computer vision. This intentionally simple setup is not designed for better detection performance, but to exactly reveal characteristics of the Transformer family in object detection as unbiased as possible.

### 2.1 Architecture

An overview of the model is depicted in Fig. 1. The change from a ViT to a YOLOS detector is simple: (1) YOLOS drops the [CLS] token for image classification and appends one hundred randomly initialized detection tokens ([DET] tokens) to the input patch embedding sequence for object detection.

---

<sup>2</sup>Salute to You Only Look Once (YOLO) [46].

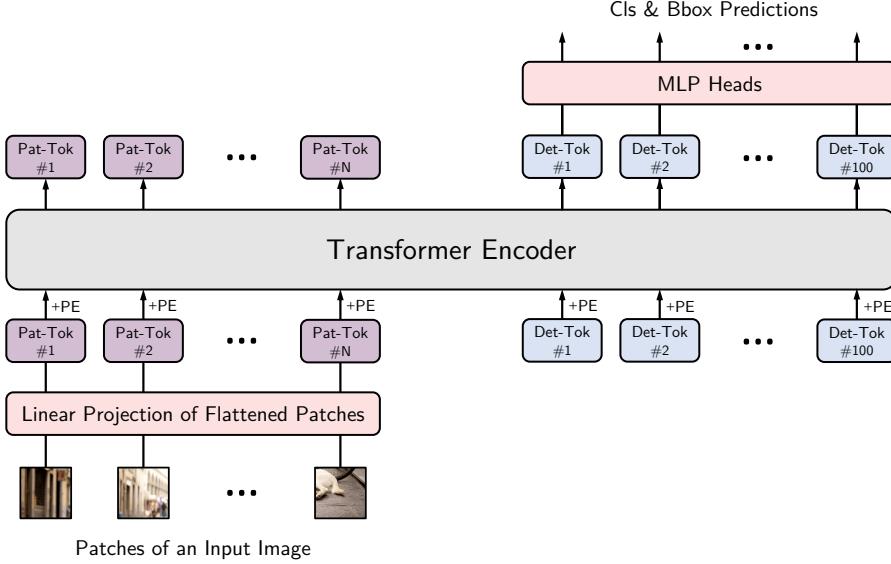


Figure 1: YOLOS overview. “Pat-Tok” refers to patch token, which is the embedding of a flattened image patch. “Det-Tok” refers to [DET] token, which is a learnable embedding for object detection predictions. “PE” refers to positional embedding. During training, YOLOS produces an optimal bipartite matching between predictions from one hundred [DET] tokens and ground truth objects. During inference, YOLOS directly outputs the final set of predictions in parallel. The figure style is inspired by Dosovitskiy et al. [20].

(2) During training, YOLOS replaces the image classification loss in ViT with the bipartite matching loss to perform object detection in a set prediction manner following Carion et al. [9]. Here we highlight the design methodology of YOLOS.

**Detection Token.** We purposefully choose randomly initialized [DET] tokens as proxies for object representations to avoid inductive biases of 2D structure and prior knowledge about the task injected during label assignment. When fine-tuning on COCO, for each forward pass, an optimal bipartite matching between predictions generated by [DET] tokens and ground truth objects is established. This procedure plays the same role as label assignment [9, 69], but is unaware of the input 2D structure, *i.e.*, YOLOS does not need to re-interpret the output sequence of ViT to an 2D feature maps for label assignment. Theoretically, it is feasible for YOLOS to perform any dimensional object detection without knowing the exact spatial structure and geometry, as long as the input is always flattened to a sequence in the same way for each pass.

**Fine-tuning at Higher Resolution.** When fine-tuning on COCO, all the parameters are initialized from ImageNet-1k pre-trained weights except for the MLP heads for classification & bounding box regression as well as one hundred [DET] tokens, which are randomly initialized. Both the classification and bounding box regression heads are implemented by MLP with two hidden layers using separate parameters. During fine-tuning, the image has a much higher resolution than pre-training, we keep the patch size the same ( $16 \times 16$ ), which results in a larger effective sequence length. While ViT can handle arbitrary sequence lengths, the positional embeddings need to adapt to the longer input sequences. We perform 2D interpolation of the pre-trained position embeddings in the same way as Dosovitskiy et al. [20].

**Inductive Bias.** We carefully design YOLOS for minimal additional inductive biases injection. The inductive biases inherent from ViT come from the patch extraction at the network stem part as well as the resolution adjustment for position embeddings. Apart from that, YOLOS adds no

non-degenerated (*e.g.*,  $3 \times 3$  or other non -  $1 \times 1$ ) convolutions upon ViT<sup>3</sup>. From the representation learning perspective, we choose to use [DET] tokens as proxies of objects for final predictions to avoid additional 2D inductive biases as well as heuristics. The performance-oriented design inspired by modern CNN architectures such as pyramidal feature hierarchy, 2D local spatial attention as well as the region-wise pooling operation is not applied. All these efforts are meant to exactly unveil the versatility and transferability of Transformer from image recognition to object detection in a pure sequence-to-sequence manner, with minimal knowledge about the input spatial structure and geometry.

**Comparisons with DETR.** The design of YOLOS is inspired by DETR [9]: YOLOS use [DET] tokens as proxies for object representations to avoid inductive biases about 2D structures and prior knowledge about the task injected during label assignment, and YOLOS is optimized in a similar way as DETR. Meanwhile, there are some key differences between the two<sup>4</sup>: (1) DETR uses a randomly initialized Transformer with an encoder-decoder architecture, while YOLOS studies the transferability of the pre-trained encoder-only ViT [20, 55]. (2) DETR uses decoder-encoder attention (cross attention) between image features and object queries with auxiliary decoding losses deeply supervised at each decoder layer, while YOLOS always looks at only one sequence for each layer, without distinguishing patch tokens and [DET] tokens in terms of operations. The quantitative comparisons are given in Sec. 3.4.

## 3 Experiments

### 3.1 Setup

**Pre-training.** We pre-train all YOLOS / ViT models on ImageNet-1k [48] dataset using the data-efficient training strategy suggested by Touvron et al. [55]. The parameters are initialized with a truncated normal distribution and optimized using AdamW [38]. The learning rate and batch size are  $1 \times 10^{-3}$  and 1024, respectively. The learning rate decay is cosine and the weight decay is 0.05. Rand-Augment [13] and random erasing [67] implemented by timm library [62] are used for data augmentation. Stochastic depth [31], Mixup [66] and Cutmix [64] are used for regularization.

**Fine-tuning.** We fine-tune all YOLOS models on COCO object detection benchmark [35] using in a similar way as Carion et al. [9]. All the parameters are initialized from ImageNet-1k pre-trained weights except for the MLP heads for classification & bounding box regression as well as one hundred [DET] tokens, which are randomly initialized. We train YOLOS on a single node with  $8 \times 12G$  GPUs. The learning rate and batch size are  $2.5 \times 10^{-5}$  and 8 respectively. The learning rate decay is cosine and the weight decay is  $1 \times 10^{-4}$ .

As for data augmentation, we use multi-scale augmentation, resizing the input images such that the shortest side is at least 256 and at most 608 pixels while the longest at most 864 for tiny models. For small and base models, we resize the input images such that the shortest side is at least 480 and at most 800 pixels while the longest at most 1333. We also apply random crop augmentations during training following Carion et al. [9]. The number of [DET] tokens are 100 and we keep the loss function as well as loss weights the same as DETR, while we don't apply dropout [51] or stochastic depth during fine-tuning since we find these regularization methods hurt performance.

**Model Variants.** With available computational resources, we study several YOLOS variants. Detailed configurations are summarized in Tab. 1. The input patch size for all models is  $16 \times 16$ . YOLOS-Ti (Tiny), -S (Small), and -B (Base) directly correspond to DeiT-Ti, -S, and -B [55]. From the model scaling perspective [19, 53, 58], the small and base models of YOLOS / DeiT can be seen as performing width scaling ( $w$ ) [29, 65] on the corresponding tiny model.

---

<sup>3</sup>We argue that it is imprecise to say Transformer do not have convolutions. All linear projection layers in Transformer are equivalent to point-wise or  $1 \times 1$  convolutions with sparse connectivity, parameter sharing, and equivalent representations properties, which can largely improve the computational efficiency compared with the “all-to-all” interactions in fully-connected design that has even weaker inductive biases [4, 23].

<sup>4</sup>The original DETR [9] relies on CNN features, but it is feasible to use ViT as a feature extractor for DETR. Therefore we do not treat this as a key difference.

Model	DeiT [55] Model	Layers (Depth)	Embed. Dim. (Width)	Pre-train Resolution	Heads	Params.	FLOPs	$\frac{f(\text{Lin.})}{f(\text{Att.})}$
YOLOS-Ti	DeiT-Ti		192		3	5.7 M	1.2 G	5.9
YOLOS-S	DeiT-S	12	384	224	6	22.1 M	4.5 G	11.8
YOLOS-B	DeiT-B		768		12	86.4 M	17.6 G	23.5
YOLOS-S ( <i>dwr</i> )	–	19	240	272	6	13.7 M	4.6 G	5.0
YOLOS-S ( <i>dwr</i> )	–	14	330	240	6	19.0 M	4.6 G	8.8

Table 1: Variants of YOLOS. “*dwr*” and “*dwr*” refer to uniform compound model scaling and fast model scaling, respectively. The “*dwr*” and “*dwr*” notations are inspired by Dollár et al. [19]. Note that all the numbers listed are for pre-training, which could change during fine-tuning, *e.g.*, the resolution, parameters and FLOPs.

Besides, we investigate two other model scaling strategies which proved to be effective in CNNs. The first one is uniform compound scaling (*dwr*) [19, 53]. In this case, the scaling is uniform w.r.t. FLOPs along all model dimensions (*i.e.*, width ( $w$ ), depth ( $d$ ) and resolution ( $r$ )). The second one is fast scaling (*dwr*) [19] that encourages primarily scaling model width ( $w$ ), while scaling depth ( $d$ ) and resolution ( $r$ ) to a lesser extent w.r.t. FLOPs. During the ImageNet-1k pre-training phase, we apply *dwr* and *dwr* scaling to DeiT-Ti ( $\sim 1.2\text{G}$  FLOPs) and scale the model to  $\sim 4.5\text{G}$  FLOPs to align with the computations of DeiT-S.

For typical CNN architectures, the model complexity or FLOPs ( $f$ ) are proportional to  $dw^2r^2$  [19]. Formally,  $f(\text{CNN}) \propto dw^2r^2$ . Different from CNN, there are two kinds of operations that contribute to the FLOPs of ViT. The first one is the linear projection (Lin.) or point-wise convolution, which fuses the information across different channels point-wisely via learnable parameters. The complexity is  $f(\text{Lin.}) \propto dw^2r^2$ , which is the same as  $f(\text{CNN})$ . The second one is the spatial attention (Att.), which aggregates the spatial information depth-wisely via computed attention weights. The complexity is  $f(\text{Att.}) \propto dwr^4$ , which grows quadratically with the input sequence length or number of pixels.

Note that the available scaling strategies are designed for architectures with complexity  $f \propto dw^2r^2$ , so theoretically the *dwr* as well as *dwr* model scaling are not directly applicable to ViT. However, during pre-training phase the resolution is relatively low, therefore  $f(\text{Lin.})$  dominates the FLOPs ( $\frac{f(\text{Lin.})}{f(\text{Att.})} > 5$ ). Our experiments indicate that some model scaling properties of ViT are consistent with CNNs when  $\frac{f(\text{Lin.})}{f(\text{Att.})}$  is large.

In this paper, we do not study larger models in Dosovitskiy et al. [20] for computing resource constraints.

### 3.2 The Effects of Pre-training

We study the effects of different pre-training strategies (label-supervised and self-supervised) when transferring ViT (*i.e.*, DeiT-Ti and DeiT-S) from ImageNet-1k to the COCO object detection benchmark via YOLOS. For object detection, the input shorter size is 512 for tiny models and is 800 for small models during inference. The results are shown in Tab. 2.

**Necessity of Pre-training.** At least under prevalent transfer learning paradigms [9, 55], the pre-training is necessary in terms of computational efficiency. For both tiny and small models, we find that pre-training on ImageNet-1k saves the total theoretical computations (total pre-training FLOPs & total fine-tuning FLOPs) compared with training on COCO from random initialization (training from scratch [27]). Models trained from scratch with hundreds of epochs still lag far behind the pre-trained ViT even if given more total FLOPs budgets. This seems quite different from typical modern CNN-based detectors, which can catch up with pre-trained counterparts quickly [27].

**Label-supervised Pre-training.** For supervised pre-training with ground truth labels, we find that different-sized models prefer different pre-training schedules: 200 epochs pre-training for YOLOS-Ti still cannot catch up with 300 epochs pre-training even with a 300 epochs fine-tuning schedule, while for the small model 200 epochs pre-training provides feature representations as good as 300 epochs pre-training for transferring to the COCO object detection benchmark.

Model	Pre-train Method	Pre-train Epochs	Fine-tune Epochs	Pre-train pFLOPs	Fine-tune pFLOPs	Total pFLOPs	ImNet Top-1	AP
YOLOS-Ti	Rand. Init.	0	600	0	$14.2 \times 10^2$	$14.2 \times 10^2$	–	19.7
	Label Sup. [55]	200		$3.1 \times 10^2$		$10.2 \times 10^2$	71.2	26.9
	Label Sup. [55]	300	300	$4.7 \times 10^2$	$7.1 \times 10^2$	$11.8 \times 10^2$	72.2	28.7
	Label Sup. () [55]	300		$4.7 \times 10^2$		$11.8 \times 10^2$	74.5	29.7
YOLOS-S	Rand. Init.	0	250	0	$5.9 \times 10^3$	$5.9 \times 10^3$	–	20.9
	Label Sup. [55]	100		$0.6 \times 10^3$		$4.1 \times 10^3$	74.5	32.0
	Label Sup. [55]	200		$1.2 \times 10^3$		$4.7 \times 10^3$	78.5	36.1
	Label Sup. [55]	300	150	$1.8 \times 10^3$	$3.5 \times 10^3$	$5.3 \times 10^3$	79.9	36.1
	Label Sup. () [55]	300		$1.8 \times 10^3$		$5.3 \times 10^3$	81.2	37.2
	DINO Self Sup. [10]	800	150	$4.7 \times 10^3$	$3.5 \times 10^3$	$8.2 \times 10^3$	–	36.2

Table 2: The effects of pre-training. “pFLOPs” refers to petaFLOPs ( $\times 10^{15}$ ). “ImNet” refers to ImageNet-1k. “” refers to the distillation strategy introduced by Touvron et al. [55].

With additional transformer-specific distillation (“”) introduced by Touvron et al. [55], the detection performance is further improved by  $\sim 1$  AP for both tiny and small models, in part because exploiting a CNN teacher [43] during pre-training helps ViT adapt to COCO better. It is also promising to directly leverage [DET] tokens to help smaller YOLOS learn from larger YOLOS on COCO during fine-tuning in a similar way as Touvron et al. [55], we leave it as a future work.

**Self-supervised Pre-training.** The success of Transformer in NLP greatly benefits from large-scale self-supervised pre-training [17, 41, 42]. In vision, pioneering works [11, 20] train self-supervised Transformers following the masked auto-encoding paradigm in NLP. Recent works [10, 12] based on siamese networks show intriguing properties as well as excellent transferability to downstream tasks. We perform a small transfer learning experiment using DINO [10] self-supervised pre-trained weights. For YOLOS-S model, the transfer performance of DINO on COCO object detection is on a par with label-supervised pre-training, suggesting great potentials of self-supervised pre-training for ViT on challenging object-level recognition tasks.

**YOLOS as a Transfer Learning Benchmark for ViT.** From the above analysis, we conclude that the ImageNet-1k pre-training results cannot precisely reflect the transfer learning performance on COCO object detection. Compared with widely used image recognition transfer learning benchmarks such as CIFAR-10/100 [33], Oxford-IIIT Pets [40] and Oxford Flowers-102 [39], YOLOS is more sensitive to the pre-train scheme and the performance is far from saturating. Therefore it is reasonable to consider YOLOS as a challenging transfer learning benchmark to evaluate different (label-supervised or self-supervised) pre-training strategies for ViT.

### 3.3 Pre-training and Transfer Learning Performance of Different Scaled Models

We study the pre-training and the transfer learning performance of different model scaling strategies, *i.e.*, width scaling ( $w$ ), uniform compound scaling ( $dwr$ ) and fast scaling ( $dwr$ ). The models are scaled from  $\sim 1.2G$  to  $\sim 4.5G$  FLOPs for pre-training. Detailed model configurations and descriptions are given in Sec. 3.1 and Tab. 1.

We pre-train all the models for 300 epochs on ImageNet-1k with input resolution determined by the corresponding scaling strategies, and then fine-tune these models on COCO for 150 epochs. Few literatures are available for resolution scaling in object detection, where the inputs are usually oblong in shape and the multi-scale augmentation [9, 26] is used as a common practice. Therefore for each model during inference, we select the smallest resolution (*i.e.*, the shorter size) ranging in [480, 800] producing the highest box AP, which is 784 for  $dwr$  scaling and 800 for all the others. The results are summarized in Tab. 3.

**Pre-training.** Both  $dwr$  and  $dwr$  scaling can improve the accuracy compared with simple  $w$  scaling, *i.e.*, the DeiT-S baseline. Other properties of each scaling strategy are also consistent with CNNs [19, 53]. *e.g.*,  $w$  scaling is the most speed friendly.  $dwr$  scaling achieves the strongest accuracy.  $dwr$  is nearly as fast as  $w$  scaling and is on a par with  $dwr$  scaling in accuracy. Perhaps the reason

Scale	Image Classification @ ImageNet-1k				Object Detection @ COCO val			
	FLOPs	$\frac{f(\text{Lin.})}{f(\text{Att.})}$	FPS	Top-1	FLOPs	$\frac{f(\text{Lin.})}{f(\text{Att.})}$	FPS	AP
-	1.2 G	5.9	1315	72.2	81 G	0.28	12.0	29.6
w	4.5 G	11.8	615	79.9	200 G	0.55	5.7	36.1
dwr	4.6 G	5.0	386	80.5	174 G	0.35	4.5	36.2
dwr	4.6 G	8.8	511	80.4	179 G	0.49	5.4	37.6

Table 3: Pre-training and transfer learning performance of different scaled models. FLOPs and FPS data of object detection are measured over the first 100 images of COCO val split during inference following Carion et al. [9]. FPS is measured with batch size 1 on a single 1080Ti GPU.

why these CNN model scaling strategies are still applicable to ViT is that during pre-training the linear projection ( $1 \times 1$  convolution) dominates the model computations.

**Transfer Learning.** The picture changes when transferred to COCO. The input resolution  $r$  is much higher so the spatial attention takes over and linear projection part is no longer dominant in terms of FLOPs ( $\frac{f(\text{Lin.})}{f(\text{Att.})} \propto \frac{w}{r^2}$ ). Canonical CNN model scaling recipes do not take spatial attention computations into account. Therefore there is some inconsistency between pre-training and transfer learning performance: Despite being strong on ImageNet-1k, the *dwr* scaling achieves similar box AP as simple *w* scaling. Meanwhile, the performance gain from *dwr* scaling on COCO cannot be clearly explained by the corresponding CNN scaling methodology that does not take  $f(\text{Att.}) \propto dwr^4$  into account. The performance inconsistency between pre-training and transfer learning calls for novel model scaling strategies for ViT that considering spatial attention complexity.

### 3.4 Comparisons with CNN-based Object Detectors

In previous sections, we treat YOLOS as a touchstone for the transferability of ViT. In this section, we consider YOLOS as an object detector and we compare YOLOS with some modern CNN detectors.

Method	Backbone	Size	AP	Params. (M)	FLOPs (G)	FPS
YOLOv3-Tiny [45]	DarkNet [45]	$416 \times 416$	16.6	8.9	5.62	330
YOLOv4-Tiny [58]	COSA [58]	$416 \times 416$	21.7	6.1	6.96	371
<b>YOLOS-Ti</b>	DeiT-Ti (⌚) [55]	$256 \times *$	23.1	6.5	3.45	103
CenterNet [68]	ResNet-18 [25]	$512 \times 512$	28.1	—	—	129
YOLOv4-Tiny (3l) [58]	COSA [58]	$320 \times 320$	28.7	—	—	252
Def. DETR [70]	FBNet-V3 [14]	$800 \times *$	27.9	12.2	12.26	35
<b>YOLOS-Ti</b>	DeiT-Ti (⌚) [55]	$432 \times *$	28.6	6.5	12.00	84

Table 4: Comparisons with some tiny-sized modern CNN detectors. All models are trained to be fully converged. “Size” refers to input resolution for inference. FLOPs and FPS data are measured over the first 100 images of COCO val split during inference following Carion et al. [9]. FPS is measured with batch size 1 on a single 1080Ti GPU.

**Comparisons with Tiny-sized CNN Detectors.** As shown in Tab. 4, The tiny-sized YOLOS model achieves impressive performance compared with well-established and highly-optimized CNN object detectors. YOLOS-Ti is strong in AP and is competitive in FLOPs & FPS even though the Transformer is not intentionally designed to optimize these factors. From the model scaling perspective [19, 53, 58], YOLOS-Ti can serve as a promising model scaling start point.

**Comparisons with DETR.** The relations and differences in model design between YOLOS and DETR are given in Sec. 2.1, here we make quantitative comparisons between the two.

As shown in Tab. 5, YOLOS-Ti still performs better than the DETR counterpart, while larger YOLOS models with width scaling become less competitive: YOLOS-S with more computations is 0.8 AP lower compared with a similar-sized DETR model. Even worse, YOLOS-B cannot beat DETR with

Method	Backbone	Epochs	Size	AP	Params. (M)	FLOPs (G)	FPS
Def. DETR [70]	FBNet-V3 [14]	150	800 × *	27.5	12.2	12.26	35
<b>YOLOS-Ti</b>	DeiT-Ti (◐) [55]	300	432 × *	28.6	6.5	12.00	84
<b>YOLOS-Ti</b>	DeiT-Ti (◐) [55]	300	528 × *	30.0	6.5	21.35	51
DETR [9]	ResNet-18-DC5 [25]		800 × *	36.9	28.7	128.9	7.4
<b>YOLOS-S</b>	DeiT-S [55]	150	800 × *	36.1	30.7	200.2	5.7
<b>YOLOS-S (dwr)</b>	DeiT-S [55] (dwr Scale [19])		704 × *	37.2	27.9	127.5	7.2
<b>YOLOS-S (dwr)</b>	DeiT-S [55] (dwr Scale [19])		784 × *	37.6	27.9	179.0	5.4
DETR [9]	ResNet-101-DC5 [25]	150	800 × *	42.5	60	253	–
<b>YOLOS-B</b>	DeiT-B (◐) [55]		800 × *	42.0	127	537	–

Table 5: Comparisons with different DETR models. Tiny-sized models are trained to be fully converged. “Size” refers to input resolution for inference. FLOPs and FPS data are measured over the first 100 images of COCO val split during inference following Carion et al. [9]. FPS is measured with batch size 1 on a single 1080Ti GPU. The “ResNet-18-DC5” implantation is from timm library [62].

over  $2\times$  parameters and FLOPs. Even though YOLOS-S with *dwr* scaling is able to perform better than the DETR counterpart, the performance gain cannot be clearly explained as discussed in Sec. 3.3.

**Meanings of the Results.** Although the performance is seemingly discouraging, the numbers are meaningful, as YOLOS is not purposefully designed for better performance, but designed to precisely reveal the transferability of ViT in object detection. *E.g.*, YOLOS-B is directly adopted from the BERT-Base architecture [17]. This 12 layers, 768 channels Transformer along with its variants have shown impressive performance on a wide range of NLP tasks. We demonstrate that with minimal modifications, this kind of architecture can also be successfully transferred (*i.e.*, AP = 42.0) to the challenging COCO object detection benchmark in computer vision from a pure sequence-to-sequence perspective. The minimal modifications from YOLOS exactly reveal the versatility and generality of Transformer.

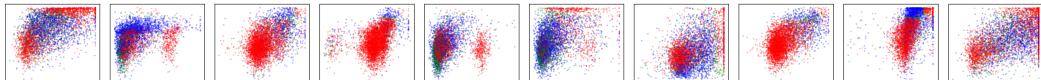


Figure 2: Visualization of all box predictions on all images from COCO val split for the first ten [DET] tokens. Each box prediction is represented as a point with the coordinates of its center normalized by each thumbnail image size. The points are color-coded so that **blue** points corresponds to small objects, **green** to medium objects and **red** to large objects. We observe that each [DET] token learns to specialize on certain regions and sizes. The visualization style is inspired by Carion et al. [9].

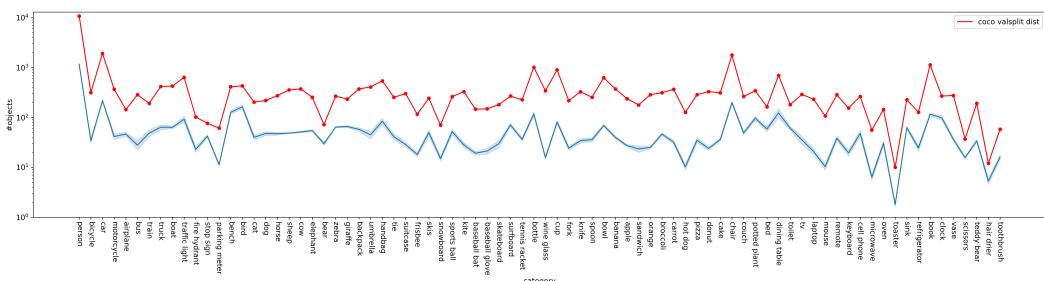


Figure 3: The statistics of all ground truth object categories (the **red** curve) and the statistics of all object category predictions from all [DET] tokens (the **blue** curve) on all images from COCO val split. The error bar of the **blue** curve represents the variability of the preference of different tokens for a given category, which is small. This suggests that different [DET] tokens are category insensitive.

**Towards Better Bigger Transformer.** Transformer blocks are more complicated than convolution kernels for there are two kinds of computations with different complexities (see Sec. 3.1 & Sec. 3.3). Moreover, in addition to macroscopic scaling dimensions such as depth, width, and resolution, there are other degrees of freedom that could be scalable in ViT, *e.g.*, the width (number of channels) of Query (Key) & Value in multi-head self-attention, the bottleneck ratio (could be inverted) in MLP, as well as the number of attention heads. From the above analysis, the tiny-sized YOLOS model is comparable with modern CNN object detectors, suggesting a promising model scaling start point. However, larger YOLOS models with simple width scaling [55] produce relatively less competitive detection results. Therefore larger models call for better scaling strategies tailored for ViT. We also hope the methodology of model scaling in computer vision can inspire the Transformer design in NLP.

**Inspecting Detection Tokens.** As an object detector, YOLOS uses [DET] tokens to represent detected objects. We find that different [DET] tokens are sensitive to object locations and sizes, while insensitive to object categories, as shown in Fig. 2 and Fig. 3.

## 4 Related Work

**Vision Transformer for Object Detection.** There has been a lot of interest in combining CNNs with forms of self-attention mechanisms [56] to improve object detection performance [8, 30, 61], while recent works trend towards augmenting Transformer with CNNs (or CNN design). Beal et al. [5] propose to use a pre-trained ViT as the feature extractor for an R-CNN [22] object detector. Despite being effective, they fail to ablate the CNN architectures, region-wise pooling operations [21, 24, 26] as well as hand-crafted components such as dense anchors [47] and NMS. Inspired by modern CNN architecture, some works [37, 57, 60, 63] introduce the pyramidal feature hierarchy and locality to Vision Transformer design, which largely boost the performance in dense prediction tasks including object detection. However, these architectures are performance-oriented and cannot reflect the properties of the naïve or vanilla Vision Transformer [20] that directly inherited from Vaswani et al. [56]. Another series of work, the DEtection TRansformer (DETR) families [9, 70], use a random initialized Transformer to encode & decode CNN features for object detection, which does not reveal the transferability of a pre-trained Transformer.

In this paper, we argue for the characteristics of a pre-trained naïve Vision Transformer [20] in object detection, which is rare in the existing literature.

**Pre-training and fine-tuning.** The textbook-style usage of Transformer [56] follows a “pre-training & fine-tuning” paradigm. In NLP, large transformer-based models are often pre-trained on large corpora and then fine-tuned for different tasks at hand [17, 41]. In computer vision, Dosovitskiy et al. [20] apply Transformer to image recognition at scale using modern vision transfer learning recipe [32]. They show that a standard Transformer encoder architecture is able to attain excellent results on mid-sized or small image recognition benchmarks (*e.g.*, ImageNet-1k [48], CIFAR-10/100 [33], *etc.*) when pre-trained at sufficient scale (*e.g.*, JFT-300M [52], ImageNet-21k [16]). Touvron et al. [55] achieves competitive Top-1 accuracy by training Transformer on ImageNet-1k only, and is also capable of transferring to smaller downstream tasks [33, 39, 40]. However, existing transfer learning literature of Transformer arrest in image-level recognition and does not touch more complex tasks in vision such as object detection, which is widely used to benchmark CNNs transferability.

Our work aims to bridge this gap. We study the performance and properties of ViT on the challenging COCO [35] object detection benchmark when pre-trained on the mid-sized ImageNet-1k dataset using different strategies.

## 5 Conclusion and Future Work

In this paper, we have explored the transferability of the vanilla ViT pre-trained on mid-sized ImageNet-1k dataset to the more challenging COCO object detection benchmark. We demonstrate that 2D object detection can be accomplished in a pure sequence-to-sequence manner with minimal additional inductive biases. The performance on COCO is promising, and these initial results are meaningful, suggesting the versatility and generality of Transformer to various downstream tasks.

There are still many challenges that remain and needed to be resolved in the future. One is the long input sequence / high input resolution in object detection as well as other dense prediction tasks. Since the self-attention operation scales quadratically with the sequence length, we are unable to touch larger models in this work. Transformers such as [6, 15] that can efficiently process thousands of tokens are urgently needed. Another one is the model scaling method tailored for ViT, for there are two kinds of operations that contribute to the FLOPs of one Transformer layer, which is different from CNN.

## 6 Appendix

### 6.1 Mathematical Formulation of YOLOS

The formulation of YOLOS basically follows Dosovitskiy et al. [20].

**Stem.** The standard ViT [20] receives a 1D sequence of embedded tokens as the input. To handle 2D image inputs, we reshape the image  $\mathbf{x} \in \mathbb{R}^{H \times W \times C}$  into a sequence of flattened 2D patch tokens  $\mathbf{x}_p \in \mathbb{R}^{N \times (P^2 \cdot C)}$ . Here,  $(H, W)$  is the resolution of the input image,  $C$  is the number of input channels,  $(P, P)$  is the resolution of each image patch, and  $N = \frac{HW}{P^2}$  is the resulting number of patches, which also serves as the effective input sequence length for YOLOS. Then we flatten the patches  $\mathbf{x}_p$  and map to  $D$  dimensions with a trainable linear projection  $\mathbf{E}$ <sup>5</sup>. We refer to the output of this projection as the patch tokens. Meanwhile, one hundred [DET] tokens  $\mathbf{x}_{\text{DET}}$  are appended to the patch tokens. Position embeddings  $\mathbf{E}_{PE}$  are added to all the input tokens to retain positional information. We use standard learnable 1D position embeddings following Dosovitskiy et al. [20]. The resulting sequence serves as the input of YOLOS. Formally:

$$\mathbf{z}_0 = [\mathbf{x}_p^1 \mathbf{E}; \mathbf{x}_p^2 \mathbf{E}; \dots; \mathbf{x}_p^N \mathbf{E}; \mathbf{x}_{\text{DET}}^1; \dots; \mathbf{x}_{\text{DET}}^{100}] + \mathbf{E}_{PE}, \quad \mathbf{E} \in \mathbb{R}^{(P^2 \cdot C) \times D}, \mathbf{E}_{PE} \in \mathbb{R}^{(N+100) \times D} \quad (1)$$

**Backbone.** The backbone of YOLOS is exactly the same as ViT, which consists of a stack of Transformer encoder layers [56]. One Transformer encoder layer consists of one multi-head self-attention (MSA) block and one MLP block. LayerNorm (LN) [2] is applied before every block, and residual connections [25] are applied after every block [3, 59]. The MLP contains two hidden layers with a GELU [28] non-linearity. Formally:

$$\begin{aligned} \mathbf{z}'_\ell &= \text{MSA}(\text{LN}(\mathbf{z}_{\ell-1})) + \mathbf{z}_{\ell-1}, & \ell &= 1 \dots L \\ \mathbf{z}_\ell &= \text{MLP}(\text{LN}(\mathbf{z}'_\ell)) + \mathbf{z}'_\ell, & \ell &= 1 \dots L \\ \mathbf{y} &= \text{LN}(\mathbf{z}_L^0) \end{aligned} \quad (2)$$

**Detector Heads.** Both the classification head and the bounding box regression head are implemented by a MLP with two hidden layers.

### 6.2 Position Embedding (PE) of YOLOS

In object detection and many other computer vision benchmarks, the image resolutions as well as the aspect ratios are usually not fixed as the image classification task. Due to the changes in input resolutions & aspect ratios (sequence length) from the image classification task to the object detection task, the position embedding (PE) in ViT / YOLOS has also to be changed and adapted<sup>6</sup>. The changes in PE could affect the model size and performance. In this work, we study two types of PE settings for YOLOS:

- Type-I adds randomly initialized PE to the input of each intermediate Transformer layer as DETR [9], and the PE is 1D learnable (considering the inputs as a sequence of patches in the raster order) as ViT [20]. For the first layer, the PE is interpolated following ViT. The size of PEs is usually smaller than the input sequence size considering the model parameters. In the paper, small- and base-sized models use this setting.
- Type-II interpolates the pre-trained 1D learnable PE to a size similar to or slightly larger than the input size, and adds no PE in intermediate Transformer layers. In the paper, tiny-sized models use this setting.

In a word, Type-I uses more PEs and Type-II uses larger PE.

---

<sup>5</sup>This set of operations is equivalent to a convolution with `kernel_size = P × P, stride = P, in_plane = P2 · C` and `out_plane = D`, as implemented by Wightman [62].

<sup>6</sup>PE for one hundred [DET] tokens is not affected.

**Type-I PE.** This setting adds PE to the input of each Transformer layer following DETR [9], and the PE considering the inputs as a sequence of patches in the raster order following ViT [20]. Specifically, during fine-tuning, the PE of the first layer is interpolated from the pre-trained one, and the PEs for the rest intermediate layers are randomly initialized and trained from scratch. In our paper, small- and base-sized models use this setting. The detailed configurations are given in Tab. 6.

Model	PE-cls to PE-det @ First Layer	Rand. Init. PE-det @ Mid. Layer	cls → det Params. (M)
YOLOS-S	$\frac{224}{16} \times \frac{224}{16} \nearrow \frac{512}{16} \times \frac{864}{16}$	$\frac{512}{16} \times \frac{864}{16}$	$22.1 \rightarrow 30.7$
YOLOS-S ( <i>dwr</i> )	$\frac{224}{16} \times \frac{224}{16} \nearrow \frac{800}{16} \times \frac{1344}{16}$	$\frac{800}{16} \times \frac{1344}{16}$	$13.7 \rightarrow 22.0$
YOLOS-S ( <i>dwr</i> )	$\frac{224}{16} \times \frac{224}{16} \nearrow \frac{512}{16} \times \frac{864}{16}$	$\frac{512}{16} \times \frac{864}{16}$	$19.0 \rightarrow 27.6$
YOLOS-B	$\frac{384}{16} \times \frac{384}{16} \nearrow \frac{800}{16} \times \frac{1344}{16}$	$\frac{800}{16} \times \frac{1344}{16}$	$86.4 \rightarrow 127.8$

Table 6: Type-I PE configurations for YOLOS models. “PE-cls  $\nearrow$  PE-det” refers to performing 2D interpolation of ImageNet-1k pre-trained PE-cls to PE-det for object detection. The PEs added in the intermediate (Mid.) layers (all the other layers of YOLOS except the first layer) are randomly initialized.

From Tab. 6, we conclude that it is expensive in terms of model size to use intermediate PEs for object detection. In other words, about  $\frac{1}{3}$  of the model weights is for providing positional information only. Despite being heavy, we argue that the randomly initialized intermediate PEs do not directly inject additional inductive biases and they learn the positional relation from scratch. Nevertheless, for multi-scale inputs during training or input with different sizes & aspect ratios during inference, we (have to) adjust the PE size via 2D interpolation on the fly<sup>7</sup>. As mentioned in Dosovitskiy et al. [20] and in the paper, this operation could introduce inductive biases.

To control the model size, these intermediate PE sizes are usually set to be smaller than the input sequence length, *e.g.*, for typical models YOLOS-S and YOLOS-S (*dwr*), the PE size is  $\frac{512}{16} \times \frac{864}{16}$ . Since the *dwr* scaling is more parameter friendly compared with other model scaling approaches, we use a larger PE for YOLOS-S (*dwr*) than other small-sized models to compensate for the number of parameters. For larger models such as YOLOS-Base, we do not consider the model size so we also choose to use larger PE.

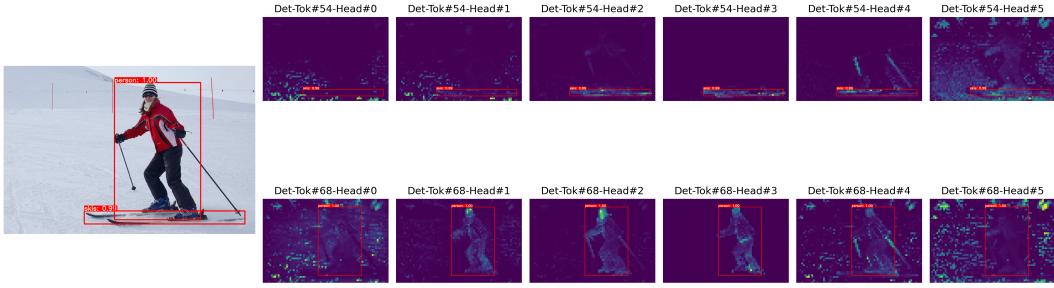
Using 2D PE can save a lot of parameters, *e.g.*, DETR uses two long enough PE (Length = 50 for regular models and Length = 100 for DC5 models) for both  $x$  and  $y$  axes. We don’t consider 2D PE in this work.

Model	PE Type	PE-cls to PE-det @ First Layer	Rand. Init. PE-det @ Rest Layer	Params. (M) cls → det	AP
YOLOS-Ti	Type-I	$\frac{224}{16} \times \frac{224}{16} \nearrow \frac{512}{16} \times \frac{864}{16}$	$\frac{512}{16} \times \frac{864}{16}$	$5.7 \rightarrow 9.9$	28.3
	Type-II	$\frac{224}{16} \times \frac{224}{16} \nearrow \frac{800}{16} \times \frac{1344}{16}$	No PE	$5.7 \rightarrow 6.5$	28.7
YOLOS-S	Type-I	$\frac{224}{16} \times \frac{224}{16} \nearrow \frac{512}{16} \times \frac{864}{16}$	$\frac{512}{16} \times \frac{864}{16}$	$22.1 \rightarrow 30.7$	36.1
	Type-II	$\frac{224}{16} \times \frac{224}{16} \nearrow \frac{960}{16} \times \frac{1600}{16}$	No PE	$22.1 \rightarrow 24.6$	36.6

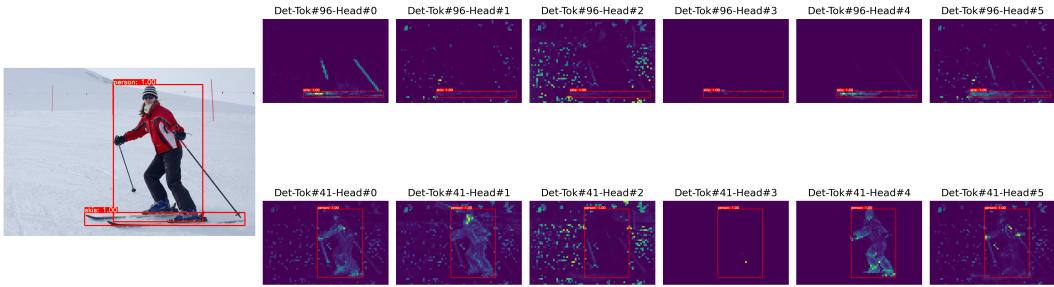
Table 7: Some instantiations of Type-II PE. They are lighter and better than Type-I counterparts.

**Type-II PE.** Later, we find that interpolating the pre-trained PE at the first layer to a size similar to or larger than the input sequence length as the only PE can provide enough positional information, and is more efficient than using more smaller-sized PEs in the intermediate layers. In other words, it is redundant to use intermediate PEs given one large enough PE in the first layer. Some instantiations are shown in Tab. 7. In the paper, tiny-sized models use this setting. This type of PE is more promising, and we will make a profound study about this setting in the future.

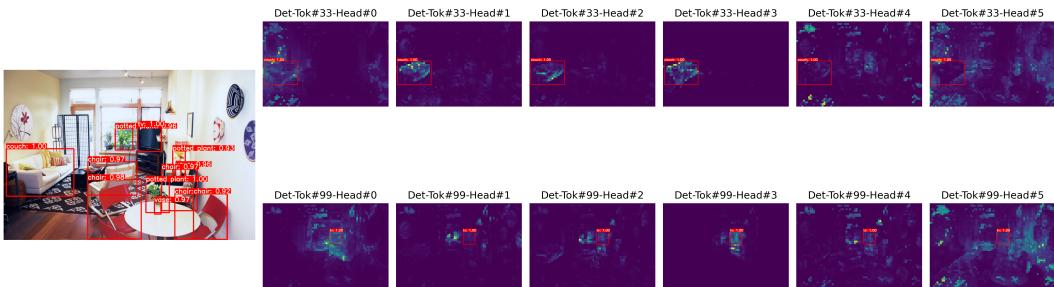
<sup>7</sup>There are some kind of data augmentations that can avoid PE interpolation, *e.g.*, large scale jittering used in Tan et al. [54], which randomly resizes images between  $0.1\times$  and  $2.0\times$  of the original size then crops to a fixed resolution. However, scale jittering augmentation usually requires longer training schedules, in part because when the original input image is resized to a higher resolution, the cropped image usually has a smaller number of objects than the original, which could weaken the supervision signal therefore needs longer training to compensate. So there is no free lunch.



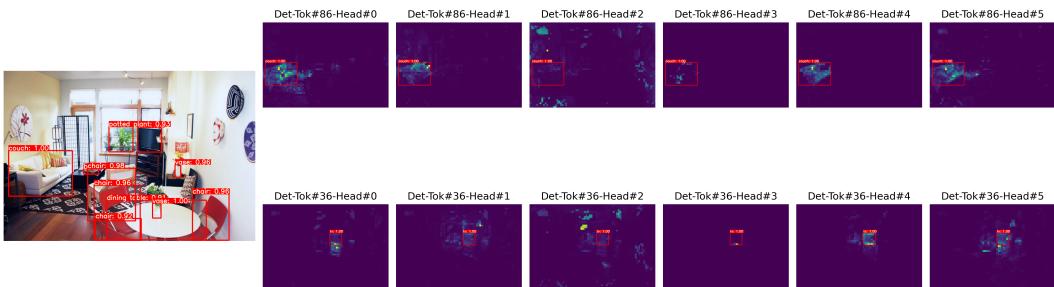
(a) YOLOS-S, 200 epochs pre-trained, COCO AP = 36.1.



(b) YOLOS-S, 300 epochs pre-trained, COCO AP = 36.1.

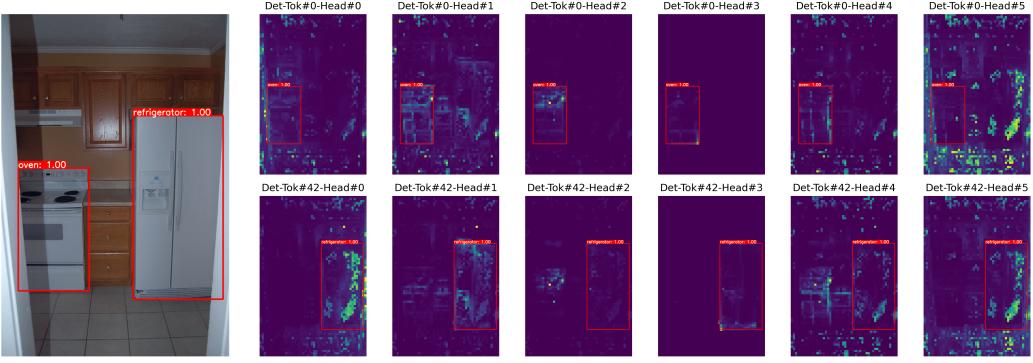


(c) YOLOS-S, 200 epochs pre-trained, COCO AP = 36.1.

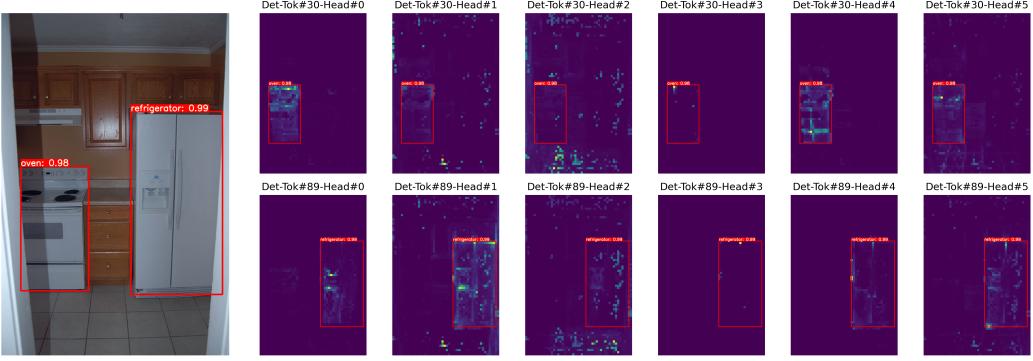


(d) YOLOS-S, 300 epochs pre-trained, COCO AP = 36.1.

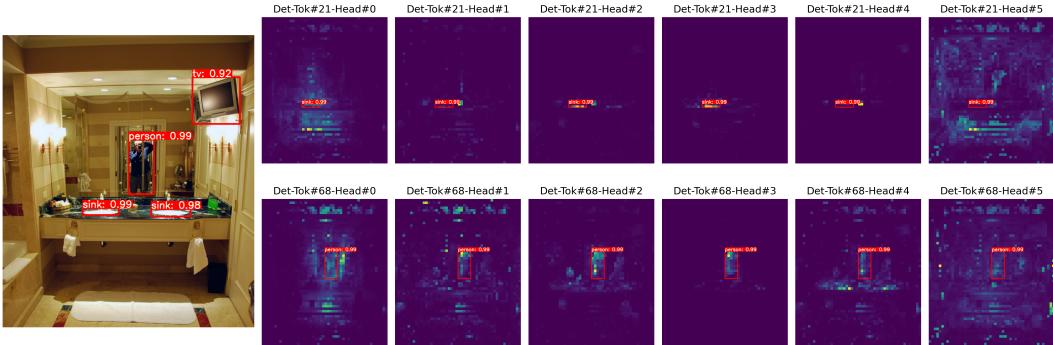
Figure 4: The self-attention map visualization of the [DET] tokens and the corresponding predictions on the heads of the last layer of two different YOLOS-S models.



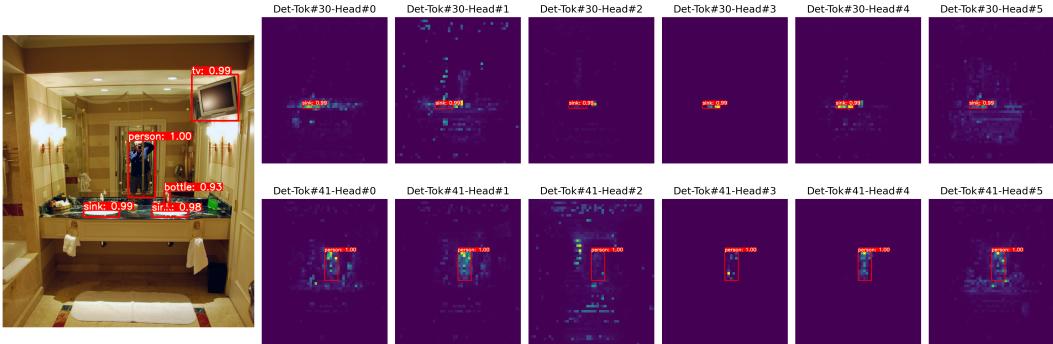
(a) YOLOS-S, 200 epochs pre-trained, COCO AP = 36.1.



(b) YOLOS-S, 300 epochs pre-trained, COCO AP = 36.1.



(c) YOLOS-S, 200 epochs pre-trained, COCO AP = 36.1.



(d) YOLOS-S, 300 epochs pre-trained, COCO AP = 36.1.

Figure 5: The self-attention map visualization of the [DET] tokens and the corresponding predictions on the heads of the last layer of two different YOLOS-S models.

### 6.3 Self-attention Maps of YOLOS

We inspect the self-attention of the [DET] tokens that related to the predictions on the heads of the last layer of YOLOS-S. The visualization pipeline follows Caron et al. [10]. The visualization results are shown in Fig. 4 & Fig. 5. We conclude that:

- For a given YOLOS model, different self-attention heads focus on different patterns & different locations. Some visualizations are interpretable while others are not.
- We study the attention map differences of two YOLOS models, *i.e.*, the 200 epochs ImageNet-1k [48] pre-trained YOLOS-S and the 300 epochs ImageNet-1k pre-trained YOLOS-S. Note that the AP of these two models is the same (AP= 36.1). From the visualization, we conclude that for a given predicted object, the corresponding [DET] token as well as the attention map patterns are usually different for different models.

## References

- [1] Edward H Adelson, Charles H Anderson, James R Bergen, Peter J Burt, and Joan M Ogden. Pyramid methods in image processing. *RCA engineer*, 1984.
- [2] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [3] Alexei Baevski and Michael Auli. Adaptive input representations for neural language modeling. *arXiv preprint arXiv:1809.10853*, 2018.
- [4] Peter W Battaglia, Jessica B Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinicius Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, et al. Relational inductive biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*, 2018.
- [5] Josh Beal, Eric Kim, Eric Tzeng, Dong Huk Park, Andrew Zhai, and Dmitry Kislyuk. Toward transformer-based object detection. *arXiv preprint arXiv:2012.09958*, 2020.
- [6] Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*, 2020.
- [7] Samuel R Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. A large annotated corpus for learning natural language inference. *arXiv preprint arXiv:1508.05326*, 2015.
- [8] Yue Cao, Jiarui Xu, Stephen Lin, Fangyun Wei, and Han Hu. Gcnet: Non-local networks meet squeeze-excitation networks and beyond. In *ICCV*, 2019.
- [9] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020.
- [10] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. *arXiv preprint arXiv:2104.14294*, 2021.
- [11] Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever. Generative pretraining from pixels. In *ICML*, 2020.
- [12] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. *arXiv preprint arXiv:2104.02057*, 2021.
- [13] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *CVPRW*, 2020.
- [14] Xiaoliang Dai, Alvin Wan, Peizhao Zhang, Bichen Wu, Zijian He, Zhen Wei, Kan Chen, Yuandong Tian, Matthew Yu, Peter Vajda, et al. Fbnetv3: Joint architecture-recipe search using neural acquisition function. *arXiv preprint arXiv:2006.02049*, 2020.
- [15] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V Le, and Ruslan Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. *arXiv preprint arXiv:1901.02860*, 2019.
- [16] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009.

- [17] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [18] William B Dolan and Chris Brockett. Automatically constructing a corpus of sentential paraphrases. In *IWP*, 2005.
- [19] Piotr Dollár, Mannat Singh, and Ross Girshick. Fast and accurate model scaling. *arXiv preprint arXiv:2103.06877*, 2021.
- [20] 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.
- [21] Ross Girshick. Fast r-cnn. In *ICCV*, 2015.
- [22] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR*, 2014.
- [23] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016.
- [24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Spatial pyramid pooling in deep convolutional networks for visual recognition. *TPAMI*, 2015.
- [25] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [26] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross B. Girshick. Mask r-cnn. In *ICCV*, 2017.
- [27] Kaiming He, Ross Girshick, and Piotr Dollár. Rethinking imagenet pre-training. In *ICCV*, 2019.
- [28] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016.
- [29] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- [30] Han Hu, Jiayuan Gu, Zheng Zhang, Jifeng Dai, and Yichen Wei. Relation networks for object detection. In *CVPR*, 2018.
- [31] Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q Weinberger. Deep networks with stochastic depth. In *ECCV*, 2016.
- [32] Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly, and Neil Houlsby. Big transfer (bit): General visual representation learning. *arXiv preprint arXiv:1912.11370*, 2019.
- [33] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [34] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *NeurIPS*, 2012.
- [35] 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 *ECCV*, 2014.
- [36] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *CVPR*, 2017.
- [37] 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. *arXiv preprint arXiv:2103.14030*, 2021.
- [38] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- [39] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In *ICVGIP*, 2008.
- [40] Omkar M. Parkhi, Andrea Vedaldi, Andrew Zisserman, and C. V. Jawahar. Cats and dogs. In *CVPR*, 2012.

- [41] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. 2018.
- [42] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI blog*, 2019.
- [43] Ilija Radosavovic, Raj Prateek Kosaraju, Ross Girshick, Kaiming He, and Piotr Dollár. Designing network design spaces. In *CVPR*, 2020.
- [44] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*, 2016.
- [45] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018.
- [46] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *CVPR*, 2016.
- [47] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *arXiv preprint arXiv:1506.01497*, 2015.
- [48] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *IJCV*, 2015.
- [49] Erik F Sang and Fien De Meulder. Introduction to the conll-2003 shared task: Language-independent named entity recognition. *arXiv preprint cs/0306050*, 2003.
- [50] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [51] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *JMLR*, 2014.
- [52] Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. Revisiting unreasonable effectiveness of data in deep learning era. In *ICCV*, 2017.
- [53] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *ICML*, 2019.
- [54] Mingxing Tan, Ruoming Pang, and Quoc V Le. Efficientdet: Scalable and efficient object detection. In *CVPR*, 2020.
- [55] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. *arXiv preprint arXiv:2012.12877*, 2020.
- [56] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *arXiv preprint arXiv:1706.03762*, 2017.
- [57] Ashish Vaswani, Prajit Ramachandran, Aravind Srinivas, Niki Parmar, Blake Hechtman, and Jonathon Shlens. Scaling local self-attention for parameter efficient visual backbones. *arXiv preprint arXiv:2103.12731*, 2021.
- [58] Chien-Yao Wang, Alexey Bochkovskiy, and Hong-Yuan Mark Liao. Scaled-yolov4: Scaling cross stage partial network. *arXiv preprint arXiv:2011.08036*, 2020.
- [59] Qiang Wang, Bei Li, Tong Xiao, Jingbo Zhu, Changliang Li, Derek F Wong, and Lidia S Chao. Learning deep transformer models for machine translation. *arXiv preprint arXiv:1906.01787*, 2019.
- [60] Wenhui 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. *arXiv preprint arXiv:2102.12122*, 2021.
- [61] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *CVPR*, 2018.
- [62] Ross Wightman. Pytorch image models. <https://github.com/rwightman/pytorch-image-models>, 2019.

- [63] Weijian Xu, Yifan Xu, Tyler Chang, and Zhuowen Tu. Co-scale conv-attentional image transformers. *arXiv preprint arXiv:2104.06399*, 2021.
- [64] Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *ICCV*, 2019.
- [65] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.
- [66] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017.
- [67] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random erasing data augmentation. In *AAAI*, 2020.
- [68] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. *arXiv preprint arXiv:1904.07850*, 2019.
- [69] Benjin Zhu, Jianfeng Wang, Zhengkai Jiang, Fuhang Zong, Songtao Liu, Zeming Li, and Jian Sun. Autoassign: Differentiable label assignment for dense object detection. *arXiv preprint arXiv:2007.03496*, 2020.
- [70] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020.