

# You Only Look One-level Feature

Qiang Chen<sup>1,2\*</sup>, Yingming Wang<sup>4</sup>, Tong Yang<sup>4</sup>, Xiangyu Zhang<sup>4</sup>, Jian Cheng<sup>1,2,3†</sup>, Jian Sun<sup>4</sup>

<sup>1</sup>NLPR, Institute of Automation, Chinese Academy of Sciences

<sup>2</sup>School of Artificial Intelligence, University of Chinese Academy of Sciences

<sup>3</sup>CAS Center for Excellence in Brain Science and Intelligence Technology

<sup>4</sup>MEGVII Technology

{qiang.chen, jcheng}@nlpr.ia.ac.cn, {wangyingming, yangtong, zhangxiangyu, sunjian}@megvii.com

## Abstract

This paper revisits feature pyramids networks (FPN) for one-stage detectors and points out that *the success of FPN is due to its divide-and-conquer solution to the optimization problem in object detection rather than multi-scale feature fusion*. From the perspective of optimization, we introduce an alternative way to address the problem instead of adopting the complex feature pyramids - utilizing only one-level feature for detection. Based on the simple and efficient solution, we present *You Only Look One-level Feature (YOLOF)*. In our method, two key components, *Dilated Encoder and Uniform Matching*, are proposed and bring considerable improvements. Extensive experiments on the COCO benchmark prove the effectiveness of the proposed model. Our YOLOF achieves comparable results with its feature pyramids counterpart RetinaNet while being 2.5× faster. Without transformer layers, YOLOF can match the performance of DETR in a single-level feature manner with 7× less training epochs. With an image size of 608 × 608, YOLOF achieves 44.3 mAP running at 60 fps on 2080Ti, which is 13% faster than YOLOv4. Code is available at <https://github.com/megvii-model/YOLOF>.

## 1. Introduction

In state-of-the-art two-stage detectors [22, 13, 3] and one-stage detectors [23, 38], feature pyramids become an essential component. The most popular way to build feature pyramids is the feature pyramid networks (FPN) [22], which mainly brings two benefits: (1) *multi-scale feature fusion*: fusing multiple low-resolution and high-resolution feature inputs to obtain better representations; (2) *divide-and-conquer*: detecting objects on different levels regarding

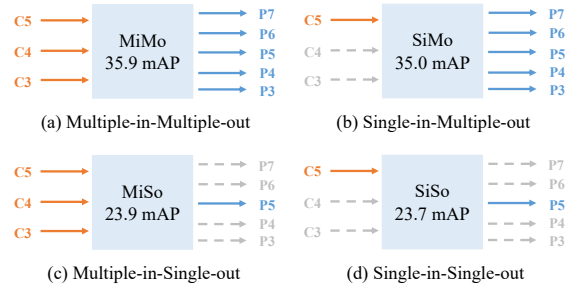


Figure 1. Comparison of box AP among the Multiple-in-Multiple-out (MiMo), Single-in-Multiple-out (SiMo), Multiple-in-Single-out (MiSo), and Single-in-Single-out (SiSo) encoders on COCO validation set. Here, we adopt the original RetinaNet [23] as our baseline model, where C3, C4, and C5 denote output features of the backbone with a downsample rate of {8, 16, 32} and P3 to P7 represent the feature levels used for final detection. All results reported in the figure use the same backbone, ResNet-50 [14]. The structure of MiMo is same as the FPN in RetinaNet [23]. A detailed illustration of the structure for all encoders can be found in the Figure 8.

objects’ scales. A common belief for FPN is that its success relies on the fusion of multiple level features, inducing a line of studies of designing complex fusion methods manually [25, 17, 28], or via Neural Architecture Search (NAS) algorithms [9, 37]. However, the belief ignores the function of the divide-and-conquer in FPN. It leads to fewer studies on how these two benefits contribute to FPN’s success and may hinder new advances.

This paper studies the influence of FPN’s two benefits in one-stage detectors. We design experiments by decoupling the *multi-scale feature fusion* and the *divide-and-conquer* functionalities with RetinaNet [23]. In detail, we consider FPN as a *Multiple-in-Multiple-out (MiMo) encoder*, which encodes multi-scale features from the backbone and provides feature representations for the decoder (the detection heads). We conduct controlled comparisons among *Multiple-in-Multiple-out (MiMo)*, *Single-in-*

\*This work is done during Qiang Chen’s internship at MEGVII Technology.

†Corresponding author.

*Multiple-out* (SiMo), *Multiple-in-Single-out* (MiSo), and *Single-in-Single-out* (SiSo) encoders in Figure 1. Surprisingly, the SiMo encoder, which only has one input feature C5 and does not perform feature fusion, can achieve comparable performance with the MiMo encoder (i.e., FPN). The performance gap is less than 1 mAP. In contrast, the performance drops dramatically ( $\geq 12$  mAP) in MiSo and SiSo encoders. These phenomena suggest two facts: (1) the C5 feature carries sufficient context for detecting objects on various scales, which enables the SiMo encoder to achieve comparable results; (2) the *multi-scale feature fusion* benefit is far away less critical than the *divide-and-conquer* benefit, thus *multi-scale feature fusion might not be the most significant benefit of FPN*, which is also demonstrated by ExFuse [50] in semantic segmentation. Thinking one step deeper, *divide-and-conquer* is related to the optimization problem in object detection. It divides the complex detection problem into several sub-problems by object scales, facilitating the optimization process.

The above analysis suggests that the essential factor for the success of FPN is its solution to the optimization problem in object detection. The *divide-and-conquer* solution is a good way. But it brings memory burdens, slows down the detectors, and make detectors' structure complex in one-stage detectors like RetinaNet [23]. Given that the C5 feature carries sufficient context for detection, we show a simple way to address the optimization problem.

We propose *You Only Look One-level Feature (YOLOF)*, which only uses one single C5 feature (with a downsample rate of 32) for detection. To bridge the performance gap between the SiSo encoder and the MiMo encoder, we first design the structure of the encoder properly to extract the multi-scale contexts for objects on various scales, compensating for the lack of multiple-level features; then, we apply a uniform matching mechanism to solve the imbalance problem of positive anchors raised by the sparse anchors in the single feature.

Without bells and whistles, YOLOF achieves comparable results with its feature pyramids counterpart RetinaNet [23] but  $2.5\times$  faster. In a single feature manner, YOLOF matches the performance of the recent proposed DETR [4] while converging much faster ( $7\times$ ). With an image size of  $608 \times 608$  and other techniques [1, 47], YOLOF achieve 44.3 mAP running at 60 fps on 2080Ti, which is 13% faster than YOLOv4 [1]. In a nutshell, the contributions of this paper are:

- We show that the most significant benefits of FPN is its *divide-and-conquer* solution to the optimization problem in dense object detection rather than the *multi-scale feature fusion*.
- We present YOLOF, which is a simple and efficient baseline without using FPN. In YOLOF, we propose

two key components, *Dilated Encoder* and *Uniform Matching*, bridging the performance gap between the SiSo encoder and the MiMo encoder.

- Extensive experiments on COCO benchmark indicates the importance of each component. Moreover, we conduct comparisons with RetinaNet [23], DETR [4] and YOLOv4 [1]. We can achieve comparable results with a faster speed on GPUs.

## 2. Related Works

**Multiple-level feature detectors.** It is a conventional technique to employ multiple features for object detection. Typical approaches to construct multiple features can be categorized into image pyramid methods and feature pyramid methods. Image pyramids based detector such as DPM [8] dominates the detection in the pre-deep learning era. In CNN-based detectors, the image pyramids method also wins some researchers' [34, 35] praise as it can achieve higher performance out of the box. However, the image pyramids method is not the only way to obtain multiple features; it is more efficient and natural to exploit feature pyramids' power in CNN models. SSD [26] first utilizes multiple-scale features and performs object detection on each scale for different scales objects. FPN [22] follows SSD [26] and UNet [33] and constructs semantic-riched feature pyramids by combining shallow features and deep features. After that, several works [17, 25, 9, 37] follow FPN and focus on how to obtain better representations. FPN becomes an essential component and dominates modern detectors. It is also applied to popular one-stage detectors, such as RetinaNet [23], FCOS [38], and their variants [48]. Another line of method to get feature pyramids is to use multi-branch and dilation convolution [20]. Different from the above works, our method is a single-level feature detector.

**Single-level feature detectors.** In early times, the R-CNN series [11, 10, 31] and R-FCN [6] only extract RoI features on a single feature, while their performances lag behind their multiple feature counterparts [22]. Also, in one-stage detectors, YOLO [29] and YOLOv2 [30] only use the last output feature of the backbone. They can be super fast but have to bear a performance decline in detection. CornerNet [19] and CenterNet [51, 7] follow this fashion and achieve competitive results while using a single feature with a downsample rate of 4 to detect all the objects. Using a high-resolution feature map for detection brings enormous memory cost and is not friendly to deployment. Recently, DETR [4] introduces the transformer [39] to detection and shows that it could achieve state-of-the-art results only use a single C5 feature. Due to the totally anchor-free mechanism and transformer learning phase, DETR needs a

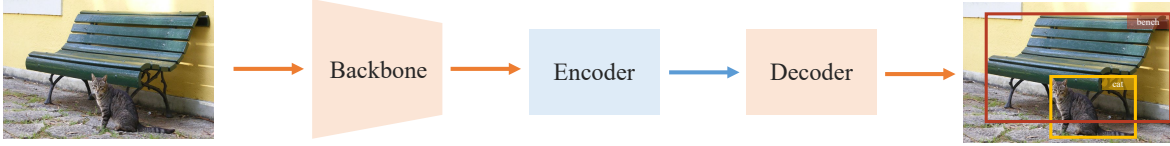


Figure 2. An illustration of the detection pipeline. In this paper, we format the detection pipeline into three parts: (1) the backbone; (2) the encoder, which receives inputs from the backbone and distributes representations for detection; (3) the decoder, which performs classification and regression tasks and generate final prediction boxes. The color for the encoder is corresponding to the one in Figure 1.

long training schedule for its convergence. The long training schedule characteristic is cumbersome for further improvements. Unlike these papers, we investigate the working mechanism of multiple-level detection. From the perspective of optimization, we provide an alternative solution to the widely used FPN. Moreover, YOLOF converges faster and achieves promising performance; thus, YOLOF can serve as a simple baseline for fast and accurate detectors.

### 3. Cost Analysis of MiMo Encoders

As mentioned in Section 1, the success of FPN in dense object detection is due to its solution to the optimization problem. However, the multi-level feature paradigm is inevitable to make detectors complex, brings memory burdens, and slows down the detector. In this section, we provide a quantitative study on the cost of MiMo encoders.

We design experiments based on RetinaNet [23] with ResNet-50 [14]. In detail, we format the pipeline for the detection task as a combination of three key parts: the backbone, the encoder, and the decoder (Figure 2). In this view, we show the FLOPs of each component in Figure 3. Compared with SiSo encoders, the MiMo encoder brings enormous memory burdens to the encoder and the decoder (134G vs. 6G) (Figure 3). Moreover, the detector with MiMo encoder runs much slower than the ones with SiSo encoders (13 FPS vs. 34 FPS) (Figure 3). The slow speed is caused by detecting objects on high-resolution feature maps in the detector with MiMo encoder, such as the C3 feature (with a downsample rate of 8). Given the above drawbacks of the MiMo encoder, we aim to find an alternative way to solve the optimization problem while keeping the detector simple, accurate, and fast simultaneously.

### 4. Method

Motivated by the above purpose and the finding that the C5 feature contains enough context for detecting numerous objects, we try to replace the complex MiMo encoder with the simple SiSo encoder in this section. But this replacement is **nontrivial** as the performance drops extensively when applying SiSo encoders according to the results in Figure 3. Given the situation, we carefully analyze the obstacles preventing SiSo encoders from getting a comparable

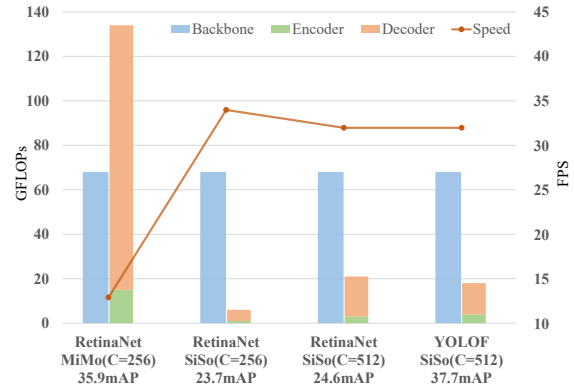


Figure 3. FLOPs, accuracy, and speed comparison between the models that adopt MiMo and SiSo encoders on COCO. As the FLOPs of the decoder is affected by the encoder’s outputs, we stack the FLOPs of the encoder and the decoder in the figure to better understanding the effects of encoders on the FLOPs. All models use the same backbone, ResNet-50. All FLOPs are measured with a shorter edge size 800 over the first 100 images of COCO val2017. The FPS is calculated with batch size 1 on 2080Ti from the total inference pure compute time reported in the Detectron2 [42]. In the figure,  $C$  represents the number of channels used in the model’s encoder and decoder.

performance with MiMo encoders. We find that two problems brought by SiSo encoders are responsible for the performance drop. The first problem is that the range of scales matching to the C5 feature’s receptive field is limited, which impedes the detection performance for objects across various scales. The second one is the imbalance problem on positive anchors raised by sparse anchors in the single-level feature. Next, we discuss these two problems in detail and provide our solutions.

#### 4.1. Limited Scale Range

Recognizing objects at vastly different scales is a fundamental challenge in object detection. One feasible solution to this challenge is to leverage multiple-level features. In detectors with MiMo or SiMo encoders, they construct multiple-level features with different receptive fields (P3-P7) and detect objects on the level with receptive field matching to their scales. However, the single-level feature setting changes the game. There is only one output feature in SiSo encoders, whose receptive field is a constant. As

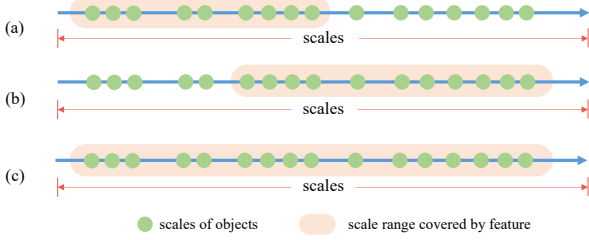


Figure 4. A toy example to illustrate the relation between the object scales and the scale range covered by the single feature. The axis in this figure denotes the scales. (a) indicates that the feature’s receptive field can only cover a limited scale range; (b) shows that the enlarged scale ranges enable the feature to cover large objects while miss covering small ones; (c) represents that all scales can be covered the feature with multiple receptive fields.

shown in Figure 4(a), the C5 feature’s receptive field can only cover a limited scale range, resulting in poor performance if the objects’ scales mismatches with the receptive field. To achieve the goal of detecting all objects with SiSo encoders, we have to find a way to generate an output feature with various receptive fields, compensating for the lack of multiple-level features.

We begin with enlarging the receptive field of the C5 feature by stacking standard and dilated convolutions [45]. Although the covered scale range is enlarged to some extent, it still can not cover all object scales as the enlarging process multiplies a factor greater than 1 to all originally covered scales. We illustrate the situation in Figure 4(b), where the whole scale range shifts to larger scales compare with the one in Figure 4(a). Then, we combine the original scale range and the enlarged scale range by adding the corresponding features, resulting in an output feature with multiple receptive fields covering all object scales (Figure 4(c)). The above operations can be easily achieved by constructing residual blocks [14] with dilations on the middle  $3 \times 3$  convolution layer.

**Dilated Encoder:** Based on the above designs, we propose our SiSo encoder in Figure 5, named as *Dilated Encoder*. It contains two main components: the *Projector* and the *Residual Blocks*. The projection layer first applies one  $1 \times 1$  convolution layer to reduce the channel dimension, then add one  $3 \times 3$  convolution layer to refine semantic contexts, which is the same as in the FPN [22]. After that, we stack four successive dilated residual blocks with different dilation rates in the  $3 \times 3$  convolution layers to generate output features with multiple receptive fields, covering all objects’ scales.

**Discussion:** Dilated convolution [45] is a common strategy to enlarge the features’ receptive field in object detection. As reviewed in the Section 2, TridentNet [20] use dilated convolution to generate multi-scale features. It deals with

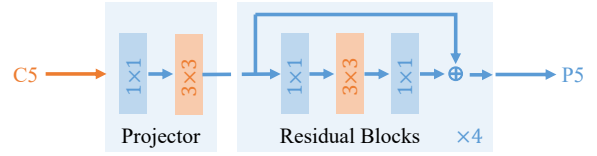


Figure 5. An illustration of the structure of *Dilated Encoder*. In the figure,  $1 \times 1$  and  $3 \times 3$  denotes  $1 \times 1$  and  $3 \times 3$  convolution layers and  $\times 4$  means four successive residual blocks. All convolution layers in Residual Blocks are followed by a batchnorm layer [15] and a ReLU layer [27], while in Projector, we only use convolution layers and batchnorm layers [15].

the scale variation problem in object detection via multi-branch structure and weight sharing mechanism, which is different from our single-level feature setting. Moreover, *Dilated Encoder* stack dilated residual blocks one by one without weight sharing. Although DetNet [21] also successively applies dilated residual blocks, its purpose is to maintain the spatial resolution of the features and keep more details in the backbone’s outputs, while ours is to generate a feature with multiple receptive fields out of the backbone. The design of *Dilated Encoder* enables us to detecting all objects on single-level feature instead of on multiple-level features like TridentNet [20] and DetNet [21].

## 4.2. Imbalance Problem on Positive Anchors

The definition of positive anchors is crucial for the optimization problem in object detection. In anchor-based detectors, strategies to define positive are dominated by measuring the IoUs between anchors and ground-truth boxes. In RetinaNet [23], if the max IoU of the anchor and ground-truth boxes is greater than a threshold 0.5, this anchor will be set as positive. We call it Max-IoU matching.

In MiMo encoders, the anchors are pre-defined on multiple levels in a dense paved fashion, and the ground-truth boxes generate positive anchors in feature levels corresponding to their scales. Given the divide-and-conquer mechanism, Max-IoU matching enables ground-truth boxes in each scale to generate a sufficient number of positive anchors. However, when we adopt the SiSo encoder, the number of anchors diminish extensively compare to the one in the MiMo encoder, from  $100k$  to  $5k$ , resulting in sparse anchors<sup>1</sup>. Sparse anchors raise a matching problem for detectors when applying Max-IoU matching, as shown in Figure 6. Large ground-truth boxes induce more positive anchors than small ground-truth boxes in natural, which cause an imbalance problem for positive anchors. This imbalance makes detectors pay attention to large ground-truth boxes while ignoring the small ones when training.

<sup>1</sup>In SiSo encoders, we simply collapse multiple anchors on multiple-level features to single-level, e.g., we construct 5 anchors with different anchor sizes of {32, 64, 128, 256, 512} on each position of the C5 feature.



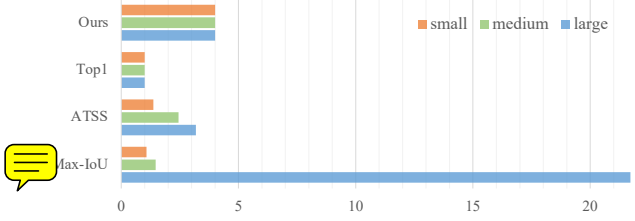


Figure 6. Distribution of the generated positive anchors in various matching methods with single feature. This figure aims to show the balancedness of the generated positive anchors. The positive anchors in the Max-IoU are dominated by large ground-truth boxes, causing huge imbalance across object scales. ATSS alleviates the imbalance problem by adaptively sampling positive anchors when training. The Top1 and Ours adopt a uniform matching, generating positive anchors in a balanced manner regardless of small, medium, and large objects.

**Uniform Matching:** To solve this imbalance problem in positive anchors, we propose an *Uniform Matching* strategy: **adopting the  $k$  nearest anchor as positive anchors for each ground-truth box, which makes sure that all ground-truth boxes can be matched with the same number of positive anchors uniformly regardless of their sizes** (Figure 6). Balance in positive samples makes sure that all ground-truth boxes participate in training and contribute equally. Besides, following Max-IoU matching [23], we set IoU thresholds in Uniform Matching to ignore large IoU ( $>0.7$ ) negative anchors and small IoU ( $<0.15$ ) positive anchors.

**Discussion: relation to other matching methods.** Applying topk in the matching process is not new. ATSS [48] first select topk anchors for each ground-truth box on  $\mathcal{L}$  feature levels, then samples positive anchors among  $k \times \mathcal{L}$  candidates by dynamic IoU thresholds. However, ATSS focuses on defining positives and negatives adaptively, while our uniform matching focuses on achieving **balance on positive samples with sparse anchors**. Although several previous methods achieve balance on positive samples, their matching processes are **not** designed for this imbalance problem. For example, YOLO [29] and YOLOv2 [30] match the ground-truth boxes with the best matching cell or anchor; DETR [4] and [36] apply Hungarian algorithm [18] for matching. These matching methods can be view as top1 matching, which is a specific case of our uniform matching. More importantly, the difference between the uniform matching and the learning-to-match methods is that: the learning-to-match methods, such as FreeAnchor [49] and PAA [16], adaptively separate anchors into positives and negatives according to the learning status, while uniform matching is **fixed** and does not evolve with training. The uniform matching is proposed to address the specific imbalance problem on positive anchors under the SiSo design. The comparison in Figure 6 and the results in Table 5e demonstrate the significance of the balance in positives in

SiSo encoders.

### 4.3. YOLOF

Based on the solutions above, we propose a fast and straightforward framework with single-level feature, denoted as YOLOF. We format YOLOF into three parts: the backbone, the encoder, and the decoder. The sketch of YOLOF is shown in Figure 9. In this section, we give a brief introduction to the main components of YOLOF.

**Backbone.** In all models, we simply adopt the ResNet [14] and ResNeXt [43] series as our backbone. All models are pre-trained on ImageNet. The output of the backbone is the C5 feature map which has 2048 channels and with a downsample rate of 32. To make a fair comparison with other detectors, all batchnorm layers in the backbone are frozen by default.

**Encoder.** For the encoder (Figure 5), we first follow FPN by adding two projection layers (one  $1 \times 1$  and one  $3 \times 3$  convolution) after the backbone, resulting in a feature map with 512 channels. Then, to enable the encoder’s output feature to cover all objects on various scales, we propose to add residual blocks, which consist of three consecutive convolutions: the first  $1 \times 1$  convolution apply channel reduction with a reduction rate of 4, then a  $3 \times 3$  convolution with dilation is used to enlarge the receptive field, at last, a  $1 \times 1$  convolution to recover the number of channels.

**Decoder.** For the decoder, we adopt the main design of RetinaNet, which consists of two parallel task-specific heads: the classification head and the regression head (Figure 9). We only add two minor modifications. The first one is that we follow the design of FFN in DETR [4] and make the number of convolution layers in two heads different. There are four convolutions followed by batch normalization layers and ReLU layers on the regression head while only have two on the classification head. The second is that we follow Autoassign [52] and add an implicit objectness prediction (without direct supervision) for each anchor on the regression head. The final classification scores for all predictions are generated by multiplying the classification output with the corresponding implicit objectness.

**Other Details.** As mentioned in the previous section, the pre-defined anchors in YOLOF **are sparse**, decreasing the match quality between anchors and ground-truth boxes. We add a random shift operation on the image to circumvent this problem. The operation shifts the image randomly with a maximum of 32 pixels in left, right, top, and bottom directions and aims to inject noises into the object’s position in the image, increasing the probability of ground-truth boxes matching with high-quality anchors. Moreover, we found that a restriction on the anchors’ center’s shift is also helpful to the final classification when using a single-level feature.

Model	schedule	$AP$	$AP_{50}$	$AP_{75}$	$AP_S$	$AP_M$	$AP_L$	#params	GFLOPs	FPS
RetinaNet [23]	1x	35.9	55.7	38.5	19.4	39.5	48.2	38M	201	13
RetinaNet-R101 [23]	1x	38.3	58.5	41.3	21.7	42.5	51.2	57M	266	11
RetinaNet+	1x	37.7	58.1	40.2	22.2	41.7	49.9	38M	201	13
RetinaNet-R101+	1x	40.0	60.4	42.7	23.2	44.1	53.3	57M	266	10
YOLOF	1x	37.7	56.9	40.6	19.1	42.5	53.2	44M	86	32
YOLOF-R101	1x	39.8	59.4	42.9	20.5	44.5	54.9	63M	151	21
YOLOF-X101	1x	42.2	62.1	45.7	23.2	47.0	57.7	102M	289	10
YOLOF-X101 <sup>†</sup>	3x	44.7	64.1	48.6	25.1	49.2	<b>60.9</b>	102M	289	10
YOLOF-X101 <sup>†‡</sup>	3x	<b>47.1</b>	<b>66.4</b>	<b>51.2</b>	<b>31.8</b>	<b>50.9</b>	60.6	102M	-	-

Table 1. Comparison with RetinaNet on the COCO2017 validation set. The top section shows the results of RetinaNet. The middle section gives the results of an improved RetinaNet (with a “+”), which is RetinaNet with GIoU [32], GN [41], and implicit objectness. The last section shows the results of various YOLOF models. In the table, the model with a suffix of R101 or X101 means it use ResNet-101 [14] or ResNeXt-101-64×4d [43] as backbone. For those not marked with suffix, they adopt ResNet-50 [14] by default. In the last two rows, we use multi-scale training and testing techniques († indicates multi-scale training and ‡ means multi-scale testing), whose settings follow HTC [5]. More details about the settings can be found in the Appendix. In the last three columns, we show models’ number of parameters (#params), GFLOPs, and inference speed. All FLOPs are measured with a shorter edge size 800 over the first 100 images of COCO val2017. Moreover, the FPS in the table is calculated with batch size 1 on 2080Ti from the total inference pure compute time reported in the Detectron2 [42].

We add a restriction that the centers’ shift for all anchors should smaller than 32 pixels.

## 5. Experiments

We evaluate our YOLOF on the MS COCO [24] benchmark and conduct comparisons with RetinaNet [23] and DETR [4]. Then, we provide a detailed ablation study of each component’s design with quantitative results and analysis. Finally, to give insights to further research on single-level detection, we provide error analysis and show the weaknesses of YOLOF compared with DETR [4]. The details are as follows.

**Implementation Details.** YOLOF is trained with synchronized SGD over 8 GPUs with a total of 64 images per mini-batch (8 images per GPU). All models are trained with an initial learning rate of 0.12. Moreover, following DETR [4], we set a smaller learning rate for the backbone, which is 1/3 of the base learning rate. To stabilize the training at the beginning, we extend the number of warmup iterations from 500 to 1500. For training schedules, as we increase the batch size, the ‘1×’ schedule setting in YOLOF is a total of 22.5k iterations and with base learning rate decreased by 10 in the 15k and the 20k iteration. Other schedules are adjusted according to the principles in Detectron2 [42]. For model inference, we employ NMS with a threshold of 0.6 to post-process the results. For other hyperparameters, we follow the settings of RetinaNet [23].

### 5.1. Comparison with previous works

**Comparison with RetinaNet:** To make a fair comparison, we align RetinaNet with YOLOF by employing generalized

IoU [32] for the box loss, adding an implicit objectness prediction, and applying group normalization layers [41] in heads (as there are only two images per GPU and both BN [15] and SyncBN [46] give poor results in RetinaNet<sup>2</sup>, we use GN [41] instead of BN [15] in the heads). The results are presented in Table 1. All ‘1×’ models are trained with a single scale that the shorter side is set as 800 pixels and the longer side is at most 1333 [23]. In the top section, we give RetinaNet baseline results trained with Detectron2 [42]. In the middle section, we present the results of the improved RetinaNet baseline (with a “+”), whose settings are aligned with YOLOF. In the last section, we show results from multiple YOLOF models. Thanks to the single-level feature, YOLOF achieves results *on par with* RetinaNet+ with a 57% flops reduction (flops for each component in YOLOF are shown in Figure 3) and a 2.5× speed up. Due to the large stride (32) of the C5 feature, YOLOF has an inferior performance (−3.1) than RetinaNet+ on small objects. However, YOLOF achieves better performance on large objects (+3.3) as we add dilated residual blocks in the encoder. The comparison between RetinaNet+ and YOLOF with a ResNet-101 [14] show similar evidence as well. Although YOLOF is inferior to RetinaNet+ on small objects when applying the same backbone, it can match small objects’ performance with a stronger backbone ResNeXt [43] while running at the same speed. Moreover, to prove that our method is compatible and complementary to current technologies in object detec-

<sup>2</sup>[https://github.com/facebookresearch/detectron2/blob/master/detectron2/modeling/meta\\_arch/retinanet.py#L532](https://github.com/facebookresearch/detectron2/blob/master/detectron2/modeling/meta_arch/retinanet.py#L532)

Model	Epochs	#params	GFLOPS/FPS	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
DETR [4]	500	41M	86/24*	42.0	62.4	44.2	20.5	45.8	61.1
DETR-R101 [4]	500	60M	152/17*	43.5	<b>63.8</b>	46.4	21.9	48.0	<b>61.8</b>
YOLOF	72	44M	86/32	41.6	60.5	45.0	22.4	46.2	57.6
YOLOF-R101	72	63M	151/21	<b>43.7</b>	62.7	<b>47.4</b>	<b>24.3</b>	<b>48.3</b>	58.9

Table 2. Comparison with DETR on the COCO2017 validation set. We conduct comparisons with backbone ResNet-50 (without suffix) and ResNet-101 (with a suffix R101). To make fair comparison, YOLOF adopts multi-scale training (same as in Table 1) with a '6×' schedule, which is roughly 72 epochs. For the FPS of DETR, \* means we follow the method in the original paper [4] and re-measure it on 2080Ti.

Model	Epochs	FPS	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
YOLOv4 [1]	273	53*	43.5	<b>65.7</b>	47.3	<b>26.7</b>	47.6	53.3
YOLOF-DC5	184	<b>60†</b>	<b>44.3</b>	62.9	<b>47.5</b>	24.0	<b>48.5</b>	<b>60.4</b>

Table 3. Comparison with YOLOv4 on the COCO *test-dev* set. We train YOLOF-DC5 with a '15×' schedule (184 epochs) and compare it with YOLOv4. In the table, † means that the FPS for YOLOF-DC5 is measured by following YOLOv4 [1]. It is different from the method used in Table 1, 2 in this paper. In YOLOv4 [1], the authors fuse the convolution layer and the batch normalization layer, then measure the inference time after converting the model to half-precision. \* represents that we get the speed for YOLOv4 on 2080Ti from the official repo <https://github.com/AlexeyAB/darknet#geforce-rtx-2080-ti>.

tion, we show results that training with multi-scale images and a longer schedule in the last two rows of Table 1. Finally, with the help of multi-scale testing, we obtain our final result of 47.1 mAP and a competitive performance of 31.8 mAP on small objects.

**Comparison with DETR.** DETR [4] is a recent proposed detector which introduces transformer [39] to object detection. It achieves surprising results on the COCO benchmark [24] and proves that by only adopting a single C5 feature, it can achieve comparable results with a multi-level feature detector (Faster R-CNN w/ FPN [22]) for the first time. Given this, one might expect that layers capture global dependencies such as transformer layers [39] are required to achieve promising results in single-level feature detection. *However, we show that a conventional network with local convolution layers can also achieve this goal.* We compare DETR with global layers and YOLOF with local convolution layers in Table 2. The results show that YOLOF matches the DETR’s performance, and YOLOF gets more benefits from deeper networks than DETR (w/ ResNet-50 (−0.4) vs. w/ ResNet-101 (+0.2)). Interestingly, we find that YOLOF outperforms DETR on small objects (+1.9 and +2.4) while lags behind DETR on large objects (−3.5 and −2.9). The finding is consistent with the local and global discussion above. More importantly, compared with DETR, YOLOF converge much faster ( $\sim 7\times$ ), making it more suitable than DETR to serve as a simple baseline for single-level detectors.

<i>Dilated Encoder</i>	<i>Uniform Matching</i>	AP	$\Delta$	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
		21.1	-16.6	8.6	31.1	34.5
✓		29.1	-8.6	9.5	32.2	50.6
	✓	33.8	-3.9	17.7	40.9	43.8
✓	✓	<b>37.7</b>	-	<b>19.1</b>	<b>42.5</b>	<b>53.2</b>

Table 4. Effect of *Dilated Encoder* and *Uniform Matching* with ResNet-50. These two components improve the original single-level detector by 16.6 mAP. Note that the result of 21.1 mAP in the table is not a bug. It perform slightly worse than the detectors with SiSo encoders in Figure 1 and Figure 3 due to the design of the decoder in YOLOF - only two convolution layers in the classification head.

**Comparison with YOLOv4.** YOLOv4 [1] is an optimal speed and accuracy multi-level feature detector. It combines many tricks to achieve state-of-the-art results. As our purpose is to build a simple and fast baseline for single-level detectors, investigation on the bag of freebie tricks is outside of the scope of this work. Thus, we do not expect a rigidly aligned comparison on performance. To compare our YOLOF with YOLOv4, we apply the data augmentation methods as YOLOv4, adopt a three-phase training pipeline, modify the training settings accordingly, and add dilations on the last stage of the backbone (YOLOF-DC5 in Table 3). More technical details about the model and the training settings are given in the Appendix. As shown in Table 3, YOLOF-DC5 can run 13% faster than YOLOv4 with a 0.8 mAP improvement on overall performance. YOLOF-DC5 achieves less competitive results on small objects than YOLOv4 (24.0 mAP vs. 26.7 mAP) while outperforms it on large objects by a large margin (+7.1 mAP). The above results indicate that single-level detectors have great potential to achieve state-of-the-art speed and accuracy simultaneously.

## 5.2. Ablation Experiments

We run a number of ablations to analyze YOLOF. We first provide an overall analysis of the two proposed components. Then, we show the ablation experiments on detailed designs of each component. Results are shown in Table 4, 5 and discussed in detail next.

$N$	AP	AP <sub>s</sub>	AP <sub>m</sub>	AP <sub>l</sub>	Dilations	AP	AP <sub>s</sub>	AP <sub>m</sub>	AP <sub>l</sub>	Dilations & Shortcut	AP	AP <sub>s</sub>	AP <sub>m</sub>	AP <sub>l</sub>
0	33.8	17.7	40.9	43.8	1,1,1,1	35.5	17.6	41.4	48.4	2,4,6,8 ✓	<b>37.7</b>	<b>19.1</b>	<b>42.5</b>	<b>53.2</b>
2	34.9	17.8	41.3	46.8	2,2,2,2	36.4	18.1	41.8	50.2	2,4,6,8 -	34.1	16.2	38.4	47.5
4	<b>35.5</b>	<b>17.6</b>	<b>41.4</b>	<b>48.4</b>	3,3,3,3	36.9	18.4	42.1	51.0	1,1,1,1 ✓	35.5	17.6	41.4	48.4
6	36.0	17.7	41.9	49.5	1,2,3,4	37.4	18.6	42.6	51.8	1,1,1,1 -	32.6	15.0	38.4	44.2
8	36.6	18.5	42.0	50.7	<b>2,4,6,8</b>	<b>37.7</b>	<b>19.1</b>	<b>42.5</b>	<b>53.2</b>					
10	36.9	18.3	42.4	50.4	3,6,9,12	37.3	18.7	42.1	52.6					

(a) **Number of ResBlocks** (ResNet-50): More residual blocks bring more gains.  $N$  represent the number of ResBlocks. To keep YOLOF simple and neat, we add 4 blocks in the encoder by default.

(b) **Different dilations** (ResNet-50-N4): 'N4' means we add 4 ResBlocks in the encoder. Dilatation in the residual block gives large gains on large objects and slightly improve the performance of small and medium objects.

(c) **Add shortcut or not** (ResNet-50): YOLOF results with shortcuts or not on various dilation settings. Shortcut brings considerable gains on all object scales and becomes more important when the dilations are adopted (+3.6 AP with dilations 2,4,6,8 vs. +2.9 AP when dilations are all ones).

$topk$	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>s</sub>	AP <sub>m</sub>	AP <sub>l</sub>	Matching Methods	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>s</sub>	AP <sub>m</sub>	AP <sub>l</sub>
top1	35.9	55.6	38.4	17.5	40.3	50.2	Max-IoU Matching [23]	29.1	45.9	29.6	9.5	32.2	50.6
top2	37.2	56.7	39.9	18.9	41.6	52.0	ATSS(topk=9) [48]	34.6	54.3	37.1	17.7	40.6	46.9
top3	37.5	<b>57.1</b>	40.2	18.6	41.9	52.5	ATSS(topk=15) [48]*	36.5	55.9	38.6	18.1	41.4	50.8
<b>top4</b>	<b>37.7</b>	56.9	<b>40.6</b>	<b>19.1</b>	<b>42.5</b>	<b>53.2</b>	Hungarian Matching [4]	35.8	55.5	38.3	18.2	39.9	50.2
top5	37.5	56.7	40.3	18.1	42	<b>53.2</b>	<b>Uniform Matching</b>	<b>37.7</b>	<b>56.9</b>	<b>40.6</b>	<b>19.1</b>	<b>42.5</b>	<b>53.2</b>

(d) **Number of positives** (ResNet-50-N4): Number of positive anchors in *Uniform Matching*. Increase the positive anchor for each ground-truth box can improve the performance while it saturates when too many positive anchors. We choose the top4 anchors in YOLOF which achieves best results.

(e) **Uniform matching vs. other matchings** (ResNet-50-N4): Comparison with other matching methods. Uniform Matching achieve balance in positive anchors and get the best results among other matching methods, which is consistent with the comparison in Figure 6. Note that '\*' represents that we get the best result for ATSS [48] when setting topk as 15. More details can be found in the Appendix.

Table 5. **Ablations.** We show ablation experiments for *Dilation Encoder* and *Uniform Matching* on COCO2017 val set with ResNet-50.

**Dilated Encoder and Uniform Matching:** Table 4 shows that both *Dilated Encoder* and *Uniform Matching* are necessary to YOLOF and bring considerable improvements. Specifically, Dilated Encoder has a significant impact on large objects (43.8 vs. 53.2) and slightly improves the results of small and medium objects. The results indicate that the *limited scale range* is a severe problem in the C5 feature (Section 4.1). Our Dilated Encoder provides a simple but effective solution to this problem. On the other side, the performance of small and medium objects drops significantly ( $\sim 10AP$ ) without uniform matching, while the large objects' performance is only lightly affected. The finding is consistent with the *imbalance problem on positive anchors* analyzed in Section 4.2. The positive anchors are dominated by large objects, resulting in poor results on small and medium objects. Finally, when we remove both Dilated Encoder and Uniform Matching, a single-level feature detector's performance drops back to  $\sim 20$  mAP like the results in Figure 1 and Figure 3.

**Number of ResBlock:** YOLOF stacks residual blocks in the SiSo encoder. The results in Table 5a shows that stacking more blocks gives extensive improvements on large objects, which is due to the increment of the feature scale range. Although we observe continuous improvements with more blocks, we choose to add four residual blocks to keep YOLOF simple and neat.

**Different dilations:** Following the analysis in Section 4.1, to enable the C5 feature to cover large scales, we replace

the standard  $3 \times 3$  convolution layer in the residual blocks with its dilated counterpart. We show the results with different dilations in the residual blocks in Table 5b. Applying dilations to residual blocks bring improvements to YOLOF, while the improvements are saturated when using too large dilations. We conjecture that the reason for this phenomenon is that dilations of 2, 4, 6, 8 are enough to match object scales in all images.

**Add shortcut or not:** Table 5c shows that shortcuts play an essential role in Dilated Encoder. The performance of all objects will drop significantly if we remove the shortcuts in residual blocks. According to Section 4.1, shortcuts combine different scale ranges. A largely and densely paved scale range covered by the feature is the critical factor for detecting all objects in a single-level feature manner.

**Number of positives:** A comparison among the number of induced positive anchors by ground-truth boxes is conducted in Table 5d. Intuitively, more positive anchors can achieve better performance as the learning will be easier when given more samples. Thus, in our uniform matching manner, we empirically increase the number of positive anchors induced by each ground-truth box. As shown in Table 5d, the hyper-parameter  $k$  is very robust for the performance when  $k$  is larger than 1, which may suggest that the most important is the uniform matching manner in YOLOF. We set *top4* for our uniform matching as it is the best choice according to the results.



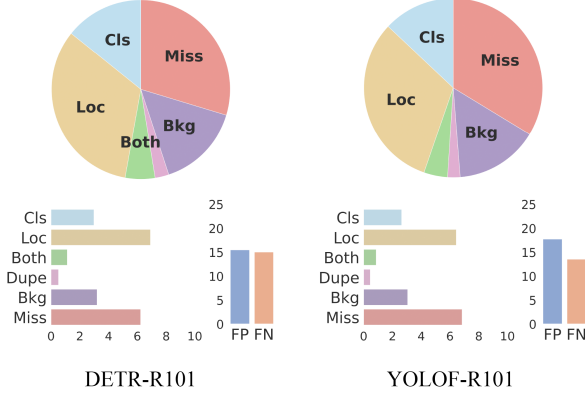


Figure 7. Error analysis for DETR-R101 and YOLOF-R101. According to TIDE [2], the figure shows the six types of errors (Cls: classification error; Loc: localization error; Both: both cls and loc error; Dupe: duplicate predictions error; Bkg: background error; Miss: missing error). The pie chart shows the relative contribution of each error, while the bar plots show their absolute contribution. FP and FN means false positive and false negative respectively.

**Uniform matching vs. other matchings:** We compare the uniform matching with other matching strategies for YOLOF and show results in Table 5e. The proposed uniform matching strategy can achieve the best results, compatible with the imbalance analysis in Figure 6. It worth noting that the Hungarian matching strategy can be roughly treated as Top1 matching (Table 5d) so that they get similar performance. The difference between them is that an anchor will only match one object in Hungarian matching while the Top1 matching does not have this constraint, and the experiments show that this is not important. The original ATSS find that top9 anchors are the best choice, while we find top15 anchors are much better in the single-level feature detector. By using top15 anchors, ATSS achieves a good result of 36.5 mAP while still lags behind our uniform matching by a 1.2 mAP gap.

### 5.3. Error Analysis

We add error analysis for YOLOF in this section to provide insights for future research in single-level feature detection. We adopt the recent proposed tool TIDE [2] to compare YOLOF with DETR [4]. As illustrated in Figure 7, DETR has a larger error in localization than YOLOF, which may be related to its regression mechanism. DETR regresses objects in a total anchor free manner and predicts the location globally in the image, which causes difficulties in localization. In contrast, YOLOF relies on pre-defined anchors, which is responsible for higher missing error than DETR [4] in the predictions. According to the analysis in Section 4.2, the anchors of YOLOF are sparsely and not flexible enough in the inference stage. Intuitively, there are situations that there are no high-quality anchors pre-defined around a ground-truth box. Thus, introducing the anchor-

free mechanism into YOLOF may help alleviate this problem, and we leave it for future work.

## 6. Conclusion

In this work, we identify that the success of FPN is due to its divide-and-conquer solution to the optimization problem in dense object detection. Given that FPN makes network structure complex, brings memory burdens, and slows down the detectors, we propose a simple but highly efficient method without using FPN to address the optimization problem differently, denoted as YOLOF. We prove its efficacy by making fair comparisons with RetinaNet and DETR. We hope our YOLOF can serve as a solid baseline and provide insight for designing single-level feature detectors in future research.

## Appendix A: More Details

**Detailed Structures of All Encoders:** In Figure 8, we illustrate the detailed processes of generating outputs in encoders. The four encoders differ in the number of input features and output features. (a) The *Multiple-in-Multiple-out* (MiMo) encoder receives three levels from the backbone and output five levels. The structure of the MiMo encoder is the same as FPN in RetinaNet [23]. (b) *Single-in-Multiple-out* (SiMo) only has one C5 feature for the input. As there are no other inputs, we remove the  $1 \times 1$  convolution layer designed for C3 and C4. (c) *Multiple-in-Single-out* (MiSo) receives three input features while only generate one output feature P5. To fully utilize the context in the input features, we adopt a structure similar to PANet [25] in MiSo. (d) In the *Single-in-Single-out* (SiSo) encoder, we remove all other convolution layers and only keep the convolution layers in the level of C5.

**Network Architecture of YOLOF** In Figure 9, we show a detailed network architecture of YOLOF. YOLOF detects objects on single-level feature, which is very simple. Our method consists of three components: the backbone, the encoder, and the decoder. The detailed design of these components are presented in Section 4.3.

**Training Time & Memory:** In this section, we compare training time and training memory among YOLOF, DETR [4], and RetinaNet+ [23]. As shown in Table 6, due to the long training schedule, DETR needs 112.5 hours to converge on COCO with eight 2080Ti GPUs, while YOLOF and RetinaNet+ only need 4.5 hours and 9.8 hours, respectively. As for training memory, YOLOF needs less memory than RetinaNet+ and DETR, which make YOLOF be trained with larger batch size and converge faster.

**More Implementation Details:** The default training settings for YOLOF is a total of 64 images per mini-batch (8 images per GPU) with an initial learning rate of 0.12.

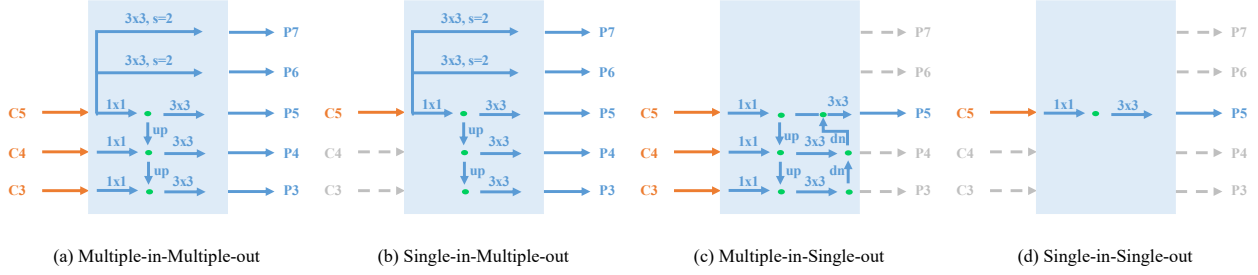


Figure 8. Detailed Structures of Multiple-in-Multiple-out (MiMo), Single-in-Multiple-out (SiMo), Multiple-in-Single-out (MiSo), and Single-in-Single-out (SiSo) encoders.

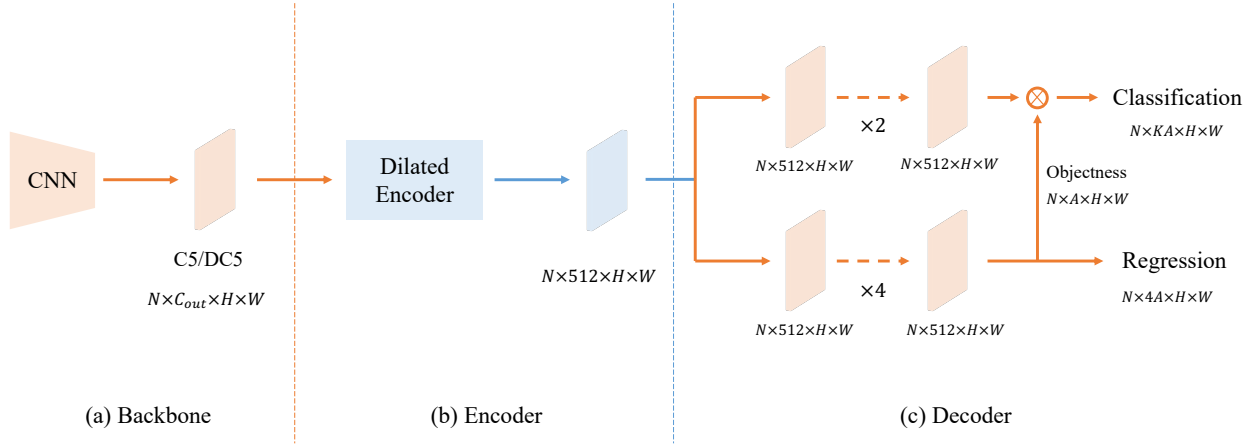


Figure 9. The sketch of YOLOF, which consists of three main components: the backbone, the encoder, and the decoder. In the figure, 'C5/DC5' represents the output feature of the backbone with downsample rate of 32/16. ' $C_{out}$ ' means the number of channels of the feature. We set the number of channels as 512 for feature maps in the encoder and the decoder.  $H \times W$  is the height and width of feature maps.

Model	Memory/Images	Training Time
YOLOF	5.3G / 8	4.5h
RetinaNet+ [23]	4.9G / 2	9.8h
DETR [4]	7.1G / 2	112.5h

Table 6. Comparison of training memory and training time among different models. All models are trained with eight 2080Ti GPUs with their default settings, i.e., we train YOLOF and RetinaNet+ [23] in a '1x' schedule, while train DETR [4] with 150 epochs on COCO2017 training set.

While for ResNeXt-101 [43], we train with 4 images per GPU (batch size 32) and set the learning rate to 0.06 following the linear rule [12]. For multi-scale training, DETR [4] apply random crop plus resize to simulate large image size during training. In YOLOF, we simply resize the image to large size. For multi-scale training, we follow HTC [5] and adopt a strategy of random sample the image size between [400, 1400] with its largest edge no greater than 1600 pixels.

**Detailed Settings to Compare with YOLOv4** To match the performance of YOLOv4, we first increase the number of dilated residual blocks in the dilated encoder from

4 to 8. We adjust the dilations of these dilated residual blocks according to experimental results. We find that the dilations [1, 2, 3, 4, 5, 6, 7, 8] give the best result. Then following YOLOv4 [1], we adopt its data augmentations, take the CSPDarkNet-53 [40] as the backbone, replace all the batch normalization layers with its synchronized counterpart, and apply LeakyReLU [44] in the encoder and the decoder instead of ReLU layers. According to the results in Table 9, YOLOF-DC5 gives better results than YOLOF. Thus we use YOLOF-DC5 as the baseline model in this section. After that, we set an initial learning rate of 0.04 for the whole model. To train the final model, we adopt a three-phase training. At first, we training YOLOF-DC5 for a '9x' schedule; then we increase the ignore threshold for negative anchors from 0.75 to 0.8 and train a '3x' schedule based on the previous model (this phase gives a 0.5 mAP gain); at last, we train another '3x' schedule by following the recipe introduced in [47]. The final result shown in Table 3 is produced by the SWA model, which is obtained by averaging 12 checkpoints (the SWA model gives a  $\sim 1$  mAP improvement).

Model	areas	sizes	ratios	$AP$	$AP_{50}$	$AP_{75}$	$AP_s$	$AP_m$	$AP_l$
YOLOF	5	1	1	<b>37.7</b>	56.9	40.6	19.1	42.5	53.2
YOLOF	5	1	3	<b>37.7</b>	57.2	40.7	19.4	42.0	52.2
YOLOF	5	3	1	35.0	52.3	37.9	15.0	40.6	52.9
YOLOF	5	3	3	35.4	52.4	38.3	14.8	41.2	52.5

Table 7. Results of YOLOF with different multiple anchors per location on COCO [24] validation set.

Model & k	5	7	9	11	13	15	17	19
YOLOF (with ATSS)	33.7	33.8	34.6	35.8	35.5	<b>36.5</b>	36.3	36.2

Table 8. An illustration of how performance changes with the variation of the hyper-parameter  $k$  in ATSS [48].

Model	FPS	$AP$	$AP_{50}$	$AP_{75}$	$AP_S$	$AP_M$	$AP_L$
YOLOF-DC5-R50	24	39.2	58.6	42.7	22.3	43.9	50.8
YOLOF-DC5-R101*	17	40.5	59.8	43.9	23.0	44.9	53.8

Table 9. Additional results of YOLOF-DC5 with different backbones on COCO *val* split. \* means that due to the limited memory of 2080Ti, we train with 4 images per GPU (batch size 32) for ResNet-101. Higher performance can be achieved if train with 8 images per GPU or apply SyncBN (BN layers in the encoder and decoder restrict the improvements).

## Appendix B: Additional Experimental Results

**Number of Anchors:** In RetinaNet [23], anchors are generated from multiple level features (P3-P7) with areas of  $32^2$  to  $512^2$ , respectively. At each level feature, RetinaNet paves anchors with sizes  $\{2^0, 2^{1/3}, 2^{2/3}\}$  and aspect ratios  $\{0.5, 1, 2\}$ . While in YOLOF, we only have a one-level feature to place anchors. To cover all objects’ scales, we add anchors with areas of  $\{32^2, 64^2, 128^2, 256^2, 512^2\}$ , size  $\{1\}$ , and aspect ratio  $\{1\}$  in the single feature map, resulting in 5 anchors in each position. Moreover, we investigate the influence of more anchors in YOLOF. Following RetinaNet, we generate 45 anchors in each position with different sizes ( $\{2^0, 2^{1/3}, 2^{2/3}\}$ ) and more aspect ratios ( $\{0.5, 1, 2\}$ ). All results are shown in Table 7. The results show that adding more aspect ratios does not change the performance of YOLOF, while the performance drops with more sizes. Thus, we choose to add a minimum of five anchors for YOLOF by default.

**Hyper-parameter of ATSS:** Here, we provide the results of using different values of  $k$  in ATSS [48] in Table 8. The results show that the choice of  $k = 9$  used in the original paper is not the best choice in YOLOF. According to the results, we choose  $k = 15$  for ATSS in this paper.

**Results with Dilated C5:** In this paper, we show that YOLOF performs well on the C5 feature. To boost the performance of YOLOF, we detect objects on a feature map with higher resolution than the C5 feature. Following DETR [4], we construct a backbone with dilation and without stride on its last stage. The backbone’s output feature is denoted as DC5, with a downsample rate of 16. In

Table 9, we show the results of YOLOF-DC5 on COCO *val* split with ResNet-50 and ResNet-101 as the backbone. YOLOF-DC5 achieves higher performance than the original YOLOF but runs at a slower speed as the feature’s resolution is larger than C5. To achieve the results, we first add a smaller anchor, resulting in 6 anchors per location ( $\{16, 32, 64, 128, 256, 512\}$ ), then we increase the *topk* from 4 to 8 and change the ignore threshold for positive anchors from 0.15 to 0.1. Other parameters are the same as before.

## Acknowledgements

This work is supported by The National Key Research and Development Program of China (No. 2017YFA0700800), Beijing Academy of Artificial Intelligence (BAAI), National Natural Science Foundation of China (No.61972396, 61876182, 61906193), National Key Research and Development Program of China (No. 2020AAA0103402), the Strategic Priority Research Program of Chinese Academy of Sciences (No. XDB32050200), and The NSFC-General Technology Collaborative Fund for Basic Research (Grant No. U1936204).

## References

- [1] Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao. Yolo4: Optimal speed and accuracy of object detection. *arXiv preprint arXiv:2004.10934*, 2020. **2, 7, 10**
- [2] Daniel Bolya, Sean Foley, James Hays, and Judy Hoffman. Tide: A general toolbox for identifying object detection errors. In *ECCV*, 2020. **9**
- [3] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delying into high quality object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6154–6162, 2018. **1**
- [4] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. *arXiv preprint arXiv:2005.12872*, 2020. **2, 5, 6, 7, 8, 9, 10, 11**
- [5] Kai Chen, Jiangmiao Pang, Jiaqi Wang, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jianping Shi, Wanli Ouyang, et al. Hybrid task cascade for instance segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4974–4983, 2019. **6, 10**
- [6] Jifeng Dai, Yi Li, Kaiming He, and Jian Sun. R-fcn: Object detection via region-based fully convolutional networks. In *Advances in neural information processing systems*, pages 379–387, 2016. **2**
- [7] Kaiwen Duan, Song Bai, Lingxi Xie, Honggang Qi, Qingming Huang, and Qi Tian. Centernet: Keypoint triplets for object detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 6569–6578, 2019. **2**

- [8] Pedro F Felzenszwalb, Ross B Girshick, David McAllester, and Deva Ramanan. Object detection with discriminatively trained part-based models. *IEEE transactions on pattern analysis and machine intelligence*, 32(9):1627–1645, 2009. 2
- [9] Golnaz Ghiasi, Tsung-Yi Lin, and Quoc V Le. Nas-fpn: Learning scalable feature pyramid architecture for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7036–7045, 2019. 1, 2
- [10] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448, 2015. 2
- [11] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 580–587, 2014. 2
- [12] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyröla, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large mini-batch sgd: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017. 10
- [13] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 1
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 1, 3, 4, 5, 6
- [15] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015. 4, 6
- [16] Kang Kim and Hee Seok Lee. Probabilistic anchor assignment with iou prediction for object detection. In *ECCV*, 2020. 5
- [17] Tao Kong, Fuchun Sun, Chuanqi Tan, Huaping Liu, and Wenbing Huang. Deep feature pyramid reconfiguration for object detection. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 169–185, 2018. 1, 2
- [18] Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955. 5
- [19] Hei Law and Jia Deng. Cornernet: Detecting objects as paired keypoints. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 734–750, 2018. 2
- [20] Yanghao Li, Yuntao Chen, Naiyan Wang, and Zhaoxiang Zhang. Scale-aware trident networks for object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 6054–6063, 2019. 2, 4
- [21] Zeming Li, Chao Peng, Gang Yu, Xiangyu Zhang, Yangdong Deng, and Jian Sun. Detnet: Design backbone for object detection. In *Proceedings of the European conference on computer vision (ECCV)*, pages 334–350, 2018. 4
- [22] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2117–2125, 2017. 1, 2, 4, 7
- [23] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017. 1, 2, 3, 4, 5, 6, 8, 9, 10, 11
- [24] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 6, 7, 11
- [25] Shu Liu, Lu Qi, Haifang Qin, Jianping Shi, and Jiaya Jia. Path aggregation network for instance segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8759–8768, 2018. 1, 2, 9
- [26] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *European conference on computer vision*, pages 21–37. Springer, 2016. 2
- [27] Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In *ICML*, 2010. 4
- [28] Jiangmiao Pang, Kai Chen, Jianping Shi, Huajun Feng, Wanli Ouyang, and Dahua Lin. Libra r-cnn: Towards balanced learning for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 821–830, 2019. 1
- [29] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016. 2, 5
- [30] Joseph Redmon and Ali Farhadi. Yolo9000: better, faster, stronger. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7263–7271, 2017. 2, 5
- [31] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015. 2
- [32] Hamid Rezaatoughi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian Reid, and Silvio Savarese. Generalized intersection over union: A metric and a loss for bounding box regression. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 658–666, 2019. 6
- [33] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 2
- [34] Bharat Singh and Larry S Davis. An analysis of scale invariance in object detection snip. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3578–3587, 2018. 2
- [35] Bharat Singh, Mahyar Najibi, and Larry S Davis. Sniper: Efficient multi-scale training. In *Advances in neural information processing systems*, pages 9310–9320, 2018. 2
- [36] Russell Stewart, Mykhaylo Andriluka, and Andrew Y Ng. End-to-end people detection in crowded scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2325–2333, 2016. 5



- [37] Mingxing Tan, Ruoming Pang, and Quoc V Le. Efficientdet: Scalable and efficient object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10781–10790, 2020. 1, 2
- [38] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. Fcos: Fully convolutional one-stage object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 9627–9636, 2019. 1, 2
- [39] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. 2, 7
- [40] Chien-Yao Wang, Hong-Yuan Mark Liao, Yueh-Hua Wu, Ping-Yang Chen, Jun-Wei Hsieh, and I-Hau Yeh. Cspnet: A new backbone that can enhance learning capability of cnn. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 390–391, 2020. 10
- [41] Yuxin Wu and Kaiming He. Group normalization. In *Proceedings of the European conference on computer vision (ECCV)*, pages 3–19, 2018. 6
- [42] Yuxin Wu, Alexander Kirillov, Francisco Massa, Wan-Yen Lo, and Ross Girshick. Detectron2. <https://github.com/facebookresearch/detectron2>, 2019. 3, 6
- [43] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1492–1500, 2017. 5, 6, 10
- [44] Bing Xu, Naiyan Wang, Tianqi Chen, and Mu Li. Empirical evaluation of rectified activations in convolutional network. *arXiv preprint arXiv:1505.00853*, 2015. 10
- [45] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*, 2015. 4
- [46] Hang Zhang, Kristin Dana, Jianping Shi, Zhongyue Zhang, Xiaogang Wang, Amrbrish Tyagi, and Amit Agrawal. Context encoding for semantic segmentation. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 7151–7160, 2018. 6
- [47] Haoyang Zhang, Ying Wang, Feras Dayoub, and Niko Sünderhauf. Swa object detection. *arXiv preprint arXiv:2012.12645*, 2020. 2, 10
- [48] Shifeng Zhang, Cheng Chi, Yongqiang Yao, Zhen Lei, and Stan Z Li. Bridging the gap between anchor-based and anchor-free detection via adaptive training sample selection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9759–9768, 2020. 2, 5, 8, 11
- [49] Xiaosong Zhang, Fang Wan, Chang Liu, Rongrong Ji, and Qixiang Ye. FreeAnchor: Learning to match anchors for visual object detection. In *Neural Information Processing Systems*, 2019. 5
- [50] Zhenli Zhang, Xiangyu Zhang, Chao Peng, Xiangyang Xue, and Jian Sun. Exfuse: Enhancing feature fusion for semantic segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 269–284, 2018. 2
- [51] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. *arXiv preprint arXiv:1904.07850*, 2019. 2
- [52] 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. 5