# YOLO-MS: Rethinking Multi-Scale Representation Learning for Real-time Object Detection

**Yuming Chen**[1], **Xinbin Yuan**[1], **Ruiqi Wu**[1], **Jiabao Wang**[1], **Qibin Hou**[1], **Ming-Ming Cheng**[1]
[1]VCIP, School of Computer Science, Nankai University.

https://github.com/FishAndWasabi/YOLO-MS

## Abstract

We aim at providing the object detection community with an efficient and performant object detector, termed YOLO-MS. The core design is based on a series of investigations on how convolutions with different kernel sizes affect the detection performance of objects at different scales. The outcome is a new strategy that can strongly enhance multi-scale feature representations of real-time object detectors. To verify the effectiveness of our strategy, we build a network architecture, termed YOLO-MS. We train our YOLO-MS on the MS COCO dataset from scratch without relying on any other large-scale datasets, like ImageNet, or pre-trained weights. Without bells and whistles, our YOLO-MS outperforms the recent state-of-the-art real-time object detectors, including YOLO-v7 and RTMDet, when using a comparable number of parameters and FLOPs. Taking the XS version of YOLO-MS as an example, with only 4.5M learnable parameters and 8.7G FLOPs, it can achieve an AP score of 43%+ on MS COCO, which is about 2%+ higher than RTMDet with the same model size. Moreover, our work can also be used as a plug-and-play module for other YOLO models. Typically, our method significantly improves the AP of YOLOv8 from 37%+ to 40%+ with even fewer parameters and FLOPs.

## 1 Introduction

Real-time object detection, exemplified by the YOLO series [43–45, 1, 50, 22, 13, 25, 51, 39, 56, 23, 38], has found significant applications in industries, particularly for edge devices, such as drones and robotics. Different from previous heavy object detectors [62, 58, 46, 3, 49, 7], real-time object detectors aim to pursue an optimal trade-off between speed and accuracy. In pursuit of this objective, a mountain of work has been proposed: from the first generation of DarkNet [45] to CSPNet [52], and then to the recent extended ELAN [51], the architectures of real-time object detectors have experienced great changes accompanied by the fast growth of performance.

Despite the impressive performance, recognizing objects at different scales remains a fundamental challenge for real-time object detectors. This motivates us to design a robust encoder architecture for learning expressive multi-scale feature representations. Specifically, we consider encoding multi-scale features for real-time object detection from two new perspectives.

- From a local view, we design a MS-Block with a simple yet effective hierarchical feature fusion strategy [12]. Inspired by Res2Net [12], we introduce multiple branches in MS-Block to perform feature extraction, but differently, we use an inverted bottleneck block with depth-wise convolution to enable the efficient use of large kernels.

- From a global view, we propose gradually increasing the kernel size of convolutions as the network goes deeper. We use small kernel convolutions in shallow layers to process the high-resolution features more efficiently. On the other hand, we adopt large kernel convolutions in deep layers to capture wide-range information.
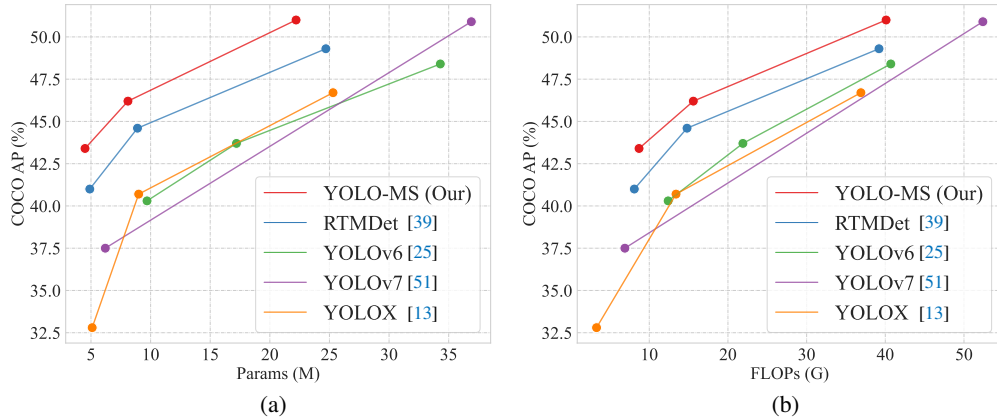
Figure 1: **Comparison with other state-of-the-art real-time object detectors on the MS COCO dataset [33].** (a) AP performance v.s. #parameters. (b) AP performance v.s. computations (FLOPs). The input size utilized to compute the FLOPs is $640 \times 640$. It is clear that our proposed YOLO-MS achieves the best trade-off between performance and computations.

Based on the above design principles, we present our real-time object detector, termed YOLO-MS. To evaluate the performance of our YOLO-MS, we conduct comprehensive experiments on the MS COCO [33] dataset. Quantitative comparisons with other state-of-the-art methods are also provided to showcase the strong performance of our method. As demonstrated in the Fig. 1, YOLO-MS outperforms other recent real-time object detectors with a better computation-performance trade-off. Specifically, the YOLO-MS-XS achieves an AP score of 43%+ on MS COCO, with only 4.5M learnable parameters and 8.7G FLOPs. YOLO-MS-S and YOLO-MS yield 46%+ AP and 51%+ AP, respectively, with 8.1M and 22.2M learnable parameters. Moreover, our work can also be used as a plug-and-play module for other YOLO models. Typically, our method significantly improves the AP of YOLOv8 from 37%+ to 40%+ with even fewer parameters and FLOPs.

## 2 Related Work

### 2.1 Real-time Object Detection

The task of object detection is to detect objects in a specific scene. Despite the outstanding performance achieved by the multi-stage detectors [15, 14, 46, 31, 18, 2, 54, 41, 4] and end-to-end detectors [62, 3, 40, 26, 34, 58], their complicated structures often hinder them from achieving real-time performance, which is a prerequisite for real-world applications for object detection. For the trade-off between speed and accuracy, much effort has gone into developing efficient detectors. As opposed to the previous two-stage detectors, most real-time object detection networks adopt the one-stage framework [32, 49, 61, 59, 28, 27]. In particular, the YOLO series [43–45, 1, 50, 22, 25, 51, 39, 56, 23, 38] is the most typical representation.

As a key factor of model performance, the architecture design is the main concern during the development of YOLO. Starting from the ancestor of the YOLO family, *i.e.*, YOLOv1 [43], the network architectures have undergone dramatic changes. Recently, YOLOLv4 [1] modifies the DarkNet with cross-stage partial connections (CSPNet) [52] for better performance. YOLOv6 [25] and PPYOLOE [56] explore the re-parameterization techniques in YOLO, which can achieve higher accuracy without extra inference costs. YOLOv7 [51] proposes Extended Efficient Layer Aggregation Networks (E-ELAN) that can learn and converge effectively by controlling the shortest longest gradient path. RTMDet [39] introduces large kernel convolution ($5 \times 5$) into networks to boost the feature extraction capability of the basic block. The larger receptive fields allow for more comprehensive contextual modeling and significantly improve accuracy.

In this paper, instead of introducing new training or optimization techniques, we focus on improving real-time object detectors by learning more expressive multi-scale feature representations. This makes our method quite different from the previous works mentioned above.

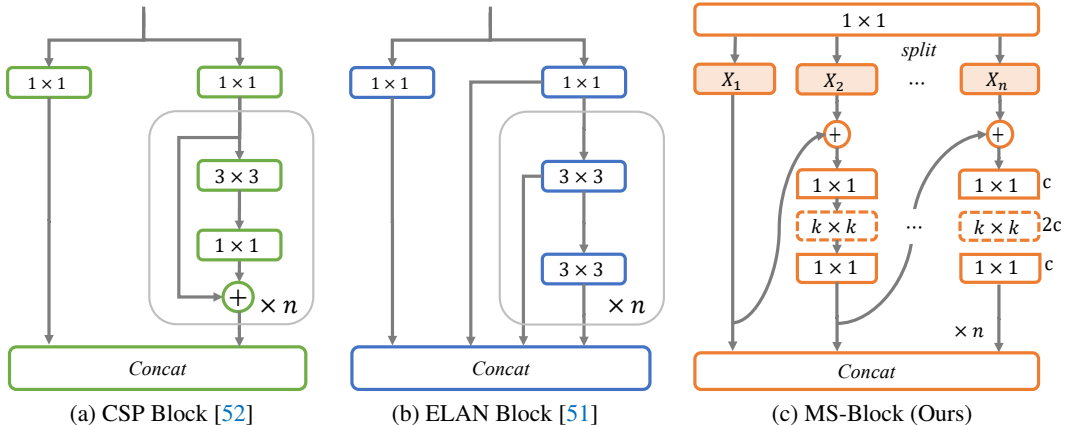| | | |
|---|---|---|
| (a) CSP Block [52] | (b) ELAN Block [51] | (c) MS-Block (Ours) |

Figure 2: **Comparisons with two popular building blocks, which are widely used in previous real-time object detectors.** The dashed box represents the depth-wise convolution. $n$ refers to the number of layers used in the block. $c$ and $k$ correspond to the number of channels and the convolution kernel size, respectively.

## 2.2 Multi-scale feature representation

Multi-scale feature representation learning in computer vision has been studied for a long history [55, 6, 11, 57, 16, 20]. A strong multi-scale feature representation capability can boost the model performance effectively, which has been demonstrated in many tasks [31, 35, 19, 60, 56], including real-time object detection [45, 1, 22, 25, 39].

**Multi-scale feature learning in real-time object detection.** Many real-time object detectors extract multi-scale features by integrating features from different feature levels in the neck part [31, 35]. For instance, YOLOv3 [45] and later YOLO series introduce FPN [31] and PAFPN [35], respectively, to capture rich multi-scale semantics. The SPP [19] module is also widely utilized to enlarge the receptive field. Additionally, multi-scale data augmentations [45] are also widely used as effective training skills. However, the mainstream of the basic building block overlooks the importance of multi-scale feature representation while focusing on how to facilitate detection efficiency or how to introduce new training techniques, especially for CSP [52] block and ELAN [51, 53] block. Differently, our method concentrates on learning richer multi-scale features.

**Large kernel convolution.** Our work is also related to large kernel convolution. Recently, large kernel convolutions have been revitalized in a depth-wise manner [9]. The wider receptive field provided by large kernel convolutions can be a powerful technique for building strong multi-scale feature representations [17, 16, 20, 36, 30, 29]. In the real-time object detection field, RTMDet [39] first attempts to adopt large kernel convolution into networks. However, the kernel size only reaches $5 \times 5$ due to speed limitation. Homogeneous block design in different stages limits the application of large-kernel convolutions. In this paper, we propose a HKS protocol based on some empirical findings. Our HKS protocol adopts convolutional layers with different kernel sizes at different stages, enabling a good trade-off between speed and accuracy with the help of large kernel convolutions.

## 3 Method

As a crucial topic in modern object detectors, multi-scale feature representation has a great influence on the detection performance [31, 35]. In this section, we analyze how to design a powerful encoder architecture that can effectively learn expressive multi-scale feature representations from two perspectives.

### 3.1 Multi-Scale Building Block Design

CSP Block is a stage-level gradient path-based network that balances the gradient combinations and computational cost [52]. It is a fundamental building block widely utilized in the YOLO series.

3

There are several variants having been proposed, including the original version in YOLOv4 [1] and YOLOv5 [22], CSPVoVNet [24] in Scaled YOLOv4 [50], ELAN [53] in YOLOv7 [51], and large kernel unit proposed in RTMDet [39]. We present the structure of the original CSP block and ELAN in the Fig. 2(a) and Fig. 2(b) respectively.

A key aspect that has been overlooked by the aforementioned real-time detectors is how to encode multi-scale features in the basic building blocks. One of the powerful design principles is Res2Net [12], which aggregates features from different hierarchies to enhance the multi-scale representations. However, this principle does not thoroughly explore the role of large kernel convolutions, which have demonstrated effective in CNN-based models for visual recognition tasks [20, 36, 9]. The main obstacle in incorporating large kernel convolutions into Res2Net is the computational overhead they introduce as the building block adopts the standard convolutions. In our method, we propose to replace the standard $3 \times 3$ convolution with the inverted bottleneck layer [47] to enjoy large-kernel convolutions.

**MS-Block**. Based on the earlier analysis, we propose a novel block with a hierarchical feature fusion strategy [12], dubbed MS-Block, to enhance the capability of real-time object detectors in extracting multi-scale features while maintaining a fast inference speed. The specifics of MS-Block are illustrated in Fig. 2(c). Suppose $X \in R^{H \times W \times C}$ refers to the input feature. After the transition via a $1 \times 1$ convolution, the channel dimension of $X$ is increased to $n \times C$. Then, we split $X$ into $n$ distinct groups, denoted as $\{X_i\}$ where $i \in 1, 2, 3, ..., n$. To lower the computational cost, we choose $n$ to be 3. Note that apart from $X_1$, each other group goes through an inverted bottleneck layer, denoted by $IB_{k \times k}(\cdot)$ where $k$ refers to kernel size, to attain $Y_i$. The mathematical representation of $Y_i$ can be described as follows:

$$Y_i = \begin{cases} X_i, & i = 1 \\ IB_{k \times k}(Y_{i-1} + X_i). & i > 1 \end{cases} \quad (1)$$

Following the formula, we do not connect the inverted bottleneck layer to $X_1$, allowing it to act as a cross-stage connection and retain information from the previous layers. Finally, we concatenate all splits and apply a $1 \times 1$ convolution to conduct an interaction among all splits, each of which encodes features at different scales. This $1 \times 1$ convolution is also used to adjust the channel numbers when the network goes deeper.

### 3.2 Heterogeneous Kernel Selection Protocol

In addition to the design of building block, we also delve into the usage of convolution from a macro perspective. Previous real-time object detectors adopt homogeneous convolutions (i.e., convolutions with the same kernel size) in different encoder stages but we argue that this is not the optimal option for extracting multi-scale semantic information.

In a pyramid architecture, high-resolution features extracted from the shallow stages of the detectors are commonly used to capture fine-grained semantics, which will be used to detect small objects. On the contrary, low-resolution features from the deeper stages of the network are utilized to capture high-level semantics, which will be used to detect large objects. If we adopt uniform small-kernel convolutions for all stages, the *Effective Receptive Field (ERF)* in the deep stage is limited, affecting the performance on large objects. Incorporating large-kernel convolutions in each stage can help address this limitation. However, large kernels with a large *ERF* can encode broader regions, which increases the probability of including contaminative information outside the small objects and decreases the inference speed as analyzed in our experiment section.



Figure 3: **Illustration of HKS protocol**. Note that the shades of color in the figure represent the kernel size $k$ used in the MS-Block.

In this work, we propose to leverage heterogeneous convolutions in different stages to assist in capturing richer multi-scale features. To be specific, we leverage the smallest-kernel convolution in
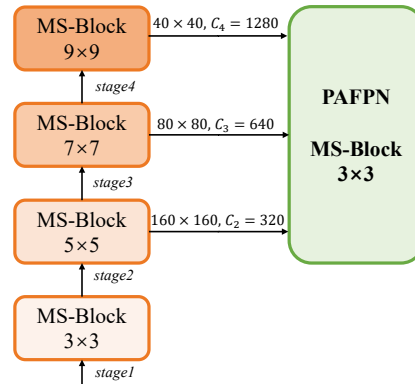
the first stage of the encoder while the largest-kernel convolution is in the last stage. Subsequently, we incrementally increase the kernel size for the middle stages, aligning it with the increment of feature resolution. This strategy allows for extracting both fine-grained and coarse-grained semantic information, enhancing the multi-scale feature representation capability of the encoder.

As presented in Fig. 3, we assign values of $k$ to 3, 5, 7, and 9 from the shallow stage to the deep stage in the encoder. We call this Heterogeneous Kernel Selection (HKS) Protocol. Our HKS protocol is able to enlarge the receptive fields in the deep stages without introducing any other impact to the shallow stages. Fig. 4 in Sec. 4 support our analysis. Additionally, HKS not only helps encode richer multi-scale features but also ensures efficient inference.

As shown in Table 1, applying large-kernel convolutions to high-resolution features is computationally expensive. However, our HKS protocol adopts large-kernel convolution on low-resolution features, thereby substantially reducing computational costs compared to only using large-kernel convolutions. In practice, we empirically found that our YOLO-MS with HKS protocol achieves nearly the same inference speed as that only using depth-wise $3 \times 3$ convolutions.

Table 1: **FPS ($\times 10^3$) of convolutions with different kernel sizes in different stages of the network.** The **gray** color indicates the kernel size used in each stage of our YOLO-MS. The computation environment of the benchmark for all models is the same.

| Stage | Input Size | #Channels | $3 \times 3$ | $5 \times 5$ | $7 \times 7$ | $9 \times 9$ |
|-------|-----------|-----------|--------------|--------------|--------------|--------------|
| #1 | $320 \times 320$ | 160 | **2.72** | 1.38 | 0.93 | 0.65 |
| #2 | $160 \times 160$ | 320 | 5.53 | **2.78** | 1.86 | 1.31 |
| #3 | $80 \times 80$ | 640 | 10.46 | 5.52 | **3.65** | 2.65 |
| #4 | $40 \times 40$ | 1280 | 14.25 | 10.73 | 7.21 | **5.21** |

## 3.3 Architecture

As shown in Fig. 3, the backbone of our model consists of four stages, each of which is followed by a $3 \times 3$ convolution with stride 2 for downsampling. We add an SPP block [19] after the third stage as done in [39]. Following [1, 52], we use PAFPN as neck to build feature pyramid [31, 35] upon our encoder. It fuses multi-scale features extracted from different stages in the backbone. The basic building blocks used in the neck are also our MS-Block, in which $3 \times 3$ depth-wise convolutions are used for fast inference.

In addition, to achieve a better trade-off between speed and accuracy, we halve the channel depth of the multi-level features from the backbone. We present three variants of YOLO-MS, namely YOLO-MS-XS, YOLO-MS-S, and YOLO-MS. The detailed configurations of different scales of our YOLO-MS are listed in Tab. 2. For other parts of YOLO-MS, we keep them the same to [39].

Table 2: **Brief configurations of the proposed YOLO-MS.** $c_i$ refers to the channel number of features of the $i$th stage. "Channel Expansion Ratio" indicates the ratio of the channel expansion process within the inverted bottleneck of the MS-Block. Based on RTMDet, we implement three variants with parameters 4.54M, 8.13M, and 22.17M, respectively.

| Model | $\{C_1, C_2, C_3, C_4\}$ | Channel Expansion Ratio | Params(M) | FLOPs(G) |
|-------|--------------------------|-------------------------|-----------|----------|
| YOLO-MS-XS | $\{32, 64, 128, 256\}$ | $c \rightarrow 2c$ | 4.54 | 8.74 |
| YOLO-MS-S | $\{43, 86, 172, 344\}$ | $c \rightarrow 2c$ | 8.13 | 15.58 |
| YOLO-MS | $\{64, 128, 256, 512\}$ | $c \rightarrow 2c$ | 22.17 | 40.09 |

# 4 Experiments

## 4.1 Experiment Setup

**Implementation details.** Our implementation is based on the MMDetection [5] framework and PyTorch [42]. All experiments are conducted with a batch size of 8 per GPU to guarantee impartiality on a machine with 8 V100 GPUs. All the scales of YOLO-MS are trained from scratch for 300

Table 3: **Effectiveness of MS-Block with different channel expansion ratio.** $r$ refers to the channel expansion ratio of the inverted bottleneck. All models are benchmarked under the same computational environment.

| $r$ | Params(M) | FLOPs(G) | FPS | AP |
|---|---|---|---|---|
| 1 | 4.13 | 8.24 | 134.1 | 42.2 |
| 2 | 4.54 | 8.74 | 130.3 | 43.4 |
| 3 | 5.12 | 9.53 | 125.9 | 44.4 |

Table 4: **Effectiveness of MS-Block with feature fusion strategy.** *fuse* denotes integrating the features from two branches through addition, as described in Equation (1).

| fuse | AP | $AP_s$ | $AP_m$ | $AP_l$ |
|---|---|---|---|---|
|  | 42.2 | 23.2 | 46.6 | 58.0 |
| ✓ | 43.4 | 23.7 | 48.3 | 60.3 |

epochs without relying on other large-scale datasets, like ImageNet [8], or pre-trained weights. More implementation details can be found in the supplementary materials.

**Datasets.** We evaluate the proposed detector on the widely used MS COCO [33] benchmark, following the standard practice in most previous works. Specifically, we use the *train2017* set, which includes $115K$ images, for training, and the *val2017* set, which includes $5K$ images, for validation. The standard COCO-style measurement, i.e., Average Precision (AP), is utilized as the primary challenge metric for evaluation. In addition, the mAP with IoU thresholds of 0.5 and 0.75, as well as AP for small, medium, and large objects, are reported as supporting metrics.

**Benchmark Settings.** Following previous works, we measure all models' frames per second (FPS) using an NVIDIA 3090 GPU in a full-precision floating-point format (FP32). During testing, we perform inference using TensorRT 8.4.3 and cuDNN 8.2.0, without the NMS post-processing step. The batch size used for the inference process is set to 1. Additionally, the FLOPs are calculated based on an input size of $640 \times 640$, using the MMDetection Framework [5].

## 4.2 Analysis of MS-Block

In this subsection, we conduct a series of ablation analysis on our MS-Block. By default, we use YOLO-MS-XS for all experiments.

**Inverted Bottleneck.** We conduct an ablation study on the channel expansion ratio denoted as $r$ in the inverted bottleneck of the MS-Block, and the results are shown in Tab. 3. The results reveal that the detector achieves the best trade-off when $r = 2$. Furthermore, the channel expansion rate significantly affects the detector's performance, with $1.2$ AP lower for $r = 1$ and $1$ AP higher for $r = 3$ than $r = 2$. This implies that channel expansion within the inverted bottleneck strengthens the representational power of the depth-wise convolution, abounding the features' semantic richness. However, while $r = 3$ yields the highest performance, the computational cost is higher. Hence, we use $r = 2$ as the default setting in all subsequent experiments to maintain high computational efficiency.

**Feature Fusion Strategy.** In general, MS-Block progressively fuses the features between adjacent branches through addition. We conduct an ablation study to evaluate the effectiveness of the feature fusion strategy. The results are presented in Tab. 4, which indicates that the feature fusion operation between branches is critical to improve the model performance. Particularly, it brings a significant improvement of +1.2% AP to YOLO-MS.

**Number of MS-Layers.** We also analyze the computational cost and inference speed of different numbers of MS-Layers, represented by $N_l$. The results are demonstrated in Tab. 5. One can see that the number of MS-Layers in MS-Block significantly impacts the speed of YOLO-MS. For instance, in the case of YOLO-MS-XS, as $N_l$ increases from 1 to 2 and then to 3, the number of parameters increases by 25.8% and 51.5%, respectively. Also, the FLOPs increase by 18.1% and 36.2%, respectively. The FPS of the inference process when $N_l = 2$ and $N_l = 3$ is also dropped by 9.2% and 16.6%, respectively. Hence, we use $N_l = 1$ as the default setting in all subsequent experiments.

**Attention Mechanism.** In accordance with RTMDet, we utilize SE attention [21] after the last $1 \times 1$ convolution to capture inter-channel correlations. We perform experiments to study the impact of channel attention here. The computational analysis is presented in Tab. 6, while the performance is

Table 5: **Computational cost analysis of the number of the layers.** $N_l$ refers to the number layers of MS-Block. All models are benchmarked under the same computational environment.

| Scale | $N_l$ | AP | Params(M) | FLOPs(G) | FPS |
|---|---|---|---|---|---|
| XS | 1 | 43.4 | 4.54 | 8.74 | 130.3 |
| XS | 2 | 44.7 | 5.71 | 10.32 | 117.9 |
| XS | 3 | 45.4 | 6.88 | 11.90 | 108.7 |
| S | 1 | 46.2 | 8.13 | 15.58 | 111.1 |
| S | 2 | - | 10.18 | 18.34 | 100.9 |
| S | 3 | - | 12.24 | 21.09 | 93 |

Table 6: **Computational cost analysis of the attention mechanisms.** All models are benchmarked under the same computational environment.

| Scale | SE | AP | Params(M) | FLOPs(G) | FPS |
|---|---|---|---|---|---|
| XS |  | 43.1 | 4.49 | 8.739 | 134.6 |
| XS | ✓ | 43.4 | 4.54 | 8.742 | 130.3 |
| S |  | 46.2 | 8.05 | 15.575 | 114.9 |
| S | ✓ | 46.2 | 8.13 | 15.579 | 111.1 |
| - |  | 51.0 | 21.99 | 40.088 | 89.0 |
| - | ✓ | 50.8 | 22.17 | 40.094 | 79.6 |

Table 7: **Ablation study of the branches number.** $N_b$ refers to the number of branches in the inverted bottleneck of MS-Block. All models are benchmarked under the same computational environment.

| $N_b$ | Params | FLOPs | FPS | AP | $AP_{50}$ | $AP_{75}$ |
|---|---|---|---|---|---|---|
| 2 | 4.63M | 8.86G | 134.8 | 42.2 | 59.1 | 46.0 |
| 3 | 4.54M | 8.74G | 130.3 | 43.4 | 60.4 | 47.6 |
| 4 | 4.36M | 8.52G | 124.9 | 43.2 | 59.9 | 47.1 |

Table 8: **Comparison with different kernel size settings.** $k_i$ represents the kernel size in the $i$th stage. All models are benchmarked under the same computational environment.

| $[k_1, k_2, k_3, k_4]$ | AP | $AP_{50}$ | $AP_{75}$ | Params | FLOPs |
|---|---|---|---|---|---|
| $[3, 3, 3, 3]$ | 42.4 | 59.6 | 46.5 | 4.44 | 8.65 |
| $[5, 5, 5, 5]$ | 42.9 | 59.8 | 46.7 | 4.47 | 8.75 |
| $[7, 7, 7, 7]$ | 42.8 | 59.8 | 46.6 | 4.52 | 8.90 |
| $[9, 9, 9, 9]$ | 42.9 | 59.8 | 46.8 | 4.58 | 9.10 |
| $[11, 11, 11, 11]$ | 43.2 | 60.1 | 47.2 | 4.61 | 9.34 |
| $[9, 7, 5, 3]$ | 42.5 | 59.5 | 46.1 | 4.47 | 8.97 |
| $[3, 7, 11, 15]$ | 42.9 | 60.1 | 46.5 | 4.68 | 8.90 |
| $[5, 7, 9, 11]$ | 42.9 | 60.0 | 46.9 | 4.56 | 8.87 |
| $[3, 5, 7, 9]$ | 43.4 | 60.4 | 47.5 | 4.54 | 8.74 |

shown in Tab. 12. Interestingly, the attention mechanism can only slightly improve the performance but slows down the inference time. Therefore, users can selectively use channel attention according to their own conditions.

**Number of branches.** Our MS-Block partitions the feature input and propagates it through multiple branches. However, increasing the branch number also leads to the increase of MS-Layers and the decrease the channel number in each branch. To investigate the impact of the number of branches, denoted as $N_b$, we conduct an ablation study. The results are reported in Tab. 7. Intriguingly, directly expanding the number of branches does not always lead to performance improvement. Specifically, when $N_b = 3$, YOLO-MS reaches the best performance of 43.4% AP, which is 1.2% and 0.2% higher than $N_b = 2$ and $N_b = 4$, respectively. Hence, we use $N_b = 3$ as the default setting in all subsequent experiments.

**Ablation study on PAFPN module.** We have performed an ablation study on the PAFPN module, and the results are listed in Tab. 9. PAFPN [35] is a popular structure that has been widely used in other YOLOs. We remove the PAFPN from YOLO-MS to further verify the effectiveness. Experimental results indicate that our proposed method with few computational costs can produce nearly the same performance as PAFPN without pre-trained weights. Besides, our proposed method also outperforms the baseline without PAFPN. Moreover, our method is orthogonal to the FPN module. We present the comparison between original PAFPN and PAFPN-MS (PAFPN with MS-Block). As demonstrated, the detector with PAFPN-MS achieves better performance (+0.2% AP) with only ~ 60% parameters and ~ 80% FLOPs.

**Analysis of image resolution.** Here, we conduct an experiment to investigate the correlation between image resolution and the design of the multi-scale building block. We apply the Test Time Augmentation during inference, performing multi-scaling transformations ($320 \times 320$, $640 \times 640$, and $1280 \times 1280$) on the image. In addition, we also use these resolutions individually for testing. Note that the image resolution we used in training is $640 \times 640$. The results are provided in Tab. 10. Experimental results demonstrate a consistent trend: as the image resolution increases, there is also

Table 9: **Ablation study of FPN module.** * refers to train with pertained models. 'PAFPN-MS' refers to the PAFPN with MS-Block.

| Model | Neck | Params | FLOPs | AP |
|---|---|---|---|---|
| RTMDet-tiny | PAFPN | 4.9 | 8.1 | 38.8 |
| RTMDet-tiny* | - | 4.1 | 8.6 | 34.8 |
| YOLO-MS-XS | - | 4.0 | 8.7 | 38.0 |
| YOLO-MS-XS | PAFPN | 6.6 | 10.7 | 43.2 |
| YOLO-MS-XS | PAFPN-MS | 4.5 | 8.7 | 43.4 |

Table 10: **Comparison with different kernel size settings.** 'Ours' represents YOLO-MS-XS. 'TTA' refers to the Test Time Augmentation.

| Model | Resolution | AP | $AP_s$ | $AP_m$ | $AP_l$ |
|---|---|---|---|---|---|
| RTMDet-tiny | $320 \times 320$ | 30.0 | 8.3 | 36.0 | 54.0 |
| Ours | | 32.8 | 10.2 | 36.0 | 55.7 |
| RTMDet-tiny | $640 \times 640$ | 41.0 | 20.7 | 45.3 | 58.0 |
| Ours | | 43.4 | 23.7 | 48.3 | 60.3 |
| RTMDet-tiny | $1280 \times 1280$ | 35.2 | 29.2 | 45.8 | 36.4 |
| Ours | | 38.2 | 30.7 | 46.5 | 39.2 |
| RTMDet-tiny | TTA | 41.9 | 28.1 | 47.1 | 55.5 |
| Ours | | 44.8 | 30.5 | 48.7 | 58.6 |

Table 11: **Integration with other YOLO detectors.** * refers to train with pertained models. The performance results of other detectors are referenced from MMDetection [5].

| Model | Epochs | AP | $AP_{50}$ | $AP_{75}$ | $AP_s$ | $AP_m$ | $AP_l$ | Params (M) | FLOPs (G) |
|---|---|---|---|---|---|---|---|---|---|
| RTMDet-tiny* [39] | 300 | 41.0 | 57.4 | 44.4 | 20.7 | 45.3 | 58.0 | 4.9 | 8.1 |
| YOLO-MS-XS | 300 | 43.4 (+2.4) | 60.4 | 47.6 | 23.7 | 48.3 | 60.3 | 4.5 | 8.7 |
| YOLOv6-tiny [25] | 400 | 41.0 | 57.4 | 43.9 | 21.2 | 45.7 | 57.7 | 9.7 | 12.4 |
| YOLOv6-MS-tiny | 400 | 43.5 (+2.5) | 57.4 | 44.4 | 26.0 | 48.3 | 57.8 | 8.1 | 9.6 |
| YOLOv8-n [23] | 500 | 37.2 | 52.7 | 40.3 | 18.9 | 40.5 | 52.5 | 3.2 | 4.4 |
| YOLOv8-MS-n | 500 | 40.3 (+3.1) | 56.4 | 43.9 | 22.0 | 44.6 | 53.7 | 2.9 | 4.4 |

an increase on $AP_s$. However, we can achieve higher $AP_l$ for low-resolution images. This also verifies the effectiveness of our HKS protocol.

**Application to other YOLOs.** Our proposed methods can be used as a plug-and-play module for other YOLO models. To demonstrate the generalization ability of our method, we apply the proposed method to other YOLO models and conducted a comprehensive comparison on MS COCO. The results are listed in Tab. 11. As shown below, the AP scores of YOLOv6 [25] and YOLOv8 [23] can be significantly increased to 43.5% (+2.5%) and 40.3% (+3.1%), with even fewer parameters and FLOPs, respectively.

### 4.3 Analysis of HKS Protocol

In this subsection, we perform experiments to evaluate the effectiveness of HKS by exploring different kernel size settings. To simplify notation, we use the format $[k_1, k_2, k_3, k_4]$, where $k_i$ represents the kernel size in stage $i$. We investigate both homogeneous kernel size settings with 3, 5, 7, 9, and 11, as well as an inverted version of HKS, *i.e.* $[9, 7, 5, 3]$.

The results of our experiments, as shown in Tab. 8, reveal interesting insights. We observe that simply increasing the convolution kernel size does not always lead to significant performance improvements. However, when we employ HKS, we achieve a substantial boost in performance (43.4% AP), which is better than all other settings with homogeneous kernel size.

Furthermore, the order in which the convolution kernels are arranged within the stages plays a crucial role. Specifically, when a large kernel is used in the shallow stage, and a small kernel is used in the deep stage, the performance drops by 0.9% AP compared to HKS. This suggests that the deep stage requires a larger receptive field to effectively capture coarse-grained information, in contrast to the shallow stage.

Considering the computational costs, our HKS stands out as it incurs the least computation overhead. This indicates that by strategically placing convolutions with different kernel sizes in suitable positions, we can maximize the efficient utilization of these convolutions.
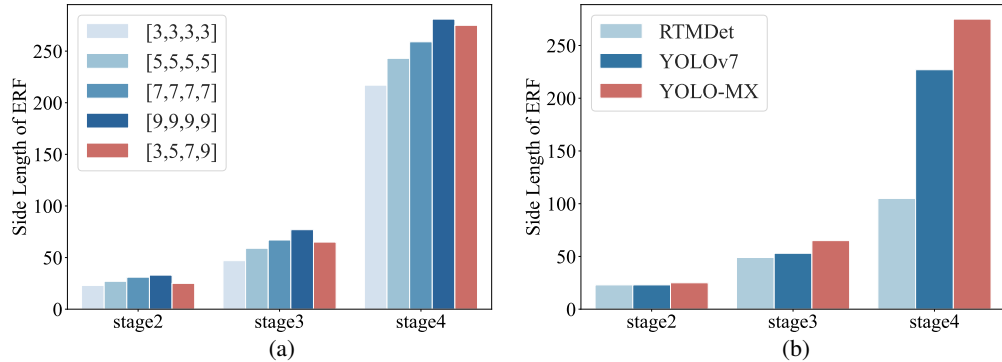
Figure 4: **Statistics analysis of the Effective Receptive Field.** (a) Comparison between different kernel size settings. (b) Comparison between different real-time detectors. The red bar represents our YOLO-MS. The legend of the left part refers to the kernel size used in different stages.
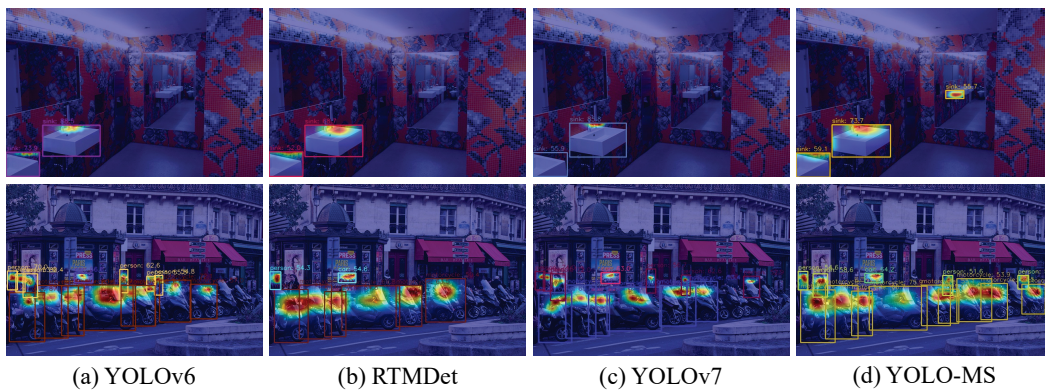


| (a) YOLOv6 | (b) RTMDet | (c) YOLOv7 | (d) YOLO-MS |

Figure 5: **Visual comparison with state-of-the-art through Grad-CAM [48].**

We also conduct experiments using two new settings, $[5, 7, 9, 11]$ and $[3, 7, 11, 15]$. The results are provided in the table below. According to the results, it intuitively indicates the setting of $[3, 5, 7, 9]$ achieves better performance with lower computational costs.

**Effective Receptive Field Analysis**. Previous studies [37, 10] have introduced the concept of the *ERF* as a metric to understand the behavior of deep convolutional neural networks (CNNs). *ERF* measures the effective region in the input space influenced by a feature representation. Here, we leverage the concept of *ERF* further to investigate the effectiveness of HKS. Specifically, we measure the side length of the area containing high-contribution pixels in the *ERF* for stages 2, 3, and 4 of the encoder. Notably, the computation of the *ERF* follows the methodology described in [10], and detailed information can be found in the supplementary materials.

The visual comparison is presented in Fig. 4. As shown in Fig. 4(a), as the kernel size increases, the area of *ERF* also becomes larger in all stages, which supports the positive correlation between kernel size and receptive fields. Moreover, in the shallow stage, the area of *ERF* is smaller than most other settings, while in the deep stage, it is the opposite. This observation indicates that the protocol effectively enlarges the receptive fields in the deep stage without compromising the shallow stage. In the Fig. 4(b), we can observe that our HKS achieves the largest *ERF* in the deep stage, allowing us to detect large objects better.

## 4.4 Comparison with State-of-the-Arts

**Visualization comparison with CAM.** To evaluate which part of the image attracts the attention of the detector, we use Grad-CAM [48] to generate class response maps. We visualize the class response maps generated from the neck of YOLOv6-tiny [25], RTMDet-tiny [39], YOLOV7-tiny [51], and YOLO-MS-XS. We also select typical images of different sizes, including small, medium, and

Table 12: **Comparison with state-of-the-art real-time object detectors.** * refers to train with pertained models. † refers to without SE attention [21]. All models are benchmarked under the same computational environment. The performance results of other detectors are referenced from MMDetection [5].

| Model | AP | $AP_{50}$ | $AP_{75}$ | $AP_s$ | $AP_m$ | $AP_l$ | Latency (ms) | Params (M) | FLOPs (G) |
|---|---|---|---|---|---|---|---|---|---|
| YOLOv5-n [22] | 28.0 | 45.9 | 29.4 | 14.0 | 31.8 | 36.6 | 4.8 | 1.9 | 2.3 |
| YOLOX-tiny [13] | 32.8 | 50.3 | 34.8 | 14.0 | 35.5 | 48.3 | 4.7 | 5.1 | 7.6 |
| YOLOv6-n [25] | 36.2 | 51.6 | 39.2 | 16.8 | 40.2 | 52.6 | 7.5 | 4.3 | 5.5 |
| YOLOv7-tiny [51] | 37.5 | 55.8 | 40.2 | 19.9 | 41.1 | 50.8 | 5.2 | 6.2 | 6.9 |
| YOLOv6-tiny [25] | 40.3 | 57.4 | 43.9 | 21.2 | 45.6 | 57.5 | 9.6 | 9.7 | 12.4 |
| RTMDet-tiny* [39] | 41.0 | 57.4 | 44.4 | 20.7 | 45.3 | 58.0 | 7.2 | 4.9 | 8.1 |
| YOLO-MS-XS | 43.4 | 60.4 | 47.6 | 23.7 | 48.3 | 60.3 | 7.6 | 4.5 | 8.7 |
| YOLO-MS-XS† | 43.1 | 60.1 | 46.8 | 24.0 | 47.8 | 59.1 | 7.4 | 4.5 | 8.7 |
| YOLOv5-s [22] | 37.7 | 57.1 | 41.0 | 21.7 | 42.5 | 48.8 | 5.4 | 7.2 | 8.3 |
| YOLOX-s [13] | 40.7 | 59.6 | 44.3 | 23.9 | 45.2 | 53.8 | 5.8 | 9.0 | 13.4 |
| YOLOv6-s [25] | 43.7 | 60.8 | 47.0 | 23.6 | 48.7 | 59.8 | 8.9 | 17.2 | 21.9 |
| RTMDet-s* [39] | 44.6 | 61.7 | 48.3 | 24.2 | 49.2 | 61.8 | 7.7 | 8.9 | 14.8 |
| YOLO-MS-S | 46.2 | 63.7 | 50.5 | 26.9 | 50.5 | 63.0 | 9.0 | 8.1 | 15.6 |
| YOLO-MS-S† | 46.2 | 63.6 | 50.3 | 27.5 | 50.6 | 62.9 | 8.7 | 8.1 | 15.6 |
| YOLOv5-m [22] | 45.3 | 64.1 | 49.4 | 28.4 | 50.8 | 57.7 | 7.4 | 21.2 | 24.5 |
| YOLOv5-l [22] | 48.8 | 67.4 | 53.3 | 33.5 | 54.0 | 61.8 | 11.0 | 46.6 | 54.6 |
| YOLOX-m [13] | 46.7 | 65.6 | 50.3 | 29.0 | 51.2 | 60.9 | 8.2 | 25.3 | 36.9 |
| YOLOv6-m [25] | 48.4 | 65.7 | 52.7 | 30.0 | 54.1 | 64.5 | 10.6 | 34.3 | 40.7 |
| YOLOv6-l [25] | 51.0 | 68.4 | 55.2 | 33.5 | 56.2 | 67.3 | 14.7 | 58.5 | 71.4 |
| RTMDet-m [39] | 49.3 | 66.9 | 53.9 | 30.5 | 53.6 | 66.1 | 10.5 | 24.7 | 39.2 |
| YOLO-MS | 51.0 | 68.6 | 55.7 | 33.1 | 56.1 | 66.5 | 12.3 | 22.2 | 40.1 |
| YOLO-MS† | 50.8 | 68.4 | 55.4 | 33.2 | 54.8 | 66.4 | 11.2 | 22.0 | 40.1 |

large objects, from the MS COCO dataset [33]. The visualization results are shown in the Fig. 5. Both YOLOv6-tiny, RTMDet-tiny, and YOLOV7-tiny fail to detect small, densely-packed objects, such as crowds of people, and ignore parts of the objects. Conversely, YOLO-MS-XS exhibits a strong response for all objects in the class response maps, indicating its remarkable multi-scale feature representation ability. Furthermore, it highlights that our detector achieves excellent detection performance for objects of different sizes and images containing objects of various densities.

**Quantitative comparison.** We compare YOLO-MS with the current state-of-the-art object detectors. From the Tab. 12, we can see YOLO-MS achieves the remarkable speed-accuracy trade-off. Compared with the second-best tiny detector, *i.e.* RTMDet [39], **YOLO-MS-XS** reaches 43.4% AP, which is 2.3% AP higher than it with the ImageNet [8] pretrained model. **YOLO-MS-S** achieves 46.2% AP which brings 5.7% AP improvement with around half parameters size compared with YOLOv6. Furthermore, **YOLO-MS** exhibits a detection performance of 51.0% AP, which is superior to state-of-the-art object detectors with similar parameters and computational complexity, even the large-scale models, *e.g.,* YOLOv6-M and YOLOv6-L. In conclusion, YOLO-MS is able to serve as a promising baseline for real-time object detection, offering powerful multi-scale feature representation.

## 5   Conclusions and Discussions

This paper proposes a high-performance real-time object detector with a reasonable computational cost. To achieve this goal, we investigate the usage of convolution with varying kernel sizes from global and local perspectives and build an encoder with a powerful ability to extract multi-scale feature representations. Our experimental studies find that the proposed HKS strategy and MS-Block significantly enhance the speed-accuracy trade-off of our detector and outperform other real-time detectors. We hope to bring new insight to the object detection community.

**Limitations.** Although YOLO-MS achieves remarkable performance with a reasonable computational cost, there is still a gap in inference speed compared to state-of-the-art real-time detectors. It is mainly caused by the inefficiency of the large kernel convolution and the hierarchy structure. This

encourages us further to optimize the usage of large kernel convolution and streamline the structure of YOLO-MS for a faster version of our YOLO-MS.

# References

[1] Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao. Yolov4: Optimal speed and accuracy of object detection, 2020. 1, 2, 3, 4, 5

[2] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delving into high quality object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 2

[3] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*, pages 213–229. Springer, 2020. 1, 2

[4] Kai Chen, Jiangmiao Pang, Jiaqi Wang, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jianping Shi, Wanli Ouyang, Chen Change Loy, and Dahua Lin. Hybrid task cascade for instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 2

[5] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jiarui Xu, Zheng Zhang, Dazhi Cheng, Chenchen Zhu, Tianheng Cheng, Qijie Zhao, Buyu Li, Xin Lu, Rui Zhu, Yue Wu, Jifeng Dai, Jingdong Wang, Jianping Shi, Wanli Ouyang, Chen Change Loy, and Dahua Lin. MMDetection: Open mmlab detection toolbox and benchmark. *arXiv preprint arXiv:1906.07155*, 2019. 5, 6, 8, 10

[6] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation, 2017. 3

[7] Jifeng Dai, Yi Li, Kaiming He, and Jian Sun. R-fcn: Object detection via region-based fully convolutional networks. *Advances in neural information processing systems*, 29, 2016. 1

[8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 6, 10

[9] Xiaohan Ding, Xiangyu Zhang, Jungong Han, and Guiguang Ding. Scaling up your kernels to 31x31: Revisiting large kernel design in cnns. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11963–11975, 2022. 3, 4

[10] Xiaohan Ding, Xiangyu Zhang, Jungong Han, and Guiguang Ding. Scaling up your kernels to 31x31: Revisiting large kernel design in cnns. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11963–11975, June 2022. 9

[11] Shanghua Gao, Zhong-Yu Li, Qi Han, Ming-Ming Cheng, and Liang Wang. Rf-next: Efficient receptive field search for convolutional neural networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–19, 2022. 3

[12] Shang-Hua Gao, Ming-Ming Cheng, Kai Zhao, Xin-Yu Zhang, Ming-Hsuan Yang, and Philip Torr. Res2net: A new multi-scale backbone architecture. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(2):652–662, 2021. 1, 4

[13] Zheng Ge, Songtao Liu, Feng Wang, Zeming Li, and Jian Sun. Yolox: Exceeding yolo series in 2021, 2021. 1, 2, 10

[14] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, December 2015. 2

[15] 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 (CVPR)*, June 2014. 2

[16] Meng-Hao Guo, Cheng-Ze Lu, Qibin Hou, Zheng-Ning Liu, Ming-Ming Cheng, and Shi min Hu. Segnext: Rethinking convolutional attention design for semantic segmentation. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022. 3

[17] Meng-Hao Guo, Cheng-Ze Lu, Zheng-Ning Liu, Ming-Ming Cheng, and Shi-Min Hu. Visual attention network, 2022. 3

[18] Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Oct 2017. 2

[19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Spatial pyramid pooling in deep convolutional networks for visual recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(9):1904–1916, 2015. 3, 5

[20] Qibin Hou, Cheng-Ze Lu, Ming-Ming Cheng, and Jiashi Feng. Conv2former: A simple transformer-style convnet for visual recognition. *arXiv preprint arXiv:2211.11943*, 2022. 3, 4

[21] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 6, 10

[22] Glenn Jocher. ultralytics/yolov5: v3.1 - Bug Fixes and Performance Improvements. https://github.com/ultralytics/yolov5, Oct. 2020. 1, 2, 3, 4, 10

[23] Glenn Jocher, Ayush Chaurasia, and Jing Qiu. YOLO by Ultralytics, Jan. 2023. 1, 2, 8

[24] Youngwan Lee, Joong-won Hwang, Sangrok Lee, Yuseok Bae, and Jongyoul Park. An energy and gpu-computation efficient backbone network for real-time object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2019. 4

[25] Chuyi Li, Lulu Li, Hongliang Jiang, Kaiheng Weng, Yifei Geng, Liang Li, Zaidan Ke, Qingyuan Li, Meng Cheng, Weiqiang Nie, Yiduo Li, Bo Zhang, Yufei Liang, Linyuan Zhou, Xiaoming Xu, Xiangxiang Chu, Xiaoming Wei, and Xiaolin Wei. Yolov6: A single-stage object detection framework for industrial applications, 2022. 1, 2, 3, 8, 9, 10

[26] Feng Li, Hao Zhang, Shilong Liu, Jian Guo, Lionel M. Ni, and Lei Zhang. Dn-detr: Accelerate detr training by introducing query denoising. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13619–13627, June 2022. 2

[27] Xiang Li, Wenhai Wang, Xiaolin Hu, Jun Li, Jinhui Tang, and Jian Yang. Generalized focal loss v2: Learning reliable localization quality estimation for dense object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11632–11641, June 2021. 2

[28] Xiang Li, Wenhai Wang, Lijun Wu, Shuo Chen, Xiaolin Hu, Jun Li, Jinhui Tang, and Jian Yang. Generalized focal loss: Learning qualified and distributed bounding boxes for dense object detection. *Advances in Neural Information Processing Systems*, 33:21002–21012, 2020. 2

[29] Yuxuan Li, Qibin Hou, Zhaohui Zheng, Ming-Ming Cheng, Jian Yang, and Xiang Li. Large selective kernel network for remote sensing object detection. *arXiv preprint arXiv:2303.09030*, 2023. 3

[30] Yuxuan Li, Xiang Li, and Jian Yang. Spatial group-wise enhance: Enhancing semantic feature learning in cnn. In *Proceedings of the Asian Conference on Computer Vision (ACCV)*, pages 687–702, December 2022. 3

[31] 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. 2, 3, 5

[32] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollar. Focal loss for dense object detection. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Oct 2017. 2

[33] 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 *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014. 2, 6, 10

[34] Shilong Liu, Feng Li, Hao Zhang, Xiao Yang, Xianbiao Qi, Hang Su, Jun Zhu, and Lei Zhang. DAB-DETR: Dynamic anchor boxes are better queries for DETR. In *International Conference on Learning Representations*, 2022. 2

[35] 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 (CVPR)*, June 2018. 3, 5, 7

[36] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11976–11986, 2022. 3, 4

[37] Wenjie Luo, Yujia Li, Raquel Urtasun, and Richard Zemel. Understanding the effective receptive field in deep convolutional neural networks. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016. 9

[38] Wenyu Lv, Shangliang Xu, Yian Zhao, Guanzhong Wang, Jinman Wei, Cheng Cui, Yuning Du, Qingqing Dang, and Yi Liu. Detrs beat yolos on real-time object detection, 2023. 1, 2

[39] Chengqi Lyu, Wenwei Zhang, Haian Huang, Yue Zhou, Yudong Wang, Yanyi Liu, Shilong Zhang, and Kai Chen. Rtmdet: An empirical study of designing real-time object detectors. *arXiv preprint arXiv:2212.07784*, 2022. 1, 2, 3, 4, 5, 8, 9, 10

[40] Depu Meng, Xiaokang Chen, Zejia Fan, Gang Zeng, Houqiang Li, Yuhui Yuan, Lei Sun, and Jingdong Wang. Conditional detr for fast training convergence. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 3651–3660, October 2021. 2

[41] 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/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 2

[42] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019. 5

[43] 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 (CVPR)*, June 2016. 1, 2

[44] Joseph Redmon and Ali Farhadi. Yolo9000: Better, faster, stronger. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.

[45] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement, 2018. 1, 2, 3

[46] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28, 2015. 1, 2

[47] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4510–4520, 2018. 4

[48] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Oct 2017. 9

[49] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. Fcos: Fully convolutional one-stage object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9627–9636, 2019. 1, 2

[50] Chien-Yao Wang, Alexey Bochkovskiy, and Hong-Yuan Mark Liao. Scaled-yolov4: Scaling cross stage partial network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13029–13038, June 2021. 1, 2, 4

[51] Chien-Yao Wang, Alexey Bochkovskiy, and Hong-Yuan Mark Liao. Yolov7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 1, 2, 3, 4, 9, 10

[52] 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. 1, 2, 3, 5

[53] Chien-Yao Wang, Hong-Yuan Mark Liao, and I-Hau Yeh. Designing network design strategies through gradient path analysis. *arXiv preprint arXiv:2211.04800*, 2022. 3, 4

[54] Jiaqi Wang, Kai Chen, Shuo Yang, Chen Change Loy, and Dahua Lin. Region proposal by guided anchoring. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 2

[55] Andrew Witkin. Scale-space filtering: A new approach to multi-scale description. In *ICASSP'84. IEEE International Conference on Acoustics, Speech, and Signal Processing*, volume 9, pages 150–153. IEEE, 1984. 3

[56] Shangliang Xu, Xinxin Wang, Wenyu Lv, Qinyao Chang, Cheng Cui, Kaipeng Deng, Guanzhong Wang, Qingqing Dang, Shengyu Wei, Yuning Du, et al. Pp-yoloe: An evolved version of yolo. *arXiv preprint arXiv:2203.16250*, 2022. 1, 2, 3

[57] Bowen Yin, Xuying Zhang, Qibin Hou, Bo-Yuan Sun, Deng-Ping Fan, and Luc Van Gool. Camoformer: Masked separable attention for camouflaged object detection. *arXiv preprint arXiv:2212.06570*, 2022. 3

[58] Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun Zhu, Lionel M Ni, and Heung-Yeung Shum. Dino: Detr with improved denoising anchor boxes for end-to-end object detection. *arXiv preprint arXiv:2203.03605*, 2022. 1, 2

[59] 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 (CVPR)*, June 2020. 2

[60] Xuying Zhang, Bowen Yin, Zheng Lin, Qibin Hou, Deng-Ping Fan, and Ming-Ming Cheng. Referring camouflaged object detection. *arXiv preprint arXiv:2306.07532*, 2023. 3

[61] Zhaohui Zheng, Ping Wang, Wei Liu, Jinze Li, Rongguang Ye, and Dongwei Ren. Distance-iou loss: Faster and better learning for bounding box regression. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, 2020. 2

[62] 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. 1, 2