

CE-FPN: Enhancing Channel Information for Object Detection

Yihao Luo, Xiang Cao, Juntao Zhang, Xiang Cao, Jingjuan Guo, Haibo Shen, Tianjiang Wang and Qi Feng

Abstract—Feature pyramid network (FPN) has been an effective framework to extract multi-scale features in object detection. However, current FPN-based methods mostly suffer from the intrinsic flaw of channel reduction, which brings about the loss of semantical information. And the miscellaneous fused feature maps may cause serious aliasing effects. In this paper, we present a novel channel enhancement feature pyramid network (CE-FPN) with three simple yet effective modules to alleviate these problems. Specifically, inspired by sub-pixel convolution, we propose a sub-pixel skip fusion method to perform both channel enhancement and upsampling. Instead of the original 1×1 convolution and linear upsampling, it mitigates the information loss due to channel reduction. Then we propose a sub-pixel context enhancement module for extracting more feature representations, which is superior to other context methods due to the utilization of rich channel information by sub-pixel convolution. Furthermore, a channel attention guided module is introduced to optimize the final integrated features on each level, which alleviates the aliasing effect only with a few computational burdens. Our experiments show that CE-FPN achieves competitive performance compared to state-of-the-art FPN-based detectors on MS COCO benchmark.

Index Terms—Object detection, Feature pyramid network, Channel enhancement, Sub-pixel convolution

I. INTRODUCTION

Object detection is a fundamental task in computer vision, which is widely applied to various applications, such as object tracking [1], [2], person re-identification [3], [4], etc. With the advances in deep convolutional networks, a number of deep detectors have been developed to achieve remarkable performance recently [5], [6], [7], [8], [9], [10], [11]. Among these detectors, FPN [6] constructs an effective framework to address the issue of scale variations, a primary challenge in object detection. In FPN, multi-scale feature maps are created by propagating the semantical information from high levels into lower levels. By fusing multi-scale features with shallow content descriptive and deep semantical features, FPN-based methods substantially improve the performance of object detection.

However, there exist two widely-held limitations in FPN [12], [8]: (1) Information decay during fusion; (2) Aliasing effects in cross-scale fusion. The existing methods such as

PAFPN [13], Libra R-CNN [12], and AugFPN [8] can alleviate these problems to some extent, but there is still the possibility of further improvement. Meanwhile, in the light of our observations, FPN-based methods also suffer from an intrinsic flaw about channel reduction. We will describe these issues as following:

Information loss of channel reduction. As illustrated in Fig. 1(a), FPN-based methods adopt 1×1 convolutional layers to reduce channel dimensions of the output feature maps C_i from the backbone, which also loses channel information. C_i generally extract thousands of channels in high-level feature maps, which are reduced to a much smaller constant in F_i (e.g. 2048 to 256).

The existed methods [12], [9] mainly add extra modules on the channel-reduced maps rather than make full use of C_i as shown in Fig. 1(b), 1(c). EfficientDet [9] develops various configurations of different FPN channels. It indicates that increasing FPN channels improves the performance with more parameters and FLOPs, so EfficientDet still adopts relatively few channels and proposes complex-connected BiFPN for better accuracy. Therefore, declining channels from the backbone outputs substantially reduces the computation consumption for subsequent prediction but also brings about the loss of information.

Information decay during fusion. The low-level and high-level information are complementary for object detection, while the semantical information would be diluted in the progress of top-down feature fusion [12]. PAFPN [13] and Libra R-CNN [12] propose innovative fusion methods to make full use of features on each level. Nevertheless, the representation ability of high-level semantical feature has not been utilized mostly for larger receptive fields. The exploitation of context information [8] is a proper approach to improve feature representation, which prevents increasing computational burden by adding deeper convolutional layers directly.

Aliasing effects in cross-scale fusion. Cross-scale fusion and skip connections are widely used to improve the performance [12], [9]. The intuitive and simple connections achieve the full use of diverse features on each level. However, there exist semantical differences in cross-scale feature maps, so that direct fusion after interpolation may cause aliasing effects [6]. And the miscellaneous integrated features might confuse the localization and recognition tasks [11]. Motivated by the refinement of Non-local attention [14] on the integrated features, more attention modules could be designed to optimize the fused aliasing features and enhance their discriminative abilities.

In this paper, we propose three novel components to deal

This work was supported in part by the National Nature Science Foundation of China under Grant 61572214. (Corresponding author: Qi Feng, fengqi@hust.edu.cn.)

Yihao Luo, Xiang Cao, Juntao Zhang, Xiang Cao, Haibo Shen, Tianjiang Wang and Qi Feng are with School of Computer Science and Technology, Huazhong University of Science and Technology, Wuhan, Hubei, 430074, China.

Jingjuan Guo is with School of Computer and Information Engineering, Henan University, Kaifeng 475004, China.

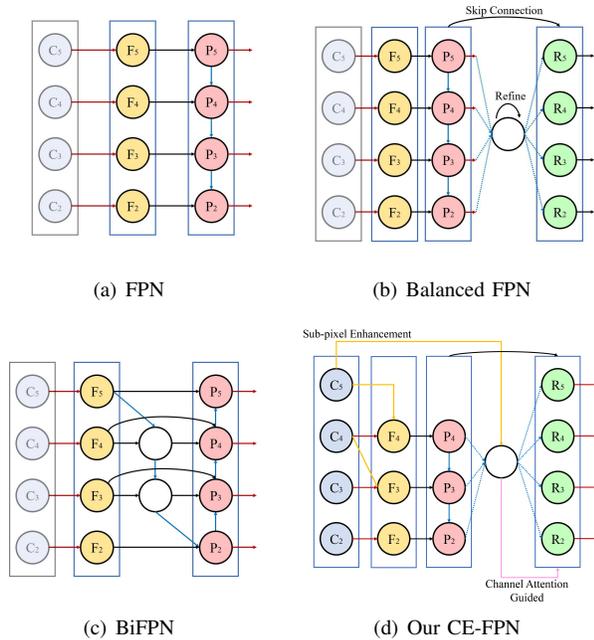


Fig. 1: Comparison of feature pyramid networks design in the case of 4 levels. The translucent nodes indicate underutilization. (a). FPN [6] introduces a top-down pathway to fuse multi-scale features. (b). Libra R-CNN [12] proposes a balanced FPN with integration and refinement. (c). EfficientDet [9] adds an extra bottom-up pathway and skip connections. (d). Our CE-FPN base on the integration framework, in which sub-pixel enhancement and attention guided modules are proposed to make the most of rich channel information.

with the issues above respectively. First, inspired by sub-pixel convolution [15] in super-resolution, we introduce a sub-pixel skip fusion method for utilizing the original cross-scale backbone outputs with rich channel information as shown in Fig. 1(d). Second, we present a sub-pixel context enhancement module for extracting and integrating diverse context information from the highest-level feature map. Sub-pixel convolution is an upsampling method that increases channel dimensions of low-resolution images first, which also brings about extra computation and unreliability. It is worth noting that high-level features in FPN have obtained adequate amounts of channels, which allows sub-pixel convolution to be employed directly. Instead of the original 1×1 convolution and upsampling, the proposed methods can alleviate channel information loss. Hence we extend the original upsampling function of sub-pixel convolution to fuse channel information, which is different from CARAFE [16]. Third, we propose a simple yet effective channel attention guided module to optimize the final integrated features on each level. The attention module alleviates the aliasing effect only with a few computational burdens. We name our whole model as Channel Enhancement Feature Pyramid Network (CE-FPN), which is flexible and generalizable for various FPN-based detectors.

Without bells and whistles, by replacing FPN with CE-FPN, we achieve 38.8 points and 40.9 points Average Precision (AP) with Faster R-CNN when using ResNet-50 and ResNet-101 [17] respectively on MS COCO [18] with the

$1 \times$ schedule in [19]. The experiments results show that our modules improve performance significantly only with slightly computational cost. Meanwhile, CE-FPN achieves competitive performance compared to state-of-the-art FPN-based methods such as Libra R-CNN [12] and AugFPN [8]. The main contributions of our work are summarized as follows:

- We propose two novel channel enhancement methods inspired by sub-pixel convolution. We extend the intrinsic upsampling function of sub-pixel convolution to integrate rich channel information in our modules.
- We introduce simple yet effective channel attention guided module to optimize the integrated features on each level.
- We evaluate the proposed framework on MS COCO and obtain significant improvements over state-of-the-art FPN-based detectors.

II. RELATED WORK

A. Deep object detectors

With the advances in deep convolutional networks, remarkable progress has been achieved in object detection. Object detectors based on deep learning are generally divided into two categories: two-stage detectors and one-stage detectors. The successful two-stage detectors generate regions of interest (ROI) firstly and then refine ROI with classifier and regressor. Faster R-CNN [20] proposes region proposal network (RPN) and develops an end-to-end framework, which significantly improves the efficiency of detectors. The proposed RPN integrates proposal generation with a single convolutional network that has been a paradigm for the two-stage detector. Numerous extended studies of this framework have been proposed which improve the performance significantly, such as Mask R-CNN [21], Cascade R-CNN [22] and CBNet [23].

On the other hand, one-stage detectors adopt a unified network to achieve locations and classifications directly with more efficiency yet less accuracy. SSD [5] handles objects of various sizes on multi-scale features. YOLO [24] makes use of the whole feature map to predict both classification confidences and bounding boxes. RetinaNet [25] follows the FPN framework and utilizes focal loss to suppress the gradients of easy negative samples, which promotes the performance significantly. Besides, there are extensive proposed one-stage detectors for enhancing the network architectures or detection process, such as YOLOv2 [26], YOLOv3 [27] and DSOD [28].

Recently, increasing researches have made remarkable progress from different concerns such as anchor-free [29], [10], multi-scale [13], [12], context extraction [30], [31], and attention module [32], [11].

B. Multi-scale feature augmentation

FPN [6] constructs an effective framework to address the issue of scale variations by merging features via a top-down pathway, which is popularly applied and further studied [8], [11], [32], [12], [9]. PANet [13] investigates an extra bottom-up pathway for further increase of the low-level information in deep layers. Libra R-CNN [12] introduces a balanced

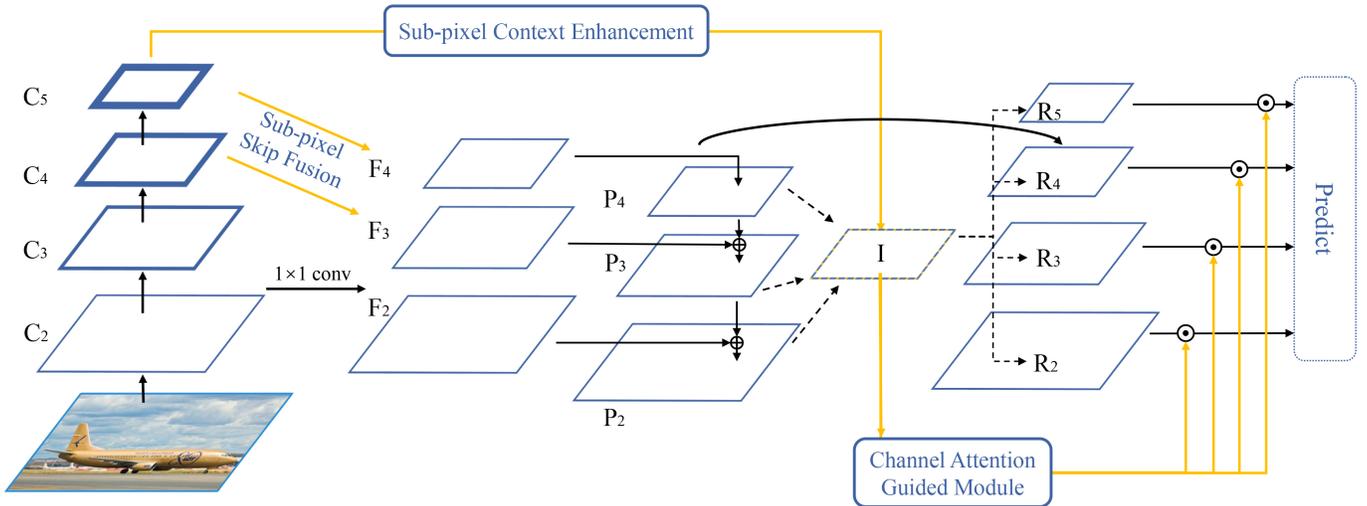


Fig. 2: An overview of our CE-FPN. We follow the framework of integration map [12]. Sub-pixel Skip Fusion (SSF) is proposed to fuse the high-level backbone outputs with the channel-reduced features, which can be performed independently without the integrated feature map I . Sub-pixel Context Enhancement (SCE) extracts and integrates diverse context information from the highest-level feature map to augment the feature representation of I . Channel Attention Guided Module (CAG) adopts a channel attention module to extract channel weights from I , and products the final integrated features respectively.

feature pyramid to pay equal attention to multi-scale features through integration and refinement. NAS-FPN [33] adopts neural architecture search for learning better fusion among all cross-scale connections. EfficientDet [9] proposes a weighted bi-directional FPN to perform easy and fast feature fusion. And AugFPN [8] proposes a series of augmentation methods for FPN.

In the aspect of feature augmentation, context information can facilitate the performance of localization and classification [34]. PSPNet [35] utilizes pyramid pooling to extract hierarchical global context. And [36] proposes a context refinement algorithm to refine each region proposals. Meanwhile, attention mechanism [37] is generally adopted to enhance feature representation in various vision tasks. CoupleNet [38] extracts the attention-related information with global and local features of the objects. MAD [39] searches for neuron activations aggressively from high-level and low-level information streams.

Based on the methods above, we focus on reducing the information loss due to channel decline in FPN construction and optimizing the final features after complicated integration.

III. PROPOSED METHODS

In this section, we introduce a Channel Enhancement Feature Pyramid Network (CE-FPN) to alleviate channel information loss and optimize the integrated features. In CE-FPN, three components are proposed: Sub-pixel Skip Fusion (SSF), Sub-pixel Context Enhancement (SCE), and Channel Attention Guided Module (CAG). We will describe them in detail following.

A. Overall

The overall network architecture is shown in Fig. 2. Following the setting of FPN [6], CE-FPN generates a 4-level

feature pyramid. We denote the output of the backbone as $\{C_2, C_3, C_4, C_5\}$, which have strides of $\{4, 8, 16, 32\}$ pixels with respect to the input image. $\{F_2, F_3, F_4\}$ are the features with the same reduced channels of 256 after 1×1 convolution. The feature pyramid $\{P_2, P_3, P_4\}$ is generated by the top-down pathway in FPN. We remove the nodes of F_5 and P_5 , which are the original highest-level feature with semantical information for FPN. Because that our proposed methods have fully utilized channel information from C_5 . Repetitive feature fusion may cause not only more serious aliasing effects, but also unnecessary computational burdens. The effects of this procedure are analyzed in Sec IV-D. The integration map I is produced through interpolation and max-pooling. And predictions are performed independently at all final results $\{R_2, R_3, R_4, R_5\}$, which corresponds to the feature pyramid of the original FPN.

B. Sub-pixel Skip Fusion

In FPN, Residual networks [17] are widely used as backbone with the output channels of $\{256, 512, 1024, 2048\}$, in which high-level features $\{C_4, C_5\}$ contain rich semantical information. As shown in Fig. 3(a), the 1×1 convolution layers are adopted to reduce channel dimensions of C_i for computation efficiency, which causes a serious loss of channel information. The further studied FPN-based methods [12], [9], [13] generally focus on developing effective modules on the feature pyramid P_i with 256 channels, while the rich channel information of C_i are underutilized.

Based on this observation, we expect that the channel-rich features $\{C_4, C_5\}$ could be developed to improve performance of the resulting feature pyramid. To this end, we introduce a direct fusion method to merge low-resolution (LR) features to high-resolution (HR) inspired by sub-pixel convolution. Sub-pixel convolution is an upsampling method [15], which augments dimensions of width and height via shuffling pixels

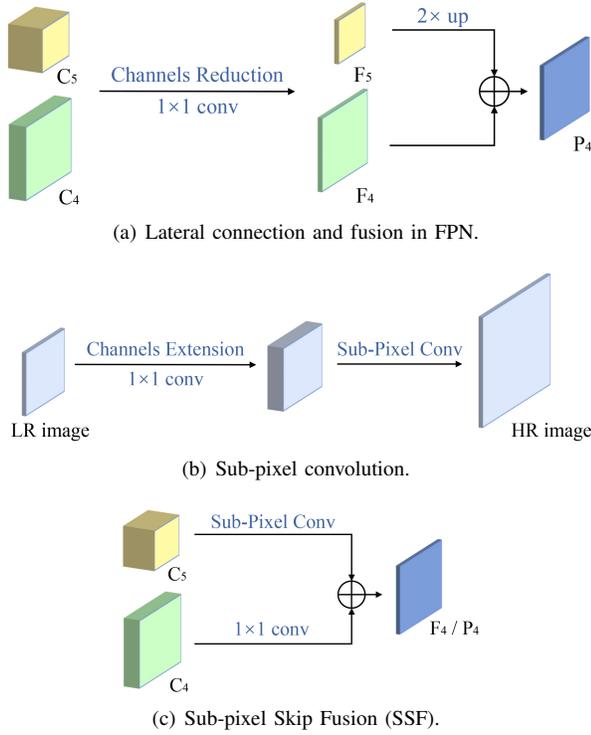


Fig. 3: The design ideas of Sub-pixel Skip Fusion (SSF) as a fusion example of C_5 . (a) In FPN, 1×1 convolution layers are adopted to reduce channel dimensions before fusion, which loses channel information. (b) The pipeline of sub-pixel convolution. Channel dimensions should be extended before upsampling. (c) In SSF, channel dimensions of C_5 would not reduce for upsampling.

on the dimensions of channel. The pixel shuffle operator rearranges the feature of shape $H \times W \times C \cdot r^2$ to $rH \times rW \times C$, which can be mathematically defined as

$$\mathcal{PS}(F)_{x,y,c} = F_{\lfloor x/r \rfloor, \lfloor y/r \rfloor, C \cdot r \cdot \text{mod}(y,r) + C \cdot \text{mod}(x,r) + c}, \quad (1)$$

where r denotes the upscaling factor, F is the input feature, and $\mathcal{PS}(F)_{x,y,c}$ denotes the output feature pixel on coordinates (x, y, c) .

As shown in Fig. 3(b), dimensions of the LR image channel need to be increased first when using sub-pixel convolution as upsampling, which brings about extra computation. And the HR image are unreliable that need additional training. Thus FPN adopts nearest neighbor upsampling for simplicity. Nevertheless, we observe that the amounts of channels in $\{C_4, C_5\}$ (1024, 2048) are sufficient to perform sub-pixel convolution. So we introduce Sub-pixel Skip Fusion (SSF) to upsampling the LR image directly without channel reduction as shown in Fig. 3(c). SSF utilizes the rich channel information of $\{C_4, C_5\}$ and merge them into F_i , which is described as

$$F_i = \begin{cases} \varphi(C_i) + \mathcal{PS}(\bar{\varphi}(C_{i+1})) & i = 3, 4 \\ \varphi(C_i) & i = 2 \end{cases}, \quad (2)$$

where φ denotes 1×1 convolution to reduce channels, and i indicates the index of pyramid levels, $\bar{\varphi}$ denotes channel transformation. And the factor r in sub-pixel convolution is adopted as 2 to double the spatial scale for fusion. $\bar{\varphi}$ adopts

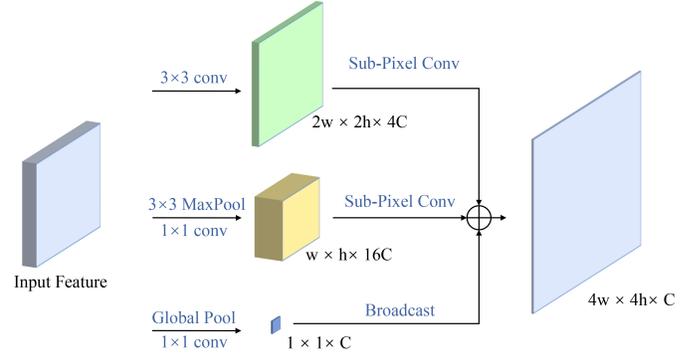


Fig. 4: Illustration of Sub-pixel Context Enhancement (SCE).

1×1 convolution or split operation to change channel dimensions for double sub-pixel upsampling, which is analyzed in Sec IV-D. And if channel dimensions fill the bill, $\bar{\varphi}$ performs identity mapping. Then the feature pyramid P_i are produced by F_i through element-wise summation and nearest-neighbor upsampling, which is the same as in FPN.

As shown in Fig. 2, SSF can be seen as two extra connections from C_5 to F_4 and C_4 to F_3 . SSF performs upsampling and channel fusion simultaneously, which utilizes the rich channel information of high-level features $\{C_4, C_5\}$ to enhance the representation ability of the feature pyramid.

C. Sub-pixel Context Enhancement

On the one hand, feature maps of lower levels are endowed with diverse context information naturally by merging the semantical information from higher levels in conventional FPN. But the highest-level feature only contains single scale context information that is not benefited from others. On the other hand, input images with higher resolution (e.g. shorter sizes of 800 pixels) require neurons with larger receptive fields to obtain more semantical information for capturing large objects [31], [11]. To alleviate the two issues, we adopt the framework of integration map [12] and introduce Sub-pixel Context Enhancement (SCE) to exploit more contextual information with a larger receptive field on C_5 . The extracted context features are merged into the integration map I . SCE follows the design ideas of SSF to utilizes the rich channel information of C_5 .

The key idea of SCE is to fuse large-field local information and global contextual information for generating more discriminative features. We assume that the shape of input feature map C_5 is $2w \times 2h \times 8C$, and the output integration map I is $4w \times 4h \times C$. C is adopted as 256. We perform three scales of context features through parallel pathways shown as Fig. 4.

First, we apply a 3×3 convolution on C_5 to extract local information. Meanwhile, it transforms channel dimensions for sub-pixel upsampling. Then we adopt sub-pixel convolution to perform double scale upsampling, which is similar to SSF.

Second, the input feature is downsampled to $w \times h$ by a 3×3 max-pooling and undergoes a 1×1 convolution layer to extend channel dimensions. Then it follows a $4 \times$ upsampling

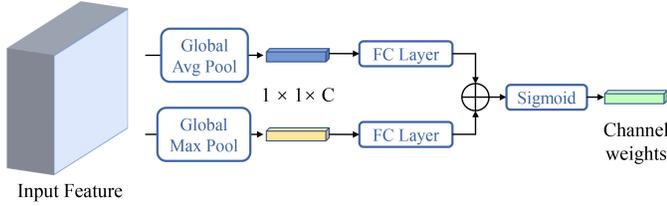


Fig. 5: Illustration of Channel Attention Guided Module (CAG).

sub-pixel convolution. This pathway obtains rich contextual information for larger receptive fields.

Third, we perform a global average pooling on C_5 for global contextual information. Afterward, the global feature of $1 \times 1 \times 8C$ is squeezed to $1 \times 1 \times C$ and broadcasted to the size of $4w \times 4h$. The first and third pathways extract local and global contextual information respectively.

At last, the three generated feature maps are aggregated to the integration map I by element-wise summation. By extending feature representations of three scales, SCE effectively enlarges the receptive field of C_5 and refines the representation ability of I . Therefore, semantical information in the highest-level feature has been fully utilized for FPN. The nodes of F_5 and P_5 are removed for simplicity.

D. Channel Attention Guided Module

There exist semantical differences in cross-scale feature maps, the miscellaneous integrated features may cause aliasing effects to confuse the localization and recognition tasks [6], [11]. In FPN, 3×3 convolution is appended on each merged feature map to generate the final feature pyramid. The proposed SSF and SCE fuse more cross-scale feature maps so that the aliasing effects are more serious than the original FPN. For mitigating the negative impacts of aliasing effects, an intuitive solution is to develop attention modules on the feature pyramid. However, performing independent attention modules on each level of the pyramid is computation-expensive, since some detectors adopt a 6-level pyramid or even more. Meanwhile, we expect that the attention mechanism of different levels can learn from the information of other levels. To this end, we propose a Channel Attention Guided Module (CAG) inspired by CBAM [40], which guides each level of the pyramid to alleviate aliasing effects. CAG extracts channel weights only through the integration map I . And then the channel weights are multiplied to each output feature.

The pipeline of CAG is shown in Fig. 5. We first employ a global average pooling and a global max pooling independently to aggregate two different spatial context information. Next, the two descriptors are forwarded to fully connected layers respectively. Finally, the output feature vectors are merged through element-wise summation and a sigmoid function. The process can be formulated as

$$\mathbf{CA}(x) = \sigma(\mathbf{fc}_1(\text{AvgPool}(x)) + \mathbf{fc}_2(\text{MaxPool}(x))), \quad (3)$$

$$R_i = \mathbf{CA}(I) \odot P_i. \quad (4)$$

where $\mathbf{CA}()$ represents the channel attention function, σ denotes the sigmoid function, and i indicates the index of pyramid levels.

CAG is simply designed to reduce the misleading of aliasing features, rather than sophisticated architecture [8] to enhance more discriminative abilities of features. So lightweight computing is central to our design and $\mathbf{CA}()$ is robust to other channel attention models.

IV. EXPERIMENTS

A. Dataset and Evaluation Metrics

We perform our experiments on MS COCO [18] benchmark. It contains 80 categories and consists of 115k images for training (train-2017) and 5k images for validation (val-2017). There are also 20k images in test-dev that have no released publicly labels. We train all the models on train-2017 and report results of ablation study on val-2017. We submit the final results to the evaluation server of test-dev for comparison. The performance metrics follow the standard COCO-style mean Average Precision (mAP) metrics under different IoU thresholds, ranging from 0.5 to 0.95 with an interval of 0.05.

B. Implementation Details

All of our experiments are implemented based on mmdetection [19]. Mmdetection has been upgraded to v2.0 with higher baseline performance than v1.0. The performance improvement compared with baseline becomes more difficult. Therefore, we re-implement the baseline on mmdetection v2.0 for fair comparisons. We train and test the detectors with the resolution of (1333, 800) on 4 NVIDIA Quadro P5000 GPUs (2 images per GPU). In the training process, $1 \times$ schedule denotes 12 epochs, and 24 epochs for $2 \times$ schedule. The learning rate dropped by 0.1 after 8 and 11 epochs respectively in $1 \times$ schedule, and 16, 22 epochs for $2 \times$ schedule. Our CE-FPN can be applied to any FPN-based detectors. Faster R-CNN [20] and RetinaNet [25] are chosen as the baseline detectors, which represent two-stage and one-stage detectors respectively. FPN [6] and RoIAlign [21] are incorporated into the naive Faster R-CNN to provide a strong baseline. The initial learning rate of Faster R-CNN and RetinaNet is set to 0.01 and 0.005 respectively. The classical networks ResNet-50, ResNet-101 [17] and ResNext101-64x4d [41] are adopted as backbones for comparative experiments. The dimension of the feature pyramid channel is set as 256. And the other settings follow the basic framework if not specifically noted.

C. Main Results

To verify the effectiveness of our method for performance improvement, we evaluate CE-FPN on COCO test-dev subset. For fair comparisons with the corresponding baselines, we report our re-implemented results. As shown in Table I, by replacing FPN with CE-FPN, Faster R-CNN using ResNet-50 and ResNet-101 as backbone achieves 38.8 and 40.9 AP, which is 1.4 and 1.5 points higher than the baseline respectively. When using ResNext101-64x4d backbone, a much more powerful feature extractor, our model achieves 43.1 AP. When $2 \times$

TABLE I: Comparison with baselines and state-of-the-art methods on COCO test-dev. The symbol '*' means our re-implemented results through mmdetection.

Method	Backbone	Schedule	AP	AP_{50}	AP_{75}	AP_S	AP_M	AP_L
Baseline:								
RetinaNet* [25]	ResNet-50	1×	36.3	55.5	38.7	20.5	40.1	47.5
Faster RCNN* [20]	ResNet-50	1×	37.4	58.2	40.4	21.2	40.8	48.1
Faster RCNN* [20]	ResNet-101	1×	39.4	60.2	43.0	22.3	43.3	49.9
Faster RCNN* [20]	ResNet-101	2×	39.9	60.3	43.2	23.0	43.8	52.9
Faster RCNN* [20]	ResNext101-64x4d	1×	41.8	64.4	45.6	24.7	46.1	54.2
State-of-the-art:								
CARAFE [16]	ResNet-50	-	38.1	60.7	41.0	22.8	41.2	46.9
RetinaNet w/ AugFPN [8]	ResNet-50	1×	37.5	58.4	40.1	21.3	40.5	47.3
Faster RCNN w/ AugFPN [8]	ResNet-50	1×	38.8	61.5	42.0	23.3	42.1	47.7
Faster RCNN w/ AugFPN [8]	ResNet-101	1×	40.6	63.2	44.0	24.0	44.1	51.0
Faster RCNN w/ AugFPN [8]	ResNet-101	2×	41.5	63.9	45.1	23.8	44.7	52.8
Faster RCNN w/ AugFPN [8]	ResNext101-64x4d	1×	43.0	65.6	46.9	26.2	46.5	53.9
Libra RetinaNet* [12]	ResNet-50	1×	37.8	57.5	40.5	21.5	40.8	47.4
Libra RCNN* [12]	ResNet-50	1×	38.6	60.0	42.0	22.4	41.3	47.7
Libra RCNN* [12]	ResNet-101	1×	40.2	61.2	44.1	22.7	43.6	52.1
Libra RCNN* [12]	ResNet-101	2×	41.0	62.0	44.7	23.3	43.9	52.6
Libra RCNN* [12]	ResNext101-64x4d	1×	43.0	64.2	46.9	25.2	45.9	54.1
Ours:								
RetinaNet w/ CEFPN	ResNet-50	1×	37.8	57.4	40.1	21.3	40.8	46.8
Faster RCNN w/ CEFPN	ResNet-50	1×	38.8	60.5	41.9	22.5	41.7	48.1
Faster RCNN w/ CEFPN	ResNet-101	1×	40.9	62.5	44.4	23.5	44.2	51.4
Faster RCNN w/ CEFPN	ResNet-101	2×	41.3	62.7	44.8	23.2	44.4	52.7
Faster RCNN w/ CEFPN	ResNext101-64x4d	1×	43.1	64.7	46.9	25.6	46.5	54.0

schedule is adopted, Faster R-CNN with ResNet-101 achieves 39.9 AP due to full training. And the AP boost of CE-FPN still gains up to 1.4. Hence, CE-FPN can consistently bring non-negligible performance improvement even with more powerful backbones. As for RetinaNet, a typical one-stage detector, the performance is boosted to 37.8 AP from 36.3 AP. In addition, from column AP_S , AP_M and AP_L (AP results for small, medium and large objects respectively), we notice that our model brings comprehensive improvement. All improvements indicate that our CE-FPN is effective.

Furthermore, we compare our CE-FPN with other state-of-the-art detectors. Noted that the baseline of mmdetection v2.0 has higher performance than before [42], there seems to be less improvement in scores. So we also re-implement Libra RCNN on mmdetection v2.0 for fair comparisons. Compared with the original report [12], the final performance of our re-implementation results are similar. As shown in Table I, CE-FPN achieves competitive performance compared to state-of-the-art Libra R-CNN and AugFPN.

We also show the comparisons of the qualitative results between FPN and our CE-FPN in Fig. 6. It can be seen that CE-FPN generates satisfactory results for small, medium, and large objects, while the typical FPN generates inferior results. The typical FPN model occasionally misses some objects since these objects may be too small or out of the receptive fields. And it may also localize wrong and aliased objects. Our CE-FPN is more discriminative and performs much better by exploring rich channel information and alleviating aliasing effects. Both models are built upon Faster R-CNN with ResNet-50 and 1× schedule. The images are chosen from COCO val-2017. We compare the detection performance with threshold = 0.5.

TABLE II: Effect of each component on COCO val-2017. **SSF**: Sub-pixel Skip Fusion, **SCE**: Sub-pixel Context Enhancement, **CAG**: Channel Attention Guided Module.

SSF	SCE	CAG	AP	AP_{50}	AP_{75}	AP_S	AP_M	AP_L
			36.1	55.4	38.5	19.8	39.7	47.1
✓			36.6	55.8	39.0	20.9	40.2	47.8
	✓		37.0	56.3	39.3	20.7	40.8	48.1
		✓	37.1	56.2	39.5	21.1	41.2	48.2
✓	✓		37.2	56.7	39.7	21.3	41.3	48.6
✓	✓	✓	37.5	57.3	40.2	21.6	41.2	48.7

D. Ablation Experiments

We also analyze the effect of each proposed component of CE-FPN on COCO val-2017 subset. The overall ablation studies are reported in Table II. We gradually add Sub-pixel Skip Fusion (SSF), Sub-pixel Context Enhancement (SCE), and Channel Attention Guided Module (CAG) on RetinaNet with the backbone of ResNet-50. And the improvements brought by the combination of SSF and SCE are also presented to demonstrate the effectiveness of our sub-pixel based modules. The training process follows 1× schedule (12 epochs). Ablation experiments are implemented with the same settings for fair comparisons.

The effect of Sub-pixel Skip Fusion.

We first implement SSF on RetinaNet without the integration map I . The results show that the naive fusion brings 0.5 points higher AP than the corresponding baseline. When performing SSF independently, the nodes of F_5 and P_5 are preserved. As mentioned before, SSF can be seen as adding two extra connections from C_5 to F_4 and C_4 to F_3 . We also implement the two connections through 1×1 convolution + linear interpolation and compare it with SSF. The simple linear upsampling may cause more serious aliasing effects, which has no sense of performance improvement. Table III proves

TABLE V: Ablation studies of different configurations of channel attention modules on COCO val-2017.

Settings	AP	AP_{50}	AP_{75}	AP_S	AP_M	AP_L
Integration Atten.	36.7	55.9	39.4	20.9	40.7	47.9
Part Atten.	37.0	56.1	39.2	20.6	40.8	48.1
Guided Atten.	37.1	56.2	39.5	21.1	41.2	48.2

TABLE VI: Comparison of FLOPs and parameters in different configurations.

Config.	AP	#GFLOPs	#params.(M)
baseline	36.1	250.34	37.74
SSF_a	36.6	252.54	39.84
SSF_c	36.6	250.34	37.74
SCE w/ $F5, P5$	37.0	272.45	65.54
SCE w/o $F5, P5$	37.0	271.28	65.01
Part Attention	37.0	250.35	37.77
Guided Attention	37.1	250.35	37.75
CE-FPN (ALL)	37.5	271.29	65.02

effectiveness of utilizing rich channel information by sub-pixel convolution.

The effect of Channel Attention Guided Module.

Based on the above methods, diverse cross-scale feature maps are fused into the final feature pyramid. To alleviate the negative impacts of aliasing, CAG takes advantage of the integration map to optimize the channel information in output features. CAG boosts the performance by 1.0 AP.

We also conduct ablation experiments to study different the effects of different attention configurations. The channel attention module based on Eq 3 is adopted in all three configurations. First, we add a self-attention mechanism only on the integration map (Integration Attention). This procedure further refines features to be more discriminative. Second, we perform channel attention on each part of output features (Part Attention). It attempts to eliminate the negative impacts of aliasing effects independently at each level. The third one is our CAG, which extracts channel weights through the integration map and then multiplies to each level (Guided Attention). Table IV and Table VI shows the performance and computational costs of these attention modules. From the results we can observe that CAG is better than other configurations.

Computational Costs.

Table VI demonstrates the computation increase of our proposed methods. It is worthy to note that SSF does not introduce extra computation and parameters. SCE brings a few computational burdens and CAG causes slight computation costs. When adding all components, CE-FPN increases slight computational costs and achieves significant improvement compared to the baseline.

E. Inference speed

We also test the inference speed of CE-FPN. The inference time is the average over COCO val-2017 split. All the runtimes are tested on NVIDIA Quadro P5000. When using ResNet-50 backbone with the input of (1333, 800), Faster R-CNN with CE-FPN can run at 9.8 fps, and the Faster RCNN with FPN can run at 10.5 fps. The inference speed is decreased by about

6.67%. Since the performance of our CE-FPN is similar to that of AugFPN, we compare its inference speed. When replacing FPN with AugFPN in Faster R-CNN, the inference speed has dropped by 17.2% [8]. The comparison result validates the speed superiority of our methods.

V. CONCLUSION

In this paper, we propose a novel channel enhancement feature pyramid network (CE-FPN) with three simple yet effective components to alleviate channel information loss and the aliasing effects. Specifically, we extend the intrinsic upsampling function of sub-pixel convolution to utilize rich channel information in Sub-pixel Skip Fusion (SSF) and Sub-pixel Context Enhancement (SCE). Then we introduce Channel Attention Guided Module (CAG) to alleviate the aliasing effects on each feature level. Our experiments demonstrate that our CE-FPN well generalizes to various FPN-based detectors and achieves significant improvement only with a few computation increases. In future work, we will verify the generalization of CE-FPN on more backbones and other vision tasks with multi-scale.

REFERENCES

- [1] B. Li, J. Yan, W. Wu, Z. Zhu, and X. Hu, "High performance visual tracking with siamese region proposal network," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8971–8980.
- [2] P. Feng, C. Xu, Z. Zhao, F. Liu, J. Guo, C. Yuan, T. Wang, and K. Duan, "A deep features based generative model for visual tracking," *Neurocomputing*, vol. 308, pp. 245–254, 2018.
- [3] C. Yuan, J. Guo, P. Feng, Z. Zhao, Y. Luo, C. Xu, T. Wang, and K. Duan, "Learning deep embedding with mini-cluster loss for person re-identification," *Multimedia Tools and Applications*, vol. 78, no. 15, pp. 21 145–21 166, 2019.
- [4] C. Yuan, J. Guo, P. Feng, Z. Zhao, C. Xu, T. Wang, G. Choe, and K. Duan, "A jointly learned deep embedding for person re-identification," *Neurocomputing*, vol. 330, pp. 127–137, 2019.
- [5] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," in *Proceedings of the European conference on computer vision (ECCV)*, 2016, pp. 21–37.
- [6] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2117–2125.
- [7] Z. Ge, Z. Jie, X. Huang, C. Li, and O. Yoshie, "Delving deep into the imbalance of positive proposals in two-stage object detection," *Neurocomputing*, vol. 425, pp. 107–116, 2021.
- [8] C. Guo, B. Fan, Q. Zhang, S. Xiang, and C. Pan, "Augfpn: Improving multi-scale feature learning for object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 12 595–12 604.
- [9] M. Tan, R. Pang, and Q. V. Le, "Efficientdet: Scalable and efficient object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10 781–10 790.
- [10] Z. Tian, C. Shen, H. Chen, and T. He, "Fcos: Fully convolutional one-stage object detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 9627–9636.
- [11] J. Cao, Q. Chen, J. Guo, and R. Shi, "Attention-guided context feature pyramid network for object detection," *arXiv preprint arXiv:2005.11475*, 2020.
- [12] J. Pang, K. Chen, J. Shi, H. Feng, W. Ouyang, and D. Lin, "Libra r-cnn: Towards balanced learning for object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 821–830.
- [13] S. Liu, L. Qi, H. Qin, J. Shi, and J. Jia, "Path aggregation network for instance segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8759–8768.

- [14] X. Wang, R. Girshick, A. Gupta, and K. He, “Non-local neural networks,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7794–7803.
- [15] W. Shi, J. Caballero, F. Huszár, J. Totz, A. P. Aitken, R. Bishop, D. Rueckert, and Z. Wang, “Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2016, pp. 1874–1883.
- [16] J. Wang, K. Chen, R. Xu, Z. Liu, C. C. Loy, and D. Lin, “Carafe: Content-aware reassembly of features,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 3007–3016.
- [17] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [18] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft coco: Common objects in context,” in *Proceedings of the European conference on computer vision (ECCV)*, 2014, pp. 740–755.
- [19] K. Chen, J. Wang, J. Pang, and et al., “MMDetection: Open mmlab detection toolbox and benchmark,” *arXiv preprint arXiv:1906.07155*, 2019.
- [20] S. Ren, K. He, R. Girshick, and J. Sun, “Faster r-cnn: towards real-time object detection with region proposal networks,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 6, pp. 1137–1149, 2016.
- [21] K. He, G. Gkioxari, P. Dollár, and R. Girshick, “Mask r-cnn,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2017, pp. 2961–2969.
- [22] Z. Cai and N. Vasconcelos, “Cascade r-cnn: Delving into high quality object detection,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 6154–6162.
- [23] Y. Liu, Y. Wang, S. Wang, T. Liang, Q. Zhao, Z. Tang, and H. Ling, “Cbnet: A novel composite backbone network architecture for object detection,” in *Proceedings of the AAAI conference on artificial intelligence*, 2020, pp. 11 653–11 660.
- [24] J. Redmon, S. K. Divvala, R. B. Girshick, and A. Farhadi, “You only look once: Unified, real-time object detection,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 779–788.
- [25] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, “Focal loss for dense object detection,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2017, pp. 2980–2988.
- [26] J. Redmon and A. Farhadi, “YOLO9000: better, faster, stronger,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 6517–6525.
- [27] —, “Yolov3: An incremental improvement,” *arXiv preprint arXiv:1804.02767*, 2018.
- [28] Z. Shen, Z. Liu, J. Li, Y. Jiang, Y. Chen, and X. Xue, “DSOD: learning deeply supervised object detectors from scratch,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2017, pp. 1937–1945.
- [29] H. Law and J. Deng, “Cornernet: Detecting objects as paired keypoints,” in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 734–750.
- [30] J. Guo, C. Yuan, Z. Zhao, P. Feng, Y. Luo, and T. Wang, “Object detector with enriched global context information,” *Multimedia Tools and Applications*, vol. 79, no. 39, pp. 29 551–29 571, 2020.
- [31] Z. Qin, Z. Li, Z. Zhang, Y. Bao, G. Yu, Y. Peng, and J. Sun, “Thundernet: Towards real-time generic object detection on mobile devices,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 6718–6727.
- [32] Q. Zhao, T. Sheng, Y. Wang, Z. Tang, Y. Chen, L. Cai, and H. Ling, “M2det: A single-shot object detector based on multi-level feature pyramid network,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 33, no. 01, 2019, pp. 9259–9266.
- [33] G. Ghiasi, T.-Y. Lin, and Q. V. Le, “Nas-fpn: Learning scalable feature pyramid architecture for object detection,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7036–7045.
- [34] S. Gidaris and N. Komodakis, “Object detection via a multi-region and semantic segmentation-aware cnn model,” in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1134–1142.
- [35] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, “Pyramid scene parsing network,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2881–2890.
- [36] Z. Chen, S. Huang, and D. Tao, “Context refinement for object detection,” in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 71–86.
- [37] F. Wang, M. Jiang, C. Qian, S. Yang, C. Li, H. Zhang, X. Wang, and X. Tang, “Residual attention network for image classification,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 3156–3164.
- [38] Y. Zhu, C. Zhao, H. Guo, J. Wang, X. Zhao, and H. Lu, “Attention couplenet: Fully convolutional attention coupling network for object detection,” *IEEE Transactions on Image Processing*, vol. 28, no. 1, pp. 113–126, 2018.
- [39] H. Li, Y. Liu, W. Ouyang, and X. Wang, “Zoom out-and-in network with map attention decision for region proposal and object detection,” *International Journal of Computer Vision*, vol. 127, no. 3, pp. 225–238, 2019.
- [40] S. Woo, J. Park, J.-Y. Lee, and I. So Kweon, “Cbam: Convolutional block attention module,” in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 3–19.
- [41] S. Xie, R. B. Girshick, P. Dollár, Z. Tu, and K. He, “Aggregated residual transformations for deep neural networks,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2017, pp. 5987–5995.
- [42] J. Ma and B. Chen, “Dual refinement feature pyramid networks for object detection,” *arXiv preprint arXiv:2012.01733*, 2020.