

SDTP: Semantic-aware Decoupled Transformer Pyramid for Dense Image Prediction

Zekun Li Yufan Liu Bing Li Weiming Hu Kebin Wu Pei Wang

(lizekun2018,yufan.liu)@ia.ac.cn, (bli,wmhu)@nlpr.ia.ac.cn, wukebin@huawei.com, wangpei@cert.org.cn

Abstract

Although transformer has achieved great progress on computer vision tasks, the scale variation in dense image prediction is still the key challenge. Few effective multi-scale techniques are applied in transformer and there are two main limitations in the current methods. On one hand, self-attention module in vanilla transformer fails to sufficiently exploit the diversity of semantic information because of its rigid mechanism. On the other hand, it is hard to build attention and interaction among different levels due to the heavy computational burden. To alleviate this problem, we first revisit multi-scale problem in dense prediction, verifying the significance of diverse semantic representation and multi-scale interaction, and exploring the adaptation of transformer to pyramidal structure. Inspired by these findings, we propose a novel Semantic-aware Decoupled Transformer Pyramid (SDTP) for dense image prediction, consisting of Intra-level Semantic Promotion (ISP), Cross-level Decoupled Interaction (CDI) and Attention Refinement Function (ARF). ISP explores the semantic diversity in different receptive space. CDI builds the global attention and interaction among different levels in decoupled space which also solves the problem of heavy computation. Besides, ARF is further added to refine the attention in transformer. Experimental results demonstrate the validity and generality of the proposed method, which outperforms the state-of-the-art by a significant margin in dense image prediction tasks. Furthermore, the proposed components are all plug-and-play, which can be embedded in other methods.

1. Introduction

Inspired by the great success of transformer [25] in natural language processing (NLP) [9], there is an increasing effort on applying it to computer vision. Following vision transformer (ViT) [10], which is the first attempt of applying transformer to vision, plenty of studies [26, 19] have adopted transformer to dense image prediction tasks, such as object detection, semantic segmentation and instance

segmentation.

However, the scale variation in dense image prediction is still the key challenge, even for transformer-based methods. While the pyramidal structure represented by Feature Pyramid Network (FPN) [14] is an effective method to tackle this problem, but few attempts are made to apply multi-scale technique in transformer. On one hand, vanilla transformer fails to explore the diversity of high-level feature’s semantic information, because of the fixed receptive field of patches and the limitations of the self-attention mechanism. On the other hand, establishing attention directly on low-level feature maps, which has large size in general, is infeasible due to quadratic computational cost, let alone building attention and interaction among different levels.

To alleviate this problem, we first revisit the multi-scale problem in dense prediction and obtain several findings. Through the decomposition experiments of FPN, we find that the semantic level (i.e., C_5) plays a significant role in dense prediction, and even the single input of C_5 can achieve comparable performance. Besides the semantic information, multi-level interaction is indispensable, which can promote the mutual learning among different levels and suppress the redundant information in low levels. In addition, we make various attempts to adapt transformer to FPN. The results show that in the proposed decoupled space, transformer saves large computation and achieves high performance.

Based on these findings, we propose a novel Semantic-aware Decoupled Transformer Pyramid (SDTP) for dense image prediction consisting of three components: Intra-level Semantic Promotion (ISP), Cross-level Decoupled Interaction (CDI) and Attention Refinement Function (ARF). ISP exploits the semantic diversity in various receptive space to amply mine the semantic information, integrating local to global information in transformer flexibly. CDI builds the global attention and interaction across different feature levels with the help of proposed decoupled technique, which also solves the problem of heavy computation. Besides, ARF module is further embedded in attention module in transformer to refine the attention map. These three components are all plugged and can be embedded in various

arXiv:2109.08963v1 [cs.CV] 18 Sep 2021

methods.

In summary, this work makes the following contributions:

- The multi-scale problem in dense image prediction is revisited with insightful findings of key factors leading to the success of multi-scale prediction: sufficient semantic information and effective interaction computation. Besides, one simple yet effective solution is presented to reduce the computational cost of multi-scale transformer.
- We propose SDTP to alleviate scale variation problem in dense image prediction, which is an enhanced drop-in replacement of FPN with transformer, for generating more representative multi-scale features.
- SDTP consists of three components: Intra-level Semantic Promotion (ISP), Cross-level Decoupled Interaction (CDI) and Attention Refinement Function (ARF). ISP makes full use of semantic information in various receptive space. CDI builds a global attention and interaction among different levels in decoupled space, and ARF is embedded in transformer block to refine the attention map. Each of these components can be separately utilized in various methods.

2. Related Works

2.1. Dense image prediction task

Object detection. Most detectors can be divided into two types: one-stage detectors (e.g., RetinaNet [15], OneNet [23]) and multi-stage detectors (e.g., Faster R-CNN [21], Cascade R-CNN [1]). Recently, inspired by the excellent performance of transformer in NLP tasks, some studies applied transformer to object detection tasks. DETR [2] is proposed to utilize transformer for end-to-end detection for the first time. Deformable DETR [30] further embeds an effective attention module, leading to better performance than DETR.

Segmentation segmentation. FCN [20] utilizes a fully convolutional network to obtain segmentation maps. Inspired by FCN, U-Net [22] is widely used in medical segmentation, fusing multi-level information for prediction. Besides, the series of DeepLab [17, 6] apply dilated convolution to obtain large receptive fields for more spatial information.

Instance segmentation. By adding a paralleled mask head, Mask R-CNN [11] extends Faster R-CNN for instance segmentation task. To refine the edge of instances, PointRend [13] provides a more sophisticated mask heads with a set of points. In HTC [4], apart from a semantic segmentation branch for contextual information, a cascade structure is used to combine the object detection and instance segmentation for multi-stage prediction. All these dense image

prediction tasks are faced with scale variation problem and need multi-scale features for precise prediction.

2.2. Method for scale variation

Scale variation of object instances is a giant obstacle in dense image prediction tasks and most current studies rely on multi-scale technique. Feature pyramid network (FPN) [14] is a classical structure that contains a top-down pathway to fuse the adjacent features. After that, a series of studies are proposed to further improve FPN. For example, PANet [18] introduces a bottom-up pathway to shorten information path among different levels and FPG [3] utilizes a deep multi-pathway feature pyramid to make the feature fusion in various directions. Different from FPN and its variants, YOLOF [7] is proposed to utilize a single level feature to alleviate the scale variation through multi-receptive field blocks. Therefore, how to deal with multi-scale and multi-receptive field features effectively is the key for scale variation.

2.3. Vision transformer

Recently, a resurgence of work in transformer has led to major advances in vision tasks. ViT [10] constructs a pure transformer backbone for image classification. Following ViT, a series of works were presented. For example, T2T-ViT [28] splits the image into tokens of overlapping patches, to strengthen the interaction of tokens. In addition, the shifted windows for self-attention are proposed in Swin transformer [19], providing local connections among different windows. Besides the methods above, some other methods explore to integrating CNN and transformer, taking advantage of both sides. For instance, PVT [26] applies the transformer to the pyramid structure used in ResNet and CVT [27] introduces the convolutions to vision transformer which combines the preponderance of each others. However, most of previous studies concentrate on building the attention in the same scales of backbone, ignoring multi-scale interaction among different levels. Contrastively, the proposed method designs a multi-scale transformer, which can build long-range relationship among different levels to alleviate scale variation in dense prediction tasks.

3. Revisiting dense multi-scale prediction

Before proposing our multi-scale transformer, in this section, we revisit FPN and transformer, exploring the factors that influence the performance of dense multi-scale prediction. Specifically, we first disassemble FPN and examine the effectiveness of each component. Then the transformer is analyzed to discuss the adaptability on multi-scale integration. Based on these analysis, a number of corresponding findings are obtained to give insight into dense multi-scale detection. Here, the experiments are implemented with Faster R-CNN and RetinaNet based on ResNet-50 in

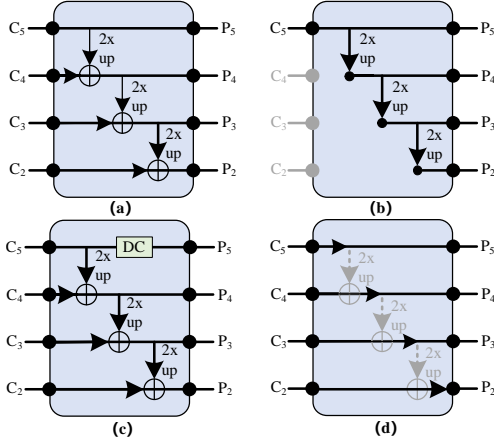


Figure 1: Illustrations of our experiments of FPN. (a) is the FPN baseline. The single input C_5 is acted as an example shown in (b). *DC* denotes the dilated convolution in (c). (d) represents the pyramidal structure without interaction.

object detection task, as shown in Fig 1. Note that $\{C_i\}_{i=2}^5$ represent the output features from different stages with a down-sampling rate of $\{4, 8, 16, 32\}$ and $\{P_i\}_{i=2}^5$ denote the features for final prediction. RetinaNet is implemented without the use of C_2 .

Finding 1: The high semantic level (C_5) plays a significant role on dense multi-scale prediction performance. Single input of C_5 can even achieve a comparable accuracy.

As shown in Fig 1 (b), we select different levels as single input of FPN and keep the multi-level output through the operations of down-sampling and up-sampling, in order to validate the importance of different level input features. The experimental results are reported in Fig 2, which shows that, the single input from all levels has different degrees of performance degradation in Faster R-CNN and RetinaNet likewise has the same phenomenon. Among them, interestingly, model with single C_5 input only has slightly accuracy degradation, which is comparable to the baseline with FPN. This observation suggests that C_5 has sufficient semantic information, which is vital for detection performance. Based on this, we further enhance C_5 by embedding multiple receptive field, which leads to high performance gain. As shown in Fig 1 (c), dilated convolutions with a rate of 3 is utilized to enhance the receptive field of C_5 . Thanks to the extension of receptive field, the model achieves better performance by 0.3 and 0.2 points on Faster R-CNN and RetinaNet. This indicates that adopting multi-receptive field can help semantic level to complementally learn abundant semantic scale information, expanding the diversity of semantic features.

Finding 2: Although FPN with input of single high semantic level achieves good performance, the interaction among multiple levels is indispensable.

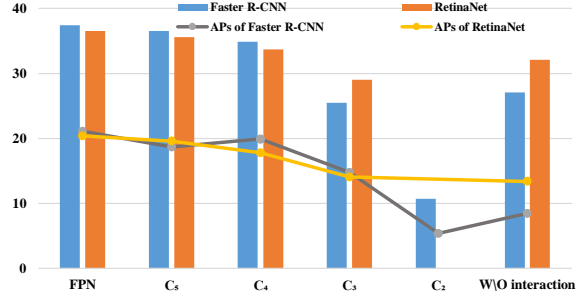


Figure 2: Experimental results in analyzing factors of FPN with Faster R-CNN and RetinaNet on COCO. AP_S denotes the the average precision of small objects.

Table 1: Results of different attempts of applying transformer to FPN.

Method	AP	AP_S	AP_M	AP_L	Flops (G)
baseline	37.4	21.2	41.0	48.1	207.07
p-MSA	-	-	-	-	-
s-MSA	37.3	19.9	41.7	49.7	229.87
d-MSA	37.7	22.1	41.5	48.5	218.46

As presented in Figure 1 (d), we directly remove the interaction among different levels to explore the necessity. The results in Fig 2 show that the detection performance has approximate 11 *mAP* reduction in Faster R-CNN and about 5 *mAP* in RetinaNet. Particularly, the small objects’ performance of both detectors has steep decrease. Furthermore, although single input of C_5 can achieve the comparable result with baseline, the performance of four levels’ input drops dramatically in both detectors without interaction. Thus, these phenomenons indicate that: (1) the interaction among different levels is indispensable to fully exploit the complementary roles of each feature level for dense multi-scale prediction; (2) there exists redundant information in low levels affecting the performance.

Finding 3: Feature interaction in the proposed decoupled space is promising to promote the efficiency of transformer with high performance.

Transformer is good at capturing the global relationship by multi-head self attention (MSA), but there is no suitable method to solve the scale variation problem in dense image prediction tasks. Here, we conduct three attempts to apply transformer to multi-scale detection: (1) primitive MSA (p-MSA): the primitive transformer block is directly applied to different levels of FPN; (2) stride-MSA (s-MSA): in order to reduce the computation, (C_2, C_3, C_4) are down-sampled with strides of (8, 4, 2) to achieve the same size with C_5 , respectively; (3) decoupled MSA (d-MSA): we propose a decoupled style of features ¹ to represent features with low flops. Given an image containing $h \times w$ patches with c chan-

¹We use pooling operator to decouple h and w .

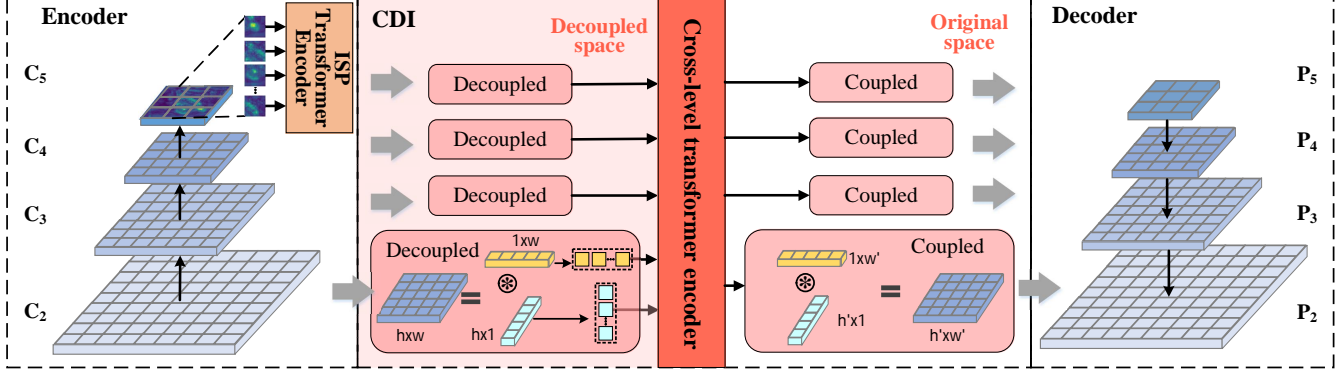


Figure 3: The framework of the proposed SDTP.

nels, the computational complexity of three attempts can be described as:

$$\mathcal{O}(p\text{-MSA}) = \sum_{i=2}^5 (4h_i w_i c_i^2 + 2(h_i w_i)^2 c_i), \quad (1)$$

$$\mathcal{O}(s\text{-MSA}) = \sum_{i=2}^5 (4 \frac{h_i w_i}{s_i^2} c_i^2 + 2(\frac{h_i w_i}{s_i^2})^2 c_i), \quad (2)$$

$$\mathcal{O}(d\text{-MSA}) = \sum_{i=2}^5 (4(h_i + w_i) c_i^2 + 2(h_i^2 + w_i^2) c_i). \quad (3)$$

As shown in Tab 1, it can be seen that the implementation of p-MSA fails due to the huge computation cost. s-MSA and d-MSA can significantly decrease the heavy computation, which verifies the efficiency of s-MSA and d-MSA. However, s-MSA performs worse than d-MSA and d-MSA, which has the smallest computation, even obtains the best performance. This may be because features of low levels are generally of large spatial size and contain redundant information, while our decoupled attempt can achieve computational savings with better performance and remove some redundancy. Thus, the decoupled-style method gives us some enlightenments to design an efficient and effective transformer to enable interaction among representative features.

4. The Proposed Method

According to our above findings, diverse and sufficient semantic information, effective multi-scale interaction and proper feature dimension reduction are all important for applying transformer in a multi-scale way. Hence, we propose a novel and effective method called SDTP as depicted in Fig 3. In the proposed SDTP framework, encoder extracts multi-scale features and decoder transmits the semantic and scale information from high level to low level. Besides, three main components, including ISP, CDI and ARF, are designed to apply multi-scale technique to transformer. In particular, ISP transformer encoder is applied to the high

semantic level (C_5) for exploring intra-level semantic diversity. CDI adopts the decoupled style of features to effectively enable the sufficient interaction among tokens from various levels by cross-level transformer encoder. Besides, ARF is embedded in the attention module of transformer block for precise correlation result.

4.1. Intra-level Semantic Promotion (ISP)

Transformer is proficient in capturing global information relying on its core module self-attention. However, self-attention generally acts on the single feature to build its long range information. This mechanism only focuses on the current features' states, which ignores to explore the feature diversity (e.g., multi-receptive field of one feature). Even though in pyramid structure, high semantic feature has strong representation, but still losses diverse semantic scale detail information. As proved in finding 1, we learn that the semantic level (i.e., C_5) contains useful context which is crucial to performance, and the enhancement of C_5 (i.e., multi-receptive enhancement) can obtain a satisfied performance gain. Thus, we design the ISP transformer encoder to explore semantic diversity information of the high semantic level in receptive space, which can flexibly integrate the local and global information in the transformer.

In particular, as shown in Fig 4, the whole process of ISP can be formulated as:

$$\hat{C}_5 = \text{Attn}_{\text{ISP}}(\text{LN}(C_5)) + C_5, \quad (4)$$

$$C_5^* = \text{MLP}(\text{LN}(\hat{C}_5)) + \hat{C}_5, \quad (5)$$

where $\text{LN}(\cdot)$ denotes the normalization layer and $\text{MLP}(\cdot)$ is multi-layer perception. $\text{Attn}_{\text{ISP}}(\cdot)$ is the key module to enhance the semantic diversity through multi-receptive tokens, which can be written as:

$$\text{Attn}_{\text{ISP}}(C_5) = \text{MMA}[\mathcal{G}(C_5)], \quad (6)$$

in which $\mathcal{G}(\cdot)$ Eq. 7 is the function to generate multi-receptive tokens and $\text{MMA}(\cdot)$ denotes the multi-head

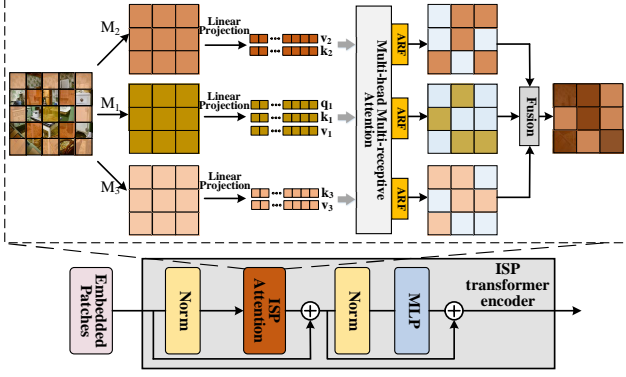


Figure 4: Illustration of ISP transformer encoder. *ARF* represents the attention refinement function module.

multi-receptive attention in Eq. 9. We explain these two components in the following:

$$\mathcal{G}(C_5) = \{q_s, k_s, v_s\}_{s=1}^S = \{Reshape(M_s)\}_{s=1}^S \in \mathbb{R}^{c \times (hw)}, \quad (7)$$

$$s.t. \quad M_s = \sum_{\Phi \in \mathcal{O}_s} \Phi(C_5), \quad s = (1, \dots, S). \quad (8)$$

$M_s \in \mathbb{R}^{c \times h \times w}$ denotes the various states of features with different receptive fields. \mathcal{O}_s represents the operation set, including 3×3 dilated convolution with different rates and the position embedding. S represents the number of dilated rates which we adopt 3 in our method. To keep the original local information without loss, we always include rate of 1. Obtaining query q , key k and value v as shown, we utilize the initial state' query q_1 to explore and search complementary semantic scale information of other states by $MMA(\cdot)$. Details can be formulated as follows:

$$MMA = \mathbb{C} \left\{ \left[Attn(q_1 w_d^q, k_s w_d^k, v_s w_d^v) \right]_{d=1}^D \right\}, \quad s = (1, 2, 3). \quad (9)$$

Here, $\mathbb{C}(\cdot)$ means concatenation and D is the number of heads. $Attn(\cdot)$ computes the token-wise correlation among the inputs. w_j^q, w_j^k, w_j^v are the linear projection parameters. Our design keeps the query q_1 of high priority to avoid the information loss. And we utilize q_1 to search the diverse semantic message in a local-to-global way, which helps mine the diversity and make up the deficiency. Besides, due to the tiny spatial size of C_5 , the computational cost is limited.

In summary, we realize intra-level diversity semantic promotion due to the ISP transformer encoder, which aggregates multi-receptive information in a local-to-global way. ISP makes full use of the superiority and diversity of high semantic level, which can flexibly explore its different scale states to search effective information in receptive space.

4.2. Cross-level Decoupled Interaction (CDI)

Transformer has a dominant performance in NLP task due to capturing long-range relationship. Hence, it is desirable to establish cross-level interaction based on transformer. However, there exist two fatal obstacles. The first one is huge computation cost. In transformer-based vision tasks, we regard image patches as tokens, and the resulted sequence length is generally large. What is worse, comparing with classification task, dense image prediction tasks like object detection need image of higher-resolution to gain precise prediction. Thus, applying transformer on features of high resolution in dense prediction tasks is difficult, let alone the interaction among different levels. Secondly, the popular interaction style in FPN is insufficient and rigid, since non-adjacent levels can not learn from each other. To address these problems, we propose the cross-level decoupled interaction module (CDI).

On the one hand, it reduces the input dimension for transformer in decoupled space, making multi-scale interaction practical and efficient. On the other hand, the proposed decoupled operation effectively decrease the redundancy in original features, which thereafter deepen the multi-scale interaction and improves the performance.

As shown in Fig 3, the proposed CDI decouples one feature map into a the vertical feature and a horizontal one with much smaller dimensions. For clearness, we present the decoupled process in Eq. 10 and Eq. 11.

$$Y_i = \mathcal{F}_{3 \times 1} \left[\sum_j^w \Psi[\mathcal{F}_{1 \times 1}(C_i)] \cdot \sum_j^w (C_i) \right] \in \mathbb{R}^{c \times h \times 1}, \quad (10)$$

$$X_i = \mathcal{F}_{1 \times 3} \left[\sum_j^h \Psi[\mathcal{F}_{1 \times 1}(C_i)] \cdot \sum_j^h (C_i) \right] \in \mathbb{R}^{c \times 1 \times w}, \quad (11)$$

in which Ψ denotes the activation function and \mathcal{F} is the convolution operations with its kernel size shown in the subscript to enhance the learnability. Y_i and X_i are the decoupled features of C_i . Considering the decoupled features from multiple levels as tokens, we leverage them to make a flexible and sufficient interaction in cross-level transformer as shown in Fig 5. Firstly, we can obtain q, k, v from Y_i and X_i :

$$\{q_{ih}, k_{ih}, v_{ih}\}_{i=2}^5 = \{Reshape(Y_i)\}_{i=2}^5 \in \mathbb{R}^{c \times h}, \quad (12)$$

$$\{q_{iw}, k_{iw}, v_{iw}\}_{i=2}^5 = \{Reshape(X_i)\}_{i=2}^5 \in \mathbb{R}^{c \times w}. \quad (13)$$

Then we design a multi-head global attention module (MGA) to make the features learn cross-level knowledge from each other. Note that the MGA for tokens from decoupled vertical and horizontal features are implemented separately. For tokens form X_i , this process can be formulated

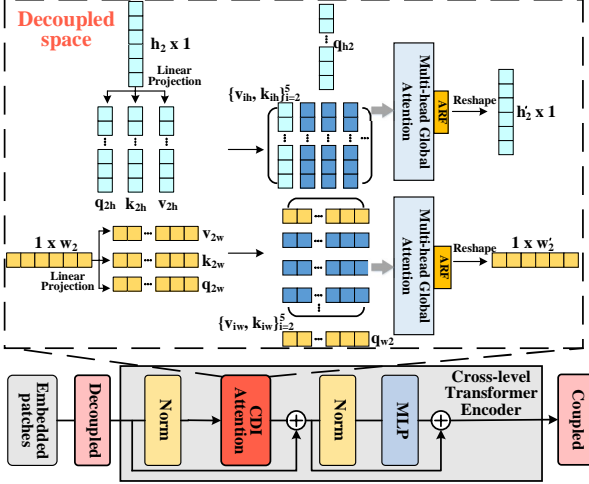


Figure 5: Illustration of cross-level transformer encoder. We use C_2 as an example to show the process.

as below, and the MGA for Y_i is implemented similarly.

$$MGA = \mathbb{C}[Attn(q_i w^q, \{k_i w^k\}_{i=2}^5, \{v_i w^v\}_{i=2}^5)]_{d=1}^D. \quad (14)$$

The output of MGA are \hat{X}_i and \hat{Y}_i with sufficient scale information from all levels owing to the sufficient interaction. Particularly, in contrast to the traditional top-down fusion style in FPN, our design is more flexible since it allows for interactions for features of any two levels. Finally, we make the re-couple of \hat{C}_i for the last process such as MLP in transformer block.

Besides, in order to avoid significant information loss, we introduce an effective loss function for decoupled process to achieve the end-to-end optimization:

$$\mathcal{L}_{dep} = \sum_{i=2}^5 \|C_i - (Y_i \otimes X_i)\|_2. \quad (15)$$

Note that \otimes is the Kronecker sum operator. The total loss can be described as:

$$\mathcal{L} = \mathcal{L}_{org} + \lambda \mathcal{L}_{dep}, \quad (16)$$

in which \mathcal{L}_{org} represent the original loss of dense prediction task. Besides, we leverage λ to balance the loss, which is set to 0.01 in our experiments.

Because of the large size and redundant information in lower levels of dense prediction tasks, applying transformer to capture long-range relationship is impractical. The proposed CDI adopts the decoupled style to obtain more representative features with low dimension, and makes tokens learn adequate cross-level message through the multi-head global attention, which proves to be efficient and effective.

4.3. Attention Refinement Function (ARF)

As the core block in transformer, self-attention module generates an attention map that depicts the correlation between different tokens. The previous studies (e.g., [10, 26]) adopt Softmax as the activation function in correlation calculation. However, Softmax takes account of all tokens together, ignoring the independence and inconsistency of various tokens. In dense prediction tasks, there are enormous tokens and the use of Softmax is likely to alter the correlation map by enhancing some relations and weakening others. However, we desire for a correlation map that can reflect the relations more accurately.

To alleviate this problem, we propose a flexible and effective activation function, which is formulated as:

$$\mathcal{U}(x) = Max\left[\frac{e^x - e^{-x}}{e^x + e^{-(x+2\tau)}}, 0\right], \quad (17)$$

where, τ is a hyper-parameter. Our proposed activation function considers each correlation individually and eliminates the redundant relations, avoiding the limitations of Softmax. Moreover, it is easy to find that Tanh is a special case of this function when the hyper-parameter τ is set to zero and the activation value is positive. Comparing with Tanh, this function can boost values that are close to zero by adjusting the hyper-parameter τ and suppress the redundant relations to zero whose activation values are negative. Rather than using Softmax as activation function, we rethink the calculation of attention in self-attention and propose ARF to refine the correlation result. ARF can be put in other transformer-based methods as a pluggable component to refine the attention.

5. Experiments

5.1. Settings

Data and evaluation. Our experiments are implemented on the MS COCO 2017 [16] for object detection and instance segmentation and ADE20K [29] for semantic segmentation. MS COCO 2017 contains 80 object categories for detection and instance segmentation tasks, which consists of 115k images for training (*train2017*) and 5k images for validation (*val2017*). We train on *train2017*, and report results on *val2017*. Then performance is evaluated by standard COCO-style Average Precision (AP) metrics on small, medium and large objects, i.e., AP_s , AP_m and AP_l . Besides, AP^b and AP^m denote the AP of bounding box and mask. ADE20K is a challenging scene parsing benchmark containing 150 fine-grained semantic categories, consisting of 20k images for training and 2k images for validation. The mean Intersection-over-Union (mIoU) is the primary metric to evaluate the performance of semantic segmentation.

Implementation details. We implement our method based on mmdetection [5]. In order to ensure the fairness of com-

Table 2: **Object Detection:** Performance comparisons with popular detectors. "SDTP" denotes our method. "√" means the baselines integrated with our transformer pyramid.

Method	Backbone	SDTP	AP^b	AP^b_S	AP^b_M	AP^b_L
RetinaNet	R50	√	36.5	20.4	40.3	48.1
	R101	√	38.1	21.8	41.8	49.1
Faster R-CNN	R50	√	38.5	21.7	42.8	50.4
	R101	√	40.0	22.5	44.1	52.1
Cascade R-CNN	R50	√	37.4	21.2	41.0	48.1
	R101	√	39.4	22.7	42.7	51.0
Cascade R-CNN	R50	√	39.4	22.4	43.7	51.1
	R101	√	40.8	23.3	44.9	54.0
Cascade R-CNN	R50	√	40.3	22.5	43.8	52.9
	R101	√	41.7	24.2	45.0	54.9
Cascade R-CNN	R50	√	42.0	23.4	45.8	55.7
	R101	√	43.2	25.3	47.1	57.3

parisons, we also re-implement baseline methods on mmdetection [5]. The results in our re-implementation are generally better than those in the reference papers. Besides, the backbone used in our experiments are pre-trained in ImageNet [8]. If not otherwise specified, all the baseline methods are equipped with FPN and the ablation studies are implemented with Faster R-CNN based on ResNet50. All other hyper-parameters in our work follow the settings in mmdetection.

5.2. Performance

To demonstrate the generality of our method, we implement our framework on three dense prediction tasks: object detection, semantic segmentation and instance segmentation.

Object detection. The results on common object detectors are shown in Tab 2. The proposed SDTP achieves consistent improvement on both single-stage and two-stage detectors. When paired with strong detectors, SDTP still shows the superiority with better performance. In particular, the results of small objects has significant improvement due to the effective utilization of high-resolution features because of the sufficient interaction among multiply levels.

Semantic segmentation. We also conduct experiments to prove the effectiveness of SDTP on semantic segmentation task. As shown in Tab 3, we compare our method with Semantic FPN [12] and PointRend. With the help of SDTP, we again outperform the baselines. Especially, the $mIoU$ of PointRend with ResNet50 increases nearly 4 points with SDTP. Since the semantic segmentation is more sensitive to multi-resolution information, the performance margin brought by SDTP is remarkable.

Instance segmentation. We continue to validate the generalization of SDTP on instance segmentation task as shown in Tab 5. Our method improves the baseline on both detection and instance segmentation task with great margin. Even with the strong method such as HTC, SDTP still has significant increase by 1.5 points. Besides, benefiting from the

Table 3: **Semantic Segmentation:** Performance comparisons with common semantic segmentation methods of different backbones on ADE20K validation set.

Method	Backbone	SDTP	$mIoU$	$mAcc$	$aAcc$
Semantic FPN	R50	√	37.48	47.57	78.02
	R101	√	38.77	49.35	79.13
PointRend	R50	√	39.35	49.43	79.19
	R101	√	41.52	51.68	80.13
PointRend	R50	√	37.63	48.14	77.80
	R101	√	41.53	51.82	79.56
PointRend	R50	√	40.01	50.56	79.09
	R101	√	42.39	53.44	80.35

diverse semantic information, the performance on large object in instance segmentation has witnessed dominant promotion.

Comparison on transformer-based method. Apart from the above experiments, we further assess the superiority and generality of SDTP on transformer-based backbone. As seen in Tab 6, we apply our method to Mask R-CNN based on backbones of two versions of PVT. SDTP still obtains better performance on two dense prediction tasks.

Comparison with State-of-the-art methods. All shown in Tab 4, SDTP achieves consistently non-negligible improvements even with more powerful backbones and more training epochs. For example, when applying ResNeXt101-32×4d and ResNeXt101-64×4d as the feature extractors of Faster R-CNN, our SDTP still improves performance by 1.1 and 0.9 points, respectively. Besides, in the comparison between HTC and its SDTP variant which are both trained with 20 epochs, SDTP still wins the comparisons. SDTP brings consistent improvements on various backbones, methods and learning schedules. This proves the generalization and superiority of SDTP.

5.3. Ablation studies

Ablation studies on each component. To analyze the importance of each module in SDTP, we gradually apply them to the model. As in Tab 7, all three parts are essential. Especially, we can find that ISP alone increases 1.3 points. The adding of CDI helps to give a boost of 1.0 point. These results are in consistency with our first two findings that both semantic diversity and interaction across multi-level features are the core of multi-scale technique. Furthermore, the proposed ARF provides more accurate attention map in building the intra-level and cross-level transformer, and exploiting ARF enables the performance grow to 39.4.

Ablation studies on different dilated rates in ISP. The results illustrated in Tab 8 show that exploring receptive space brings improvements to SDTP. In detail, the combination of 1,3,6 has the best performance. Another observation is that the performances of the first two rows are lower, which indicates that the selected dilation rates should be of large

Table 4: **Comparisons with the state-of-the-art methods:** The symbol “*” means our re-implemented results on mmdetection.

Method	Backbone	Schedule	AP	AP_{50}	AP_{75}	AP_S	AP_M	AP_L
Faster R-CNN*	ResNeXt101-32×4d	12	41.2	62.1	45.1	24.0	45.5	53.5
Faster R-CNN*	ResNeXt101-64×4d	12	42.1	63.0	46.3	24.8	46.2	55.3
Mask R-CNN*	ResNeXt101-32×4d	12	41.9	62.5	45.9	24.4	46.3	54.0
Cascade R-CNN*	ResNeXt101-32×4d	12	43.7	62.3	47.7	25.1	47.6	57.3
DETR[2]	ResNet50	500	42.0	62.4	44.2	20.5	45.8	61.1
DETR[2]	ResNet101	500	43.5	63.8	46.4	21.9	48.0	61.8
Deformable DETR[30]	ResNet50	50	43.8	62.6	47.7	26.4	47.1	58.0
Sparse R-CNN[24]	ResNet101	36	44.1	62.1	47.2	26.1	46.3	59.7
Cascade Mask R-CNN*	ResNeXt101-32×4d	20	45.0	63.2	49.1	26.7	48.9	59.0
HTC*	ResNet101	20	44.8	63.3	48.8	25.7	48.5	60.2
SDTP Faster R-CNN(ours)	ResNeXt101-32×4d	12	42.3	63.7	46.1	25.3	46.6	54.8
SDTP Faster R-CNN(ours)	ResNeXt101-64×4d	12	43.0	63.9	46.6	25.3	46.9	55.8
SDTP Mask R-CNN(ours)	ResNeXt101-32×4d	12	43.2	64.3	47.1	25.9	47.1	56.5
SDTP Cascade R-CNN(ours)	ResNeXt101-32×4d	12	44.6	63.8	48.6	26.1	48.7	57.8
SDTP Sparse R-CNN(ours)	ResNet101	36	44.6	63.1	48.3	27.3	47.6	60.0
SDTP Cascade Mask R-CNN (ours)	ResNeXt101-32×4d	20	45.7	64.5	49.8	27.1	49.1	59.6
SDTP HTC* (ours)	ResNet101	20	45.8	64.9	49.8	27.3	49.3	60.5

Table 5: **Instance Segmentation:** Performance comparisons with strong instance segmentation methods.

Method	Backbone	SDTP	AP^b	AP_S^b	AP^m	AP_L^m
Mask R-CNN	R50	✓	38.2	21.9	34.7	47.2
	R101	✓	40.0	22.8	36.2	53.1
PointRend	R50	✓	40.0	22.6	36.1	49.5
	R101	✓	41.6	24.5	37.2	54.7
HTC	R50	✓	38.4	22.8	36.3	48.5
	R101	✓	40.9	25.5	38.0	50.8
HTC	R50	✓	42.3	23.7	37.4	51.7
	R101	✓	43.8	25.7	38.9	57.2

Table 6: **Comparison with transformer-based backbone:** Performance comparisons paired with Mask R-CNN.

Method	Backbone	SDTP	AP^b	AP_S^b	AP^m	AP_L^m
Mask R-CNN	PVT-Tiny	✓	36.7	21.6	35.1	48.5
	PVT-Small	✓	38.3	23.5	36.2	54.2
Mask R-CNN	PVT-Small	✓	40.4	22.9	37.8	53.6
	PVT-Small	✓	41.4	23.3	38.3	58.0

Table 7: **Effectiveness of each component.**

Method	AP	AP_S	AP_M	AP_L
baseline	37.4	21.2	41.0	48.1
baseline+ISP	38.7	22.3	42.5	50.8
baseline++CDI	38.4	21.3	42.0	49.9
baseline+ISP+CDI	39.0	22.6	42.6	50.7
baseline+ISP+CDI+ARF	39.4	22.7	42.7	51.0

difference so that the semantic features can be diverse.

Ablation studies on τ in ARF. Experimental results related with the settings of ARF are presented in Tab 9. Paired with our proposed function, the suppression of irrelevant infor-

Table 8: **Different settings of dilated rates in ISP.**

Rates	AP	AP_S	AP_M	AP_L
1,2,3	38.1	21.5	41.6	49.6
1,2,4	37.9	21.4	41.2	49.1
1,3,6	38.7	22.3	42.5	50.8
2,4,6	38.4	21.9	42.2	49.1
3,6,12	38.3	21.4	41.6	50.0

mation helps for better performance. We choose 2 for τ which achieves best performance.

Table 9: **Comparisons among different τ of ARF.**

Method	τ	AP	AP_S	AP_M	AP_L
Softmax	-	39.0	22.6	42.6	50.7
Tanh	-	39.1	22.4	42.6	51.0
Proposed	1	39.0	22.1	42.5	50.9
	2	39.4	22.7	42.7	51.0
	3	39.1	22.4	42.2	50.8
	4	39.2	22.6	42.5	51.4

6. Conclusion

In this paper, we focus on dealing with the scale variation problem in dense image prediction tasks with the aid of multi-scale and transformer techniques. To begin with, we revisit the dense multi-scale prediction, and obtain important insights that semantic diversity and interaction among different levels are the key elements. Based on these findings, we propose a novel semantic-aware decoupled transformer pyramid, which includes three simple yet effective components, i.e., Intra-level Semantic Promo-

tion, Cross-level Decoupled Interaction and Attention Refinement Function. SDTP has shown the generality and effectiveness on various dense image prediction tasks. Additionally, SDTP and its three key components can be easily extended to other methods.

References

- [1] 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*, pages 6154–6162, 2018.
- [2] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European Conference on Computer Vision*, pages 213–229. Springer, 2020.
- [3] Kai Chen, Yuhang Cao, Chen Change Loy, Dahua Lin, and Christoph Feichtenhofer. Feature pyramid grids. *arXiv preprint arXiv:2004.03580*, 2020.
- [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 *IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [5] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jiarui Xu, et al. Mmdetection: Open mmlab detection toolbox and benchmark. *arXiv preprint arXiv:1906.07155*, 2019.
- [6] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2017.
- [7] Qiang Chen, Yingming Wang, Tong Yang, Xiangyu Zhang, Jian Cheng, and Jian Sun. You only look one-level feature. *arXiv preprint arXiv:2103.09460*, 2021.
- [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.
- [9] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [11] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017.
- [12] Alexander Kirillov, Ross Girshick, Kaiming He, and Piotr Dollár. Panoptic feature pyramid networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6399–6408, 2019.
- [13] Alexander Kirillov, Yuxin Wu, Kaiming He, and Ross Girshick. Pointrend: Image segmentation as rendering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9799–9808, 2020.
- [14] 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.
- [15] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.
- [16] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, pages 740–755, 2014.
- [17] Chenxi Liu, Liang-Chieh Chen, Florian Schroff, Hartwig Adam, Wei Hua, Alan L Yuille, and Li Fei-Fei. Auto-deeplab: Hierarchical neural architecture search for semantic image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 82–92, 2019.
- [18] Shu Liu, Lu Qi, Haifang Qin, Jianping Shi, and Jiaya Jia. Path aggregation network for instance segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8759–8768, 2018.
- [19] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. *arXiv preprint arXiv:2103.14030*, 2021.
- [20] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on*

computer vision and pattern recognition, pages 3431–3440, 2015.

- [21] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015.
- [22] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015.
- [23] Peize Sun, Yi Jiang, Enze Xie, Zehuan Yuan, Changhu Wang, and Ping Luo. Onenet: Towards end-to-end one-stage object detection. *arXiv preprint arXiv:2012.05780*, 2020.
- [24] Peize Sun, Rufeng Zhang, Yi Jiang, Tao Kong, Chenfeng Xu, Wei Zhan, Masayoshi Tomizuka, Lei Li, Zehuan Yuan, Changhu Wang, et al. Sparse r-cnn: End-to-end object detection with learnable proposals. *arXiv preprint arXiv:2011.12450*, 2020.
- [25] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.
- [26] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. *arXiv preprint arXiv:2102.12122*, 2021.
- [27] Haiping Wu, Bin Xiao, Noel Codella, Mengchen Liu, Xiyang Dai, Lu Yuan, and Lei Zhang. Cvt: Introducing convolutions to vision transformers. *arXiv preprint arXiv:2103.15808*, 2021.
- [28] Ze Yang, Shaohui Liu, Han Hu, Liwei Wang, and Stephen Lin. Reppoints: Point set representation for object detection. *arXiv preprint arXiv:1904.11490*, 2019.
- [29] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 633–641, 2017.
- [30] 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.