

# ISTR: End-to-End Instance Segmentation with Transformers

Jie Hu<sup>1</sup>, Liujuan Cao<sup>1</sup>, Yao Lu<sup>1</sup>, ShengChuan Zhang<sup>1</sup>, Yan Wang<sup>2</sup>,  
Ke Li<sup>3</sup>, Feiyue Huang<sup>3</sup>, Ling Shao<sup>4</sup>, and Rongrong Ji<sup>1,5</sup>.

<sup>1</sup>Media Analytics and Computing Lab, Department of Artificial Intelligence,  
School of Informatics, Xiamen University, <sup>2</sup>Pinterest, <sup>3</sup>Tencent Youtu Lab,

<sup>4</sup>Inception Institute of Artificial Intelligence, <sup>5</sup>Institute of Artificial Intelligence, Xiamen University.

## Abstract

*End-to-end paradigms significantly improve the accuracy of various deep-learning-based computer vision models. To this end, tasks like object detection have been upgraded by replacing non-end-to-end components, such as removing non-maximum suppression by training with a set loss based on bipartite matching. However, such an upgrade is not applicable to instance segmentation, due to its significantly higher output dimensions compared to object detection. In this paper, we propose an instance segmentation Transformer, termed ISTR, which is the first end-to-end framework of its kind. ISTR predicts low-dimensional mask embeddings, and matches them with ground truth mask embeddings for the set loss. Besides, ISTR concurrently conducts detection and segmentation with a recurrent refinement strategy, which provides a new way to achieve instance segmentation compared to the existing top-down and bottom-up frameworks. Benefiting from the proposed end-to-end mechanism, ISTR demonstrates state-of-the-art performance even with approximation-based suboptimal embeddings. Specifically, ISTR obtains a 46.8/38.6 box/mask AP using ResNet50-FPN, and a 48.1/39.9 box/mask AP using ResNet101-FPN, on the MS COCO dataset. Quantitative and qualitative results reveal the promising potential of ISTR as a solid baseline for instance-level recognition. Code has been made available at: <https://github.com/hujiecpp/ISTR>.*

## 1. Introduction

A growing trend in the recent development of computer vision is to remove the handcrafted components to enable end-to-end training and inference, which has demonstrated significant improvement in multiple fields. However, this end-to-end paradigm still lacks applications for instance segmentation that aims to jointly detect and segment each object in an image. Existing instance segmentation approaches either need a manually-designed post-processing

step called non-maximum suppression (NMS) to remove duplicate predictions [16, 28, 19, 4, 22, 44], or are early trials on small datasets and lack evaluation against modern baselines [31, 33]. Popular approaches also rely on a top-down or bottom-up framework that decomposes instance segmentation into several dependent tasks, preventing them from being end-to-end.

Besides instance segmentation, object detection also faces similar challenges. Recent studies enable end-to-end object detection by introducing a set prediction loss [18, 3, 37, 43, 56], with optional use of Transformers [42]. The set prediction loss enforces bipartite matching between labels and predictions to penalize redundant outputs, thus avoiding NMS during inference. However, enabling end-to-end instance segmentation is not as trivial as adding a mask branch and changing the loss. We conducted a proof-of-concept experiment by adapting the end-to-end object detection approach to instance segmentation. The results in Table 1a show the inferior performance of doing so.

We argue that the reason behind the failure is the insufficient number of samples for learning the mask head. On the one hand, the dimensions of masks are much higher than those of classes and boxes. For example, a mask usually has a  $28 \times 28$  or higher resolution on the COCO dataset, while a bounding box only needs two coordinates to represent. Therefore, the mask head requires more samples for training. On the other hand, the proposal bounding boxes obtained by the bipartite matching are usually on a small scale, which also raises the problem of sparse training samples. For example, Mask R-CNN [16] uses 512 proposal bounding boxes to extract the region of interest (RoI) features for training the mask head, while the number of proposal bounding boxes, *i.e.*, ground truths per image on the COCO dataset, is only 7.7 on average after bipartite matching. The gap between the demand and supply of training samples makes the approach prone to failing.

While we could blindly augment the ground truth samples to alleviate the problem at the cost of longer training time, we argue that there might be a smarter way. Not all

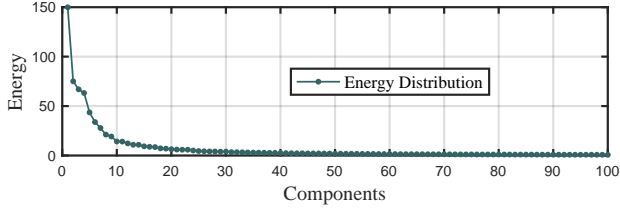


Figure 1: **Component analysis** of masks, ranking the Top@100 components by energy. The majority of the mask information is embedded into the first few components.

$28 \times 28$  entries are likely to appear as a mask. The distribution of the natural masks may lie in a low-dimensional manifold instead of being uniformly scattered. Based on this intuition, we carried out several dimension reduction experiments on the masks from the training data, and surprisingly found even linear methods, such as principal component analysis (PCA), can do a decent job. The energy distribution of different components is shown in Fig. 1. We observe that the first few components can represent the majority of the mask information. Therefore, in this paper, we propose to achieve end-to-end instance segmentation by regressing low-dimensional embeddings instead of raw masks, which enables the training to be effectively conducted with a small number of matched samples. We also extend the definition of bipartite matching cost based on the mask embeddings. Furthermore, regressing with the embeddings enables us to design a recurrent refinement strategy that can process detection and segmentation concurrently. This provides a new way of instance segmentation compared to the top-down and bottom-up frameworks, and boosts the performance.

Specifically, we propose a new end-to-end instance segmentation framework built upon a Transformer, termed ISTR. ISTR predicts low-dimensional mask embeddings, and then matches them with ground truth mask embeddings for the set loss. With the recurrent refinement strategy, ISTR updates the query boxes and refines the set of predictions. Benefiting from the proposed end-to-end mechanism, we find that even with the suboptimal mask embeddings obtained by the closed-form solution of PCA, ISTR can achieve state-of-the-art performance. With a single 1080Ti GPU, ISTR obtains a 46.8/38.6 box/mask AP with 13.8 fps using ResNet50-FPN, and a 48.1/39.9 box/mask AP with 11.0 fps using ResNet101-FPN on the `test-dev` split of the COCO dataset [26]. Our contributions are summarized as follows:

- We propose a new framework, termed instance segmentation Transformer (ISTR). For the first time, we demonstrate the potential of using Transformers in end-to-end instance segmentation.
- ISTR predicts low-dimensional mask embeddings instead of high-dimensional masks, which facilitates the

training with a small number of samples and inspires the design of a bipartite matching cost for masks.

- With a recurrent refinement strategy, ISTR concurrently detects and segments instances, providing a new perspective for achieving instance segmentation compared to the bottom-up and top-down frameworks.
- Without bells and whistles, ISTR demonstrates accuracy and run-time performance on par with the state-of-the-art methods on the challenging COCO dataset.

## 2. Related Work

**Instance Segmentation:** Instance segmentation requires instance-level and pixel-level predictions. Existing works can be summarized into three categories. Top-down methods [1, 16, 23, 28, 19, 7, 4, 22, 48] first detect and then segment the objects. Bottom-up methods [52, 8, 27, 12] view instance segmentation as a label-then-cluster problem, learning to classify each pixel and then clustering them into groups for each object. The latest work, SOLO [44, 46], deals with instance segmentation without dependence on box detection. The proposed ISTR provides a new perspective that directly predicts a set of bounding boxes and mask embeddings, which avoids decomposing instance segmentation into dependent tasks. Note that the idea of regressing mask embeddings is also investigated in MEInst [53]. However, with redundant predictions in each pixel, MEInst obtains suboptimal performance compared to ISTR.

**End-to-End Instance-Level Recognition:** Recent studies have revealed the great potential of end-to-end object detection [18, 3, 37, 43, 56, 55]. As such, the bipartite matching cost has become an essential component for achieving end-to-end object detection. For instance segmentation, the works [31, 33] explored the end-to-end mechanism with recurrent neural networks. However, these early trials were only evaluated on small datasets and not against current baselines. In contrast, ISTR uses the similarity metric of mask embeddings as the bipartite matching cost for masks, and, for the first time, incorporates Transformers [42] to improve end-to-end instance segmentation.

**Transformers in Computer Vision:** The breakthroughs of Transformers [42] in natural language processing have sparked great interest in the computer vision community. The critical component of Transformer is the multi-head attention, which can significantly enhance the capacity of models [14, 20]. So far, Transformers have been successfully used for image recognition [11, 41], object detection [3, 56, 37], segmentation [54, 51], image super-resolution [49], video understanding [36, 13], image generation [5, 45], visual question answering [38, 35], and several other tasks [21, 10, 50]. With the sequence information between frames, the contemporary work [47] achieves end-to-end video instance segmentation with a Transformer. With-

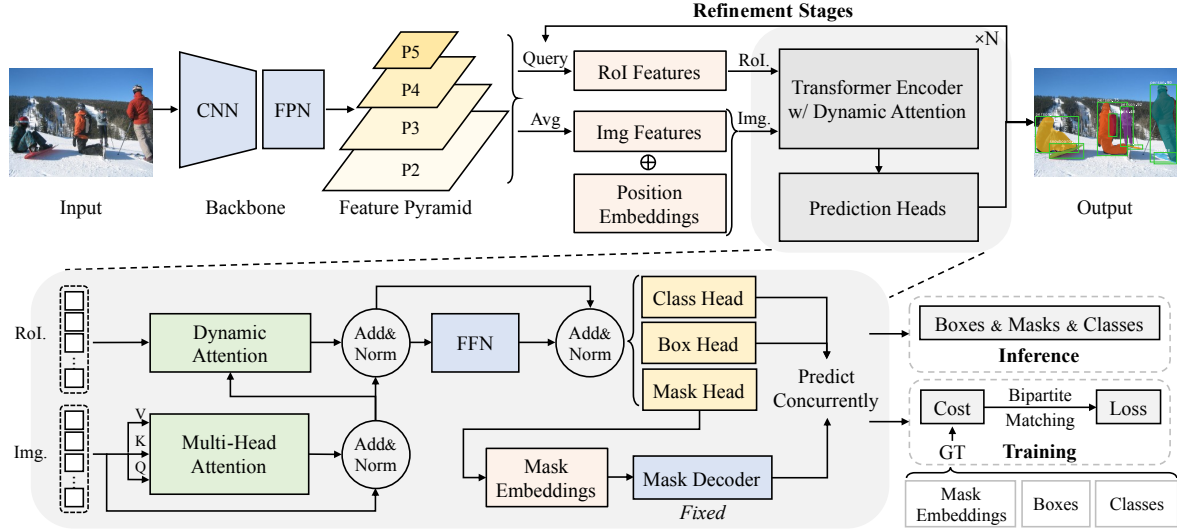


Figure 2: **Framework** of ISTR. Top: Overview of the pipeline. Input images are sent to a convolutional neural network (CNN) with a feature pyramid network (FPN) [24] to produce the feature pyramid. The feature maps from the feature pyramid are cropped and aligned by learnable query boxes with RoIAlign [16] to get the RoI features. Image features are obtained by summing and averaging the feature maps. Then, a Transformer encoder with dynamic attention fuses the image and RoI features for prediction heads. The predicted bounding boxes, classes, and masks are recurrently refined in  $N$  stages by updating the query boxes. During training, the predictions are matched with ground truth labels to calculate the set loss. During inference, the predictions are directly used as the final results without NMS. Bottom: Details of the modules. FFN denotes the feed-forward network, and the mask decoder is pre-learned and fixed during training.

out continuous frames, our study aims to segment instances for a single image, making its design entirely different from the framework of [47].

**Multi-Task Learning:** The benefit of learning detection and segmentation jointly was first studied in the work of [15]. After that, Mask R-CNN [16] also demonstrated that bounding box detection could benefit from multi-task learning. Recent works [6, 28, 34, 2] have provided more complex mechanisms to improve the performance for multi-tasks. In our work, we also observe a performance boost when concurrently processing detection and segmentation.

### 3. Proposed Method

ISTR aims to directly predict a set of mask embeddings, classes, and bounding boxes for each instance. To this end, we first introduce a generalized formulation to extract embeddings for representing and reconstructing the masks in Section 3.1. A bipartite matching cost and a set loss are introduced in Section 3.2 to pair and regress the predictions with ground truth labels. Finally, a model that predicts a set of outputs and learns their relations is proposed in Section 3.3. The overall framework of ISTR is shown in Fig. 2.

#### 3.1. Mask Embeddings

To provide a formulation that effectively extracts mask embeddings, we constrain the mutual information between the original and reconstructed masks:

$$\max \mathcal{I}(\mathbf{M}, f(g(\mathbf{M}))), \quad (1)$$

where  $\mathcal{I}(\cdot, \cdot)$  denotes the mutual information between two random variables,  $\mathbf{M}$  denotes a set of masks  $\{\mathbf{m}_i \in \mathbb{R}^{s^2} | i = 1, \dots, n\}$ ,  $s^2$  is the dimension of masks,  $g(\cdot)$  denotes the mask encoder for extracting embeddings and  $f(\cdot)$  denotes the mask decoder for reconstructing masks. Eq. 1 guarantees that the encoding and decoding phases have minimal information loss, which implicitly encourages the embeddings to represent the masks. After derivation, we have a generalized objective function for the mask embeddings:

$$\min \sum_{i=1}^n \|\mathbf{m}_i - f(\mathbf{r}_i)\|_2^2, \quad (2)$$

where  $\mathbf{r}_i = g(\mathbf{m}_i)$  denotes the mask embeddings, and  $\|\cdot\|_2$  is the L2-norm. By making the functions of the encoder and decoder simple linear transformations via a matrix  $\mathbf{D} \in \mathbb{R}^{s^2 \times l}$ , i.e.,  $f(g(\mathbf{m}_i)) = \mathbf{D}\mathbf{D}^T \mathbf{m}_i$  and  $\mathbf{D}\mathbf{D}^T = \mathbf{I}_l$ , the

objective function becomes:

$$D^* = \arg \min_D \sum_{i=1}^n \|\mathbf{m}_i - DD^T \mathbf{m}_i\|_2^2, \quad (3)$$

where  $l$  is the dimension of the mask embeddings, and  $I_l$  denotes the  $l \times l$  unit matrix. Eq. 3 has the same formulation as the objective function of PCA, which provides a closed-form solution to learn the transformation. Note that the objective function in Eq. 2 can also be optimized by other models, such as an autoencoder.

### 3.2. Matching Cost and Prediction Loss

After obtaining the encoder and decoder for mask embeddings, we define a bipartite matching cost and a set prediction loss for end-to-end instance segmentation. Let us denote the ground truth bounding boxes, classes, and masks as  $\mathbf{Y} = \{\mathbf{b}_i, \mathbf{c}_i, \mathbf{m}_i | i = 1, \dots, n\}$ . The predicted bounding boxes, classes, and mask embeddings are denoted as  $\tilde{\mathbf{Y}} = \{\tilde{\mathbf{b}}_i, \tilde{\mathbf{c}}_i, \tilde{\mathbf{r}}_i | i = 1, \dots, k\}$ , where  $k > n$ .

**Bipartite Matching Cost:** For the bipartite matching, we search for a permutation of  $n$  non-repeating integers  $\sigma \in \{1, 2, \dots, k\}$  with the lowest cost, as:

$$\sigma^* = \arg \min_{\sigma} \sum_{i=1}^n (\mathcal{C}_{box}(\mathbf{b}_i, \tilde{\mathbf{b}}_{\sigma(i)}) + \mathcal{C}_{cls}(\mathbf{c}_i, \tilde{\mathbf{c}}_{\sigma(i)}) + \mathcal{C}_{mask}(\mathbf{m}_i, \tilde{\mathbf{r}}_{\sigma(i)}). \quad (4)$$

Inspired by [3, 37], we define the matching cost for bounding boxes as:

$$\mathcal{C}_{box}(\mathbf{b}_i, \tilde{\mathbf{b}}_{\sigma(i)}) = \lambda_{L1} \cdot \mathcal{C}_{L1}(\mathbf{b}_i, \tilde{\mathbf{b}}_{\sigma(i)}) + \lambda_{giou} \cdot \mathcal{C}_{giou}(\mathbf{b}_i, \tilde{\mathbf{b}}_{\sigma(i)}), \quad (5)$$

and the matching cost for classes as:

$$\mathcal{C}_{cls}(\mathbf{c}_i, \tilde{\mathbf{c}}_{\sigma(i)}) = -\lambda_{cls} \cdot \tilde{p}_{\sigma(i)}(\mathbf{c}_i), \quad (6)$$

where  $\lambda$  denotes the hyperparameters that balance the costs,  $\mathcal{C}_{L1}(\cdot, \cdot)$  denotes the L1 cost,  $\mathcal{C}_{giou}(\cdot, \cdot)$  denotes the generalized IoU [32] cost, and  $\tilde{p}_{\sigma(i)}(\mathbf{c}_i)$  is the probability of class  $\mathbf{c}_i$ . Instead of directly matching the high-dimensional masks, we use the similarity metric between mask embeddings to match them, which is defined as:

$$\mathcal{C}_{mask}(\mathbf{m}_i, \tilde{\mathbf{r}}_{\sigma(i)}) = -\frac{1}{2} \lambda_{mask} \cdot \left( \frac{\tilde{\mathbf{r}}_{\sigma(i)} \cdot \mathbf{g}(\mathbf{m}_i)}{\|\tilde{\mathbf{r}}_{\sigma(i)}\|_2 \cdot \|\mathbf{g}(\mathbf{m}_i)\|_2} + 1 \right), \quad (7)$$

where the mask embeddings are L2 normalized, and the dot product between two normalized vectors is used to calculate the cosine similarity. We add 1 to the result and divide it by 2 to guarantee that the values are in the range of [0, 1].

---

### Algorithm 1 Instance Segmentation Transformer

---

// ——— Training Phase ———

**Input:** Images and ground truth labels.

**Output:** Learned ISTR model.

1: Learn the mask encoder and decoder via Eq. 2.

2: Initialize the learnable query boxes  $\tilde{\mathbf{B}}^0$ .

3: **Repeat**

4: **For**  $i = 1, 2, \dots, N$  stage:

5: Obtain RoI features by RoIAlign with  $\tilde{\mathbf{B}}^{i-1}$ .

6: Predict via ISTR encoder and heads.

7: Match predictions with labels via Eq. 4.

8: Calculate loss via Eq. 8 and train ISTR.

9: Update  $\tilde{\mathbf{B}}^i$  from  $\tilde{\mathbf{B}}^{i-1}$  with the predicted boxes.

10: **Until** scheduled epochs.

// ——— Inference Phase ———

**Input:** Images to be processed.

**Output:** Detected and segmented objects.

1: **For**  $i = 1, 2, \dots, N$  stage:

2: Obtain RoI features by RoIAlign with  $\tilde{\mathbf{B}}^{i-1}$ .

3: Obtain predictions via ISTR encoder and heads.

4: Update  $\tilde{\mathbf{B}}^i$  from  $\tilde{\mathbf{B}}^{i-1}$  with the predicted boxes.

5: Output the set of predictions in the final stage.

---

**Set Prediction Loss:** For the set prediction loss, we use the matched predictions to regress the ground truth targets. The set prediction loss is defined as:

$$\mathcal{L}_{set}(\mathbf{Y}, \tilde{\mathbf{Y}}, \sigma^*) = \frac{1}{n} \sum_{i=1}^n (\mathcal{L}_{box}(\mathbf{b}_i, \tilde{\mathbf{b}}_{\sigma^*(i)}) + \mathcal{L}_{cls}(\mathbf{c}_i, \tilde{\mathbf{c}}_{\sigma^*(i)}) + \mathcal{L}_{mask}(\mathbf{m}_i, \tilde{\mathbf{r}}_{\sigma^*(i)}), \quad (8)$$

where  $\mathcal{L}_{box}(\cdot, \cdot)$  is defined the same as  $\mathcal{C}_{box}(\cdot, \cdot)$ , and  $\mathcal{L}_{cls}(\cdot, \cdot)$  is the focal loss [25] for classification. For masks, we add the dice loss [30] to improve the learned embeddings for reconstructing the masks. The mask loss is defined as:

$$\mathcal{L}_{mask}(\mathbf{m}_i, \tilde{\mathbf{r}}_{\sigma^*(i)}) = \lambda_{mask} \cdot \left( \mathcal{L}_{L2}(g(\mathbf{m}_i), \tilde{\mathbf{r}}_{\sigma^*(i)}) + \mathcal{L}_{dice}(\mathbf{m}_i, f(\tilde{\mathbf{r}}_{\sigma^*(i)})) \right), \quad (9)$$

where  $\mathcal{L}_{L2}(\cdot, \cdot)$  denotes the L2 loss and  $\mathcal{L}_{dice}(\cdot, \cdot)$  denotes the dice loss.

### 3.3. Instance Segmentation Transformer

The architecture of ISTR is depicted in Fig. 2. It contains four main components: a CNN backbone with FPN [24] to extract features for each instance, a Transformer encoder with dynamic attention to learn the relations between objects, a set of prediction heads that conduct detection and segmentation concurrently, as well as the  $N$ -step recurrent update for refining the set of predictions.

	$AP^m$	$AP_{50}^m$	$AP_{75}^m$	$AP^b$	$AP_{50}^b$	$AP_{75}^b$
<i>mask</i>	31.6	53.3	33.0	40.2	58.4	43.6
$l=40$	33.7	<b>55.6</b>	35.5	<b>41.1</b>	58.9	<b>44.8</b>
$l=60$	<b>34.2</b>	<b>55.6</b>	<b>36.4</b>	41.0	58.9	44.4
$l=80$	33.8	55.3	35.9	40.6	58.6	44.1
$l=784$	31.4	55.2	31.8	<b>41.1</b>	<b>60.0</b>	44.4

(a) **Mask vs. Mask Embeddings:** Regression with mask embeddings instead of masks brings better performance to the mask APs. The performance improves when the embedding dimension  $l=60$ , and saturates when the dimension  $l=80$ . Directly expanding the mask as embeddings, *i.e.*,  $l=784$ , has worse performance.

	$AP^m$	$AP_{50}^m$	$AP_{75}^m$	$AP_S^m$	$AP_M^m$	$AP_L^m$
dice	32.6	55.0	33.5	17.3	34.8	47.6
L2	33.8	55.5	35.5	17.3	36.1	49.5
L2+dice	<b>34.2</b>	<b>55.6</b>	<b>36.4</b>	<b>17.6</b>	<b>36.5</b>	<b>50.6</b>

(c) **Loss Functions:** Learning masks with both a pixel-level dice loss and an embedding-level L2 loss yields gains in mask APs.

	img.	type	pos.	$AP^m$	$AP_{50}^m$	$AP_{75}^m$	$AP_S^m$	$AP_M^m$	$AP_L^m$	$AP^b$	$AP_{50}^b$	$AP_{75}^b$	$AP_S^b$	$AP_M^b$	$AP_L^b$
<i>features</i>	✓	max		33.3	54.6	35.2	17.4	35.6	48.5	40.0	57.6	43.5	24.2	42.6	52.6
	✓	avg		33.6	55.3	35.5	17.3	36.4	49.5	40.6	58.8	44.1	23.9	43.2	54.1
	-	-	✓	34.0	55.3	36.3	17.4	36.4	50.5	40.9	58.8	44.3	24.3	43.2	55.0
	✓	avg	✓	<b>34.2</b>	<b>55.6</b>	<b>36.4</b>	<b>17.6</b>	<b>36.5</b>	<b>50.6</b>	<b>41.0</b>	<b>58.9</b>	<b>44.4</b>	<b>24.8</b>	<b>43.3</b>	<b>55.3</b>

(e) **Pooling Type:** The global average pooling yields better performance than the global max-pooling. Combining the image features with position embeddings increases the performance in both mask and box APs.

	$AP^m$	$AP_{50}^m$	$AP_{75}^m$	$AP^b$	$AP_{50}^b$	$AP_{75}^b$
w/o	33.8	55.5	35.7	40.7	58.8	44.1
dice	33.9	55.4	35.8	40.8	58.8	44.3
L1	34.0	55.5	35.6	40.9	<b>58.9</b>	44.3
cosine	<b>34.2</b>	<b>55.6</b>	<b>36.4</b>	<b>41.0</b>	<b>58.9</b>	<b>44.4</b>

(b) **Mask Cost Functions:** Matching with the dice loss between the predicted and ground truth masks performs slightly better than w/o the mask cost. The L1 loss between the predicted and encoded mask embeddings also has a slight improvement. Using cosine similarity as the mask cost function brings expected gains.

	$AP^m$	$AP_{50}^m$	$AP_{75}^m$	$AP^b$	$AP_{50}^b$	$AP_{75}^b$
multi-head	22.0	43.9	19.6	31.2	50.5	32.5
dynamic	<b>34.2</b>	<b>55.6</b>	<b>36.4</b>	<b>41.0</b>	<b>58.9</b>	<b>44.4</b>
	+12.2	+11.7	+16.8	+9.8	+8.4	+11.9

(d) **Attention Type:** Dynamic attention brings significant gains compared with multi-head attention in fusing the ROI and image features.

Table 1: **Ablations.** We train on the COCO `train2017` split, and report *mask* as well as *box* APs on the `val2017` split.

**Backbone:** We use a CNN backbone with FPN to extract the features ranging from P2 to P5 level of the feature pyramid. Then,  $k$  learnable query boxes  $\tilde{B}^0 = \{\tilde{b}_i^0 | i = 1, \dots, k\}$  initially covering the whole images are used to extract  $k$  ROI features  $U^0 \in \mathbb{R}^{k \times d \times t \times t}$  via RoIAlign [16]. Image features  $P \in \mathbb{R}^{k \times d}$  are first extracted by averaging and summing the features from P2 to P5, and then expand the first dimension to  $k$  for each ROI feature. Learnable position embeddings  $E \in \mathbb{R}^{k \times d}$  are initialized randomly.

**Transformer Encoder and Dynamic Attention:** The sum of image features and position embeddings is first transformed by three learnable weight matrices to obtain the inputs  $Q = (P + E)W_Q$ ,  $K = (P + E)W_K$ ,  $V = (P + E)W_V$  for the self-attention module defined as:

$$Z = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V. \quad (10)$$

The multi-head attention comprises multiple self-attention blocks, *e.g.*, eight in the original Transformer [42], to encapsulate multiple complex relationships amongst different features. Inspired by [37], we add a dynamic attention module for better fusing the ROI and image features, which is defined as the attention conditioned on the ROI features  $U^i$  in the  $i$ -th step:

$$O^i = U^i \cdot fc(Z), \quad (11)$$

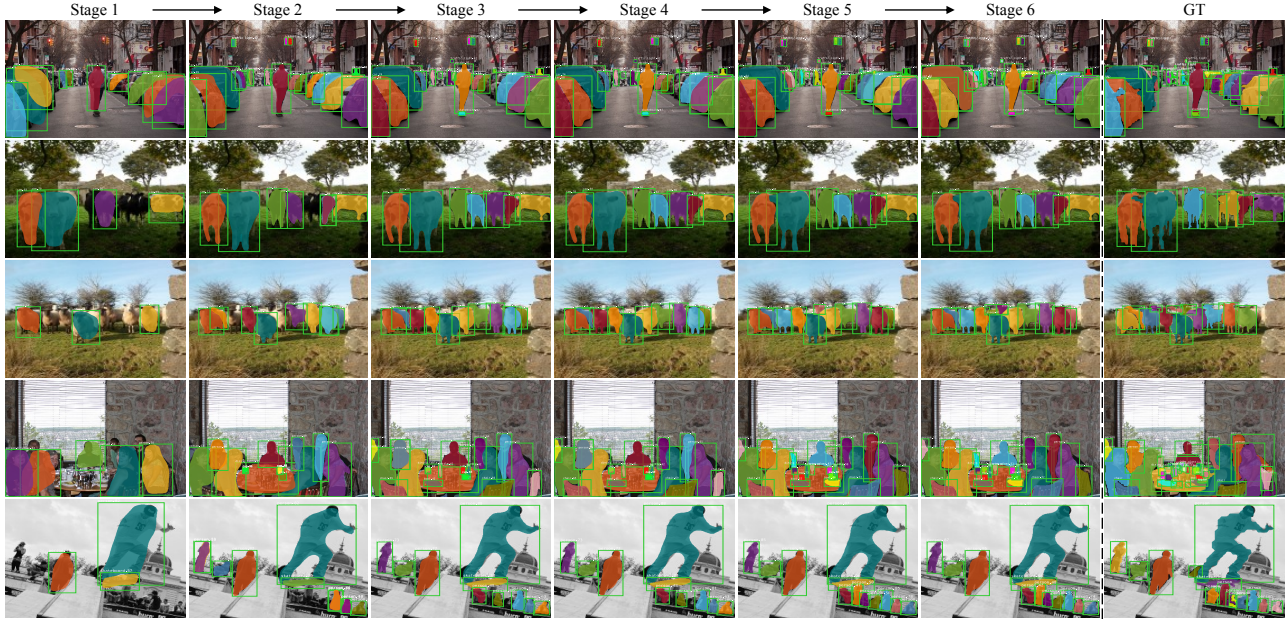
where  $fc(\cdot)$  denotes a fully connected layer to generate the dynamic parameters. The obtained features  $O^i$  are then used in the prediction heads to produce the set of outputs.

**Prediction Heads:** The set of predictions is computed by the heads, including a class head, a box head, a mask head, and a fixed mask decoder. The box head predicts the residual values of normalized center coordinates, height, and width for updating the query boxes  $\tilde{B}^i$  in the  $i$ -th step, and the class head predicts the classes using a softmax function. The mask head outputs the mask embeddings, which are then reconstructed to predict masks using the pre-learned mask decoder.

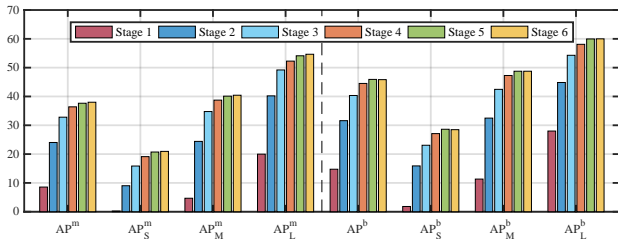
**Recurrent Refinement Strategy:** The query boxes  $\tilde{B}^i$  are recurrently updated by the predicted boxes, which refines the predictions and makes it possible to process the detection and segmentation concurrently. The overall process is summarized in Alg. 1.

## 4. Experiments

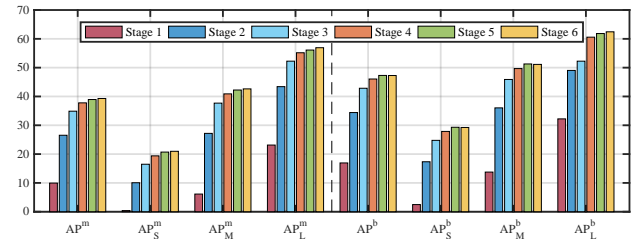
**Dataset and Evaluation Metrics:** Our experiments are performed on the MS COCO dataset [26], which contains 123K images with 80-class instance labels. Our models are trained on the `train2017` split (118K images), and the ablation study is carried out on the `val2017` split



(a) **Qualitative Results:** Masks and bounding boxes are from ISTR using ResNet101-FPN on the COCO val2017 split, with a threshold of 0.4.



(b) AP results on the COCO val2017 split with **ResNet50-FPN**.



(c) AP results on the COCO val2017 split with **ResNet101-FPN**.

Figure 3: **Stage Analysis.** We report quantitative and qualitative results of ISTR for  $N = 6$  stages. Both results show that the bounding boxes and masks are refined stage by stage, and the performance saturates in the last two stages.

(5K images). Final results are reported on the `test-dev` split, which has no public labels and is evaluated on the server. We report the standard COCO metrics including AP (*i.e.*, averaged over IoU thresholds),  $AP_{50}$ ,  $AP_{75}$ , and  $AP_S$ ,  $AP_M$ ,  $AP_L$  (*i.e.*, AP at different scales) for both boxes and masks, denoted as  $AP^b$  and  $AP^m$ , respectively.

**Training Details:** We follow the strategy in [53, 39] to encode the ground truth mask embeddings. ResNet50 and ResNet101 [17] pre-trained on ImageNet [9] are used as our backbone networks. FPN [24] is used to extract the feature pyramid. For the ablation study, all the models are trained over 12 epochs with learning rate decay, dividing by 10 at epoch 9 and 11, respectively. All results in the ablation study are tested with ResNet50-FPN and reported on the COCO val2017 split. The training schedule for the final models is 36 epochs with the learning rate divided by 10 at epoch 27 and 33, respectively. The mini-batch contains 16 images, and the models are trained with eight GPUs. Following [37], we use the AdamW [29] optimizer

and an initial learning rate of 0.000025. The input images are resized such that the shortest side is at least 480 and at most 800 pixels, while the longest side is at most 1333. The number of predictions, *i.e.*,  $k$ , is set to 300, and the number of self-attention blocks in the multi-head attention is set to 8. The number of recurrent refinement stages is set to 6. Following [3, 56], we set  $\lambda_{cls}$ ,  $\lambda_{L1}$ , and  $\lambda_{giou}$  to 2, 5, and 2, respectively. The hyperparameter  $\lambda_{mask}$  is set to 2.

#### 4.1. Ablation Experiments

To analyze ISTR, we conduct ablation studies on the choices of mask embeddings, cost functions, loss functions, the effect of position embeddings, and pooling types. Results are shown in Table 1 and discussed in detail next.

**Mask vs. Mask Embeddings:** Table 1a shows the models with various types of mask representations, including the original masks with  $28 \times 28$  dimensions and the mask embeddings with different dimensions  $l$ . We also expand the original masks to vectors to test the performance. Us-

method	backbone	Epochs	AP <sup>m</sup>	AP <sup>m</sup> <sub>S</sub>	AP <sup>m</sup> <sub>M</sub>	AP <sup>m</sup> <sub>L</sub>	AP <sup>b</sup>	AP <sup>b</sup> <sub>S</sub>	AP <sup>b</sup> <sub>M</sub>	AP <sup>b</sup> <sub>L</sub>	FPS	Time	GPU
Mask R-CNN [16]	ResNet50-FPN	36	37.5	21.1	39.6	48.3	41.3	24.2	43.6	51.7	15.3	65.6	1080Ti
MEInst [53]	ResNet50-FPN	36	33.5	19.3	35.7	42.1	42.5	25.6	45.1	52.2	15.0	66.8	1080Ti
CondInst [40]	ResNet50-FPN	36	37.8	21.0	40.3	48.7	42.1	25.1	44.5	52.1	15.4	65.0	1080Ti
BlendMask [4]	ResNet50-FPN	36	37.8	18.8	40.9	53.6	43.0	25.3	45.4	54.0	15.0	66.8	1080Ti
SOLOv2 [46]	ResNet50-FPN	36	38.2	16	<b>41.2</b>	<b>55.4</b>	40.7	18.4	43.5	57.6	10.5	95.5	1080Ti
DETR [3]	ResNet50	500	-	-	-	-	42.0	20.5	45.8	<b>61.1</b>	-	-	-
Sparse R-CNN [37]	ResNet50-FPN	36	-	-	-	-	44.5	26.9	47.2	59.5	-	-	-
<b>ISTR, ours</b>	ResNet50-FPN	36	<b>38.6</b>	<b>22.1</b>	40.4	50.6	<b>46.8</b>	<b>27.8</b>	<b>48.7</b>	59.9	13.8	72.5	1080Ti
Mask R-CNN [16]	ResNet101-FPN	36	38.8	21.8	41.4	50.5	43.1	25.1	46.0	54.3	11.8	85.0	1080Ti
MEInst [53]	ResNet101-FPN	36	35.3	20.4	37.8	44.5	44.5	26.8	47.3	54.9	11.2	89.3	1080Ti
CondInst [40]	ResNet101-FPN	36	39.1	21.5	41.7	50.9	43.5	25.8	46.0	54.1	12.0	83.2	1080Ti
BlendMask [4]	ResNet101-FPN	36	39.6	22.4	42.2	51.4	44.7	26.6	47.5	55.6	11.5	86.6	1080Ti
SOLOv2 [46]	ResNet101-FPN	36	39.7	17.3	<b>42.9</b>	<b>57.4</b>	42.6	22.3	46.7	56.3	9.0	111.6	1080Ti
DETR [3]	ResNet101	500	-	-	-	-	43.5	21.9	48.0	<b>61.8</b>	-	-	-
Sparse R-CNN [37]	ResNet101-FPN	36	-	-	-	-	45.6	<b>28.3</b>	48.3	61.6	-	-	-
<b>ISTR, ours</b>	ResNet101-FPN	36	<b>39.9</b>	<b>22.8</b>	41.9	52.3	<b>48.1</b>	<b>28.7</b>	<b>50.4</b>	61.5	11.0	91.3	1080Ti

Table 2: **Quantitative Results** of ISTR on the COCO test-dev split. All the models are learned by multi-scale training. The results of FPS are measured with a single 1080Ti GPU. The performance of Mask R-CNN is the result of the modified version with implementation details in TensorMask [7].

ing raw masks for segmentation is implemented with the mask head from Mask R-CNN in a top-down manner, and other settings are the same as ISTR. We obtain the following results. First, predicting mask embeddings brings better performance to the mask APs than predicting masks. The results verify our concern that the high-dimensional masks cannot be effectively learned with a small number of matched samples. In contrast, the mask embeddings can be well regressed, as their dimensions are much lower than those of masks, *e.g.*, 40, 60, 80 vs. 784. Second, the performance of mask embeddings improves when the dimension  $l=60$  and saturates when the dimension  $l=80$ . Finally, regressing with the expanded masks, *i.e.*,  $l=784$ , as embeddings has worse performance in mask APs.

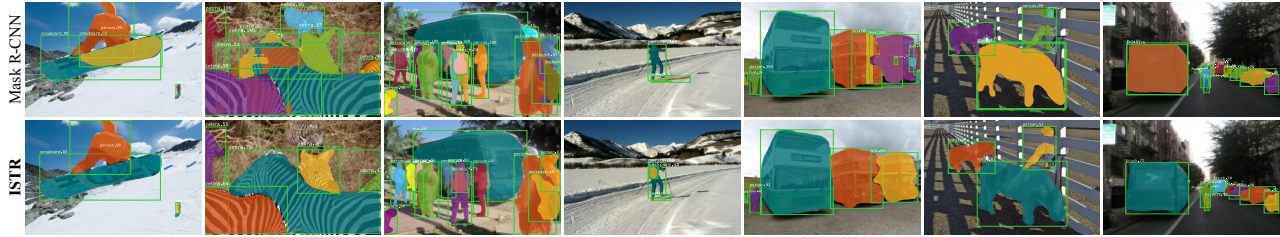
**Cost Functions:** Appropriate cost functions can match high-quality predictions with ground truth labels for the set loss. In Table 1b, we compare the performance with different choices of cost functions for the masks. Matching the predicted masks with ground truth masks by the dice loss does not produce expected gains compared to matching without the mask cost function. Using the L1 loss between embeddings slightly improves performance, and using the cosine similarity between embeddings as the mask cost function brings expected gains.

**Loss Functions:** We investigate two types of loss functions for ISTR: the dice loss at the pixel-level and the L2 loss at the embedding-level. The dice loss is calculated using the masks reconstructed by the mask decoder. As shown in Table 1c, only calculating the dice loss at the pixel-level without constraining the mask embeddings has an inferior performance. Although the L2 loss between the predicted

and encoded mask embeddings improves the performance, the learned mask embeddings are slightly suboptimal for reconstructing the original masks. Therefore, training with both the pixel-level dice loss and the embedding-level L2 loss produces better results.

**Attention Type:** We next study the effect of the dynamic attention module, which is used to fuse the RoI and image features, by replacing it with the multi-head attention module. As shown in Table 1d, the dynamic attention module performs much better than the multi-head attention module. We believe this may be because the multi-projection in the multi-head attention complicates the fusion of RoI and image features, which is essential for learning the relations between objects. From the results, we infer that a single projection is more effective for learning such relations.

**Pooling Type:** In Table 1e, we evaluate various strategies for obtaining the image features by extracting different information from input images. As can be seen, the global max-pooling does not obtain a high score, while the global average pooling performs better. We believe this is because the max-pooling extracts features from the highest activated pixel in the feature map, which usually corresponds to a single activated object. In contrast, the global average pooling extracts features that contain information about the whole image. We also find that the position embeddings are essential for achieving high results. By summing the position embeddings with averaged image features, the performance is significantly improved.



(a) Mask R-CNN [16] (top) vs. ISTR (bottom) using ResNet101-FPN. Mask R-CNN suffers inferior segmentation when bad duplicate removal occurs.



(b) More visualization results of ISTR, using ResNet101-FPN and running at 11.0 fps on a 1080Ti GPU, with 39.9 mask AP (Table 2).

Figure 4: **Qualitative Results** of ISTR on the COCO `test-dev` split. Predictions are shown with a threshold of 0.4

## 4.2. Stage Analysis

One of the essential components of ISTR is the recurrent refinement strategy, which provides a new way to achieve instance segmentation compared to the bottom-up and top-down strategies. The prediction heads infer the classes, bounding boxes, and masks for each instance using the query boxes updated in each stage. We investigate the performance of the recurrent refinement stages both quantitatively and qualitatively in Fig. 3. From the visualization of masks and bounding boxes shown in Fig. 3a, we can see that both masks and boxes are refined from coarse to fine. The mask and box APs shown in Fig. 3b and Fig. 3c using different backbones also demonstrate that the results are gradually refined step by step.

## 4.3. Main Results

We compare ISTR with the state-of-the-art instance segmentation methods as well as the latest end-to-end object detection methods to demonstrate its superior performance.

**Quantitative Results:** From Table 2, we can see that ISTR performs well, especially on small objects. For example, the  $AP_S^m$  of ISTR based on ResNet101-FPN outperforms SOLOv2 based on ResNet101-FPN by 5.5 points. We believe this is because the bipartite matching cost does not filter small objects for training. MEInst also uses mask em-

beddings for instance segmentation. However, the performance of MEInst suffers significantly due to the redundant predictions of mask embeddings. For example, the  $AP^m$  of ISTR based on ResNet101-FPN outperforms MEInst based on ResNet101-FPN by a large margin of 4.6 points. Besides, we also find performance gains of ISTR in detection when comparing the results with the state-of-the-art end-to-end object detection methods. We find that the  $AP^b$  of ISTR outperforms DETR and sparse R-CNN based on ResNet101-FPN by 4.6 and 2.5 points, respectively. Overall, it is surprising that, despite the suboptimal mask embeddings from PCA, ISTR can still obtain such a good result. This demonstrates the strength of the proposed end-to-end mechanism and shows the potential of concurrently conducting detection and segmentation with Transformers.

**Qualitative Results:** We show some examples comparing ISTR with Mask R-CNN in Fig. 4a. As can be seen, Mask R-CNN suffers inferior performance when NMS does not remove the duplicate predictions. More visualization results in Fig. 4b suggest that, although ISTR obtains state-of-the-art mask APs, there is still room for further improvement by learning finer masks. We leave this for future work.



## 5. Conclusion

In this paper, we propose a new framework, termed instance segmentation Transformer (ISTR), to explore the end-to-end mechanism for instance segmentation. ISTR predicts low-dimensional mask embeddings instead of high-dimensional masks, which inspires the design of a mask matching cost and facilitates the regression. Besides, ISTR concurrently conducts detection and segmentation via a recurrent refinement strategy, which provides a new perspective to achieve end-to-end instance segmentation and boosts the performance of both tasks. On the challenging COCO dataset, the strong performance of ISTR demonstrates its potential for instance-level recognition tasks.

## References

- [1] Daniel Bolya, Chong Zhou, Fanyi Xiao, and Yong Jae Lee. Yolact: Real-time instance segmentation. In *Proceedings of the IEEE International Conference on Computer Vision*, 2019. 2
- [2] Jiale Cao, Yanwei Pang, and Xuelong Li. Triply supervised decoder networks for joint detection and segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 3
- [3] 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*, 2020. 1, 2, 4, 6, 7
- [4] Hao Chen, Kunyang Sun, Zhi Tian, Chunhua Shen, Yongming Huang, and Youliang Yan. Blendmask: Top-down meets bottom-up for instance segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020. 1, 2, 7
- [5] Hanqing Chen, Yunhe Wang, Tianyu Guo, Chang Xu, Yiping Deng, Zhenhua Liu, Siwei Ma, Chunqing Xu, Chao Xu, and Wen Gao. Pre-trained image processing transformer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021. 2
- [6] Liang-Chieh Chen, Alexander Hermans, George Papandreou, Florian Schroff, Peng Wang, and Hartwig Adam. Masklab: Instance segmentation by refining object detection with semantic and direction features. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 3
- [7] Xinlei Chen, Ross Girshick, Kaiming He, and Piotr Dollár. Tensormask: A foundation for dense object segmentation. In *Proceedings of the IEEE International Conference on Computer Vision*, 2019. 2, 7
- [8] Bert De Brabandere, Davy Neven, and Luc Van Gool. Semantic instance segmentation with a discriminative loss function. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshop*, 2017. 2
- [9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2009. 6
- [10] Carl Doersch, Ankush Gupta, and Andrew Zisserman. Crosstransformers: spatially-aware few-shot transfer. In *Advances in Neural Information Processing Systems*, 2020. 2
- [11] 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*, 2020. 2
- [12] Naiyu Gao, Yanhu Shan, Yupei Wang, Xin Zhao, Yanan Yu, Ming Yang, and Kaiqi Huang. Ssap: Single-shot instance segmentation with affinity pyramid. In *Proceedings of the IEEE International Conference on Computer Vision*, 2019. 2
- [13] Rohit Girdhar, Joao Carreira, Carl Doersch, and Andrew Zisserman. Video action transformer network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 2
- [14] Kai Han, Yunhe Wang, Hanqing Chen, Xinghao Chen, Jianyuan Guo, Zhenhua Liu, Yehui Tang, An Xiao, Chunqing Xu, Yixing Xu, et al. A survey on visual transformer. *Arxiv preprint*, 2020. 2
- [15] Bharath Hariharan, Pablo Arbeláez, Ross Girshick, and Jitendra Malik. Simultaneous detection and segmentation. In *European Conference on Computer Vision*, 2014. 3
- [16] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE International Conference on Computer Vision*, 2017. 1, 2, 3, 5, 7, 8
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016. 6
- [18] Han Hu, Jiayuan Gu, Zheng Zhang, Jifeng Dai, and Yichen Wei. Relation networks for object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 1, 2
- [19] Zhaojin Huang, Lichao Huang, Yongchao Gong, Chang Huang, and Xinggang Wang. Mask scoring r-cnn. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 1, 2
- [20] Salman Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Shahbaz Khan, and Mubarak Shah. Transformers in vision: A survey. *Arxiv preprint*, 2021. 2
- [21] Manoj Kumar, Dirk Weissenborn, and Nal Kalchbrenner. Colorization transformer. In *International Conference on Learning Representations*, 2021. 2
- [22] Youngwan Lee and Jongyoul Park. Centermask: Real-time anchor-free instance segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020. 1, 2
- [23] Yi Li, Haozhi Qi, Jifeng Dai, Xiangyang Ji, and Yichen Wei. Fully convolutional instance-aware semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 2
- [24] 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*, 2017. 3, 4, 6
- [25] 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*, 2017. 4
- [26] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, 2014. 2, 5
- [27] Shu Liu, Jiaya Jia, Sanja Fidler, and Raquel Urtasun. Sgn: Sequential grouping networks for instance segmentation. In *Proceedings of the IEEE International Conference on Computer Vision*, 2017. 2
- [28] 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*, 2018. 1, 2, 3
- [29] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2018. 6
- [30] Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi. V-net: Fully convolutional neural networks for volumetric medical image segmentation. In *International Conference on 3D Vision*, 2016. 4
- [31] Mengye Ren and Richard S Zemel. End-to-end instance segmentation with recurrent attention. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 1, 2
- [32] Hamid Rezaatoughi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian Reid, and Silvio Savarese. Generalized intersection over union: A metric and a loss for bounding box regression. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 4
- [33] Bernardino Romera-Paredes and Philip Hilaire Sean Torr. Recurrent instance segmentation. In *European Conference on Computer Vision*, 2016. 1, 2
- [34] Yunhang Shen, Rongrong Ji, Yan Wang, Yongjian Wu, and Liujuan Cao. Cyclic guidance for weakly supervised joint detection and segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 3
- [35] Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. Vi-bert: Pre-training of generic visual-linguistic representations. In *International Conference on Learning Representations*, 2020. 2
- [36] Chen Sun, Austin Myers, Carl Vondrick, Kevin Murphy, and Cordelia Schmid. Videobert: A joint model for video and language representation learning. In *Proceedings of the IEEE International Conference on Computer Vision*, 2019. 2
- [37] 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. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021. 1, 2, 4, 5, 6, 7
- [38] Hao Tan and Mohit Bansal. Lxmert: Learning cross-modality encoder representations from transformers. In *Conference on Empirical Methods in Natural Language Processing*, 2019. 2
- [39] Zhi Tian, Tong He, Chunhua Shen, and Youliang Yan. Decoders matter for semantic segmentation: Data-dependent decoding enables flexible feature aggregation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 6
- [40] Zhi Tian, Chunhua Shen, and Hao Chen. Conditional convolutions for instance segmentation. In *European Conference on Computer Vision*, 2020. 7
- [41] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. *Arxiv preprint*, 2020. 2
- [42] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, 2017. 1, 2, 5
- [43] Jianfeng Wang, Lin Song, Zeming Li, Hongbin Sun, Jian Sun, and Nanning Zheng. End-to-end object detection with fully convolutional network. *Arxiv preprint*, 2020. 1, 2
- [44] Xinlong Wang, Tao Kong, Chunhua Shen, Yuning Jiang, and Lei Li. Solo: Segmenting objects by locations. In *European Conference on Computer Vision*, 2020. 1, 2
- [45] Xinpeng Wang, Chandan Yeshwanth, and Matthias Nießner. Sceneformer: Indoor scene generation with transformers. *Arxiv preprint*, 2020. 2
- [46] Xinlong Wang, Rufeng Zhang, Tao Kong, Lei Li, and Chunhua Shen. Solov2: Dynamic and fast instance segmentation. In *Advances in Neural Information Processing Systems*, 2020. 2, 7
- [47] Yuqing Wang, Zhaoliang Xu, Xinlong Wang, Chunhua Shen, Baoshan Cheng, Hao Shen, and Huaxia Xia. End-to-end video instance segmentation with transformers. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021. 2, 3
- [48] Enze Xie, Peize Sun, Xiaoge Song, Wenhai Wang, Xuebo Liu, Ding Liang, Chunhua Shen, and Ping Luo. Polarmask: Single shot instance segmentation with polar representation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020. 2
- [49] Fuzhi Yang, Huan Yang, Jianlong Fu, Hongtao Lu, and Bainig Guo. Learning texture transformer network for image super-resolution. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020. 2
- [50] Han-Jia Ye, Hexiang Hu, De-Chuan Zhan, and Fei Sha. Few-shot learning via embedding adaptation with set-to-set functions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020. 2
- [51] Linwei Ye, Mrigank Rochan, Zhi Liu, and Yang Wang. Cross-modal self-attention network for referring image segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 2
- [52] Jialin Yuan, Chao Chen, and Li Fuxin. Deep variational instance segmentation. In *Advances in Neural Information Processing Systems*, 2020. 2

- [53] Rufeng Zhang, Zhi Tian, Chunhua Shen, Mingyu You, and Youliang Yan. Mask encoding for single shot instance segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020. [2](#), [6](#), [7](#)
- [54] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei Fu, Jianfeng Feng, Tao Xiang, Philip HS Torr, et al. Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021. [2](#)
- [55] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. *Arxiv preprint*, 2019. [2](#)
- [56] 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*, 2020. [1](#), [2](#), [6](#)