

# PPMN: Pixel-Phrase Matching Network for One-Stage Panoptic Narrative Grounding

Zihan Ding\*  
dingzihan737@gmail.com  
Institute of Artificial Intelligence,  
Beihang University  
Hangzhou Innovation Institute,  
Beihang University

Junshi Huang  
Meituan

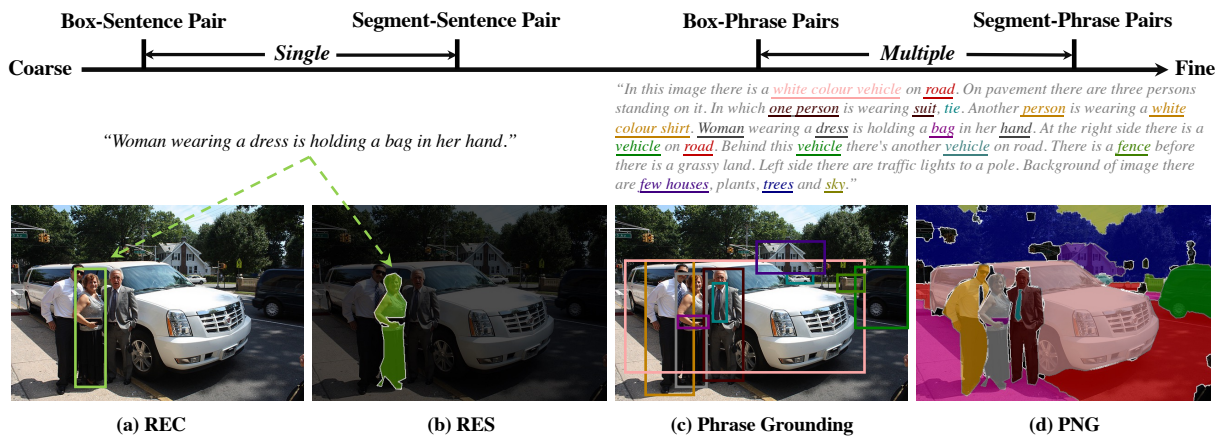
Zi-han Ding\*  
zihanding819@gmail.com  
Institute of Artificial Intelligence,  
Beihang University  
Hangzhou Innovation Institute,  
Beihang University

Xiaoming Wei  
Meituan

Tianrui Hui†  
huitianrui@gmail.com  
Institute of Information Engineering,  
Chinese Academy of Sciences  
School of Cyber Security, University  
of Chinese Academy of Sciences

Xiaolin Wei  
Meituan

Si Liu  
Institute of Artificial Intelligence,  
Beihang University  
Hangzhou Innovation Institute,  
Beihang University



**Figure 1: Illustration of different visual grounding tasks. (a) Referring Expression Comprehension (REC) aims to locate the object referred by a sentence. (b) Referring Expression Segmentation (RES) aims to segment the object referred by a sentence. (c) Phrase Grounding aims to locate multiple objects referred by noun phrases. (d) Panoptic Narrative Grounding (PNG) aims to generate a panoptic segmentation according to dense narrative captions. For better viewing of all figures in this paper, please see original zoomed-in color pdf file.**

\*Equal contribution.

†Corresponding author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

MM '22, October 10–14, 2022, Lisboa, Portugal

© 2022 Association for Computing Machinery.

ACM ISBN 978-1-4503-9203-7/22/10...\$15.00

<https://doi.org/10.1145/3503161.3548086>

## ABSTRACT

Panoptic Narrative Grounding (PNG) is an emerging task whose goal is to segment visual objects of *things* and *stuff* categories described by dense narrative captions of a still image. The previous two-stage approach first extracts segmentation region proposals by an off-the-shelf panoptic segmentation model, then conducts coarse region-phrase matching to ground the candidate regions for each noun phrase. However, the two-stage pipeline usually suffers from the performance limitation of low-quality proposals in the first stage and the loss of spatial details caused by region feature pooling, as well as complicated strategies designed for *things* and *stuff* categories separately. To alleviate these drawbacks, we propose a one-stage end-to-end Pixel-Phrase Matching Network (PPMN), which

directly matches each phrase to its corresponding pixels instead of region proposals and outputs panoptic segmentation by simple combination. Thus, our model can exploit sufficient and finer cross-modal semantic correspondence from the supervision of densely annotated pixel-phrase pairs rather than sparse region-phrase pairs. In addition, we also propose a Language-Compatible Pixel Aggregation (LCPA) module to further enhance the discriminative ability of phrase features through multi-round refinement, which selects the most compatible pixels for each phrase to adaptively aggregate the corresponding visual context. Extensive experiments show that our method achieves new state-of-the-art performance on the PNG benchmark with 4.0 absolute Average Recall gains<sup>1</sup>.

## CCS CONCEPTS

• **Computing methodologies** → **Scene understanding; Image segmentation.**

## KEYWORDS

Panoptic Narrative Grounding, Pixel-Phrase Matching

### ACM Reference Format:

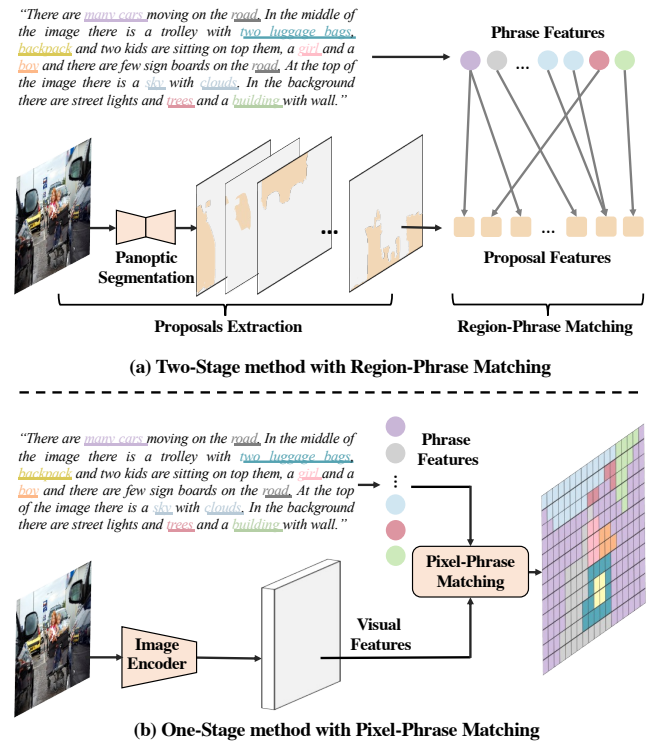
Zihan Ding, Zi-han Ding, Tianrui Hui, Junshi Huang, Xiaoming Wei, Xiaolin Wei, and Si Liu. 2022. PPMN: Pixel-Phrase Matching Network for One-Stage Panoptic Narrative Grounding. In *Proceedings of the 30th ACM International Conference on Multimedia (MM '22)*, October 10–14, 2022, Lisboa, Portugal. ACM, Lisbon, Portugal, 10 pages. <https://doi.org/10.1145/3503161.3548086>

## 1 INTRODUCTION

Panoptic Narrative Grounding (PNG) [18], which aims to segment visual objects of both *things* and *stuff* (i.e., panoptic) categories in an image grounded by dense narrative captions, is a newly proposed and more general formulation of the natural language visual grounding problem. In Figure 1 we illustrate the difference between PNG and traditional visual grounding tasks. Referring expression comprehension and segmentation [10, 25, 34, 41, 44, 49] ground the subject of a single sentence in the form of a single box and a single segment respectively, while phrase grounding [27, 75] extends them by locating multiple objects grounded by multiple noun phrases in the caption. Among these grounding tasks, PNG achieves the finest cross-modal alignment between multiple noun phrases and multiple segments. Compared to Panoptic Segmentation [32], PNG also possesses more flexibility as it requires the model to produce a panoptic segmentation according to free-form narrative captions instead of a fixed set of pre-determined object categories. With spatially finer and denser cross-modal alignment between noun phrases and segmentation masks, PNG is able to benefit many downstream applications including text-driven image manipulation [79], visual question answering [2], and intelligent robots [80], etc.

Along with the PNG benchmark, González *et al.* [18] also provide a strong **two-stage** baseline upon the Cross-Modality Relevance (CMR) [82] model (Figure 2(a)). Concretely, they first extract the region feature of each segmentation proposal from an off-the-shelf panoptic segmentation model [31]. Then, they match (i.e., ground) candidate segmentation regions conditioning on the affinity matrix calculated between features of all region proposals and noun phrases. Although this baseline has achieved good performance,

<sup>1</sup><https://github.com/dzh19990407/PPMN>



**Figure 2: Comparison between the two-stage baseline [18] and our one-stage Pixel-Phrase Matching Network (PPMN). (a) The two-stage baseline first extracts segmentation proposals using an off-the-shelf panoptic segmentation model. Then, it conducts coarse Region-Phrase Matching to select grounding results for each noun phrase. (b) Our one-stage PPMN conducts finer pixel-phrase matching by calculating the matching scores between features of all pixels and all phrases to directly generate panoptic segmentation.**

some limitations may still exist: 1) The matching performance in the second stage is restricted by the quality of segmentation proposals and the recall rate of the panoptic segmentation model in the first stage. 2) Region feature of each segmentation proposal is pooled into a single vector, in which detailed spatial information may be lost, incurring a coarse cross-modal matching process. 3) The two-stage pipeline is complicated and involves many manually designed rules, such as different feature extraction strategies for *things* and *stuff* categories, and the extra post-processing for plurals noun phrases in the inference phase.

In order to alleviate the above problems, we propose a simple yet effective **one-stage** approach for PNG, which enjoys end-to-end optimization. Specifically, instead of conducting indirect and coarse Region-Phrase Matching with an off-the-shelf panoptic segmentation model, our one-stage approach is realized via a direct and fine Pixel-Phrase Matching Network (PPMN). As illustrated in Figure 2(b), our PPMN matches each phrase to its corresponding pixels by directly calculating the matching matrix between features of all phrases and all pixels, so that each phrase can obtain a response map derived from the matching matrix. The panoptic segmentation is then generated by simply combining the results of all phrases.

It is interesting to note that multiple phrases can correspond to the same ground-truth mask in the task setting of PNG. Therefore, each pixel is actually classified to multiple labels (*i.e.*, phrases) when applying binary cross-entropy loss on all response maps, where responses on positive labels are forced to be higher and otherwise on negative ones. Overall, the merits of our PPMN contain three aspects: 1) One-stage framework can exploit sufficient cross-modal semantic correspondence from the supervision of densely annotated pixel-phrase pairs rather than sparse region-phrase pairs. 2) By avoiding proposal pooling, detailed spatial information can be preserved as well for more accurate and complete segmentation results. 3) Our PPMN can better tap the category priors contained intrinsically in the natural language to distinguish *things/stuff* and *singulars/plurals*, which removes complicated post-processing and largely simplifies the whole pipeline.

To further enhance the discriminative ability of each phrase feature, we propose a Language-Compatible Pixel Aggregation (LCPA) module. In detail, we select the most compatible pixels for each phrase according to matching scores and adaptively aggregate their features via a multi-head cross-modal attention mechanism. By this means, each phrase is aware of the corresponding visual contextual information. Our LCPA module is applied for multiple rounds to gradually refine phrase features with adaptive visual clues for more accurate pixel-phrase matching.

The main contributions of our paper are summarized as follows: 1) We propose a novel one-stage end-to-end method termed Pixel-Phrase Matching Network (PPMN) for the emerging PNG task, where each phrase is directly matched with its corresponding pixels instead of pre-generated region proposals to produce finer panoptic segmentation. 2) For more accurate pixel-phrase matching, we also propose a Language-Compatible Pixel Aggregation (LCPA) module to gradually enhance the discriminative ability of phrase features with adaptive visual clues. 3) Extensive experiments show that our one-stage method outperforms the previous two-stage baseline on the PNG benchmark with a significant gain of 4.0 overall Average Recall. Our code will be made publicly available to facilitate further research in this field.

## 2 RELATED WORK

### 2.1 Panoptic Segmentation

Panoptic segmentation [32], which aims to assign a semantic label and an instance id to each pixel, has received much attention recently. Mainstream methods [8, 31, 37, 60, 69] are usually composed of two specifically designed branches for segmenting foreground *things* (*i.e.*, instance segmentation) and background *stuff* (*i.e.*, semantic segmentation) respectively. For example, Kirillov *et al.* [31] propose to extend Mask R-CNN [19] with a semantic segmentation branch using a shared FPN [42] backbone. With the rise of Transformer [62], some methods [6, 9, 39, 63, 81] have emerged with an end-to-end set prediction objective, and generate panoptic masks through attention blocks. Differently, PNG aims to generate the panoptic segmentation for a still image according to dense narrative captions. To tackle this task, we propose a simple yet effective one-stage method that directly matches each phrase with its corresponding pixels to segment both *things* and *stuff* categories in a unified way.

### 2.2 Referring Expression Comprehension

Referring Expression Comprehension (REC) aims to predict the bounding box of the object referred by the referring expression. Existing methods can be roughly divided into two categories: two-stage methods [7, 27, 45, 48, 58, 65, 70, 76, 78] and one-stage methods [11, 40, 71–73]. Conventional two-stage methods follow the propose-and-match paradigm and design sophisticated ways of cross-modal interaction for accurate vision-language alignment (*e.g.*, fine-grained context modeling [45, 76, 78] and graph attention [27, 65, 70]). Recently, some methods [7, 48, 58] leverage the powerful modeling ability of BERT [12] to learn better vision-language representations on large-scale image-text datasets. By introducing language into one-stage detectors [57, 83], one-stage approaches [40, 71–73] are gradually taking over the mainstream as they achieve a good balance between effectiveness and efficiency. Moreover, TransVG [11] take a further step for directly regressing the box coordinates while using Transformers [58] to conduct cross-modal alignment.

### 2.3 Referring Expression Segmentation

Referring expression segmentation (RES) requires models to predict foreground pixels for the object described by the input referring expression. Hu *et al.* [22] first propose an one-stage framework for RES, where they use FCN [47] and LSTM [21] to extract visual feature maps and sentence features respectively. Then, they fuse them by concatenation to form the cross-modal features, on which they apply deconvolution layers to generate the segmentation mask. To exploit different types of informative words in the expression and accurately align the two modalities, CMPC [23] first perceives all possible entities under the guidance of *entity* and *attribute* clues of the input expression and then utilizes *relation* words to filter out irrelevant ones. Motivated by the powerful ability of Transformers [62] for capturing long-range dependencies, some methods [13, 17, 26, 36, 74] design complex cross-attention mechanisms to model the semantic relationship between vision and language modalities. Moreover, RES can be applied on consecutive video frames for segmenting the queried actors [14, 15, 24]. Differently, PNG grounds multi noun phrases in the narrative caption for panoptic categories, and in this paper we formulate PNG as a direct pixel-phrase matching process to fully mine the cross-modal correspondence from densely annotated pixel-phrase pairs.

### 2.4 Phrase Grounding

Phrase grounding aims to localize multiple regions with bounding boxes in an image referred by noun phrases in natural language descriptions. Early methods [1, 29, 55, 64] used to first extract region and phrase embeddings independently, and then learn the semantic correspondence between phrase-region pairs in a shared embedding subspace via manually designed loss functions. In order to exploit visual and textual context, some methods begin to adopt different ways of cross-modal interaction [4, 16, 46, 52, 77]. Mu *et al.* [52] propose to distinguish various contexts via motif-aware graph learning, which can achieve fine-grained contextual fusion. Yu *et al.* [77] explore how to capture omni-range dependencies through both multi-level and multi-modal interaction. Moreover, recent research shows that phrase grounding can benefit a lot from

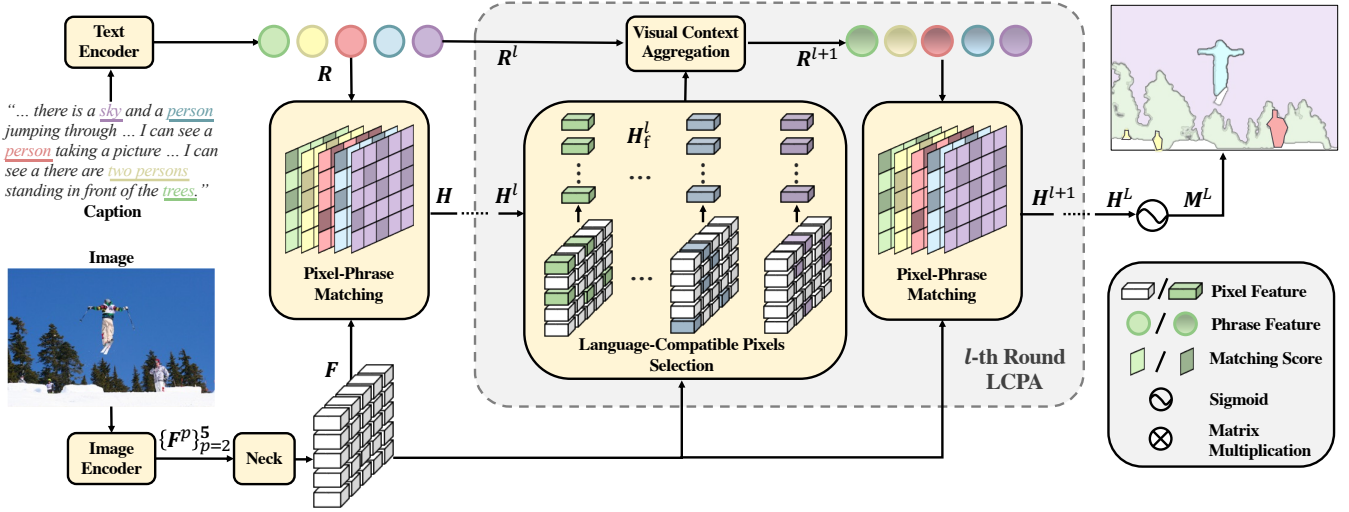


Figure 3: Overview of our Pixel-Phrase Matching Network (PPMN). We use an image encoder to extract multi-scale image features  $\{F^p\}_{p=2}^5$  and aggregate them into a single feature map  $F$ . As for the linguistic modality, we use a text encoder to extract noun phrase features  $R$ . We obtain the raw matching maps  $H$  via Pixel-Phrase Matching. To elaborate on the LCPA approach, we take the  $l$ -th round as an example. Firstly, we use  $H^l$  to select language-compatible pixels  $H_f^l$  for each noun phrase from  $F$ . After, we aggregate the specific visual context using each phrase feature in  $R^l$  as the query to get the refined phrase features  $R^{l+1}$ . And we conduct Pixel-Phrase Matching to obtain the refined matching maps  $H^{l+1}$ . After multiple rounds of refinement, we apply sigmoid on  $H^L$  from the last round to get response maps  $M^L$ , from which we generate panoptic segmentation.

the large-scale vision-language pretraining [28, 35]. Compared with phrase grounding, PNG is spatially finer and more general as it includes segmentation annotations of all panoptic categories (*i.e.*, *things* and *stuff*). In this paper, we propose a one-stage Pixel-Phrase Matching Network (PPMN) to better mine the fine-grained semantic richness of a visual scene for the PNG task.

### 3 METHOD

In this section, we first introduce the feature extraction processes for visual and linguistic modalities (§3.1). Then, we describe how the panoptic narrative grounding task can be formulated as a pixel-phrase matching problem (§3.2). Finally, we introduce the proposed Language-Compatible Pixel Aggregation (LCPA) module which enhances the discriminative ability of phrase features by aggregating the corresponding visual context for each noun phrase to achieve accurate pixel-phrase matching (§3.3). The overall pipeline of our Pixel-Phrase Matching Network (PPMN) is illustrated in Figure 3.

#### 3.1 Feature Extraction

For the visual modality, we use FPN [42] with a ResNet-101 [20] backbone as the image encoder to extract multi-scale feature maps  $F^p \in \mathbb{R}^{H^p \times W^p \times C_v}$ ,  $p \in \{2, 3, 4, 5\}$ , where  $H^p = \frac{H^0}{2^p}$ ,  $W^p = \frac{W^0}{2^p}$ , and  $C_v$  are the height, width, and channel number of the  $p$ -th visual features respectively,  $H^0$  and  $W^0$  are the original scale of the input image. To enhance the positional information, we add the sinusoids positional encoding [62] with  $F^5$ . Next, we send  $\{F^p\}_{p=2}^5$  into the semantic FPN neck [31] to obtain the final visual feature

map  $F \in \mathbb{R}^{H \times W \times C_v}$  with strong semantic representation and low-level local details, where  $H = \frac{H^0}{8}$  and  $W = \frac{W^0}{8}$  are the height and width. For the linguistic modality, we use the “base-uncased” version of BERT [12] as the text encoder to encode each word of the narrative caption to a real-valued vector, from which we extract features of noun phrases as  $R \in \mathbb{R}^{N \times C_r}$ , where  $N$  is the maximum number of noun phrases and  $C_r$  is the channel number.

#### 3.2 Pixel-Phrase Matching Formulation

For pixel-phrase matching, we directly use each noun phrase to group its corresponding pixels based on response values calculated between representations of all pixels and all noun phrases. Concretely, we first project the visual features  $F$  and phrase features  $R$  to the same  $C$ -dimensional subspace by linear layers:

$$\hat{F} = FW_1, \hat{R} = RW_2, \quad (1)$$

where  $W_1 \in \mathbb{R}^{C_v \times C}$  and  $W_2 \in \mathbb{R}^{C_r \times C}$  are projection parameters,  $\hat{F} \in \mathbb{R}^{H \times W \times C}$  and  $\hat{R} \in \mathbb{R}^{N \times C}$  are projected features. Next, we first reshape  $\hat{F}$  to  $\mathbb{R}^{HW \times C}$  and then conduct matrix multiplication between  $\hat{F}$  and  $\hat{R}$  to obtain response maps between all pixels and all noun phrases as follows:

$$H = \hat{R}\hat{F}^T, \quad (2)$$

$$M = \sigma(H), \quad (3)$$

where  $\hat{F}^T \in \mathbb{R}^{C \times HW}$  is the transpose of  $\hat{F}$ ,  $H \in \mathbb{R}^{N \times HW}$  is the raw matching maps, and  $\sigma$  denotes sigmoid function. After, we reshape  $M \in \mathbb{R}^{N \times HW}$  to  $\mathbb{R}^{N \times H \times W}$  and  $M^n \in \mathbb{R}^{H \times W}$  is the response map of the  $n$ -th noun phrase, based on which we can generate the segmentation result.

It is straight-forward to train a pixel-phrase matching network: given ground-truth binary masks of all phrases  $Y \in \mathbb{R}^{N \times H \times W}$ , we apply the binary cross-entropy (BCE) loss on the response maps  $\mathbf{M}$ . The operation can be written as:

$$\mathcal{L}_{\text{bce}}(\mathbf{M}^{n,i}, Y^{n,i}) = Y^{n,i} \log(\mathbf{M}^{n,i}) + (1 - Y^{n,i}) \log(1 - \mathbf{M}^{n,i}), \quad (4)$$

$$\bar{\mathcal{L}}_{\text{bce}}(\mathbf{M}, Y) = -\frac{1}{NHW} \sum_{n=1}^N \sum_{i=1}^{HW} \mathcal{L}_{\text{bce}}(\mathbf{M}^{n,i}, Y^{n,i}). \quad (5)$$

It is worth noting that Eq. 5 can be rewritten as follows:

$$\bar{\mathcal{L}}_{\text{bce}}(\mathbf{M}, Y) = -\frac{1}{HW} \sum_{i=1}^{HW} \frac{1}{N} \sum_{n=1}^N \mathcal{L}_{\text{bce}}(\mathbf{M}^{n,i}, Y^{n,i}), \quad (6)$$

$$= -\frac{1}{HW} \sum_{i=1}^{HW} \mathcal{L}_{\text{multi-cl}}(\mathbf{M}^{:,i}, Y^{:,i}),$$

which shows that the pixel-phrase matching is equivalent to performing the multi-label classification process on each pixel if we consider  $\mathbf{M}^{:,i} \in \mathbb{R}^{1 \times N}$  as a probability distribution over all possible  $N$  noun phrases (*i.e.*, labels) for  $i$ -th pixel. In this view, it forces the model to produce high responses to positive phrases and otherwise to negative ones.

Since BCE loss treats each pixel separately, it can not handle the foreground-background sample imbalance problem. We apply Dice loss [51] to alleviate this issue following previous works [59, 66, 81]:

$$\mathcal{L}_{\text{dice}}(\mathbf{M}^n, Y^n) = 1 - \frac{2 \sum_{i=1}^{HW} \mathbf{M}^{n,i} Y^{n,i}}{\sum_{i=1}^{HW} \mathbf{M}^{n,i} + \sum_{i=1}^{HW} Y^{n,i}}, \quad (7)$$

$$\bar{\mathcal{L}}_{\text{dice}}(\mathbf{M}, Y) = \frac{1}{N} \sum_{n=1}^N \mathcal{L}_{\text{dice}}(\mathbf{M}^n, Y^n). \quad (8)$$

Overall, the final training loss function  $\mathcal{L}_{\text{ppm}}$  of our PPMN can be formulated as:

$$\mathcal{L}_{\text{ppm}} = \lambda_{\text{bce}} \bar{\mathcal{L}}_{\text{bce}} + \lambda_{\text{dice}} \bar{\mathcal{L}}_{\text{dice}}, \quad (9)$$

where  $\lambda_{\text{bce}}$  and  $\lambda_{\text{dice}}$  are hyperparameters used to balance these two losses. We empirically find  $\lambda_{\text{bce}} = 1$  and  $\lambda_{\text{dice}} = 1$  work best.

### 3.3 Language-Compatible Pixel Aggregation

Without the explicit cross-modal interaction between visual and linguistic modalities, our model can only generate sub-optimal panoptic segmentation results using the limited category priors implied in noun phrases. To endow the phrase features with stronger discriminative ability, we propose a Language-Compatible Pixel Aggregation (LCPA) module to refine phrase features for multiple rounds. This process is formulated as:

$$\mathbf{R}^{l+1}, \mathbf{H}^{l+1} = \text{LCPA}^l(\mathbf{F}, \mathbf{R}^l, \mathbf{H}^l), l = 0, 1, \dots, L-1, \quad (10)$$

while  $\mathbf{R}^0 = \mathbf{R}$ ,  $\mathbf{H}^0 = \mathbf{H}$  and  $L$  is the number of multiple rounds.

In the  $l$ -th round, we obtain indexes  $\mathbf{H}_{\text{index}}^l \in \mathbb{R}^{N \times S \times 2}$  of the  $S$  most compatible pixels for all phrases from  $\mathbf{H}^l$ :

$$\mathbf{H}_{\text{index}}^l = \text{MaxPool}(\mathbf{H}^l, S), \quad (11)$$

where  $\text{MaxPool}(\cdot, S)$  is an adaptive max pooling layer that returns  $S$  max pooling indexes instead of the values in our implementation. Afterwards, we use  $\mathbf{H}_{\text{index}}^l$  to sample compatible pixel features  $\mathbf{H}_f^l \in \mathbb{R}^{N \times S \times C}$  from  $\mathbf{F}$ .

Next, we enhance the discriminative ability of each noun phrase by aggregating the visual context of its most compatible pixel features. To this end, we follow the implementation practice of Transformer [62] and revise it to a multi-head cross-modal attention mechanism  $\text{MCA}(\cdot)$ . Details of this process for the  $n$ -th noun phrase feature  $(\mathbf{R}^l)^n \in \mathbb{R}^{1 \times C}$  can be formulated as follows:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{W}_5(\mathbf{K}\mathbf{W}_6)^T}{\sqrt{C}}\right)\mathbf{V}\mathbf{W}_7, \quad (12)$$

$$\text{MCA} = \mathbb{C}[\{\text{Attention}((\mathbf{R}^l)^{n,d}, (\mathbf{H}_f^l)^{n,d}, (\mathbf{H}_f^l)^{n,d})\}_{d=1}^D], \quad (13)$$

where  $\mathbb{C}(\cdot)$  denotes concatenation and  $D$  is the number of heads.  $\mathbf{W}_5$ ,  $\mathbf{W}_6$ , and  $\mathbf{W}_7$  are projection parameters.  $\mathbf{Q}$ ,  $\mathbf{K}$ ,  $\mathbf{V}$  are the query, key, and value respectively.  $(\mathbf{R}^l)^{n,d} \in \mathbb{R}^{1 \times \frac{C}{d}}$  and  $(\mathbf{H}_f^l)^{n,d} \in \mathbb{R}^{1 \times \frac{C}{d}}$  are input features of  $d$ -th head. The output of  $\text{MCA}(\cdot)$  is the refined  $n$ -th phrase feature, denoted as  $(\hat{\mathbf{R}}^l)^n \in \mathbb{R}^{1 \times C}$ . After applying  $\text{MCA}(\cdot)$  on each phrase separately, we concatenate all refined phrase features and feed them to a standard Feed-Forward Network (FFN) to obtain the refined phrase features  $\mathbf{R}^{l+1} \in \mathbb{R}^{N \times C}$ :

$$\hat{\mathbf{R}}^l = \mathbb{C}[\{(\hat{\mathbf{R}}^l)^n\}_{n=1}^N] + \mathbf{R}^l, \quad (14)$$

$$\mathbf{R}^{l+1} = \text{LN}(\text{FFN}(\text{LN}(\hat{\mathbf{R}}^l))) + \hat{\mathbf{R}}^l, \quad (15)$$

where LN is the LayerNorm [3].

At last, we project  $\mathbf{F}$  and  $\mathbf{R}^{l+1}$  to the same subspace. Specifically, we project the visual features  $\mathbf{F}$  with a  $1 \times 1$  convolution layer followed by GroupNorm [67] and ReLU activation [53]. As for phrase features, we apply a multi-layer perceptron (MLP) with 3 hidden layers to  $\mathbf{R}^{l+1}$ . Then we conduct matrix multiplication between them to acquire  $\mathbf{H}^{l+1} \in \mathbb{R}^{N \times H \times W}$  and apply sigmoid function on  $\mathbf{H}^{l+1}$  to obtain response maps  $\mathbf{M}^{l+1} \in \mathbb{R}^{N \times H \times W}$ . In the training phase, we apply  $\mathcal{L}_{\text{ppm}}$  to all rounds of  $\{\mathbf{M}^l\}_{l=0}^{L-1}$  for sufficient intermediate supervisions on the learning of our LCPA module:

$$\mathcal{L} = \sum_{l=0}^{L-1} \mathcal{L}_{\text{ppm}}(\mathbf{M}^l, Y). \quad (16)$$

In the inference phase, we employ a threshold of 0.5 to obtain the grounding results from the last round response maps  $\mathbf{M}^L$ .

## 4 EXPERIMENTS

### 4.1 Dataset and Evaluation Criteria

We evaluate our PPMN on the Panoptic Narrative Grounding benchmark [18], which is extended from MS COCO [43] dataset. It includes 726,445 noun phrases from the whole Localized Narratives annotations [56] that are matched with 659,298 unique segments from MS COCO panoptic segmentation annotations. Each narrative caption contains an average of 11.3 noun phrases, of which 5.1 noun phrases are grounded. Following González *et al.* [18], we adopt the Average Recall to evaluate our PPMN. Concretely, we first calculate the Intersection over Union (IoU) between segmentation predictions and ground-truth masks for all evaluated noun phrases. Then, we compute recall at different IoU thresholds to obtain a curve where recall approaches one at very low IoU values and decreases at higher IoU values. The Average Recall refers to the area under the curve described above. As for plural noun phrases, of which

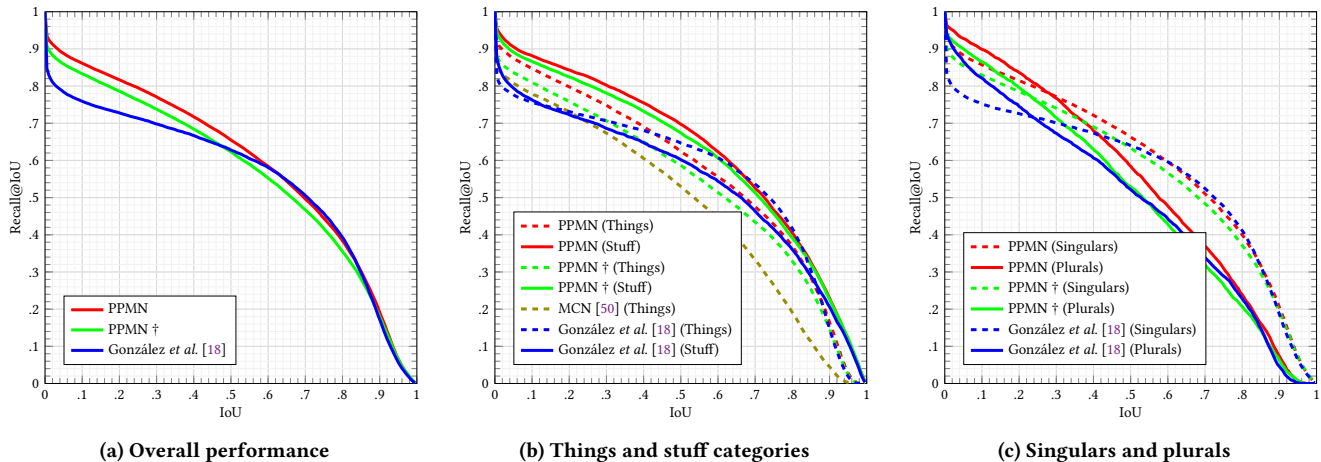


Figure 4: Average Recall Curve for our PPMN method performance (a) compared to the state-of-the-art methods, and disaggregated into (b) things and stuff categories, and (c) singulars and plurals noun phrases.

Table 1: Comparison with state-of-the-art methods on the PNG benchmark, disaggregated into (a) things and stuff categories, and (b) singulars and plurals noun phrases. † denotes training without COCO panoptic segmentation annotations.

(a) Things and stuff categories.				(b) Singulars and plurals noun phrases.			
Method	overall	Average Recall things	stuff	Method	overall	Average Recall singulars	plurals
González et al. [18]	55.4	56.2	54.3	González et al. [18]	55.4	56.2	48.8
MCN [50]	-	48.2	-	PPMN †	56.7	57.4	49.8
PPMN †	56.7	53.4	61.1	PPMN	<b>59.4 (+4.0)</b>	<b>60.0 (+3.8)</b>	<b>54.0 (+5.2)</b>
PPMN	<b>59.4 (+4.0)</b>	<b>57.2 (+1.0)</b>	<b>62.5 (+8.2)</b>				

each phrase is annotated with multiple instances, we aggregate all corresponding ground truth masks into a single segmentation.

## 4.2 Implementation Details

Consistent with González et al. [18], we use FPN [42] with a ResNet-101 [20] backbone pre-trained with Panoptic Feature Pyramid Network [31] on MS COCO [43] with 3x schedule using the official implementation [68]. We fix parameters in FPN. The input image is resized with a shorter side to 800 and a longer side up to 1333, without changing the aspect ratio. For the linguistic input, we adopt the pretrained “base-uncased” BERT model [12] to convert each word in the narrative caption into a 768-dimensional vector. The maximum length of the input caption is 230, of which at most 30 different noun phrases need to be grounded. Adam [30] is utilized as the optimizer. We implement our proposed PPMN in PyTorch [54] and train it with batch size 12 for 14 epochs on 4 NVIDIA A100 GPUs. The initial learning rate is set to  $1e^{-4}$ , which is divided by 2 for every 2 epochs started from the 10-th epoch, while the learning rate for the text encoder is set to  $1e^{-5}$  constantly. The default numbers of multiple rounds  $L$  and language-compatible pixels  $S$  are set to 3 and 200, respectively. In the inference phase, we average matching maps of all words included in each noun phrase following the setup of the two-stage baseline [18].

## 4.3 Comparison with State-of-the-Art Methods

We conduct experiments on the PNG benchmark to compare our PPMN with state-of-the-art methods. Since PNG is an emerging and challenging task, there are only two methods to compare at present. Quantitative results are shown in Table 1 and Figure 4.

Compared to the two-stage baseline [18], our PPMN achieves significant performance boosts of 4.0/1.0/8.2/3.8/5.2 on the Average Recall metric for *overall/things/stuff/singulars/plurals* splits (Table 1), indicating the superiority of our one-stage Pixel-Phrase Matching approach. From Figure 4 we can observe that the area between the red (i.e., PPMN) and blue curves (i.e., González et al. [18]) is relative large when the IoU is less than 0.5, which means that our proposed PPMN can ground much more panoptic objects than the two-stage baseline. Besides, our PPMN achieves comparable performances when the IoU is close to one, verifying that our PPMN can generate accurate and complete segmentation results without extra manually designed rules and post-processing. As for the computational overhead, our pixel-phrase matching and LCPA module occupy only 1.14 GFLOPs and 2.03 GFLOPs with negligible parameters respectively, which are not computationally heavy. To further demonstrate the effectiveness of our proposed PPMN, we directly use the ResNet-101 [20] pretrained on ImageNet [33] dataset as the image encoder, whose performances are shown in the row “PPMN †” of Table 1 and the green curves in Figure 4. It can

be seen that its Average Recall values exceed the two-stage baseline on all but one (*i.e.*, *things*) of splits, proving that our PPMN is strong enough to learn discriminative visual and linguistic features by end-to-end training, even without the panoptic priors obtained from pretraining on COCO panoptic annotations. It is worth noting that the two-stage baseline is similar to a state-of-the-art phrase grounding method, *i.e.*, GLIP [35], where they separately encode visual and phrase features and conduct cross-modal fusion before matching candidate regions with noun phrases. Thus, the performance gap between our PPMN and the two-stage baseline indicates that directly adapting phrase grounding methods to PNG can not yield compelling results. And we believe our PPMN can be a better baseline to promote research in PNG.

Multi-task Collaborative Network (MCN) [50] is a state-of-the-art visual grounding method that achieves joint learning of REC and RES. Considering that REC and RES only include objects belonging to the *things* category, we only evaluate its performance on the *things* split of the PNG benchmark [18]. As shown in Table 1a and Figure 4b, there is a big performance gap between MCN and methods designed for PNG. We claim that this phenomenon is caused by the sparse and coarse annotations (*i.e.*, box/segment-sentence pair in Figure 1) of previous grounding tasks. In those settings, models are not forced to mine the fine-grained cross-modal semantic relationship between pixels and phrases, while our proposed method achieves this by the Pixel-Phrase Matching approach.

#### 4.4 Ablation Studies

To test if our LCPA module confers benefits, we conduct ablation studies on the PNG benchmark to evaluate its different designs.

**Table 2: The number of language-compatible pixels  $S$ . We set the number of multiple rounds  $L = 1$  in this experiment.**

$S$	Average Recall				
	overall	singulars	plurals	things	stuff
0	51.6	51.9	48.7	49.5	54.6
100	58.5	59.1	53.8	56.2	61.8
200	<b>58.7</b>	<b>59.2</b>	<b>53.9</b>	<b>56.4</b>	<b>61.9</b>
300	57.4	57.9	52.6	55.0	60.7

**Number of Language-Compatible Pixels.** We evaluate different numbers of language-compatible pixels  $S$  in the Table 2. As shown in the 1-st and 2-nd rows, the performance boosts significantly when noun phrases are aware of their corresponding visual context, verifying that this cross-modal interaction can enhance the discriminative ability of noun phrase features. The best performance is achieved when  $S = 200$ . Qualitative analysis of the language-compatible pixels in different rounds is shown in §4.5.

**Number of Multiple Rounds.** We demonstrate the influence of the number of multiple rounds  $L$  in Table 3, where the Average Recall values on all evaluated splits grows gradually as  $L$  increases. Considering the balance between computational overhead and performance, we choose  $L = 3$  in our final model.

**Visual Context Aggregation.** Here we demonstrate different ways of aggregating the visual context of the language-compatible pixels for each noun phrase (Table 4). MUTAN [5] is a conventional multi-modal fusion scheme originally designed for VQA [2],

**Table 3: The number of multiple rounds  $L$ . We set the number of language-compatible pixels  $S = 200$  in this experiment.**

$L$	Average Recall				
	overall	singulars	plurals	things	stuff
1	58.7	59.2	53.9	56.4	61.9
2	59.2	59.8	53.9	56.9	62.4
3	59.4	60.0	<b>54.0</b>	<b>57.2</b>	62.5
4	<b>59.5</b>	<b>60.1</b>	<b>54.0</b>	57.1	<b>62.7</b>

**Table 4: Different implementation techniques of the visual context aggregation part in LCPA.**

Method	Average Recall				
	overall	singulars	plurals	things	stuff
MUTAN [5]	59.1	59.6	53.9	56.8	62.1
SKNet [38]	59.2	59.8	53.8	56.8	<b>62.5</b>
MCA (§3.3)	<b>59.4</b>	<b>60.0</b>	<b>54.0</b>	57.2	<b>62.5</b>

which is able to model fine and rich cross-modal interactions between visual and language modalities based on a Tucker decomposition [61]. Besides, we evaluate another implementation technique for the visual context aggregation by modifying the selective kernel convolution proposed in SKNet [38]. Concretely, we first conduct average pooling on sampled compatible pixel features to obtain the integrated visual context for each noun phrase. Then we fuse all noun phrase features and their corresponding visual context via element-wise multiplication and generate the channel-wise addition weights using a linear transformation. After, we use these weights to adaptively fuse features of language-compatible pixels and noun phrases. From Table 4 we can see that our LCPA is robust to the specific implementation techniques of the visual context aggregation part, and we use MCA (§3.3) in this paper for its best performances on all evaluated metrics.

#### 4.5 Qualitative Analysis

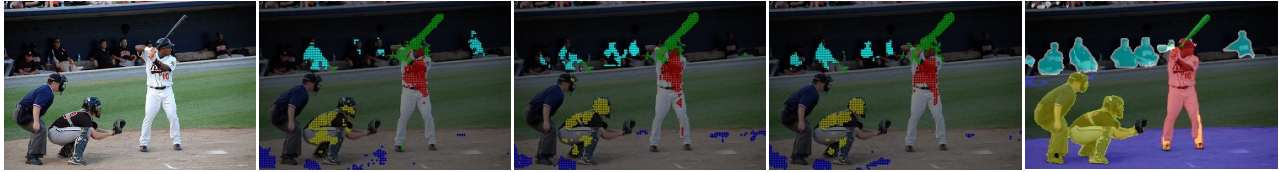
In Figure 5, we visualize locations of the language-compatible pixels in each round. Taking the 1-st row as an example, the language-compatible pixels for the plurals noun phrase “few people” (*i.e.*, light blue points) can adaptively attend to regions of different people who are sitting on the chair in different rounds. By this means, the noun phrase can aggregate sufficient visual context for all grounded people and segment them accurately and completely.

Figure 6 shows some qualitative results of our proposed PPMN compared to the ground truth masks. Surprisingly our proposed PPMN can produce better segmentation results than annotations. In the 1-st row, our PPMN hardly generates false positive predictions on the rock when grounding the noun phrase “sea”. Furthermore, our PPMN segments finer and smoother boundaries between the trees and sky in the 2-nd row and between the grass and red soil field in the 3-rd row.

## 5 CONCLUSION AND DISCUSSION

In this paper, we propose a simple yet effective one-stage framework that can be end-to-end optimized for Panoptic Narrative Grounding (PNG), called Pixel-Phrase Matching Network (PPMN). The fine-grained Pixel-Phrase Matching strategy encourages our model

In this picture there is a man holding a bat. Two persons at abc are in squat position. Some grass is visible on the ground. There are few people sitting on the bench and bottles are kept on the bench.



In this image we can see a chairs on the sand, and umbrella on it. Here is the water, and at above here is the sky.

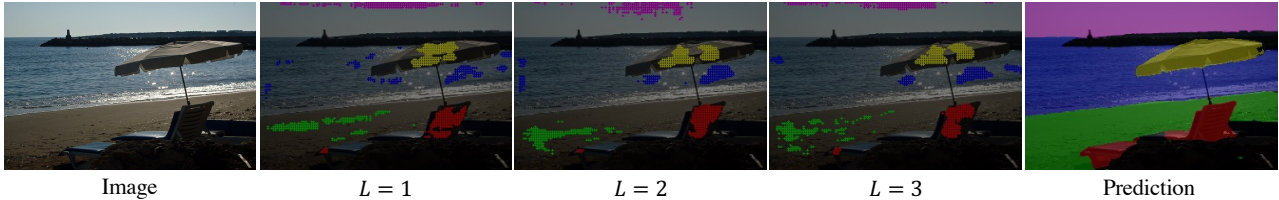


Figure 5: Visualization of the language-compatible pixels' locations in different rounds (2-nd to 4-th column) and segmentation results (5-th column).

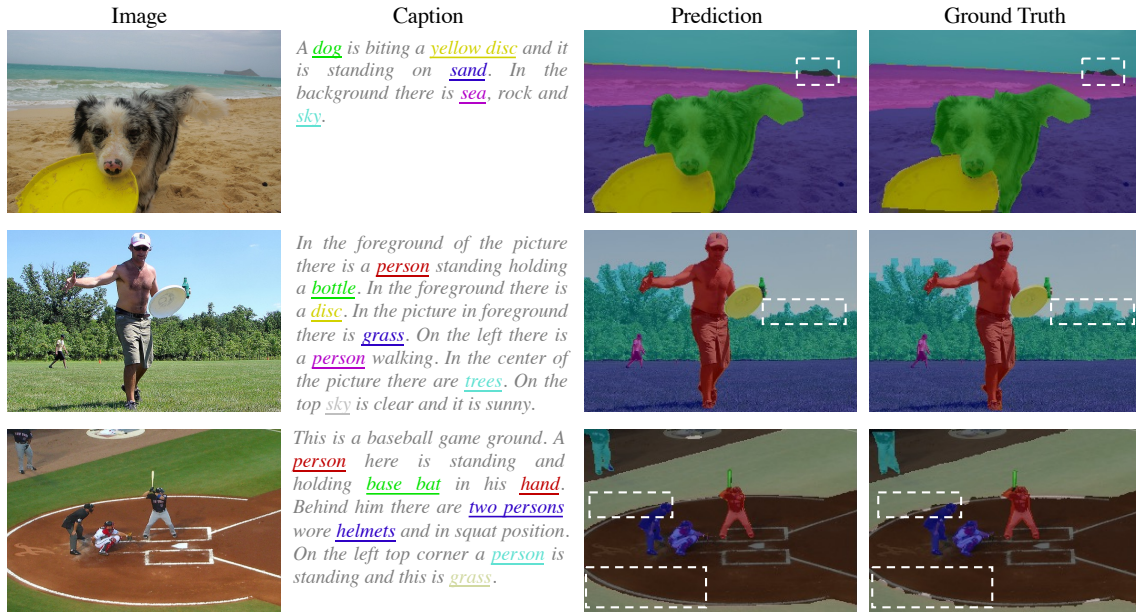


Figure 6: Qualitative analysis for our proposed PPMN. The regions where our PPMN produces better segmentation results than ground truth masks are highlighted within the white dashed boxes.

to explore sufficient cross-modal semantic relevance from densely annotated pixel-phrase pairs of the PNG benchmark. Moreover, it can liberate our model from the performance bottleneck caused by the indirect two-stage baseline and cumbersome manually designed training/inference pipelines. To further enhance the discriminative ability of noun phrase representations, we propose a Language-Compatible Pixel Aggregation (LCPA) module to adaptively aggregate representative visual context to each noun phrase from its corresponding language-compatible pixels for multiple rounds. Experiments show that our PPMN outperforms previous methods

by a large margin. In the future, we plan to extend our PPMN to other visual grounding tasks, such as phrase grounding, to further verify its generalization ability.

## 6 ACKNOWLEDGMENTS

This work was supported in part by the National Natural Science Foundation of China under Grant 62122010 and Grant 61876177, in part by the Fundamental Research Funds for the Central Universities, and in part by the Key Research and Development Program of Zhejiang Province under Grant 2022C01082.



## REFERENCES

- [1] Hassan Akbari, Svebor Karaman, Surabhi Bhargava, Brian Chen, Carl Vondrick, and Shih-Fu Chang. 2019. Multi-level multimodal common semantic space for image-phrase grounding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12476–12486.
- [2] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. 2015. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*. 2425–2433.
- [3] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. 2016. Layer normalization. *arXiv preprint arXiv:1607.06450* (2016).
- [4] Mohit Bajaj, Lanjun Wang, and Leonid Sigal. 2019. G3raphground: Graph-based language grounding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 4281–4290.
- [5] Hedi Ben-Younes, Rémi Cadene, Matthieu Cord, and Nicolas Thome. 2017. Mutan: Multimodal tucker fusion for visual question answering. In *Proceedings of the IEEE international conference on computer vision*. 2612–2620.
- [6] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. 2020. End-to-end object detection with transformers. In *European conference on computer vision*. Springer, 213–229.
- [7] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. 2020. Uniter: Universal image-text representation learning. In *European conference on computer vision*. Springer, 104–120.
- [8] Bowen Cheng, Maxwell D Collins, Yukun Zhu, Ting Liu, Thomas S Huang, Hartwig Adam, and Liang-Chieh Chen. 2020. Panoptic-DeepLab: A Simple, Strong, and Fast Baseline for Bottom-Up Panoptic Segmentation. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 12472–12482.
- [9] Bowen Cheng, Alex Schwing, and Alexander Kirillov. 2021. Per-pixel classification is not all you need for semantic segmentation. *Advances in Neural Information Processing Systems* 34 (2021).
- [10] Ying Cheng, Ruize Wang, Jiashuo Yu, Rui-Wei Zhao, Yuejie Zhang, and Rui Feng. 2021. Exploring Logical Reasoning for Referring Expression Comprehension. In *Proceedings of the 29th ACM International Conference on Multimedia*. 5047–5055.
- [11] Jiajun Deng, Zhengyuan Yang, Tianlang Chen, Wengang Zhou, and Houqiang Li. 2021. Transvg: End-to-end visual grounding with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 1769–1779.
- [12] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [13] Henghui Ding, Chang Liu, Suchen Wang, and Xudong Jiang. 2021. Vision-language transformer and query generation for referring segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 16321–16330.
- [14] Zihan Ding, Tianrui Hui, Junshi Huang, Xiaoming Wei, Jizhong Han, and Si Liu. 2022. Language-Bridged Spatial-Temporal Interaction for Referring Video Object Segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4964–4973.
- [15] Zihan Ding, Tianrui Hui, Shaofei Huang, Si Liu, Xuan Luo, Junshi Huang, and Xiaoming Wei. 2021. Progressive multimodal interaction network for referring video object segmentation. *The 3rd Large-scale Video Object Segmentation Challenge* (2021), 7.
- [16] Pelin Dogan, Leonid Sigal, and Markus Gross. 2019. Neural sequential phrase grounding (seqground). In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4175–4184.
- [17] Guang Feng, Zhiwei Hu, Lihe Zhang, and Huchuan Lu. 2021. Encoder Fusion Network with Co-Attention Embedding for Referring Image Segmentation. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 15501–15510.
- [18] Cristina González, Nicolás Ayobi, Isabela Hernández, José Hernández, Jordi Pont-Tuset, and Pablo Arbeláez. 2021. Panoptic Narrative Grounding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 1364–1373.
- [19] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. 2017. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*. 2961–2969.
- [20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [21] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9, 8 (1997), 1735–1780.
- [22] Ronghang Hu, Marcus Rohrbach, and Trevor Darrell. 2016. Segmentation from natural language expressions. In *European Conference on Computer Vision*. Springer, 108–124.
- [23] Shaofei Huang, Tianrui Hui, Si Liu, Guanbin Li, Yunchao Wei, Jizhong Han, Luoqi Liu, and Bo Li. 2020. Referring Image Segmentation via Cross-Modal Progressive Comprehension. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE Computer Society, 10485–10494.
- [24] Tianrui Hui, Shaofei Huang, Si Liu, Zihan Ding, Guanbin Li, Wenguan Wang, Jizhong Han, and Fei Wang. 2021. Collaborative spatial-temporal modeling for language-queried video actor segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4187–4196.
- [25] Tianrui Hui, Si Liu, Shaofei Huang, Guanbin Li, Sansi Yu, Faxi Zhang, and Jizhong Han. 2020. Linguistic structure guided context modeling for referring image segmentation. In *European Conference on Computer Vision*. Springer, 59–75.
- [26] Yang Jiao, Zequn Jie, Weixin Luo, Jingjing Chen, Yu-Gang Jiang, Xiaolin Wei, and Lin Ma. 2021. Two-stage Visual Cues Enhancement Network for Referring Image Segmentation. In *Proceedings of the 29th ACM International Conference on Multimedia*. 1331–1340.
- [27] Chenchen Jing, Yuwei Wu, Mingtao Pei, Yao Hu, Yunde Jia, and Qi Wu. 2020. Visual-semantic graph matching for visual grounding. In *Proceedings of the 28th ACM International Conference on Multimedia*. 4041–4050.
- [28] Aishwarya Kamath, Mannat Singh, Yann LeCun, Gabriel Synnaeve, Ishan Misra, and Nicolas Carion. 2021. MDETR-modulated detection for end-to-end multimodal understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 1780–1790.
- [29] Andrej Karpathy, Armand Joulin, and Li F Fei-Fei. 2014. Deep fragment embeddings for bidirectional image sentence mapping. *Advances in neural information processing systems* 27 (2014).
- [30] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [31] Alexander Kirillov, Ross Girshick, Kaiming He, and Piotr Dollár. 2019. Panoptic feature pyramid networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6399–6408.
- [32] Alexander Kirillov, Kaiming He, Ross Girshick, Carsten Rother, and Piotr Dollár. 2019. Panoptic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9404–9413.
- [33] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2012. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems* 25 (2012).
- [34] Liuwu Li, Yuqi Bu, and Yi Cai. 2021. Bottom-Up and Bidirectional Alignment for Referring Expression Comprehension. In *Proceedings of the 29th ACM International Conference on Multimedia*. 5167–5175.
- [35] Liunian Harold Li, Pengchuan Zhang, Haotian Zhang, Jianwei Yang, Chunyuan Li, Yiyu Zhong, Lijuan Wang, Lu Yuan, Lei Zhang, Jenq-Neng Hwang, et al. 2021. Grounded Language-Image Pre-training. *arXiv preprint arXiv:2112.03857* (2021).
- [36] Muchen Li and Leonid Sigal. 2021. Referring transformer: A one-step approach to multi-task visual grounding. *Advances in Neural Information Processing Systems* 34 (2021).
- [37] Qizhu Li, Xiaojuan Qi, and Philip HS Torr. 2020. Unifying training and inference for panoptic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 13320–13328.
- [38] Xiang Li, Wenhai Wang, Xiaolin Hu, and Jian Yang. 2019. Selective kernel networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 510–519.
- [39] Zhiqi Li, Wenhai Wang, Enze Xie, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, Tong Lu, and Ping Luo. 2021. Panoptic SegFormer. *arXiv preprint arXiv:2109.03814* (2021).
- [40] Yue Liao, Si Liu, Guanbin Li, Fei Wang, Yanjie Chen, Chen Qian, and Bo Li. 2020. A real-time cross-modality correlation filtering method for referring expression comprehension. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10880–10889.
- [41] Yue Liao, Aixi Zhang, Zhiyuan Chen, Tianrui Hui, and Si Liu. 2022. Progressive Language-customized Visual Feature Learning for One-stage Visual Grounding. *IEEE Transactions on Image Processing* (2022).
- [42] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. 2017. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2117–2125.
- [43] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *European conference on computer vision*. Springer, 740–755.
- [44] Si Liu, Tianrui Hui, Shaofei Huang, Yunchao Wei, Bo Li, and Guanbin Li. 2021. Cross-modal progressive comprehension for referring segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2021).
- [45] Xihui Liu, Zihao Wang, Jing Shao, Xiaogang Wang, and Hongsheng Li. 2019. Improving referring expression grounding with cross-modal attention-guided erasing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1950–1959.
- [46] Yongfei Liu, Bo Wan, Xiaodan Zhu, and Xuming He. 2020. Learning cross-modal context graph for visual grounding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 11645–11652.
- [47] Jonathan Long, Evan Shelhamer, and Trevor Darrell. 2015. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 3431–3440.
- [48] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *Advances in neural information processing systems* 32 (2019).

- [49] Gen Luo, Yiyi Zhou, Rongrong Ji, Xiaoshuai Sun, Jinsong Su, Chia-Wen Lin, and Qi Tian. 2020. Cascade grouped attention network for referring expression segmentation. In *Proceedings of the 28th ACM International Conference on Multimedia*. 1274–1282.
- [50] Gen Luo, Yiyi Zhou, Xiaoshuai Sun, Liujuan Cao, Chenglin Wu, Cheng Deng, and Rongrong Ji. 2020. Multi-Task Collaborative Network for Joint Referring Expression Comprehension and Segmentation. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 10031–10040.
- [51] Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi. 2016. V-net: Fully convolutional neural networks for volumetric medical image segmentation. In *2016 fourth international conference on 3D vision (3DV)*. IEEE, 565–571.
- [52] Zongshen Mu, Siliang Tang, Jie Tan, Qiang Yu, and Yueting Zhuang. 2021. Disentangled motif-aware graph learning for phrase grounding. In *Proc 35th AAAI Conf on Artificial Intelligence*.
- [53] Vinod Nair and Geoffrey E Hinton. 2010. Rectified linear units improve restricted boltzmann machines. In *Icml*.
- [54] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems* 32 (2019).
- [55] Bryan A Plummer, Paige Kordas, M Hadi Kiapour, Shuai Zheng, Robinson Piramuthu, and Svetlana Lazebnik. 2018. Conditional image-text embedding networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 249–264.
- [56] Jordi Pont-Tuset, Jasper Uijlings, Soravit Changpinyo, Radu Soricut, and Vittorio Ferrari. 2020. Connecting vision and language with localized narratives. In *European Conference on Computer Vision*. Springer, 647–664.
- [57] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. 2016. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 779–788.
- [58] Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. 2019. VL-BERT: Pre-training of Generic Visual-Linguistic Representations. In *International Conference on Learning Representations*.
- [59] Zhi Tian, Chunhua Shen, and Hao Chen. 2020. Conditional convolutions for instance segmentation. In *European Conference on Computer Vision*. Springer, 282–298.
- [60] Zhi Tian, Bowen Zhang, Hao Chen, and Chunhua Shen. 2022. Instance and panoptic segmentation using conditional convolutions. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2022).
- [61] Ledyard R Tucker. 1966. Some mathematical notes on three-mode factor analysis. *Psychometrika* 31, 3 (1966), 279–311.
- [62] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
- [63] Huiyu Wang, Yukun Zhu, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. 2021. Max-deeplab: End-to-end panoptic segmentation with mask transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5463–5474.
- [64] Liwei Wang, Yin Li, Jing Huang, and Svetlana Lazebnik. 2018. Learning two-branch neural networks for image-text matching tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 41, 2 (2018), 394–407.
- [65] Peng Wang, Qi Wu, Jiewei Cao, Chunhua Shen, Lianli Gao, and Anton van den Hengel. 2019. Neighbourhood watch: Referring expression comprehension via language-guided graph attention networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1960–1968.
- [66] Xinlong Wang, Rufeng Zhang, Tao Kong, Lei Li, and Chunhua Shen. 2020. Solov2: Dynamic and fast instance segmentation. *Advances in Neural information processing systems* 33 (2020), 17721–17732.
- [67] Yuxin Wu and Kaiming He. 2018. Group normalization. In *Proceedings of the European conference on computer vision (ECCV)*. 3–19.
- [68] Yuxin Wu, Alexander Kirillov, Francisco Massa, Wan-Yen Lo, and Ross Girshick. 2019. Detectron2. <https://github.com/facebookresearch/detectron2>
- [69] Yuwen Xiong, Renjie Liao, Hengshuang Zhao, Rui Hu, Min Bai, Ersin Yumer, and Raquel Urtasun. 2019. Upsnet: A unified panoptic segmentation network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8818–8826.
- [70] Sibe Yang, Guanbin Li, and Yizhou Yu. 2019. Dynamic graph attention for referring expression comprehension. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 4644–4653.
- [71] Sibe Yang, Guanbin Li, and Yizhou Yu. 2020. Propagating over phrase relations for one-stage visual grounding. In *European Conference on Computer Vision*. Springer, 589–605.
- [72] Zhengyuan Yang, Tianlang Chen, Liwei Wang, and Jiebo Luo. 2020. Improving one-stage visual grounding by recursive sub-query construction. In *European Conference on Computer Vision*. Springer, 387–404.
- [73] Zhengyuan Yang, Boqing Gong, Liwei Wang, Wenbing Huang, Dong Yu, and Jiebo Luo. 2019. A fast and accurate one-stage approach to visual grounding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 4683–4693.
- [74] Zhao Yang, Jiaqi Wang, Yansong Tang, Kai Chen, Hengshuang Zhao, and Philip HS Torr. 2021. LAVT: Language-Aware Vision Transformer for Referring Image Segmentation. *arXiv preprint arXiv:2112.02244* (2021).
- [75] Jiabo Ye, Xin Lin, Liang He, Dingbang Li, and Qin Chen. 2021. One-Stage Visual Grounding via Semantic-Aware Feature Filter. In *Proceedings of the 29th ACM International Conference on Multimedia*. 1702–1711.
- [76] Licheng Yu, Zhe Lin, Xiaohui Shen, Jimei Yang, Xin Lu, Mohit Bansal, and Tamara L Berg. 2018. MAttNet: Modular Attention Network for Referring Expression Comprehension. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*. IEEE, 1307–1315.
- [77] Tianyu Yu, Tianrui Hui, Zhihao Yu, Yue Liao, Sansi Yu, Faxi Zhang, and Si Liu. 2020. Cross-modal omni interaction modeling for phrase grounding. In *Proceedings of the 28th ACM International Conference on Multimedia*. 1725–1734.
- [78] Hanwang Zhang, Yulei Niu, and Shih-Fu Chang. 2018. Grounding referring expressions in images by variational context. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 4158–4166.
- [79] Tianhao Zhang, Hung-Yu Tseng, Lu Jiang, Weilong Yang, Honglak Lee, and Irfan Essa. 2021. Text as neural operator: Image manipulation by text instruction. In *Proceedings of the 29th ACM International Conference on Multimedia*. 1893–1902.
- [80] Weixia Zhang, Chao Ma, Qi Wu, and Xiaokang Yang. 2020. Language-guided navigation via cross-modal grounding and alternate adversarial learning. *IEEE Transactions on Circuits and Systems for Video Technology* 31, 9 (2020), 3469–3481.
- [81] Wenwei Zhang, Jiangmiao Pang, Kai Chen, and Chen Change Loy. 2021. K-net: Towards unified image segmentation. *Advances in Neural Information Processing Systems* 34 (2021).
- [82] Chen Zheng, Quan Guo, and Parisa Kordjamshidi. 2020. Cross-modality relevance for reasoning on language and vision. *arXiv preprint arXiv:2005.06035* (2020).
- [83] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. 2019. Objects as points. *arXiv preprint arXiv:1904.07850* (2019).