

Context Does Matter: End-to-end Panoptic Narrative Grounding with Deformable Attention Refined Matching Network

Yiming Lin, Xiao-Bo Jin[†], Qiufeng Wang
School of Advanced Technology
Xi'an Jiaotong-Liverpool University, Suzhou, China
{yiming.lin21@student, xiaobo.jin@, qiufeng.wang@}xjtlu.edu.cn

Kaizhu Huang
Data Science Research Center
Duke Kunshan University, Suzhou, China
kaizhu.huang@dukekunshan.edu.cn

Abstract—Panoramic Narrative Grounding (PNG) is an emerging visual grounding task that aims to segment visual objects in images based on dense narrative captions. The current state-of-the-art methods first refine the representation of phrase by aggregating the most similar k image pixels, and then match the refined text representations with the pixels of the image feature map to generate segmentation results. However, simply aggregating sampled image features ignores the contextual information, which can lead to phrase-to-pixel mis-match. In this paper, we propose a novel learning framework called Deformable Attention Refined Matching Network (DRMN), whose main idea is to bring deformable attention in the iterative process of feature learning to incorporate essential context information of different scales of pixels. DRMN iteratively re-encodes pixels with the deformable attention network after updating the feature representation of the top- k most similar pixels. As such, DRMN can lead to accurate yet discriminative pixel representations, purify the top- k most similar pixels, and consequently alleviate the phrase-to-pixel mis-match substantially. Experimental results show that our novel design significantly improves the matching results between text phrases and image pixels. Concretely, DRMN achieves new state-of-the-art performance on the PNG benchmark with an average recall improvement 3.5%. The codes are available in: <https://github.com/JaMesLiMers/DRMN>.

Index Terms—Visual Grounding, Panoptic Narrative Grounding, One-stage Method

I. INTRODUCTION

Panoptic Narrative Grounding (PNG) [1], one emerging visual grounding task, has recently drawn great attention in data mining and computer vision including grounded context recognition [2], visual question answering [3], and visual-language model pre-training [4]. Given an image and its associated dense narrative caption, the goal of PNG is to segment the visuals of things and stuff based on the visuals mentioned in the caption (see illustration in Fig. 1). In contrast to other related tasks, PNG extends the grounding range from the bounding box of the foreground class (called “object”) to a segmentation mask containing both foreground

and background classes (named “object” and “Stuff”), thus defining the finest-grained alignment between multiple noun phrases and segments. A detailed comparison between PNG and other related vision-based tasks can be seen in Sect. II.



Fig. 1: Illustration of the PNG problem: Given an image (left) and corresponding caption (middle), the goal is to generate a panoptic segmentation (right) based on all visual objects contained in the caption (i.e., labeling each object and its associated segmented region with the same color).

In general, there are two families of methods for PNG. The first type of methods typically exploits a two-stage pipeline [1], which matches by computing the affinity matrix between object proposals (extracted by off-the-shelf models) and noun phrases. As such, the object proposal model, i.e., the off-the-shelf model will limit the performance ceiling. On the other hand, the one-stage or end-to-end methods [5], [6] alleviate this problem by directly generating a response map between all noun phrases and image pixels. To better fuse information from different modalities, Ding et al. [5] propose a language-compatible pixel aggregation (LCPA) module to aggregate the most compatible features from images to noun phrases. Namely, taking each noun phrase as a query feature, LCPA samples top- k image-compatible features, which are then used as key and value features. Finally the multi-head cross-modal attention is adopted to aggregate visual features.

Albeit its promising performance, LCPA simply aggregates sampled image features without taking into account the contextual information, which could lead to serious phrase-to-pixel mis-match. Concretely, LCPA tends to push the phrase feature towards the center of top- k sampled image features. This algorithm works well when the top- k sampled features are related to the target visual object with high similarity.

This work was partially supported by “Qing Lan Project” in Jiangsu universities, NSFC under No. 62106081, 62376113, 92370119, RDF with No. RDF-22-01-020, Jiangsu Science and Technology Programme under No. BE2020006-4 and Natural Science Foundation of the Jiangsu Higher Education Institutions of China under No. 22KJB520039

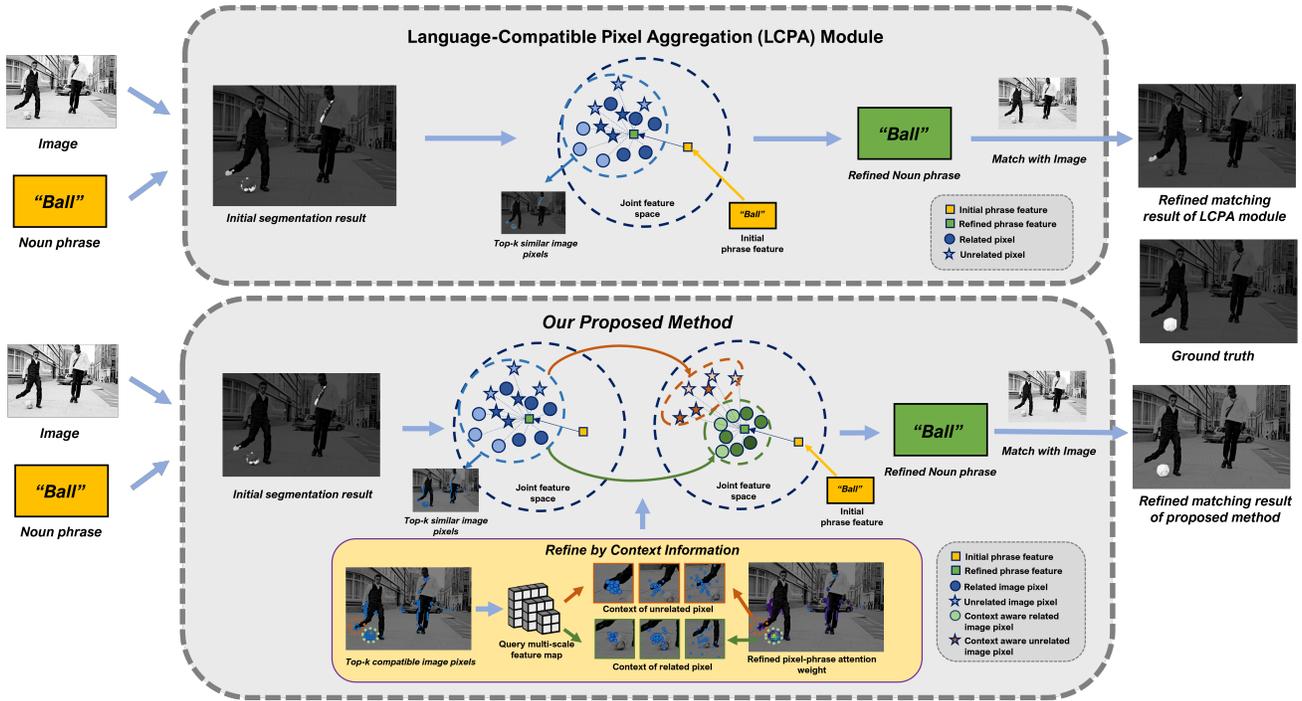


Fig. 2: Insight of our proposed method. The upper part addresses the limitation of LCPA with a hard example. We introduce the essential context information in the multi-scale feature map as a cue to refine the sampled top- k image feature. By sharpening the representation of different visual objects points, our method can filter out the sampled points by purifying irrelevant visual objects, which further enhances the final segmentation result.

However, such strategy may also inevitably introduce unrelated image pixels. As illustrated in a hard example of Fig. 2, the unrelated object’s pixels (for “feet”) fused with the related object’s pixels (for “ball”) dominate the top- k sampled features with high similarity, thus inducing a serious mis-match. To alleviate this problem, we argue that relevant context information is crucial for differentiating and purifying the top- k sampled points. Essentially, integration of context information would enable more accurate and discriminate pixel representation, since related pixels enjoy more similar context whilst unrelated ones may share distinctive context. In other words, integration of context information could sharpen the representation of different visual pixels, thereby providing potentials to improve the segmentation performance.

Motivated from the above observations and inspired by the current object detection method [7], we design a novel deformable attention mechanism to extract essential pixel contextual information in multi-scale feature maps in an iterative way, resulting in a simple yet effective end-to-end model, called Deformable Attention Fine Matching Network (DRMN). An overview of its framework is shown in Fig. 3. Similar to Transformer, DRMN is a multi-scale encoding-decoding method that offers a full range of context information at different scales, including global information and local information. For feature extraction, we engage the deformable attention to encode image features at different levels, which will generate a cross-fused multi-scale word-pixel matching

matrix to obtain initial word-pixel matching results. In addition, in the feature aggregation stage, we further incorporate word embedding representations into pixel encoding representations. Specifically, we follow insights from the DETR decoder [8] to refine the sampled features: for each word vector, our model queries its nearest k pixels and applies an attention mechanism to encode them after updating the vector representations of these pixels with the word vector.

Our contributions are four-fold:

- From the perspective of multimodal information fusion, we design a deformable attention model with multi-scale encoding and decoding functions in the aggregation process of pixel features to better encode the context information around the pixel, effectively alleviating the phrase-to-pixel mis-match problem.
- From the perspective of model structure, unlike the existing DETR-based models [8], [9], our novel design leverages directly multimodal transformers for highly intertwined feature aggregation. Inheriting insights from DETR feature aggregation, our new approach offers a more sparse and interpretable way of exploiting DETR for vision-based tasks.
- From an algorithmic perspective, we simplify the multi-round pixel feature refinement process into an iterative process of two subproblems: the fuzzy K-means clustering subproblem and the multi-objective assignment subproblem, the latter of which can be efficiently solved

by online gradient descent.

- From an experimental point of view, the results of multiple categories and the overall results on the public PNG benchmark show the superiority of our method compared to previous methods, where the average recall rate is 3.5% higher than the second-ranked method.

II. RELATED WORK

In this section, we overview PNG in contrast to different related vision-based tasks. Overall, Table I shows the comparison granularities of related vision-based tasks, among which the PNG task provides the most fine-grained alignment between different types of nouns and segmentation.

TABLE I: Comparison of different granularities of related vision-based tasks. Considering the typical segmentation categories in computer vision tasks between things (countable objects) and stuff (amorphous regions of similar texture), the datasets of the other tasks mainly focus on things categories.

Grounding Task	Language Granularity	Visual Granularity	Semantic Generality
REC [10]	Short phrase	Bounding box	Things
PG [11]	Noun phrase	Bounding box	Mainly Things
RES [12]	Short phrase	Segmentation	Things
PNG [1]	Noun phrase	Segmentation	Things + Stuff

A. Visual Grounding with Bounding Box Regression

The goal of Referent Expression Comprehension (REC) task is to predict the corresponding bounding box in an image for a given referring expression. Current methods can be categorized into two-stage and one-stage approaches. Two-stage methods [13] first propose bounding box proposals in the image, then match the proposal-referring expression pairs. Inspired by object detection techniques [14], the one-stage methods [15] directly generate results based on the input textual information without explicit matching. Recently, some methods have explored multi-modal pre-training models [16] in REC [17], taking architectures similar to BERT, to obtain joint representations of images and texts.

The phrase grounding task aims to find the corresponding bounding box in an image for multiple noun phrases mentioned in an input caption. Early methods [11], [18] adopted representation learning, which first project region proposal and phrase embeddings onto a same subspace, then learn semantic similarity between them. In recent years, researchers have explored various methods [19]–[21] for fusing and learning multi-modal features. It is worth mentioning that recent large-scale visual-language pre-training models [22], [23] have adopted weakly supervised phrase grounding [24] loss to align image-noun phrase pairs.

Recently, some methods [8], [9], [25] have modified the transformer-based object detection framework to address the aforementioned bounding box regression problems. TransVG [25] first proposed a pure transformer framework for visual grounding tasks. Furthermore, some methods [8], [9]

drew inspirations from the DETR object detection framework. MDETR [8] employed a transformer encoder-decoder structure, where the transformer simultaneously extracted features from both the image and text in the encoder, and introduced QA-specific queries in the decoder for visual grounding-related task decoding. Dynamic MDETR [9] engaged the idea of deformable attention in the decoder to reduce computation. It is worth noting that our approach differs from the aforementioned methods. Instead of using a transformer to simultaneously encode image and text information, we handle the multi-modal feature interaction in the decoder through top- k sampling. In particular, in the decoder, we consider the features of the top- k image positions as object queries, and design the deformable attention mechanism to extract object-relevant features and then aggregate the extracted object query features into the textual features.

B. Visual Grounding with Segmentation

The task of Referent Expression Segmentation (RES) is to generate a segmentation map of the referred object according to the input referring expression. The first proposed method [12] on this task is a one-stage model that first concatenated textual features and global image feature, then decoded the segmentation mask through deconvolution layers. Recently, inspired by multi-modal transformers, various fine-grained modeling approaches [26], [27] have been proposed to facilitate interactions between different modalities. For example, SHNET [28] concatenated textual features with different levels of image features as joint input of the transformer, then adopted language features to guide the information exchange between different levels of image features. The LAVT model [27] developed the PWAM module, which directly used attention to expand textual information to the size of the image feature map for pixel-word feature fusion.

PNG aims to segment corresponding things or stuff in an image based on the multiple noun phrases mentioned in the image caption. This task was initially proposed with a two-stage method by González et al. [1], along with a dataset. They extracted segmentation proposals from off-the-shelf models which were matched with the extracted noun phrase features. Later, some work explored the one-stage paradigm. For example, PPMN [5] achieved feature fusion between different modalities through a sampling strategy end-to-end. EPNG [6] further optimized the inference speed and achieved real-time segmentation effects while sacrificing little accuracy.

III. MAIN METHOD

In this section, we first introduce the process of feature extraction for image text (III-A). Then, we describe how initial segmentation results are generated (III-C). Subsequently, we present our proposed Multi-round Visual-Language Aggregation Module, which selectively aggregates image features into textual features to enhance the model’s performance (III-D). Finally, we detail the loss function and introduce the training process of the model (III-E). The overall workflow of our

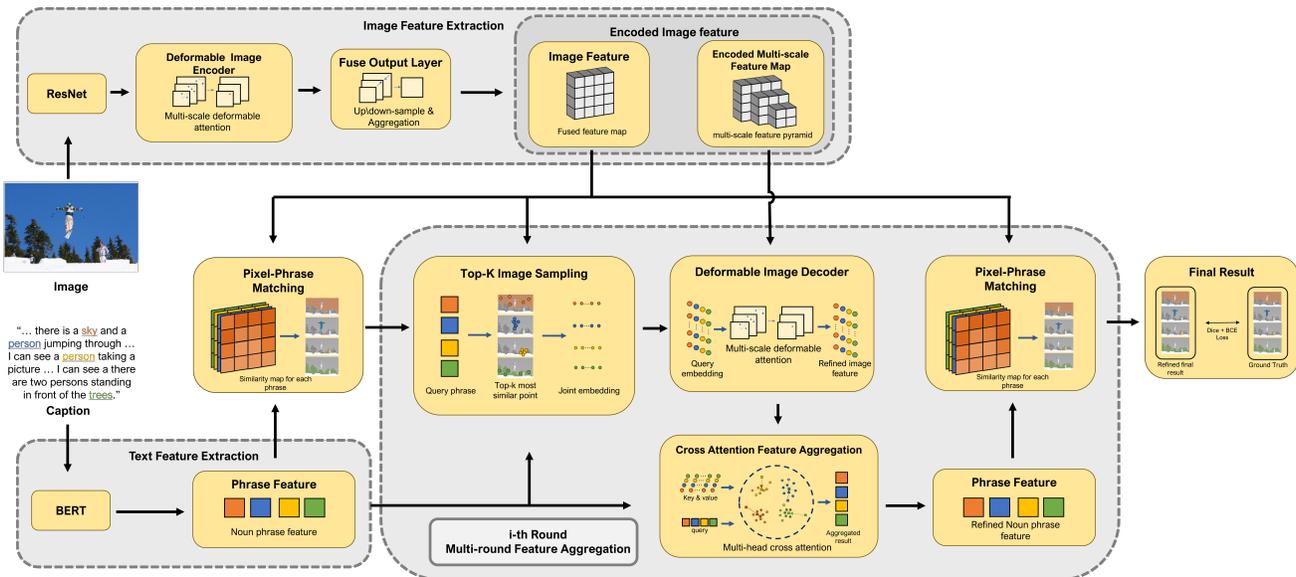


Fig. 3: Overview of our model. We integrate essential image context information in feature extraction and multi-round feature aggregation phases with deformable attention. First, we utilize BERT to encode textual features and employ deformable attention to encode multi-scale image feature maps. Furthermore, we generate initial image-text matching results based on textual and image features. Finally, in the multi-round feature aggregation, we aggregate the top- k image features into text feature based on the matching results. The model utilizes deformable attention to refine sampled image features further, then aggregates the refined features into textual features through a cross-attention mechanism to generate improved matching results.

model is shown in Fig. 3 and the pseudo code of the algorithm is shown in Alg. 1.

A. Feature Extraction

In the feature extraction stage, we can employ off-the-shelf methods to extract features from visual and linguistic modalities.

For the text modality, we leverage BERT to extract features for each word in the image caption. Specifically, we focus on extracting features from the noun phrase part T of the title $\mathcal{G} = \text{BERT}(T) \in R^{n \times d}$, where n represents the maximum number of words in all input noun phrases and d denotes the dimensionality of the textual feature embedding representation.

As for image modality, given any image I , we use ResNet as the backbone to extract multi-scale feature pyramid $S = \{F_2, F_3, F_4, F_5\}$ such as

$$F_l = \text{flatten}(\text{ResNet}(I, l)), \quad l = 2, 3, 4, 5. \quad (1)$$

Here I is an RGB image with height h and width w and the l -th scale output will be a matrix of size $\frac{hw}{4^l} \times c$, where c is the number of channels in the output map.

In the subsequent feature extraction stage, first we normalize the coordinates x 's and y 's of all points in the feature map of each scale to the range $[0, 1]$. We then grid them to obtain the reference point matrix with the same size $p_l = \text{flatten}(\text{grid}(M_l)) \in R^{\frac{hw}{4^l} \times 2}$, where the size of the tensor M is $\frac{h}{2} \times \frac{w}{2} \times 2$ and its elements $M_l(i, j) = (i, j)$.

Meanwhile, we add the feature maps of multiple scales obtained by the feature pyramid and their position codes and straighten them to get the matrix of $\frac{hw}{4^l} \times c$

$$\hat{F}_l = F_l + \text{pos}(F_l), \quad l = 2, 3, 4, 5, \quad (2)$$

where $\text{pos}(\cdot)$ represents the positional encoding function. We take these feature representations and concatenate them row by row into a matrix as the input of the DeformLayer function in the initial stage

$$\hat{F} = \text{catrow}(\{\hat{F}_l\}). \quad (3)$$

B. Deformable Layer

In order to better integrate feature information at different scales, we use multiple deformable attention layers to aggregate information from various levels of the feature pyramid according to the position and characteristics of the input feature points, where the input of deformable attention is the query features $Q = \hat{F}$, the reference positions $P = \{p_2, p_3, p_4, p_5\}$ corresponding to the feature point in the image and multi-scale feature map $S = \{F_2, F_3, F_4, F_5\}$.

Similar to Transformer, on multiple feature maps of the pyramid, each feature point will be re-expressed as a linear combination of other feature points. The difference is that its value V is based on the feature map with a bilinear sampling operation instead of itself. Moreover, both the sampling offset

and the self-attention correlation coefficient depend on the query Q , specifically described as follows

$$\Delta p_l = \hat{F}W_l^p \quad (4)$$

$$V_l = \phi_{\text{bilin}}(F_l, p_l + \Delta p_l) \quad (5)$$

$$V = \text{catrow}(\{V_l\}) \quad (6)$$

$$\hat{V} = \text{softmax}(\hat{F}W^a)VW^v, \quad (7)$$

where $l = 2, 3, 4, 5$ means transformation on multiple scales, W_l^p , W^v and W^a represent the linear mapping to be learned, and $\phi_{\text{bilin}}(F_l, p_l + \Delta p_l)$ indicates that pixels are sampled by bilinear interpolation on the position $p_l + \Delta p_l$ of the feature map F_l . Similarly, multiple heads are introduced to obtain image feature representations with multiple attentions, and these representations are concatenated and linearly mapped to obtain a multi-scale refined representation \mathcal{V} .

Subsequently, we can construct a deformable coding layer with deformable attention representation

$$\mathcal{F} = \text{FFN}(\text{norm}(\mathcal{V} + \text{dropout}(\mathcal{V}))), \quad (8)$$

where FFN, norm and dropout indicate feedforward network layer, normalization layer, and dropout layer respectively.

For convenience, we represent the above whole process as the following function

$$\mathcal{F} = \text{DeformLayer}(\hat{F}, \{F_l\}, \{p_l\}). \quad (9)$$

C. Image and Text Matching

Below we describe how to match text embedding \mathcal{G} with multiple feature maps \mathcal{F}_l of different scales to obtain multiple similarity matrices.

In the initial stage, we concatenate multi-scale feature maps by row as DeformLayer function to get their attention representation, and then we restore them into multiple feature maps \mathcal{F}_l (line 8-12 in Alg. 1). Then, using the third layer as a benchmark, all the other layers are down-sampled or up-sampled to the same size as the feature map of the third layer

$$\bar{F}_l = \phi_{\text{sampling}}(\text{reshape}(\mathcal{F}_l), 2^{l-3}), \quad l = 2, 3, 4, 5. \quad (10)$$

In this way, we can fuse the feature map output with the same scale, and convert the fused image output into vectors to facilitate matching with the text vector (line 13 in Alg. 1)

$$\mathbb{F} = \text{vect} \left(\frac{1}{4} \sum_{l=2}^5 \bar{F}_l \right), \quad (11)$$

where $\text{vect}(\cdot)$ means pulling a three-dimensional tensor $\mathbb{F} \in R^{(h/8) \times (w/8) \times c}$ into a two-dimensional vector $\mathbb{F} \in R^{(hw/8^2) \times c}$.

Now, we map the representation G of phrases to the feature space of pixels to compute their similarity matrix (line 15 in Alg. 1)

$$\hat{G} = GV^g, \quad H = \text{sigmoid}(\hat{G}\mathbb{F}^T), \quad (12)$$

where V^g is a projection matrix and sigmoid is the sigmoid function.

It is worth noting that during the iterative process, we directly use the latest representation \hat{F} of the pixel to calculate the similarity matrix (line 25 in Alg. 1)

$$H = \text{sigmoid}(\hat{G}\hat{F}^T). \quad (13)$$

D. Multi-round Feature Aggregation Module

Algorithm 1 Multi-round Feature Aggregation

```

1: Input: Image  $I$  and caption  $T$ 
2:  $\mathcal{G} = \text{BERT}(T)$ 
3: for  $l = 2, 3, 4, 5$  do
4:    $p_l = \text{flatten}(\text{grid}(M_l))$ 
5:    $F_l = \text{flatten}(\text{ResNet}(I, l))$ 
6:    $\hat{F} = F_l + \text{pos}(F_l)$ 
7: end for
8:  $\hat{F} = \text{catrow}(\{\hat{F}_l\})$ 
9: for  $t = 1, 2, \dots, T$  do
10:   $\hat{F} = \text{DeformLayer}(\hat{F}, \{F_l\}, \{p_l\})$ 
11: end for
12:  $\{\mathcal{F}_l\} = \text{splitrow}(\hat{F})$ 
13:  $\hat{F} = \text{avg}(\{\mathcal{F}_l\})$ 
14:  $\hat{\mathcal{G}} = \mathcal{G}V^g$ 
15:  $H = \text{sigmoid}(\hat{\mathcal{G}}\hat{F}^T)$ 
16:  $\mathcal{H} = []$ 
17: for  $i = 1, 2, \dots, I$  do
18:   $S = \text{topk}(H, k)$ 
19:  for  $j = 1, 2, \dots, n$  do
20:     $s = S[j, :]$ 
21:     $\hat{F}[s, :] = \hat{F}[s, :] + \text{pos}(\hat{F}[s, :]) + \hat{\mathcal{G}}[j, :]$ 
22:     $\hat{F}[s, :] = \text{DeformLayer}(\hat{F}[s, :], \{\mathcal{F}_l\}, \{p_l\})$ 
23:     $\hat{\mathcal{G}}[j, :] = \text{CrossAttention}(\hat{\mathcal{G}}[j, :], \hat{F}[s, :], \hat{F}[s, :])$ 
24:  end for
25:   $H = \text{sigmoid}(\hat{\mathcal{G}}\hat{F}^T)$ 
26:   $\mathcal{H}.\text{append}(\phi_{\text{sampling}}(H, 2^{-3}))$ 
27: end for
28: return  $\mathcal{H}$ 

```

Alg. 1 shows the entire process of our multi-round feature aggregation: initially establish the relationship between text and pixels, and then refine these relationships through continuous iteration.

First, we select k pixels with the highest similarity for each row on the multi-scale similarity matrix $S = \text{topk}(H, k)$, where S is a matrix of dimension $n \times k$.

In the refinement phase (lines 12-24), we first update the embedded representations of the k nearest image pixels each time with the text's representation based on the previous iteration (line 20-21 in Alg. 1). Next, we apply **DeformLayer** again to regenerate the multi-scale representation of pixels for the top- k image positions (line 22 in Alg. 1). Notably, we re-use the multi-scale output $\{\mathcal{F}_l\}$ of **DeformLayer** from the initial stage as its input.

Subsequently, we update the representation of noun phrases using the weighted sum of the current top- k image features

(line 23 in Alg. 1). Below we will describe its implementation in detail.

Given a query Q , a key K and a value V , through the attention we can get Q 's updated weighted representation

$$\text{attn}(Q, K, V) = \text{softmax} \left(\frac{QW_q(KW_k)^T}{\sqrt{c}} \right) VW_v. \quad (14)$$

Here W_q , W_k and W_v represent the projection matrices, which project a row vector to the c -dimensional space.

We treat each row of \hat{G} as a query, \hat{F} as key and value, and split them into M blocks along the dimension of representation. The multi-head attention representation of $\hat{G}[j, :]$ can be computed ($s = S[j, :]$)

$$G[j, :] = \text{catcol}(\{\text{attn}(\hat{G}[j, \text{ids}(i)], \hat{F}[s, \text{idx}(i)], \hat{F}[s, \text{idx}(i)])\}).$$

where $j = 1, 2, \dots, n$ and catcol represents the concatenation along the column and $\text{idx}(i)$ represents the column index set of the i -th sub-block. Then we sequentially perform addition, dropout, norm, and FFN operations on G and \hat{G}

$$\bar{G} = \text{dropout}(G + \hat{G}), \quad (15)$$

$$\hat{G} = \text{FFN}(\text{norm}(\hat{G} + \bar{G})). \quad (16)$$

E. Loss Function

Once we have a series of predicted values \mathcal{H} of the correlation coefficient of text and pixel, we can define the optimized loss function based on Binary Cross Entropy (BCE) and Dice loss according to the true value Y

$$\mathcal{L}(\mathcal{H}, Y) = \sum_{i=1}^I \lambda_{bce} \mathcal{L}_{bce}(\mathcal{H}_i, Y) + \lambda_{dice} \mathcal{L}_{dice}(\mathcal{H}_i, Y), \quad (17)$$

where λ_{bce} and λ_{dice} are the weight coefficients of the loss which are both set to 1 in our experiments. Specifically, BCE loss is the average loss of all text-pixel pairs

$$\mathcal{L}_{bce}(\mathcal{H}_i, Y) = \frac{1}{nhw} \sum_{j=1}^n \sum_{k=1}^{hw} \text{CE}(Y(j, k), \mathcal{H}_i(j, k)), \quad (18)$$

where CE is the cross entropy loss.

In general, the goal of BCE loss is to compute a binary classification loss for all pixels, but this loss does not consider the problem of class imbalance. To alleviate this problem, we introduce the Dice loss as an additional loss

$$\mathcal{L}_{dice}(\mathcal{H}_i, Y) = \frac{1}{n} \sum_{j=1}^n \left(1 - \frac{2 \sum_{k=1}^{hw} \mathcal{H}_i(j, k) Y(j, k)}{\sum_{k=1}^{hw} \mathcal{H}_i(j, k) + Y(j, k)} \right).$$

To provide sufficient intermediate supervision during the encoding stage, we follow the setup of [5], which applies the loss \mathcal{L} to the predicted value H of all refinement stages. For the inference, we obtain grounded results from the previous round of response maps with a threshold of 0.5.

F. Discussion on Multi-round Feature Aggregation

In our task, we are given an embedded representation t_j of n words, and the goal is to assign all m pixels x_i in the image to these n noun phrases. For convenience, we assume that x_i and t_j are located in a common vector space. We then iteratively optimize the representations of x_i s and t_j s through the objective function $\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^n u_{ij}^2 \|x_i - t_j\|^2$, where u_{ij} represents the probability that x_i belongs to t_j .

1) *Solving x_i 's for known t_j 's*: We can define the following loss function

$$\min_{u, x} \quad \mathcal{L}(u, x) = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^n u_{ij}^2 \|x_i - t_j\|^2, \quad (19)$$

$$\text{s.t.} \quad \sum_{i=1}^m u_{ij} = k, \quad u_{ij} \in \{0, 1\}, \quad k < n. \quad (20)$$

We add the constraint $\sum_{i=1}^m u_{ij} = k$ and $k < n$ to avoid trivial solutions. Obviously, if $m = n$, then $x_j = t_j$.

Assume that for each target point t_j , its k closest points are $x_{r(j,1)}, x_{r(j,2)}, \dots, x_{r(j,k)}$. If we fix x_i to find the optimal point of u_{ij} , we have

$$u_{ij} = \begin{cases} 1, & i \in \{r(j, 1), r(j, 2), \dots, r(j, k)\}, \\ 0, & \text{otherwise.} \end{cases} \quad (21)$$

Next, if we fix u_{ij} , we then calculate the gradient of $\mathcal{L}(u, x)$ with respect to x to get

$$\frac{\partial \mathcal{L}}{\partial x_i} = \sum_{j=1}^n u_{ij}^2 x_i - \sum_{j=1}^n u_{ij}^2 t_j. \quad (22)$$

Hence, we get the update formula of x_i

$$x_i = x_i - \alpha \frac{\partial \mathcal{L}}{\partial x_i} = (1 - \alpha \sum_{j=1}^n u_{ij}^2) x_i + \sum_{j=1}^n \alpha u_{ij}^2 t_j, \quad (23)$$

where α ($0 < \alpha < 1$) is a step size. Note that the above equation is a batch processing method for n targets, and its online update method for target t_j can be given as

$$x_i = (1 - \alpha u_{ij}^2) x_i + \alpha u_{ij}^2 t_j, \quad (24)$$

which can be further simplified to

$$x_i = \begin{cases} (1 - \alpha) x_i + \alpha t_j, & x_i \in \text{topk}(t_j), \\ x_i, & \text{otherwise.} \end{cases} \quad (25)$$

Here $x_i \in \text{topk}(t_j)$ means $i \in r(j, 1), r(j, 2), \dots, r(j, k)$. Note that u_{ij} is fixed, the optimization problem is a strictly convex optimization problem about x . Therefore, setting an appropriate step size can ensure that the function value decreases after each gradient descent.

In line 21 of Alg. 1, we update x_i with the following formula

$$x_i = \begin{cases} f(x_i) + t_j, & x_i \in \text{topk}(t_j), \\ x_i, & \text{otherwise.} \end{cases} \quad (26)$$

Here $f(x_i)$ represents the encoded representation of x_i . At the same time, when we solve the nearest k points from t_j , we exploit the predicted correlation coefficient instead of the Euclidean distance.

2) *Solving t_i 's for known x_i 's*: Drawing the idea of fuzzy K-means, we define the following loss function

$$\min_{u,t} \quad \mathcal{L}(u,t) = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^n u_{ij}^2 \|x_i - t_j\|^2, \quad (27)$$

$$s.t. \quad \sum_{j=1}^n u_{ij} = 1, \quad i = 1, 2, \dots, m, \quad (28)$$

whose Lagrangian function is

$$\mathcal{J}(u,t,\lambda) = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^n u_{ij}^2 \|x_i - t_j\|^2 + \sum_{i=1}^m \lambda_i \left(\sum_{j=1}^n u_{ij} - 1 \right).$$

According to KKT, we obtain the optimal u_{ij} and t_j satisfying

$$u_{ij} = \frac{1/\|x_i - t_j\|^2}{\sum_{k=1}^n 1/\|x_i - t_k\|^2}, \quad (29)$$

$$t_j = \frac{\sum_{i=1}^n u_{ij}^2 x_i}{\sum_{i=1}^n u_{ij}^2}. \quad (30)$$

The above updating formulas show that t_j is the weighted average of x_i , and the weight of each item is inversely proportional to the distance, or proportional to the similarity.

In line 23 of our algorithm, we apply a multi-head attention mechanism to represent each t_j as an adaptive weighting sum of top-k x_i , where the weight of x_i with respect to t_j is expressed as a normalized dot product.

IV. EXPERIMENTS

A. Dataset and Evaluation Criteria

We compare the performance of our proposed method with the other methods on the benchmark PNG dataset, which matches noun phrase annotations in the Localized Narrative dataset [29] with panoptic segmentation annotations in the MS COCO dataset [30]. As the only publicly available benchmark in PNG, this dataset contains 726,445 noun phrases matched to segments involving 659,298 unique segments, and it covers 47.5% of the segmentation annotations in the MS COCO panoptic segmentation dataset and 45.1% of the noun phrases in the Localized Narrative dataset. On average, each title in the dataset has 5.1 noun phrases. The train and validation splits contains 133,103 and 8533 localized narratives, respectively.

We adopt the average recall as the evaluation metric for model performance following the previous practice. It calculates the recall for different intersection over (IoU) thresholds between the segmentation result and the ground truth, then draws curves based on different thresholds. The area under the curve represents the average recall value, i.e., for plural noun phrases. All ground truth annotations are merged into a single segment to compute the IoU result.

B. Implementation Details

Our backbone configuration is consistent with the PPMN baseline model [5], where we utilize official pre-trained [31] ResNet101 model (with 3x schedule) on the MS COCO dataset [30] as the image backbone. For the text input, we use the pre-trained ‘‘base-uncased’’ BERT model [32] to convert each word

in the narrative captions into a 768-dimensional vector. The longest caption contains 230 characters, with up to 30 different noun phrases that need be localized. We do not update the image and text pre-trained backbone models during training.

Furthermore, we only apply image size augmentation to the input image, which is resized to a resolution between 800 and 1,333 pixels while maintaining the aspect ratio. We implement our proposed model using PyTorch and train it with a batch size 10 for 20 epochs on three NVIDIA 3090 GPUs. The Adam optimizer is used with a fixed learning rate of 10^{-4} . During inference, we obtain segmentation results following the configuration of the two-stage model [1], which averages the matching graphs of all words in each noun phrase.

C. Experimental Results

To validate the effectiveness of the context information we introduced, we compare the performance of our proposed model with other methods on the PNG dataset. The main results are shown in Table II. We also compare the recall curves of these methods in Fig. 4. It is worth noting that our best model do not update the image and text backbones during training, and the results obtained using the same training strategy are labeled with PPMN†.

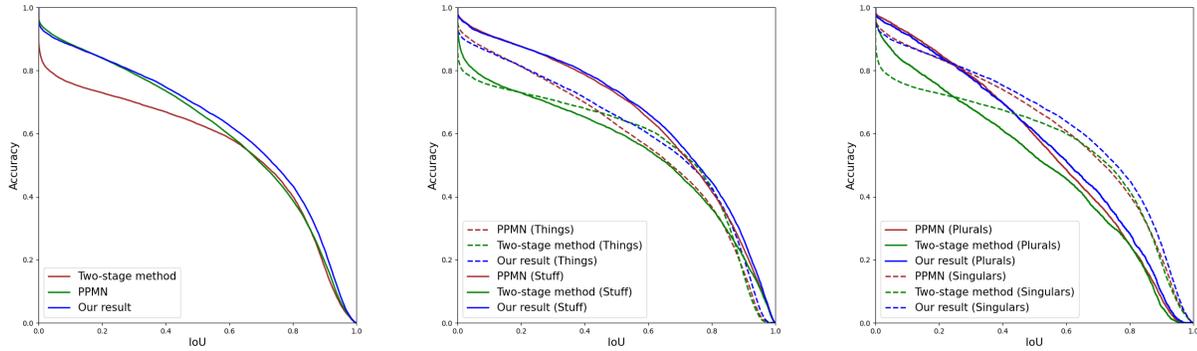
Compared to the current state-of-the-art methods on the PNG dataset, our proposed model achieves an average 3.5% (from 2.7% to 3.9%) improvements in average recall across various metrics. Specifically, our method achieves 3.5/3.1/3.9/3.6/2.7 improvements in whole/thing/thing/singular/plural categories, validating the effectiveness of our proposed method.

Fig. 4 depicts recall values for different classes at different IOU thresholds. In Fig. 4a, when the IoU threshold is larger than 0.3, our method (blue curve) consistently outperforms the baseline model (green curve), showing that the image context information can indeed benefit the segmentation results. Furthermore, in Fig. 4b, we can see that our method exhibits significant performance gains in object categories compared to the baseline model, even approaching the accuracy of the two-stage method. This further demonstrates that context information could enhance the representation ability of the aggregated text feature, leading to better segmentation results. In Fig. 4c, we further investigate the detailed performance of our method on object categories (stuff and singular): our method still improves the segmentation results on both categories, which indicates that the essential context information may benefit all object categories’ results.

D. Ablation Studies

To validate the effectiveness of our proposed method on different components, we conduct ablation experiments on the PNG task and compare the results under different parameter settings.

1) *Number of deformable encoder layer*: In Table III, we show how various deformable encoder layers affect the model’s performance. The results show that combining multi-scale image context information can improve segmentation



(a) Overall performance

(b) Things and stuff performance

(c) Singulars and plurals performance

Fig. 4: Average recall curves for the PNG dataset: (a) overall performance compared to other state-of-the-art methods, (b) curves for things and categories of things, and (c) curves for singular and plural noun phrases.

TABLE II: Results of our method for the panoptic narrative base task, compared with state-of-the-art methods.

Method model	Average Recall				
	<i>overall</i>	<i>singulars</i>	<i>plurals</i>	<i>things</i>	<i>stuff</i>
PNG [1]	55.4	56.2	48.8	56.2	54.3
PPMN† [5]	56.7	57.4	49.8	53.4	61.1
EPNG [6]	58.0	58.6	52.1	54.8	62.4
PPMN [5]	59.4	60.0	54.0	57.2	62.5
DRMN(Our)	62.9 (+3.5)	63.6 (+3.6)	56.7 (+2.7)	60.3 (+3.1)	66.4 (+3.9)

performance, whereas too many encoder layers may lead to performance degradation.

TABLE III: Ablation study on the number of encoder layers.

Num of Encoder Layers	Average Recall				
	<i>overall</i>	<i>singulars</i>	<i>plurals</i>	<i>things</i>	<i>stuff</i>
0	61.5	62.2	55.4	58.7	65.4
1	62.4	62.9	57.1	60.1	65.6
2	62.9	63.6	56.7	60.3	66.4
3	62.6	63.3	55.8	60.2	65.9

2) Number of rounds for multi-round feature aggregation:

We also examine the performance of different rounds of feature aggregation on the model in Table IV where the results show that introducing multi-round feature aggregation suddenly improves model performance. We also observe that the performance of the Singular and Stuff categories gradually improves as the number of stages increases. However, there are some fluctuations in the performance of the plurals category. Since we find some incomplete annotations in the PNG dataset for the plurals category during training (as shown in Fig. 5), we believe it is reasonable for the model to be slightly unstable in testing on this category.

3) Number of sample points for multi-round feature aggregation module: We conduct further studies to evaluate our proposed multi-round feature aggregation module by examining the impact of different numbers of sampling points on model performance. The results are reported in Table V. Since the context information covers a large extent on image during

TABLE IV: Ablation study on the number of feature aggregation round.

Number of Rounds	Average Recall				
	<i>overall</i>	<i>singulars</i>	<i>plurals</i>	<i>things</i>	<i>stuff</i>
0	59.6	60.2	54.5	57.1	63.0
1	62.7	62.7	57.2	60.4	65.9
2	62.7	63.5	56.2	60.1	66.4
3	62.9	63.6	56.7	60.3	66.4

top-k image refinement stage, as further shown in Fig. 8, our model can perform well even with a small number of sample points, This indicates a small set for the context information refined top-k image features may be enough to cover the object information. We observe that increasing the number of sampling points improves the stuff category more. We attribute this to the mask of stuff category which usually covers more space in the image. Hence increasing the sampling points may help to preserve the semantics of different parts of the ground truth mask information.

TABLE V: Ablation on the number of sampled image points.

Sampling Points (<i>k</i>)	Average Recall				
	<i>overall</i>	<i>singulars</i>	<i>plurals</i>	<i>things</i>	<i>stuff</i>
10	62.8	63.4	56.9	60.4	66.1
50	62.9	63.6	56.7	60.4	66.4
100	62.9	63.6	56.7	60.3	66.4
400	62.7	63.3	56.5	60.2	66.2

E. Qualitative Analysis

We illustrate the qualitative results of our proposed model for text paragraphs in Fig. 5. It is observed that comparing to the baseline model, our model predicts more complete segmentation results (“doors”, “windows” and “refrigerator” in the first row, “person”, “group of people” in the second row), indicating that refinement of top- k sampled image features benefits to cluster more related pixels. It is worth mentioning that the model is even able to locate more complete ground truth annotation (the “few bowls” result in the first row, the “two persons” result in the third row).

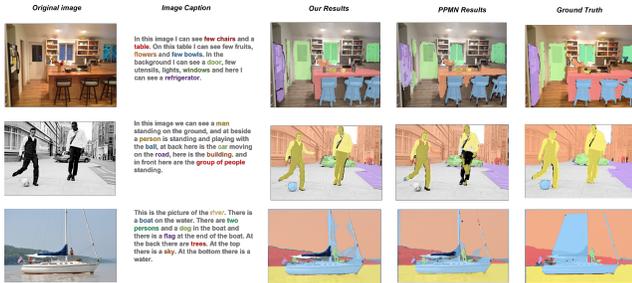


Fig. 5: Qualitative results for Panoptic Narrative Grounding. The segmentation masks in the image correspond one-to-one to the colors mentioned in the text.

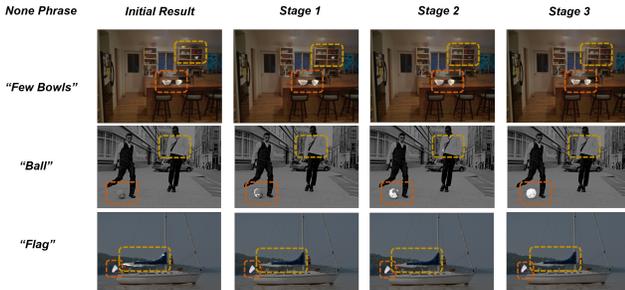


Fig. 6: Refinement results in each stage for specific visual objects. Distracting objects and target objects are highlighted in dashed yellow and red boxes, respectively.

To visualize how the context information benefits segmentation, we show the instructive results during the multi-round feature aggregation process in Fig. 6. In the first and third rows, the model gradually filters out irrelevant results in each round of refinement, improving segmentation results despite irrelevant or similar objects in the initial matching. The example in the second row shows how a segmented object goes from a low response matching result to an almost correct matching during refinement. Compared to the example in Fig. 2, such results validate that context information does matter in alleviate the phrase-to-pixel mis-match and thus improve the performance in PNG.

We further visualize the top- k image locations most similar to the text during each round of refinement and corresponding weights in the cross-attention mechanism in Fig. 7. As seen in

the figure, our proposed method generally puts the weight on the most relevant objects and gradually filters out the impact of irrelevant objects.

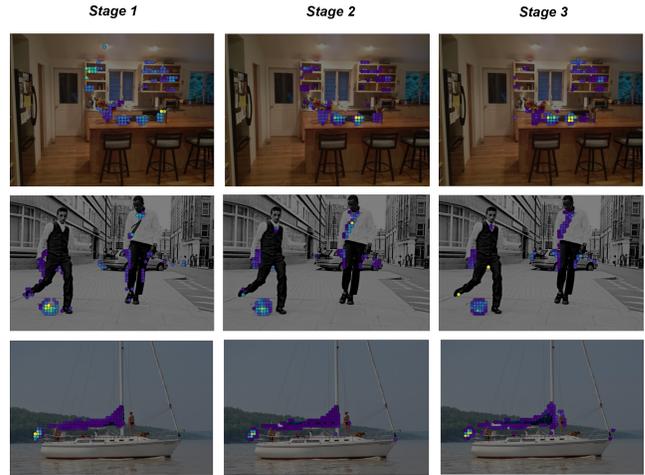


Fig. 7: Attention weights for top- k image locations. Weights are averaged over all heads in a multi-head cross-attention layer. Lighter colors indicate greater weight and vice versa.

To better visualize what context information is introduced by deformable attention during the aggregation stage, we further visualize the 50 most important offset points obtained by sampling the top- k most relevant points in the last round of the deformable attention mechanism for different layers in Fig. 8. We can find that in the primary layers containing relatively detailed low-level information, the offset point is around the target object, which may refine more detailed information for the target object. In the latter layers that contain relatively high-level information, the offset points attend to the general context around the target object. The above observations indicate our proposed refinement method effectively concentrates on both detailed and general context information for the top- k sampled image points.



Fig. 8: Visualization of the most important offset points in the deformable attention layer, according to the top- k query of phrase “few bowl”. We visualize the top-50 most important points based on the attention weights for each layer of the multi-scale feature map.

V. CONCLUSION

In this paper, we propose a novel one-stage model named Deformable-Attention Refined Matching Network (DRMN) for Panoptic Narrative Grounding (PNG) task. Built upon the end-to-end one-stage model architecture, we integrate the essential context information of multi-scale image features

in the multi-modal information fusion module as an addition cue to enhance the feature discriminative ability. Furthermore, we employ a clustering framework to interpret our proposed module and validate our method through experiments on the benchmark PNG dataset. The results demonstrate that our proposed model can achieve new state-of-the-art performance, with a 3.5% improvement on the average recall metric.

REFERENCES

- [1] C. González, N. Ayobi, I. Hernández, J. Hernández, J. Pont-Tuset, and P. Arbeláez, “Panoptic narrative grounding,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 1364–1373.
- [2] S. Pratt, M. Yatskar, L. Weihs, A. Farhadi, and A. Kembhavi, “Grounded situation recognition,” in *Proceedings of the European Conference on Computer Vision (ECCV)*. Springer, 2020, pp. 314–332.
- [3] D. A. Hudson and C. D. Manning, “Gqa: A new dataset for real-world visual reasoning and compositional question answering,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6700–6709.
- [4] L. H. Li, P. Zhang, H. Zhang, J. Yang, C. Li, Y. Zhong, L. Wang, L. Yuan, L. Zhang, J.-N. Hwang *et al.*, “Grounded language-image pre-training,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 10965–10975.
- [5] Z. Ding, Z.-h. Ding, T. Hui, J. Huang, X. Wei, X. Wei, and S. Liu, “Ppmm: Pixel-phrase matching network for one-stage panoptic narrative grounding,” in *Proceedings of the 30th ACM International Conference on Multimedia*, 2022, pp. 5537–5546.
- [6] H. Wang, J. Ji, Y. Zhou, Y. Wu, and X. Sun, “Towards real-time panoptic narrative grounding by an end-to-end grounding network,” *arXiv preprint arXiv:2301.03160*, 2023.
- [7] X. Zhu, W. Su, L. Lu, B. Li, X. Wang, and J. Dai, “Deformable detr: Deformable transformers for end-to-end object detection,” *arXiv preprint arXiv:2010.04159*, 2020.
- [8] A. Kamath, M. Singh, Y. LeCun, G. Synnaeve, I. Misra, and N. Carion, “Mdetr-modulated detection for end-to-end multi-modal understanding,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 1780–1790.
- [9] F. Shi, R. Gao, W. Huang, and L. Wang, “Dynamic mdetr: A dynamic multimodal transformer decoder for visual grounding,” *arXiv preprint arXiv:2209.13959*, 2022.
- [10] H. De Vries, F. Strub, S. Chandar, O. Pietquin, H. Larochelle, and A. Courville, “Guesswhat?! visual object discovery through multi-modal dialogue,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2017, pp. 5503–5512.
- [11] B. A. Plummer, L. Wang, C. M. Cervantes, J. C. Caicedo, J. Hockenmaier, and S. Lazebnik, “Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2015, pp. 2641–2649.
- [12] R. Hu, M. Rohrbach, and T. Darrell, “Segmentation from natural language expressions,” in *Proceedings of the European Conference on Computer Vision (ECCV)*. Springer, 2016, pp. 108–124.
- [13] H. Zhang, Y. Niu, and S.-F. Chang, “Grounding referring expressions in images by variational context,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4158–4166.
- [14] R. Girshick, “Fast r-cnn,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2015, pp. 1440–1448.
- [15] Z. Yang, B. Gong, L. Wang, W. Huang, D. Yu, and J. Luo, “A fast and accurate one-stage approach to visual grounding,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 4683–4693.
- [16] Y.-C. Chen, L. Li, L. Yu, A. El Kholi, F. Ahmed, Z. Gan, Y. Cheng, and J. Liu, “Uniter: Universal image-text representation learning,” in *Proceedings of the European Conference on Computer Vision (ECCV)*. Springer, 2020, pp. 104–120.
- [17] S. Kazemzadeh, V. Ordonez, M. Matten, and T. Berg, “Referitgame: Referring to objects in photographs of natural scenes,” in *In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*, 2014, pp. 787–798.
- [18] B. A. Plummer, A. Mallya, C. M. Cervantes, J. Hockenmaier, and S. Lazebnik, “Phrase localization and visual relationship detection with comprehensive image-language cues,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2017, pp. 1928–1937.
- [19] P. Dogan, L. Sigal, and M. Gross, “Neural sequential phrase grounding (seqground),” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4175–4184.
- [20] Z. Mu, S. Tang, J. Tan, Q. Yu, and Y. Zhuang, “Disentangled motif-aware graph learning for phrase grounding,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 15, 2021, pp. 13587–13594.
- [21] T. Yu, T. Hui, Z. Yu, Y. Liao, S. Yu, F. Zhang, and S. Liu, “Cross-modal omni interaction modeling for phrase grounding,” in *Proceedings of the 28th ACM international conference on multimedia*, 2020, pp. 1725–1734.
- [22] S. Liu, Z. Zeng, T. Ren, F. Li, H. Zhang, J. Yang, C. Li, J. Yang, H. Su, J. Zhu *et al.*, “Grounding dino: Marrying dino with grounded pre-training for open-set object detection,” *arXiv preprint arXiv:2303.05499*, 2023.
- [23] H. Zhang, P. Zhang, X. Hu, Y.-C. Chen, L. Li, X. Dai, L. Wang, L. Yuan, J.-N. Hwang, and J. Gao, “Glipv2: Unifying localization and vision-language understanding,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 36067–36080, 2022.
- [24] Q. Wang, H. Tan, S. Shen, M. W. Mahoney, and Z. Yao, “Maf: Multi-modal alignment framework for weakly-supervised phrase grounding,” *arXiv preprint arXiv:2010.05379*, 2020.
- [25] J. Deng, Z. Yang, T. Chen, W. Zhou, and H. Li, “Transvg: End-to-end visual grounding with transformers,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 1769–1779.
- [26] H. Ding, C. Liu, S. Wang, and X. Jiang, “Vision-language transformer and query generation for referring segmentation,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 16321–16330.
- [27] Z. Yang, J. Wang, Y. Tang, K. Chen, H. Zhao, and P. H. Torr, “Lavt: Language-aware vision transformer for referring image segmentation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 18155–18165.
- [28] K. Jain and V. Gandhi, “Comprehensive multi-modal interactions for referring image segmentation,” *arXiv preprint arXiv:2104.10412*, 2021.
- [29] J. Pont-Tuset, J. Uijlings, S. Changpinyo, R. Soricut, and V. Ferrari, “Connecting vision and language with localized narratives,” in *Proceedings of the European Conference on Computer Vision (ECCV)*. Springer, 2020, pp. 647–664.
- [30] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft coco: Common objects in context,” in *Proceedings of the European Conference on Computer Vision (ECCV)*. Springer, 2014, pp. 740–755.
- [31] Y. Wu, A. Kirillov, F. Massa, W.-Y. Lo, and R. Girshick, “Detectron2,” <https://github.com/facebookresearch/detectron2>, 2019.
- [32] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” *arXiv preprint arXiv:1810.04805*, 2018.