# Rethinking the Reference-based Distinctive Image Captioning

Yangjun Mao*
Zhejiang University
maoyj0119@zju.edu.cn

Long Chen*
Columbia University
zjuchenlong@gmail.com

Zhihong Jiang
Zhejiang University
zju_jiangzhihong@zju.edu.cn

Dong Zhang
Hong Kong University of Science and
Technology, dongz@ust.hk

Zhimeng Zhang, Jian Shao
Zhejiang University
zhimeng@zju.edu.cn,jshao@zju.edu.cn

Jun Xiao†
Zhejiang University
junx@zju.edu.cn

## ABSTRACT

Distinctive Image Captioning (DIC) — generating distinctive captions that describe the unique details of a target image — has received considerable attention over the last few years. A recent DIC work proposes to generate distinctive captions by comparing the target image with a set of semantic-similar reference images, *i.e.*, reference-based DIC (Ref-DIC). It aims to make the generated captions can tell apart the target and reference images. Unfortunately, reference images used by existing Ref-DIC works are easy to distinguish: *these reference images only resemble the target image at scene-level and have few common objects, such that a Ref-DIC model can trivially generate distinctive captions even without considering the reference images.* For example, if the target image contains objects "towel" and "toilet" while all reference images are without them, then a simple caption "A bathroom with a towel and a toilet" is distinctive enough to tell apart target and reference images. To ensure Ref-DIC models really perceive the unique objects (or attributes) in target images, we first propose two new Ref-DIC benchmarks. Specifically, we design a two-stage matching mechanism, which strictly controls the similarity between the target and reference images at object-/attribute- level (vs. scene-level). Secondly, to generate distinctive captions, we develop a strong Transformer-based Ref-DIC baseline, dubbed as **TransDIC**. It not only extracts visual features from the target image, but also encodes the differences between objects in the target and reference images. Finally, for more trustworthy benchmarking, we propose a new evaluation metric named **DisCIDEr** for Ref-DIC, which evaluates both the accuracy and distinctiveness of the generated captions. Experimental results demonstrate that our TransDIC can generate distinctive captions. Besides, it outperforms several state-of-the-art models on the two new benchmarks over different metrics.

## CCS CONCEPTS

• **Computing methodologies → Computer vision**.

## KEYWORDS

Image Captioning, Distinctiveness, Benchmark, Transformer

## 1 INTRODUCTION

Image captioning, *i.e.*, generating natural language descriptions to summarize the salient contents of a target image, has drawn much attention from the multimedia community. It has great impacts on many downstream applications, such as helping blind people and developing navigation systems. However, as revealed in [11, 12], conventional image captioning models tend to generate over-generic captions or even identical captions when input images are similar. Obviously, these generic captions neglect the unique details of the target image. Recent captioning works [12, 25, 26, 42] begin to make generated captions more distinctive and ask these captions describe more unique details of each target image, called Distinctive Image Captioning (DIC).

Currently, mainstream DIC works follow the same setting as plain image captioning: using one single image as input, and generating distinctive captions for each image, dubbed as Single-image DIC (**Single-DIC**). In this setting, they tend to generate a totally distinctive caption. By "totally", we mean the generated caption is asked to distinguish its corresponding image from all images in the dataset, *i.e.*, dataset-level distinctiveness. To this end, they always resort to reinforcement learning and develop different distinctive rewards [25, 26]. However, this Single-DIC setting has two inherent issues: 1) It is difficult (or impossible) to generate a totally distinctive caption for the target image unless we describe all the details in the image. 2) Even for our humans, we still need some reference images when generating distinctive captions. For example in Figure 1 (b), without any reference images, people won't know what should be emphasized in the *target image*, and may simply predict "A bathroom with a towel" for the image. In contrast, they will focus on the unique colors of "towel" and "shower curtain", and predict "A bathroom with a pink towel and a blue shower curtain" when they use the "white shower curtain" in *ref-img1* and "yellow towel" in *ref-img2* as references.

For human-like distinctive captioning, a recent work proposes to study the DIC task based on a group of semantic-similar reference images, dubbed Reference-based DIC (**Ref-DIC**). Different

**Figure 1: (a): An example of constructed reference image group used in existing Ref-DIC work [43]. (b): Selected reference images for the same target image using our two-stage matching mechanism. We use same colors to denote the same object categories in the different images (*e.g.*, "`towel`" is with yellow box, "`shower curtain`" is with red box, and "`toilet`" is with green box). We only show two reference images here.**

from Single-DIC, they use the target image and all reference images as input and these reference images will inform DIC models which parts of the target image should be emphasized. Compared to Single-DIC, the generated captions are only asked to distinguish the target image from the group of reference images, *i.e.*, group-level distinctiveness. Unfortunately, the reference images used in existing Ref-DIC works [43] can be trivially distinguished: *these reference images only resemble the target image at the scene-level and have few common objects, thus Ref-DIC models can simply generate distinctive captions even without considering the reference images.* For example in Figure 1 (a), *target* and *reference images* have no object in common (*e.g.*, "`towel`", "`shower curtain`", or "`toilet`"), each object in *target image* is unique, such that the Ref-DIC model can trivially generate "`a bathroom with a towel`" to tell the target and reference images apart.

As mentioned above, reference images are crucial when defining the unique details in the target image. In this paper, we propose two new benchmarks for the Ref-DIC task: **COCO-DIC** and **Flickr30K-DIC**. To strictly control the unique details between target and reference images, we propose a two-stage matching mechanism, which can measure image similarity at the object-/attribute- level (vs. scene-level in [43]), and deliberately make target and reference images have some common objects. Under this mechanism, Ref-DIC models can learn to focus on the unique attributes and objects in target image. As the example in Figure 1 (b), compared to *ref-img1*, *target image* has the unique attribute "`blue`" of "`shower curtain`" and the unique object "`towel`".

To achieve group-level distinctiveness, we propose to emphasize both unique attributes and objects in target image. Thus, we also propose a new Transformer-based captioning model named **TransDIC**, which directly gives each region in target image (target

regions) some region references when generating captions. Specifically, we firstly find similar regions from reference images (reference regions) for each target region. Then, we send target region and its corresponding reference regions into the *Two-Flow Encoder* Module, which consists of a Target flow and a Target-Reference flow. The Target flow aims to encode target image features through self-attention blocks [38]. And the Target-Reference flow enables cross-image interactions between target and reference images through a multi-layer co-attention. Different from the existing Ref-DIC work [43] which proposes an attention module to focus on unique objects in target image, our TransDIC directly enables the feature interactions between target and reference images.

Furthermore, to fully take advantage of ground-truth captions of reference images, we propose a new CIDEr-based [40] metric termed as **DisCIDEr**. According to our definition of group-level distinctiveness, we believe frequently-used n-grams in ground-truth captions of reference images should be given less weight at evaluation time. The metric can not only directly evaluate the distinctiveness, but also preserve the accuracy advantage of CIDEr. Extensive experimental results on multiple Ref-DIC benchmarks (COCO-DIC and Flickr30K-DIC) have demonstrated the effectiveness of our proposed TransDIC.

In summary, we make three contributions in this paper:

(1) We proposed two new benchmarks for Ref-DIC, constructed groups in those benchmarks can inform the model which parts in target image should be emphasized.

(2) We developed a new model TransDIC for Ref-DIC, which achieves better performance than the state-of-the-art Ref-DIC models in terms of both accuracy and distinctiveness.

(3) We proposed a new metric named DisCIDEr which considers both the distinctiveness and accuracy of the generated caption.

## 2 RELATED WORK

**Image Captioning.** Most modern image captioning models typically employ an encoder-decoder framework for caption generation [15, 16, 22, 41, 44]. Within this framework, many efforts have been made to improve the architecture, including attention mechanisms [1, 5, 7, 14, 46], graph convolution networks [8, 50, 51], and transformer-based models [10, 17]. Meanwhile, another series of works explore different training objectives at the training stage. For example, Dai *et al.* [11] and Shetty *et al.* [35] leverage Generative Adversarial Network (GAN) to the improve the diversity of generated captions. Some recent captioning works apply Reinforcement Learning (RL) to captioning and achieve great success [24, 32, 34, 48]. These models directly optimize non-differentiable evaluation metrics (*e.g.*, BLEU [27], CIDEr [40]), which boost the caption generation procedure at the sentence-level.

**Distinctive Image Captioning (DIC).** Compared with conventional image captioning, DIC is a more challenging task, which tends to generate more informative and descriptive captions. According to the stage they take effect, existing solutions can be coarsely divided into two categories: *Inference-based* and *Training-based* methods. Inference-based models mainly modify the caption decoding procedure at inference time and thus can be applied to any captioning architecture [39, 45]. In contrast, training-based methods, resort to different training objects [25, 26] or the progressive
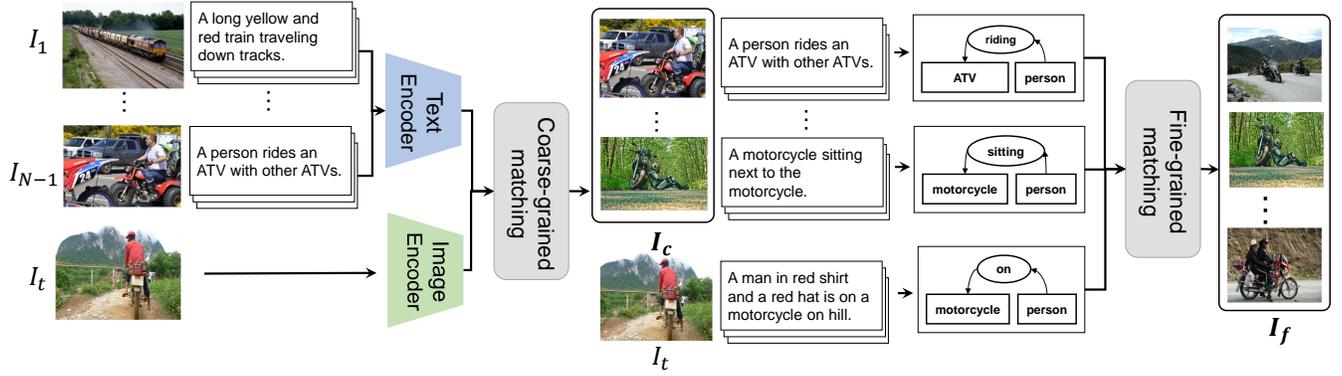
**Figure 2: The pipeline of our two-stage matching procedure. In the first stage, we calculate image-text similarity scores between target image $I_t$ and the captions of all other images in $\mathcal{D}$ through CLIP [31] and construct a coarse-grained group $\mathcal{I}_c$ for target image $I_t$. In the second stage, we leverage scene graphs to calculate the object and attribute overlaps between images in $\mathcal{I}_c$ and $I_t$. We rearrange $\mathcal{I}_c$ according to their similarity scores with $I_t$, and finally get the fine-grained reference image group $\mathcal{I}_f$.**

training procedure [23]. Recently, some works begin to study the DIC task based on semantic-similar reference images. Specifically, Wang *et al.* [42] propose to assign higher weights to distinctive ground-truth captions at the training stage, Wang *et al.* [43] use multiple images as input to emphasize distinctive objects. In this paper, we propose a co-attention based model to directly enable the feature interactions between target and reference images.

To evaluate the distinctiveness of generated captions in DIC, several new evaluation metrics are developed. SPICE-U [45] is designed for Single-DIC. CIDErBtw [42] measures the distinctiveness of caption at sentence-level similarity. DisWordRate [43] directly evaluates the occurrences of distinctive words. In this paper, we develop a new metric named DisCIDEr for Ref-DIC. Compared to existing metrics, our metric fully explore the distinctiveness of each individual n-gram in ground-truth captions of target image.
**Multi-input Image Captioning.** Several captioning settings need multiple images as input. According to the number of input images, they can be divided into two categories: *Two-image based* and *Group-based* captioning. Two-image based captioning tends to describe the common [36] or different [28, 30, 37, 49] parts between the two images. Thus, the two images in their settings always have strong correlations. For example, the change captioning task takes before and after images as input and describes the changes between them. In contrast, group-based captioning typically uses a group of images as references to investigate certain properties of the target image. For example, Chen *et al.* [4] firstly model the relevance and diversity between target and reference images and aim to generate diverse captions for the target image. Li *et al.* [19] tend to describe a group of target images using another group of semantically similar images as references.

## 3 PROPOSED BENCHMARKS

In this section, we firstly formally define the Ref-DIC task. Then we describe our solution for Ref-DIC benchmarks construction. Finally, we provide details of our proposed COCO-DIC and Flickr30K-DIC benchmarks for Ref-DIC.

### 3.1 Task Definition: Reference-based DIC

Given a **target image** $I_t$ and a group of $K$ **reference images** $\mathcal{I}_r = \{I_i\}_{i=1}^K$ which are semantic-similar to $I_t$, Ref-DIC models aim to generate a natural language sentence $S = \{w_1, w_2, \ldots, w_T\}$. The generated sentence $S$ should not only correctly describe the target image $I_t$, but also contain sufficient details about $I_t$, so it can tell apart target and reference images. For example in Figure 1 (b), given the target and reference images, Ref-DIC models aim to generate a distinctive caption "a bathroom with a pink towel, a blue shower curtain and a toilet". The detail "pink towel" is helpful to distinguish target image from *ref-img2* because the "towel" in *ref-img2* is "white". On the contrary, predicting "a bathroom with a towel and a toilet" fails to meet the requirements because it is suitable for both target and reference images.

### 3.2 Ref-DIC Benchmarks Construction

Given a conventional image captioning dataset $\mathcal{D}$, suppose it contains $N$ images and each one has $M$ corresponding ground-truth captions. We build new Ref-DIC benchmarks based on $D$, by coupling each image (target image) with several semantic-similar reference images. Specifically, each image in $\mathcal{D}$ will be regarded as a target image $I_t$, and all remaining $N-1$ images are termed as its **candidate reference images**. For each target image $I_t$, our goal is to retrieve $K$ reference images from its candidate reference images to construct the reference image group $\mathcal{I}_r$.

To achieve group-level distinctiveness, $I_t$ and retrieved $\mathcal{I}_r$ should have some common objects, such that $\mathcal{I}_r$ will inform the model to focus on the unique details in $I_t$. To this end, we design a two-stage matching mechanism. In the first stage, we construct a **coarse-grained group** $\mathcal{I}_c$ based on the image-text similarity score for each target image. Then in the second stage, we investigate fine-grained details of $I_t$ and $\mathcal{I}_c$, and construct a **fine-grained group** $\mathcal{I}_f$ based on $\mathcal{I}_c$. Finally, we select $K$ images out of $\mathcal{I}_f$ to construct the $\mathcal{I}_r$. We detailed introduce our two-stage matching pipeline below.

*3.2.1 Coarse-grained Group Construction.* Following [43], we use an image-text retrieval model to calculate similarity scores between
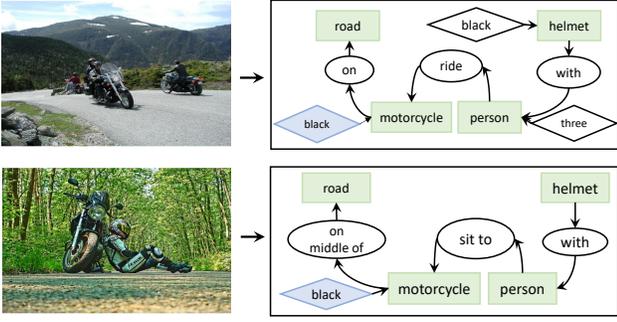
**Figure 3: An example of parsed scene graphs for the ground-truth captions of two images. Two graphs have four object overlaps: "helmet", "people", "motorcycle", and "road" (green), and one attribute overlap "black" (blue).**

**Table 1: Statistical summary of the COCO-DIC, Flickr30K-DIC, and existing Ref-DIC benchmark [43]. "#overlaps" denotes the number of object/attribute overlap in each dataset.**

| Datasets | images | | | #overlaps in a group | | |
|---|---|---|---|---|---|---|
| | Train | Val | Test | Train | Val | Test |
| Wang *et al.* [43] | 133,980 | 5,562 | 5,538 | 3.8 | 3.7 | 3.7 |
| COCO-DIC | 123,287 | 5,000 | 5,000 | 5.0 | 4.9 | 4.9 |
| Flickr30K-DIC | 29,000 | 1,014 | 1,000 | 5.3 | 5.3 | 5.3 |

images and texts. Specifically, as shown in Figure 2 (left), we use a pre-trained CLIP [31] to firstly extract the visual feature of target image $I_t$ and text features of ground-truth captions from candidate reference images. Then, we perform the cosine similarity between visual feature and all text features to get $(N-1) \times M$ scores. Finally, we select $||\mathcal{I}_c||$ captions with the highest scores, and their corresponding images are used as the coarse group $\mathcal{I}_c$ for $I_t$.

The CLIP-based matching mechanism can effectively filter out some obviously unrelated candidate images. However, since it encodes texts at the sentence level, several fine-grained details may be neglected when computing similarity. For example in Figure 2, ground-truth captions of candidate image $I_1$ contain the "train", thus $I_1$ will be removed due to a low scene-level similarity score to $I_t$. Meanwhile, $I_{N-1}$ is considered similar to $I_t$ by the CLIP because both of them describe the scene of "someone is riding a vehicle". However, they resemble each other only at the sense-level and do not contain any common objects (*e.g.*, "motorcycle" in $I_t$ and "ATV" in $I_{N-1}$ are totally different objects).

*3.2.2 Fine-grained Group Construction.* To overcome the shortcoming of the coarse-grained group, we propose a fine-grained matching mechanism that directly uses object and attribute overlaps between two images as the similarity measurement. Firstly, for $I_t$ or any image in $\mathcal{I}_c$, we parse all its ground-truth captions into one scene graph. Then, we extract objects and attributes from scene graphs [6, 18, 21, 47] of two images to calculate overlaps. Specifically, objects from the two graphs will be compared according to their categories. However, two attributes should firstly correspond to the same objects and then compare to each other. For example in

Figure 3, when calculating object overlaps, "black helmet" (top) and "helmet" (bottom) from two graphs denote one-time object overlap (*e.g.*, common object "helmet"). As for attribute overlaps, two graphs have both "black motorcycle" in common and denote one-time attribute overlap[1]. In this paper, we take the sum over object and attribute overlaps as the final similarity score for two images. And we sort all images in $\mathcal{I}_c$ according to their similarity scores to $I_t$ to construct the fine-grained group $\mathcal{I}_f$. It is worth noting that we select $K$ images but not the top-K from $\mathcal{I}_f$ to construct $I_r$. The reason for this choice is that we believe the most similar images from $\mathcal{I}_f$ may contain the identical objects and attributes as $I_t$, thus they won't help to emphasize any unique details in $I_t$. More detailed discussion about the top-K selection are left in Table 5.

### 3.3 Benchmarks: COCO-DIC & Flickr30K-DIC

We apply our matching mechanism to widely-used captioning benchmarks MS-COCO [9] and Flickr30K [29] to construct **COCO-DIC** and **Flickr30K-DIC**, respectively. Some basic statistics about our proposed benchmarks are reported in Table 1. Different from the construction procedure proposed in [43], they avoid image reuse (or overlap) among different constructed groups. Thus, some images which are not similar enough may be forced to construct a group. In contrast, we find $K$ reference images independently to ensure the similarity within a group.

## 4 PROPOSED APPROACH

### 4.1 Preliminaries

*4.1.1 Transformer-based Image Captioning.* Transformer [38] follows the standard **Encoder-Decoder** architecture. It employs the self-attention mechanism to explore the internal correlation within the sequential data, which has been widely adopted by numerous image captioning models [10, 17]. For a given image $I$, they use proposal features [1] extracted by an object detector as input: $X = \{x_i\}_{i=1}^N$ where $x_i \in \mathbb{R}^d$ is the feature vector for $i$-th proposal in $I$ and $N$ is the number of proposals. They employ multiple self-attention layers as **Encoder**, and the outputs of the $l$-th layer is calculated as follows:

$$H_{l-1} = \mathbf{LN}\left(O_{l-1} + \mathbf{MH}(O_{l-1}, O_{l-1}, O_{l-1})\right), \qquad (1)$$

$$O_l = \mathbf{LN}\left(H_{l-1} + \mathbf{FFN}(H_{l-1})\right), \qquad (2)$$

where $O_0$ refers to input proposal features $X$, and $O_{l-1}$ is outputs of the $(l-1)$-th layer. $\mathbf{LN}(\cdot)$ denotes the layer normalization [2], $\mathbf{FFN}(\cdot)$ denotes the feed forward network, and $\mathbf{MH}(\cdot)$ denotes the multi-head attention [38]. Encoded visual features are fed into the **Decoder** for caption generation. It generates a word probability distribution $P_t = P(w_t|w_{1:t-1}, I)$ at each time step t conditioning on the previously generated words $\{w_1, \ldots, w_{t-1}\}$ and image $I$.

*4.1.2 Model Optimization.* Mainstream captioning works typically resort to a two-stage procedure for model optimization [10, 34]. Given an image $I$ and its ground-truth captions $C = \{c_i\}_{i=1}^M$. They firstly apply a cross-entropy loss (XE) to pre-train the model and then employ reinforcement learning (RL) to finetune sequence generation.

---

[1]Note that the "black helmet" and "black motorcycle" contain no attribute overlap because the attribute "black" belongs to different objects.
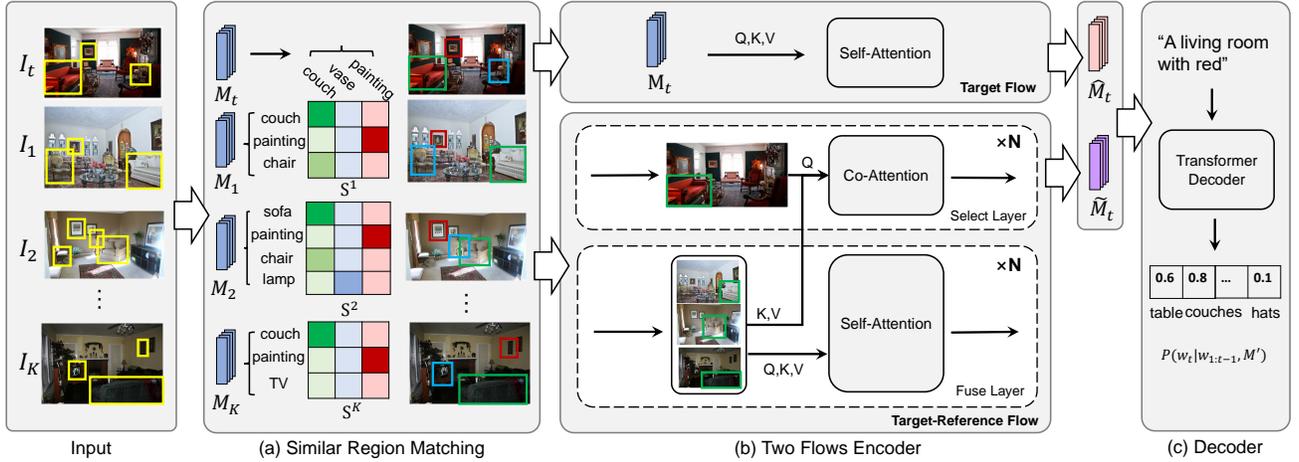
**Figure 4: Overview of our proposed TransDIC model. It consists of three parts: (a) A Similar Region Matching module that couples each target region with reference regions. (b) A Two-Flow Encoder module that encodes both target and reference images. (c) A plain captioning decoder. In module (a), we use the same colors to denote the same objects. We firstly send target and reference images into the similar region matching module to construct Target-Reference tuples, *e.g.*, all regions marked with green boxes contain the object "couch", thus form a Target-Reference tuple. Then we send constructed tuples into the Two-Flow Encoder for target image and cross-image features extraction. Finally, both kinds of features are sent into the decoder for caption generation.**

When training with XE, for a ground-truth caption $c_i = \{w_t\}_{t=1}^T$, they ask the model to minimize the following cross-entropy loss:

$$L_{xe} = -\sum_{t=1}^T \log P(w_t|w_{1:t-1}, I), \tag{3}$$

where $P(w_t|w_{1:t-1}, I)$ denotes the predicted probability of word $w_t$.

When training with reinforcement learning, they firstly generate top-n captions $\hat{C} = \{\hat{c}_i\}_{i=1}^n$ through beam search, and then optimize the following RL loss [10][2]:

$$L_{rl} = -\frac{1}{n}\sum_{i=1}^n ((r(\hat{c}_i, C) - b) \log p(\hat{c}_i)), \tag{4}$$

where $r(\cdot, \cdot)$ is the reward function computed between $\hat{c}_i$ and $C$, and $b = (\sum_{i=1}^n r(\hat{c}_i, C))/n$ is the baseline, calculated as the mean of the rewards obtained by the generated captions.

## 4.2 TransDIC: Transformer-based Ref-DIC

Given a target image $I_t$, we term all regions (or proposals) in it as **target regions** $R^t = \{r_n^t\}_{n=1}^N$. Our proposed model tends to give each target region $r_n^t$ some region references when generating captions. To this end, we firstly couple each $r_n^t$ with semantic-similar regions from $\mathcal{I}_r$ as the reference, *i.e.*, **reference regions**. Then we send each target region and its reference regions into the model for distinctive caption generation.

Specifically, our TransDIC consists of three components: 1) Similar Region Matching module in Section 4.2.1. 2) Two-Flow Encoder module. The module contains two parallel data flows to extract both target image and cross-image features in Section 4.2.2. 3)

A plain Transformer-based Captioning Decoder for caption generation. An overview of our model is shown in Figure 4.

*4.2.1 Similar Region Matching.* For each target region, we retrieve regions from reference images with the highest similarity scores to it as its reference regions. Given a target image $I_t$ and its corresponding $\mathcal{I}_r$, their region features are firstly projected into the memory space through an MLP layer. We denote memory features for $I_t$ and the K images in $\mathcal{I}_r$ as $M_t = \{m_j^t\}_{j=1}^N$ and $M_k = \{m_i^k\}_{i=1}^N, k = \{1 \ldots K\}$, respectively.

Then, we calculate the cosine similarity scores between features in $M_t$ and $M_k$, *i.e.*,

$$S_{ij}^k = cos(m_i^k, m_j^t), \tag{5}$$

where $m_j^t$ represents the $j$-th region in $I_t$, and $m_i^k$ represents the $i$-th region in reference image $I_k$. We apply max operation to get the most similar region for $r_j^t$ according to the calculated $S^k$:

$$\hat{r}_j^k = \arg\max_i(\{S_{ij}^k\}_{i=1}^N), \tag{6}$$

where $\hat{r}_j^k$ denotes the most similar region from image $I_k$ for $r_j^t$, *i.e.*, reference region. As an example in Figure 4 (a), reference image $I_K$ contains regions "couch", "painting" and "TV". For the region "couch" in $I_t$, we can learn from the similarity matrix $S^K$: the region is similar to the "couch" in $I_K$ (deep green) while is different from the "TV" or "painting" (light green) in $I_K$. Max operation is then token along each column of $S^K$, and "couch" in $I_k$ is selected as the reference region for "couch" in $I_t$.

Finally, for each target region $r_n^t$, we gather K reference regions, one for each, from K reference images. We put these K + 1 regions together and term them as a **Target-Reference region tuple**:

$$T_n = \{r_n^t, \hat{r}_n^1, \hat{r}_n^2, \ldots, \hat{r}_n^K\}. \quad n = \{1, \ldots, N\} \tag{7}$$

For example in Figure 4 (a), all regions marked with green boxes form a Target-Reference tuple (for "couch").

*4.2.2 Two-Flow Encoder.* Our proposed module takes $M_t$ and $M_k$ as input, and extracts target image features and cross-image features through the **Target flow** and the **Target-Reference flow**, respectively. An overview of the module is shown in 4 (b).

**Target flow.** The flow enables the region interactions within the target image $I_t$. Same as the standard transformer-based captioning model, it sends memory features $M_t$ into multiple self-attention layers, and finally outputs encoded features $\hat{M}_t = \{\hat{m}_i^t\}_{i=1}^N$ for $I_t$.

**Target-Reference flow.** This data flow consists of **select layers** and **fuse layers**. Given a Target-Reference region tuple $T_n$, we denote the memory features for target region and reference regions in $T_n$ as $\bar{m}_n^t$ and $\bar{M}_n = \{\bar{m}_n^i\}_{i=1}^K$, respectively. The flow takes in those two kinds of features and generates cross-image features through **select** and **fuse** layers:

*Fuse layer*. The goal of the fuse layer is to enable the interactions among memory features within $\bar{M}_n$. We stack multiple fuse layers, and the $l$-th fuse layer is calculated as follows:

$$U_l = \mathbf{MH}(U_{l-1}, U_{l-1}, U_{l-1}), \qquad (8)$$

where $U_0$ refers to $\bar{M}_n$ and $U_{l-1}$ is the outputs of the $(l-1)$-th fuse layer. Because all the features in $\bar{M}_n$ are semantic-similar, the model can learn to capture the primary concepts they are describing.

*Select layer*. The select layer builds on the co-attention mechanism. We set features of target region as query, features of reference regions as key and value in multi-head attention. Multiple co-attention layers are stacked, the $l$-th select layer is computed as:

$$V_l = \mathbf{MH}(V_{l-1}, U_{l-1}, U_{l-1}), \qquad (9)$$

where $V_0$ refers to $\bar{m}_n^t$, $V_{l-1}$ and $U_{l-1}$ are the outputs of the $(l-1)$-th select and fuse layer, respectively. By the residual connection in self-attention blocks, feature $\bar{m}_n^t$ will gradually select useful information from reference images while preserving the original information from $I_t$.

As an example shown in Figure 4 (b), our model can learn to focus on the unique attributes and objects in $I_t$. For unique attributes, we send all reference regions of "couch" (green boxes) into fuse layers, the model will be informed they are describing the concept "couch". Since target region also describes "couch", the select layer learns to focus on the unique color "red" of the "couch" in $I_t$. When predicting unique objects, for "vase" region in $I_t$, because all selected reference regions for it (blue boxes) do not contain the same concept, the select layer learns that "vase" is a unique object in $I_t$.

We use the outputs of the last select layer as final refined target feature $\tilde{m}_n^t$ for region $r_n^t$ in $I_t$. For N Target-Reference tuples in $I_t$, we can get $\tilde{M}_t = \{\tilde{m}_i^t\}_{i=1}^N$ as the outputs of Target-Reference flow. Finally, we concatenate the outputs of Target flow and Target-Reference flow as the final outputs of the Two-Flow Encoder:

$$M_t' = [\hat{M}_t; \tilde{M}_t], \qquad (10)$$

where $[\cdot; \cdot]$ denotes concatenation operation, and $M_t'$ will be sent into decoder for caption generation.

# 5 EXPERIMENTS

In this section, we describe the datasets used for experiments and introduce a new distinctiveness-based evaluation metric **DisCIDEr**.
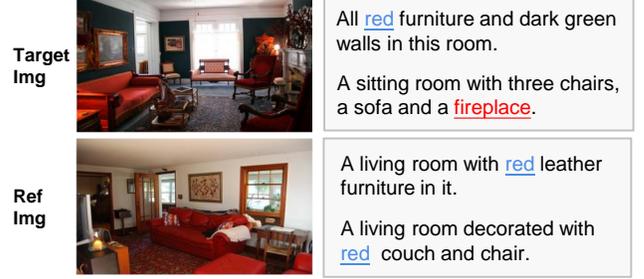


**Figure 5: An example of the intuition behind DisCIDEr. N-gram "red" appears in both target and reference images, thus should be given less attention (blue). In contrast, we should pay more attention to "fireplace", because it appears only in target image (red).**

We conduct extensive experiments and ablation studies to reveal the superiority of our proposed model, as well as our proposed benchmarks for Ref-DIC.

## 5.1 Datasets

We developed the **COCO-DIC** and **Flickr30K-DIC** based on the MS-COCO [9] and Flickr30K [29]. They contain 123287 and 31014 images, respectively. Each image is annotated with 5 ground-truth captions. For both datasets, we followed the splits provided by [16], and constructed reference image groups within the training, validation, and test splits. For completeness, we also reported results on Wang *et al.* [43]'s dataset for Ref-DIC.

## 5.2 Evaluation Metrics

We applied two kinds of metrics to evaluate the accuracy and distinctiveness of generated captions. For accuracy evaluation, we calculated four commonly used evaluation metrics: BLEU-N (B-N) (1- to 4-grams) [27], ROUGE-L (R) [20], METEOR (M) [3], and CIDEr (C) [40]. For distinctiveness evaluation, we developed a new metric named **DisCIDEr** (DisC). We introduce the DisCIDEr below.

**DisCIDEr.** All existing metrics designed for Ref-DIC task fail to fully explore the distinctiveness of each individual n-gram in GT captions of target image. To solve this, we assign these n-grams with different weights according to group-level distinctiveness: if an n-gram occurs frequently in ground-truth captions of reference images, it is less distinctive. As an example in Figure 5, both target and reference images are describing "red sofa", so we should assign lower weights to the word "red" in ground-truth captions of $I_t$ at evaluation time. Instead, since only image $I_t$ contains the object "fireplace", we should put more weight on it.

To realize this intuition, we modifies the n-gram weighting procedure of CIDEr by adding a re-weight term. For a target image $I_t$, we denote its generated and ground-truth captions as $c$ and $S_t = \{s_t^i\}_{i=1}^M$, respectively. Similarly, we denote ground-truth captions for reference images $\mathcal{I}_r$ as $S_r = \{s_r^j\}_{j=1}^M, r = \{1 \dots K\}$. The number of times an n-gram $\omega_d$ occurs in $s_t^j$ is denoted as $h_d(s_t^j)$. We modify CIDEr by adding an **Inverse reference frequency** term after it when calculating $g_d(s_t^j)$ for the n-grams in ground-truth

**Table 2: Comparison of captions accuracy on COCO family with state-of-the-art image captioning models.**

| Model | B-1 | B-4 | M | R | C | DisC |
|---|---|---|---|---|---|---|
| **Dataset: MS-COCO** | | | | | | |
| UpDown [1] | 79.8 | 36.3 | 27.7 | 56.9 | 120.1 | — |
| AoANet [14] | 80.2 | 38.9 | 29.2 | 58.8 | 129.8 | — |
| Transformer [38] | 80.0 | 38.2 | 28.9 | 58.2 | 127.3 | 98.7 |
| $M^2$ Transformer [10] | 80.8 | 39.1 | 29.2 | 58.6 | 131.2 | — |
| DiscCap [26] | — | 36.1 | 27.4 | 57.3 | 114.3 | — |
| CIDErBtwCap [42] | — | 38.5 | 29.1 | 58.8 | 127.8 | — |
| **Dataset: COCO-DIC** | | | | | | |
| GdisCap [43] | 80.0 | 37.3 | 28.4 | 57.5 | 125.8 | 96.6 |
| CAGC [19] | 80.7 | 38.1 | 28.7 | 57.9 | 127.9 | 98.0 |
| **TransDIC (Ours)** | **81.6** | **39.3** | **29.2** | **58.5** | **132.0** | **102.2** |
| **Dataset: Wang _et al._ [43]** | | | | | | |
| GdisCap [43] | 80.2 | 37.7 | 28.3 | 57.3 | 126.6 | 97.7 |
| CAGC [19] | 80.4 | 37.7 | 28.7 | 57.6 | 127.2 | 98.0 |
| **TransDIC (Ours)** | **81.0** | **38.8** | **29.1** | **58.2** | **130.8** | **101.9** |

**Table 3: Comparison of captions accuracy on Flickr30K family with state-of-the-art image captioning models.**

| Model | B-1 | B-4 | M | R | C | DisC |
|---|---|---|---|---|---|---|
| **Dataset: Flickr30K** | | | | | | |
| NIC [41] | 66.3 | 18.3 | — | — | — | — |
| Xu _et al._ [46] | 66.9 | 19.9 | 18.5 | — | — | — |
| Transformer [38] | 70.7 | 27.7 | 21.4 | 49.0 | 61.2 | 39.1 |
| **Dataset: Flickr30K-DIC** | | | | | | |
| GdisCap [43] | 71.7 | 29.0 | 22.1 | 49.6 | **65.6** | 41.2 |
| CAGC [19] | 72.9 | 29.1 | 21.9 | 50.1 | 62.2 | 39.0 |
| **TransDIC** | **73.2** | **30.1** | **22.5** | **50.3** | 65.1 | **41.4** |

**Table 4: Ablation study of Two-Flow Encoder on COCO-DIC. "Fuse" and "Select" denote the fuse layer and select layer in the Target-Reference flow, respectively.**

| Fuse | Select | B-1 | B-4 | M | R | C | DisC |
|---|---|---|---|---|---|---|---|
| ✗ | ✗ | 80.0 | 38.2 | 28.9 | 58.2 | 127.3 | 98.7 |
| ✗ | ✓ | 81.3 | 38.4 | 29.0 | 58.1 | 130.0 | 100.2 |
| ✓ | ✗ | **81.6** | 39.1 | **29.1** | 58.4 | 131.3 | 101.6 |
| ✓ | ✓ | **81.6** | **39.3** | **29.1** | **58.5** | **132.0** | **102.2** |

captions of $I_t$:

$$g_d(s_t^i) = \frac{h_d(s_t^i)}{\sum_{w_l \in \Omega} h_l(s_t^i)} \underbrace{\log(\frac{|I|}{\max(1, \sum_{I_p \in I} \min(1, \sum_{s_p^q \in I_p} h_d(s_p^q)))})}_{\text{Inverse document frequency}}$$

$$\underbrace{\log(\frac{m + K}{n + \sum_{S_u \in S_{1:K}} \min(1, \sum_{s_u^v \in S_u} h_d(s_u^v))})}_{\text{Inverse reference frequency}}, \quad (11)$$

where $\Omega$ is the vocabulary of all n-grams and $I$ is the set of all images. $m$ and $n$ are two parameters. In this way, DisCIDEr can evaluate the group-level distinctiveness while preserve advantage of n-gram based metric. We refer reader to [40] for more details.

## 5.3 Implementation Details

Following [1], we used the region proposal features extracted by the Faster R-CNN [33] with dimension 2048, and the memory space was of dimension 512. The number of self-attention blocks in Target flow, select and fuse layers in Target-Reference flow were set to 3. In the two-stage matching procedure, the size of $\mathcal{I}_c$ was 500. For $\mathcal{I}_r$, we used images with top-p to top-(p+K-1) highest similarity scores, where p is an adjustable parameter for group similarity and both benchmarks set p to 3, K to 5. Parameters $m$ and $n$ in DisCIDEr were set to 0.8 and 5.0.

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

We reported our results on two kinds of datasets: 1) Our constructed COCO-DIC and Flickr30K-DIC datasets. 2) The dataset proposed in [43]. We compared our TransDIC model with three kinds of state-of-the-art models: 1) **NIC** [41], **Xu _et al._** [46], **UpDown** [1], **AoANet** [14], **Transformer** [38], $M^2$ **Transformer** [10]. They only aim to generate captions with high accuracy. 2) **DiscCap** [26], **CIDErBtwCap** [42]. They are designed for the Single-DIC. 3) **Gdis-Cap** [43] that is designed for the Ref-DIC. We also compared **CAGC** [19] which use multiple images as input.

**Results on COCO-family Benchmarks.** From Table 2, we can observe: 1) For accuracy evaluation, our proposed TransDIC achieves the best performance on all conventional metrics at both COCO-DIC and Wang _et al._ [43] (_e.g._, 132.0 vs. 125.8 in GdisCap on CIDEr). Meanwhile, our model outperforms some strong state-of-the-art models (_e.g._, 132.0 vs. 131.2 in $M^2$ Transformer on CIDEr) in terms of accuracy-based metrics. 2) For distinctiveness evaluation, our model gets the highest scores on DisCIDEr in the two datasets.

**Results on Flickr30K-family Benchmarks.** From Table 3, we can observe: 1) For accuracy evaluation, our TransDIC achieves the largest performance gains on most metrics while it is narrowly defeated by GdisCap on CIDEr (65.1 vs. 65.6). 2) For distinctiveness-based metrics, our model outperforms GdisCap by 0.2 on DisCIDEr (41.4 vs. 41.2), despite it having a lower CIDEr score.

## 5.5 Ablation Studies

We conducted extensive experiments to verify the influences of the proposed Two-Flow Encoder module and group similarity.

_5.5.1 Influence of Two-Flow Encoder._ To measure the influence of each component in our proposed Target-Reference flow, we trained an ordinary transformer as the baseline and three variants of our model. 1) Target-Reference flow only contains the select layer: stacked select layers always take the original reference features as input. 2) Target-Reference flow only contains the fuse layer: outputs of the last fuse layer are directly used as the outputs of the Target-Reference flow. 3) Transformer with complete select and fuse layers. All these models were trained on COCO-DIC and the results were shown in Table 4.

**Results**. As can be observed in rows 2 and 3, two additional components can improve captioning performance consistently in terms of both accuracy and distinctiveness. Above all, our complete model achieves the most promising performance in all metrics.
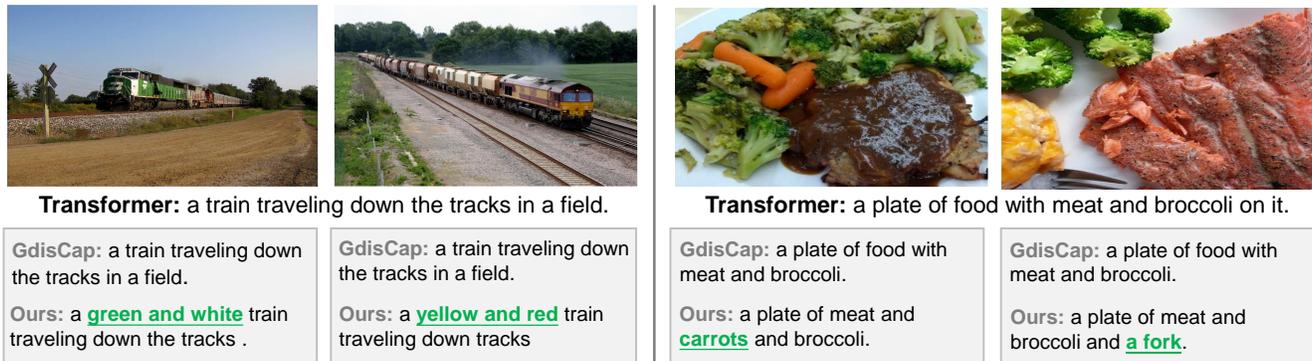
**Transformer:** a train traveling down the tracks in a field.

> **GdisCap:** a train traveling down the tracks in a field.
>
> **Ours:** a **green and white** train traveling down the tracks .

> **GdisCap:** a train traveling down the tracks in a field.
>
> **Ours:** a **yellow and red** train traveling down tracks

**Transformer:** a plate of food with meat and broccoli on it.

> **GdisCap:** a plate of food with meat and broccoli.
>
> **Ours:** a plate of meat and **carrots** and broccoli.

> **GdisCap:** a plate of food with meat and broccoli.
>
> **Ours:** a plate of meat and broccoli and **a fork**.

**Figure 6: Examples of generated captions for two similar images. The green words indicate the unique details in the images.**

**Table 5: Ablation study of group similarity on COCO-DIC and Flickr30K-DIC.**

(a) Comparison of different group trained with TransDIC on COCO-DIC

| TransDIC at COCO-DIC | | | | |
|---|---|---|---|---|
| Group | B-1 | B-4 | C | DisC |
| top1-5 | **81.6** | 39.1 | **132.2** | 101.4 |
| top2-6 | 81.4 | **39.3** | 131.5 | 101.4 |
| top3-7 | **81.6** | **39.3** | 132.0 | **102.2** |
| top4-8 | 81.0 | 38.7 | 130.7 | 101.5 |

(b) Comparison of different group trained with GdisCap on COCO-DIC

| GdisCap at COCO-DIC | | | | |
|---|---|---|---|---|
| Group | B-1 | B-4 | C | DisC |
| top1-5 | **80.4** | **38.0** | **126.9** | 97.0 |
| top2-6 | 79.8 | 37.5 | 125.1 | 96.2 |
| top3-7 | 80.0 | 37.3 | 125.8 | 96.6 |
| top4-8 | 80.2 | 37.6 | 125.3 | **97.1** |

(c) Comparison of different group trained with TransDIC on Flickr30K-DIC

| TransDIC at Flickr30K-DIC | | | | |
|---|---|---|---|---|
| Metric | B-1 | B-4 | C | DisC |
| top1-5 | 72.5 | 29.4 | 64.5 | 40.6 |
| top2-6 | 70.8 | 28.8 | 62.4 | 39.3 |
| top3-7 | 73.2 | 30.1 | 65.1 | 41.6 |
| top4-8 | **73.5** | **30.9** | **66.7** | **42.0** |

*5.5.2 Influence of Group Similarity.* To quantify the influence of group similarity, we used different p when choosing K most similar images from $\mathcal{I}_f$, *e.g.*, top-2 to top-6 **(top2-6)** group when setting p to 2. The results were reported in Table 5.

**Results**. From Table 5 (a), we can observe: Top1-5 group is surpassed by top3-7 group on DisCIDEr (101.4 vs. 102.2), despite it having a marginal improvement on CIDEr (132.2 vs. 132.0). The results indicate that the most similar reference group is not always helpful to group-level distinctiveness.

## 5.6 Qualitative Results

We illustrated the qualitative results of our proposed TransDIC and compared it with the Transformer and SOTA Ref-DIC model GdisCap [43] in Figure 6. Naive captioning models generate identical captions for similar images. In contrast, our TransDIC can describe the unique attributes and objects in the target image. For unique attributes, as shown in Figure 6 (left), TransDIC precisely captures the unique attributes "green and white" and "yellow and red" for the two trains, respectively. In Figure 6 (right), for unique objects, TransDIC captures unique objects "carrots" and "a fork" for each individual image. The results demonstrate that TransDIC can generate distinctive captions in terms of unique objects and attributes.

## 6 LIMITATIONS

One possible limitation of our work is that if original human-annotated captions omit some objects or attributes, it will lead to:

1) our proposed two-stage matching mechanism may fail to collect object-/attribute- level similarity reference images. 2) our proposed DisCIDEr may degrade to existing CIDEr. We believe this omitting problem is due to the natural defect of datasets (COCO/Flickr30K): Since these datasets are annotated for general captioning task, human annotators may tend to simply describe the objects while ignoring its attributes (*e.g.*, color) when there is no reference image.

## 7 CONCLUSIONS

In this paper, we argued that all existing DIC works fail to achieve group-level distinctiveness. To solve this problem, firstly, we introduced a two-stage matching mechanism and proposed two new benchmarks for Ref-DIC. Then, we developed a Transformer-based model for Ref-DIC. Finally, we came up with a new evaluation metric termed DisCIDEr. We conducted extensive experiments to verify the effectiveness of TransDIC. Moving forward, we are going to 1) extend our Ref-DIC into video domains, or 2) design stronger Ref-DIC models.

# REFERENCES

[1] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. 2018. Bottom-up and top-down attention for image captioning and visual question answering. In *CVPR*.

[2] J. L. Ba, J. R. Kiros, and G. E. Hinton. 2016. Layer Normalization. (2016).

[3] Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *ACL workshop*.

[4] Fuhai Chen, Rongrong Ji, Xiaoshuai Sun, Yongjian Wu, and Jinsong Su. 2018. Groupcap: Group-based image captioning with structured relevance and diversity constraints. In *CVPR*. 1345–1353.

[5] Long Chen, Zhihong Jiang, Jun Xiao, and Wei Liu. 2021. Human-like controllable image captioning with verb-specific semantic roles. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 16846–16856.

[6] Long Chen, Hanwang Zhang, Jun Xiao, Xiangnan He, Shiliang Pu, and Shih-Fu Chang. 2019. Counterfactual critic multi-agent training for scene graph generation. In *ICCV*.

[7] Long Chen, Hanwang Zhang, Jun Xiao, Liqiang Nie, Jian Shao, Wei Liu, and Tat-Seng Chua. 2017. Sca-cnn: Spatial and channel-wise attention in convolutional networks for image captioning. In *CVPR*.

[8] Shizhe Chen, Qin Jin, Peng Wang, and Qi Wu. 2020. Say As You Wish: Fine-grained Control of Image Caption Generation with Abstract Scene Graphs. In *CVPR*.

[9] X. Chen, H. Fang, T. Y. Lin, R. Vedantam, S. Gupta, P. Dollar, and C. L. Zitnick. 2015. Microsoft COCO Captions: Data Collection and Evaluation Server. *arXiv* (2015).

[10] M. Cornia, M. Stefanini, L. Baraldi, and R. Cucchiara. 2020. Meshed-Memory Transformer for Image Captioning. In *CVPR*.

[11] Bo Dai, Sanja Fidler, Raquel Urtasun, and Dahua Lin. 2017. Towards diverse and natural image descriptions via a conditional gan. In *ICCV*.

[12] Bo Dai and Dahua Lin. 2017. Contrastive learning for image captioning. In *NeurIPS*.

[13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *CVPR*.

[14] Lun Huang, Wenmin Wang, Jie Chen, and Xiao-Yong Wei. 2019. Attention on attention for image captioning. In *ICCV*. 4634–4643.

[15] Justin Johnson, Andrej Karpathy, and Li Fei-Fei. 2016. Densecap: Fully convolutional localization networks for dense captioning. In *CVPR*. 4565–4574.

[16] Andrej Karpathy and Li Fei-Fei. 2015. Deep visual-semantic alignments for generating image descriptions. In *CVPR*.

[17] Guang Li, Linchao Zhu, Ping Liu, and Yi Yang. 2019. Entangled transformer for image captioning. In *ICCV*. 8928–8937.

[18] Lin Li, Long Chen, Yifeng Huang, Zhimeng Zhang, Songyang Zhang, and Jun Xiao. 2022. The Devil is in the Labels: Noisy Label Correction for Robust Scene Graph Generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 18869–18878.

[19] Zhuowan Li, Quan Tran, Long Mai, Zhe Lin, and Alan L Yuille. 2020. Context-aware group captioning via self-attention and contrastive features. In *CVPR*. 3440–3450.

[20] Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *ACL workshop*.

[21] An-An Liu, Hongshuo Tian, Ning Xu, Weizhi Nie, Yongdong Zhang, and Mohan Kankanhalli. 2021. Toward region-aware attention learning for scene graph generation. *IEEE Transactions on Neural Networks and Learning Systems* (2021).

[22] An-An Liu, Yingchen Zhai, Ning Xu, Weizhi Nie, Wenhui Li, and Yongdong Zhang. 2021. Region-Aware Image Captioning via Interaction Learning. *IEEE Transactions on Circuits and Systems for Video Technology* (2021).

[23] Lixin Liu, Jiajun Tang, Xiaojun Wan, and Zongming Guo. 2019. Generating Diverse and Descriptive Image Captions Using Visual Paraphrases. In *ICCV*.

[24] Siqi Liu, Zhenhai Zhu, Ning Ye, Sergio Guadarrama, and Kevin Murphy. 2017. Improved image captioning via policy gradient optimization of spider. In *ICCV*. 873–881.

[25] Xihui Liu, Hongsheng Li, Jing Shao, Dapeng Chen, and Xiaogang Wang. 2018. Show, tell and discriminate: Image captioning by self-retrieval with partially labeled data. In *ECCV*.

[26] Ruotian Luo, Brian Price, Scott Cohen, and Gregory Shakhnarovich. 2018. Discriminability objective for training descriptive captions. In *CVPR*.

[27] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: a method for automatic evaluation of machine translation. In *ACL*.

[28] Dong Huk Park, Trevor Darrell, and Anna Rohrbach. 2019. Robust change captioning. In *ICCV*. 4624–4633.

[29] B. A. Plummer, L. Wang, C. M. Cervantes, J. C. Caicedo, and S. Lazebnik. 2016. Flickr30k Entities: Collecting Region-to-Phrase Correspondences for Richer Image-to-Sentence Models. In *ICCV*.

[30] Yue Qiu, Shintaro Yamamoto, Kodai Nakashima, Ryota Suzuki, Kenji Iwata, Hirokatsu Kataoka, and Yutaka Satoh. 2021. Describing and Localizing Multiple Changes with Transformers. In *ICCV*. 1971–1980.

[31] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *ICML*. 8748–8763.

[32] Marc'Aurelio Ranzato, Sumit Chopra, Michael Auli, and Wojciech Zaremba. 2015. Sequence level training with recurrent neural networks. *arXiv* (2015).

[33] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2016. Faster r-cnn: Towards real-time object detection with region proposal networks. *IEEE TPAMI* (2016).

[34] Steven J Rennie, Etienne Marcheret, Youssef Mroueh, Jerret Ross, and Vaibhava Goel. 2017. Self-critical sequence training for image captioning. In *CVPR*.

[35] Rakshith Shetty, Marcus Rohrbach, Lisa Anne Hendricks, Mario Fritz, and Bernt Schiele. 2017. Speaking the same language: Matching machine to human captions by adversarial training. In *ICCV*.

[36] Alane Suhr, Stephanie Zhou, Ally Zhang, Iris Zhang, Huajun Bai, and Yoav Artzi. 2019. A corpus for reasoning about natural language grounded in photographs. (2019).

[37] Hao Tan, Franck Dernoncourt, Zhe Lin, Trung Bui, and Mohit Bansal. 2019. Expressing visual relationships via language. (2019).

[38] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *NeurIPS*.

[39] Ramakrishna Vedantam, Samy Bengio, Kevin Murphy, Devi Parikh, and Gal Chechik. 2017. Context-aware captions from context-agnostic supervision. In *CVPR*.

[40] Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. 2015. Cider: Consensus-based image description evaluation. In *CVPR*.

[41] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. 2015. Show and tell: A neural image caption generator. In *CVPR*.

[42] Jiuniu Wang, Wenjia Xu, Qingzhong Wang, and Antoni B Chan. 2020. Compare and Reweight: Distinctive Image Captioning Using Similar Images Sets. In *ECCV*.

[43] Jiuniu Wang, Wenjia Xu, Qingzhong Wang, and Antoni B Chan. 2021. Group-based distinctive image captioning with memory attention. In *ACM MM*. 5020–5028.

[44] Yanhui Wang, Ning Xu, An-An Liu, Wenhui Li, and Yongdong Zhang. 2021. High-Order Interaction Learning for Image Captioning. *IEEE Transactions on Circuits and Systems for Video Technology* (2021).

[45] Zeyu Wang, Berthy Feng, Karthik Narasimhan, and Olga Russakovsky. 2020. Towards Unique and Informative Captioning of Images. In *ECCV*.

[46] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. 2015. Show, attend and tell: Neural image caption generation with visual attention. In *ICML*.

[47] Ning Xu, An-An Liu, Yongkang Wong, Weizhi Nie, Yuting Su, and Mohan Kankanhalli. 2020. Scene graph inference via multi-scale context modeling. *IEEE Transactions on Circuits and Systems for Video Technology* (2020), 1031–1041.

[48] Ning Xu, Hanwang Zhang, An-An Liu, Weizhi Nie, Yuting Su, Jie Nie, and Yongdong Zhang. 2019. Multi-level policy and reward-based deep reinforcement learning framework for image captioning. *IEEE Transactions on Multimedia* (2019), 1372–1383.

[49] An Yan, Xin Eric Wang, Tsu-Jui Fu, and William Yang Wang. 2021. L2C: Describing Visual Differences Needs Semantic Understanding of Individuals. (2021).

[50] Xu Yang, Kaihua Tang, Hanwang Zhang, and Jianfei Cai. 2019. Auto-encoding scene graphs for image captioning. In *CVPR*. 10685–10694.

[51] Ting Yao, Yingwei Pan, Yehao Li, and Tao Mei. 2018. Exploring visual relationship for image captioning. In *ECCV*. 684–699.

## Appendix. A  USER STUDY INTERFACE

We conduct user studies to validate the effectiveness of our proposed **DisCIDEr**. Specifically, we invite 5 experts and give them a picture and two captions that describe it. They are asked to choose which one better matches the image in terms of accuracy and distinctiveness and the user interface for this is shown in Fig. 7. The caption which got more than 3 votes is regarded as human judgment. We randomly select 100 images from the test split of MS-COCO and ask the experts to give their judgments. A metric agrees with human judgment only if it gives a higher score to the same caption that experts choose. We use the agreement counts between human judgment and the metric as the effectiveness measurement.
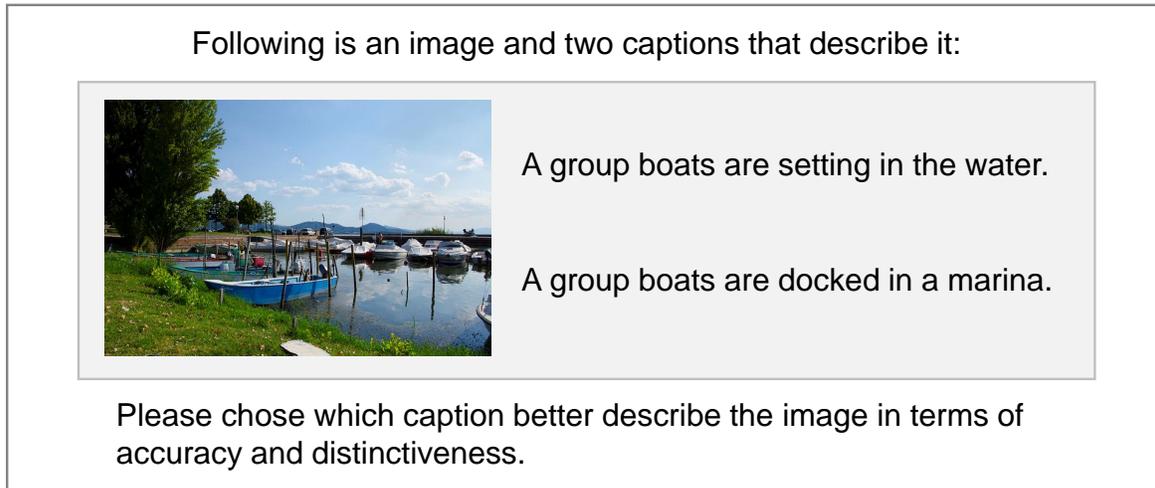


**Figure 7: We display an image and two captions generated by Transformer [38] and our proposed TransDIC, respectively. The users are asked to choose which caption better describe the image in terms of accuracy and distinctiveness.**

**User Studies on DisCIDEr.** We conducted a user study to validate the effectiveness of **DisCIDEr** with 5 experts. We randomly selected 100 images (100 trials) from the test set. In each trial, an image and two captions (generated by Transformer and TransDIC, respectively) were displayed and these experts are asked to choose the better caption in terms of distinctiveness and accuracy. These captions which got more than 3 votes were regarded as human judgment. We calculated the agreements between human judgments and different metrics (*i.e.*, whether humans and metrics give higher scores to a same caption). Results were reported in Table 6. From Table 6, we can observe that DisCIDEr achieves better agreement than both the accuracy-based metric CIDEr and the distinctiveness-based metric DisWordRate.

**Table 6: User study of different metrics**

| Metric | CIDEr [40] | DisWordRate [43] | DisCIDEr |
|---|---|---|---|
| Agreements | 58 | 62 | 64 |

## Appendix. B  DETAILS OF THESE COMPARED BASELINES

In this section, we describe implementation details of two compared state-of-the-art baselines: GdisCap [43] and CAGC [19]. These two group-based captioning baselines are both built on top of the widely-used Transformer [38] architecture, which takes a sequence of visual features as input and encode them through multiple self-attention layers. We introduce the two baselines in detail.

(1) **GdisCap** [43]: It develops a Group-based Memory Attention (GMA) module which assigns higher weights to the distinctive regions and two special losses: DisLoss and MemCls loss, which directly encourage the model to generate distinctive words. We re-implement their proposed GMA module and DisLoss.

(2) **CAGC** [19]: It is designed to describe a group of target images in the context of another group of related reference images. To this end, they extract visual features for all images in target and reference groups from the ResNet50 network [13] (after pool5 layer). Then, these Visual features of images are sent into multiple self-attention layers to enable the interactions within each group (target image interacts with other images only within target group, and same for reference group). Finally, they construct group representations and contrastive representations for both groups and send them into a decoder for captions generation. To make a fair comparison, we replace the visual features with mean pooled bottom-up features [1] and re-implement their visual content encoder.

## Appendix. C    MORE QUALITATIVE RESULTS

We report more qualitative results in Fig. 8 and Fig. 9. In Fig. 8, we compare our TransDIC with Transformer and GdisCap. Remarkably, our model can precisely capture the unique objects and attributes in the image. For example in Fig. 8 (row2, column 3), both Transformer and GdisCap wrongly mix the tie and shirt into "blue shirt" while TransDIC correctly matches the "shirt" with "black" and "tie" with "blue". In terms of numerals, in Fig. 8 (row2, column 1), only TransDIC predicts the right number, whereas others mistake it for "one". In Fig. 9, we test all models with two similar images, as can be seen on the left side, TransDIC correctly captures the "red and white" and "green and red" for the two buses, respectively.
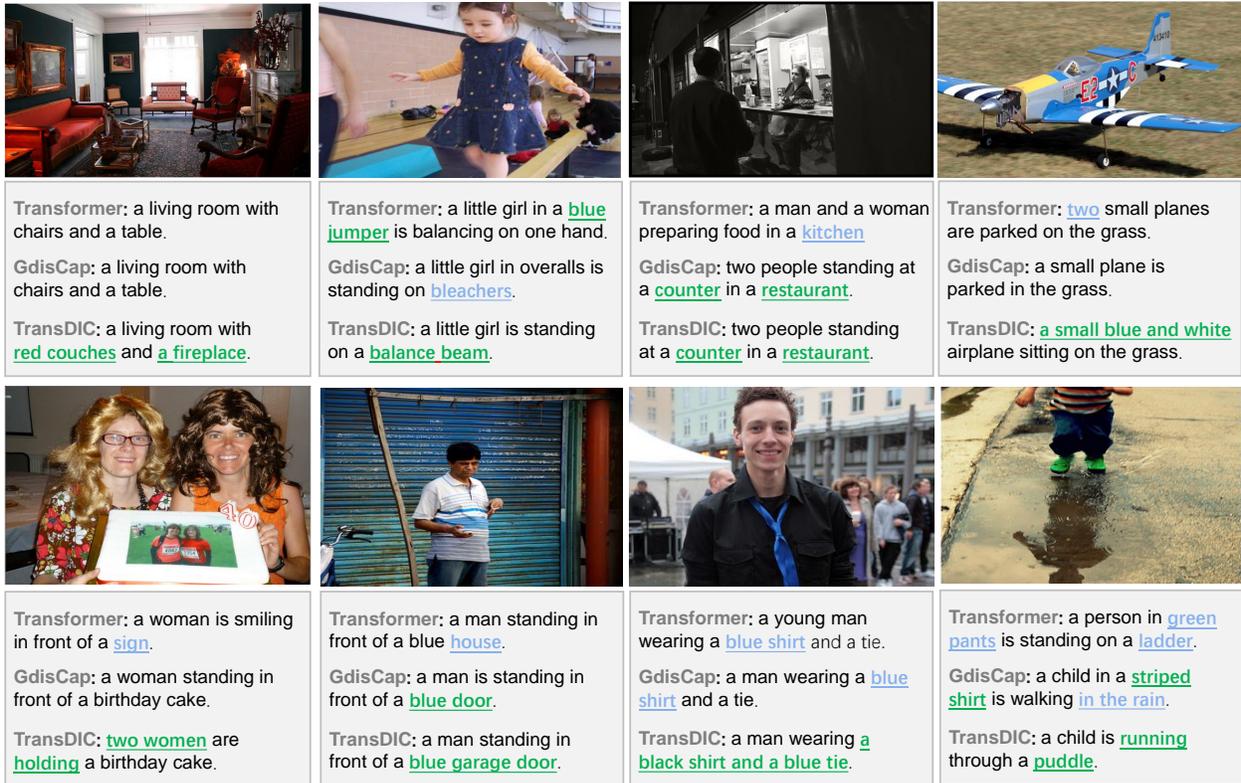


**Figure 8: Examples of generated captions for two similar images using Transformer, GdisCap and TransDIC. The green words indicate the unique details in the images while blue denote the mistakes in the generated captions.**
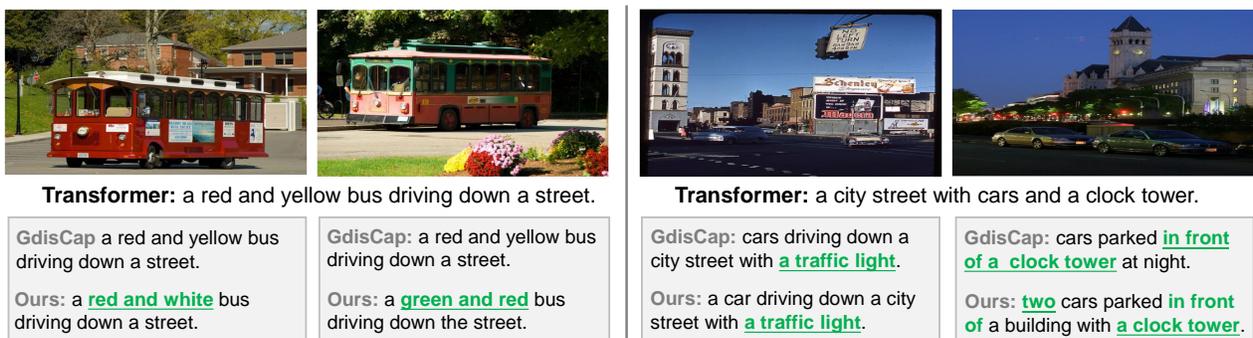


**Figure 9: Examples of generated captions for two similar images. The green words indicate the unique details in the images.**