# Compare and Reweight: Distinctive Image Captioning Using Similar Images Sets

Jiuniu Wang[1,2,3], Wenjia Xu[2,3,4], Qingzhong Wang[1], and Antoni B. Chan[1]

[1] Department of Computer Science, City University of Hong Kong
{jiuniwang2-c,qingzwang2-c}@my.cityu.edu.hk, abchan@cityu.edu.hk
[2] Aerospace Information Research Institute, Chinese Academy of Sciences
[3] University of Chinese Academy of Sciences
xuwenjia16@mails.ucas.ac.cn
[4] Max Planck Institute for Informatics

**Abstract** A wide range of image captioning models has been developed, achieving significant improvement based on popular metrics, such as BLEU, CIDEr, and SPICE. However, although the generated captions can accurately describe the image, they are generic for similar images and lack distinctiveness, *i.e.*, cannot properly describe the uniqueness of each image. In this paper, we aim to improve the distinctiveness of image captions through training with sets of similar images. First, we propose a distinctiveness metric — between-set CIDEr (CIDErBtw) to evaluate the distinctiveness of a caption with respect to those of similar images. Our metric shows that the human annotations of each image are not equivalent based on distinctiveness. Thus we propose several new training strategies to encourage the distinctiveness of the generated caption for each image, which are based on using CIDErBtw in a weighted loss function or as a reinforcement learning reward. Finally, extensive experiments are conducted, showing that our proposed approach significantly improves both distinctiveness (as measured by CIDErBtw and retrieval metrics) and accuracy (*e.g.*, as measured by CIDEr) for a wide variety of image captioning baselines. These results are further confirmed through a user study. Project page: https://wenjiaxu.github.io/ciderbtw/.

## 1 Introduction

Image captioning is attracting increasing attention from researchers in the fields of computer vision and natural language processing. It is promising in various applications such as human-computer interaction and medical image understanding [36,24,11,41,3,37]. Currently, the limitation of image captioning models is that the generated captions tend to consist of common words so that many images have similar or even the same captions (see Figure 1). The distinctive concepts in images are ignored, which limits the application of image captioning. Although, auxiliary information such as where, when and who takes the picture could be used to generate personalized captions [4,26], many images do not have such information. In terms of the quality of generated captions, [21] summarizes four attributes that encourage auto-generated captions to resemble human

| CIDErBtw | Human Ground-truth Captions: |
|---|---|
| 53.5 | 1: A living room with a big table next to a book shelf. |
| 40.2 | 2: The large room has a wooden table with chairs and a couch. |
| 54.2 | 3: A living room decorated with a modern theme. |
| 73.0 | 4: A living room with wooden floors and furniture. |
| | Machine generated captions: |
| 141.5 | Baseline: A living room with a couch and a table. |
| 68.5 | Ours: A living room filled with wooden table and a large window. |

| CIDErBtw | Human Ground-truth Captions: |
|---|---|
| 55.6 | 1: A living room filled with nice furniture and a persian rug. |
| 55.9 | 2: An image of a living room setting with furniture and curtains. |
| 38.9 | 3: An open living room with brown walls and beige carpeting. |
| 31.2 | 4: A large tan living room bathed in sunlight. |
| | Machine generated captions: |
| 174.2 | Baseline: A living room with a couch and a table. |
| 88.3 | Ours: A living room with a white couch and a painting. |

**Figure 1:** The human ground-truth captions of a target image and a semantically similar image contain both common words (highlighted in green) and distinctive words (highlighted in red for the target, and blue for the similar image). The baseline model, Transformer [33] trained with MLE and SCST, generates the same caption for both images, while our model generates distinctive captions with words unique to each image. The distinctiveness is measured using CIDErBtw, the CIDEr metric between the target caption and the GT captions of the similar images set, where lower values mean more distinctive.

language: fluency, relevance, diversity, and descriptiveness. Various models and metrics have been proposed to improve the fluency and relevance of the captions so as to obtain accurate results. However, these captions are poor at mimicking the inherent characteristics of human language: *distinctiveness*, which refers to the specific and detailed aspects of the image that distinguish it from other similar images.

Some recent works have focused on generating more diverse and descriptive captions, with techniques such as conditional generative adversarial networks (GANs) [5,30], self-retrieval [6,23,35] and two-stage LSTM [21]. Some works propose metrics for evaluating the *diversity* of a set of generated captions for a single image, based on the percentage of unique n-grams or novel sentences [30] or the similarity between pairs of captions [38]. However, only encouraging the diversity, such as using synonyms or changing word order, may not help with generating distinctive captions among multiple similar images. For instance, the human caption in Figure 1 "*an image of a living room setting with furniture and curtains*" is telling the same story as "*a living room with furniture and curtains*". Although the two sentences have different syntax and the first sentence is more diverse according to some metrics, the distinctiveness is not improved. In this paper, we mainly focus on promoting the *distinctiveness* of image captioning, where the caption should describe the important and specific aspects of an image that can distinguish it from other similar images. To evaluate distinctiveness, the retrieval metric is generally employed in recent works [6,22,23,21]. However, using self-retrieval in captioning models could lead to repetition problem [38,35], *i.e.*, the generated captions could repeat distinct words, which hurts language

fluency. Also, its result may vary when choosing different retrieval models or candidate images pool. In this work, we propose a general metric for distinctiveness, Between-Set CIDEr (CIDErBtw), by measuring the semantic distance between an image's caption and captions from a set of similar images. If the caption is distinct, i.e., captures unique concepts in its image, then it should have less overlap with its similar image set, i.e., lower CIDErBtw. We found that the human annotations of each image are not equivalent based on distinctiveness. Consider the example image and caption pairs shown in Figure 1, some ground-truth captions contain more distinct concepts (e.g., *bathed in sunlight*) and detailed description that can distinguish the image from its similar image (e.g., *wooden floor* and *brown walls*). However, traditional training objectives such as maximum likelihood estimation (MLE) and reinforcement learning (RL) treat every ground-truth caption equally. Thus, one possible method for improving distinctiveness is to give more weight to the distinctive ground-truth captions during training. In this way, the captioning model learns to focus on important visual objects or properties, and generate distinctive words instead of generic ones. In summary, the contributions of our paper are three-fold:

- We propose a novel metric CIDErBtw to evaluate the distinctiveness of captions within similar image sets. Experiments show that our metric aligns with human judgment for distinctiveness.
- We use CIDErBtw as guidance for training, encouraging the model to learn from more distinctive captions. Experiments show that training with CIDErBtw is generic and yields consistent improvement for many baseline models.
- Based on the transformer network trained with SCST (self-critical sequence training) [29] and CIDErBtw strategies, we generate distinctive captions while maintaining state-of-the-art performance according to evaluation metrics such as CIDEr and BLEU. Both automatic metrics and human evaluation demonstrate that our captions are more accurate and more distinctive.

## 2   Related work

**Captioning models.** A wide range of image captioning models have been developed [36,24,11,41,3,37], achieving satisfying results as measured by popular metrics, such as BLEU [25], CIDEr [34] and SPICE [1]. Generally, an image captioning model is composed of three modules: 1) visual feature extractor, 2) language generator, and 3) the connection between vision and language. Convolutional neural networks (CNNs) [31,14] are widely used as visual feature extractors. Recently, object-level features extracted by Faster-RCNN [28] have also been introduced into captioning models [2], significantly improving the performance of image captioning models. [42] proposed a hierarchy parsing model to fuse multi-level image features extracted by mask-RCNN [13], which improves the performance of the baseline models. In terms of language generators, LSTMs [15] and its variants are the most popular, while some works [3,37] use CNNs as the decoder since LSTMs cannot be trained in parallel. More recently, transformers [33,27,9] show improved performance in both language gen-

eration and language understanding, where the multi-head attention plays the most important role and the receptive field is much larger than CNNs. Stacking multi-head attention layers could mitigate the long-term dependency problem in LSTMs. Hence, the transformer model could handle much longer texts. For vision-language connection, attention mechanisms [41,29,2,16] are used to reveal the co-occurrence between concepts and objects in the images.

**Distinctive image captioning.** Previous works [6,5,38] reveal that training the captioning model with MLE loss or CIDEr reward result in over-generic captions, since the captioning models try to predict an "average" caption that is close to all ground-truth captions. These captions lack distinctiveness, *i.e.*, they describe images with similar semantic content using the same caption. Recently, various works aim to solve this problem. In summary, they propose three aspects to consider: (1) *diversity*: describe one image with notably different expressions every time like humans [5], or use rich and diverse wording [38] to generate captions; (2) *discriminability*: describe an image by referring to the important and detailed aspects of the image, which is accurate, and informative [22,23,21,35]; (3) *distinctiveness*: describe the important and specific aspects of an image that can distinguish the image from other similar images [6,21]. In our paper, we focus on the last aspect, distinctiveness.

To promote diversity, some works [5,30] employ GANs, where an evaluator distinguishes the generated captions from human annotations, encouraging the captions to be similar to human annotations. Instead of using generative models, VisPara-Cap [21] employs two-stage LSTM and visual paraphrases to improve diversity and discriminability, where the two-stage model is trained with a pair of ground-truth image captions from an image — the first caption is less complex, and the next one with rich information is more distinctive. In contrast, our method is based on weighting all the ground-truth captions according to their distinctiveness, which retains more information for training. During inference, VisPara-Cap [21] first generates a simple caption and then paraphrases it into a more distinctive caption, which is a two-stage model and time-consuming. Another drawback of the model is that it cannot be trained in SCST [29] manner, and therefore the performance based on BLEU [25], CIDEr [34], and SPICE [1] is limited. In contrast, our method is able to improve both traditional metric scores and distinctiveness, and it can be applied to any image captioning model.

Contrastive learning [6] and self-retrieval [22,23,35] are introduced into captioning models to improve the distinctiveness of the generated captions. DiscCap [23], CL [6] and PSST [35] employ image retrieval to optimize the contrastive loss, which aims at pushing the generated caption far from other images in the training batch. On one hand, image retrieval encourages a model to generate distinctive words, while on the other hand, it hurts the accuracy and caption quality — weighting too much on image retrieval could lead a model to repeat the distinctive words [38]. In contrast, we encourage the generated caption to learn from its own ground-truth captions, giving more weights to captions that are distinct from other similar images, and disregard those generic captions. Thus both accuracy and distinctiveness are promoted in our model.

**Metrics for distinctiveness.** Traditional metrics such as BLEU [25], ME-TEOR [7], ROUGE-L [19], CIDEr [34] and SPICE [1] normally consider the overlap between a generated caption and the ground-truth captions. These metrics treat all ground-truth equally, and thus a generated caption that only uses common words could obtain high scores, reflecting the statistics of human annotations. Some works aim to generate multiple captions to cover more concepts in an image [5,30,8,39] and several diversity metrics are proposed, such as the number of novel captions, the number of distinct n-grams [40], mBLEU [30], local and global word recall [32], and self-CIDEr [38]. However, these metrics only encourage the diversity and discriminability and do not explicitly evaluate distinctiveness. Although generating multiple captions could cover distinctive concepts, it is difficult to summarize them into one human-like description.

Currently, the retrieval approach is the most popular evaluation metric for distinctiveness. A generated caption is used as the query and a pre-trained image-text embedding model, *e.g.*, VSE++ [10], is employed to rank the given images, with recall at $K$ (R@$K$) normally used to measure the distinctiveness of captions. Ideally, a correct and distinctive caption should retrieve the image that was used to generate the caption. The drawback of retrieval-based metrics is that they are time-consuming, since it requires using a deep retrieval model. Moreover, different trained models could result in different R@$K$. In contrast, our proposed CIDErBtw metric for distinctiveness is fast and easy to implement, allowing it to be incorporated into various training protocols and captioning models.

## 3   Methodology

In this paper, we aim to obtain a distinct caption that describes the important, specific, and detailed aspects of an image. To achieve this goal, we train the captioning model to focus on important details that would distinguish the target image from semantically similar images. Our work involves two main components, the Between-Set CIDEr (CIDErBtw) that measures the distinctiveness of an image caption from those of similar images, and several strategies for training distinctive models based on CIDErBtw.

The image captioning model aims to generate a sentence $c^*$ to describe the semantics of the target image $I_0$. In the image caption dataset, the image $I_0$ is provided with $N$ annotated ground-truth captions $C^0 = \{c_1^0, c_2^0, \ldots, c_N^0\}$. We first find $K$ similar images $\{I_1, I_2, \ldots, I_K\}$ that are semantically similar to $I_0$, and then calculate the CIDErBtw values of $C^0$ using these similar images. During training process, CIDErBtw can be used as an indicator of which ground-truth captions deserve more attention, or as a part of the reward in reinforcement learning (RL). This will train the model to generate a caption different from those of the similar images. Moreover, CIDErBtw can work as an evaluation metric to measure distinctiveness.

### 3.1   Similar images set

According to the split of the training, validation, and testing dataset, we measure the similarity of the target image $I_0$ to every image within the same split. For each image $I_0$ in the dataset, we find the top $K$ images $\{I_1, I_2, \ldots, I_K\}$ with the highest semantic similarity to form a *similar images set*. Similar images sets in the training split are used when calculating the loss and the reward during training, while those in the validation and test split are used to evaluate the distinctiveness of generated captions.

Given every target image, we generate its similar images set according to an image-to-caption retrieval process. We use VSE++ [10] to encode images and captions into a joint semantic space, and obtain similar images sets via retrieval. Given target image $I_0$, we obtain a set of closest captions $\{c'_1, c'_2, \ldots, c'_{N'}\}$ in the joint space by image-to-caption retrieval, where $N' = N(K + 1)$ to ensure that at least $K+1$ images are obtained to construct the similar images set. The top $K$ images corresponding to this caption set are considered as similar to the target image $I_0$. When using the retrieval method, the similarity of $I_i$ to $I_j$ denoted as $S(I_i, I_j)$ can be expressed like

$$S(I_i, I_j) = \max_{k \in \{1, \cdots, N\}} g_r(I_i, c_k^j), \quad g_r(I_i, c_k^j) = \frac{\phi(I_i)^T \theta(c_k^j)}{\|\phi(I_i)\|\|\theta(c_k^j)\|}, \quad (1)$$
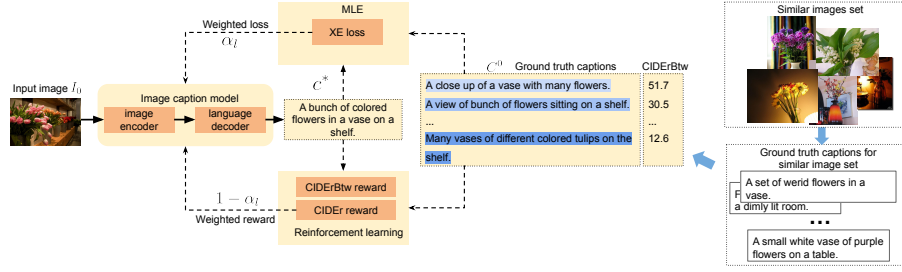
where $g_r(I_i, c_k^j)$ represents the retrieval score between the target image $I_i$ and the $k$-th ground-truth caption of $I_j$, and $\phi(\cdot)$ and $\theta(\cdot)$ are the image and caption encoders.

### 3.2   Between-set CIDEr (CIDErBtw)

Next, we introduce the definition of Between-set CIDEr (CIDErBtw) and its applications. In this paper, we mainly apply CIDErBtw in the following three aspects. During training, CIDErBtw is used to reweight the cross entropy (XE) loss and the reinforcement learning (RL) reward for each ground-truth caption. The CIDErBtw metric is also used directly as part of the reward to guide RL. During inference, CIDErBtw is used as a metric to measure the distinctiveness of a generated caption.

**CIDErBtw definition.** CIDErBtw reflects the distinctiveness of a caption $c$ by measuring the similarity of $c$ to the captions of similar images $C^{(s)}$. Specifically, given a caption $c$ for image $I_0$, the similar images set $\{I_1, I_2, \ldots, I_K\}$ retrieved in Section 3.1 and their ground-truth captions $C^{(s)} = \{c_n^k\}_{n=1, k=1}^{N, K}$, we define the CIDErBtw score of $c$ as

$$CIDErBtw(c) = \frac{1}{KN} \sum_{k=1}^{K} \sum_{n=1}^{N} g_c(c, c_n^k), \quad (2)$$

**Figure 2:** The framework of our CIDErBtw image captioning model. $\alpha_l$ is a hyperparameter that controls the weight of the two optimization modules. The solid and dashed lines represent the forward and backward process. $c^*$ and $C^0$ indicate the generated caption and the ground-truth captions. With CIDErBtw, we reweight the ground-truth captions when calculating the XE loss and reward. The shade of blue shows the CIDErBtw weight $w_i$ for each caption.

where $N$ is the number of ground-truth captions provided for each image, $g_c(c, c_n^k)$ represents the CIDEr value between $c$ and $c_n^k$. Actually, the methodology could be extended to use any caption metric to measure between-set similarity. Here we use CIDEr because it focuses more on low frequency words (through TF-IDF vectors) that could be more distinctive, is efficient to compute, and is the most frequently used metric to evaluate performance of image captioning models.

**CIDErBtw weight.** For conventional training strategies such as MLE and reinforcement learning, we maximize the likelihood or reward for the given ground-truth captions $C^0 = \{c_1^0, c_2^0, \ldots, c_N^0\}$. In previous methods, each ground-truth caption $c_i^0$ is treated equally, whereas these ground-truth might have different distinctiveness. In this work, we focus more attention to distinctive ground-truth captions by reweighting the training loss. For every training image $I_0$, we provide its $N$ ground-truth captions $C^0$ with different weights $W = \{w_1, w_2, \ldots, w_N\}$, according to their CIDErBtw scores $V = \{v_1, v_2, \ldots, v_N\}$,

$$v_i = CIDErBtw(c_i^0), \quad w_i = \lambda_w - \alpha_w \frac{v_i}{\max_i(v_i)}, \tag{3}$$

where $\lambda_w$ and $\alpha_w$ are hyperparameters. Here $w_i$ indicates the contribution of the $i$-th ground-truth caption during model training. More distinctive captions will have lower $v_i$, leading to higher weight $w_i$.

### 3.3 CIDErBtw training strategies

Figure 2 shows the overall framework of our CIDErBtw Image Caption model. The model is composed of a image encoder and language decoder. These two modules can generate a caption $c^*$ for input image $I_0$. There are two criteria to update the parameters of our image captioning model, the XE loss $\mathcal{L}_{XE}$ and RL

reward $\mathcal{L}_{RL}$. We apply a hyperparameter $\alpha_l$ to control the weight of these two criteria,

$$\mathcal{L} = \alpha_l \mathcal{L}_{XE} + (1 - \alpha_l)\mathcal{L}_{RL}. \tag{4}$$

Following SCST (self-critical sequence training) [29], the training process of our model can be divided into two steps. The first step only trains with $\mathcal{L}_{XE}$, setting $\alpha_l = 1$, and the second step only trains with $\mathcal{L}_{RL}$, setting $\alpha_l = 0$.

**Reweighting XE loss.** Given the words in a ground-truth caption $c_i^0 = \{d_1, d_2, \ldots, d_T\}$, XE loss can be expressed as

$$L_{XE}(c_i^0) = -\sum_{t=1}^{T} \log p_\theta(d_t | d_{1:t-1}, I_0), \tag{5}$$

where $p_\theta(d_t | d_{1:t-1}, I_0)$ denotes the probability of the word $d_t$ given the word sequence $d_1, \ldots, d_{t-1}$ and image $I_0$. The CIDErBtw weighted XE loss is then

$$\mathcal{L}_{XE} = \sum_{i=1}^{N} w_i L_{XE}(c_i^0). \tag{6}$$

**Reweighting RL reward.** For RL, we reweight the CIDEr reward according to the CIDErBtw to focus more on distinctive captions, resulting in a new reward,

$$\tilde{R}(c^*) = \frac{1}{N} \sum_{i=1}^{N} w_i g_c(c^*, c_i^0), \tag{7}$$

where $g_c(c^*, c_i^0)$ is the CIDEr value between $c^*$ and ground-truth $c_i^0$.

**CIDErBtw reward.** Finally, when performing RL, our CIDErBtw can also be used as a part of the reward related to distinctiveness. We combine the CIDErBtw score with the prevous reward $\tilde{R}(c^*)$ and obtain the final RL reward $R(c^*)$ and RL loss $\mathcal{L}_{RL}$ as

$$R(c^*) = \tilde{R}(c^*) - \alpha_r CIDErBtw(c^*), \quad \mathcal{L}_{RL} = -\mathbb{E}_{c^* \sim p_\theta}[R(c^*)], \tag{8}$$

where $CIDErBtw(c^*)$ represents CIDErBtw score of the generated caption $c^*$ defined in (2), $\alpha_r$ is a hyperparameter controlling the relative contributions, and the greedy sampling is used as the RL policy $p_\theta$.

**CIDErBtw evaluation metric.** CIDEr measures the similarity between the generated caption $c^*$ and its ground-truth captions $C^0$, and has become an important evaluation metric in image captioning. We believe the distinctiveness should also be measured when evaluating the quality of generated captions.

Thus we propose to use CIDErBtw as a complementary evaluation metric for image captioning models. We hope that the caption $c^*$ generated by the model is closer to the semantics of target image $I_0$, while far from the semantics of other $K$ similar images $\{I_1, I_2, \ldots, I_K\}$. Therefore, the $c^*$ generated by a more distinctive image captioning model will have a lower CIDErBtw. Note that for evaluation, the similar image sets are computed using the validation or test split. Note that CIDEtBtw requires human annotations to evaluate the generated captions, which is similar to other captioning evaluations, e.g., CIDEr [34], BLEU [25], METEOR [7], and ROUGE [19]. Although VSE++ does not require human annotation for evaluation, it still needs ground-truth captions in the training phase, and the performance is highly related to the training data.

## 4 Experiments

In this section, we conduct extensive experiments to evaluate the effectiveness of CIDErBtw in generating distinctive captions. Note that our motivation is to generate distinctive captions as well as achieve high caption quality.

### 4.1 Implementation details

**Dataset.** We use the MSCOCO dataset [20] with Karpathy spliting [17]. The numbers of images are 113,287 for training, 5,000 for validation, and 5,000 for testing. There are five annotated captions for each image.

**Models.** For the image encoder, following Luo *et al.* [23], we use two types of features in the experiments, *i.e.*, the FC features and the spatial features. The FC features are extracted from Resnet-101 [14], and each image is encoded as a vector of dimension $2,048$. The spatial features are extracted from the output of a Faster-RCNN [28] following UpDown [2].

Our experiments are performed using four baseline models, *i.e.*, FC [29], Att2in [29], UpDown [2], and Transformer [33]. FC model only uses the FC features, Att2in and Transformer only use the spatial features, and UpDown uses both types of features. Each model is trained using four methods: 1) MLE with standard XE loss, denoted as "*model*"; 2) MLE with CIDErBtw-weighted XE loss in (6), denoted as "*model*+CIDErBtw"; 3) SCST [29], which trains with standard XE loss first, and then switches to RL with CIDEr reward, denoted as "*model*+SCST"; 4) SCST using weighted XE loss and weighted RL reward in (7), denoted as "*model*+SCST+CIDErBtw".

**Training details.** We set $\lambda_w$ as 1.5, $\alpha_w$ between 0.25 to 1.25 when reweighting the loss and the reward. $\alpha_r$ is set to 0.4 when using CIDErBtw reward, and 0 otherwise. We use Adam [18] to optimize the training parameters with an initial learning rate $5 \times 10^{-4}$ and a decay factor 0.8 every three epochs. During test time, we apply beam search with size five to generate captions.

**Metrics.** For evaluation we consider two groups of metrics. The first group includes language quality metrics CIDEr, BLEU3, BLEU4, METEOR, ROUGE-L, and SPICE for evaluating the accuracy and quality of generated captions. The second group assesses the distinctiveness of captions, and includes our CIDErBtw metric and retrieval metrics (*i.e.*, R@1, R@5, R@10). When calculating CIDErBtw, we collect $K = 5$ similar images for each target image, so the CIDErBtw score measures the similarity between the generated caption and 25 captions from the similar images set, with lower values indicating more distinctiveness. Similar images sets are generated using a pre-trained VSE++ [10] to perform the caption-to-image retrieval (see Section 3.1). For the retrieval metrics, we follow the protocol in [21,23,6]. Given a generated caption, images are retrieved in the joint semantic space of the pre-trained VSE++, with the goal to retrieve the original image. Recall at $K$ (R@$K$) is used to measure the retrieval performance, where a higher recall represents a better distinctiveness.

### 4.2   Experiment results

In this section, we present the experiment results to show the effectiveness of CIDErBtw training strategies at improving caption distinctiveness. Due to space constraints, the ablation study is presented in the supplemental.

**Effect of CIDErBtw strategies.** The main results are presented in the top and middle of Table 1. All baseline models obtain better performances when using CIDErBtw weighting in training process, for both MLE or SCST, which suggests that our method is widely applicable to many existing models. Specifically, our method both reduces the CIDErBtw score and improves other accuracy metrics, such as CIDEr. This shows that the generated captions become more similar to ground-truth captions, while more distinctive from other images' captions since redundancy is suppressed. Among the four baseline models, CIDErBtw reweighted loss and reward have the largest effect on Transformer [33]. Most likely the multi-head attention and larger receptive field of Transformer allow it better extract details and context from the image that is distinctive.

Next we apply all three of our CIDErBtw reward strategies together on Transformer+SCST, which is denoted as "+CIDErBtwReward" in Table 1. Compared to only using reweighted loss and reward (Transformer+SCST+CIDErBtw), adding the CIDErBtw reward in RL improves both the CIDErBtw and retrieval metrics significantly (i.e., improves distinctiveness), at the expense of a small decrease in accuracy (CIDEr).

Finally, we examine the disadvantage of SCST that directly optimizing CIDEr reward improves the fluency of captions but also leads to common and generic words. Consistent with [38,21], the baseline models trained with SCST obtain higher CIDEr but also perform worse in CIDErBtw and R@$K$, compared with models trained only with MLE. Optimizing the model with CIDErBtw weighted reward will relieve this problem, and the distinctness of captions will be promoted, while maintaining or even improving the overall quality of the captions.

| Method | CIDEr↑ | CIDErBtw↓ | BLEU3↑ | BLEU4↑ | METEOR↑ | ROUGE-L↑ | SPICE↑ | R@1↑ | R@5 ↑ | R@10↑ |
|---|---|---|---|---|---|---|---|---|---|---|
| FC [29] | 97.90 | 83.35 | 41.81 | 31.58 | 25.22 | 53.34 | 17.99 | 15.44 | 40.36 | 55.08 |
| FC+CIDErBtw (ours) | 98.82 | 83.22 | 42.03 | 31.79 | 25.46 | 53.48 | 18.29 | 16.24 | 41.54 | 56.64 |
| Att2in [29] | 110.04 | 83.19 | 46.36 | 35.75 | 26.79 | 56.18 | 19.91 | 17.44 | 43.88 | 58.02 |
| Att2in+CIDErBtw (ours) | 110.97 | 82.42 | 46.63 | 36.0 | 27.03 | 56.30 | 20.01 | 17.98 | 44.72 | 58.62 |
| UpDown [2] | 111.25 | 79.46 | 45.64 | 35.93 | 27.54 | 56.24 | 20.54 | 20.10 | 47.58 | 61.92 |
| UpDown+CIDErBtw (ours) | 112.77 | 78.34 | 46.35 | 36.10 | 27.69 | 56.36 | 20.68 | 20.92 | 49.72 | 63.98 |
| Transformer [33] | 110.13 | 80.98 | 44.80 | 34.46 | 26.98 | 55.30 | 20.18 | 21.52 | 49.88 | 64.70 |
| Transformer+CIDErBtw (ours) | 112.44 | **75.35** | 45.44 | 35.01 | 27.59 | 55.66 | 20.74 | 21.84 | 50.48 | 65.04 |
| FC+SCST [29] | 104.43 | 90.09 | 43.10 | 31.59 | 25.46 | 54.33 | 18.67 | 11.44 | 33.16 | 48.04 |
| FC+SCST+CIDErBtw (ours) | 104.76 | 89.41 | 43.25 | 31.72 | 25.60 | 54.35 | 18.58 | 11.74 | 33.62 | 48.32 |
| Att2in+SCST [29] | 117.96 | 87.40 | 47.22 | 35.31 | 27.17 | 56.92 | 20.57 | 16.00 | 41.55 | 56.66 |
| Att2in+SCST+CIDErBtw (ours) | 118.48 | 87.21 | 47.33 | 35.41 | 27.27 | 56.94 | 20.77 | 16.82 | 42.26 | 57.72 |
| UpDown+SCST [2] | 121.94 | 86.82 | 48.82 | 36.12 | 27.95 | 57.61 | 21.29 | 18.50 | 46.34 | 61.70 |
| UpDown+SCST+CIDErBtw (ours) | 123.02 | 86.42 | 48.98 | 36.39 | 28.12 | 57.78 | 21.44 | 19.68 | 47.30 | 62.78 |
| Transformer+SCST [33] | 125.13 | 86.68 | 50.26 | 38.04 | 27.96 | 58.60 | 22.30 | 23.38 | 54.34 | 68.44 |
| Transformer+SCST+CIDErBtw (ours) | **128.11** | 84.70 | **51.29** | **39.0** | **29.12** | **59.24** | 22.92 | 24.46 | 55.22 | 69.02 |
| +CIDErBtwReward (ours) | 127.78 | 82.74 | 50.97 | 38.52 | 29.09 | 58.82 | **22.96** | **26.46** | **57.98** | **71.28** |
| Stack-Cap [12] | 120.4 | 88.7 | 47.9 | 36.1 | 27.4 | 56.9 | 20.9 | 21.9 | 49.7 | 63.7 |
| DiscCap [23] | 120.1 | 89.2 | 48.5 | 36.1 | 27.7 | 57.8 | 21.4 | 21.6 | 50.3 | 65.4 |
| VisPara-Cap [21] | 86.9 | - | - | 27.1 | - | - | 21.1 | 26.3 | 57.2 | 70.8 |
| CL-Cap [6] | 114.2 | 81.3 | 46.0 | 35.3 | 27.1 | 55.9 | 19.7 | 24.1 | 52.5 | 67.5 |
| PSST [35] | 111.9 | - | - | 32.2 | 26.4 | 54.4 | 20.6 | 45.3† | 79.4† | 89.9† |

**Table 1:** Comparison of caption accuracy and distinctiveness on MSCOCO test split: (top) baseline models trained with MLE using standard or our weighted XE loss; (middle) models trained with SCST using standard or our weighted loss/reward; (bottom) SOTA methods for generating distinctive/discriminative captions. CIDEr, BLEU3/4, METEOR, ROUGE-L, and SPICE measure caption accuracy, while CIDErBtw and R@$K$ measure distinctiveness. ↑ or ↓ show whether higher or lower scores are better for each metric. CIDErBtw could not be computed for some models because the captions are not publicly available. Our self-retrieval results (R@$K$) and those of [12,23,21,6] use the pre-trained VSE++ model and the same protocol. † Note that [35] reports self-retrieval results using a different retrieval model/protocol – they use their own model for retrieval – which makes it not directly comparable.

**Reasons for improving CIDEr.** Results in Table 1 show that models trained with CIDErBtw obtain better performance for *both* distinctiveness metrics and accuracy metrics. Given that our training method puts more weight on distinct ground-truth captions, it is expected that we will obtain lower CIDErBtw and higher R@$K$ scores. However, the reason why our method also improves caption accuracy (CIDEr) is less obvious, especially for SCST, which *directly* optimizes CIDEr using RL. Note that CIDEr is based on the cosine similarity between TFIDF vectors, and thus low-frequency words (with higher IDF weights) will have higher impact on the CIDEr score. Since rare words are also distinct, their usage in a caption should increase CIDEr. If using distinct words can increase CIDEr, then why does RL with CIDEr reward not use distinct words? We speculate that RL gets stuck in a local minimum of models that only use frequent words because of two reasons: 1) equal weighting of an image's ground-truth captions encourages the model to predict the common words that match all captions; and 2) regularization encourages models to use smaller vocabularies – using less words means less non-zero weights in the network, and lower model

complexity. By reweighting the reward with CIDErBtw, more reward is obtained when using diverse words, which effectively moves the learning process out of this local minimum.

**Comparison with state-of-the-art.** We list the performance of state-of-the-art captioning models that focus on distinctiveness at the bottom of Table 1. Compared to these models, our model (Transformer+SCST+CIDErBtw, and +CIDErBtwReward) generally achieves superior results in both accuracy and distinctiveness — our model obtains both a high CIDEr score and low CIDErBtw score (or high retrieval score) at the same time. Specifically, Stack-Cap [12] and DiscCap [23] have lower accuracy (CIDEr 120) and less distinctiveness (CIDErBtw 89, R@1 22), compared to our model. VisPara-Cap [21] has high distinctiveness by using visual paraphrases, slightly worse than our model (+CIDErBtwReward), while the accuracy (CIDEr 86.9) is much lower than our model. CL-Cap [6] and PSST [35] directly optimize the retrieval loss, aiming to identify the input image among a set of randomly-chosen distractor images, which improves the distinctiveness. CL-Cap has similar distinctiveness as our method, obtaining worse $R@K$ than ours, but better CIDErBtw.[5] However, directly optimizing the training parameter with retrieval loss results in low-quality captions, lowering the accuracy (CIDEr 114.2 and 111.9) compared to our model.

### 4.3   User Study

To fairly evaluate the quality of generated sentences and verify the consistency between the metrics and human perspective, we conducted two user studies. Firstly, we performed a user study on image retrieval to assess distinctiveness, following the protocol in [23]. The task involves displaying the target image, a semantically similar image which is retrieved following the method in Section 3.1, and a generated caption describing the target image. The users are asked to choose the image that more closely matches the caption.

In the second experiment, we compare captions generated from a baseline model trained with and without CIDErBtw. In each trial, an image and two captions are displayed, and the user is asked to choose the better caption with respect to two criteria: distinctiveness and accuracy. In each experiment, we randomly sample 50 similar images pair from the test split. We perform the experiment on four captioning models: UpDown [2] and Transformer [33] trained by SCST with and without CIDErBtw (denoted as UD, UD+CIDErBtw, TF and TF+CIDErBtw). Twenty people participated in the user study, and we collected about $6,000$ responses in total. See the supplemental for more details.

The results for the image retrieval user study are shown in Table 2. Compared to the baseline model, our method increases the accuracy of image retrieval by $5.6\%$ and $14.4\%$. This user study is consistent with the automatic image retrieval results ($R@K$), and indicates that captions generated by our model are more
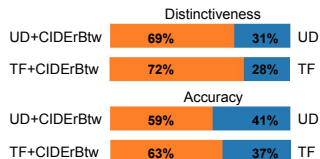
---

[5] We could not compare distinctiveness with PSST since their captions are not publicly available, and they use a different evaluation protocol for $R@K$.

distinctive in terms of both machine and human perception, than those of the baseline models.

The result for the distinctiveness/accuracy user study are shown in Figure 3. From human perspective, captions from our models are more distinctive than the baseline models (our captions are selected 69% and 72% of the time). The improvement of accuracy is not as much (our model selected 59% and 63% of the time), since the baseline models already generate captions that are accurate. Again this is consistent with the observations from the machine-based metrics (CIDErBtw and CIDEr).

| Method | image retrieval |
|---|---|
| UD | 68.7% ** |
| UD+CIDErBtw | **74.3%** ** |
| TF | 75.2% * |
| TF+CIDErBtw | **79.6%** * |



**Table 2:** User study on image retrieval to assess caption distinctiveness. Our models trained with CIDErBtw generated more distinctive captions, enabling the user to more accurately select the correct image, compared with the baselines (2-sample z-test on proportions, * p<0.05, ** p<0.01).

**Figure 3:** User study comparing captions generated from models trained with and without our CIDErBtw. Users selected our models trained with CIDErBtw more frequently when assessing accuracy and distinctiveness (Chi-Square test, p<0.001 for each pair).

### 4.4 Qualitative Results

We next show qualitative results for the baseline model Transformer+SCST, and our model Transformer+SCST+CIDErBtw in Figure 4. The baseline model generates captions that accurately describe the main object, but are quite generic and monotonous. Intuitively, in order to increase a caption's distinctiveness, the model should focus on more properties that would distinguish the image from others, such as color, numbers, or other objects/background in the image. Our method focuses on more of these aspects and generates accurate results. Our captions describe more properties of the main object, such as "*black suit*", "*red tie*" and "*a man and a child*". We also describe backgrounds that are distinctive, such as "*pictures on the wall*" and "*city street at night*".

In order to show the distinctiveness of our model, we present a similar images set with the same semantic meaning in Figure 5. The baseline model generates captions that follow generic templates, e.g. "*train on the track*" or "*at a train station*". Although the captions are correct, it is hard to tell the images apart according to the captions. Our model enriches the description by mentioning the colors, e.g. the "*green and yellow*" and "*yellow and black*" distinguishing the first two images, and the background environment, e.g. "*under a bridge*" and "*in a forest*". Furthermore, our model is more sensitive to the relative positions of objects, e.g. "*next to each other on the tracks*". However, a more descriptive caption may also lead to some errors. For instance, the train in the third image is not actually "*under a bridge*". More details are in the supplementary material.

| Baseline: | (86.7) A man in a suit and tie. | (59.0) A living room with a television and on the television | (66.8) A stop sign on the side of a road. | (90.2) A man standing on the beach with a surfboard. |
| Ours: | (49.0) A man in a black suit wearing a red tie. | (56.6) A living room with a television and pictures on the wall. | (58.3) A stop sign on the side of a city street at night. | (69.7) A man and a child standing on the beach with a surfboard. |

**Figure 4:** Example captions from the baseline model and our model. The distinctive words are highlighted. The number in parenthesis is the CIDErBtw score, with lower values meaning more distinctive.



| Baseline: | (122.3) A yellow train on the tracks at a train station. | (110.7) A yellow train on the tracks of a track. | (104.9) A yellow train on the tracks at a train station. | (157.8) A train is sitting on the tracks. | (92.8) Two trains on the tracks at a train station. |
| Ours: | (93.6) A green and yellow train is on the tracks at a train station. | (103.2) A yellow and black train is on the tracks. | (79.4) A yellow train is on the tracks under a bridge. | (141.7) A train is on the tracks in a forest. | (22.3) Two red trains parked next to each other on the tracks. |

**Figure 5:** Example captions for a set of similar images.

## 5  Conclusion

In this paper, we consider an important property, *distinctiveness* of image captions, and proposed a metric CIDErBtw to evaluate distinctiveness, which can be calculated quickly and easily implemented. We found that human annotations for each image vary in distinctiveness based on CIDErBtw. To improve the distinctiveness of generated captions, we developed a novel training strategy, where each human ground-truth annotation is assigned a weight based on its distinctiveness computed by CIDErBtw. Thus, during training the model pays more attention to the captions that are more distinctive. We also consider using CIDErBtw directly as part of the reward in RL. Extensive experiments were conducted, and we showed that our method is widely applicable to many captioning models. Experimental results demonstrate that our training strategy is able to improve both accuracy and distinctiveness, achieving state-of-the-art performance on CIDEr, CIDErBtw and retrieval metrics (R@$K$).

# References

1. Anderson, P., Fernando, B., Johnson, M., Gould, S.: SPICE: Semantic propositional image caption evaluation. In: ECCV (2016) 3, 4, 5
2. Anderson, P., He, X., Buehler, C., Teney, D., Johnson, M., Gould, S., Zhang, L.: Bottom-up and top-down attention for image captioning and visual question answering. In: CVPR. pp. 6077–6086 (2018) 3, 4, 9, 11, 12
3. Aneja, J., Deshpande, A., Schwing, A.: Convolutional image captioning. In: CVPR (2018) 1, 3
4. Chunseong Park, C., Kim, B., Kim, G.: Attend to you: Personalized image captioning with context sequence memory networks. In: CVPR. pp. 895–903 (2017) 1
5. Dai, B., Fidler, S., Urtasun, R., Lin, D.: Towards diverse and natural image descriptions via a Conditional GAN. In: ICCV (2017) 2, 4, 5
6. Dai, B., Lin, D.: Contrastive learning for image captioning. In: NeurIPS (2017) 2, 4, 10, 11, 12
7. Denkowski, M., Lavie, A.: Meteor universal: Language specific translation evaluation for any target language. In: EACL Workshop (2014) 5, 9
8. Deshpande, A., Aneja, J., Wang, L., Schwing, A.G., Forsyth, D.: Fast, diverse and accurate image captioning guided by part-of-speech. In: CVPR (2019) 5
9. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: BERT: Pre-training of deep bidirectional transformers for language understanding. NAACL (2018) 3
10. Faghri, F., Fleet, D.J., Kiros, J.R., Fidler, S.: VSE++: Improving visual-semantic embeddings with hard negatives. BMVC (2018) 5, 6, 10
11. Farhadi, A., Hejrati, M., Sadeghi, M.A., Young, P., Rashtchian, C., Hockenmaier, J., Forsyth, D.: Every picture tells a story: Generating sentences from images. In: ECCV (2010) 1, 3
12. Gu, J., Cai, J., Wang, G., Chen, T.: Stack-captioning: Coarse-to-fine learning for image captioning. In: AAAI (2018) 11, 12
13. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. In: CVPR. pp. 2961–2969 (2017) 3
14. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR. pp. 770–778 (2016) 3, 9
15. Hochreiter, S., Schmidhuber, J.: Long short-term memory. Neural computation 9(8), 1735–1780 (1997) 3
16. Huang, L., Wang, W., Chen, J., Wei, X.Y.: Attention on attention for image captioning. In: ICCV. pp. 4634–4643 (2019) 4
17. Karpathy, A., Fei-Fei, L.: Deep visual-semantic alignments for generating image descriptions. In: CVPR. pp. 3128–3137 (2015) 9
18. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: ICLR (2015) 9
19. Lin, C.Y.: Rouge: A package for automatic evaluation of summaries. In: ACL Workshop (2004) 5, 9
20. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft COCO: Common objects in context. In: ECCV (2014) 9
21. Liu, L., Tang, J., Wan, X., Guo, Z.: Generating diverse and descriptive image captions using visual paraphrases. In: CVPR. pp. 4240–4249 (2019) 1, 2, 4, 10, 11, 12
22. Liu, X., Li, H., Shao, J., Chen, D., Wang, X.: Show, tell and discriminate: Image captioning by self-retrieval with partially labeled data. In: ECCV. pp. 338–354 (2018) 2, 4

23. Luo, R., Price, B., Cohen, S., Shakhnarovich, G.: Discriminability objective for training descriptive captions. In: CVPR. pp. 6964–6974 (2018) 2, 4, 9, 10, 11, 12
24. Mao, J., Xu, W., Yang, Y., Wang, J., Huang, Z., Yuille, A.: Deep captioning with multimodal recurrent neural networks (m-rnn). In: ICLR (2015) 1, 3
25. Papineni, K., Roukos, S., Ward, T., Zhu, W.J.: BLEU: a method for automatic evaluation of machine translation. In: ACL (2002) 3, 4, 5, 9
26. Park, C.C., Kim, B., Kim, G.: Towards Personalized Image Captioning via Multimodal Memory Networks. In: IEEE TPAMI (2018) 1
27. Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I.: Language models are unsupervised multitask learners. OpenAI Blog **1**(8),  9 (2019) 3
28. Ren, S., He, K., Girshick, R., Sun, J.: Faster R-CNN: Towards real-time object detection with region proposal networks. In: NeurIPS. pp. 91–99 (2015) 3, 9
29. Rennie, S.J., Marcheret, E., Mroueh, Y., Ross, J., Goel, V.: Self-critical sequence training for image captioning. In: CVPR. pp. 7008–7024 (2017) 3, 4, 8, 9, 11
30. Shetty, R., Rohrbach, M., Hendricks, L.A.: Speaking the same language: Matching machine to human captions by adversarial training. In: ICCV (2017) 2, 4, 5
31. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. In: ICLR (2015) 3
32. Van Miltenburg, E., Elliott, D., Vossen, P.: Measuring the diversity of automatic image descriptions. In: COLING. pp. 1730–1741 (2018) 5
33. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. In: NeurIPS. pp. 5998–6008 (2017) 2, 3, 9, 10, 11, 12
34. Vedantam, R., Lawrence Zitnick, C., Parikh, D.: CIDEr: Consensus-based image description evaluation. In: CVPR. pp. 4566–4575 (2015) 3, 4, 5, 9
35. Vered, G., Oren, G., Atzmon, Y., Chechik, G.: Joint optimization for cooperative image captioning. In: CVPR. pp. 8898–8907 (2019) 2, 4, 11, 12
36. Vinyals, O., Toshev, A., Bengio, S., Erhan, D.: Show and tell: A neural image caption generator. In: CVPR (2015) 1, 3
37. Wang, Q., Chan, A.B.: CNN+CNN: Convolutional decoders for image captioning. CVPR Workshop (2018) 1, 3
38. Wang, Q., Chan, A.B.: Describing like humans: on diversity in image captioning. In: CVPR (2019) 2, 4, 5, 10
39. Wang, Q., Chan, A.B.: Towards diverse and accurate image captions via reinforcing determinantal point process. arXiv (2019) 5
40. Xu, J., Ren, X., Lin, J., Sun, X.: Diversity-promoting GAN: A cross-entropy based generative adversarial network for diversified text generation. In: EMNLP. pp. 3940–3949 (2018) 5
41. Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A., Salakhudinov, R., Zemel, R., Bengio, Y.: Show, attend and tell: Neural image caption generation with visual attention. In: ICML (2015) 1, 3, 4
42. Yao, T., Pan, Y., Li, Y., Mei, T.: Hierarchy parsing for image captioning. In: ICCV. pp. 2621–2629 (2019) 3