

Direct Metric Optimization for Image Captioning through Reward-Weighted Augmented Data Utilization

Anonymous ACL submission

Abstract

While image captioning is an essential field of vision language models (VLM), a lack of continuity between the learning objective and final performance metrics of VLMs complicates their training and optimization. Reinforcement learning (RL) can directly optimize such metrics, but it is accompanied by a significant computational cost, making it difficult to apply to recent large-scale VLMs. In this paper, we propose *Direct Metric Optimization* (DMO), which is a lightweight final-metric-optimizing training method. We replace the computationally expensive exploration process in RL with an offline, diverse text data augmentation and show that self-supervised training on reward-weighted augmented data leads to direct and stable metric optimization. Our experiments demonstrate that DMO achieves performance comparable to those of the state-of-the-art RL method while saving hundreds of times more model forwarding iterations and greater amounts of computation time. This suggests that DMO constitutes a promising alternative for metric optimization in the era of large-scale VLMs.

1 Introduction

With the advent of CLIP (Radford et al., 2021), the boundaries between vision and language modalities in machine learning have been dissolved, leading to rapid advancements in research involving these areas. Furthermore, the rise of large language models (LLM) has led to the emergence of large-scale vision language models (VLM), extending their influence to practical applications. For example, models such as ChatGPT (Achiam et al., 2023) and Gemini (Team et al., 2023) generate detailed natural language descriptions from visual information. With the increasing prevalence of VLMs, methods for customizing and fine-tuning these models for specific domains or individuals are attracting significant interest and attention (Sun et al., 2023; Zhao et al., 2023).

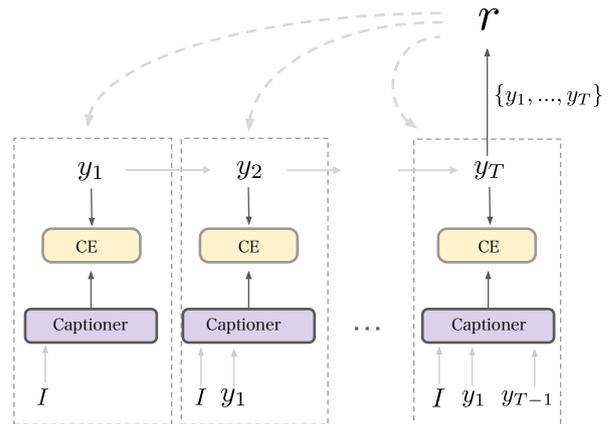


Figure 1: An overview of our *Direct Metric Optimization* (DMO). Image and tokens are denoted as I and y_i and CE stands for cross entropy function. Precomputed rewards r assign different weights to each sample in a textually augmented dataset, effectively enhancing the targeted performance metrics.

Recent standard captioning models adopt self-supervised learning for training purposes (Wang et al., 2022; Yu et al., 2022; Alayrac et al., 2022; Li et al., 2023). This method treats the ground truth captions both as inputs and labels, and the model predicts only the next token from the given image and preceding tokens. Specifically, recent transformer-based encoder-decoder models can conduct the next token prediction of each step in parallel, significantly enhancing their computational efficiency. However, this approach is subject to certain limitations. Typically, the performance of image captioning is evaluated using metrics such as BLEU (Papineni et al., 2002) or CIDEr (Vedantam et al., 2015); however, self-supervised learning of language modeling does not necessarily optimize those metrics. To target final metrics directly, reinforcement learning (RL) methods have been employed (Ranzato et al., 2015; Zhang et al., 2017b; Rennie et al., 2017; Gao et al., 2019). Reinforce-

063 ment learning is a powerful method capable of op- 113
 064 timizing even non-differentiable metrics; however, 114
 065 it has certain drawbacks, such as learning instabil- 115
 066 ity and significant time and computational costs. 116
 067 With the growing trend of using large pre-trained
 068 models, those challenges have become increasingly
 069 serious. Conducting RL with models containing
 070 billions of parameters demands extensive computa-
 071 tional time and resources, making the application
 072 of RL methods impractical.

073 To bypass the prohibitive computational cost of 118
 074 RL, we propose to replace the expensive explo- 119
 075 ration process in RL with diverse text data aug- 120
 076 mentation and reduce RL to simple importance- 121
 077 weighted self-supervised learning. The approach 122
 078 that utilizes previously collected data for RL is 123
 079 known as offline-RL (Levine et al., 2020). Partic- 124
 080 ularly in our approach, datasets are augmented by 125
 081 various methods and the augmentation diversity 126
 082 brings a variety of samples of different rewards, 127
 083 enabling the efficient estimation of the optimal cap-
 084 tion for the image. We call this metric-optimizing
 085 self-supervised training *Direct Metric Optimization*
 086 (DMO). Our experiments demonstrate that DMO
 087 achieves performance on par with state-of-the-art
 088 (SOTA) RL methods in standard image captioning
 089 metrics while retaining lightweight computational
 090 efficiency and learning stability. This highlights
 091 DMO’s significant practical advantages in metric
 092 optimization, especially considering the increasing
 093 need to tune and customize large-scale VLMs.

094 2 Preliminaries

095 2.1 Self-Supervised Learning for Image 096 Captioning

097 The standard approach for image captioning in re-
 098 cent years has been to employ an encoder-decoder
 099 model, where the encoder maps the image into the
 100 latent space and extracts features from the image,
 101 and the text decoder autoregressively generates to-
 102 kens for the next step from extracted features and
 103 previously generated tokens. In the self-supervised
 104 learning of language models (LM), the model is of-
 105 ten trained by teacher-forcing (Williams and Zipser,
 106 1989). This method increases the likelihood of
 107 ground-truth sentences by aligning the model’s
 108 conditional distribution $p_\theta(y_i|I, y_{<i})$ with the la-
 109 bel distribution $q(y_i)$ using cross-entropy (CE) for
 110 each step $i \in \{1, \dots, T\}$. Here, y_i represents
 111 a text token at step i and I is the given image.
 112 The label distribution $q(y_i)$ is a one-hot vector or

label-smoothed vector (Szegedy et al., 2016) from
 ground truth label y_i . The objective function of self-
 supervised learning for language modeling \mathcal{L}_{LM} is
 expressed as follows:

$$\mathcal{L}_{\text{LM}} = \sum_i^T \text{CE}(q(y_i), p_\theta(y_i|I, y_{<i})). \quad (1)$$

Especially in the recent Transformer-based archi-
 tecture, the predictions of the next tokens at each
 time step can be performed in parallel (Vaswani
 et al., 2017). Thus this training method is extremely
 time and computationally efficient because it does
 not require recursive operations, as is the case with
 conventional RNN-based methods (Vinyals et al.,
 2015; Xu et al., 2015).

2.2 Reinforcement Learning for Image Captioning

Typically, the performance of captioning models
 is assessed using metrics such as BLEU or CIDEr.
 However, self-supervised learning does not neces-
 sarily optimize these metrics. Furthermore, these
 evaluation metrics are often non-differentiable,
 making it impossible to apply the gradient descent
 directly. One approach for directly optimizing
 those non-differentiable metrics is to employ RL.
 Recent studies of image captioning have applied
 the various RL algorithms including REINFORCE
 and Actor-Critic (Ranzato et al., 2015; Zhang et al.,
 2017b; Liu et al., 2017; Rennie et al., 2017; Zhang
 et al., 2021), to captioning tasks by regarding the
 captioning models as agents and the final evalu-
 ation metrics (such as CIDEr) as rewards. The
 objective function of the captioning model in the
 RL framework is expressed as follows:

$$\mathcal{L}_{\text{RL}} = -\mathbb{E}_{w_i \sim p_\theta} [r(\{w_1, \dots, w_T\})], \quad (2)$$

where w_i is the token sampled from the model’s
 distribution p_θ at the time step i and T is the total
 length of the token sequence. The partial deriva-
 tives of \mathcal{L}_{RL} can be determined using the REIN-
 FORCE algorithm (Williams and Zipser, 1989).
 The calculation of the expected values is approx-
 imated by Monte Carlo sampling within a mini-
 batch as follows:

$$\begin{aligned} \nabla_\theta \mathcal{L}_{\text{RL}} &\approx -r(\{w_1, \dots, w_T\}) \cdot \\ &\nabla_\theta \log p_\theta(\{w_1, \dots, w_T\}|I). \end{aligned} \quad (3)$$

In RL, the reward metric need not be confined to automatic evaluation scores. It can accept any user-defined score including ambiguous evaluations from humans (Christiano et al., 2017; Ouyang et al., 2022), thus demonstrating significant versatility and adaptability to various tasks.

Notwithstanding these advantages, the RL methodology presents significant challenges. First, for the method to be effective, extensive exploration is needed, which makes RL time-consuming. Second, the learning process tends to be unstable, especially during the early stage of training, when the model poorly samples the high-rewarded sequences. (Ranzato et al., 2015; Rennie et al., 2017). In addition, the sampling process in RL is computationally inefficient because the gradient computation requires the last token w_T (see Equation 3), but token generation is done in a left-to-right manner, undermining the computational parallelism of Transformer architecture.

2.3 Text Data Augmentation (TDA)

Data augmentation is a traditional yet effective method that is used to enhance a language model’s performance (Li et al., 2023; Fan et al., 2023; Yang et al., 2023). Previous studies have shown that text data augmentation (TDA) strategies can be broadly categorized into three types: paraphrasing, noising and sampling (Li et al., 2022). Paraphrasing is a method that generates data that convey very similar information as the original data with restrained changes. Noising adds noise to datasets to improve the robustness of the model. Sampling produces new novel data from the data distribution master. The primary objective of these strategies is to introduce diversity into the dataset. This is particularly crucial in scenarios with limited datasets, where models are prone to overfitting. Augmenting the dataset and smoothing the distribution can effectively prevent this overfitting.

In the fields of image and audio processing, data augmentation has traditionally achieved significant success (DeVries and Taylor, 2017; Zhang et al., 2017a). In contrast, it has not been explored as extensively in the field of natural language processing. This disparity can be attributed to the challenges inherent in text augmentation. Unlike images and audio, which comprise continuous data, tokenized text data is discrete and even minor alterations can lead to significant semantic shifts. Implementing superficial changes while controlling these seman-

tic variations is not straightforward, and universally effective methods for achieving this are yet to be established (Feng et al., 2021).

3 Proposed Framework

3.1 Overview of Method

In RL, the sequences are sampled from the model’s distribution p_θ , but this raises problems of its high computational cost and instability in the early stages of training. We propose the following perspective shift: What if we were to consider the sequences drawn from the given dataset as the sequences sampled from the model itself? This approach allows us to obtain the gradients of RL objective function with ground truth data in a self-supervised manner. Furthermore, since sequences are pre-sampled, the bottleneck in RL, specifically the recursive generation process, is resolved. This significantly enhances computational efficiency.

The approach that utilizes previously collected data for RL is known as "offline-RL" (Levine et al., 2020), and is commonly used to bypass the computationally expensive exploration in RL (Chen et al., 2021; Jang et al., 2022; Shi et al., 2023; Baheti et al., 2023). Applying the offline-RL to image captioning, however, is not straightforward. First, because of the nature of major captioning metrics (such as BLEU, METEOR, and ROUGE-L) that measure the overlap of n-grams, words or subsequences with a set of ground-truth captions, ground truth captions always receive rewards of 1 in the offline-RL framework. Because the rewards are indicators of the quality of samples, receiving a constant value of rewards gives no clue about how good each caption is, and consequently, there is no advantage compared with standard self-supervised training. The second problem is data suboptimality, a common challenge in offline-RL (Levine et al., 2020). The reliance on limited static data restricts exposure to high-reward samples, thereby capping the model’s performance improvements. We address those obstacles by introducing diverse text data augmentation (TDA). With TDA, we expose models to a variety of expressions with different rewards outside the original dataset, providing a greater number of clues about the optimal caption for the images. Furthermore, substituting TDA for exploration improves the stability of the learning process in the early stage. This is because, unlike RL, TDA can consistently provide reasonable quality samples and training does not rely on

the model’s capability of sampling high-reward sequences. This metric-optimizing self-supervised training on textually augmented datasets, which we call *Direct Metric Optimization*, offers the following two significant advantages.

1. It allows direct optimization of metrics in a self-supervised manner, significantly enhancing computational efficiency.
2. Training is stable even at an early stage because it does not rely on the model’s capability of generating captions of high rewards.

3.2 Direct Metric Optimization

In the proposed DMO method, sequences are sampled from textually augmented dataset D_{aug} . Because the dataset is known, scores for each ground truth sample can be calculated in advance. Let the score function (for example, the BLEU and CIDEr scorer) be the reward function $r(\cdot)$. Once ground truth data $d = \{y_1, \dots, y_T\}$ from dataset D_{aug} is sampled, the gradient of our DMO objective \mathcal{L}_{DMO} is defined as follows:

$$\nabla_{\theta} \mathcal{L}_{\text{DMO}}(d) = -r(\{y_1, \dots, y_T\}) \cdot \nabla_{\theta} \log p_{\theta}(\{y_1, \dots, y_T\} | I) \quad (4)$$

$$= -r(d) \nabla_{\theta} \sum_i^T \log p_{\theta}(y_i | I, y_{<i}) \quad (5)$$

$$= r(d) \nabla_{\theta} \sum_i^T \text{CE}(q(y_i), p_{\theta}(y_i | I, y_{<i})) \quad (6)$$

$$= r(d) \nabla_{\theta} \mathcal{L}_{\text{LM}}(d), \quad (7)$$

where $q(y_i)$ is a one-hot vector from label y_i . This can be interpreted as a reward-weighted gradient of self-supervised learning loss. This approach eliminates the bottleneck inherent in RL, specifically the recursive generation process, through the utilization of pre-sampled sequences. Consequently, it enables the model to leverage the parallel computational capabilities of the Transformer architecture, resulting in a substantial enhancement of computational efficiency.

This reward-weighted self-supervised training on augmented datasets is related to noise/similarity-aware supervised training that adaptively assigns different weights to each sample (Atliha and Šešok, 2020; Yang et al., 2023; Kang et al., 2023). However, there are notable differences. First, while those noise-aware methods often focus on large-scale pre-training from noisy datasets and mitigate

the effect of noisy samples, our approach features the finetuning stage with relatively small and clean datasets, and deliberately augments datasets to introduce the diversity of samples. Second, while similarity-aware methods often utilize CLIP/BERT scores or custom weights (Ding et al., 2019; Atliha and Šešok, 2020; Yang et al., 2023), we directly employ target metrics for sample weighting. While the CLIP/BERT score is useful for denoising or filtering, training with these measures does not directly lead to the optimization of the final metrics. With these perspectives, our method enables more effective optimization of the target metrics.

4 Experiment Implementations

This section describes the detailed experiment implementations for the evaluation of our DMO training.

4.1 Datasets and Captioning Models

We validate our method with the MS-COCO dataset (Lin et al., 2014) and Flickr8K (Hodosh et al.), which are commonly used in image captioning research. Both datasets have 5 captions per image. For the evaluation, the datasets are split into training, validation, and testing sets according to the Karpathy method (Karpathy and Fei-Fei, 2015) so that the numbers of images in the training, validation and test datasets become 6091/1000/1000 for Flickr8k and 113287/5000/5000 for MS-COCO. For captioning models, we employ GIT-base/large (Wang et al., 2022) and BLIP2-2.7b (Li et al., 2023), which have different sizes of parameters and architectures. GIT has a simplified architecture of one image encoder and one text decoder and the base model has 178M parameters while the large model has 390M parameters for each. BLIP2-2.7B has 2.7B parameters and it employs large pre-trained frozen models for its vision encoder and text decoder. Both models are pre-trained on datasets that include COCO (Lin et al., 2014), Visual Genome (Krishna et al., 2017), CC3M (Sharma et al., 2018), and other large image-text pair datasets. In our experiments, all models are finetuned for 3 epochs with learning rate 1.0×10^{-5} and batch size 960, using a fixed single random seed. Further details are explained in Appendix A.1.

4.2 Text Data Augmentation Strategy

Based on Section 2.3, we adopt the following three augmentation methods accordingly. From each of the following three methods, two augmented captions are randomly sampled for each image and added to the original training dataset.

- **Back-translation:** The En-Fr translation model from MarianNMT (Junczys-Dowmunt et al., 2018) is adopted. Back-translation is applied to each ground truth caption and the same number of back-translated captions as original captions are created.
- **Pre-trained VLM Sampling:** Using the COCO-pre-trained BLIP2-6.7B model, five captions are generated from each image in the training dataset with temperature 1.0.
- **Paraphrasing by LLM:** We employ Llama2-7b-chat (Touvron et al., 2023) to paraphrase captions. The detailed prompt text is presented in Appendix B.

We do not explicitly adopt a noising strategy, as sufficient semantic noise is introduced by each TDA method.

4.3 Metrics and Rewards

We evaluate the performance of models by CIDEr, BLEU-4, METEOR (Banerjee and Lavie, 2005), ROUGE-L (Lin, 2004) and SPICE (Anderson et al., 2016). Since our method requires each sample to be scored by a reward function, we directly use the metrics above as the rewards in the training. As the scoring of each metric demands a set of ground truth captions as references, we employ the original training dataset as the reference dataset.

5 Results

We evaluate our proposed method in terms of metric optimization performance, learning stability and computational efficiency. We further investigate how reward-weighting architecture facilitates robust metric optimization by comparing DMO with standard LM training under noisy data and limited data settings.

5.1 Evaluating Metric Optimization Performance

5.1.1 Does DMO Enhance Final Metrics?

First, we examine whether our method effectively improves the targeted metrics for image captioning. We use a textually augmented Flickr8k dataset

Training Metric	Evaluation Metrics			
	CIDEr	B4	MET.	ROU.
CIDEr	97.0	33.3	27.7	57.5
BLEU-4	96.8	33.5	27.7	57.5
METEOR	96.1	32.8	27.5	57.0
ROUGE-L	96.0	33.0	27.6	57.5
standard LM	95.1	33.4	27.0	57.1

Table 1: Performance evaluation of GIT-base model optimized for each metric by DMO on the textually augmented Flickr8k dataset. 'B4', 'MET.' and 'ROU.' refer to BLEU-4, METEOR and ROUGE-L, respectively.

(TDA-Flickr8k) and apply DMO to the GIT-base so that each CIDEr, BLEU-4, METEOR, and ROUGE-L is optimized respectively. We then evaluate whether DMO improves these metrics compared to training with the standard Language Model (LM) loss. The result is presented in table 1. We find that when optimized for each metric, there is an improvement in each metric compared to training with the standard LM loss. This result implies that our method can effectively enhance the target metrics. Interestingly, optimizing for CIDEr or BLEU-4 leads to improved scores in other metrics as well. This can be attributed to the similarities in the way each evaluation metric is measured. In the following experiments, we use CIDEr as the target metric because CIDEr-optimizing DMO leads to general improvements in scores across other metrics.

5.1.2 Comparison of DMO and LM Training

We compare the performance of DMO training with LM training for different models and dataset settings. We use three models, GIT-base, GIT-large, and BLIP2-2.7b and two datasets, the Flickr8k and the COCO dataset. We apply CIDEr-optimizing DMO to each model and compare CIDEr, BLEU-4, METEOR, ROUGE-L, and SPICE scores with the models trained by standard LM loss. The results are presented in Table 2. We observe that DMO results in significant performance improvements in almost all models and datasets compared with models trained with LM loss without TDA. On the Flickr8k dataset, all models exhibit score improvements across all metrics with DMO. A similar trend is observed when fine-tuning GIT-base on the COCO dataset. This suggests that DMO consistently enhances the scores beyond standard LM training across various models and datasets.

In the analysis comparing with LM training on

Captioning Model	Optimization Method	Flickr8k					MS-COCO				
		CIDEr	B4	MET.	ROU.	SPICE	CIDEr	B4	MET.	ROU.	SPICE
GIT-base	LM	95.1	33.4	27.0	57.1	21.4	135.6	41.0	30.5	60.4	23.6
GIT-base	LM w/ TDA	93.6	33.2	26.9	57.0	21.0	132.1	39.4	29.8	59.9	23.5
GIT-base	DMO	99.6	35.4	27.9	58.1	22.3	137.4	41.5	30.5	60.9	24.0
GIT-large	LM	96.3	33.3	26.9	56.9	21.2	140.9	42.5	31.3	61.3	24.3
GIT-large	LM w/ TDA	101.1	34.9	27.9	58.0	22.2	134.7	39.8	30.3	60.4	23.9
GIT-large	DMO	110.7	37.6	29.1	60.2	23.2	140.6	42.0	31.1	61.5	24.3
BLIP2-2.7b	LM	101.3	33.8	28.6	58.5	23.4	132.2	39.1	29.6	59.5	23.3
BLIP2-2.7b	LM w/ TDA	100.2	32.7	28.3	58.2	22.8	132.9	38.7	29.9	59.7	23.6
BLIP2-2.7b	DMO	103.7	33.8	28.7	58.4	23.7	138.3	41.1	30.4	60.7	24.0

Table 2: Evaluation of three models trained by standard LM training (with and without TDA) and CIDEr-optimizing DMO on Flickr8k and COCO datasets. The performance metrics include CIDEr, BLEU-4 (B4), METEOR (MET.), ROUGE-L (ROU.) and SPICE.

TDA datasets, we observe that LM training on TDA datasets causes a decline in the performance in certain scenarios, such as training GIT-base/BLIP2 on Flickr8k and GIT-base/large on the COCO. This implies that the TDA datasets possess excessive noise and this noise leads to a deterioration in the performance of the models trained with standard LM loss. In contrast, DMO training, which utilizes the TDA dataset, exhibits rather enhanced performance. BLIP2 trained with DMO on Flickr8k and GIT-base DMO-trained on the COCO show improved scores for almost all metrics while those trained with LM loss show worse performance by introducing TDA. These findings indicate that our DMO can effectively leverage even noisy datasets that would deteriorate the performance of regular LM training. Moreover, when GIT-large is trained on TDA-COCO, a reduction in performance is observed for both LM and DMO training. However, the decline in performance is significantly different: 6.1 points for LM training compared with 0.3 points for DMO training, highlighting DMO’s robustness under noisy dataset conditions.

5.1.3 Does TDA-Diversity Matter?

We hypothesize that diversifying the augmentation techniques serves as a replacement for exploration, enhancing the performance of DMO. To validate this hypothesis, we conduct an ablation study and evaluate the performance of models trained with DMO on datasets augmented by a single method and on datasets augmented by multiple methods, respectively. We denote the datasets augmented solely by the back-translation, pre-trained BLIP2 sampling, and Llama2 paraphrasing as D_{bktrs} , D_{blip2} , and D_{llama} respectively. For a fair comparison, each dataset is adjusted to

Dataset Setting	Evaluation Metrics				
	CIDEr	B4	MET.	ROU.	SPICE
D_{bktrs}	93.6	32.8	27.3	56.9	22.1
D_{blip2}	98.2	33.0	28.0	57.9	22.1
D_{llama}	99.5	34.3	28.0	58.0	21.9
D_{all}	99.6	35.4	27.9	58.1	22.3
baseline	95.1	33.4	27.0	57.1	21.4

Table 3: Scores of GIT-base model trained with CIDEr-optimizing DMO on each dataset setting. Baseline is the score of LM training without any TDA. B4: BLEU-4, MET: METEOR, ROU: ROUGE-L.

have approximately the same number of image-caption pairs. For datasets D_{bktrs} , D_{blip2} , D_{llama} , we increase 5 captions per image by augmentation. Note that the dataset D_{all} is constructed by sampling two augmented captions from each augmentation method and adding them to the original dataset. We use the Flickr8k dataset and train GIT-base by CIDEr-optimizing DMO. The results are presented in Table 3. While all data augmentation methods except for back-translation improve performance over the baseline, which is the score of LM training without any TDA, the highest performance is achieved with the dataset that combines all augmentation methods, suggesting that exposing the model to a variety of expressions from diverse augmentation techniques yields the most significant performance improvement. Further analysis of the advantages of combining multiple TDA methods is discussed in Appendix C.

5.2 Stability and Computational Efficiency

DMO replaces the exploration with TDA and resolves the learning instability and computational bottleneck of RL. We examine DMO’s stability and

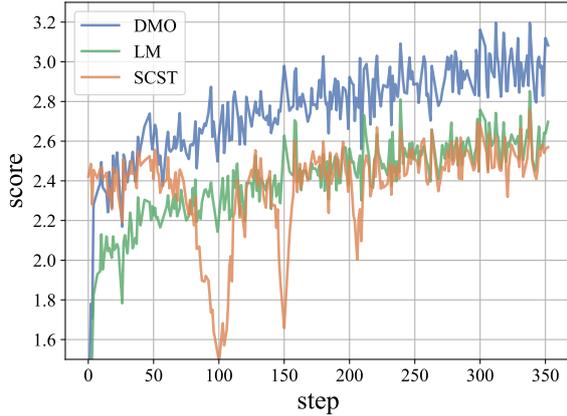


Figure 2: Transition of CIDEr scores for the GIT-base model trained using three different methods: DMO, standard LM training, and SCST. The scores reflect the CIDEr values of sequences greedily generated by each model for images in the mini-batches of the Flickr8k training dataset.

490 efficiency by comparing them with those of RL and
 491 standard LM training. We use the GIT-base model
 492 trained on the Flickr8k dataset. As an RL method,
 493 we employ SCST (Rennie et al., 2017), which is
 494 one of the most prominent reinforcement learning
 495 methods for image captioning, utilized in training
 496 many SOTA models (Wang et al., 2022; Xu et al.,
 497 2023). We assess the stability of training by con-
 498 ducting experiments with five seeds and calculating
 499 the average and variance of the final scores at the
 500 3rd and 20th epoch. To equalize the number of
 501 parameter updates, the number of explorations in
 502 RL is set to equal the number of image-text pairs
 503 in the training dataset. SCST requires intricate hyper-
 504 parameter tuning and we present the detailed con-
 505 figuration for SCST training in Appendix A.2. The
 506 result is presented in Table 4. With 3-epoch train-
 507 ing, our DMO produces the highest score. SCST is
 508 still unstable at 3 epochs, with a standard deviation
 509 (SD) of 4.1, which is significantly larger than LM’s
 510 SD of 0.74 or DMO’s SD of 1.26. Figure 2 shows
 511 the CIDEr score progression on training data up to
 512 5 epochs for DMO, LM, and SCST, respectively.
 513 While DMO and LM show steady score improve-
 514 ments, SCST exhibits large fluctuations. This insta-
 515 bility in SCST is attributed to its poor capability of
 516 sampling high-reward sequences in early training.
 517 On the other hand, DMO utilizes samples from
 518 TDA throughout the entire training, which makes
 519 training more stable. By the end of the 20th epoch,
 520 SCST achieves the highest score, owing to SCST’s

Method	3 epochs	20 epochs
LM	94.4 ± 0.74	94.2 ± 1.65
SCST	93.6 ± 4.10	99.8 ± 1.50
DMO	97.7 ± 1.26	98.0 ± 1.36

Table 4: The average and standard deviation of CIDEr scores when the model was trained for 3 epochs and 20 epochs with five different seeds by each method.

Optimization Method	Forwarding iterations	Execution time
LM	1.68×10^3	13 sec
SCST	1.50×10^5	9971 sec
DMO	1.68×10^3	14 sec

Table 5: The number of model forwarding iterations and execution time for the GIT-base to complete 3 epochs on Flickr8k. The time is measured during the loss computation and the number of iterations is measured by counting the number of batch model forwarding.

521 ability to continually explore and obtain new sam-
 522 ples. DMO displays minimal score improvements
 523 from epoch 3 to epoch 20, suggesting that training
 524 for only 3 epochs may suffice for model optimiza-
 525 tion in DMO training while SCST requires at least
 526 20 epochs. This indicates that DMO optimizes the
 527 performance of the model more rapidly than SCST.

528 Additionally, we measure the number of model
 529 forwarding iterations and the training time required
 530 for each method to complete 3 epochs. To elimi-
 531 nate the impact of differences in implementation
 532 and hardware differences on timing, we specifically
 533 measure the duration between feeding the data to
 534 the model and obtaining the loss values. The re-
 535 sult is presented in Table 5. With respect to com-
 536 putational efficiency, we find that LM and DMO
 537 have the same number of batch forwardings, while
 538 SCST requires approximately 100 times more (the
 539 rationale behind the results is explained in Ap-
 540 pendix A.1). In terms of execution time, while
 541 DMO training takes approximately the same du-
 542 ration as LM training, SCST requires about 1000
 543 times longer. The slight increase in time for DMO
 544 training compared to LM training is due to the ne-
 545 cessity of reward-weighting operations. On the
 546 other hand, SCST requires the recursive genera-
 547 tion process, resulting in a number of forwardings
 548 approximately 100 times greater and a duration ap-
 549 proximately 1000 times longer, compared to LM

Noise ratio	LM		DMO	
0%	95.1	(−0%)	98.6	(−0%)
20%	82.1	(−14%)	97.9	(−1%)
40%	69.6	(−27%)	98.0	(−1%)
60%	58.6	(−38%)	90.9	(−8%)
80%	5.1	(−95%)	81.9	(−17%)

Table 6: CIDEr scores of GIT-base model trained with noisy dataset. Noise ratio is the ratio of original ground-truth captions replaced by irrelevant random captions. The numbers in parentheses represent the percentage decrease from the score at the 0% noise ratio.

Dataset Size	LM		DMO	
100%	95.1	(−0%)	98.6	(−0%)
80%	92.1	(−3%)	97.4	(−1%)
60%	93.8	(−1%)	97.6	(−1%)
40%	90.4	(−5%)	95.4	(−3%)
20%	86.5	(−9%)	95.0	(−4%)

Table 7: CIDEr scores of GIT-base model trained with the limited number of data. Dataset size is the volume of the available training data. The numbers in parentheses represent the percentage decrease from the score at the 100% dataset volume.

and DMO training. These results emphasize the DMO’s substantially greater computational efficiency compared to that of SCST.

5.3 Noise Robustness and Data Efficiency

We experimentally show that DMO can train models robustly even in data-noisy or low-resource settings, by effectively leveraging the reward. To evaluate the pure effect of reward utilization, we compare LM training and DMO training without TDA.

5.3.1 Evaluation on Extremely Noisy Dataset

In this experiment, we examine how robustly our method can train models on the noisy dataset. To simulate the case where the training dataset is extremely noisy, we construct datasets in which a certain percentage of the ground truth captions in the Flickr8k dataset are replaced with entirely irrelevant captions (we randomly sampled from the COCO dataset). With these datasets, we train GIT-base both with DMO and LM loss and observe how training is affected by those toxic samples. In DMO training, we use the original clean dataset for the reference dataset to ensure that the quality of each sample is accurately scored. We increase the noise ratio from 0% to 80% and evaluate the CIDEr scores of the model trained by LM and DMO. The results are presented in Table 6. Compared to the baseline, both LM and DMO training exhibit a decline in performance; however, while LM training experiences a significant performance drop, DMO manages to minimize this reduction. This indicates that, by utilizing scores as cues, our proposed method can effectively discern samples that should be learned from samples that should be ignored, enabling robust learning even from noisy datasets.

5.3.2 Evaluation under Low-Resource Setting

We simulate a scenario where training is constrained by limited data samples due to low computational resources, as is typical in edge device training that is aimed at minimizing time and battery consumption. We construct small datasets to evaluate how effectively data can be utilized under conditions of low resource availability. We reduce the amount of training data progressively from 20% to 80%. For the scoring in DMO training, we use the full original dataset as the reference dataset. The results are presented in Table 7. While LM training exhibits a 9% drop in scores as the data size decreases, DMO demonstrates robust learning even with limited data, exhibiting a smaller decrease of 4%. This result demonstrates that our method can efficiently learn even from a small number of data samples by leveraging the importance score of each sample.

6 Conclusion

In this paper, we present *Direct Metric Optimization* (DMO), which is a lightweight final-metric-optimizing training method. We hypothesize that diverse text augmentation can substitute the exploration in RL, and show that self-supervised training on reward-weighted augmented data leads to direct and stable metric optimization. Our experiments demonstrate that DMO can directly optimize evaluation metrics across models of various architectures and parameter sizes, and stably achieves performance comparable to the SOTA RL method while saving hundreds of times more model forwarding iterations and greater amounts of computation time. With these practical advantages of stable and lightweight cost of tuning, DMO emerges as a new promising choice for metric optimization in the era of large-scale VLMs.

622 Limitations

623 Although our experiments yield promising results,
624 it is important to acknowledge the limitations of our
625 method. The first limitation is the quality subopti-
626 mality of TDA. Our approach substitutes the explo-
627 ration phase in RL with diverse data augmentation.
628 However, in theory, data augmentation is distinct
629 from exploration because it does not actively pur-
630 sue higher rewards. Consequently, over extended
631 training periods, RL methods, which consistently
632 seek new and higher-quality samples, can outper-
633 form our DMO which relies on a fixed dataset.
634 However, considering that RL methods for VLM
635 are often very sensitive to hyperparameters and
636 challenging to optimize, our DMO offers distinct
637 practical advantages such as learning stability and
638 the straightforward training process without the
639 need for intricate hyperparameter tuning—benefits
640 that are absent in most of RL approaches.

641 Another limitation is the data augmentation over-
642 head. While DMO avoids the computationally ex-
643 pensive exploration process in RL, data augmenta-
644 tion still necessitates a certain computational cost.
645 Therefore, considering the data preparation phase
646 in addition to the training phase, the computational
647 costs required for DMO increase and DMO’s supe-
648 riority over RL methods in terms of computational
649 costs is diminished. However, a key distinction
650 from exploration is that TDA can be conducted on
651 separate machines (e.g., cloud servers) from the
652 one the target VLM is deployed on. This aspect
653 becomes particularly beneficial for model tuning in
654 scenarios where resources such as time, memory,
655 and battery of devices are constrained, as is typical
656 in edge device training. Collecting augmented data
657 on different servers enables models on the resource-
658 constrained device to bypass the data augmentation
659 overhead, making DMO a genuinely lightweight
660 metric optimization method.

661 References

662 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama
663 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
664 Diogo Almeida, Janko Altenschmidt, Sam Altman,
665 Shyamal Anadkat, et al. 2023. Gpt-4 technical report.
666 *arXiv preprint arXiv:2303.08774*.

667 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc,
668 Antoine Miech, Iain Barr, Yana Hasson, Karel
669 Lenc, Arthur Mensch, Katherine Millican, Malcolm
670 Reynolds, et al. 2022. Flamingo: a visual language

model for few-shot learning. *Advances in Neural
Information Processing Systems*, 35:23716–23736.

Peter Anderson, Basura Fernando, Mark Johnson, and
Stephen Gould. 2016. Spice: Semantic propositional
image caption evaluation. In *ECCV*.

Viktat Atliha and Dmitriy Šešok. 2020. Text augmen-
tation using bert for image captioning. *Applied Sci-
ences*, 10(17):5978.

Ashutosh Baheti, Ximing Lu, Faeze Brahman, Ronan Le
Bras, Maarten Sap, and Mark Riedl. 2023. Improving
language models with advantage-based offline policy
gradients. *arXiv preprint arXiv:2305.14718*.

Satanjeev Banerjee and Alon Lavie. 2005. **METEOR:
An automatic metric for MT evaluation with im-
proved correlation with human judgments**. In *Pro-
ceedings of the ACL Workshop on Intrinsic and Ex-
trinsic Evaluation Measures for Machine Transla-
tion and/or Summarization*, pages 65–72, Ann Arbor,
Michigan. Association for Computational Linguis-
tics.

Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee,
Aditya Grover, Misha Laskin, Pieter Abbeel, Aravind
Srinivas, and Igor Mordatch. 2021. Decision trans-
former: Reinforcement learning via sequence mod-
eling. *Advances in neural information processing
systems*, 34:15084–15097.

Paul F Christiano, Jan Leike, Tom Brown, Miljan Mar-
tic, Shane Legg, and Dario Amodei. 2017. **Deep
reinforcement learning from human preferences**. In
Advances in Neural Information Processing Systems,
volume 30. Curran Associates, Inc.

Terrance DeVries and Graham W Taylor. 2017. Im-
proved regularization of convolutional neural net-
works with cutout. *arXiv preprint arXiv:1708.04552*.

Guiguang Ding, Minghui Chen, Sicheng Zhao, et al.
2019. **Neural image caption generation with
weighted training and reference**. *Cognitive Com-
putation*, 11:763–777.

Lijie Fan, Dilip Krishnan, Phillip Isola, Dina Katabi,
and Yonglong Tian. 2023. Improving clip training
with language rewrites. In *NeurIPS*.

Steven Y Feng, Varun Gangal, Jason Wei, Sarath Chan-
dar, Soroush Vosoughi, Teruko Mitamura, and Ed-
uard Hovy. 2021. A survey of data augmentation ap-
proaches for nlp. *arXiv preprint arXiv:2105.03075*.

Junlong Gao, Shiqi Wang, Shanshe Wang, Siwei Ma,
and Wen Gao. 2019. Self-critical n-step training for
image captioning. In *Proceedings of the IEEE/CVF
Conference on Computer Vision and Pattern Recog-
nition*, pages 6300–6308.

Micah Hodosh, Peter Young, and Julia Hockenmaier.
Flickr8k dataset.

723	Youngsoo Jang, Jongmin Lee, and Kee-Eung Kim. 2022.	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	779
724	Gpt-critic: Offline reinforcement learning for end-to-	Carroll Wainwright, Pamela Mishkin, Chong Zhang,	780
725	end task-oriented dialogue systems. In <i>10th Inter-</i>	Sandhini Agarwal, Katarina Slama, Alex Ray, et al.	781
726	<i>national Conference on Learning Representations,</i>	2022. Training language models to follow instruc-	782
727	<i>ICLR 2022. International Conference on Learning</i>	tions with human feedback. <i>Advances in Neural</i>	783
728	<i>Representations, ICLR.</i>	<i>Information Processing Systems</i> , 35:27730–27744.	784
729	Marcin Junczys-Dowmunt, Roman Grundkiewicz,	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	785
730	Tomasz Dwojak, Hieu Hoang, Kenneth Heafield,	Jing Zhu. 2002. Bleu: a method for automatic evalu-	786
731	Tom Neckermann, Frank Seide, Ulrich Germann,	ation of machine translation. In <i>Proceedings of the</i>	787
732	Alham Fikri Aji, Nikolay Bogoychev, André F. T.	<i>40th annual meeting of the Association for Computa-</i>	788
733	Martins, and Alexandra Birch. 2018. Marian: Fast	<i>tional Linguistics</i> , pages 311–318.	789
734	neural machine translation in C++ . In <i>Proceedings of</i>	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya	790
735	<i>ACL 2018, System Demonstrations</i> , pages 116–121,	Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sas-	791
736	Melbourne, Australia. Association for Computational	try, Amanda Askell, Pamela Mishkin, Jack Clark,	792
737	Linguistics.	et al. 2021. Learning transferable visual models from	793
738	Wooyoung Kang, Jonghwan Mun, Sungjun Lee, and	natural language supervision. In <i>International confer-</i>	794
739	Byungseok Roh. 2023. Noise-aware learning from	<i>ence on machine learning</i> , pages 8748–8763. PMLR.	795
740	web-crawled image-text data for image captioning.	Marc’Aurelio Ranzato, Sumit Chopra, Michael Auli,	796
741	In <i>Proceedings of the IEEE/CVF International Con-</i>	and Wojciech Zaremba. 2015. Sequence level train-	797
742	<i>ference on Computer Vision</i> , pages 2942–2952.	ing with recurrent neural networks. <i>arXiv preprint</i>	798
743	Andrej Karpathy and Li Fei-Fei. 2015. Deep visual-	<i>arXiv:1511.06732</i> .	799
744	semantic alignments for generating image descrip-	Steven J Rennie, Etienne Marcheret, Youssef Mroueh,	800
745	tions. In <i>Proceedings of the IEEE conference on</i>	Jerret Ross, and Vaibhava Goel. 2017. Self-critical	801
746	<i>computer vision and pattern recognition</i> , pages 3128–	sequence training for image captioning. In <i>Proceed-</i>	802
747	3137.	<i>ings of the IEEE conference on computer vision and</i>	803
748	Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin John-	<i>pattern recognition</i> , pages 7008–7024.	804
749	son, Kenji Hata, Joshua Kravitz, Stephanie Chen,	Piyush Sharma, Nan Ding, Sebastian Goodman, and	805
750	Yannis Kalantidis, Li-Jia Li, David A Shamma, et al.	Radu Soricut. 2018. Conceptual captions: A cleaned,	806
751	2017. Visual genome: Connecting language and vi-	hypernymed, image alt-text dataset for automatic im-	807
752	sion using crowdsourced dense image annotations.	age captioning . In <i>Proceedings of the 56th Annual</i>	808
753	<i>International journal of computer vision</i> , 123:32–73.	<i>Meeting of the Association for Computational Lin-</i>	809
754	Sergey Levine, Aviral Kumar, George Tucker, and Justin	<i>guistics (Volume 1: Long Papers)</i> , pages 2556–2565,	810
755	Fu. 2020. Offline reinforcement learning: Tutorial,	Melbourne, Australia. Association for Computational	811
756	review, and perspectives on open problems. <i>arXiv</i>	<i>Linguistics</i> .	812
757	<i>preprint arXiv:2005.01643</i> .	Ruizhe Shi, Yuyao Liu, Yanjie Ze, Simon S Du, and	813
758	Bohan Li, Yutai Hou, and Wanxiang Che. 2022. Data	Huazhe Xu. 2023. Unleashing the power of pre-	814
759	augmentation approaches in natural language pro-	trained language models for offline reinforcement	815
760	cessing: A survey . <i>AI Open</i> , 3:71–90.	learning. <i>arXiv preprint arXiv:2310.20587</i> .	816
761	Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi.	Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu,	817
762	2023. Blip-2: Bootstrapping language-image pre-	Chunyuan Li, Yikang Shen, Chuang Gan, Liang-	818
763	training with frozen image encoders and large lan-	Yan Gui, Yu-Xiong Wang, Yiming Yang, et al. 2023.	819
764	guage models. <i>arXiv preprint arXiv:2301.12597</i> .	Aligning large multimodal models with factually aug-	820
765	Chin-Yew Lin. 2004. ROUGE: A package for auto-	mented rlhf. <i>arXiv preprint arXiv:2309.14525</i> .	821
766	matic evaluation of summaries . In <i>Text Summariza-</i>	Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe,	822
767	<i>tion Branches Out</i> , pages 74–81, Barcelona, Spain.	Jon Shlens, and Zbigniew Wojna. 2016. Rethinking	823
768	Association for Computational Linguistics.	the inception architecture for computer vision. In	824
769	Tsung-Yi Lin, Michael Maire, Serge J. Belongie,	<i>Proceedings of the IEEE conference on computer</i>	825
770	Lubomir D. Bourdev, Ross B. Girshick, James Hays,	<i>vision and pattern recognition</i> , pages 2818–2826.	826
771	Pietro Perona, Deva Ramanan, Piotr Doll’ar, and	Gemini Team, Rohan Anil, Sebastian Borgeaud,	827
772	C. Lawrence Zitnick. 2014. Microsoft COCO: com-	Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu,	828
773	mon objects in context . <i>CoRR</i> , abs/1405.0312.	Radu Soricut, Johan Schalkwyk, Andrew M Dai,	829
774	Siqi Liu, Zhenhai Zhu, Ning Ye, Sergio Guadarrama,	Anja Hauth, et al. 2023. Gemini: a family of	830
775	and Kevin Murphy. 2017. Improved image caption-	highly capable multimodal models. <i>arXiv preprint</i>	831
776	ing via policy gradient optimization of spider. In	<i>arXiv:2312.11805</i> .	832
777	<i>Proceedings of the IEEE international conference on</i>		
778	<i>computer vision</i> , pages 873–881.		

833	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrubhi Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .	889
834		890
835		891
836		892
837		
838		
839	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>Advances in neural information processing systems</i> , 30.	893
840		894
841		895
842		896
843		897
844	Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. 2015. Cider: Consensus-based image description evaluation. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 4566–4575.	898
845		899
846		900
847		901
848		902
849	Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. 2015. Show and tell: A neural image caption generator. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 3156–3164.	903
850		904
851		905
852		906
853		907
854	Jianfeng Wang, Zhengyuan Yang, Xiaowei Hu, Linjie Li, Kevin Lin, Zhe Gan, Zicheng Liu, Ce Liu, and Lijuan Wang. 2022. Git: A generative image-to-text transformer for vision and language. <i>arXiv preprint arXiv:2205.14100</i> .	908
855		909
856		910
857		911
858		912
859	Ronald J Williams and David Zipser. 1989. A learning algorithm for continually running fully recurrent neural networks. <i>Neural computation</i> , 1(2):270–280.	913
860		914
861		915
862	Haiyang Xu, Qinghao Ye, Ming Yan, Yaya Shi, Jiabo Ye, Yuanhong Xu, Chenliang Li, Bin Bi, Qi Qian, Wei Wang, et al. 2023. mplug-2: A modularized multi-modal foundation model across text, image and video. <i>arXiv preprint arXiv:2302.00402</i> .	916
863		917
864		918
865		
866		
867	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 <i>International conference on machine learning</i> , pages 2048–2057. PMLR.	919
868		920
869		921
870		922
871		923
872		924
873	Kaicheng Yang, Jiankang Deng, Xiang An, Jiawei Li, Ziyong Feng, Jia Guo, Jing Yang, and Tongliang Liu. 2023. Alip: Adaptive language-image pre-training with synthetic caption. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 2922–2931.	925
874		926
875		927
876		928
877		929
878		930
879	Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. 2022. Coca: Contrastive captioners are image-text foundation models. <i>arXiv preprint arXiv:2205.01917</i> .	931
880		932
881		933
882		
883	Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. 2017a. mixup: Beyond empirical risk minimization. <i>arXiv preprint arXiv:1710.09412</i> .	934
884		935
885		936
886	Le Zhang, Yanshuo Zhang, Xin Zhao, and Zexiao Zou. 2021. Image captioning via proximal policy optimization. <i>Image and Vision Computing</i> , 108:104126.	937
887		
888		
	Li Zhang, Flood Sung, Feng Liu, Tao Xiang, Shaogang Gong, Yongxin Yang, and Timothy M Hospedales. 2017b. Actor-critic sequence training for image captioning. <i>arXiv preprint arXiv:1706.09601</i> .	
	Zhiyuan Zhao, Bin Wang, Linke Ouyang, Xiaoyi Dong, Jiaqi Wang, and Conghui He. 2023. Beyond hallucinations: Enhancing lvlms through hallucination-aware direct preference optimization. <i>arXiv preprint arXiv:2311.16839</i> .	
	A Experiments Setting in Detail	
	A.1 Hardware Environment	
	We use two RTX A6000 Ada GPUs for GIT-base/large and four H100 GPUs for BLIP2-2.7B. Because each GPU has a different size of VRAM, we adopt the gradient accumulation method so that the batch size becomes 960 regardless of the size of the VRAM of the GPU that is used in each experiment. We present the hyperparameters for GIT-base/large and BLIP2-2.7B in Table 8. This configuration explains the number of model iterations in LM/DMO training presented in Table 5. The training uses 6091 images and 11 captions per image (including augmented captions) across 3 epochs. Dividing this by a mini-batch size of 60 and 2 GPUs results in approximately 1.68×10^3 iterations. For SCST, the maximum token length is set to 128 and the model recursively generates tokens up to this length, resulting in nearly 100 times more model forwarding iterations compared to LM/DMO training.	
	A.2 Configuration for SCST Training	
	Due to the instability of training with SCST, it necessitates pre-fine-tuning through standard self-supervised training (Rennie et al., 2017). Therefore, we initially fine-tune the model for three epochs with a learning rate of 1.0×10^{-5} before applying SCST. Given that fine-tuning has already been completed, we reduce the learning rate to 5.0×10^{-6} and only update the parameters of the text decoder to stabilize training. Moreover, during the sampling process in SCST, we opt for a temperature of 0.1. This is because we observe that higher temperatures, such as 0.5 or 1.0, often lead the model to generate random, meaningless sequences of words, which ultimately results in model collapse.	
	B Prompt for Llama2 Paraphrasing	
	For the Llama2-paraphrasing method, we employ the same prompting method proposed in LaCLIP (Fan et al., 2023). We present the prompt that	

Model	Learning Rate	Batch Size	Mini-Batch Size	Grad. Acc. Step	GPU
GIT-base	1.0×10^{-5}	960	60	8	RTX A6000 Ada \times 2
GIT-large	1.0×10^{-5}	960	30	16	RTX A6000 Ada \times 2
BLIP2-2.7B	1.0×10^{-5}	960	30	8	H100 \times 4

Table 8: Hyperparameters for GIT-base/large and BLIP2-2.7B. The number of batch sizes is equal to the product of the mini-batch size, the number of gradient accumulation steps (Grad. Acc. Step) and the number of GPUs.

Rewrites the following image descriptions:

A young child splashes in a green and yellow wading pool. => A little boy is splashing around in the water in a kiddie pool.

A piece of banana, some strawberries, shish-kabobs, and a muffin are on a tray. => A tray of fruit and a muffin on a table.

Three people are sitting on a bus stop bench in between two tropical palm trees. => Three individuals are seated on a bench at a bus stop, flanked by two tropical palm trees.

<target captions> =>

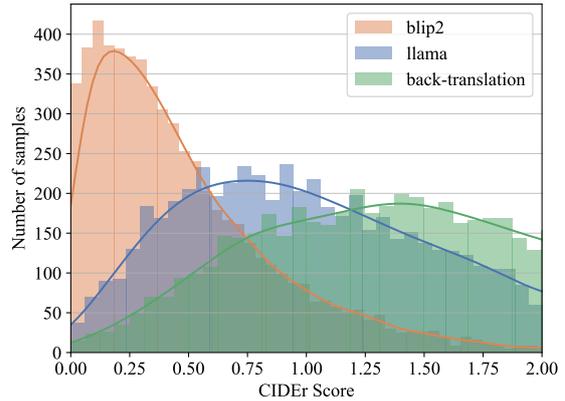


Figure 4: CIDEr score distributions of captions augmented by back-translation, BLIP2-sampling, and Llama2-paraphrasing.

Figure 3: The prompt that is used in Llama2-paraphrasing. The part '<target caption>' in the prompt text is replaced by the caption to be augmented.

is used in our experiments in Figure 3. The part '<target captions>' in the prompt text is replaced by the caption to be paraphrased. In the prompt, three examples of paraphrasing are provided. The first and second examples are constructed by regarding the two captions for the same image in each Flickr8K and COCO dataset as the captions before and after paraphrasing. The third paraphrasing example is made by feeding ChatGPT4 a caption from Flickr8k with the prompt "rewrite this image caption".

C Distributions Difference by TDA Method

In this section, we explore how the quality of samples varies across different TDA methods. We present the distribution of scores of samples generated by each TDA method in Figure 4. Scores are CIDEr scores based on the training dataset of Flickr8k. The back-translation tends to yield higher scores while BLIP2-sampling tends to pro-

duce samples of lower scores. Examples of generated captions by each TDA method are shown in Figure 5. We find that captions generated by the back-translation show little change compared to the ground truth (GT) captions. On the other hand, BLIP2 sampling generates captions that are significantly different from GT in terms of style and level of detail. Back-translation receives the GT captions to augment captions. Thus captions generated by back-translation closely resemble GT captions. On the other hand, BLIP2-sampling generates captions solely from images. Therefore, captions generated by BLIP2-sampling often deviate from the GT captions and sometimes include incorrect descriptions (e.g., BLIP2 misidentified the person in the image as "woman" in Figure 5). Moreover, because the paraphrasing by Llama2 takes the GT captions as a prompt, the generated captions by Llama2 are semantically close to the GT captions. However, the changes from GT captions are greater than those of captions augmented by back-translation, owing to the prompt which encourages Llama2 to change expressions (e.g., 'surf' is paraphrased as 'ride a wave' in Figure 5).

Interestingly, considering the result that DMO training on a dataset solely augmented by back-

958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983



Ground truth

- a man surfing in the ocean.

Back-translation

- a man surfing the ocean.

BLIP2-sampling

- a woman wearing a wet suit riding a single white surfboard on a large wave in fresh ocean water.

Llama2-paraphrasing

- a man rides a wave in the ocean.

Figure 5: Examples of captions generated by each augmentation method, back-translation, BLIP2-sampling, and Llama2-paraphrasing.

984 translation does not improve scores (shown in Ta-
985 ble 3), TDA that often produces samples of high
986 scores may not provide an advantage for DMO if
987 it hardly alters the expression of the original GT
988 data. Rather, TDA methods that produce samples
989 of diverse expressions and structures can improve
990 the performance of DMO by introducing various in-
991 formation that is not present in the original dataset.
992 With the analysis above, we emphasize that the di-
993 versity of samples is important for DMO training
994 and especially diversifying the augmentation tech-
995 niques themselves is an effective approach because
996 utilizing multiple augmentation methods generates
997 a diverse set of samples across various distribu-
998 tions.

999 **D Examples of Captions Generated by**
1000 **DMO, LM Training and SCST**

1001 Figure 6 shows examples of images and corre-
1002 sponding ground truth captions and captions gen-
1003 erated by a GIT-base model trained by DMO, stan-
1004 dard LM training and SCST. Models are trained for
1005 three epochs on the textually augmented Flickr8k
1006 dataset. Expressions which seem to be unique to
1007 the image are underlined. Those examples show
1008 that the model trained by DMO captures distinct-
1009 ive information and objects within the image and
1010 depicts them in the generated captions.



Ground truth: a girl wearing red and blue clothing poses for a man kneeling to take her picture.

DMO: a man is kneeling down to take a picture of a girl in a red and blue dress.

LM: a man and a girl are playing in a field.

SCST: a man and a little girl are playing in the grass.



Ground truth: a parade of people is marching in santa costumes on a city street.

DMO: a group of people dressed in santa costumes walk down a city street.

LM: a group of people dressed in costumes walk down a street.

SCST: a group of men in red costumes are walking down the street.



Ground truth: some children are playing soccer in front of a large wooden door.

DMO: three children are playing soccer in front of a large wooden door.

LM: three boys are playing soccer in front of a building.

SCST: three children are playing soccer in front of a building.



Ground truth: a woman sings into a microphone while playing guitar.

DMO: a woman is singing into a microphone while holding a guitar.

LM: a woman is singing a guitar.

SCST: a woman is playing a guitar.



Ground truth: a man and a police officer are smiling at the photographer.

DMO: a man in a suit and tie is standing next to a police officer.

LM: a man in a suit and tie is standing in front of a building with a man in a suit.

SCST: two men in uniform are smiling at the camera.

Figure 6: Examples of ground truth captions and captions generated by GIT-base model trained by DMO, standard LM training and SCST. Each model is trained on the textually augmented Flickr8k dataset for 3 epochs. Expressions which seem to be unique to the image are underlined.