

# Long-Tail Classification for Distinctive Image Captioning: A Simple yet Effective Remedy for Side Effects of Reinforcement Learning

Anonymous ACL submission

## Abstract

Distinctiveness is a desirable feature of image captions. Captions should cover the characteristic details of input images. However, recent high-performing captioning models that are trained with reinforcement learning (RL) tend to generate overly generic captions despite their high performance in various other criteria. Interestingly, it has also been reported that their outputs are composed of a limited number of common words and rarely contain tail-class words, *i.e.*, low-frequency words in the training corpus. Vocabulary size is closely related to distinctiveness as it is difficult for a model to describe details beyond its vocabulary. Based on this insight, we hypothesize that the limited vocabulary of RL models is the major factor limiting their distinctiveness. We recast distinctive image captioning as a simpler task of long-tail classification to increase the vocabulary and then propose lightweight fine-tuning methods to encourage tail-class word generation. The experimental results demonstrate that our methods significantly enhance the distinctiveness of existing RL models as well as their vocabulary size, without sacrificing quality. Our methods also outperform previous distinctiveness-aware methods with a small computational cost of minor modifications to pre-trained RL models.<sup>1</sup>

## 1 Introduction

Image captioning plays a fundamental role at the intersection of computer vision and natural language processing by converting the information in images into natural language descriptions. Generated captions can be used in various downstream tasks, such as aiding visually impaired users (Gurari et al., 2020), visual question answering on images and videos (Fisch et al., 2020; Kim et al., 2020), visual dialogue (White et al., 2021), and news generation (Zhang et al., 2021b).

<sup>1</sup>The code will be made available on our website.

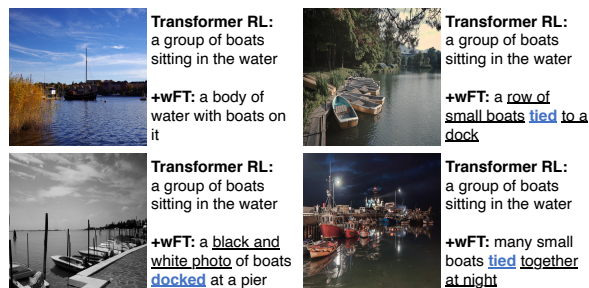


Figure 1: Caption examples in the MS COCO validation set. **Transformer RL** is a Transformer captioning model trained with RL and **+wFT** is our fine-tuning method. Transformer RL generates exactly the same caption for the four images. The underlined words indicate the characteristic information that are not mentioned by Transformer RL, and the blue words are those that have never appeared in the outputs of the model.

For those downstream tasks, the generated captions should be **distinctive**: captions should cover the characteristic and important details of the input images. However, current captioning models tend to generate overly generic captions (Dai and Lin, 2017; Dai et al., 2017; Wang and Chan, 2019; Wang et al., 2020c). For example, a high-performing captioning model based on Transformer (Vaswani et al., 2017) generates exactly the same caption for the four different images shown in Figure 1, ignoring the other salient details of each image.

To address the problem of overly generic captions, some studies have been conducted on **distinctive image captioning**, which is also called descriptive image captioning or discriminative image captioning. Previous research has created new rewards regarding distinctiveness or new model architectures to enhance distinctiveness. These approaches improved the performance with regard to distinctiveness and other evaluation metrics; however, their models come with additional computations and require training from scratch.

Instead of creating or paying those computational costs, we first analyze the cause of the current

overly generic captions to explore ways to improve the distinctiveness of *pre-trained, existing models*. In particular, we focus on high-performing captioning models that are trained with the standard **reinforcement learning (RL)** (Rennie et al., 2017), which is the *de facto* standard training method in current image captioning (Stefanini et al., 2021). Those models have greater room to improve distinctiveness as they unexpectedly perform poor in distinctiveness despite the significant advantages in various other criteria (Liu et al., 2019; Wang et al., 2020a). Interestingly, some previous studies have reported that RL decreased the vocabulary size of output captions (Wang and Chan, 2019; Liu et al., 2019; Wang et al., 2020a). Vocabulary size is closely related to distinctiveness as it is difficult for a model to describe details beyond its vocabulary. Based on this insight, we hypothesize that the limited vocabulary of RL models is the major factor limiting their distinctiveness.

To directly increase the vocabulary of RL models, we recast distinctive image captioning as a simpler task of **long-tail classification**. Unlike previous approaches, our methods do not require any distinctiveness reward, new model architecture, or training from scratch. Our methods focus on generating **tail-class words**, *i.e.*, low-frequency words in the training corpus. Owing to their simplicity, our methods can be realized by single-epoch fine-tuning of pre-trained, existing RL models.

The experimental results confirm our hypothesis by revealing that our methods significantly boost both vocabulary size and distinctiveness from existing RL models. We also demonstrate that our methods outperform previous distinctiveness-aware methods with a small computational cost of minor modifications to pre-trained RL models.

## 2 RL Model Distinctiveness and Limited Vocabulary

Currently, RL is the *de facto* standard training method for models used in image captioning because it significantly improves the performance in various evaluation metrics (Stefanini et al., 2021). However, it does not improve distinctiveness and may even decrease it (Liu et al., 2019; Wang et al., 2020a). In this section, we examine the cause of overly generic captions generated by RL models and hypothesize that their limited vocabulary hinders their distinctiveness.

### 2.1 RL in Image Captioning

We provide a brief overview of the standard RL algorithm used in image captioning. It was proposed by Ranzato et al. (2015) and refined by Rennie et al. (2017). Their goal was to directly optimize non-differentiable test-time metrics by minimizing the negative expected reward:

$$\mathcal{L}_{\text{RL}}(\theta) = -\mathbb{E}_{w^s \sim p_\theta(w^s | I)}[r(w^s)], \quad (1)$$

where  $w^s = (w_1^s, \dots, w_T^s)$  is a sequence sampled from a policy  $p_\theta$ ,  $I$  is the input image, and  $r(\cdot)$  is a reward function that returns a reward for  $w^s$ . To compute the gradient of  $\mathcal{L}(\theta)$ , Ranzato et al. (2015) applied the REINFORCE algorithm (Williams, 1992) to text generation. In their algorithm, the model parameters  $\theta$  are updated with the following gradient:

$$\nabla_\theta \mathcal{L}_{\text{RL}}(\theta) \approx -(r(w^s) - b) \nabla_\theta \log p_\theta(w^s | I). \quad (2)$$

Here,  $b$  is a baseline reward that reduces the variance in the gradient. Typically, the reward function  $r(\cdot)$  is CIDEr (Vedantam et al., 2015), and the baseline reward  $b$  is a reward for a sequence sampled with greedy decoding (Rennie et al., 2017).

### 2.2 RL Results in Limited Vocabulary

Despite its effectiveness, RL has been found to decrease distinctiveness and the number of unique n-grams in output captions (Liu et al., 2019; Wang et al., 2020a). As the relation between RL and those negative effects is not obvious, it was just considered a curious case.

However, recent research on the weaknesses of RL has revealed the relation between RL and a limited vocabulary. Recently, Choshen et al. (2020) and Kiegeland and Kreutzer (2021) empirically showed that RL makes the output distribution peaky. As shown in Section 2.1, RL samples sequences from policy  $p_\theta$ . Typically, policy  $p_\theta$  is computed using a captioning model pre-trained with the **Cross-Entropy (CE)** loss on ground-truth captions. However, text-generation models in general are known to output skewed distributions. Specifically, the distributions tend to be skewed towards **head-class words**, *i.e.*, high-frequency words in the training corpus (Nguyen and Chiang, 2018; Raunak et al., 2020; Demeter et al., 2020; Holtzman et al., 2020). Thus, RL can sample and reward head-class words but cannot sample or reward tail-class words during training.

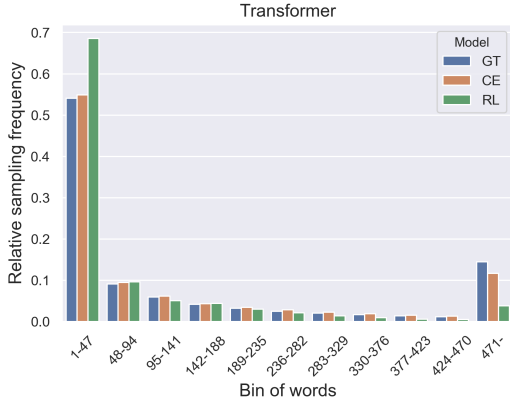


Figure 2: Relative frequency of the words in the sequences sampled for the training images. Five sequences were sampled for each image. The words (9,486 unique words excluding an out-of-vocabulary token  $\langle \text{unk} \rangle$ ) are sorted by their frequency in ground-truth captions and divided into 200 bins. We show the first 10 bins and the sum of the rest. GT is the ground-truth caption of the training images, CE is the output of a captioning model trained with the CE loss, and RL is the output of a captioning model trained with RL. Here, we used the Transformer captioning model.

This imbalance results in shifts of the probability mass from tail-class words to head-class words, further limiting the vocabulary to head-class words.

Figure 2 confirms this phenomenon in image captioning by plotting the relative frequency of the words sampled for the training images. The words are sorted by their frequency in ground-truth captions and divided into 200 bins. Compared to the ground-truth captions and sequences sampled with a CE model, the sequences sampled with an RL model are clearly limited to the head-class words, forming a peaky distribution<sup>2</sup>.

### 2.3 Limited Vocabulary Results in Overly Generic Captions

Standard encoder–decoder captioning models generate captions using sequential vocabulary-size classification. However, the actual vocabulary a model can generate is much smaller than the entire vocabulary as the output distribution is highly skewed towards head-class words. If the actual vocabulary cannot cover the details of an image, the model is forced to avoid those details and output only the information that can be described by

<sup>2</sup>Although Figure 2 shows only the results obtained with the Transformer captioning model, we also confirmed that other models output peaky distributions (Rennie et al., 2017; Anderson et al., 2018). See Appendix A for the details.

head-class words. For example, the blue words in Figure 1 are not in the actual vocabulary of the RL model; these words have never been generated by the RL model. As a result, the RL model had to ignore the characteristic relations *tied* and *docked* and ended up describing exactly the same information for all four images.

Based on the above observations, we hypothesize that the limited vocabulary of RL models hinders their distinctiveness. To directly address this limitation, we recast distinctive image captioning as a simpler task of increasing the actual vocabulary.

## 3 Long-Tail Classification to Remedy the Side Effects of RL

RL results in the limited vocabulary because it steals the probability mass from tail-class words of ground-truth captions. Thus, those tail-class words are the key to addressing the limitation. Wang and Chan (2019) jointly optimized both the RL loss and the CE loss on ground-truth captions so that the tail-class words in ground-truth captions would be more likely to be sampled during RL training. However, this approach still relies on the sampling from a skewed policy and requires training from scratch.

To increase the actual vocabulary more effectively and efficiently, we propose two fine-tuning methods based on long-tail classification. Our methods are designed to directly encourage the generation of tail-class words with only single-epoch fine-tuning on pre-trained, existing RL models.

### 3.1 Simple Fine-Tuning

The first method is a **simple fine-tuning (sFT)** method for ground-truth captions. It is based on a decoupled two-stage training (Kang et al., 2020), which is a current strong baseline model for long-tail classification (Tang et al., 2020; Menon et al., 2020; Wang et al., 2020b). Kang et al. (2020) decoupled the learning procedure into representation learning and classification, and then found that classification is critical for long-tail classification. They decoupled the classification model  $f_{\theta}(\cdot)$  into an encoder  $g_{\theta_e}(\cdot)$  and a classifier consisting of weight and bias parameters:  $f_{\theta}(x) = \mathbf{W}^{\top} g_{\theta_e}(x) + \mathbf{b}$ . In the first stage of training, they trained the entire classification model  $f_{\theta}(\cdot)$  on a full training dataset. In the second stage, they fixed the encoder parameters  $\theta_e$  and adjusted only the classifier parameters

ters. For the second-stage adjustment, they applied class-balanced sampling to encourage learning on tail-class labels.

Following Kang et al. (2020), we decouple a captioning model into an encoder and a classifier. In image captioning, the first-stage training of Kang et al. (2020) corresponds to RL training on the full training dataset. Likewise, the second-stage training corresponds to adjusting the classifier parameters on the *vocabulary-balanced* sequences. However, sampling from the skewed policy of text-generation models cannot provide sequences containing tail-class words (Section 2.2). Thus, we use ground-truth captions as relatively vocabulary-balanced samples. sFT simply fine-tunes the classifier parameters of a pre-trained RL captioning model by minimizing the CE loss on ground-truth captions:

$$\mathcal{L}_{\text{CE}}(\hat{\theta}) = -\frac{1}{T} \sum_{t=1}^T \log p_{\hat{\theta}}(w_t^g | w_{<t}^g, I), \quad (3)$$

where  $w^g = (w_1^g, \dots, w_T^g)$  is a ground-truth caption of image  $I$ , and the model parameters  $\hat{\theta}$  are initialized with RL training. The conditional probability  $p_{\theta}(w_t^g | w_{<t}^g, I)$  is computed using the following softmax function:

$$p_{\theta}(w_t^g | w_{<t}^g, I) = \frac{\exp(\beta z_{w_t^g})}{\sum_{w_i \in \mathcal{W}} \exp(\beta z_{w_i})}, \quad (4)$$

$$z = \mathbf{W}^{\top} g_{\theta_e}(w_{<t}^g, I) + \mathbf{b}, \quad (5)$$

where  $z \in \mathbb{R}^{|\mathcal{W}|}$ ,  $\mathbf{W} \in \mathbb{R}^{d \times |\mathcal{W}|}$ , and  $\mathbf{b} \in \mathbb{R}^{|\mathcal{W}|}$ .  $\mathcal{W}$  is the entire vocabulary, and  $d$  is the dimension of the hidden states of an encoder  $g_{\theta_e}(\cdot)$ .  $z_{w_i}$  indicates the element of  $z$  at the index of a word  $w_i \in \mathcal{W}$ .  $\beta$  is an inverse-temperature hyperparameter that controls the steepness of the softmax distribution. We use LSTM (Hochreiter and Schmidhuber, 1997) or Transformer (Vaswani et al., 2017) for  $g_{\theta_e}(\cdot)$ . During fine-tuning, only the classifier parameters  $\{\mathbf{W}, \mathbf{b}\} \in \hat{\theta}$  are updated with the gradients  $\nabla_{\mathbf{W}} \mathcal{L}_{\text{CE}}(\hat{\theta})$  and  $\nabla_{\mathbf{b}} \mathcal{L}_{\text{CE}}(\hat{\theta})$ , respectively.

### 3.2 Weighted Fine-Tuning

Ground-truth captions contain more tail-class words than sampled sequences, but some tail-class words are still difficult to learn because of their low frequency. Our second method is **weighted fine-tuning (wFT)**, which further pursues vocabulary balance by rebalancing the loss for head-class words and tail-class words in ground-truth captions.

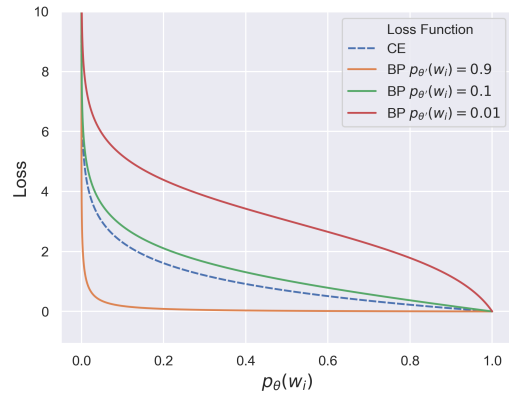


Figure 3: Visualization of the CE loss  $-\log p_{\theta}(w_i)$  and BP loss  $-\log p_{\theta, \theta'}(w_i)$ . To compute the BP loss, we need the entire distribution of  $\{p_{\theta}(w_i)\}_{w_i \in \mathcal{W}}$  and  $\{p_{\theta'}(w_i)\}_{w_i \in \mathcal{W}}$ . Here, we set the index  $i$  to 1 and assigned  $\frac{1}{5}(1 - p_{\theta}(w_1))$  to the words of the next five indices,  $w_2, \dots, w_6$ . This is because we observed that the five most probable words occupied 99% of the probability in the output distribution of the RL models. We assumed that the five most probable words were the same between  $p_{\theta}$  and  $p_{\theta'}$  as the parameters were initialized with the same RL model. Thus, we assigned  $\frac{1}{5}(1 - p_{\theta'}(w_1))$  to the words of the next five indices,  $w_2, \dots, w_6$ , likewise  $p_{\theta}$ . Here,  $\beta$  was set to 1.

To rebalance the loss, we exploit the head-class bias of RL models: RL models overly assign probability to head-class words, but not to tail-class words. Based on the head-class bias, a ground-truth word that an RL model is confident of should be a head-class word that the model is refrained from further learning, whereas a ground-truth word that an RL model is not confident of should be a tail-class word for the model to learn intensely. wFT incorporates these heuristics by modifying the probability  $p_{\theta}$  of  $\mathcal{L}_{\text{CE}}$  to the probability of the **bias product (BP)** (Clark et al., 2019; He et al., 2019),  $p_{\theta, \theta'}$ , as follows:

$$p_{\theta, \theta'}(w_t^g | w_{<t}^g, I) = \frac{\exp(s_{\theta}^t(w_t^g) + s_{\theta'}^t(w_t^g))}{\sum_{w_i \in \mathcal{W}} \exp(s_{\theta}^t(w_i) + s_{\theta'}^t(w_i))}, \quad (6)$$

where

$$s_{\theta}^t(w_i) = \log p_{\theta}(w_i | w_{<t}^g, I), \quad (7)$$

$$s_{\theta'}^t(w_i) = \log p_{\theta'}(w_i | w_{<t}^g, I). \quad (8)$$

By inserting  $p_{\theta, \theta'}$  into  $\mathcal{L}_{\text{CE}}$ , we define the objective of wFT as follows:

$$\mathcal{L}_{\text{BP}}(\hat{\theta}) = -\frac{1}{T} \sum_{t=1}^T \log p_{\hat{\theta}, \hat{\theta}'}(w_t^g | w_{<t}^g, I). \quad (9)$$



Similar to sFT, both parameters  $\hat{\theta}$  and  $\hat{\theta}'$  are initialized with a captioning model pre-trained with RL. The difference is that, although the classifier parameters of  $\hat{\theta}$  are updated, all the parameters of  $\hat{\theta}'$  are fixed during fine-tuning. Figure 3 shows the change in the BP loss compared to the CE loss. The BP severely suppresses the loss when the head-class-biased policy  $p_{\theta'}$  is confident, and largely increases the loss when  $p_{\theta'}$  is not confident. In this way, the BP allows models to unlearn the head-class bias learned with RL. As with sFT, only the classifier parameters  $\{\mathbf{W}, \mathbf{b}\} \in \hat{\theta}$  are updated with the gradients  $\nabla_{\mathbf{W}} \mathcal{L}_{\text{BP}}(\hat{\theta})$  and  $\nabla_{\mathbf{b}} \mathcal{L}_{\text{BP}}(\hat{\theta})$ , respectively.

We follow Clark et al. (2019) and He et al. (2019) at the evaluation stage, too. We use the probability  $p_{\theta}$  with updated parameters rather than the BP probability  $p_{\theta, \theta'}$  to avoid incorporating the head-class bias of  $p_{\theta'}$  into the predictions.

## 4 Experiments

### 4.1 Setup

**Dataset and Metrics.** We used the MS COCO captioning dataset<sup>3</sup> (Lin et al., 2014; Chen et al., 2015) with Karpathy splitting (Karpathy and Fei-Fei, 2015). After preprocessing, the entire vocabulary size  $|\mathcal{W}|$  was 9,487<sup>4</sup>. In the evaluation, the captions were decoded using a beam search of size 5 and evaluated using various evaluation metrics<sup>5</sup>: BLEU (Papineni et al., 2002), ROUGE (Lin, 2004), METEOR (Denkowski and Lavie, 2014), CIDEr (Vedantam et al., 2015), SPICE (Anderson et al., 2016), and RefCLIP (Hessel et al., 2021). Following the previous studies (Liu et al., 2019; Wang et al., 2020a; Shi et al., 2021b), we evaluated distinctiveness with **R@K** scores: the percentage of captions with which an image–text retrieval model<sup>6</sup> (Faghri et al., 2018) could correctly retrieve the original images from the entire validation/test images within the rank of  $K \in \{1, 5, 10\}$ .

<sup>3</sup>This dataset was intended for image captioning, which is consistent with our use. The dataset was created with the instruction to anonymize people’s proper names (Chen et al., 2015). The dataset was licensed under CC BY 4.0. Each split of training/validation/test set contained 113,287/5,000/5,000 images, and each image had five ground-truth captions.

<sup>4</sup>The words that occur less than five times in the training captions were converted to  $\langle \text{unk} \rangle$  token.

<sup>5</sup>We used the following library, and all the hyperparameters were set to the default values: <https://github.com/jmhessel/pycocoevalcap> (All Rights Reserved)

<sup>6</sup>Following Liu et al. (2019), we used a pre-trained model, `coco_vse+_resnet_restval_finetune`, which is available at <https://github.com/fartashf/vsepp> (Apache License, Version 2.0)

A higher R@K indicates that the model captures more characteristic information of images and generates more distinctive captions. Evaluation was conducted in a single run for each model.

**Comparison Models.** Following Wang et al. (2020a), we used Att2in (Rennie et al., 2017), UpDown (Anderson et al., 2018), and Transformer (Vaswani et al., 2017) as the baseline models. The models were pre-trained with the standard RL (Rennie et al., 2017) and are publicly available<sup>7</sup>. In addition to the baseline models, we compared our models with state-of-the-art distinctiveness-aware models: **CIDErBtw** (Wang et al., 2020a), **NLI** (Shi et al., 2021b), **DiscCap** (Luo et al., 2018), and **Visual Paraphrase** (Liu et al., 2019). The first three created new distinctiveness rewards to be optimized with RL. Visual Paraphrase introduced a new model architecture to paraphrase simpler captions to more complex captions. As we mentioned in the beginning of Section 3, the CE loss on ground-truth captions can be utilized in a different way from our methods. We report the results of jointly optimizing the RL loss and CE loss (**Joint CE** (Wang and Chan, 2019; Edunov et al., 2018)), and also those of solely optimizing the CE loss (**Only CE**) as the baseline without using RL<sup>8</sup>.

**Hyperparameters.** Our models used the same hyperparameters as the baseline models, except for the epoch size, learning rate, and  $\beta$  in Eq. 4. We set the epoch size for fine-tuning to 1 and searched for the best learning rate from  $\{1e-3, 1e-4, 1e-5, 1e-6\}$ . We set  $\beta$  to 1 for  $p_{\theta}$  and searched for the best  $\beta$  of the fixed policy  $p_{\theta'}$  from  $\{0.1, 1\}$ . The best hyperparameters were chosen according to the R@1 scores in the validation set. See Appendix B for the best hyperparameters.

All the models except Visual Paraphrase had the same parameter size as their baseline models<sup>9</sup>.

<sup>7</sup><https://github.com/ruotianluo/self-critical.pytorch> (MIT License): {Att2in, UpDown, Transformer}+self\_critical models.

<sup>8</sup>We optimized  $\mathcal{L}_{\text{Joint}}(\theta) = \lambda \mathcal{L}_{\text{RL}}(\theta) + (1 - \lambda) \mathcal{L}_{\text{CE}}(\theta)$  during RL training. We explored  $\lambda \in \{0.2, 0.5, 0.8\}$ .  $\lambda = 0.8$  for Transformer and  $\lambda = 0.2$  for the others achieved the best R@1 scores in the validation set.  $\lambda = 0$  for Only CE. As with our models, all hyperparameters were set to the same as the baseline models except for the  $\lambda$  and scheduled sampling (Bengio et al., 2015). We disabled scheduled sampling for the CE loss to strictly separate it from the RL loss.

<sup>9</sup>The exact number of parameters was 14,451,985 for Att2in, 52,125,025 for UpDown, and 57,474,832 for Transformer. Note that the fixed parameters  $\theta'$  were not included because they were neither trained nor used in the prediction. Visual Paraphrase has double decoders of Att2in.

	Vocabulary			Standard Evaluation					Distinctiveness				
	Unique-1	Unique-S	Length	BLEU-4	METEOR	ROUGE-L	CIDEr	SPICE	RefCLIP	R@1	R@5	R@10	
Att2in	<b>Att2in RL</b>	445	2,524	9.3	35.3	27.1	56.7	117.4	20.5	79.7	16.3	41.9	57.2
	+ sFT (Ours)	880	3,156	9.0	35.6	27.0	56.5	115.4	20.4	<b>80.3</b>	20.1	48.0	62.8
	+ wFT (Ours)	<b>1,091</b>	<b>3,749</b>	9.0	32.6	26.4	54.9	108.6	19.9	<b>80.3</b>	<b>21.7</b>	<b>50.8</b>	65.2
	CIDErBtw	470	2,630	9.3	35.7	27.2	56.9	<b>119.0</b>	20.7	79.8	17.2	44.1	58.7
	NLI	465	2,626	9.2	35.7	27.2	<b>57.0</b>	<b>119.0</b>	20.6	79.9	17.6	44.4	59.8
	DiscCap <sup>†</sup>		3,093	9.3	<b>36.1</b>			114.2	<b>21.0</b>		21.6	50.3	<b>65.4</b>
	Joint CE	700	2,907	9.1	36.0	<b>27.3</b>	56.4	111.7	19.9	80.0	19.1	46.7	61.5
	Only CE	689	2,845	9.2	35.7	27.1	56.1	110.7	20.1	79.9	19.0	46.6	61.1
-----													
UpDown	<b>UpDown RL</b>	577	3,103	9.5	<b>36.7</b>	27.9	57.6	<b>122.7</b>	<b>21.5</b>	80.5	21.1	49.9	64.6
	+ sFT (Ours)	1,190	3,788	9.2	35.7	27.5	56.5	115.9	21.0	<b>80.9</b>	25.0	56.8	71.2
	+ wFT (Ours)	<b>1,227</b>	<b>4,263</b>	9.3	32.0	26.5	54.3	107.9	20.4	<b>80.9</b>	<b>25.5</b>	<b>58.0</b>	<b>72.6</b>
	CIDErBtw	582	3,108	9.4	<b>36.7</b>	<b>28.0</b>	<b>57.7</b>	122.4	21.4	80.7	21.9	50.9	65.9
	NLI	575	3,144	9.4	<b>36.7</b>	<b>28.0</b>	<b>57.7</b>	122.4	21.4	80.6	21.5	50.7	65.6
	Joint CE	857	3,120	9.4	35.4	27.6	56.0	111.8	20.5	80.2	21.8	51.2	65.2
	Only CE	878	3,126	9.4	34.2	27.3	55.5	109.2	20.1	80.0	21.8	49.9	64.5
Transformer	<b>Transformer RL</b>	753	3,433	9.2	<b>39.0</b>	28.7	<b>58.7</b>	127.7	22.5	81.3	26.6	56.2	70.5
	+ sFT (Ours)	1,458	3,959	9.1	36.9	28.2	57.2	118.7	21.7	<b>81.5</b>	30.6	62.3	75.7
	+ wFT (Ours)	<b>1,776</b>	<b>4,274</b>	9.1	31.3	26.2	53.0	103.1	20.0	81.2	<b>32.5</b>	<b>64.5</b>	<b>77.1</b>
	CIDErBtw	837	3,609	9.5	38.6	28.8	58.6	128.2	22.6	81.2	27.7	57.6	71.6
	NLI	876	3,744	9.5	38.9	28.9	58.5	<b>129.1</b>	<b>23.0</b>	<b>81.5</b>	29.8	59.9	73.4
	Joint CE	1,083	3,491	9.3	38.6	<b>29.0</b>	58.3	123.8	21.9	81.2	27.3	57.2	70.8
	Only CE	935	3,599	9.4	35.0	27.7	56.0	112.2	20.8	80.9	26.5	55.8	69.7

Table 1: Comparison with the baseline models and state-of-the-art distinctiveness-aware models. Automatic evaluation results on the MS COCO test set. *Unique-1* and *Unique-S* indicate the number of unique unigrams and sentences, respectively. *Length* is the average length of the output captions. Scores with † were reported by Liu et al. (2019). Other scores were reproduced by us. The results of our models are colored in gray.

Our fine-tuning was completed in approximately 10 minutes for each model using a single GPU of 16 GB memory.

## 4.2 Comparison with Baseline Models and Distinctiveness-Aware Models

Table 1 shows the results compared with those obtained with the baseline models and state-of-the-art distinctiveness-aware models.

**Vocabulary.** First, we observed that our methods (sFT and wFT) successfully increased the actual vocabulary size: both of them considerably increased Unique-1 compared to all the baseline models. wFT increased the vocabulary more than sFT, indicating that rebalancing the loss further encouraged tail-class word generation. The increased vocabulary resulted in the captions more specific to each image: Unique-S also increased significantly. Consistent with previous studies (Wang and Chan, 2019; Liu et al., 2019; Wang et al., 2020a), the models trained with the CE loss (Joint CE and Only CE) achieved the larger vocabulary than the baseline RL models. The improvement of our methods were even larger than these CE models. Despite the significant increase in the vocabulary size, our method kept the captions concise: the average sentence length was similar to that of the baseline models.

**Distinctiveness.** Our goal was to address the limited vocabulary of RL models in order to increase their distinctiveness. As expected, our methods increased distinctiveness compared to the baseline models: the R@K scores of our models were considerably higher than those of all the baseline models. Corresponding to the better improvement in vocabulary size, wFT increased distinctiveness more than sFT. These results confirm our hypothesis that the major bottleneck for distinctiveness is the limited vocabulary of RL models.

Among the Att2in-based models, Visual Paraphrase achieved the highest distinctiveness. However, this model is not directly comparable because it increases the parameters for its specialized model architecture. DiscCap performed comparably with our models, but its reward requires high computational costs. CIDErBtw and NLI proposed more lightweight rewards to be applicable to larger models, but they still need training from scratch. Among the larger models (UpDown and Transformer), our models achieved the highest distinctiveness despite the small computational cost.

**Standard Evaluation.** Our models degraded the performance in the text-based evaluation metrics (BLEU-4, METEOR, ROUGE-L, CIDEr, and SPICE), but they rather outperformed the baseline

	Distinctiveness	Correctness	Fluency
<b>Transformer RL</b>	3.00	4.42	4.83
+ wFT (Ours)	<b>3.34**</b>	4.45	<b>4.84</b>
NLI	3.18**	<b>4.54</b>	4.76

Table 2: Human evaluation results on the subset of the MS COCO test set. The distinctiveness score of Transformer RL was fixed at 3.00 because we set it as the baseline model. \*\* indicates a model scored higher than the baseline model with statistical significance (t-test with  $p < 0.01$ ): one-sample t-test for distinctiveness and independent two-sample t-test for the other criteria. The results of our model are colored in gray.

models in the text-and-image-based evaluation metric (RefCLIP). Because RefCLIP correlates with human evaluation better than text-based evaluation metrics (Hessel et al., 2021; Kasai et al., 2021), the high performance in RefCLIP demonstrates that our methods do not sacrifice the quality of captions. We found that our models tended to generate tail-class words that were correct but not covered by ground-truth captions, which unfairly lowered the scores in the text-based evaluation metrics. Appendix C shows those underrated captions.

### 4.3 Human Evaluation

We conducted human evaluations using Amazon Mechanical Turk (AMT) on three criteria: distinctiveness, correctness, and fluency. Correctness and fluency are absolute scores: we instructed workers to give a maximum score 5 to captions that *did not contain* incorrect information (ungrammatical or unnatural expressions) in terms of correctness (fluency). In contrast, distinctiveness is designed as a relative score because it is difficult to set an absolute standard for distinctiveness; unlike correctness or fluency, we cannot perfectly define distinctive captions across images. Following Wang et al. (2020a), we instructed the workers to determine the distinctiveness of a caption by comparing the caption with that of a baseline model<sup>10</sup>.

We evaluated the Transformer-based models, which performed the best in the automatic evaluation. We randomly selected 50 images from the MS COCO test set and assigned five workers to each image. See Appendix D for more details on the AMT instruction. Table 2 shows the re-

<sup>10</sup>If a target caption describes the same information as a baseline caption, the workers give the target caption a score of 3; if the target caption describes more (less) characteristic information than the baseline caption, the workers give the target caption a score of 4 or 5 (1 or 2).

sults. wFT, which had the highest R@K scores, also achieved the highest distinctiveness here. wFT did not achieve the highest correctness or fluency but achieved the same or higher correctness and fluency than the baseline model. This is consistent with the scores of RefCLIP, confirming again that our methods do not degrade the quality of captions.

### 4.4 Qualitative Analysis

Figure 1 shows the caption examples in the MS COCO validation set. The blue words are those that have never appeared in the output captions of the baseline model. The number of blue words indicates that our model successfully increased the vocabulary of the baseline model. Furthermore, we observed that these blue words were utilized in the description of the characteristic information of the images. See Appendix E for more examples.

## 5 Related Work

**Image Captioning** is the task of describing images in natural languages. The quality of captions has been remarkably improved by recent advances such as the encoder-decoder captioning model (Vinyals et al., 2015), attention mechanism (Xu et al., 2015), RL training (Ranzato et al., 2015; Rennie et al., 2017), attention over bounding box features (Anderson et al., 2018), large-scale pre-training (Li et al., 2020b), and large-scale captioning datasets (Young et al., 2014; Lin et al., 2014; Chen et al., 2015; Krishna et al., 2017; Sharma et al., 2018). Despite these advancements, current captioning models generate overly generic captions (Dai and Lin, 2017; Dai et al., 2017; Wang and Chan, 2019; Wang et al., 2020c).

**Distinctive Image Captioning** has been explored to generate more informative captions. Sadovnik et al. (2012) were the first to study it. They considered the more concise and more informative captions as those that describe information distinctive from *distractor images*, *i.e.*, images similar to an input image. Andreas and Klein (2016) proposed neural listener and speaker models that cooperate to generate distinctive captions for abstract scenes. Monroe et al. (2017) adapted the models to single-colored images. Vedantam et al. (2017) and Cohn-Gordon et al. (2018) extended the domain to real images and improved inference efficiency. Recently, Wang et al. (2021) proposed a memory attention network to describe objects that are unique among distractor images.



511 These approaches require selecting distractor im- 562  
512 ages for inference. Luo et al. (2018) and Liu et al. 563  
513 (2018) proposed the methods that do not require 564  
514 this step. Their models learn to generate distinctive 565  
515 captions by optimizing the R@K scores for 566  
516 sampled captions using RL (Rennie et al., 2017). 567  
517 The R@K scores are computed with a pre-trained 568  
518 image–text retrieval model (Faghri et al., 2018) 569  
519 over images in a mini batch. Vered et al. (2019) 570  
520 proposed a method to jointly train the image–text 571  
521 retrieval model and captioning model. Despite their 572  
522 effectiveness, R@K scores are associated with high 573  
523 computational costs and require a large batch size. 574  
524 Recently, Wang et al. (2020a) and Shi et al. (2021b) 575  
525 achieved state-of-the-art distinctiveness with more 576  
526 lightweight rewards. They weighted the contribu- 577  
527 tion of ground-truth captions for the CIDEr re- 578  
528 ward according to their differences from similar 579  
529 but different captions (Wang et al., 2020a) or their 580  
530 entailment scores against other ground-truth cap- 581  
531 tions (Shi et al., 2021b). Another approach ex- 582  
532 ploited unrelated captions as negative examples and 583  
533 trained caption generators with contrastive learn- 584  
534 ing (Dai and Lin, 2017) or GAN (Dai and Lin, 585  
535 2017; Goodfellow et al., 2014).

536 Liu et al. (2019) and Wu et al. (2021) are related 586  
537 to our work in that they exploited low-frequency n- 587  
538 grams to enhance distinctiveness. Liu et al. (2019) 588  
539 divided ground-truth captions into two subsets ac- 589  
540 cording to n-gram TF-IDF scores and proposed a 590  
541 new model architecture to paraphrase low TF-IDF 591  
542 captions into high TF-IDF ones. Wu et al. (2021) 592  
543 proposed the use of n-gram TF-IDF scores as an 593  
544 additional reward to a variant of R@K reward.

545 Different from above approaches, our objective 594  
546 is set to remedy the low distinctiveness of exist- 595  
547 ing RL models. Our models can be achieved with 596  
548 single-epoch fine-tuning of pre-trained RL mod- 597  
549 els, without requiring either drastic changes in the 598  
550 model architecture (Liu et al., 2019), additional 599  
551 computational cost of rewards (Wu et al., 2021), or 600  
552 training of a model from scratch.

553 **Diverse Image Captioning** is the task of gener- 601  
554 ating a set of diverse captions for a given image 602  
555 (Wang et al., 2016). Diverse image captioning 603  
556 is aimed at enumerating various pieces of infor- 604  
557 mation with a set of captions, whereas distinctive 605  
558 image captioning aims to concisely describe the 606  
559 most characteristic information with a single cap- 607  
560 tion. Similar to this study, some studies utilized 608  
561 captions that contained more tail-class words, such

as ground-truth captions (Wang and Chan, 2019; 562  
Luo and Shakhnarovich, 2020) or captions sampled 563  
from CE models (Shi et al., 2021a). Their models 564  
learn to generate these captions in addition to the 565  
captions sampled from RL models. However, these 566  
approaches still rely on sampling from skewed poli- 567  
cies and require training of a model from scratch. 568

**Long-Tail Classification** has been studied exten- 569  
sively in various tasks as label imbalance is preva- 570  
lent across datasets (Zhang et al., 2021a; Li et al., 571  
2020a). In text-generation tasks, label imbalance 572  
exists in the frequency of words. Previous ap- 573  
proaches have addressed this imbalance by nor- 574  
malizing classifier weights (Nguyen and Chiang, 575  
2018; Raunak et al., 2020) or using variants of Fo- 576  
cal loss (Raunak et al., 2020; Gu et al., 2020; Jiang 577  
et al., 2019; Wu et al., 2020; Lin et al., 2017). In 578  
contrast to these approaches, we adopted long-tail 579  
classification to mitigate the side effects of RL in 580  
the context of distinctive image captioning. We also 581  
tried these approaches and found that our methods 582  
performed the best. See Appendix F for the details. 583

## 6 Limitations and Risks 584

Our experiments were limited to the MS COCO 585  
captioning dataset, which is the standard dataset for 586  
image captioning. The images belong to the gener- 587  
al domain (real images of common objects) and 588  
the captions are in English only. The dataset con- 589  
tains social biases and captioning models have the 590  
risk of amplifying those biases (Zhao et al., 2021, 591  
2017; Hendricks et al., 2018). Our methods are not 592  
free from the risk too as they are not designed to 593  
reduce those biases from existing models. 594

## 7 Conclusion 595

In this study, we have investigated the problem of 596  
overly generic captions of RL captioning models 597  
with the hypothesis that their limited vocabulary is 598  
the major hindrance to distinctiveness. We recast 599  
distinctive image captioning as a simpler task of 600  
long-tail classification to increase the vocabulary 601  
and then propose lightweight fine-tuning methods 602  
to encourage tail-class word generation. The exper- 603  
imental results confirm our hypothesis by demon- 604  
strating that our methods significantly enhance the 605  
distinctiveness of existing RL models as well as 606  
their vocabulary size. Our methods also outper- 607  
form previous distinctiveness-aware methods with 608  
a small computational cost of minor modifications 609  
to pre-trained RL models. 610



611  
612  
613  
614  
  
615  
616  
617  
618  
619  
  
620  
621  
622  
  
623  
624  
625  
626  
  
627  
628  
629  
630  
631  
  
632  
633  
634  
  
635  
636  
637  
638  
  
639  
640  
641  
642  
  
643  
644  
645  
  
646  
647  
  
648  
649  
650  
  
651  
652  
653  
654  
  
655  
656  
657  
658  
  
659  
660  
661  
662

## References

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

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

Jacob Andreas and Dan Klein. 2016. Reasoning about pragmatics with neural listeners and speakers. In *EMNLP*.

Samy Bengio, Oriol Vinyals, Navdeep Jaitly, and Noam Shazeer. 2015. Scheduled sampling for sequence prediction with recurrent neural networks. In *NeurIPS*.

Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. 2015. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*.

Leshem Choshen, Lior Fox, Zohar Aizenbud, and Omri Abend. 2020. On the weaknesses of reinforcement learning for neural machine translation. In *ICLR*.

Christopher Clark, Mark Yatskar, and Luke Zettlemoyer. 2019. Don't take the easy way out: Ensemble based methods for avoiding known dataset biases. In *EMNLP-IJCNLP*.

Reuben Cohn-Gordon, Noah Goodman, and Christopher Potts. 2018. Pragmatically informative image captioning with character-level inference. In *NAACL-HLT*.

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

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

David Demeter, Gregory Kimmel, and Doug Downey. 2020. Stolen probability: A structural weakness of neural language models. In *ACL*.

Michael Denkowski and Alon Lavie. 2014. Meteor universal: Language specific translation evaluation for any target language. In *The Ninth Workshop on Statistical Machine Translation*.

Sergey Edunov, Myle Ott, Michael Auli, David Grangier, and Marc'Aurelio Ranzato. 2018. Classical structured prediction losses for sequence to sequence learning. In *NAACL-HLT*.

Fartash Faghri, David J Fleet, Jamie Ryan Kiros, and Sanja Fidler. 2018. Vse++: Improving visual-semantic embeddings with hard negatives. In *BMVC*.

Adam Fisch, Kenton Lee, Ming-Wei Chang, Jonathan H Clark, and Regina Barzilay. 2020. Capwap: Captioning with a purpose. In *EMNLP*.

Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. In *NeurIPS*.

Shuhao Gu, Jinchao Zhang, Fandong Meng, Yang Feng, Wanying Xie, Jie Zhou, and Dong Yu. 2020. Token-level adaptive training for neural machine translation. In *EMNLP*.

Danna Gurari, Yinan Zhao, Meng Zhang, and Nilavra Bhattacharya. 2020. Captioning images taken by people who are blind. In *ECCV*.

He He, Sheng Zha, and Haohan Wang. 2019. Unlearn dataset bias in natural language inference by fitting the residual. In *EMNLP-IJCNLP*.

Lisa Anne Hendricks, Kaylee Burns, Kate Saenko, Trevor Darrell, and Anna Rohrbach. 2018. Women also snowboard: Overcoming bias in captioning models. In *ECCV*.

Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. 2021. Clipscore: A reference-free evaluation metric for image captioning. In *EMNLP*.

Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.

Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The curious case of neural text degeneration. In *ICLR*.

Shaojie Jiang, Pengjie Ren, Christof Monz, and Maarten de Rijke. 2019. Improving neural response diversity with frequency-aware cross-entropy loss. In *WWW*.

Bingyi Kang, Saining Xie, Marcus Rohrbach, Zhicheng Yan, Albert Gordo, Jiashi Feng, and Yannis Kalantidis. 2020. Decoupling representation and classifier for long-tailed recognition. In *ICLR*.

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

Jungo Kasai, Keisuke Sakaguchi, Lavinia Dunagan, Jacob Morrison, Ronan Le Bras, Yejin Choi, and Noah A Smith. 2021. Transparent human evaluation for image captioning. *arXiv preprint arXiv:2111.08940*.

Samuel Kiegeland and Julia Kreutzer. 2021. Revisiting the weaknesses of reinforcement learning for neural machine translation. In *NAACL-HLT*.

Hyoungun Kim, Zineng Tang, and Mohit Bansal. 2020. Dense-caption matching and frame-selection gating for temporal localization in videoqa. In *ACL*.

716	Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, Michael Bernstein, and Li Fei-Fei. 2017. Visual genome: Connecting language and vision using crowdsourced dense image annotations. <i>IJCV</i> , 123(1):32–73.	Marc’Aurelio Ranzato, Sumit Chopra, Michael Auli, and Wojciech Zaremba. 2015. Sequence level training with recurrent neural networks. In <i>ICLR</i> .	768
717			769
718			770
719		Vikas Raunak, Siddharth Dalmia, Vivek Gupta, and Florian Metzger. 2020. On long-tailed phenomena in neural machine translation. In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> .	771
720			772
721			773
722			774
723	Xiaoya Li, Xiaofei Sun, Yuxian Meng, Junjun Liang, Fei Wu, and Jiwei Li. 2020a. Dice loss for data-imbalanced nlp tasks. In <i>ACL</i> .	Steven J Rennie, Etienne Marcheret, Youssef Mroueh, Jerret Ross, and Vaibhava Goel. 2017. Self-critical sequence training for image captioning. In <i>CVPR</i> .	776
724			777
725			778
726	Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. 2020b. Oscar: Object-semantic aligned pre-training for vision-language tasks. In <i>ECCV</i> .	Amir Sadovnik, Yi-I Chiu, Noah Snively, Shimon Edelman, and Tsuhan Chen. 2012. Image description with a goal: Building efficient discriminating expressions for images. In <i>CVPR</i> .	779
727			780
728			781
729			782
730			
731	Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In <i>Text Summarization Branches Out</i> .	Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. 2018. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In <i>ACL</i> .	783
732			784
733			785
734	Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. 2017. Focal loss for dense object detection. In <i>ICCV</i> .	Jiahe Shi, Yali Li, and Shengjin Wang. 2021a. Partial off-policy learning: Balance accuracy and diversity for human-oriented image captioning. In <i>ICCV</i> .	787
735			788
736			789
737	Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In <i>ECCV</i> .	Zhan Shi, Hui Liu, and Xiaodan Zhu. 2021b. Enhancing descriptive image captioning with natural language inference. In <i>ACL</i> .	790
738			791
739			792
740			
741	Lixin Liu, Jiajun Tang, Xiaojun Wan, and Zongming Guo. 2019. Generating diverse and descriptive image captions using visual paraphrases. In <i>ICCV</i> .	Matteo Stefanini, Marcella Cornia, Lorenzo Baraldi, Silvia Cascianelli, Giuseppe Fiameni, and Rita Cucchiara. 2021. From show to tell: A survey on image captioning. <i>arXiv preprint arXiv:2107.06912</i> .	793
742			794
743			795
744	Xihui Liu, Hongsheng Li, Jing Shao, Dapeng Chen, and Xiaogang Wang. 2018. Show, tell and discriminate: Image captioning by self-retrieval with partially labeled data. In <i>ECCV</i> .	Kaihua Tang, Jianqiang Huang, and Hanwang Zhang. 2020. Long-tailed classification by keeping the good and removing the bad momentum causal effect. In <i>NeurIPS</i> .	797
745			798
746			799
747			800
748	Ruotian Luo, Brian Price, Scott Cohen, and Gregory Shakhnarovich. 2018. Discriminability objective for training descriptive captions. In <i>CVPR</i> .	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In <i>NeurIPS</i> .	801
749			802
750			803
751	Ruotian Luo and Gregory Shakhnarovich. 2020. Analysis of diversity-accuracy tradeoff in image captioning. <i>arXiv preprint arXiv:2002.11848</i> .		804
752			
753			
754	Aditya Krishna Menon, Sadeep Jayasumana, Ankit Singh Rawat, Himanshu Jain, Andreas Veit, and Sanjiv Kumar. 2020. Long-tail learning via logit adjustment. In <i>ICLR</i> .	Ramakrishna Vedantam, Samy Bengio, Kevin Murphy, Devi Parikh, and Gal Chechik. 2017. Context-aware captions from context-agnostic supervision. In <i>CVPR</i> .	805
755			806
756			807
757			808
758	Will Monroe, Robert XD Hawkins, Noah D Goodman, and Christopher Potts. 2017. Colors in context: A pragmatic neural model for grounded language understanding. <i>TACL</i> , 5:325–338.	Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. 2015. Cider: Consensus-based image description evaluation. In <i>CVPR</i> .	809
759			810
760			811
761			
762	Toan Q Nguyen and David Chiang. 2018. Improving lexical choice in neural machine translation. In <i>NAACL-HLT</i> .	Gilad Vered, Gal Oren, Yuval Atzmon, and Gal Chechik. 2019. Joint optimization for cooperative image captioning. In <i>ICCV</i> .	812
763			813
764			814
765	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In <i>ACL</i> .	Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. 2015. Show and tell: A neural image caption generator. In <i>CVPR</i> .	815
766			816
767			817
		Jiuniu Wang, Wenjia Xu, Qingzhong Wang, and Antoni B Chan. 2020a. Compare and reweight: Distinctive image captioning using similar images sets. In <i>ECCV</i> .	818
			819
			820
			821

822 Jiuniu Wang, Wenjia Xu, Qingzhong Wang, and An-  
823 toni B Chan. 2021. Group-based distinctive image  
824 captioning with memory attention. In *ACM MM*.

825 Qingzhong Wang and Antoni B Chan. 2019. Describ-  
826 ing like humans: on diversity in image captioning.  
827 In *CVPR*.

828 Xudong Wang, Long Lian, Zhongqi Miao, Ziwei Liu,  
829 and Stella Yu. 2020b. Long-tailed recognition by  
830 routing diverse distribution-aware experts. In *ICLR*.

831 Zeyu Wang, Berthy Feng, Karthik Narasimhan, and  
832 Olga Russakovsky. 2020c. Towards unique and in-  
833 formative captioning of images. In *ECCV*.

834 Zhuhao Wang, Fei Wu, Weiming Lu, Jun Xiao, Xi Li,  
835 Zitong Zhang, and Yueting Zhuang. 2016. Diverse  
836 image captioning via grouptalk. In *IJCAI*.

837 Julia White, Gabriel Poesia, Robert Hawkins, Dorsa  
838 Sadigh, and Noah Goodman. 2021. Open-domain  
839 clarification question generation without question  
840 examples. In *EMNLP*.

841 Ronald J Williams. 1992. Simple statistical gradient-  
842 following algorithms for connectionist reinforce-  
843 ment learning. *Machine learning*, 8(3):229–256.

844 Jie Wu, Tianshui Chen, Hefeng Wu, Zhi Yang,  
845 Guangchun Luo, and Liang Lin. 2021. Fine-  
846 grained image captioning with global-local discrimi-  
847 native objective. *IEEE Transactions on Multimedia*,  
848 23:2413–2427.

849 Qingyang Wu, Lei Li, Hao Zhou, Ying Zeng, and Zhou  
850 Yu. 2020. Importance-aware learning for neural  
851 headline editing. In *AAAI*.

852 Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho,  
853 Aaron Courville, Ruslan Salakhudinov, Rich Zemel,  
854 and Yoshua Bengio. 2015. Show, attend and tell:  
855 Neural image caption generation with visual atten-  
856 tion. In *ICML*.

857 Peter Young, Alice Lai, Micah Hodosh, and Julia Hock-  
858 enmaier. 2014. From image descriptions to visual  
859 denotations: New similarity metrics for semantic in-  
860 ference over event descriptions. *TACL*, 2:67–78.

861 Yifan Zhang, Bingyi Kang, Bryan Hooi, Shuicheng  
862 Yan, and Jiashi Feng. 2021a. Deep long-tailed learn-  
863 ing: A survey. *arXiv preprint arXiv:2110.04596*.

864 Zhongping Zhang, Yiwen Gu, and Bryan A Plum-  
865 mer. 2021b. Show and write: Entity-aware news  
866 generation with image information. *arXiv preprint*  
867 *arXiv:2112.05917*.

868 Dora Zhao, Angelina Wang, and Olga Russakovsky.  
869 2021. Understanding and evaluating racial biases in  
870 image captioning. In *ICCV*.

871 Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Or-  
872 donez, and Kai-Wei Chang. 2017. Men also like  
873 shopping: Reducing gender bias amplification using  
874 corpus-level constraints. In *EMNLP*.



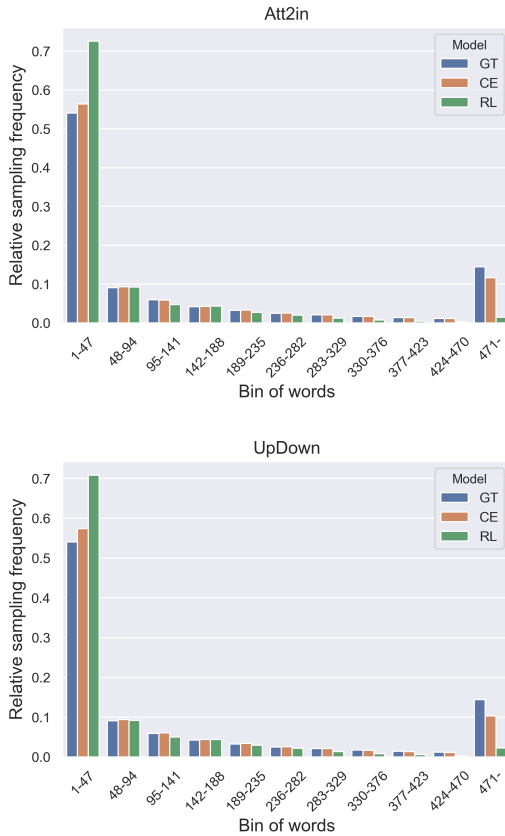


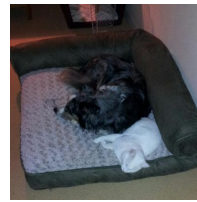
Figure 4: Relative frequency of the words in the sequences sampled for the training images. Five sequences were sampled for each image. The words (9,486 unique words excluding an out-of-vocabulary token  $\langle \text{unk} \rangle$ ) are sorted by their frequency in ground-truth captions and divided into 200 bins. We show the first 10 bins and the sum of the rest. GT is the ground-truth caption of the training images, CE is the output of a captioning model trained with the CE loss, and RL is the output of a captioning model trained with RL.

## A Peaky Distributions in Other Models

Figure 4 shows the results of the plotting in Figure 2 for the LSTM-based models: Att2in (Rennie et al., 2017) and UpDown (Anderson et al., 2018). Similar to the Transformer model, the sequences sampled with the LSTM-based RL models are clearly limited to head-class words, forming the peaky distributions.

## B Best Hyperparameters

As described in Section 4.1, we searched for the best hyperparameters for the learning rate (LR) from  $\{1e-3, 1e-4, 1e-5, 1e-6\}$ , and the inverse-temperature hyperparameter  $\beta$  of the fixed policy  $p_{\theta'}$  from  $\{0.1, 1\}$ . Note that sFT does not use  $p_{\theta'}$ . The best hyperparameters were as follows.



**Transformer RL:**  
a dog laying on top of a couch  
CIDEr: 133.3, RefCLIP: 77.7

**+wFT:**  
a dog curled up asleep on a cushion  
CIDEr: 38.7, RefCLIP: 79.2

**Human:**

an adorable dog laying down on a dog bed  
a dog and cat sleeping together on a dog bed  
a dog laying in a doggy bed with a cat  
a black down lounging on its pet bed  
black and white dog laying down on bed



**Transformer RL:**  
a person flying a kite in the ocean  
CIDEr: 60.2, RefCLIP: 79.8

**+wFT:**  
a man kiteboarding on top of a body of water  
CIDEr: 3.5, RefCLIP: 79.2

**Human:**

a person windsurfing with a grey sky in the background  
a person parasailing in the middle of the ocean  
a person riding a parachute surf board  
a man surfing alone on the ocean waters  
a man parasailing in the ocean all by himself



**Transformer RL:**  
a vase filled with yellow flowers on a table  
CIDEr: 216.7, RefCLIP: 78.5

**+wFT:**  
a clear vase filled with multi colored flowers  
CIDEr: 94.0, RefCLIP: 82.0

**Human:**

several flowers in a glass jar with water in it near a unpainted wall  
a vase filled with yellow and purple flowers  
some colorful flowers sitting on vase on the wall  
a vase of flowers on a table  
an arrangement of flowers in a clear glass canning jar haging on a wall

Figure 5: Underrated captions in the MS COCO validation set. The blue words are those that have never appeared in the output captions of the baseline model (Transformer RL). Human shows the full reference captions (ground-truth captions) of each image.

Att2in RL + sFT: LR =  $1e-4$ , 890

Att2in RL + wFT: LR =  $1e-4$ ,  $\beta = 1$ , 891

UpDown RL + sFT: LR =  $1e-4$ , 892

UpDown RL + wFT: LR =  $1e-4$ ,  $\beta = 1$ , 893

Transformer RL + sFT: LR =  $1e-5$ , 894

Transformer RL + wFT: LR =  $1e-5$ ,  $\beta = 0.1$ . 895

## C Examples of Underrated Captions

Figure 5 shows caption examples, reference captions, and their automatic evaluation scores. It is clear that our +wFT model correctly described all three images with diverse vocabulary. However, the CIDEr scores were quite low compared with those of the baseline model, Transformer RL. The cause 897  
898  
899  
900  
901  
902

Caption-A and Caption-B are the captions of the following image. Please rate the captions using the sliders below.



**Caption-A:** a cat laying on top of a red chair

**Caption-B:** a cat curled up asleep on a red chair

- How distinctive is **Caption-B**?
  - 5: **Caption-B** describes **more** characteristic information than **Caption-A**
  - 3: **Caption-B** describes **the same** information as **Caption-A**
  - 1: **Caption-B** describes **less** characteristic information than **Caption-A**

○

- How correct is **Caption-A**?
  - 5: Correct
  - 3: Slightly incorrect, but correct in the most salient contents
  - 1: Totally incorrect

○

- How correct is **Caption-B**?
  - 5: Correct
  - 3: Slightly incorrect, but correct in the most salient contents
  - 1: Totally incorrect

○

- How fluent is **Caption-A**?
  - 5: Fluent
  - 3: Slightly ungrammatical or unnatural, but understandable
  - 1: Totally ungrammatical or unnatural

○

- How fluent is **Caption-B**?
  - 5: Fluent
  - 3: Slightly ungrammatical or unnatural, but understandable
  - 1: Totally ungrammatical or unnatural

○

Submit

Figure 6: A screenshot of our AMT interface.

of this underrating is the small coverage of the reference captions: the reference captions rarely contain the tail-class words colored in blue, probably due to their low frequency. Text-based evaluation metrics such as CIDEr cannot evaluate the expressions that are correct but not covered by reference captions. In contrast, RefCLIP incorporates image features and can consider information that is not covered by reference captions. We observed that the RefCLIP scores were more plausible in these examples.

## D Details of Human Evaluation

We show our AMT interface in Figure 6. Each image was evaluated with the five questions in the discrete 5-point scale. We required workers to satisfy the following qualifications: being an AMT Master and living in the U.S. Workers were notified that this experiment was intended to evaluate caption quality. We paid \$0.1 for each image, and the median of the actual working time was 41 seconds per image. The hourly reward was estimated as






(a)		<p><b>Transformer RL:</b> a tower with a clock on top of it</p> <p><b>+wFT:</b> a clock tower with a <b>weather vane</b> on top</p> <p><b>NLI:</b> a tower with a clock on the top of it</p> <p><b>Human:</b> a weather vane atop a cathedral clock tower</p>
(b)		<p><b>Transformer RL:</b> a group of birds standing in the water</p> <p><b>+wFT:</b> a large group of <b>flamingos</b> stand in <b>shallow</b> water</p> <p><b>NLI:</b> a group of pink umbrellas are standing in the water</p> <p><b>Human:</b> a flock of pink flamingos standing in shallow water</p>
(c)		<p><b>Transformer RL:</b> a black cat wearing a hat on top of a table</p> <p><b>+wFT:</b> a cat wears a <b>funny</b> hat while <b>staring straight ahead</b></p> <p><b>NLI:</b> a black cat wearing a hat sitting on a table</p> <p><b>Human:</b> the cute black cat is wearing a bee's hat</p>
(d)		<p><b>Transformer RL:</b> a group of people riding motorcycles on a road</p> <p><b>+wFT:</b> a group of people <b>racing</b> motorcycles on a race track</p> <p><b>NLI:</b> a group of people riding motorcycles on a race track</p> <p><b>Human:</b> people are racing motorcycles on a race track</p>
(e)		<p><b>Transformer RL:</b> a dog next to a cup of coffee</p> <p><b>+wFT:</b> a dog is <b>sniffing</b> a cup of coffee</p> <p><b>NLI:</b> a dog standing next to a coffee cup on a table</p> <p><b>Human:</b> a squinting dog on a brick patio sniffs a cup of coffee</p>

Figure 7: Caption examples in the MS COCO validation set. The blue words are those that have never appeared in the output captions of the baseline model (Transformer RL). *Human* shows a ground-truth caption of each image.

\$8.78, which is higher than the minimum wage in the U.S., \$7.25 per hour.

## E Detailed Qualitative Analysis

Figure 7 shows more caption examples in the MS COCO validation set. The blue words are those that have never appeared in the output captions of the baseline model. We observed that these blue words expressed various types of characteristic information of the images. Here, *weather vane* and *flamingos* are characteristic objects of the images (a) and (b); *shallow*, *funny*, and *staring straight ahead* are characteristic attributes of the images (b) and (c); and *racing* and *sniffing* are characteristic relations in the images (d) and (e). These examples further support our hypothesis that the limited vocabulary of RL models hinders their distinctiveness.

	Vocabulary			Standard Evaluation						Distinctiveness		
	Unique-1	Unique-S	Length	BLEU-4	METEOR	ROUGE-L	CIDEr	SPICE	RefCLIP	R@1	R@5	R@10
<b>Att2in RL</b>	445	2,524	9.3	35.3	<b>27.1</b>	<b>56.7</b>	<b>117.4</b>	<b>20.5</b>	79.7	16.3	41.9	57.2
+ sFT (Ours)	880	3,156	9.0	<b>35.6</b>	27.0	56.5	115.4	20.4	<b>80.3</b>	20.1	48.0	62.8
+ wFT (Ours)	<b>1,091</b>	<b>3,749</b>	9.0	32.6	26.4	54.9	108.6	19.9	<b>80.3</b>	<b>21.7</b>	<b>50.8</b>	<b>65.2</b>
+ $\tau$ -norm	437	2,414	9.1	35.4	27.0	<b>56.7</b>	117.3	20.4	79.7	15.4	40.7	55.8
+ FL	902	3,236	9.0	35.2	27.0	56.4	114.6	20.3	<b>80.3</b>	20.3	48.4	63.4
+ AFL	886	3,104	9.0	35.4	27.0	56.6	115.2	20.4	<b>80.3</b>	19.6	47.5	62.7
<b>UpDown RL</b>	577	3,103	9.5	36.7	<b>27.9</b>	57.6	<b>122.7</b>	<b>21.5</b>	80.5	21.1	49.9	64.6
+ sFT (Ours)	1,190	3,788	9.2	35.7	27.5	56.5	115.9	21.0	<b>80.9</b>	25.0	56.8	71.2
+ wFT (Ours)	<b>1,227</b>	<b>4,263</b>	9.3	32.0	26.5	54.3	107.9	20.4	<b>80.9</b>	<b>25.5</b>	<b>58.0</b>	<b>72.6</b>
+ $\tau$ -norm	576	2,967	9.3	<b>37.0</b>	27.7	<b>57.7</b>	122.6	21.3	80.5	19.6	48.1	63.4
+ FL	1,208	3,838	9.2	35.4	27.4	56.3	114.8	20.8	<b>80.9</b>	25.3	57.2	71.0
+ AFL	1,168	3,746	9.2	35.9	27.5	56.7	116.4	20.9	<b>80.9</b>	24.7	56.5	70.5
<b>Transformer RL</b>	753	3,433	9.2	<b>39.0</b>	<b>28.7</b>	<b>58.7</b>	<b>127.7</b>	<b>22.5</b>	81.3	26.6	56.2	70.5
+ sFT (Ours)	1,458	3,959	9.1	36.9	28.2	57.2	118.7	21.7	81.5	30.6	62.3	75.7
+ wFT (Ours)	<b>1,776</b>	<b>4,274</b>	9.1	31.3	26.2	53.0	103.1	20.0	81.2	<b>32.5</b>	<b>64.5</b>	<b>77.1</b>
+ $\tau$ -norm	1,027	3,483	9.2	38.5	28.4	58.3	124.4	22.1	81.2	26.1	55.8	69.7
+ FL	1,523	4,018	9.1	36.1	28.0	56.6	116.5	21.4	81.5	31.2	63.1	76.3
+ AFL	1,402	3,908	9.1	37.4	28.3	57.5	120.5	21.9	<b>81.6</b>	30.0	62.1	75.9

Table 3: Comparison with the other long-tail classification methods. Automatic evaluation results on the MS COCO test set. *Unique-1* and *Unique-S* indicate the number of unique unigrams and sentences, respectively. *Length* is the average length of output captions. The results of our models are colored in gray.

## F Comparison with Other Long-Tail Classification Methods

We adapted the long-tail classification method of Kang et al. (2020) to remedy the side effects of RL and proposed sFT and wFT. Both methods were carefully designed for RL models, but these were not the only way to employ long-tail classification methods. In this section, we discuss the other possible adaptations based on Raunak et al. (2020).

Raunak et al. (2020) explored ways to employ long-tail classification methods to machine translation. Their first method was  $\tau$ -normalization ( $\tau$ -norm), which directly adopted the method of Kang et al. (2020). To simply make the output distributions flatter, they normalized the classifier weight  $\mathbf{W}$  as follows:

$$\widetilde{\mathbf{W}}_{w_i} = \frac{\mathbf{W}_{w_i}}{\|\mathbf{W}_{w_i}\|^\tau}, \quad (10)$$

where  $\mathbf{W}_{w_i} \in \mathbb{R}^d$  indicates a vector at the index of a word  $w_i$  and  $\tau$  is a temperature hyperparameter that controls the degree of the normalization.

The other methods of Raunak et al. (2020) were Focal loss (FL) and Anti-Focal loss (AFL). AFL is a variant of FL (Lin et al., 2017), which was aimed at reweighting the loss according to the confidence of the model predictions. Let  $p_\theta^t = p_\theta(w_i^g | w_{<t}^g, I)$ . FL and AFL in image captioning are then

written as follows:

$$\mathcal{L}_{\text{FL}}(\theta) = -\frac{1}{T} \sum_{t=1}^T (1 - p_\theta^t)^\gamma \log p_\theta^t, \quad (11)$$

$$\mathcal{L}_{\text{AFL}}(\theta) = -\frac{1}{T} \sum_{t=1}^T (1 + \alpha p_\theta^t)^\gamma \log p_\theta^t, \quad (12)$$

where  $\gamma$  and  $\alpha$  are hyperparameters that control the degree of the reweighting. Other work also explored ways to employ long-tail classification methods to text generation, but those approaches can be categorized as either  $\tau$ -norm (Nguyen and Chiang, 2018) or the variants of FL (Gu et al., 2020; Jiang et al., 2019; Wu et al., 2020), which we already explored above.

We compared our methods (sFT and wFT) with  $\tau$ -norm, FL, and AFL. In our experiments, we normalized the bias term  $\mathbf{b}^{11}$  in addition to the weight term  $\mathbf{W}$  as we found it performed better than normalizing the weight term only. For a fair comparison with our methods, we applied FL and AFL at the fine-tuning of RL models. That is, we optimized  $\mathcal{L}_{\text{FL}}(\hat{\theta})$  and  $\mathcal{L}_{\text{AFL}}(\hat{\theta})$ , where  $\hat{\theta}$  were initialized with the pre-trained RL models. We used the best hyperparameters reported in Raunak et al. (2020):  $\tau = 0.2$ ,  $\gamma = 1$ , and  $\alpha = 1$ . Similar to our models, other hyperparameters were set to the same values as the baseline models, except for the epoch size and learning rate. We explored the same values

<sup>11</sup> $\widetilde{\mathbf{b}} = \frac{\mathbf{b}}{\|\mathbf{b}\|^\tau}$ , where the value of the hyperparameter  $\tau$  was set to the same as that of  $\widetilde{\mathbf{W}}$ .



	Vocabulary			Standard Evaluation						Distinctiveness		
	Unique-1	Unique-S	Length	BLEU-4	METEOR	ROUGE-L	CIDEr	SPICE	RefCLIP	R@1	R@5	R@10
<b>Att2in RL</b>	435	2,583	9.3	<b>35.0</b>	<b>27.0</b>	<b>56.7</b>	<b>116.5</b>	<b>20.3</b>	79.8	16.2	42.5	57.0
+ sFT (Ours)	874	3,189	9.0	<b>35.0</b>	<b>27.0</b>	56.3	113.7	20.1	80.3	19.2	47.9	62.9
+ wFT (Ours)	<b>1,092</b>	<b>3,806</b>	9.0	32.5	26.4	54.8	107.2	19.7	<b>80.4</b>	<b>20.6</b>	<b>50.7</b>	<b>65.3</b>
<b>UpDown RL</b>	563	3,161	9.5	<b>36.7</b>	<b>27.9</b>	<b>57.7</b>	<b>122.3</b>	<b>21.3</b>	80.6	20.6	50.2	65.7
+ sFT (Ours)	1,222	3,805	9.2	35.4	27.4	56.7	115.3	20.7	80.9	24.6	56.2	70.9
+ wFT (Ours)	<b>1,230</b>	<b>4,311</b>	9.3	31.8	26.4	54.3	106.5	20.2	<b>81.0</b>	<b>25.9</b>	<b>58.2</b>	<b>73.6</b>
<b>Transformer RL</b>	713	3,432	9.2	<b>38.9</b>	<b>28.7</b>	<b>58.7</b>	<b>126.4</b>	<b>22.1</b>	81.2	25.4	56.3	69.8
+ sFT (Ours)	1,496	3,953	9.1	37.5	28.3	57.4	118.4	21.4	<b>81.5</b>	30.2	62.7	75.8
+ wFT (Ours)	<b>1,836</b>	<b>4,268</b>	9.1	31.1	26.3	53.3	102.2	19.8	81.3	<b>32.2</b>	<b>64.3</b>	<b>76.8</b>

Table 4: Automatic evaluation results on the MS COCO *validation* set. *Unique-1* and *Unique-S* indicate the number of unique unigrams and sentences, respectively. *Length* is the average length of output captions. The results of our models are colored in gray.

for these hyperparameters as our models: we set the epoch size for fine-tuning to 1 and searched for the best learning rate from  $\{1e-3, 1e-4, 1e-5, 1e-6\}$ . The best learning rates were chosen according to the R@1 scores in the validation set<sup>12</sup>. Note that we did not explore the learning rate for  $\tau$ -norm because it does not require training.

Table 3 shows the results. In contrast to the results reported in machine translation (Raunak et al., 2020; Nguyen and Chiang, 2018),  $\tau$ -norm models performed lower than the baseline models. These results indicate that simply flattening the output distributions does not work in image captioning. Although FL and AFL increased the vocabulary size and distinctiveness, the gains were smaller than those of wFT.

To analyze the cause of the difference between the FL, AFL, and the BP loss (wFT), we visualized the losses in Figure 8. FL suppresses the loss when a model is confident, whereas AFL increases the loss when a model is moderately confident. Compared with these losses, BP changes the loss more drastically. When the head-class-biased policy  $p_{\theta'}$  is highly confident, BP strictly suppresses the loss to prevent further learning on that word; when  $p_{\theta'}$  is not confident, BP highly increases the loss to encourage the learning on that word. This drastic rebalancing of the loss resulted in the larger vocabulary size and higher distinctiveness of wFT.

## G Validation Performance for Reproduction

Table 4 shows the performance of our models on the MS COCO validation set. We report these results for the future reproduction of our experiments.

<sup>12</sup>The best learning rates were  $1e-4$  for Att2in RL + FL/AFL,  $1e-4$  for UpDown RL + FL/AFL, and  $1e-5$  for Transformer RL + FL/AFL.

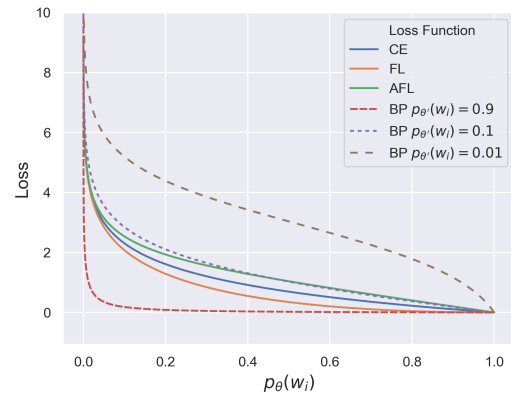


Figure 8: Visualization of the losses:  $CE - \log p_{\theta}(w_i)$ ,  $BP - \log p_{\theta, \theta'}(w_i)$ ,  $FL (1 - p_{\theta}(w_i))^{\gamma} \log p_{\theta}(w_i)$ , and  $AFL (1 + \alpha p_{\theta}(w_i))^{\gamma} \log p_{\theta}(w_i)$ . Here, we set  $\beta = 1$ ,  $\gamma = 1$ , and  $\alpha = 1$ .

The code will be also available on our website for the reproduction.

1025

1026