

Prefix-diffusion: A Lightweight Diffusion Model for Diverse Image Captioning

Anonymous ACL submission

Abstract

While impressive performance has been achieved in image captioning, the limited diversity of the generated captions and the large parameter scale remain major barriers to the real-world application of these systems. In this work, we propose a lightweight image captioning network in combination with continuous diffusion, called Prefix-diffusion. To achieve diversity, we design an efficient method that injects prefix image embeddings into the denoising process of the diffusion model. In order to reduce trainable parameters, we employ a pre-trained model to extract image features and further design an extra mapping network. Prefix-diffusion is able to generate diverse captions with relatively less parameters, while maintaining the fluency and relevance of the captions benefiting from the generative capabilities of the diffusion model. Our work paves the way for scaling up diffusion models for image captioning, and achieves promising performance compared with recent approaches.¹

1 Introduction

Image captioning, which combines computer vision (CV) and natural language processing (NLP), focuses mainly on producing a description of an image. Existing works on image captioning typically employ an encoder-decoder architecture (Vinyals et al., 2015; Anderson et al., 2018; Zhou et al., 2020) to generate captions word-by-word. However, such models require large trainable parameters to bridge the visual and textual representations. By utilizing the powerful representation capability of pre-trained models like CLIP (Radford et al., 2021), recent methods (Lovenia et al., 2022; Zhu et al., 2022; Mokady et al., 2021) map visual semantic information to language space for image captioning. Although autoregressive models have become the typical approach for image captioning,

¹Code will be released upon publication.



Figure 1: The diverse captions generated by Prefix-diffusion. The model is trained on the COCO dataset. More examples will be given in the supplementary material.

their left-to-right generative manner leads to cumulative errors. Moreover, human-like captions not only maintain fluency and relevance properties, but also contain diverse wordings and rich expressions.

Recently, the popular diffusion model (Sohl-Dickstein et al., 2015), which generates samples through an iterative denoising process, has provided a promising path to generate tokens in parallel and inherently increase the diversity of captions. Diffusion models (Sohl-Dickstein et al., 2015) have become an active area of research owing to their ability to generate comparable results with GANs (Goodfellow et al., 2020) on computer vision tasks. The strength of diffusion models trained on vast image databases has led to an almost ubiquitous fascination among researchers in producing highly typical content, such as image generation and editing (Nichol et al., 2021; Balaji et al., 2022; Kim et al., 2022; Gal et al., 2022). Nevertheless, the path is blocked by the discreteness of texts and the gap between different modals.

For the continuous diffusion models (Ho et al.,

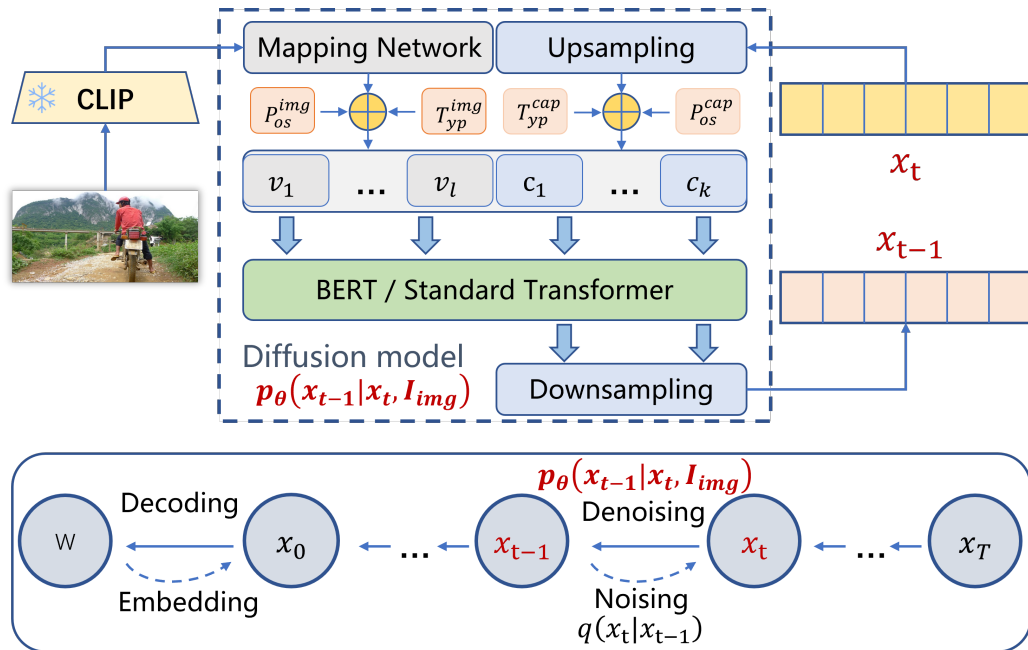


Figure 2: Illustration of Prefix-diffusion. The bottom lies the diffusion process. The reverse process is defined by $p_\theta(x_{t-1} | x_t, I_{img})$ and the diffusion model is depicted in the upper dashed box. We use the frozen CLIP to extract image features and train a lightweight mapping network to connect the image space and the text space.

2020; Nichol and Dhariwal, 2021; Song et al., 2020), they only work on continuous data but yield inferior results in generating text and image captioning, especially compared to the results of the autoregressive models. To effectively benefit from continuous diffusion, Diffusion-LM (Li et al., 2022) extends the the standard diffusion process with an embedding step followed by a rounding step, generating the high-quality text under six control targets. The discreteness of texts has been overcome, whereas the gap between different modals stays unsolved. For image captioning with continuous diffusion, it is a more challenging task, which further requires the fusion of the image information.

In this paper, we propose a lightweight captioning model based on the continuous diffusion, namely Prefix-diffusion. The model tackles three key problems in image caption generation. Firstly, we utilize diffusion models to solve the limited diversity of the generated captions. Noticing that diffusion models have the powerful generative capabilities but few research applied them to image captioning. Secondly, different from image captioning models that have a large number of parameters and are computationally expensive, our framework saves computing resources with the pre-trained CLIP model to extract image features. Last

but not least, our method is able to generate more accurate captions in parallel, since it injects prefix image embeddings into the denoising process of the diffusion model. This essentially solves the problem of sequential error accumulation.

Figure 1 shows the captions generated by Prefix-diffusion, where the captions accurately describe the content of the image with fluency. Different from the method of beam search, our method can cover all distributions of the training datasets and generate diverse captions.

The overall contributions of our work are:

- We propose a lightweight method Prefix-diffusion to generate diverse captions. Our work tackles the multi-modal issue for the diffusion model and paves the way for scaling it up for image captioning.
- Prefix-diffusion generates diverse captions in a variety of forms, which is specifically reflected in the increase of Dist-3 and vocabulary usage by 6.3 and 3.1 compared with the baselines, respectively.
- Prefix-diffusion reduces more than 38% trainable parameters compared with existing CLIP-based methods(Nukrai et al., 2022; Mokady et al., 2021), while achieving comparable or even better results in newer metrics.

2 Related Work

2.1 Image Captioning

The autoregressive models achieve promising performance on image captioning. The next token of the caption is conditioned on the former tokens. To generate more neural captions, (Lu et al., 2018) predicts the slot locations that are explicitly tied to image regions. GET (Ji et al., 2021) captures a more comprehensive global representation by using a novel transformer architecture, to guide the caption generation. Similarly, (Li et al., 2019; Luo et al., 2021) use transformer to leverage the image information efficiently. Thanks to the powerful multi-modal representation capability of CLIP (Radford et al., 2021), (Mokady et al., 2021; Galatolo et al., 2021) take an image embedding as the input which is encoded by the CLIP visual encoder. Then they use the GPT-2 (Radford et al., 2019) model to produce a sequence of words that describe the content of the input image. But autoregressive models suffer from the limitation of generation speed and the accumulation of errors.

Non-autoregressive models have recently attracted attention due to their fast inference speed and generation quality. (Gao et al., 2019) randomly masks the input sequences with certain ratios to train a masked language model, and generates captions parallelly during inference. Considering non-autoregressive image captioning as a cooperative multi-agent problem, (Guo et al., 2020) proposes a novel counterfactuals-critical multi-agent learning algorithm to improved the inference speed. (Fei, 2020) proposes a non-autoregressive image captioning approach based on the idea of iterative back modification, which refines the output in a limited number of steps. To determine the length of the image caption, (Deng et al., 2020) designs a non-autoregressive decoder for length-controllable image captioning.

2.2 Diffusion Model

Diffusion models (Sohl-Dickstein et al., 2015) have demonstrated impressive capabilities in creative applications. For text-to-image generation, a task of generating a corresponding image from a description, (Balaji et al., 2022; Nichol et al., 2021; Rombach et al., 2022; Gu et al., 2022) apply discrete diffusion models to produce high-resolution images conditioned on the text prompts. DiffSound (Yang et al., 2022) proposes a novel decoder based on the diffusion model to generate high-quality

sound. Similarly, ProDiff (Huang et al., 2022) studies on diffusion parameterization for text-to-speech and achieves superior sample quality and diversity. In the text generation domain, Diffusion-LM (Li et al., 2022) starts with a sequence of Gaussian noise vectors and denoises them incrementally into vectors corresponding to words. Diffusion-LM enables efficient gradient-based methods for controllable generation, achieving promising results in the new forms of complex fine-grained control tasks. Moreover, (Gong et al., 2022; Strudel et al., 2022) extend vanilla diffusion models to learn conditional text generation. However, few research applies the diffusion model to image captioning, because of the cross-modal challenge and the discreteness of texts.

3 Methodology

As illustrated in Figure 2, we propose Prefix-diffusion for injecting image features to learn image captioning. Different from image generating, our method requires to map discrete texts to a continuous space by a word embedding. For the conditioned image, we first extract its features by the CLIP image encoder, and then input them to the mapping network to obtain the prefix image embeddings. We then concatenate the prefix image embeddings and the caption embeddings in the denoising process of the diffusion model. The concatenated vectors are fed into a deep neural network (e.g. BERT (Kenton and Toutanova, 2019) or the standard transformer). Since our work merely trains a mapping network and a neural network, the trainable parameter scale is reduced significantly.

Forward process. Following Diffusion-LM (Li et al., 2022), we adopt an embedding function $EMB(W)$ to map a discrete word into a continuous space. Define a caption W with k words. Through the embedding function, we have $EMB(W) = [EMB(\omega_1), \dots, EMB(\omega_k)] \in \mathbb{R}^{k \times d_1}$, where d_1 is the dimension of the vector. In our experiments, we find that the value of d_1 works well at 48. Reducing the dimension will decrease the performance, while increasing the dimension will enlarge the computational burden.

For the forward process, diffusion models (Ho et al., 2020; Nichol and Dhariwal, 2021; Song et al., 2020) add noise progressively to training a sample according to a variance schedule β_1, \dots, β_T . The forward process has no learnable parameters and

we get x_t by the following equation:

$$x_t = \sqrt{1 - \beta_t}x_{t-1} + \sqrt{\beta_t}\epsilon \quad (1)$$

where $\epsilon \sim N(0, 1)$ and $\beta_t : 0.01 \rightarrow 0.03$ are hyperparameters representing the variance schedule across diffusion steps. We have tried different noise methods, with the truncation linear noise schedule method being the best. We validate this observation in section 4.3.3.

Reverse process. The reverse process generates new samples from $x_T \sim N(0, I)$. The data is sampled using the following reverse diffusion process:

$$p_\theta(x_{t-1} | x_t, I_{img}) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, I_{img}), \sigma(t)^2 I) \quad (2)$$

where I_{img} denotes the visual information from CLIP.

In order to learn the reverse process, neural networks are trained to predict μ_θ and $\sigma(t)^2$ is a fixed variance.

$$\mu_\theta(x_t, I_{img}) = \frac{\sqrt{\bar{\alpha}_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}x_t + \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t}x_0 \quad (3)$$

$$\sigma(t)^2 = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}\beta_t. \quad (4)$$

In order to get μ_θ , we compute x_0 with the following equation:

$$x_0 = \frac{1}{\sqrt{\bar{\alpha}_t}}(x_t - \sqrt{1 - \bar{\alpha}_t}\tilde{z}) \quad (5)$$

where \tilde{z} can be obtained by deep neural networks (e.g. transformer).

$$\tilde{z} = \Phi(x_t, I_{img}, t). \quad (6)$$

Here Φ denotes the neural network which is depicted in the dashed box in the Figure 2. Since the transformer architecture has been shown to outperform many other architectures on a wide range of text generation tasks, we explored two different transformer architectures as the neural network: BERT and the standard transformer. Different from other continuous diffusion approaches, we inject image features into the transformer architectures. This process changes the original mean in the caption space, as illustrated in Figure 3.

In the following, we will explain in detail how to inject the image information into the model. Firstly we use CLIP to encode image and receive its image

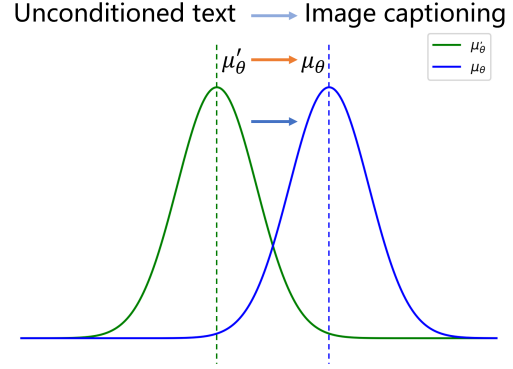


Figure 3: After we concatenate the image features in the reverse process, the original mean μ'_θ is changed to μ_θ in the caption space. Hence, the unconditioned text is converted to an image caption.

features I'_{img} . Then we train a mapping network F on I'_{img} and obtain the visual prefix I^m_{img} of length l :

$$\begin{cases} I_{img} = CLIP(image) \\ I^m_{img} = \{v'_1, v'_2, \dots, v'_l\} = F(I_{img}) \end{cases} \quad (7)$$

We specifically formulate $I^m_{img} \in \mathbb{R}^{l \times d_2}$ as $\{v'_1, v'_2, \dots, v'_l\}$ for the convenience of subsequent expression. To save the computation cost, we employ a simple Multi-Layer Perceptron (MLP) as the mapping network. Through an upsampling network, a sequence embedding x_t has the same dimension as I^m_{img} , denoted as $\{c'_1, c'_2, \dots, c'_k\} \in \mathbb{R}^{k \times d_2}$. k is the length of the caption and d_2 is the dimension of the embedding.

Before concatenating the visual prefix embedding and the caption embedding, we add positional embedding P_{os} and type embedding T_{yp} to it:

$$\{c_1, c_2, \dots, c_k\} = \{c'_1, c'_2, \dots, c'_k\} + P_{os}^{cap} + T_{yp}^{cap} \quad (8)$$

$$\{v_1, v_2, \dots, v_l\} = \{v'_1, v'_2, \dots, v'_l\} + P_{os}^{img} + T_{yp}^{img}. \quad (9)$$

The positional embedding indicates the model where the feature is located, which is essential information. Similarly, the type embedding tells the model where the image features lie. Then the visual prefix and the caption embedding are concatenated into a sequence $\{v_1, \dots, v_l, t_1, \dots, t_k\}$, and processed by a standard transformer or BERT

Method	Common Metrics \uparrow						Similarity Score \uparrow			Diversity \uparrow			
	B@1	B@3	M	R-L	C	S	CLIP-S	Ref-CLIP	P-Bert	D@2	D@3	Voc-u	
MTIC	80.8	50.9	29.2	58.6	131.2	22.6	60.3	68.6	94.0	7.9	16.3	8.3	
DLCT	81.1	51.1	29.4	58.9	133.1	22.8	60.6	69.0	94.1	8.1	17.1	8.3	
Frozen Clip Feature	CapDec	68.3	36.6	25.2	51.2	91.7	18.3	60.4	67.8	93.4	8.3	14.9	1.9
	ClipCap	73.6	42.3	26.7	54.4	105.8	<u>19.8</u>	60.8	68.6	93.8	<u>11.3</u>	21.7	2.6
	Ours(T)	<u>77.7</u>	<u>43.4</u>	25.8	<u>55.8</u>	<u>106.3</u>	19.4	<u>63.4</u>	<u>70.9</u>	93.2	11.2	<u>25.9</u>	<u>4.7</u>
	Ours(B)	78.1	44.2	<u>26.6</u>	56.1	109.3	20.4	63.7	71.2	<u>93.7</u>	12.7	28.0	5.7

Table 1: The results of image captioning on COCO. For all the metrics, the higher the better. We use boldface to indicate the best performance. The second best result is underlined. Ours(T) and Ours(B) use a standard transformer and BERT respectively. The values of vocabulary usage are reported at percentage (%).

Method	Common Metrics \uparrow						Similarity Score \uparrow			Diversity \uparrow		
	B@1	B@3	M	R-L	C	S	CLIP-S	Ref-CLIP	P-Bert	D@2	D@3	Voc-u
CapDec	57.6	27.9	20.0	44.5	42.0	14.3	58.0	61.4	<u>92.8</u>	15.5	25.2	1.3
ClipCap	67.0	<u>35.2</u>	22.5	<u>49.0</u>	<u>60.8</u>	16.5	60.9	65.0	93.0	20.9	34.5	1.77
Ours(T)	<u>68.7</u>	34.9	20.1	48.7	53.8	14.2	<u>61.6</u>	<u>66.3</u>	92.2	<u>23.1</u>	<u>41.0</u>	<u>3.6</u>
Ours(B)	71.0	36.2	<u>21.1</u>	49.3	61.4	<u>15.2</u>	64.7	68.6	92.0	27.6	46.0	4.0

Table 2: The results of image captioning on Flickr30k. For all the metrics, the higher the better. We use boldface to indicate the best performance. The second best result is underlined.

network:

$$\{y_1, y_2, \dots, y_l, y_{l+1}, \dots, y_{l+k}\} = \text{Network}(\text{concat}(v_1, \dots, v_l, c_1, \dots, c_k)). \quad (10)$$

We split y_i and use $\{y_{l+1}, \dots, y_{l+k}\}$ as the input of the downsampling, yielding the output $x_{t-1} \in \mathbb{R}^{k \times d_1}$ of the diffusion model.

Decoding process. In the decoding process, we strengthen the similarity of images and captions with CLIP scores. The benefit of CLIP in the current work is that it can provide a cosine similarity score between numerous texts and an image. Utilizing the CLIP embedding of an image, we calculate the cosine similarity between the image and the n candidate captions. We then choose the most relevant captions. The similarity is computed as follows:

$$\text{similarity}(I_{img}, W_{txt}^n) = \frac{I_{img} \cdot W_{txt}^n}{|I_{img}| \cdot |W_{txt}^n|} \quad (11)$$

where I_{img} is the image features extracted by CLIP and W_{txt}^n is the features of the n candidate captions. This is a retrieval-base (Ramos et al., 2022; Zhao et al., 2020) technique that picks the best appropriate caption from a set of candidate captions. We use this approach based on the advantage of Prefix-diffusion: our model can generate diverse captions with different Gaussian noises. We verify

the effectiveness of this retrieval-base method in section 4.3.3.

4 Experiment

In this section, we conduct quantitative and qualitative experiments to evaluate our approach. We first introduce the implementation details in subsection 4.1 and 4.2. Then we compare the performance of our approach with the others on various evaluation metrics (subsection 4.3.1 and 4.3.2). Finally, the ablation experiments (subsection 4.3.3) are also presented to analyze the significance of our design.

4.1 Dataset and Evaluation Metric

We use COCO (Lin et al., 2014) and Flickr30k (Plummer et al., 2015) as the datasets for image captioning. We split the datasets for training, validation, and test according to the Karpathy et al (Karpathy and Fei-Fei, 2015), where the test sets of the two datasets contain 5000 images and 1000 images respectively. To evaluate the generalization ability of our model, we train the model on one dataset while evaluating on the other.

In this paper, we adopt automatic evaluation to appraise the generated captions. In addition to the common metrics and similarity, we consider two metrics to evaluate the diversity of the generated captions.

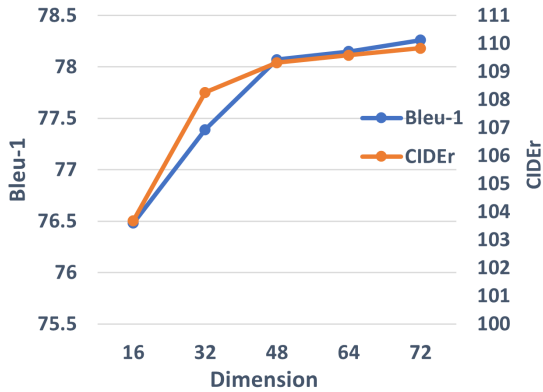


Figure 4: The performance effect of the word dimension on COCO. We report the metrics of Bleu-1 and CIDEr.

- *Common Metrics.* Following the common practice in the literatures, we perform evaluation using BLEU(B@N)(Papineni et al., 2002), METEOR(M)(Denkowski and Lavie, 2014), ROUGE-L(R-L)(Lin and Och, 2004), CIDEr(C)(Vedantam et al., 2015), SPICE(S)(Anderson et al., 2016).
- *Similarity.* We evaluate the generation by newer metrics: CLIP-S and RefCLIPScore (Ref-CLIP)(Hessel et al., 2021), BERTScore (P-Bert)(Zhang et al., 2020), which achieve higher correlation with human judgments.
- *Diversity.* Diversity (Li et al., 2016) is a metric that evaluates the diversity of the generated captions. We report Dist-2(D@2) and Dist-3(D@3) by measuring the diversity of bigrams and trigrams in the generation.
- *Vocabulary usage.* To analyze the diversity of the generated captions, according to (Dai et al., 2018), we compute vocabulary usage(Voc-u), which accounts for the percentage of words in the vocabulary that are used in the generated captions.

4.2 Baseline

We adopt the previous competitive image captioning approaches to serve as the baseline models:

MTIC (Cornia et al., 2020): MITC is a transformer-based architecture for image captioning. Its image features extracted are by ResNet (denoted as grid-based features).

DLCT (Luo et al., 2021): DLCT achieves the complementarity of region and grid features for

Method	Human Evaluation [↑]			Paras (M) [↓]
	Fluency	Sim	Div	
MTIC	3.65	3.63	3.52	38.44
DLCT	3.70	3.25	3.43	63.04
Capdec	3.53	2.95	3.29	178.03
ClipCap	3.83	3.38	3.67	155.91
Ours(T)	3.79	3.84	3.95	38.25
Ours(B)	4.07	3.95	4.12	94.83

Table 3: The results of human evaluation and the number of trainable parameters for different methods.

image captioning. To extract visual features, DLCT uses the pretrained Faster-RCNN (Ren et al., 2015).

CapDec (Nukrai et al., 2022): CapDec is a simple and intuitive approach to learning a captioning model based on CLIP.

ClipCap (Mokady et al., 2021): ClipCap leverages powerful vision-language pre-trained models (CLIP) to simplify the captioning process. And we utilize the MLP mapping network and fine-tunes the language model. All the hyper-parameters are set following its original paper.

Since CapDec and ClipCap use CLIP to extract the same image features and freeze CLIP as our model, we use these methods as the primary baselines. We train our model for 200000 steps, with a batch size of 128. The dimension of word embedding is set to 48 and the diffusion steps $T = 1000$. All the experiments are run on NVIDIA Tesla V100 GPUs. In the decoding process, we configure the value of the candidate sentences with $n = 5$. Specifically, during the evaluation, we set the denoising steps $T = 50$, which greatly reduces the generation time.

4.3 Results

4.3.1 Image Captioning

We compare Prefix-diffusion to several baselines with different evaluation metrics, as is shown in Table 1. Our model outperforms all baselines on CLIP-S and Ref-CLIP metrics, and achieves comparable results on P-Bert score, indicating that the effectiveness of the continuous diffusion on image captioning. Not only that, we have a significant improvement on some diversity metrics (such as the D@2 and D@3). Furthermore, Prefix-diffusion covers the largest percentage of words, observed from the vocabulary used to generate captions. It implies that captions generated by Prefix-

Method	Common Metrics \uparrow						Similarity Score \uparrow			Diversity \uparrow		
	B@1	B@3	M	R-L	C	S	CLIP-S	Ref-CLIP	P-Bert	D@2	D@3	Voc-u
<i>COCO</i> \implies <i>Flickr30k</i>												
CapDec	57.2	23.9	17.1	40.3	30.3	10.8	54.4	58.7	92.1	18.5	29.4	1.2
ClipCap	64.6	29.3	18.9	44.3	44.4	12.5	56.5	61.2	92.5	19.7	32.7	1.3
Ours(B)	69.5	31.2	19.3	46.6	46.8	13.0	61.2	65.3	91.9	19.4	37.0	3.0
<i>Flickr30k</i> \implies <i>COCO</i>												
CapDec	44.1	15.2	15.7	36.4	25.7	8.6	47.7	51.4	90.4	5.5	10.4	2.0
ClipCap	55.7	23.5	19.2	42.0	51.3	12.2	54.9	60.0	91.1	11.3	21.3	3.5
Ours(B)	57.2	22.4	17.5	42.5	49.3	11.3	57.5	62.8	90.4	13.6	29.9	6.6

Table 4: The results of cross-domain captioning. *COCO* \implies *Flickr30k* means model trained on *COCO* while evaluated on *Flickr30k*, and so is *Flickr30k* \implies *COCO*. We use boldface to indicate the best performance.

diffusion contain diverse wordings and rich expressions. Our model can generate high-quality captions compared with captioning approaches that extract image feature with CLIP. Prefix-diffusion performs worse than MTIC and DLCT (who not use freeze features for image captioning) on the common metrics, partially due to the proven limitations of word-overlapping-based metrics across various domains (Hessel et al., 2021; Zhang et al., 2020), and also because our generation is more diverse in expression and correctly describe the visual content, which can be observed from similarity score and diversity metrics.

We also conduct experiments on dataset of *Flickr30k*, as presented in Table 2, from which we can draw similar conclusions with the dataset of *COCO*. Our model achieves impressive performance in the image captioning task compared to the baseline models. In detail, from the results of diversity metrics, we notice that the metrics of Dist-3 and vocabulary usage increase by more than 6.0 and 3.0, respectively. Additionally, we also observe an improvement of 2.6 and 2.8 in CLIP-S and Ref-CLIP metrics, respectively. This indicates that the diffusion model can effectively improve the caption diversity while ensuring coherence and relevance in the generated captions. To generate diverse captions, existing methods tend to generate different captions via top-k sampling. Intuitively, such methods may ignore syntactic diversity and semantic diversity that humans are really interested in. Unlike existing methods, Prefix-diffusion seeks to generate multiple captions with rich expressions from different Gaussian noises.

Figure 1 shows the captions generated by Prefix-diffusion. It is observed that the generated captions

are pretty consistent with the image as well as keeping the qualified fluency. Meanwhile, our model is able to generate diverse captions that are more like human-generated.

Furthermore, we conduct human evaluation and report the number of trainable parameters to validate the applicability of our method. As is shown in Table 3, our model only requires a small number of model parameters. It brings potential advantages of saving memory storage space and computing costs, and thus being much more useful in practice. For human evaluation, we randomly selected 20 samples and presented them in a shuffled manner to 20 annotators. The annotators rated the fluency, similarity (Sim), and diversity (Div) of the captions on a scale from 1 to 5, with higher scores indicating better quality. From the human evaluation results, We can draw similar conclusions with the automatic evaluation. Our model outperforms the baselines in diversity while holding better fluency and relevance.

The dimension of word embeddings is an important hyper-parameter. The higher dimension leads to more training time and memory usage. To further study the effect of embedding dimension in Prefix-diffusion, we conduct experiments by training with different dimensions. As is shown in Figure 4, the metrics of Bleu-1 and CIDEr are improved as the embedding dimension increases. The reason is that a word embedding becomes richer with semantic information due to the higher dimension. However, there is a performance bottleneck when we continue to increase the dimension of word embeddings. It is observed that the performance trends to be stable when the dimension goes beyond 48.

n	Common Metrics \uparrow						Similarity Score \uparrow			Diversity \uparrow		
	B@1	B@3	M	R-L	C	S	CLIP-S	Ref-CLIP	P-Bert	D@2	D@3	Voc-u
1	77.2	43.6	26.0	55.6	105.2	19.5	60.4	68.6	93.1	11.9	26.4	5.4
5	78.1	44.2	26.6	56.1	109.3	20.4	63.7	71.2	93.7	12.7	28.0	5.7
10	78.3	43.8	26.6	56.0	109.1	20.3	65.3	72.2	93.4	13.1	28.8	5.8
15	78.2	43.4	26.5	55.8	108.5	20.3	66.0	72.6	93.4	13.4	29.3	5.9

Table 5: The effect of different values of candidate captions. $n = 1$ means no cosine similarity calculation in the decoding process.

Noise Schedule	Metrics \uparrow			
	B@1	CLIP-S	Ref-CLIP	P-Bert
Square	70.5	66.8	72.2	92.6
Linear	70.4	65.9	71.6	92.3
Cosine	70.5	66.5	72.0	92.5
T-Cosine	72.5	66.5	72.3	92.9
T-Linear	78.1	63.7	71.2	93.7

Table 6: The analysis of different noise schedule in the forward process. T-Linear and T-Cosine means truncation linear noise schedule and truncation cosine noise schedule respectively.

4.3.2 Cross-domain Captioning

We also conduct experiments on cross-domain captioning to evaluate the generalization capability of Prefix-diffusion. The results of the cross-domain evaluation are shown in Table 4. We train the model on the dataset of a source domain while evaluating it on another dataset. From the results of COCO \implies Flickr30k, Prefix-diffusion achieves excellent performance over all compared approaches, with the results on the common metrics being the best. In addition, it acquires significant improvements on both Dist-3 and vocabulary usage metrics. This is due to the powerful generative ability of the diffusion model. When we train on flickr30k while evaluating on COCO, the results also show that our approach has strong capability in the cross-domain scenario. By comparing the two results, we find that Prefix-diffusion works even better when trained on a larger dataset, implying the better generalization ability.

4.3.3 Ablation

We perform an ablation study on the dataset of COCO to quantify the contribution of each module in Prefix-diffusion.

Table 5 presents the effect on the number of candidate captions. From the two groups of exper-

iments, $n = 1$ and $n = 5$, it can be seen that this selection strategy improves the performance of image captioning. We observe a significant increase in the CIDEr metric, which boosts the CIDEr score from 105.2 to 109.3. It confirms the function of calculating the similarity between the image and the candidate captions and choosing the highest. But too many candidate captions lead to a reduction in the performance of the caption fluency. This is because we use the CLIP score as the only similarity selection metric, which may neglect the fluency of captions.

As presented in Table 6, We investigate the performance of different noise schedules. Observing the results, we conclude that truncated linear noise schedule is able to generate more precise and descriptive captions. We also conclude that the semantic information is corrupted by the complicated noise schedule in the forward process, leading to a more difficult learning problem in the denoising process.

5 Conclusion and Future Work

In this paper, we propose a lightweight network for image captioning in combination with continuous diffusion, called Prefix-diffusion. Experiments and further analysis demonstrate that it can generate diverse captions while maintaining the fluency and relevance of the captions. By trained on one dataset but evaluated on the other, Prefix-diffusion presents remarkable generalization ability. Besides, our model requires a small number of training parameters, which is more applicable in reality. We also conduct ablation experiments to show the effect of the selection strategy and noise schedules. The empirical results verify that Prefix-diffusion has powerful generative ability for image captioning. For future work, we will continue to explore the potential impact of diffusion models on image captioning.

539 Limitations

540 As presented in Table 1 and Table 2, though Prefix-
541 diffusion can generate diverse captions with rela-
542 tively less parameters, it is inferior to MTIC and
543 DLCT on the common metrics. But it performs
544 well on newer metrics which have been shown
545 higher correlation with human generation. The
546 reason is that our generated captions have a rich
547 expression that is inconsistent with the reference
548 text, but still convey the same underlying semantics.
549 The length is an important property as it reflects the
550 amount of information carried by a caption. Since
551 our model is a non-autoregressive model, we can-
552 not control the length of the generated text, leading
553 to a less accurate description of the image. We
554 leave this part of exploration for future work.

555 Ethics Statement

556 Since the proposed Prefix-diffusion can be used to
557 generate captions. With the advantages of being
558 accurate, diverse and descriptive, its generation is
559 more like human-generated. This would benefit im-
560 age captioning applications on downstream tasks,
561 such as chatting robots and automatic voice guide
562 system. On the other hand, the large number of
563 image captions will make it difficult to distinguish
564 human-wrote from machine-generated. Hence, ex-
565 ploring adversarial attacks on image captioning is
566 necessary. Moreover, excellent captions should
567 involve a variety of words and rich expressions,
568 which prevents them from being too dull or tedious.
569 The diffusion model generates new samples from
570 different noises. Therefore, Prefix-diffusion can be
571 used to improve the diversity of the captions.

572 References

573 Peter Anderson, Basura Fernando, Mark Johnson, and
574 Stephen Gould. 2016. [Spice: Semantic propositional
575 image caption evaluation](#). In *European conference
576 on computer vision*, pages 382–398. Springer.

577 Peter Anderson, Xiaodong He, Chris Buehler, Damien
578 Teney, Mark Johnson, Stephen Gould, and Lei Zhang.
579 2018. [Bottom-up and top-down attention for image
580 captioning and visual question answering](#). In *Pro-
581 ceedings of the IEEE conference on computer vision
582 and pattern recognition*, pages 6077–6086.

583 Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vah-
584 dat, Jiaming Song, Karsten Kreis, Miika Aittala,
585 Timo Aila, Samuli Laine, Bryan Catanzaro, et al.
586 2022. [ediffi: Text-to-image diffusion models with
587 an ensemble of expert denoisers](#). *arXiv preprint
588 arXiv:2211.01324*.

Marcella Cornia, Matteo Stefanini, Lorenzo Baraldi, 589
and Rita Cucchiara. 2020. [Meshed-memory trans- 590
former for image captioning](#). In *Proceedings of the 591
IEEE/CVF conference on computer vision and pat- 592
tern recognition*, pages 10578–10587. 593

Bo Dai, Sanja Fidler, and Dahua Lin. 2018. [A neural 594
compositional paradigm for image captioning](#). *Ad- 595
vances in Neural Information Processing Systems*, 596
31:656–666. 597

Chaorui Deng, Ning Ding, Mingkui Tan, and Qi Wu. 598
2020. [Length-controllable image captioning](#). In *Eu- 599
ropean Conference on Computer Vision*, pages 712– 600
729. Springer. 601

Michael Denkowski and Alon Lavie. 2014. [Meteor 602
universal: Language specific translation evaluation 603
for any target language](#). In *Proceedings of the ninth 604
workshop on statistical machine translation*, pages 605
376–380. 606

Zhengcong Fei. 2020. [Iterative back modification for 607
faster image captioning](#). In *Proceedings of the 28th 608
ACM International Conference on Multimedia*, pages 609
3182–3190. 610

Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, 611
Amit H Bermano, Gal Chechik, and Daniel Cohen- 612
Or. 2022. [An image is worth one word: Personaliz- 613
ing text-to-image generation using textual inversion](#). 614
arXiv preprint arXiv:2208.01618. 615

Federico A Galatolo, Mario GCA Cimino, and Gigliola 616
Vaglini. 2021. [Generating images from caption and 617
vice versa via clip-guided generative latent space 618
search](#). *arXiv preprint arXiv:2102.01645*. 619

Junlong Gao, Xi Meng, Shiqi Wang, Xia Li, Shan- 620
she Wang, Siwei Ma, and Wen Gao. 2019. [Masked 621
non-autoregressive image captioning](#). *arXiv preprint 622
arXiv:1906.00717*. 623

Shansan Gong, Mukai Li, Jiangtao Feng, Zhiyong Wu, 624
and LingPeng Kong. 2022. [Diffuseq: Sequence to 625
sequence text generation with diffusion models](#). 626

Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, 627
Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron 628
Courville, and Yoshua Bengio. 2020. [Generative 629
adversarial networks](#). *Communications of the ACM*, 630
63(11):139–144. 631

Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, 632
Bo Zhang, Dongdong Chen, Lu Yuan, and Baining 633
Guo. 2022. [Vector quantized diffusion model for text- 634
to-image synthesis](#). In *Proceedings of the IEEE/CVF 635
Conference on Computer Vision and Pattern Recog- 636
nition*, pages 10696–10706. 637

Longteng Guo, Jing Liu, Xinxin Zhu, Xingjian He, Jie 638
Jiang, and Hanqing Lu. 2020. [Non-autoregressive 639
image captioning with counterfactuals-critical multi- 640
agent learning](#). *arXiv preprint arXiv:2005.04690*. 641

752	Alec Radford, Jeffrey Wu, Rewon Child, David Luan,	Luowei Zhou, Hamid Palangi, Lei Zhang, Houdong Hu,	807
753	Dario Amodei, Ilya Sutskever, et al. 2019. Language	Jason Corso, and Jianfeng Gao. 2020. Unified vision-	808
754	models are unsupervised multitask learners . <i>OpenAI</i>	language pre-training for image captioning and vqa .	809
755	<i>blog</i> , 1(8):9.	In <i>Proceedings of the AAAI Conference on Artificial</i>	810
		<i>Intelligence</i> , volume 34, pages 13041–13049.	811
756	Rita Ramos, Bruno Martins, Desmond Elliott, and Yova	Wanrong Zhu, An Yan, Yujie Lu, Wenda Xu, Xin Eric	812
757	Kementchedjieva. 2022. Smallcap: Lightweight	Wang, Miguel Eckstein, and William Yang Wang.	813
758	image captioning prompted with retrieval augmenta-	2022. Visualize before you write: Imagination-	814
759	tion . In <i>Proceedings of the IEEE/CVF Conference</i>	guided open-ended text generation . <i>arXiv preprint</i>	815
760	<i>on Computer Vision and Pattern Recognition</i> , pages	<i>arXiv:2210.03765</i> .	816
761	2840–2849.		
762	Shaoqing Ren, Kaiming He, Ross Girshick, and Jian		
763	Sun. 2015. Faster r-cnn: Towards real-time object		
764	detection with region proposal networks . <i>Advances</i>		
765	<i>in neural information processing systems</i> , 28.		
766	Robin Rombach, Andreas Blattmann, Dominik Lorenz,		
767	Patrick Esser, and Björn Ommer. 2022. High-		
768	resolution image synthesis with latent diffusion mod-		
769	els . In <i>Proceedings of the IEEE/CVF Conference</i>		
770	<i>on Computer Vision and Pattern Recognition</i> , pages		
771	10684–10695.		
772	Jascha Sohl-Dickstein, Eric Weiss, Niru Mah-		
773	eswaranathan, and Surya Ganguli. 2015. Deep un-		
774	supervised learning using nonequilibrium thermody-		
775	namics . In <i>International Conference on Machine</i>		
776	<i>Learning</i> , pages 2256–2265. PMLR.		
777	Jiaming Song, Chenlin Meng, and Stefano Ermon. 2020.		
778	Denoising diffusion implicit models . In <i>International</i>		
779	<i>Conference on Learning Representations</i> .		
780	Robin Strudel, Corentin Tallec, Florent Althé, Yilun		
781	Du, Yaroslav Ganin, Arthur Mensch, Will Grathwohl,		
782	Nikolay Savinov, Sander Dieleman, Laurent Sifre,		
783	and Rémi Leblond. 2022. Self-conditioned embed-		
784	ding diffusion for text generation .		
785	Ramakrishna Vedantam, C Lawrence Zitnick, and Devi		
786	Pariikh. 2015. Cider: Consensus-based image de-		
787	scription evaluation . In <i>Proceedings of the IEEE</i>		
788	<i>conference on computer vision and pattern recogni-</i>		
789	<i>tion</i> , pages 4566–4575.		
790	Oriol Vinyals, Alexander Toshev, Samy Bengio, and		
791	Dumitru Erhan. 2015. Show and tell: A neural image		
792	caption generator . In <i>Proceedings of the IEEE con-</i>		
793	<i>ference on computer vision and pattern recognition</i> ,		
794	pages 3156–3164.		
795	Dongchao Yang, Jianwei Yu, Helin Wang, Wen Wang,		
796	Chao Weng, Yuexian Zou, and Dong Yu. 2022. Diff-		
797	sound: Discrete diffusion model for text-to-sound		
798	generation . <i>arXiv preprint arXiv:2207.09983</i> .		
799	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Wein-		
800	berger, and Yoav Artzi. 2020. Bertscore: Evaluating		
801	text generation with bert . In <i>International Confer-</i>		
802	<i>ence on Learning Representations</i> .		
803	Shanshan Zhao, Lixiang Li, Haipeng Peng, Zihang		
804	Yang, and Jiaxuan Zhang. 2020. Image caption		
805	generation via unified retrieval and generation-based		
806	method . <i>Applied Sciences</i> , 10(18):6235.		