

LaDiC: Are Diffusion Models Really Inferior to Autoregressive Counterparts for Image-to-text Generation?

Anonymous ACL submission

Abstract

Diffusion models have demonstrated remarkable capabilities in text-to-image generation. However, their performance in image-to-text generation, specifically image captioning, has trailed behind Auto-Regressive (AR) models, casting doubts on their suitability for such tasks. In this work, we reexamine diffusion models, highlighting their capacity for holistic context modeling and parallel decoding. These advantages address the inherent limitations of AR methods, such as slow inference speed, error propagation, and unidirectional constraints. Additionally, We identify the lack of an effective latent space for image-text alignment and the discordance between continuous diffusion processes and discrete textual data in previous works limit their performance. In response, we introduce a novel architecture, LaDiC, featuring a split BERT to create a dedicated latent space for captions and a regularization module to manage varying text lengths. Our framework further incorporates a diffuser for semantic image-to-text conversion and a Back&Refine technique to enhance token interactivity during inference. LaDiC achieves a state-of-the-art performance for diffusion-based methods on the MS COCO dataset with a BLEU@4 score of 38.2 and a CIDEr score of 126.2, demonstrating exceptional performance without pretraining or ancillary modules. This indicates strong competitiveness with AR models, revealing the previously untapped potential of diffusion models in image-to-text generation.

1 Introduction

Recently, we have witnessed a multitude of impressive and exciting applications of diffusion models in text-to-image generation tasks (OpenAI, 2023; Podell et al., 2023; Dai et al., 2023). Nevertheless, the inverse process of image-to-text generation remains less explored. Some pioneering efforts (Li et al., 2022b; Yuan et al., 2022) have

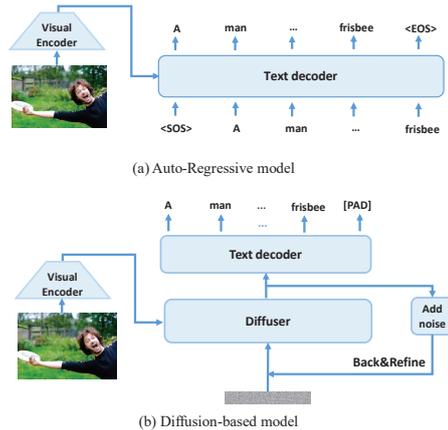


Figure 1: (a) Token-by-token generation manner of AR-based image captioning model. (b) Gradually denoising generation manner of diffusion-based model (Ours).

aimed to integrate diffusion models into text generation or Seq2Seq tasks, and they have largely followed the traditional Encoder-Decoder framework in NLP, utilizing the diffusion model as a text decoder. However, their scope has been limited to handling unimodal data. Although subsequent research (He et al., 2023b; Liu et al., 2023a) which focused on the image-to-text task introduces visual capability into this paradigm by treating visual embedding as a special token or encoded hidden states, their performance has consistently trailed behind that of Auto-Regressive (AR) models. Only through intricate architecture (Luo et al., 2022) or external data (Zhu et al., 2022) can they barely achieve comparable results, raising doubts about whether diffusion models have inherent limitations, potentially making them less suitable for the image-to-text task.

In this study, we aim to dispel this doubt by deeply reexamining the diffusion-based image-to-text paradigm and unveiling its distinct benefits. Unlike conventional AR approaches that sequentially generate captions token by token (Fig. 1a), diffusion-based models take Gaussian noise as in-

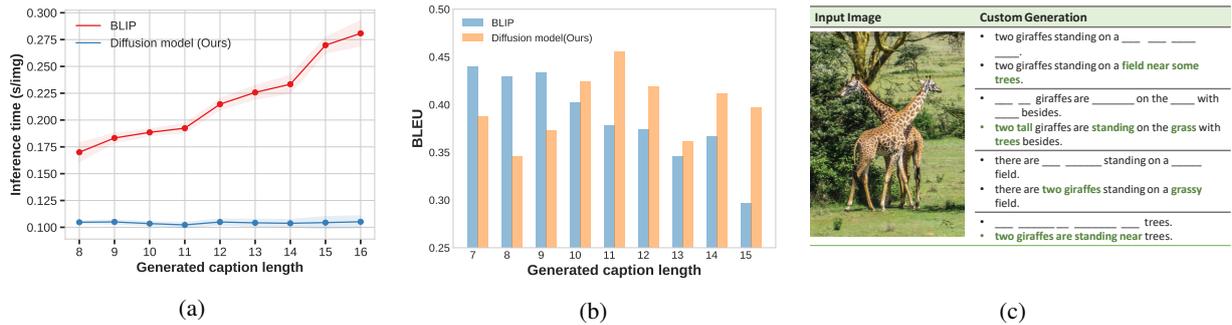


Figure 2: (a) Inference time of AR model (BLIP) and our diffusion model (LaDiC) as generated caption length increases. (b) BLEU score of BLIP and LaDiC with increasing generated caption length. (c) LaDiC’s ability of custom generation.

put and iteratively denoise it under image guidance to simultaneously produce the entire caption (Fig. 1b). Thus our diffusion-based model exhibits three key advantages: **(1) Parallel Decoding:** Diffusion-based models emit all tokens in parallel, significantly reducing inference time for lengthy target captions. As illustrated in Fig. 2a, the inference time of AR models like BLIP (Li et al., 2022a) proliferates as text length grows, while our model can set a maximum length in advance and emit all tokens concurrently. For instance, when the caption length reaches 16, our model is approximately 3× faster than BLIP. **(2) Holistic Context Consideration:** Unlike the single-directional information flow of AR models (left to right), diffusion-based models can consider more holistic contexts, mitigating error accumulation. As depicted in Fig. 2b, the BLEU metric of BLIP-generated captions declines rapidly with increasing text length, whereas our diffusion-based model maintains performance. **(3) Flexible Generation:** AR models adhere to a fixed unidirectional generation manner, whereas our model demonstrates much greater flexibility. We can custom generate captions based on tokens in nearly any position, as shown in Fig. 2c, a capability challenging for AR image captioning models.

Hence, we are convinced that diffusion-based image-to-text generation offers unique advantages and merits further exploration. Upon examining prior diffusion-based models, we deduce that their unsatisfactory performance primarily stems from two factors: **(I)** Two significant gaps exist in translating between images and text, namely the gap between visual information and textual representation, and the gap between high-level text semantics and specific words. Simultaneously addressing both gaps within the previous paradigm as shown in Fig. 3, proves to be a challenging task for dif-

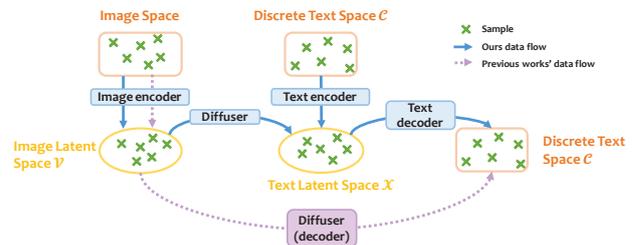


Figure 3: Comparison of the pipeline between our LaDiC and that of previous diffusion-based models.

fusion models. **(II)** Substantial discrepancies exist between text and other continuous modalities like images or audio. For instance, classical continuous diffusion models naturally align with the pixel space but struggle to transition directly to the discrete text space. Additionally, generated images have a fixed size, while caption lengths vary, presenting another challenge for diffusion models in determining the boundaries of generated captions. Given these considerations, we meticulously design a novel architecture LaDiC, a **Latent Diffusion-based Captioner**, for further amplifying the capability of diffusion models in image-to-text generation. As depicted in Fig. 3, rather than directly generating text from image representation, we treat the diffuser as an interface translating image information to high-level text representation. This approach alleviates the diffusion model’s burden, enabling it to leverage its powerful generation capabilities in high-level semantic spaces (Ramesh et al., 2022), while the decoder retains its ability to generate discrete tokens from latent space. During training, a text encoder is employed to generate ground-truth text latent codes, and during inference, it can be safely discarded.

In detail, we leverage a pre-trained language model like BERT (Devlin et al., 2019) to gener-

131 ate the text latent space, benefiting from its ability
132 to capture semantic context for creating a more
133 fluent caption. Recognizing the higher informa-
134 tion density in text compared to images, we parti-
135 tioned BERT into two parts, namely the main body
136 of the text encoder and a Non-Auto-Regressive
137 (NAR) decoder. Subsequently, we diffused on its
138 middle layer characterized by a lower-level repre-
139 sentation of text, to align more effectively with
140 images. To regulate this latent space, we pro-
141 pose a post-processing submodule after the text
142 encoder, including normalization and reassignment
143 procedures for addressing problems like variable
144 length of text. Furthermore, the diffuser serves as
145 a bridge between image and text, aiming to fit the
146 distribution of the text latent space defined above
147 conditioned on the image, wherein we utilize a
148 cross-attention mechanism for better modality fu-
149 sion. Lastly, during inference, inspired by the self-
150 conditioning (Chen et al., 2022a) which enhances
151 temporal dimension interaction, we propose the
152 Back&Refine technique to provide more interac-
153 tion between tokens in the spatial dimension, com-
154 pensating for the information loss caused by the
155 relatively independent prediction of each token in
156 the diffusion-based model.

157 We conducted experiments mainly on the COCO
158 dataset (Lin et al., 2014) to validate our model’s
159 capabilities. Remarkably, without pretraining or
160 external modules, our model achieves a BLEU@4
161 score of 38.2 and a CIDEr score of 126.2, sur-
162 passing both diffusion-based methods and tradi-
163 tional NAR models significantly. In addition to
164 the unique advantages discussed earlier, our model
165 also matches the performance of well-established
166 pretrained AR models and outperforms BLIP in
167 image paragraph captioning. These results under-
168 score the potent generative ability and immense
169 potential of diffusion models in image-to-text gen-
170 eration. We aspire that our work offers a fresh
171 perspective, fostering future research on diffusion
172 models for image-to-text generation or even other
173 text-centered multimodal generation tasks.

174 2 Related Works

175 2.1 Diffusion Models and their Applications

176 Diffusion models have recently emerged as pow-
177 erful generative models, with representative foun-
178 dational architectures such as DDPM (Ho et al.,
179 2020b) and DDIM (Song et al., 2020). These
180 methods gradually transform samples into Gaus-

181 sian noise and train a model to recover them, pre-
182 senting a simple and stable learning objective for
183 addressing issues like posterior and mode collapse
184 that challenge prior models like VAE (Kingma and
185 Welling, 2013) and GAN (Goodfellow et al., 2014).

186 The impressive generative capabilities of diffu-
187 sion models have led to their application across a
188 spectrum of fields, including image (Ramesh et al.,
189 2022; Dai et al., 2023), audio (Liu et al., 2023b;
190 Shen et al., 2023), video (Blattmann et al., 2023;
191 Girdhar et al., 2023), 3D (Poole et al., 2022), and
192 human avatar (He et al., 2023a; Hu et al., 2023),
193 among others. Yet, their adaptation to discrete
194 text spaces is an ongoing challenge. Existing ap-
195 proaches generally fall into two categories: (1) **dis-**
196 **crete diffusion models** (Austin et al., 2021; Reid
197 et al., 2022; He et al., 2022) that directly corrupt
198 text with [MASK] tokens; and (2) **continuous dif-**
199 **fusion models** (Li et al., 2022b; Gong et al., 2022;
200 Dieleman et al., 2022; Yuan et al., 2022; Lin et al.,
201 2022), which use continuous embeddings to repre-
202 sent each token. However, both approaches are con-
203 fined to unimodal and may omit high-level overall
204 semantics to some extent. Furthermore, we notice
205 the work (Lovelace et al., 2023), which explores
206 the concept of a text latent space. Yet its diffu-
207 sion model, designed for predicting BART’s (Lewis
208 et al., 2019) hidden states, still relies on an AR gen-
209 eration mechanism.

210 2.2 Image-to-text Generation

211 Image-to-text generation, or its most representa-
212 tive and general task, image captioning, aims to
213 describe the content of an image in natural lan-
214 guage. Other variants include dense captioning,
215 which illustrates each object in the picture (Johnson
216 et al., 2016), and paragraph captioning which gen-
217 erates a detailed, lengthy paragraph (Krause et al.,
218 2016) and so on. Early AR approaches for cap-
219 tioning (Karpathy and Fei-Fei, 2017; Vinyals et al.,
220 2014) employed an encoder-decoder architecture
221 with a CNN to encode images and an RNN to gener-
222 ate captions. Attention mechanisms were later
223 introduced (Huang et al., 2019; Xu et al., 2015)
224 and concurrently some researchers explored the use
225 of semantic attributes (You et al., 2016; Yao et al.,
226 2017). With the advent of Transformer (Vaswani
227 et al., 2017) and large-scale pretraining methods,
228 pretrained vision-language models (Li et al., 2022a;
229 Zhang et al., 2021; Li et al., 2020) emerged and
230 achieved high performance.

231 In contrast to the unidirectional generation of

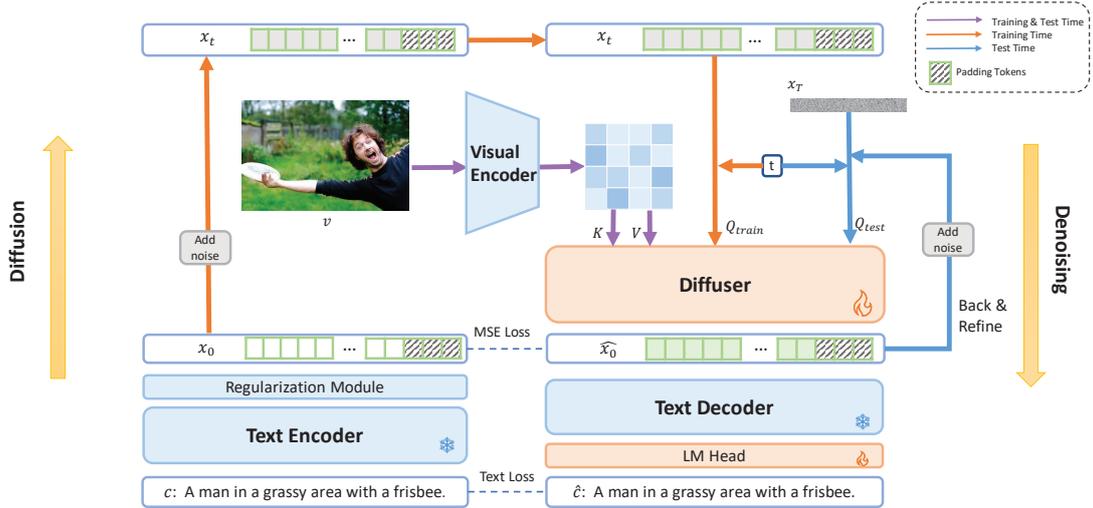


Figure 4: An overview of our LaDiC model mainly consisting of Encoder, Diffuser and Decoder. On the left is the diffusion process, and on the right is denoising process.

AR models, NAR models generate entire captions in parallel. MNIC (Gao et al., 2019) introduced the mask token strategy, and NAIC (Guo et al., 2020) employed reinforcement learning. A special class of NAR methods, diffusion-based models has recently emerged. Most models (Xu, 2022; He et al., 2023b; Liu et al., 2023a) follow the paradigm utilized in continuous diffusion models mentioned above. Additionally, Bit Diffusion (Chen et al., 2022a) encodes captions into binary bits, and DD-Cap (Zhu et al., 2022) applies a discrete diffusion model to captioning. SCD-Net (Luo et al., 2022) is the state-of-the-art diffusion-based model with a semantic-conditional diffusion process. However, its cascaded architecture is relatively complex and requires an external retrieval module, limiting its further extension. Our work reexamines the diffusion-based paradigm and proposes a novel, compact architecture with improved performance.

3 Methodology

In this section, we introduce our diffusion-based image captioning model, LaDiC. In § 3.1, we present the overall architecture of LaDiC, including its training and inference pipeline. Subsequently, from § 3.2 to § 3.4, we offer a detailed illustration.

3.1 Overview

At a macroscopic level, depicted in Fig. 3, we employ a text encoder to convert the discrete text space \mathcal{C} into a continuous latent space \mathcal{X} . Subsequently, a diffuser is trained to map the image representation space \mathcal{V} to the latent space \mathcal{X} . Specifically, in the

context of Fig. 4, given paired data (v, c) — an image and its corresponding caption, we encode the caption $c \in \mathcal{C}$ into the latent space, yielding the latent code $x_0 \in \mathcal{X}$. To model the distribution of \mathcal{X} , we adopt the diffusion models’ diffusion-denoising procedure. Initially, various levels of noise (represented by t) are introduced to x_0 to generate a noisy version x_t (left panel). Subsequently, the diffuser acts as a denoiser, recovering x_0 conditioned on the images v (right panel). Once the diffuser is sufficiently trained, a robust function $f : x_t \xrightarrow{v} x_0$ is established, connecting the image space \mathcal{V} and the latent space \mathcal{X} . During inference, given an image v^* , x_t is replaced with pure Gaussian noise $x_T \sim N(\mathbf{0}, \mathbf{I})$ and iteratively denoised by f , resulting in $x_T \xrightarrow{v^*} \hat{x}_0^*$. Finally, the decoder converts the acquired latent code back into discrete text $\hat{c}^* \in \mathcal{C}$.

3.2 Latent Space Tailored for Text

As discussed in § 1, the latent space \mathcal{X} serves as a crucial bridge between image representation \mathcal{V} and discrete text \mathcal{C} , significantly alleviating the burden on diffusion models. Therefore, it is paramount to design a latent space that incorporates rich semantic information and easy for the diffuser to adapt its distribution. Earlier studies like (He et al., 2023b) predominantly translate discrete text into continuous space using an embedding matrix, completely overlooking overall semantics, posing a challenge in aligning images with these independent token embeddings. In contrast, we utilize pre-trained language models such as BERT (Devlin et al., 2019) to construct a high-level semantic latent space through

contextual embedding methods and meanwhile harness abundant inherent knowledge from the pre-trained corpus.

Moreover, it is acknowledged that a significant information density gap exists between vision and language (He et al., 2021). Large portions of image pixels tend to be redundant while in natural language, the majority of tokens convey rich semantic information. To address this mismatch, we split the BERT model into two parts: the lower part serves as the main body of the text encoder, and the upper functions as the decoder. Through setting text latent space as the middle layer of BERT, which contains lower-level features of the text, we observed that it better aligns vision and language, thereby enhancing performance. In addition, to improve the decoder’s ability to reconstruct the text space, we make the parameters in the language model head trainable.

However, this latent space is still deemed unsatisfactory, prompting the addition of a postprocess submodule comprising two procedures: normalization and reassignment. Given that the embeddings in the BERT space vary dramatically, it is unreasonable to add the same scale of noise to various norms of embeddings. Thus we collect a subset of all captions in the dataset and calculate the mean and standard deviation of their corresponding latent codes $\hat{\mu}(x), \hat{\sigma}(x)$. During training, these statistics are used to regularize the space as follows $\text{norm}(x) = [x - \hat{\mu}(x)] / [\hat{\sigma}(x) + \epsilon]$. During inference, an unnorm module is applied to the predicted \hat{x}_0 before feeding it to the decoder. Moreover, a discrepancy between applying the diffusion model to text and image is the variable length of text, which forces the model to implicitly learn this supervised signal. In LaDiC, we extract all positions of special tokens like [CLS], [SEP], [PAD], whose representations will be messy in contextual embeddings, forming a set \mathcal{S} . We then reassign what we call an empty token to the latent code in these locations, namely pasting vectors with all 0s, as demonstrated in Equation 1. Here, x_i^{final} represents the i -th position of the final latent codes.

$$x_i^{final} = \begin{cases} [\text{norm}(x)]_i & i \notin \mathcal{S} \\ \mathbf{0}, & i \in \mathcal{S} \end{cases} \quad (1)$$

Through this technique, for short captions with pad tokens at the end, the diffuser can quickly identify this repeated pattern and easily recover these unified zero vectors, implicitly learning sentence

boundaries. This approach avoids the need for an additional module for predicting sentence length, as seen in DDCap (Zhu et al., 2022). Furthermore, despite a fixed length given during inference, the token forecasted as a pad will be mapped to the empty token defined above, and can be easily erased by postprocessing.

3.3 Diffuser Mapping Image to Text

The caption diffuser serves as an interface transforming the vision space \mathcal{V} into the text latent space \mathcal{X} . To fit the distribution of space \mathcal{X} by classical diffusion models, firstly we sample x_t , the noisy version of the latent code x_0 , as $x_t|x_0 \sim \mathcal{N}(\sqrt{\bar{\alpha}_t}x_0, \sqrt{1 - \bar{\alpha}_t}\mathbf{I})$, where $\beta_t \in (0, 1)$ is the variance schedule and $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i = \prod_{i=1}^t (1 - \beta_i)$. A notable property of this setting is that as $T \rightarrow \infty$, x_T is equivalent to an isotropic Gaussian distribution, aligning with the starting state of inference. Then for the denoising process, we use a Transformer encoder and predict the original x_0 based on the image directly, denoted as $\hat{x}_0 = f_\phi(x_t, v, t)$, where ϕ represents the parameters of the diffuser. A rigorous mathematical explanation of the diffusion model can be found in App. D if necessary.

Now, let’s delve into the architecture and training method of f . In contrast to some previous approaches that inject image information by appending the [CLS] token of the vision encoder to text (Xu, 2022; He et al., 2023b), our LaDiC model adopts the cross-attention mechanism, treating text as the query to extract information from related image patches. We hypothesize that this approach will inject vision information more effectively. Additionally, we adapt classifier-free guidance (Ho and Salimans, 2022) to this task by randomly zeroing out some images and feeding them into the model together with normal training samples. During inference, a linear combination of the conditional and unconditional estimates is performed: $\hat{x}_0 = (1 + w)f_\phi(x_t, v, t) - wf_\phi(x_t, \emptyset, t)$ where w is a predefined hyperparameter.

Regarding the loss function, a natural component is the Mean Squared Error (MSE) loss between \hat{x}_0 and x_0 . Moreover, considering that our framework has already predicted \hat{x}_0 , a loss based on the softmax distribution can be calculated as $L_{text} = \prod_{i=1}^n p_\theta(w^i|\hat{x}_0^i)$, where θ represents the parameters, to evaluate the distance of the predicted result to the ground truth in the discrete caption space \mathcal{C} . This approach makes the output

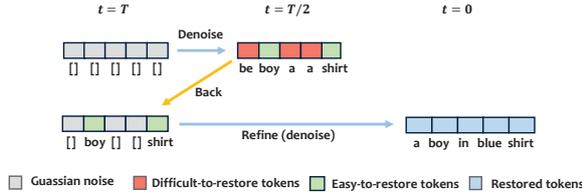


Figure 5: Illustration of Back&Refine technique.

of caption diffuser shrink faster, sharing the same intuition with XE loss in (Luo et al., 2022) and anchor loss in (Gao et al., 2022). Meanwhile, it also helps adjust the language model head in the decoder. Therefore, the loss utilized to train the caption diffuser in LaDiC is summarized as below, where $\hat{x}_0 = f_\phi(x_t, v, t)$, and λ is a hyperparameter:

$$L = \|f_\phi(x_t, v, t) - x_0\| + \lambda \prod_{i=1}^n p_\theta(w^i | \hat{x}_0^i), \quad (2)$$

3.4 Back & Refine Technique during Inference

We observe that the diffusion model exhibits a certain degree of independence in both spatial and temporal dimensions during the inference process. In the temporal aspect, (Chen et al., 2022a) found that the previously estimated \hat{x}_0 is simply discarded when estimating x_0 from a new time step. They propose self-conditioning technique, utilizing the previously generated result to improve the sample quality. However, there is little exploration in the spatial dimension, i.e., the positions of each word in a sentence. In contrast to AR models with explicit sequential dependencies across tokens, the diffusion model emits all tokens in parallel. Undoubtedly, this approach boosts the inference speed but partially loses the information flow between tokens. Considering that some tokens are easily recovered, such as the main objects in the picture, adding the same scale of noise to these well-restored tokens as the others is somewhat unreasonable and wasteful. On the contrary, we should leverage these informative tokens. Therefore, we propose a technique named Back&Refine. As illustrated in Fig. 5, let’s say we want to predict a sentence with a sequence length L and a sampling step T . Then at time $T/2$, several tokens are considered good enough, measured by the confidence scores of our model. We rank these scores and label tokens that fall in the lagging half. For these $L/2$ tokens that the model

is not currently confident about, we try to reproduce them by noising them with complete Gaussian noise, while the others remain unchanged as information. Then we set the current $t = T$ and start a brand new denoising procedure.

4 Experiments

4.1 Experimental Settings

Dataset and Metrics We conduct our experiments on MS COCO Karpathy split (Lin et al., 2014; Karpathy and Fei-Fei, 2014), which comprises 113,287 training images, 5,000 validation images, and 5,000 test images. Each image is associated with 5 reference captions. For evaluating model performance, we use several metrics including BLEU@4 (Papineni et al., 2002), CIDEr-D (Vedantam et al., 2014), METEOR (Banerjee and Lavie, 2005), ROUGE-L (Lin, 2004), and SPICE (Anderson et al., 2016). Additionally, we employ two model-based metrics: CLIP Score (Hessel et al., 2021) to assess semantic alignment between generated captions and images, and BERT Score (Zhang et al., 2019) to evaluate text quality.

Implementation Details In our LaDiC model, the encoder and decoder are frozen, except for the LM-head. The initial weights are taken from the bottom 6 layers and top 6 layers of BERT_{base} for the encoder and decoder, respectively. For the diffusion forward process, we employ the widely used cosine β schedule and adopt the noise factor (Gao et al., 2022). The diffuser consists of 12 transformer encoder blocks with additional cross-attention layer in each block and the weights are randomly initialized. To extract image features, we use the pretrained image encoder from BLIP_{base} (Li et al., 2022a), which employs ViT-B/16, for fair comparison with BLIP. The model is trained on an 8-V100 node for 60 epochs with a peak learning rate of 5e-5 and a warmup ratio of 0.1. Further details can be found in App. C.

4.2 Quantitative Analysis

We benchmark our LaDiC model against prior baselines, encompassing auto-regressive, traditional non-autoregressive, and diffusion-based models, leveraging the COCO dataset (refer to Tab. 1). Our model achieves state-of-the-art performance across various metrics for both diffusion-based and traditional NAR models. Specifically, LaDiC attains a BLEU@4 score of 38.2 and a CIDEr score of 126.2, marking improvements of 0.9 and 8.2, respectively,

Model	# Images	B@4	C	M	S	R	CLIP-score	BERT-score
<i>Autoregressive</i>								
Show and Tell (Vinyals et al., 2014)	-	31.4	97.2	25.0	18.1	53.1	69.7	93.4
CLIPCap (Mokady, 2021)	-	33.5	113.1	27.5	21.1	-	-	-
OSCAR† (Li et al., 2020)	7M	36.5	123.7	30.3	23.1	-	-	-
ViTCap† (Fang et al., 2021)	4M	36.3	125.2	29.3	22.6	58.1	-	-
VinVL† (Zhang et al., 2021)	6M	38.2	129.3	30.3	23.6	60.9	76.6	88.5
BLIP† (Li et al., 2022a)	129M	39.7	133.3	-	-	-	77.4	94.4
GIT† (Wang et al., 2022)	4M	40.4	131.4	30.0	23.0	-	-	-
<i>Traditional Non-autoregressive</i>								
NAIC _{KD} (Guo et al., 2020)	0.1M	28.5	98.2	23.6	18.5	52.3	-	-
MNIC (Gao et al., 2019)	0.1M	31.5	108.5	27.5	21.1	55.6	-	-
FNIC (Fei, 2019)	0.1M	36.2	115.7	27.1	20.2	55.3	-	-
<i>Diffusion model based</i>								
DiffCap (He et al., 2023b)	0.1M	31.6	104.3	26.5	19.6	55.1	73.6*	92.2*
Bit Diffusion (Chen et al., 2022b)	0.1M	34.7	115.0	-	-	58.0	-	-
DDCap (Zhu et al., 2022)	0.1M	35.0	117.8	28.2	21.7	57.4	74.1*	93.4*
SCD Net (Luo et al., 2022)	0.1M	37.3	118.0	28.1	21.6	58.0	74.5*	93.4*
LaDiC (ours, step 5)	0.1M	35.1	115.2	27.4	21.3	56.7	77.1	93.8
LaDiC (ours, step 30)	0.1M	38.2	126.2	29.5	22.4	58.7	77.3	94.4

Table 1: Comparison results on COCO dataset, where B@4, M, R, C denote BLEU@4, METEOR, ROUGE-L, CIDEr and SPICE scores. † indicates pretrained models and we gray them out. * represents results of models trained by ourselves. For a fair comparison, all models will not incorporate results by CIDEr optimization.

483 compared to the previous state-of-the-art method,
484 SCD-Net. Remarkably, a variant of our model, uti-
485 lizing only 5 inference steps, even outperforms all
486 prior diffusion-based models in both CLIP-Score
487 and BERT-Score. Moreover, in addition to its dis-
488 tinctive advantages over AR models, it is notewor-
489 thy that LaDiC exhibits comparable performance
490 with well-established pretraining auto-regressive
491 frameworks such as ViTCap and VinVL, despite
492 being trained on significantly less data.

493 To evaluate our model’s capacity for considering
494 holistic context, we tackle the task of image para-
495 graph captioning (Krause et al., 2016) to generate a
496 multi-sentence description of an image. Our model
497 seamlessly adapts to paragraph captioning by ex-
498 tending the predefined length without additional
499 special designs. Training our model on the dataset
500 from (Krause et al., 2016) yields a BLEU@4 score
501 of 7.3, surpassing finetuned BLIP’s 6.1 and high-
502 lighting our model’s advantage in mitigating error
503 accumulation (refer to App. B.1 for more details).
504 All these quantitative indicators above substantiate
505 the accuracy and high quality of the captions
506 generated by our model.

507 4.3 Case Studies and Human Evaluation

508 We conduct a case study to illustrate the faithful-
509 ness and diversity of the captions generated by
510 LaDiC. As depicted in Fig. 6, the generated cap-
511 tions are not only reasonable and fluent but also

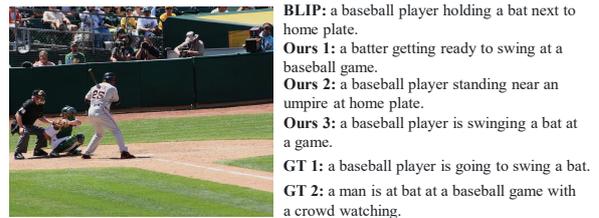


Figure 6: An example generated by our model.

512 exhibit inherent diversity due to the varied sam-
513 pling noises introduced at the start of inference.
514 Additional examples can be found in App. A.1. In
515 the context of image paragraph captioning gener-
516 ation, Fig. 7 reveals a notable difference. While each
517 sentence in the captions generated by BLIP demon-
518 strates good quality, they tend to appear somewhat
519 independent of each other, with many initiating
520 with ‘the man’ and occasionally featuring repeti-
521 tions. Conversely, by leveraging a broader context,
522 our model produces sentences with a more cohesive
523 logical relationship.

524 We conduct user studies to evaluate the gener-
525 ated captions of LaDiC, inviting volunteers to rate
526 captions on a five-point scale (1-5) for accuracy,
527 conciseness and fluency. The results, presented in
528 Tab. 2, demonstrate that our model surpasses the
529 previous diffusion-based state-of-the-art SCD-Net
530 in both aspects and achieves comparable results
531 with BLIP. Details can be found in App. B.2.



Finetuned BLIP:
 a man playing tennis. the man is wearing a white shirt and black shorts. the man is holding a tennis racket in his hand. the man is wearing a white shirt and black shorts.

Ours:
 a man playing tennis is standing on a tennis court. there is a green tennis ball above him. he is wearing white shirt, and black shorts. there is a white line on the court.

Figure 7: An example generated by fine-tuned BLIP model and ours in image paragraph captioning.

Model	SCD-Net	BLIP	Ours
Fluency	2.8	4.9	4.5
Accuracy	3.3	4.2	4.4
Conciseness	3.4	4.4	4.7

Table 2: Results of user study.

4.4 Unleashing the Speed of Diffusion Model

Despite their powerful generative capabilities, diffusion models are notorious for their slow inference speed. Most previous works require more than 50 inference steps, significantly slower than traditional NAR methods, which typically involve around ten refinement procedures. However, as shown in Tab. 1, our model achieves remarkable performance even with just 5 steps. We attribute this surprising convergence speed to specific techniques employed in our LaDiC model. Firstly, the direct prediction of x_0 and the definition of text loss enable the model to rapidly learn the distribution of discrete caption text, akin to the consistency model (Song et al., 2023). Secondly, the carefully selected noise schedule and noise factor significantly enhance the learning process of diffusion models. Regarding observed latency, the results in Tab. 3 (measured on a single A40 GPU with a batch size of 256) and Fig. 2a demonstrate that our model showcases a rapid inference speed, excelling not only in the domain of diffusion-based models but also when compared to auto-regressive models.

4.5 Customizing the Generation Process

In contrast to the unidirectional generation manner of AR models, our LaDiC model adeptly fills in empty words at almost any position within a sentence, harnessing its capability to capture more holistic information, as demonstrated in Fig. 2c. Technically, when provided with a caption containing blanks, we extract contextual embeddings of the given tokens and mask the blank tokens with Gaussian noise. The standard denoising process

Model	DiffCap	DDCap	Ours
Inference latency(s/img)	0.625	0.113	0.049

Table 3: Inference latency of diffusion-based models.

#Row	Cross-attention	Text loss	PLM	Norm-Reass	Split	B&R	B@4	C
a							15.4	46.3
b	✓						20.3	59.1
c	✓	✓					22.8	76.3
d	✓	✓	✓				26.9	91.8
e	✓	✓	✓	✓			31.6	103.5
f	✓	✓	✓	✓	✓		33.4	110.0
g	✓	✓	✓	✓	✓	✓	34.1	113.4

Table 4: Ablation on COCO dataset.

is then applied, with the exception of reinserting the embeddings of predefined tokens back to their respective positions after each inference step, ensuring that the given information is retained. Through this method, our model functions as a customized generator based on the provided tokens. Additional results can be found in App. A.2.

4.6 Ablation study

To validate the effectiveness of our core designs, we conduct ablation studies on COCO with 30 training epochs. We begin with a simple baseline that appends only the [CLS] token of the image feature to the end of text embeddings and then trains the diffuser to recover them. Subsequently, we progressively incorporate our proposed techniques to evaluate their performance. As depicted in Tab. 4, all modules exhibit performance gains. The use of PLM (BERT) and regularization in this space significantly enhance performance, emphasizing the importance of a refined latent space. Techniques aimed at better capturing visual information, such as cross-attention and splitting the BERT, also play pivotal roles in improving performance.

5 Conclusion

In this paper, we reexamine the diffusion-based image-to-text paradigm and introduce a novel architecture, denoted as LaDiC. Our model attains state-of-the-art performance among diffusion-based methods and demonstrates comparable capabilities with some pre-trained AR models. Moreover, our extensive experiments reveal the exciting advantages of diffusion models over AR models in considering more holistic contexts and emitting all tokens in parallel. Consequently, we posit that diffusion models hold substantial potential for image-to-text generation and we anticipate that our work will open new possibilities in this field.

602 Limitations

603 For simplicity and focus, this paper concentrates
604 on the main research topic of image-to-text gen-
605 eration. Nevertheless, we observe that our model
606 can be readily adapted to other modalities or even
607 pure text generation with minimal modifications.
608 We leave these potential extensions for future work,
609 and meanwhile, we hope this paper will inspire
610 confidence among researchers engaging in text-
611 centered multimodal generation tasks with diffu-
612 sion models and look forward to exciting future
613 works in this area. Furthermore, due to resource
614 constraints, the model parameters and datasets em-
615 ployed in our study are not extensive. Considering
616 the remarkable emergent abilities demonstrated by
617 scaling up autoregressive models like GPT, it be-
618 comes an intriguing and worthwhile exploration to
619 investigate whether our model or general diffusion
620 models, can exhibit similar scalability.

621 **Risk Consideration:** As a generative model, our
622 model may inadvertently produce results that are
623 challenging to distinguish from human-written con-
624 tent, raising concerns about potential misuse. Em-
625 ploying text watermark techniques could be benefi-
626 cial in mitigating this issue. Additionally, diffusion
627 models typically demand substantial computational
628 resources for training, leading to increased carbon
629 dioxide emissions and environmental impact.

630 References

631 Peter Anderson, Basura Fernando, Mark Johnson, and
632 Stephen Gould. 2016. [Spice: Semantic propositional
633 image caption evaluation](#). *ArXiv*, abs/1607.08822.

634 Jacob Austin, Daniel D. Johnson, Jonathan Ho, Daniel
635 Tarlow, and Rianne van den Berg. 2021. [Structured
636 denoising diffusion models in discrete state-spaces](#).
637 In *Advances in Neural Information Processing Sys-
638 tems*.

639 Satantjeev Banerjee and Alon Lavie. 2005. [Meteor: An
640 automatic metric for mt evaluation with improved
641 correlation with human judgments](#). In *IEEValua-
642 tion@ACL*.

643 A. Blattmann, Tim Dockhorn, Sumith Kulal, Daniel
644 Mendelevitch, Maciej Kilian, and Dominik Lorenz.
645 2023. [Stable video diffusion: Scaling latent
646 video diffusion models to large datasets](#). *ArXiv*,
647 abs/2311.15127.

648 Ting Chen, Ruixiang Zhang, and Geoffrey Hinton.
649 2022a. [Analog bits: Generating discrete data us-
650 ing diffusion models with self-conditioning](#). *arXiv
651 preprint arXiv:2208.04202*.

Ting Chen, Ruixiang Zhang, and Geoffrey E. Hinton. 652
2022b. [Analog bits: Generating discrete data us- 653
ing diffusion models with self-conditioning](#). *ArXiv*, 654
abs/2208.04202. 655

Xiaoliang Dai, Ji Hou, Chih-Yao Ma, Sam Tsai, Jialiang 656
Wang, Rui Wang, Peizhao Zhang, Simon Vanden- 657
hende, Xiaofang Wang, Abhimanyu Dubey, Matthew 658
Yu, Abhishek Kadian, Filip Radenovic, Dhruv Ma- 659
hajan, Kunpeng Li, Yue Zhao, Vladan Petrovic, 660
Mitesh Kumar Singh, Simran Motwani, Yi Wen, Yi- 661
wen Song, Roshan Sumbaly, Vignesh Ramanathan, 662
Zijian He, Peter Vajda, and Devi Parikh. 2023. [Emu:
663 Enhancing image generation models using photo-
664 genic needles in a haystack](#). 665

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and 666
Kristina Toutanova. 2019. [Bert: Pre-training of deep
667 bidirectional transformers for language understand-
668 ing](#). *ArXiv*, abs/1810.04805. 669

Sander Dieleman, Laurent Sartran, Arman Roshan- 670
nai, Nikolay Savinov, Yaroslav Ganin, Pierre H. 671
Richemond, A. Doucet, Robin Strudel, Chris Dyer, 672
Conor Durkan, Curtis Hawthorne, Rémi Leblond, 673
Will Grathwohl, and Jonas Adler. 2022. [Con-
674 tinuous diffusion for categorical data](#). *ArXiv*,
675 abs/2211.15089. 676

Zhiyuan Fang, Jianfeng Wang, Xiaowei Hu, Lin Liang, 677
Zhe Gan, Lijuan Wang, Yezhou Yang, and Zicheng 678
Liu. 2021. [Injecting semantic concepts into end-to-
679 end image captioning](#). *2022 IEEE/CVF Conference
680 on Computer Vision and Pattern Recognition (CVPR)*,
681 pages 17988–17998. 682

Zhengcong Fei. 2019. [Fast image caption generation
683 with position alignment](#). *ArXiv*, abs/1912.06365. 684

Junlong Gao, Xi Meng, Shiqi Wang, Xia Li, Shan- 685
she Wang, Siwei Ma, and Wen Gao. 2019. [Masked non-autoregressive image captioning](#). *ArXiv*,
686 abs/1906.00717. 687
688

Zhujin Gao, Junliang Guo, Xuejiao Tan, Yongxin Zhu, 689
Fang Zhang, Jiang Bian, and Linli Xu. 2022. [Dif-
690 former: Empowering diffusion model on embedding
691 space for text generation](#). *ArXiv*, abs/2212.09412. 692

Rohit Girdhar, Mannat Singh, Andrew Brown, Quentin 693
Duval, Samaneh Azadi, Sai Saketh Rambhatla, Akbar 694
Shah, Xi Yin, Devi Parikh, and Ishan Misra. 2023. [Emu video: Factorizing text-to-video generation by
695 explicit image conditioning](#). *ArXiv*, abs/2311.10709. 696
697

Shansan Gong, Mukai Li, Jiangtao Feng, Zhiyong Wu, 698
and LingPeng Kong. 2022. [Diffuseq: Sequence to
699 sequence text generation with diffusion models](#). 700

Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, 701
Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. 702
Courville, and Yoshua Bengio. 2014. [Generative
703 adversarial nets](#). In *NIPS*. 704

705	Longteng Guo, Jing Liu, Xinxin Zhu, Xingjian He, Jie Jiang, and Hanqing Lu. 2020. Non-autoregressive image captioning with counterfactuals-critical multi-agent learning . In <i>International Joint Conference on Artificial Intelligence</i> .	758
706		759
707		760
708		761
709		
710	Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Doll'ar, and Ross B. Girshick. 2021. Masked autoencoders are scalable vision learners . <i>2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 15979–15988.	762
711		763
712		
713		764
714		765
715	Tianyu He, Junliang Guo, Runyi Yu, Yuchi Wang, Jialiang Zhu, Kaikai An, Leyi Li, Xu Tan, Chunyu Wang, Han Hu, HsiangTao Wu, Sheng Zhao, and Jiang Bian. 2023a. Gaia: Zero-shot talking avatar generation .	766
716		767
717		768
718		769
719		770
720	Yufeng He, Zefan Cai, Xu Gan, and Baobao Chang. 2023b. Diffcap: Exploring continuous diffusion on image captioning . <i>ArXiv</i> , abs/2305.12144.	771
721		772
722		773
723	Zhengfu He, Tianxiang Sun, Kuan Wang, Xuanjing Huang, and Xipeng Qiu. 2022. Diffusionbert: Improving generative masked language models with diffusion models . In <i>Annual Meeting of the Association for Computational Linguistics</i> .	774
724		775
725		776
726		777
727		778
728	Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. 2021. Clipscore: A reference-free evaluation metric for image captioning . <i>ArXiv</i> , abs/2104.08718.	779
729		780
730		781
731		782
732	Jonathan Ho, Ajay Jain, and P. Abbeel. 2020a. Denoising diffusion probabilistic models . <i>ArXiv</i> , abs/2006.11239.	783
733		784
734		785
735	Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020b. Denoising diffusion probabilistic models . <i>Neural Information Processing Systems, Neural Information Processing Systems</i> .	786
736		787
737		788
738		789
739	Jonathan Ho and Tim Salimans. 2022. Classifier-free diffusion guidance .	790
740		791
741	Li Hu, Xin Gao, Peng Zhang, Ke Sun, Bang Zhang, and Liefeng Bo. 2023. Animate anyone: Consistent and controllable image-to-video synthesis for character animation . <i>arXiv preprint arXiv:2311.17117</i> .	792
742		793
743		794
744		795
745	Lun Huang, Wenmin Wang, Jie Chen, and Xiao-Yong Wei. 2019. Attention on attention for image captioning . <i>2019 IEEE/CVF International Conference on Computer Vision (ICCV)</i> , pages 4633–4642.	796
746		797
747		798
748		799
749	Justin Johnson, Andrej Karpathy, and Li Fei-Fei. 2016. Densecap: Fully convolutional localization networks for dense captioning . In <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition</i> .	800
750		801
751		802
752		803
753		804
754	Andrej Karpathy and Li Fei-Fei. 2014. Deep visual-semantic alignments for generating image descriptions . <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 39:664–676.	805
755		806
756		807
757		808
	Andrej Karpathy and Li Fei-Fei. 2017. Deep visual-semantic alignments for generating image descriptions . <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , page 664–676.	809
		810
		811
		812
	Diederik P. Kingma and Max Welling. 2013. Auto-encoding variational bayes . <i>CoRR</i> , abs/1312.6114.	
	Jonathan Krause, Justin Johnson, Ranjay Krishna, and Li Fei-Fei. 2016. A hierarchical approach for generating descriptive image paragraphs . <i>2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 3337–3345.	
	Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Fei-Fei Li. 2016. Visual genome: Connecting language and vision using crowdsourced dense image annotations .	
	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdel rahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension . In <i>Annual Meeting of the Association for Computational Linguistics</i> .	
	Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022a. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation .	
	XiangLisa Li, John Thickstun, Ishaan Gulrajani, Percy Liang, and TatsunoriB. Hashimoto. 2022b. Diffusion-lm improves controllable text generation .	
	Xiujun Li, Xi Yin, Chunyuan Li, Xiaowei Hu, Pengchuan Zhang, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, Yejin Choi, and Jianfeng Gao. 2020. Oscar: Object-semantics aligned pre-training for vision-language tasks . In <i>European Conference on Computer Vision</i> .	
	Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries . In <i>Annual Meeting of the Association for Computational Linguistics</i> .	
	Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. 2014. Microsoft coco: Common objects in context . In <i>European Conference on Computer Vision</i> .	
	Zheng-Wen Lin, Yeyun Gong, Yelong Shen, Tong Wu, Zhihao Fan, Chen Lin, Nan Duan, and Weizhu Chen. 2022. Text generation with diffusion language models: A pre-training approach with continuous paragraph denoise . In <i>International Conference on Machine Learning</i> .	
	Guisheng Liu, Yi Li, Zhengcong Fei, Haiyan Fu, Xi-angyang Luo, and Yanqing Guo. 2023a. Prefix-diffusion: A lightweight diffusion model for diverse image captioning . <i>ArXiv</i> , abs/2309.04965.	

813	Haohe Liu, Zehua Chen, Yiitan Yuan, Xinhao Mei,	Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi	867
814	Xubo Liu, Danilo P. Mandic, Wenwu Wang, and	Parikh. 2014. Cider: Consensus-based image descrip-	868
815	MarkD . Plumbly. 2023b. Audioldm: Text-to-	tion evaluation . <i>2015 IEEE Conference on Computer</i>	869
816	audio generation with latent diffusion models . <i>ArXiv</i> ,	<i>Vision and Pattern Recognition (CVPR)</i> , pages 4566–	870
817	abs/2301.12503 .	4575.	871
818	Justin Lovelace, Varsha Kishore, Chao Wan, Eliot	Oriol Vinyals, Alexander Toshev, Samy Bengio, and	872
819	Shekhtman, and Kilian Q. Weinberger. 2023. La-	D. Erhan. 2014. Show and tell: A neural image	873
820	tent diffusion for language generation .	caption generator . <i>2015 IEEE Conference on Com-</i>	874
821	Jianjie Luo, Yehao Li, Yingwei Pan, Ting Yao, Jian-	<i>puter Vision and Pattern Recognition (CVPR)</i> , pages	875
822	lin Feng, Hongyang Chao, and Tao Mei. 2022.	3156–3164.	876
823	Semantic-conditional diffusion networks for image	Jianfeng Wang, Zhengyuan Yang, Xiaowei Hu, Lin-	877
824	captioning .	jie Li, Kevin Lin, Zhe Gan, Zicheng Liu, Ce Liu,	878
825	Ron Mokady. 2021. Clipcap: Clip prefix for image	and Lijuan Wang. 2022. Git: A generative image-	879
826	captioning . <i>ArXiv</i> , abs/2111.09734 .	to-text transformer for vision and language . <i>ArXiv</i> ,	880
827	OpenAI. 2023.	abs/2205.14100 .	881
828	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	Ke Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho,	882
829	Jing Zhu. 2002. Bleu: a method for automatic evalu-	Aaron C. Courville, Ruslan Salakhutdinov, Richard S.	883
830	ation of machine translation . In <i>Proceedings of the</i>	Zemel, and Yoshua Bengio. 2015. Show, attend and	884
831	<i>40th Annual Meeting of the Association for Comput-</i>	tell: Neural image caption generation with visual	885
832	<i>ational Linguistics</i> , pages 311–318, Philadelphia,	attention . In <i>International Conference on Machine</i>	886
833	Pennsylvania, USA. Association for Computational	<i>Learning</i> .	887
834	Linguistics.	Shitong Xu. 2022. Clip-diffusion-lm: Apply diffusion	888
835	Dustin Podell, Zion English, Kyle Lacey, Andreas	model on image captioning .	889
836	Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna,	Ting Yao, Yingwei Pan, Yehao Li, Zhaofan Qiu, and	890
837	and Robin Rombach. 2023. Sdxl: Improving latent	Tao Mei. 2017. Boosting image captioning with at-	891
838	diffusion models for high-resolution image synthesis .	tributes . In <i>2017 IEEE International Conference on</i>	892
839	Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben	<i>Computer Vision (ICCV)</i> .	893
840	Mildenhall. 2022. Dreamfusion: Text-to-3d using 2d	Quanzeng You, Hailin Jin, Zhaowen Wang, Chen Fang,	894
841	diffusion . <i>arXiv</i> .	and Jiebo Luo. 2016. Image captioning with seman-	895
842	Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey	tic attention . <i>2016 IEEE Conference on Computer</i>	896
843	Chu, and Mark Chen. 2022. Hierarchical text-	<i>Vision and Pattern Recognition (CVPR)</i> , pages 4651–	897
844	conditional image generation with clip latents . <i>ArXiv</i> ,	4659.	898
845	abs/2204.06125 .	Hongyi Yuan, Zheng Yuan, Chuanqi Tan, Fei Huang,	899
846	Machel Reid, VincentJ. Hellendoorn, and Graham Neu-	and Songfang Huang. 2022. Seqdiffuseq: Text dif-	900
847	big. 2022. Diffuser: Discrete diffusion via edit-based	fusion with encoder-decoder transformers . <i>ArXiv</i> ,	901
848	reconstruction .	abs/2212.10325 .	902
849	Olaf Ronneberger, Philipp Fischer, and Thomas Brox.	Pengchuan Zhang, Xiujun Li, Xiaowei Hu, Jianwei	903
850	2015. U-net: Convolutional networks for biomedical	Yang, Lei Zhang, Lijuan Wang, Yejin Choi, and Jian-	904
851	image segmentation . <i>ArXiv</i> , abs/1505.04597 .	feng Gao. 2021. Vinvl: Revisiting visual representa-	905
852	Kai Shen, Zeqian Ju, Xu Tan, Yanqing Liu, Yichong	tions in vision-language models . In <i>2021 IEEE/CVF</i>	906
853	Leng, Lei He, Tao Qin, Sheng Zhao, and Jiang Bian.	<i>Conference on Computer Vision and Pattern Recog-</i>	907
854	2023. Naturalspeech 2: Latent diffusion models are	<i>nition (CVPR)</i> .	908
855	natural and zero-shot speech and singing synthesizers .	Tianyi Zhang, Varsha Kishore, Felix Wu, KilianQ. Wein-	909
856	<i>ArXiv</i> , abs/2304.09116 .	berger, and Yoav Artzi. 2019. Bertscore: Evaluat-	910
857	Jiaming Song, Chenlin Meng, and Stefano Ermon. 2020.	ing text generation with bert . <i>Cornell University -</i>	911
858	Denoising diffusion implicit models . <i>arXiv: Learn-</i>	<i>arXiv, Learning</i> .	912
859	ing , <i>arXiv: Learning</i> .	Zixin Zhu, Yixuan Wei, Jianfeng Wang, Zhe Gan, Zheng	913
860	Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya	Zhang, Le Wang, Gang Hua, Lijuan Wang, Zicheng	914
861	Sutskever. 2023. Consistency models . <i>ArXiv</i> ,	Liu, and Han Hu. 2022. Exploring discrete diffusion	915
862	abs/2303.01469 .	models for image captioning .	916
863	Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob		
864	Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz		
865	Kaiser, and Illia Polosukhin. 2017. Attention is all		
866	you need . In <i>NIPS</i> .		

A Additional Results

A.1 Generated Samples from COCO Dataset

Additional examples generated by our LaDiC model are presented in Fig. 10. It is shown that our model adeptly captures the main objects and their relationships in the depicted images. Simultaneously, the generated captions exhibit a high level of fluency.

A.2 Custom Generation

Utilizing the partially adding noise technique described in § 4.5, we observed that, unlike the unidirectional generation approach of AR models, our LaDiC model can effectively insert words into almost any position within a sentence. Fig. 11 offers additional examples to illustrate the generalization ability of this method.

A.3 Gradual Denosing Procedure during Inference

The inference of diffusion models involves gradually removing noise. To illustrate this process, we selected an image and showcased its caption generated at different time steps, as depicted in Fig. 8. Notably, the main objects initially emerge, and subsequently, more details are incrementally added, resulting in increasingly fluent sentences. This characteristic also serves as inspiration for our Back&Refine Technique, as discussed in § 3.4.

B Additional Details in Experiments

B.1 Details about Experiments on Image Paragraph Captioning

The objective of image paragraph captioning is to generate comprehensive paragraphs that describe images, providing detailed and cohesive narratives. This concept was initially introduced in (Krause et al., 2016), where the authors proposed a dataset comprising 19,551 images from MS COCO (Lin et al., 2014) and Visual Genome (Krishna et al., 2016), each annotated with a paragraph description. An illustrative example is presented in Fig. 9.

To assess our model’s ability to consider holistic context, we compare the performance of our model and BLIP on this task. For our model, we extend the predefined length to 60 and conduct training over 120 epochs. For BLIP, we fine-tune from BLIP_{base} using the same number of epochs and an initial learning rate of 1e-5. Subsequently, we evaluate the results using BLEU on the test set. In

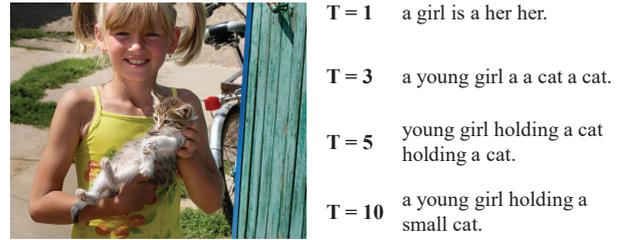


Figure 8: Gradual denosing process of diffusion models.



Sentences

- 1) A girl is eating donuts with a boy in a restaurant
- 2) A boy and girl sitting at a table with doughnuts.
- 3) Two kids sitting a coffee shop eating some frosted donuts
- 4) Two children sitting at a table eating donuts.
- 5) Two children eat doughnuts at a restaurant table.

Paragraph

Two children are sitting at a table in a restaurant. The children are one little girl and one little boy. The little girl is eating a pink frosted donut with white icing lines on top of it. The girl has blonde hair and is wearing a green jacket with a black long sleeve shirt underneath. The little boy is wearing a black zip up jacket and is holding his finger to his lip but is not eating. A metal napkin dispenser is in between them at the table. The wall next to them is white brick. Two adults are on the other side of the short white brick wall. The room has white circular lights on the ceiling and a large window in the front of the restaurant. It is daylight outside.

Figure 9: An example from image paragraph captioning dataset.

the case of BLIP, the maximum length is set to 60, and the number of beams is 5 during inference.

B.2 Human Evaluation

As a generative task, in addition to automatic metrics, it is imperative to assess results through human subjective evaluation. To this end, we utilize MOS (Mean Opinion Score) as our metric and enlist the feedback of 20 experienced volunteers, who were tasked with rating results on a scale of 1-5. They evaluated the results from three perspectives: fluency, accuracy, and conciseness. Fluency gauges the quality of generated captions in terms of language, accuracy assesses whether the main objects and actions in the caption accurately reflect the pic-

978 ture, and conciseness evaluates the extent to which
 979 generative captions are informative and succinct,
 980 avoiding unnecessary details.

981 To ensure evaluation quality, we randomly sam-
 982 pled 10 pictures from the COCO dataset and gener-
 983 ated corresponding captions for SCD-Net, BLIP¹,
 984 and our LaDiC model. Subsequently, we shuffled
 985 the three captions and required volunteers to rate
 986 them. To guarantee the reliability of the evaluation,
 987 we randomly selected 2 evaluators and calculated
 988 their correlation on each metric. This procedure
 989 was repeated 5 times, and all results were found to
 990 be satisfactory.

991 As depicted in Tab. 2, our model surpasses the
 992 previous diffusion-based state-of-the-art SCD-Net
 993 in all aspects, achieving comparable results with
 994 BLIP. A slight decrease in text quality compared to
 995 BLIP may be attributed to the substantial training
 996 data used in BLIP’s pretraining.

997 C More Hyperparameters

998 We list more hyperparameters for LaDiC model in
 999 Tab. 5.

1000 D Mathematical Details for Diffusion 1001 Models

1002 The training flow of the diffusion models is di-
 1003 vided into two phases: the forward diffusion pro-
 1004 cess and the backward denoising process. Given
 1005 a data point sampled from a real data distribu-
 1006 tion $x_0 \sim q(x)^2$, we define a forward diffusion
 1007 process in which Gaussian noise is incrementally
 1008 added to the sample, generating a sequence of
 1009 noisy samples x_1, \dots, x_T . The noise scales are
 1010 controlled by a variance schedule $\beta_t \in (0, 1)$,
 1011 and the density is expressed as $q(x_t|x_{t-1}) =$
 1012 $\mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t\mathbf{I})$. Based on the reparam-
 1013 eterization trick (Ho et al., 2020a), a nice property
 1014 of the above process is that we can sample at any

¹For BLIP, we utilized the following page for convenient inference: <https://replicate.com/salesforce/blip>.

²We follow the notation and derivation process of <https://lilianweng.github.io/posts/2021-07-11-diffusion-models>.

Hyperparameters	Values
<i>Training</i>	
Batch size	64*8(GPUs)
Epoch	60
Peak Learning rate	5e-5
Learning rate schedule	Linear
Warmup ratio	0.1
Optimizer	AdamW
β_1	0.9
β_2	0.999
<i>Inference</i>	
Method	DDIM
Sampling Criterion	Minimum Bayes Risk
<i>Diffusion Process</i>	
Diffusion steps	1000
β minimum	0.0001
β maximum	0.02
β schedule	Cosine
Classifier free probability	0.1
Classifier free weight	1
Self-conditioning probability	0.5
<i>Loss</i>	
λ	0.2
Loss type	l_2
<i>Image Encoder</i>	
Image size	224
Image Encoder	BLIP _{base}
<i>Diffuser Module</i>	
Sequence length	24
Hidden size	768
Layers	12
FFN size	3072
Attention heads	16

Table 5: More hyperparameters of our LaDiC model.

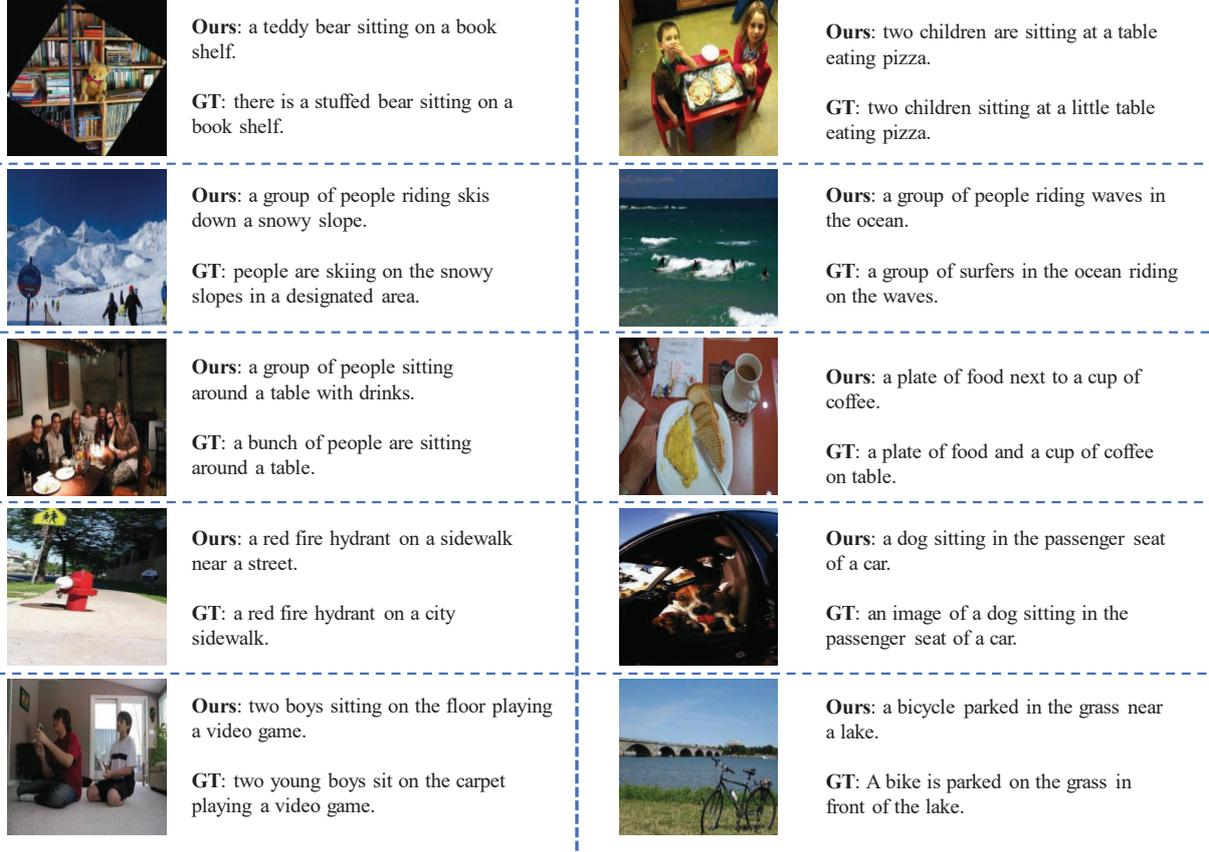


Figure 10: More examples generated by our model on COCO datasets.

arbitrary time step in a closed form:

$$\begin{aligned}
 x_t &= \sqrt{a_t}x_{t-1} + \sqrt{1 - \alpha_t}\epsilon_{t-1} \\
 &= \sqrt{a_t}(\sqrt{a_{t-1}}x_{t-2} + \sqrt{1 - \alpha_{t-1}}\epsilon_{t-2}) \\
 &\quad + \sqrt{1 - \alpha_t}\epsilon_{t-1} \\
 &= \sqrt{a_t a_{t-1}}x_{t-2} + (\sqrt{a_t(1 - \alpha_{t-1})}\epsilon_{t-2} \\
 &\quad + \sqrt{1 - \alpha_t}\epsilon_{t-1}) \\
 &= \sqrt{a_t a_{t-1}}x_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}}\bar{\epsilon}_{t-2} \\
 &= \dots \\
 &= \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon.
 \end{aligned}$$

where $\alpha_t = 1 - \beta_t$ and $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$. Thus:

$$q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, \sqrt{1 - \bar{\alpha}_t}\mathbf{I}), \quad (3)$$

Furthermore, from this equation, it becomes evident that as $T \rightarrow \infty$, x_T converges to an isotropic Gaussian distribution, aligning with the initial condition during inference.

However, obtaining the closed form of the reversed process $q(x_{t-1}|x_t)$ is challenging. Notably, if β_t is sufficiently small, the posterior will also be Gaussian. In this context, we can train a model

$p_\theta(x_{t-1}|x_t)$ to approximate these conditional probabilities:

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t)),$$

where $\mu_\theta(x_t, t)$ and $\Sigma_\theta(x_t, t)$ are parameterized by a denoising network f_θ like U-Net (Ronneberger et al., 2015) or Transformer (Vaswani et al., 2017). Similar to VAE (Kingma and Welling, 2013), we can derive the variational lower bound to optimize the negative log-likelihood of input x_0 (Ho et al., 2020b), :

$$\begin{aligned}
 \mathcal{L}_{\text{vib}} &= \mathbb{E}_q[\underbrace{D_{\text{KL}}(q(x_t|x_0)||p_\theta(x_T))}_{\mathcal{L}_T}] - \underbrace{\log p_\theta(x_0|x_1)}_{\mathcal{L}_0} \\
 &\quad + \mathbb{E}_q[\sum_{t=2}^T \underbrace{D_{\text{KL}}(q(x_{t-1}|x_t, x_0)||p_\theta(x_{t-1}|x_t))}_{\mathcal{L}_{t-1}}].
 \end{aligned}$$

With an additional condition on x_0 , the posterior of the forward process $q(x_{t-1}|x_t, x_0)$ can be calculated using Bayes theorem. Then in (Ho et al.,



Input: there is a boy [UNK]
[UNK] [UNK] cows
Output: there is a boy standing
by several cows

Input: [UNK] [UNK] [UNK]
[UNK] [UNK] on the grass
Output: An old blue car parked
on the grass

Input: a [UNK] [UNK] is
holding a [UNK] in her hand
Output: a young girl is
holding a cat in her hand.

Input: [UNK] [UNK] [UNK] [UNK]
[UNK] in front of a computer.
Output: a cup of coffee sitting in
front of a computer.

Figure 11: More examples of custom generation.

2020b) they derive:

$$\begin{aligned}
 L_t &= \mathbb{E}_{x_0, \epsilon} \left[\frac{1}{2 \|\Sigma_\theta(x_t, t)\|_2^2} \|\tilde{\mu}_t(x_t, x_0) - \mu_\theta(x_t, t)\|^2 \right] \\
 &= \mathbb{E}_{x_0, \epsilon} \left[\frac{1}{2 \|\Sigma_\theta(x_t, t)\|_2^2} \left\| \frac{1}{\sqrt{a_t}} \left(x_t - \frac{\beta_t}{\sqrt{1-a_t}} \epsilon_t \right) \right. \right. \\
 &\quad \left. \left. - \frac{1}{\sqrt{a_t}} \left(x_t - \frac{\beta_t}{\sqrt{1-a_t}} \epsilon_\theta(x_t, t) \right) \right\|^2 \right] \\
 &= \mathbb{E}_{x_0, \epsilon} \left[\frac{\beta_t^2}{2 \alpha_t (1 - \bar{\alpha}_t) \|\Sigma_\theta\|_2^2} \|\epsilon_t - \epsilon_\theta(x_t, t)\|^2 \right] \\
 &= \mathbb{E}_{x_0, \epsilon} \left[\frac{\beta_t^2}{2 \alpha_t (1 - \bar{\alpha}_t) \|\Sigma_\theta\|_2^2} \times \right. \\
 &\quad \left. \|\epsilon_t - \epsilon_\theta(\sqrt{\alpha_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon_t, t)\|^2 \right]
 \end{aligned}$$

Removing the coefficients, a much more simple DDPM learning objective can be obtained:

$$\mathcal{L}_{\text{simple}} = \sum_{t=1}^T \mathbb{E}_q \left[\|\epsilon_t(x_t, x_0) - \epsilon_\theta(x_t, t)\|^2 \right],$$

where ϵ_t is the noise added in original data x_0 . Applied to textual data, (Li et al., 2022b) introduces an even simpler architecture to train a network to predict x_0 directly, with the loss function defined as $L = \|x_0 - f_\theta(x_t, t)\|$.

During inference, the reverse process commences by sampling noise from a Gaussian distribution $p(x_T) = \mathcal{N}(x_T; 0, \mathbf{I})$ and iteratively denoising it using $p_\theta(x_{t-1}|x_t)$ until reaching x_0 . In DDIM (Song et al., 2020), a general form is derived from Equation 3.

$$\begin{aligned}
 x_{t-1} &= \sqrt{\bar{\alpha}_{t-1}} x_0 + \sqrt{1 - \bar{a}_{t-1}} \epsilon_{t-1} \\
 &= \sqrt{\bar{\alpha}_{t-1}} x_0 + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2} \epsilon_t \\
 &\quad + \sigma_t \epsilon \\
 &= \sqrt{\bar{\alpha}_{t-1}} x_0 + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2} \\
 &\quad \left(\frac{x_t - \sqrt{a_t} x_0}{\sqrt{1 - \bar{\alpha}_t}} \right) + \sigma_t \epsilon
 \end{aligned}$$

$$\begin{aligned}
 q_\sigma(x_{t-1}|x_t, x_0) &= \mathcal{N}(x_{t-1}; \sqrt{a_{t-1}} x_0 + \\
 &\quad \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2} \left(\frac{x_t - \sqrt{a_t} x_0}{\sqrt{1 - \bar{\alpha}_t}} \right), \sigma_t^2 \mathbf{I}).
 \end{aligned}$$

where $\sigma_t^2 = \eta \tilde{\beta}_t = \eta \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t$, allowing us to adjust η as a hyperparameter to control the sampling stochasticity. The special case of $\eta = 0$ renders the sampling process deterministic. This model is referred to as the denoising diffusion implicit model (DDIM). It is noteworthy that DDIM shares the same marginal distribution as DDPM. Consequently, during generation, we can sample only a subset of diffusion steps τ_1, \dots, τ_S , and the inference process becomes:

$$\begin{aligned}
 q_{\sigma, \tau}(\mathbf{x}_{\tau_i-1} | \mathbf{x}_{\tau_t}, \mathbf{x}_0) &= \mathcal{N}(\mathbf{x}_{\tau_i-1}; \sqrt{\bar{\alpha}_{\tau_i-1}} \mathbf{x}_0 \\
 &\quad + \sqrt{1 - \bar{\alpha}_{\tau_i-1} - \sigma_{\tau_i}^2} \frac{\mathbf{x}_{\tau_t} - \sqrt{\bar{\alpha}_{\tau_t}} \mathbf{x}_0}{\sqrt{1 - \bar{\alpha}_{\tau_t}}}, \sigma_{\tau_i}^2 \mathbf{I})
 \end{aligned}$$

which, significantly reduces inference latency.