

GALIP: Generative Adversarial CLIPs for Text-to-Image Synthesis

Ming Tao¹ Bing-Kun Bao^{1,2*} Hao Tang³ Changsheng Xu^{2,4,5}

¹Nanjing University of Posts and Telecommunications ²Peng Cheng Laboratory

³CVL, ETH Zürich ⁴University of Chinese Academy of Sciences

⁵NLPR, Institute of Automation, CAS

Abstract

Synthesizing high-fidelity complex images from text is challenging. Based on large pretraining, the autoregressive and diffusion models can synthesize photo-realistic images. Although these large models have shown notable progress, there remain three flaws. 1) These models require tremendous training data and parameters to achieve good performance. 2) The multi-step generation design slows the image synthesis process heavily. 3) The synthesized visual features are difficult to control and require delicately designed prompts. To enable high-quality, efficient, fast, and controllable text-to-image synthesis, we propose Generative Adversarial CLIPs, namely GALIP. GALIP leverages the powerful pretrained CLIP model both in the discriminator and generator. Specifically, we propose a CLIP-based discriminator. The complex scene understanding ability of CLIP enables the discriminator to accurately assess the image quality. Furthermore, we propose a CLIP-empowered generator that induces the visual concepts from CLIP through bridge features and prompts. The CLIP-integrated generator and discriminator boost training efficiency, and as a result, our model only requires about 3% training data and 6% learnable parameters, achieving comparable results to large pretrained autoregressive and diffusion models. Moreover, our model achieves $\sim 120\times$ faster synthesis speed and inherits the smooth latent space from GAN. The extensive experimental results demonstrate the excellent performance of our GALIP. Code is available at <https://github.com/tobran/GALIP>.

1. Introduction

Over the last few years, we have witnessed the great success of generative models for various applications [4, 45]. Among them, text-to-image synthesis [3, 15, 18–21, 25, 28, 29, 33, 46, 48–50, 57] is one of the most appealing applications. It generates high-fidelity images according to given

*Corresponding Author

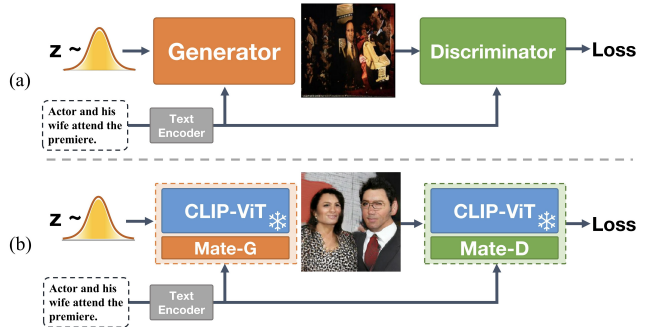


Figure 1. (a) Existing text-to-image GANs conduct adversarial training from scratch. (b) Our proposed GALIP conducts adversarial training based on the integrated CLIP model.

language guidance. Owing to the convenience of language for users, text-to-images synthesis has attracted many researchers and become an active research area.

Based on a large scale of data collections, model size, and pretraining, recently proposed large pretrained autoregressive and diffusion models, e.g., DALL-E [33] and LDM [35], show the impressive generative ability to synthesize complex scenes and outperform the previous text-to-image GANs significantly. Although these large pretrained generative models have achieved significant advances, they still suffer from three flaws. First, these models require tremendous training data and parameters for pretraining. The large data and model size brings an extremely high computing budget and hardware requirements, making it inaccessible to many researchers and users. Second, the generation of large models is much slower than GANs. The token-by-token generation and progressive denoising require hundreds of inference steps and make the generated results lag the language inputs seriously. Third, there is no intuitive smooth latent space as GANs, which maps meaningful visual attributes to the latent vector. The multi-step generation design breaks the synthesis process and scatters the meaningful latent space. It makes the synthesis process require delicately designed prompts to control.

To address the above limitations, we rethink Generative Adversarial Networks (GAN). GANs are much faster than autoregressive and diffusion models and have smooth latent

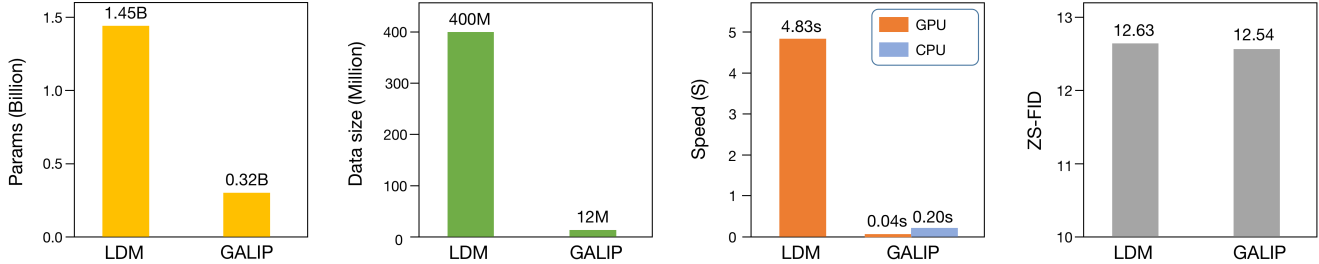


Figure 2. Comparing with Latent Diffusion Models (LDM) [35], our GALIP achieves comparable zero-shot Fréchet Inception Distance (ZS-FID) with measly 320M parameters (0.08B trainable parameters + 0.24B frozen CLIP parameters) and 12M training data. Furthermore, our GALIP only requires 0.04s to synthesize one image which is $\sim 120\times$ faster than LDM. Speed is calculated on NVIDIA 3090 GPU and Intel Xeon Silver 4314 CPU.

space, which enables more controllable synthesis. However, GAN models are known for potentially unstable training and less diversity in the generation [6]. It makes current text-to-image GANs suffer from unsatisfied synthesis quality under complex scenes.

In this work, we introduce the pretrained CLIP [30] into text-to-image GANs. The large pretraining of CLIP brings two advantages. First, it enhances the complex scene understanding ability. The pretraining dataset has many complex images under different scenes. Armed with the Vision Transformer (ViT) [8], the image encoder can extract informative and meaningful visual features from complex images to align the corresponding text descriptions after adequate pretraining. Second, the large pretraining dataset also enables excellent domain generalization ability. It contains various kinds of images, *e.g.*, photos, drawings, cartoons, and sketches, collected from a variety of publicly available sources. The various images make the CLIP model can map different kinds of images to the shared concepts and enable impressive domain generalization and zero-shot transfer ability. These two advantages of CLIP, complex scene understanding and domain generalization ability, motivate us to build a more powerful text-to-image model.

We propose a novel text-to-image generation framework named Generative Adversarial CLIPs (GALIP). As shown in Figure 1, the GALIP integrates the CLIP model [30] in both the discriminator and generator. To be specific, we propose the CLIP-based discriminator and CLIP-empowered generator. The CLIP-based discriminator inherits the complex scene understanding ability of CLIP [30]. It is composed of a frozen ViT-based CLIP image encoder (CLIP-ViT) and a learnable mate-discriminator (Mate-D). The Mate-D is mated to the CLIP-ViT for adversarial training. To retain the knowledge of complex scene understanding in the CLIP-ViT, we freeze its weights and collect the predicted CLIP image features from different layers. Then, the Mate-D further extracts informative visual features from collected CLIP features to distinguish the synthesized and real images. Based on the complex scene understanding ability of CLIP-ViT and the continuous analysis of Mate-

D, the CLIP-based discriminator can assess the quality of generated complex images more accurately.

Furthermore, we propose the CLIP-empowered generator, which exerts the domain generalization ability of CLIP [30]. It is hard for the generator to synthesize complex images directly. Some works employ sketch [10] and layout [20, 22] as bridge domains to alleviate the difficulty. However, such a design requires additional labeled data. Different from these works, the excellent domain generalization of CLIP [30] motivates us that there may be an implicit bridge domain, which is easier to synthesize but can be mapped to the same visual concepts through the CLIP-ViT. Thus, we design the CLIP-empowered generator. It is composed of a frozen CLIP-ViT and a learnable mate-generator (Mate-G). The Mate-G first predicts the implicit bridge features from text and noise. Then the bridge feature will be mapped to the visual concepts through CLIP-ViT. Furthermore, we add some text-conditioned prompts to the CLIP-ViT for task adaptation. The predicted visual concepts close the gap between text features and target images which enhances the complex image synthesis ability.

As shown in Figure 2, the proposed GALIP achieves $\sim 120\times$ faster synthesis speed and comparable synthesis ability based on significantly smaller trainable parameters and training data.

Overall, our contributions can be summarized as follows:

- We propose an efficient, fast, and more controllable model for text-to-image synthesis that can synthesize high-quality complex images.
- We propose the CLIP-based discriminator, which assesses the quality of generated complex images more accurately.
- We propose the CLIP-empowered generator, which synthesizes images based on text features and predicted CLIP visual features.
- Extensive experiments demonstrate that the proposed GALIP can achieve comparable performance with large pertaining models based on significantly smaller computational costs.

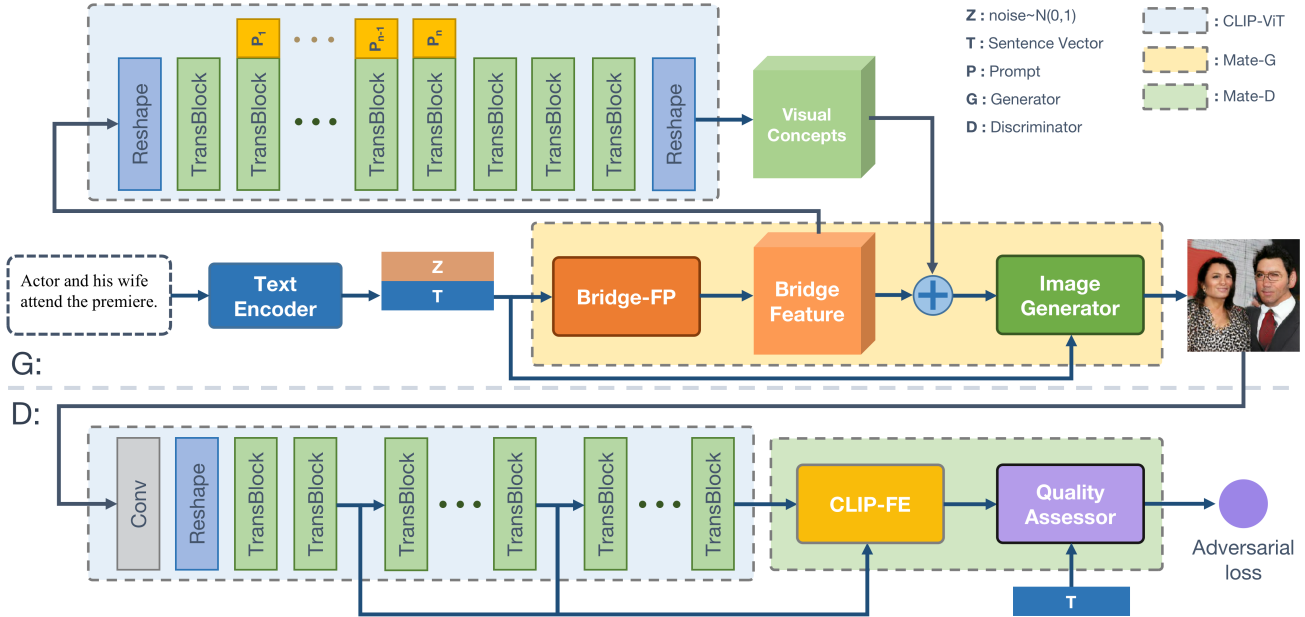


Figure 3. The architecture of the proposed GALIP for text-to-image synthesis. Armed with the CLIP-based discriminator and CLIP-empowered generator, our model can synthesize more realistic complex images.

2. Related Work

Text-to-Image GANs. GAN-INT-CLS [34] first adopted conditional GANs to synthesize images from text descriptions. To enable higher resolution synthesis, the StackGAN [54, 55], AttnGAN [48], and DM-GAN [57] stacks multiple generators and discriminators. Tao *et al.* [42] proposed a simpler yet effective text-to-image framework called DF-GAN that enables one-stage high-resolution generation. LAFITE [56] introduces CLIP text-image contrastive loss for text-to-image training and shows large improvements on CC3M [40].

Text-to-Image Large Models. Recently, large pretrained autoregressive and diffusion models have shown impressive performance on text-to-image synthesis. DALL-E [33], CogView [6], and M6 [23] leverage VQ-VAE [43] or VQ-GAN [9] to tokenize the images into discrete image tokens. Then they take the word tokens and image tokens together to pre-train a large unidirectional transformer for an autoregressive generation. Parti [51] proposes a sequence-to-sequence autoregressive model to treat text-to-image synthesis as a translation task. Cogview2 [7] employs hierarchical transformers and local parallel autoregressive generation for faster autoregressive image generation. Some works try to employ the diffusion model [5, 13, 14, 26, 41] to overcome the slow generation defect of the autoregressive model. VQ-Diffusion [11] combines the VQ-VAE [43] and diffusion model [14, 26] to eliminate the unidirectional bias and avoids accumulated prediction errors. GLIDE [27] applies guided diffusion to the problem of text-conditional image synthesis. DALL-E2 [32] combines the CLIP represen-

tation and diffusion model to make a CLIP decoder. Latent Diffusion Models (LDM) [35] apply the diffusion model in the latent space to enable the training on limited computational resources while retaining image quality. The particular text-to-image LDM is Stable Diffusion [36], which is a favorite open-source project and provides an easy-to-use interface. Imagen [38] introduces the large language model [31] to provide high-quality text features and proposes an Efficient U-Net for diffusion models.

3. Generative Adversarial CLIPs

In this paper, we propose a novel framework for text-to-image synthesis named Generative Adversarial CLIPs (GALIP). To synthesize high-quality complex images, we propose: (i) a novel CLIP-based discriminator that inherits the complex scene understanding ability of CLIP [30] for more accurate image quality assessment. (ii) a novel CLIP-empowered generator that exerts the domain generalization ability of CLIP [30] and induces the CLIP visual concepts to close the gap between text and image features. In the following of this section, we first present the overall structure of our GALIP. Then, we introduce the CLIP-based discriminator and CLIP-empowered generator in detail.

3.1. Model Overview

As shown in Figure 3, the proposed GALIP is composed of a CLIP text encoder, a CLIP-based discriminator, and a CLIP-empowered generator. The pretrained CLIP text encoder takes the text description and yields a global sentence vector T . After the text-encoder is the CLIP-empowered

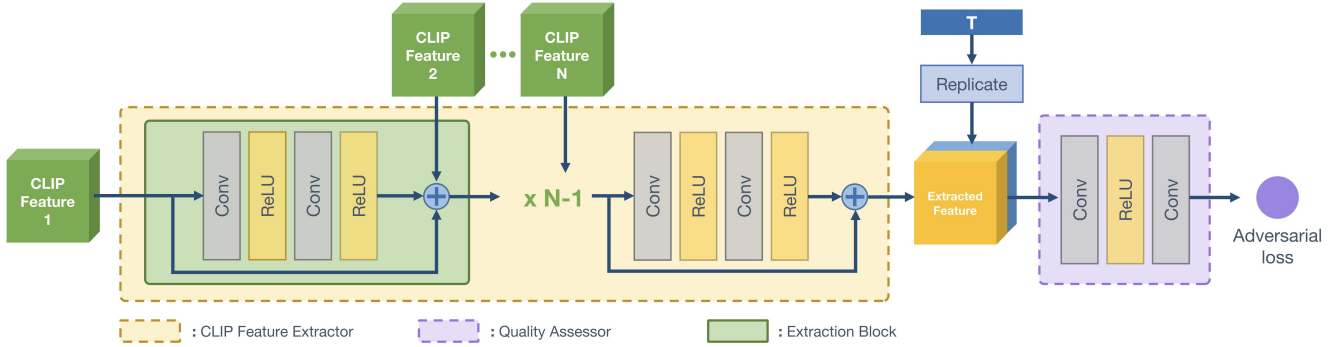


Figure 4. The architecture of the proposed Mate-D for text-to-image synthesis. It further extracts informative visual features from collected CLIP features and assesses the image quality more accurately.

generator and CLIP-based discriminator under the GAN framework. The CLIP-empowered generator is composed of a frozen CLIP-ViT and a mate generator (Mate-G). There are three main modules in the Mate-G, the bridge feature predictor (Bridge-FP), the prompt predictor, and the image generator. The CLIP-empowered generator has two inputs, the sentence vector T encoded from the text encoder and the noise vector Z sampled from the Gaussian distribution. The noise vector ensures the diversity of the synthesized images. In the CLIP-empowered generator, the sentence vector and noise are first fed into the bridge feature predictor. The bridge feature predictor translates the sentence vector and noise to the bridge feature for the CLIP-ViT. Furthermore, we add several text-conditioned prompts to the transformer blocks (TransBlock) in CLIP-ViT for task adaptation. Finally, the image generator takes the predicted visual concepts, bridge features, sentences, and noise vectors to synthesize high-quality images.

The CLIP-based discriminator is composed of a frozen CLIP-ViT and a mate discriminator (Mate-D). The CLIP-ViT converts images into image feature through a convolution layer and a series of transformer blocks. The CLIP feature extractor (CLIP-FE) in Mate-D collects the image features from different layers in CLIP-ViT. Then it further extracts informative visual features from collected CLIP features for the quality assessor. Lastly, an adversarial loss will be predicted by the quality assessor based on the extracted informative features and sentence vectors. By distinguishing synthesized images from real ones, the discriminator promotes the generator to synthesize higher-quality images.

3.2. CLIP-based Discriminator

In this section, we detailed the proposed CLIP-based discriminator, which is composed of a frozen CLIP-ViT and a Mate-D. The CLIP-based discriminator inherits the complex scene understanding ability from the frozen CLIP-ViT. Furthermore, we propose the Mate-D, which is mated to the CLIP-ViT to further extract informative visual features and distinguish real and synthesized images. The CLIP-ViT and Mate-D enable the discriminator to assess the quality of

generated complex images more accurately.

As shown in Figure 4, the Mate-D consists of a CLIP-FE and a quality assessor. To fully utilize the knowledge of complex scene understanding in CLIP-ViT, the CLIP-FE takes the CLIP image features from multilayers. There are N CLIP features collected for the CLIP-FE. We name them CLIP Feature 1 to N , which are collected from shallow to deep layers in CLIP-ViT. To further extract informative visual features from these CLIP features, we design a CLIP-FE. It contains a sequence of extraction blocks, and each block contains two convolution layers and two ReLU active functions. And the extracted image feature is summed with the shortcut and the next CLIP feature. There are $N - 1$ extraction blocks stacked in CLIP-FE. Since the CLIP feature N is only added to the processed image features in the last extraction block. To fuse the CLIP feature N , we append two convolution layers without the CLIP feature addition behind. The CLIP-FE extracts informative visual features for the quality assessor. Then the sentence vector is replicated and concatenated with the extracted image features. An adversarial loss is predicted by two convolution layers to evaluate the image quality. Furthermore, to stabilize the adversarial learning process of Mate-D, we apply the matching-aware gradient penalty (MAGP) [42] on the collected CLIP features and corresponding text features.

Based on the complex scene understanding ability of CLIP-ViT, the CLIP-based discriminator can extract more informative visual features from complex images. The higher-quality extracted visual features make it easier for the discriminator to detect unreal image parts, which improves the discriminative efficiency, thus prompting the generator to generate more realistic images.

3.3. CLIP-empowered Generator

In this section, we detail the proposed CLIP-empowered generator, which is composed of a frozen CLIP-ViT and a Mate-G. The CLIP-empowered generator exerts the domain generalization ability of the CLIP-ViT. Furthermore, we propose the Mate-G, which is mated to the CLIP-ViT to induce useful visual features from the CLIP-ViT and gen-

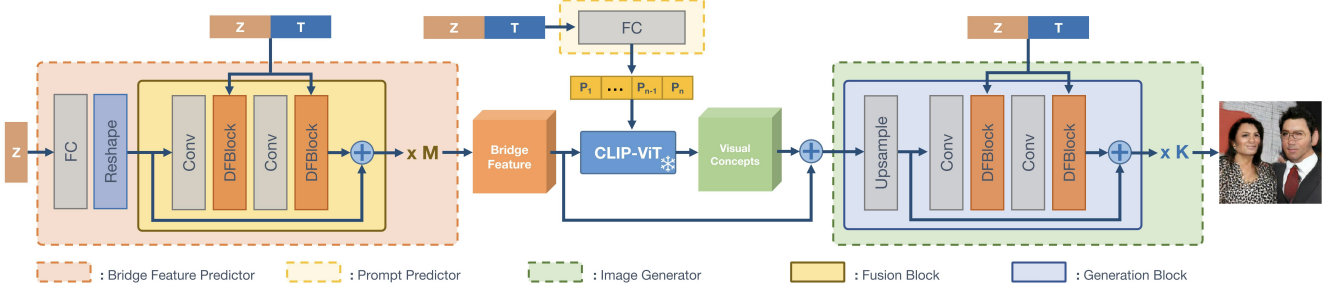


Figure 5. The architecture of the proposed CLIP-empowered generator for text-to-image synthesis. Armed with bridge feature predictor and prompt predictor, it can induce meaningful visual concepts from the frozen CLIP-ViT for image synthesis.

erate images from text and induced visual features. The Mate-G consists of a Bridge Feature Predictor (Bridge-FP), a prompt predictor, a frozen CLIP-ViT, and an image generator (see Figure 3). We detail them next.

Bridge Feature Predictor. The structure of the Bridge-FP is shown in Figure 5, as highlighted by the red dashed box. The Bridge-FP consists of an FC (Fully-Connected) layer and M fusion blocks (F-BLKs). The input noise is fed into the FC layer and reshaped to $(7, 7, 64)$ as an initial bridge feature. The initial bridge feature output by the FC layer still contains a lot of noise. Therefore, we apply a sequence of F-BLKs to fuse text information and make it more meaningful. The F-BLK is composed of two convolution layers (Conv) and two deep text-image fusion blocks (DFBlock) [42]. The DFBlock has shown its effectiveness in fusing text and image features through stacked affine transformations. Thus, we adopt it to fuse text features and intermediate bridge features. There is a shortcut addition in F-BLK for effective information propagation and gradient back-propagation. Through the Bridge-FP, the sentence and noise vectors will be translated to the bridge feature, which is adjusted to induce meaningful visual concepts from CLIP-ViT.

Prompt Predictor. The CLIP-ViT is pretrained to predict visual features from image data. There is a large gap between text and image data. To alleviate the difficulty of bridge feature translation from text features, we employ prompt tuning [16], which has shown effectiveness on domain transferring for ViT. We design a prompt predictor, which predicts prompts based on sentence and noise vectors through an FC layer. The predicted text-conditioned prompts are appended behind the visual patch embeddings in CLIP-ViT. Furthermore, we find that it is better not to add prompts to the last few layers in CLIP-ViT. The last few layers summarize the visual features and output the last image representations. The prompt predicted from text and noise in the last few layers may defect its performance.

Image Generator. The image generator consists of K generation blocks (G-BLKs). We sum the predicted visual concepts and bridge features through shortcut addition for effective information propagation and gradient back-

propagation. The image generator receives the summed visual features as input and fuses sentence and noise vectors through the DFBlocks [42] in each G-BLK. The intermediate image features grow larger during the generation process by the upsample layers. Finally, the image features are converted into high-resolution RGB images.

3.4. Objective Functions

To stabilize the training process of adversarial learning, we employ the hinge loss [52] and one-way discriminator [42]. Finally, the whole formulation of our GALIP is shown as follows:

$$\begin{aligned}
 L_D &= -\mathbb{E}_{x \sim \mathbb{P}_r}[\min(0, -1 + D(C(x), e))] \\
 &\quad - (1/2)\mathbb{E}_{G(z, e) \sim \mathbb{P}_g}[\min(0, -1 - D(C(G(z, e)), e))] \\
 &\quad - (1/2)\mathbb{E}_{x \sim \mathbb{P}_{mis}}[\min(0, -1 - D(C(x), e))] \\
 &\quad + k\mathbb{E}_{x \sim \mathbb{P}_r}[(\|\nabla_{C(x)} D(C(x), e)\| + \|\nabla_e D(C(x), e)\|)^p], \\
 L_G &= -\mathbb{E}_{G(z, e) \sim \mathbb{P}_g}[D(C(G(z, e)), e)] \\
 &\quad - \lambda\mathbb{E}_{G(z, e) \sim \mathbb{P}_g}[S(G(z, e), e)],
 \end{aligned} \tag{1}$$

where z is the noise vector sampled from Gaussian distribution; e is the sentence vector; G is the CLIP-empowered generator; D is the Mate-D; C is the frozen CLIP-ViT in CLIP-based discriminator; S represents the cosine similarity between the encoded visual and text features of CLIP; k and p are two hyper-parameters of gradient penalty; λ is the coefficients of the text-image similarity; \mathbb{P}_g , \mathbb{P}_r , \mathbb{P}_{mis} denote the synthetic data distribution, real data distribution, and mismatching data distribution, respectively.

4. Experiments

In this section, we introduce the datasets, training details, and evaluation metrics employed in our experiments, then evaluate our proposed GALIP and its variants quantitatively and qualitatively.

Datasets. We conduct experiments on four challenging datasets: CUB bird [44], COCO [24], CC3M [40], and CC12M [2]. For the CUB bird dataset, there are 11,788 images belonging to 200 bird species, with each image corresponding to ten language descriptions. The train and vali-



Figure 6. Examples of images synthesized by LAFITE [56], VQ-Diffusion [11], and our proposed GALIP conditioned on text descriptions from the test set of CUB and COCO datasets.

dation splits of the CUB bird dataset are implied as previous works did [42, 48, 54, 55, 57]. Since there are various shapes, colors, and postures of birds in the CUB dataset, it is always employed to evaluate the performance of fine-grained content synthesis. For COCO dataset, it contains 80k images for training and 40k images for testing. Each image corresponds to 5 language descriptions. The image in the COCO dataset is complex and always contains multiple objects under different scenes. The COCO dataset is always employed in recent works to evaluate the performance of complex image synthesis. For CC3M and CC12M datasets, they are two large datasets that contain about 3 and 12 million text-image pairs. It is always adopted for pretraining and to evaluate the zero-shot performance of the text-to-image model.

Training and Evaluation Details. We choose the ViT-B/32 [30] model as the CLIP model in our GALIP. In the CLIP-based discriminator, the CLIP-FE collects the CLIP feature from 2nd, 5th, 9th layers in CLIP-ViT. There are two extraction blocks stacked in CLIP-FE. In the CLIP-empowered generator, the Bridge-FP contains 4 Fusion Blocks, and the image generator contains 6 generation blocks for 224×224 image synthesis. The prompt predictor predicts 8 prompts for TransBlocks 2 to 10 in CLIP-ViT. We conduct some ablation studies on these designs. The hyper-parameters of the discriminator k and p are set to 2 and 6 as [42]. The hyper-parameters of the generator λ are set to 4 for all the datasets. Furthermore, we employ the Adam optimizer [17] with $\beta_1=0.0$ and $\beta_2=0.9$ to train our model. According to the two timescale update rule (TTUR) [12], the learning rate is set to 0.0001 for the generator and 0.0004 for the discriminator. Following the previous text-to-image works [42, 47, 48, 57], we adopt the Fréchet Inception Distance (FID) [12] and CLIPSIM [47] to evaluate the image fidelity and text-image semantic consistency. All GALIP models are trained on 8×3090 GPUs. We train our GALIP for 0.5, 1.5, 2, and 3 days on CUB, COCO, CC3M, and CC12M datasets, respectively.

Table 1. The results of FID and CLIPSIM (CS) compared with the state-of-the-art methods on the test set of CUB and COCO.

Model	CUB		COCO	
	FID ↓	CS ↑	FID ↓	CS ↑
DM-GAN [57]	16.09	-	32.64	-
XMC-GAN [53]	-	-	9.30	-
DAE-GAN [37]	15.19	-	28.12	-
DF-GAN [42]	14.81	0.2920	19.32	0.2972
LAFITE [56]	14.58	0.3125	8.21	0.3335
VQ-Diffusion [11]	10.32	-	13.86	-
GALIP (Ours)	10.08	0.3164	5.85	0.3338

4.1. Quantitative Evaluation

To evaluate the performance of our GALIP, we compare the proposed model with several state-of-the-art methods [11, 37, 42, 53, 56, 57], which have achieved impressive results in text-to-image synthesis. The results are shown in Table 1. Compared with other leading models, our GALIP has a significant improvement on both CUB and COCO datasets. Especially, compared with the recently proposed LAFITE [56], which employs CLIP text-image contrastive loss for text-to-image training, our GALIP decreases the FID metric from 14.58 to 10.08 and improves the CLIPSIM (CS) from 0.3125 to 0.3164 on the CUB dataset. Furthermore, our GALIP decreases the FID of COCO from 8.21 to 5.85 significantly. Compared with VQ-diffusion [11], which adopts diffusion models for text-to-image synthesis, our GALIP also decreases FID from 10.32 to 10.08 on the CUB dataset and decreases the FID of COCO from 13.86 to 5.85 remarkably. The quantitative comparisons on CUB and COCO datasets demonstrate that our GALIP is more effective in synthesizing high-fidelity images, especially for complex image generation.

Moreover, we evaluate the zero-shot text-to-image synthesis ability of our GALIP. The results are shown in Table 2. Compared with LAFITE [56] trained on CC3M, our GALIP (CC3M) decreases FID from 26.94 to 16.12 signif-



Figure 7. Text-to-Image samples from GALIP (CC12M) and Latent Diffusion (LAION-400M) [35, 36]. We sample 16 images from each given text description, and randomly select one as the final generation result

Table 2. We compare the performance of large pretrained autoregressive models (AR), diffusion models (DF), and GANs under zero-shot setting on the COCO test dataset.

Model	Type	Param [B]	Data size [M]	ZS-FID ↓
DALL-E [33]	AR	12	250	27.5
Cogview [6]	AR	4	30	27.1
Cogview2 [7]	AR	6	30	24.0
Parti-350M [51]	AR	0.35	>800	14.10
Parti-20B [51]	AR	20	>800	7.23
GLIDE [27]	DF	5	250	12.24
LDM [35]	DF	1.45	400	12.63
DALL-E 2 [32]	DF	6.5	250	10.39
Imagen [38]	DF	7.9	860	7.27
eDiff-I [1]	DF	9.1	1000	6.95
LAFITE [56]	GAN	0.15+0.08	3	26.94
GALIP (CC3M)	GAN	0.24+0.08	3	16.12
GALIP (CC12M)	GAN	0.24+0.08	12	12.54

icantly. It demonstrates that integrating the CLIP model in the generator and discriminator is more effective than only introducing the CLIP loss for the GAN model. Compared with autoregressive models (AR) and diffusion models (DF) which are pretrained with much larger model sizes and datasets, our GALIP also achieves competitive performance. Especially, compared with LDM [35] which is one of the most important open-source large pretrained models, our GALIP achieves better performance even with much smaller model parameters and data size. Furthermore, as shown in Figure 2, our GALIP only requires 0.04s to generate one image which is $\sim 120\times$ faster than LDM [35]. Besides, our GALIP can be inference on the CPU fastly without other acceleration settings. This significantly reduces the hardware requirements of users. In addition, the computational cost to pretrain our GALIP is quite less than these large pretrained autoregressive and diffusion models. The GALIP of CC12M is only pretrained on 8×3090 GPUs for 3 days. But these models require hundreds of GPUs and many weeks to pre-train.

4.2. Qualitative Evaluation

To evaluate the visual quality of synthesized images, we first compare the images synthesized by LAFITE [56], VQ-Diffusion [11], and our GALIP which are trained on COCO

in Figure 6. Then, we compare our GALIP (CC12M) with LDM (LAION-400M) [35, 36] in Figure 7.

As shown in the 1st, 2nd, 4th and 5th columns of Figure 6, the birds synthesized by LAFITE [56] and VQ-Diffusion [11] contain break or wrong shapes. Moreover, both LAFITE [56] and VQ-Diffusion [11] lose some fine-grained visual features (e.g., 1st, 2nd, 5th and 6th columns), which makes the synthesized images lack details and look unreal. However, the images synthesized by our GALIP have correct object shapes and clear fine-grained contents.

The superiority is more obvious in complex COCO images, which contain various shapes and multiple objects. As the results are shown in the 7th, 8th, 9th, 10th columns of Figure 7, the LAFITE [56] and VQ-Diffusion [11] models cannot synthesize the right shape of “train”, “children”, “woman”, and “stuffed bear”. Furthermore, they also cannot synthesize the right visual concept of “showing off toy cell phone” and “sitting on a book shelf”. However, armed with the proposed CLIP-based D and CLIP-empowered G, our GALIP can cope with more strict visual requirements and synthesize various shapes of different objects (see 8th, 9th, 10th and 12th columns) and present the right visual concepts in synthesized images. We also observe that LAFITE [56] and VQ-Diffusion [11] also can not synthesize correct human facial features. For example, as shown in the 8th, 9th, 12th, they can not synthesize realistic human faces. But our GALIP can synthesize these features correctly.

Moreover, we compare the images synthesized by the LDM (LAION-400M) [35, 36] and our GALIP (CC12M) in Figure 7. As the results are shown in the 1st, 4th, 5th, 8th, 11th columns of Figure 7, the LDM does not generate the objects (“ghost”, “teddy bear”, “modem”, “person”, “model”) described in the texts, but our GALIP can synthesize these objects correctly. Also, our model can generate correct visual features such as “shining eyes”, “Blue Lighthouse”, “smiling statue”, and “surprised girl” in the 3rd, 6th, 7th, 10th columns. Furthermore, as shown in the 9th, 10th, and 12th columns of Figure 7, our GALIP keeps the superior performance of human face synthesis. The extensive quantitative evaluation results demonstrate the superi-

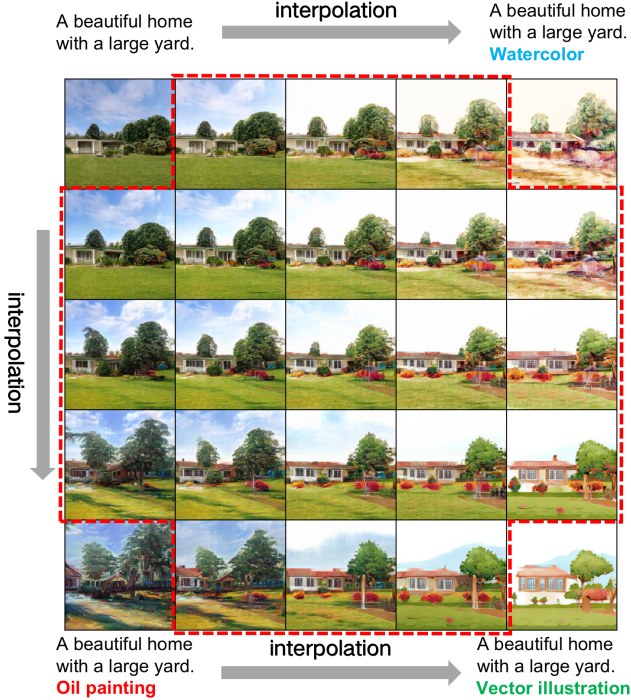


Figure 8. Images synthesized by interpolating four-sentence embeddings. Our GALIP supports gradual changes when interpolating sentence embeddings describing different image styles. It makes the degree of stylization of the image controllable and creates new styles by blending different styles.

ority and effectiveness of our proposed GALIP, which is able to generate high-fidelity, creative and complex images with various shapes and multiple objects.

Additionally, we conduct some experiments to show the smooth latent space of our GALIP. Current autoregressive and diffusion models are sensitive to input sentences. This instability makes users need to try a lot of prompts to get satisfied images. Differently, our GALIP inherits the smooth latent space from GAN, it enables gradual and smooth changes along with text changes. As shown in Figure 8, there is a smooth transition of synthesized images from top to bottom, left to right. The smooth latent space makes the degree of stylization of the image controllable. The users can fine-tune synthesized image styles like a style knob, and it also enables the users to create new styles by blending different image styles, as highlighted by the red dashed lines.

4.3. Ablation Study

To verify the effectiveness of each component in the proposed GALIP, we conduct ablation studies on the test set of the COCO dataset. The components being evaluated in this subsection include CLIP-based D (CD) and CLIP-empowered G (CG). We also further conduct ablation studies on Bridge-FP (BFP) and Prompt Predictor (PP) in CLIP-empowered G, and CLIP-FE (CFE) in CLIP-based D. Furthermore, we compare our CLIP-FE with CCM&CSM of

Table 3. The performance of different components of our model on the test set of COCO.

Architecture	FID ↓	CS ↑
Baseline	17.31	0.2996
Baseline w/ CD w/ CFE	7.92	0.3221
Baseline w/ CD w/ CCM&CSM	10.77	0.3123
Baseline w/ CD w/ BFP	6.52	0.3301
Baseline w/ CD w/ BFP w/ PP (GALIP)	5.85	0.3338
GALIP w/ CFE (2^{nd})	13.41	0.3015
GALIP w/ CFE (5^{th})	8.60	0.3145
GALIP w/ CFE (12^{th})	10.72	0.3104
GALIP w/ CFE ($2^{nd}, 5^{th}$)	6.70	0.3301
GALIP w/ CFE ($2^{nd}, 5^{th}, 12^{th}$)	6.61	0.3305
GALIP w/ CFE ($2^{nd}, 5^{th}, 9^{th}$)	5.85	0.3338
GALIP w/ CFE ($2^{nd}, 5^{th}, 8^{th}, 9^{th}$)	6.01	0.3305
GALIP w/ PP ($1^{st}-12^{th}$)	6.24	0.3320
GALIP w/ PP ($1^{st}-9^{th}$)	5.85	0.3338
GALIP w/ PP ($1^{st}-6^{th}$)	6.40	0.3310
GALIP w/ PP ($1^{st}-3^{th}$)	6.52	0.3305

Projected GAN [39], which yields a U-Net architecture to enable multi-scale feedback. In addition, we investigate the layer choice strategy of CLIP-FE and Prompt Predictor. The results on the COCO dataset are shown in Table 3.

Baseline. Our baseline is a one-stage text-to-image GAN [42]. It is composed of a CLIP text encoder and CNN-based generator and discriminator. And it generates complex images from sentence vectors directly.

Effect of CLIP-based D and CLIP-FE. The CLIP-based D decreases FID from 17.31 to 7.92 and improves CLIP-SIM (CS) from 0.2996 to 0.3221. The results demonstrate that the complex scene understanding ability of CLIP-ViT promotes the complex image synthesis ability significantly. Furthermore, we compared our CLIP-FE (CFE) with CCM&CSM [39]. Our CLIP-FE achieves better FID and CLIP-SIM. It shows that our CLIP-FE is more effective in extracting informative visual features from CLIP-ViT.

Effect of CLIP-empowered G and Bridge-FP. The CLIP-empowered G with Bridge-FP further decreases FID from 7.92 to 6.52 and improves CLIP-SIM from 0.3221 to 0.3301. It demonstrates that predicted bridge features and CLIP-ViT can enhance the complex image synthesis ability effectively.

Effect of Prompt Predictor. The proposed Prompt Predictor (PP) also decreases FID from 6.52 to 5.85 and improves CLIP-SIM from 0.3301 to 0.3338. The result demonstrates that the Prompt Predictor makes the CLIP-ViT more suitable for generation tasks and induces more meaningful features from CLIP-ViT to improve the generative ability.

CLIP Layer Selection. We find that the last few layers of CLIP-ViT defect the performance of CLIP-based D. The reason may be that the first layers of CLIP-ViT extract useful visual features and understand complex images, and the last layers focus on generalization ability to align



Figure 9. Illustration of failure cases. It is still hard for current GALIP to synthesize some imaginary images. Enlarging the model size and training data may improve image quality.

with high-level concepts in text features. The generalization ability may defect the performance of CLIP-based D because it reduces the differences between synthetic and real images and weakens the discriminator. Conversely, since CLIP-empowered G requires the generalization ability to map the bridge feature to meaningful visual features, adding prompts in the last few layers may defect the generalization ability. So we extract the CLIP features from 2^{nd} , 5^{th} , and 9^{th} layers in CLIP-based D, and add prompts to 1^{st} - 9^{th} layers. And we find that extracting more CLIP features does not lead to better performance.

4.4. Limitations

Our GALIP shows superiority in text-to-image synthesis, but some limitations should be considered in future studies. First, our model employs the CLIP model to provide text features for the generator and discriminator. However, current models [38] show that the generic large language models [31] (e.g., T5) improve the performance of text-to-image synthesis effectively. Replacing the CLIP text encoder with T5 may further improve the performance. Second, the model size and pretraining dataset are much smaller than other large pretrained models [1, 32, 35, 38, 51], it limits the synthesis ability of imaginary images (see Figure 9). Pretraining on a larger dataset with a larger model size may benefit the performance. We will try to address these limitations in our future work.

5. Conclusion

In this paper, we propose a novel Generative Adversarial CLIPs (GALIP) for text-to-image synthesis. Compared with previous models, our GALIP can synthesize higher-quality complex images. Moreover, we propose a CLIP-based discriminator and CLIP-empowered generator, which exerts the complex scene understanding and domain generalization ability of CLIP. Our GALIP achieves significant improvements on challenging datasets. Furthermore, current large models are pretrained for generative or understanding tasks. In this work, we integrate the understanding model (CLIP-ViT) into a generative model and achieve impressive results. It shows that there are some commonalities between understanding and generative models. This may be enlightening for building a general large model.

References

- [1] Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, Bryan Catanzaro, et al. ediffi: Text-to-image diffusion models with an ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*, 2022. 7, 9
- [2] Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12M: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *CVPR*, 2021. 5
- [3] Jun Cheng, Fuxiang Wu, Yanling Tian, Lei Wang, and Dapeng Tao. Rifegan: Rich feature generation for text-to-image synthesis from prior knowledge. In *CVPR*, 2020. 1
- [4] Wen-Huang Cheng, Sijie Song, Chieh-Yun Chen, Shintami Chusnul Hidayati, and Jiaying Liu. Fashion meets computer vision: A survey. *ACM CSUR*, 54(4):1–41, 2021. 1
- [5] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *NeurIPS*, 2021. 3
- [6] Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, et al. Cogview: Mastering text-to-image generation via transformers. In *NeurIPS*, 2021. 2, 3, 7
- [7] Ming Ding, Wendi Zheng, Wenyi Hong, and Jie Tang. Cogview2: Faster and better text-to-image generation via hierarchical transformers. *arXiv preprint arXiv:2204.14217*, 2022. 3, 7
- [8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021. 2
- [9] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *CVPR*, 2021. 3
- [10] Oran Gafni, Adam Polyak, Oron Ashual, Shelly Sheynin, Devi Parikh, and Yaniv Taigman. Make-a-scene: Scene-based text-to-image generation with human priors. In *ECCV*, 2022. 2
- [11] Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and Baining Guo. Vector quantized diffusion model for text-to-image synthesis. In *CVPR*, 2022. 3, 6, 7
- [12] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *NeurIPS*, 2017. 6
- [13] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *NeurIPS*, 2020. 3
- [14] Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. *JMLR*, 23:47–1, 2022. 3
- [15] Seunghoon Hong, Dingdong Yang, Jongwook Choi, and Honglak Lee. Inferring semantic layout for hierarchical text-to-image synthesis. In *CVPR*, 2018. 1

- [16] Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and Ser-Nam Lim. Visual prompt tuning. In *ECCV*, 2022. 5
- [17] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 6
- [18] Bowen Li, Xiaojuan Qi, Thomas Lukasiewicz, and Philip Torr. Controllable text-to-image generation. In *NeurIPS*, 2019. 1
- [19] Ruifan Li, Ning Wang, Fangxiang Feng, Guangwei Zhang, and Xiaojie Wang. Exploring global and local linguistic representation for text-to-image synthesis. *IEEE TMM*, 2020. 1
- [20] Wenbo Li, Pengchuan Zhang, Lei Zhang, Qiuyuan Huang, Xiaodong He, Siwei Lyu, and Jianfeng Gao. Object-driven text-to-image synthesis via adversarial training. In *CVPR*, 2019. 1, 2
- [21] Jiadong Liang, Wenjie Pei, and Feng Lu. Cpgan: Content-parsing generative adversarial networks for text-to-image synthesis. In *ECCV*, 2020. 1
- [22] Jiadong Liang, Wenjie Pei, and Feng Lu. Layout-bridging text-to-image synthesis. *arXiv preprint arXiv:2208.06162*, 2022. 2
- [23] Junyang Lin, Rui Men, An Yang, Chang Zhou, Ming Ding, Yichang Zhang, Peng Wang, Ang Wang, Le Jiang, Xianyan Jia, et al. M6: A chinese multimodal pretrainer. *arXiv preprint arXiv:2103.00823*, 2021. 3
- [24] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. 5
- [25] Bingchen Liu, Kunpeng Song, Yizhe Zhu, Gerard de Melo, and Ahmed Elgammal. Time: Text and image mutual-translation adversarial networks. In *AAAI*, 2021. 1
- [26] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *ICML*, 2021. 3
- [27] Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. In *ICML*, 2022. 3, 7
- [28] Tingting Qiao, Jing Zhang, Duanqing Xu, and Dacheng Tao. Learn, imagine and create: Text-to-image generation from prior knowledge. In *NeurIPS*, 2019. 1
- [29] Tingting Qiao, Jing Zhang, Duanqing Xu, and Dacheng Tao. Mirrorgan: Learning text-to-image generation by redescription. In *CVPR*, 2019. 1
- [30] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021. 2, 3, 6
- [31] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, Peter J Liu, et al. Exploring the limits of transfer learning with a unified text-to-text transformer. *JMLR*, 21(140):1–67, 2020. 3, 9
- [32] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. 3, 7, 9
- [33] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *ICML*, 2021. 1, 3, 7
- [34] Scott Reed, Zeynep Akata, Xinchun Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee. Generative adversarial text to image synthesis. In *ICML*, 2016. 3
- [35] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 1, 2, 3, 7, 9
- [36] Robin Rombach and Patrick Esser. Stable diffusion v1-4. <https://huggingface.co/CompVis/stable-diffusion-v1-4>. 3, 7
- [37] Shulan Ruan, Yong Zhang, Kun Zhang, Yanbo Fan, Fan Tang, Qi Liu, and Enhong Chen. Dae-gan: Dynamic aspect-aware gan for text-to-image synthesis. In *ICCV*, 2021. 6
- [38] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022. 3, 7, 9
- [39] Axel Sauer, Kashyap Chitta, Jens Müller, and Andreas Geiger. Projected gans converge faster. *NeurIPS*, 2021. 8
- [40] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, 2018. 3, 5
- [41] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *ICML*, 2015. 3
- [42] Ming Tao, Hao Tang, Fei Wu, Xiao-Yuan Jing, Bing-Kun Bao, and Xu Changsheng. Df-gan: A simple and effective baseline for text-to-image synthesis. In *CVPR*, 2022. 3, 4, 5, 6, 8
- [43] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *NeurIPS*, 2017. 3
- [44] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The Caltech-UCSD Birds-200-2011 Dataset. Technical Report CNS-TR-2011-001, California Institute of Technology, 2011. 5
- [45] Zhihao Wang, Jian Chen, and Steven CH Hoi. Deep learning for image super-resolution: A survey. *IEEE TPAMI*, 43(10):3365–3387, 2020. 1
- [46] Zihao Wang, Wei Liu, Qian He, Xinglong Wu, and Zili Yi. Clip-gen: Language-free training of a text-to-image generator with clip. *arXiv preprint arXiv:2203.00386*, 2022. 1
- [47] Chenfei Wu, Jian Liang, Lei Ji, Fan Yang, Yuejian Fang, Daxin Jiang, and Nan Duan. Nüwa: Visual synthesis pre-training for neural visual world creation. In *ECCV*, 2022. 6

- [48] Tao Xu, Pengchuan Zhang, Qiuyuan Huang, Han Zhang, Zhe Gan, Xiaolei Huang, and Xiaodong He. Attngan: Fine-grained text to image generation with attentional generative adversarial networks. In *CVPR*, 2018. 1, 3, 6
- [49] Zipeng Xu, Tianwei Lin, Hao Tang, Fu Li, Dongliang He, Nicu Sebe, Radu Timofte, Luc Van Gool, and Errui Ding. Predict, prevent, and evaluate: Disentangled text-driven image manipulation empowered by pre-trained vision-language model. In *CVPR*, 2022. 1
- [50] Guojun Yin, Bin Liu, Lu Sheng, Nenghai Yu, Xiaogang Wang, and Jing Shao. Semantics disentangling for text-to-image generation. In *CVPR*, 2019. 1
- [51] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gungjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2022. 3, 7, 9
- [52] Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. Self-attention generative adversarial networks. In *ICML*, 2019. 5
- [53] Han Zhang, Jing Yu Koh, Jason Baldridge, Honglak Lee, and Yinfei Yang. Cross-modal contrastive learning for text-to-image generation. In *CVPR*, 2021. 6
- [54] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris N Metaxas. Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. In *ICCV*, 2017. 3, 6
- [55] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris N Metaxas. Stackgan++: Realistic image synthesis with stacked generative adversarial networks. *IEEE TPAMI*, 41(8):1947–1962, 2018. 3, 6
- [56] Yufan Zhou, Ruiyi Zhang, Changyou Chen, Chunyuan Li, Chris Tensmeyer, Tong Yu, Jiuxiang Gu, Jinhui Xu, and Tong Sun. Towards language-free training for text-to-image generation. In *CVPR*, 2022. 3, 6, 7
- [57] Minfeng Zhu, Pingbo Pan, Wei Chen, and Yi Yang. Dm-gan: Dynamic memory generative adversarial networks for text-to-image synthesis. In *CVPR*, 2019. 1, 3, 6