A Diffusion-based Method for Multi-turn Compositional Image Generation

Chao Wang Toronto AI Lab LG Electronics chao2.wang@lge.com

Abstract

Multi-turn compositional image generation (M-CIG) is a challenging task that aims to iteratively manipulate a reference image given a modification text. While most of the existing methods for M-CIG are based on generative adversarial networks (GANs), recent advances in image generation have demonstrated the superiority of diffusion models over GANs. In this paper, we propose a diffusion-based method for M-CIG named conditional denoising diffusion with image compositional matching (CDD-ICM). We leverage CLIP as the backbone of image and text encoders, and incorporate a gated fusion mechanism, originally proposed for question answering, to compositionally fuse the reference image and the modification text at each turn of M-CIG. We introduce a conditioning scheme to generate the target image based on the fusion results. To prioritize the semantic quality of the generated target image, we learn an auxiliary image compositional match (ICM) objective, along with the conditional denoising diffusion (CDD) objective in a multi-task learning framework. Additionally, we also perform ICM guidance and classifier-free guidance to improve performance. Experimental results show that CDD-ICM achieves state-of-the-art results on two benchmark datasets for M-CIG, i.e., CoDraw and i-CLEVR.

1. Introduction

Image generation is a hot topic in computer vision, which has many applications in a wide range of areas, such as art, education, and entertainment. The generation of an image often needs to follow a text prompt. Additionally, sometimes the generation also needs to be based on an existing image rather than starting from scratch. Combining the above two requirements brings about *compositional image generation* (*CIG*), which is to generate a target image by changing a reference image according to a modification text. Addressing this cross-modal task is useful in computer-aided design (CAD), as it enables a computer system to generate images given verbal instructions from users.



Figure 1. A three-turn example of multi-turn compositional image generation (M-CIG).

In this paper, we focus on *multi-turn compositional image* generation (*M*-CIG), which is to perform CIG in an iterative manner. As shown in Figure 1, M-CIG can be described as a sequence of CIG turns, where the initial reference image is a background canvas, and the target image generated at each turn will be used as the reference image at the next turn. Compared with CIG, M-CIG is more challenging due to the iterative setting. Meanwhile, M-CIG is also more practical than CIG, as in the real world, a user usually needs to go through a series of incremental interactions with a computer system before achieving a final goal.

To the best of our knowledge, the existing methods for M-CIG [12, 13, 34] are mostly based on generative adversarial networks (GANs) [15], which are currently the dominant family of techniques in image generation. According to some theoretical and empirical studies [2, 4, 6, 36, 54], although GANs can generate high-quality images, they are usually difficult to train, and the diversity of the generated images is also limited. Recently, diffusion models [16, 58, 59, 62], which are another family of generative modeling techniques, have gained great popularity in image generation. Compared with GANs, diffusion models are easier to train due to the straightforward definition of objectives, and can also generate more diverse images due to the explicit modeling of data distribution. As for the quality of the generated images, it has been demonstrated that diffusion models are comparable to or even better than GANs [8, 17, 40, 63]. Therefore, we apply diffusion models to M-CIG.

Diffusion models rely on a denoising diffusion mechanism [16] to generate images, which can be conditional so that only desired images are generated. Therefore, the key to addressing M-CIG using diffusion models is to learn conditional denoising diffusion (CDD), where the condition for generating the target image at each turn comes from the reference image and the modification text. However, this raises the following two problems:

The lack of an appropriate conditioning scheme. Although many conditioning schemes have been proposed for diffusion models, most of these works only deal with unimodal cases, where the condition comes from either an image [49, 51, 53] or a text [39, 47, 52]. The conditioning scheme proposed by [22] is aimed at a multi-modal case, where the condition comes from an image-text pair, but this work assumes that the text just describes the semantics of the image rather than the desired change to it. In a word, the above conditioning schemes cannot support the application of diffusion models to M-CIG.

The concern about the semantic quality of the generated target image. For the generated target image, we are not only concerned with its visual quality, but also its semantic quality, which refers to whether it contains the desired objects and whether the contained objects constitute the desired topology. Actually, we believe that the semantic quality deserves more concern than the visual quality. The reason is two-fold. On the one hand, a high semantic quality implies a high visual quality, but the reverse is not true. On the other hand, due to the iterative nature of M-CIG, semantic mistakes are likely to accumulate from turn to turn, which may corrupt the rear turns.

To solve these problems, we propose a diffusion-based method for M-CIG named *conditional denoising diffusion with image compositional matching (CDD-ICM)*, which features a novel conditioning scheme equipped with a multi-task learning framework. Specifically, we use CLIP [45] as the backbone to encode images and texts. On this basis, we borrow a gated fusion mechanism from a question answering (QA) method [67] to perform compositional fusion between the reference image and the modification text at each turn of M-CIG, and use the result as the condition of the denoising diffusion mechanism to generate the target image. To guarantee the semantic quality of the generated target image, we learn image compositional matching (ICM) as an auxiliary objective of CDD to explicitly enhance the condition, where



Figure 2. An overview of CDD-ICM. Colored components are trainable, and those of the same color share their trainable parameters.

the compositional fusion result is aligned with the representation of the target image through contrastive learning. Moreover, we also perform ICM guidance and classifier-free guidance [18] to boost performance. Experimental results show that CDD-ICM achieves state-of-the-art (SOTA) performance on two benchmark datasets for M-CIG, namely CoDraw [24] and i-CLEVR [12].

The contribution of this paper is three-fold. First, we creatively apply diffusion models to M-CIG, where a novel conditioning scheme is developed to handle a compositional image-text pair, integrating the denoising diffusion mechanism with CLIP and a gated fusion mechanism. Second, to prioritize the semantic quality of the generated target image, we establish a multi-task learning framework for the conditioning scheme, where ICM serves as an auxiliary objective of CDD to explicitly enhance the condition. Third, our diffusion-based method outperforms the existing GAN-based methods on two M-CIG benchmark datasets.

2. Method

In this section, we elaborate CDD-ICM, which is our diffusion-based method for M-CIG. We begin by providing a task formulation of M-CIG, which is followed by a detailed introduction to the design of CDD-ICM.

2.1. Task Formulation

Given an initial reference image $a^{(1)}$, which is a background canvas, and a sequence of k modification texts $\{m^{(1)}, \ldots, m^{(k)}\}\)$, which describe the desired changes to be successively made to $a^{(1)}$, M-CIG is a k-turn iterative process, where at each turn $i \in \{1, \ldots, k\}\)$, it is required to generate a target image $z^{(i)}$ by changing the current reference image $a^{(i)}$ according to $m^{(i)}$, and if $i < k, z^{(i)}$ will be used as the next reference image $a^{(i+1)}$.

2.2. Encoding

Considering the cross-modal nature of M-CIG, we map images and texts into a joint representation space through encoding. As shown in Figure 2, we include an image encoder and a text encoder in CDD-ICM, where the former is used to encode reference images and target images, and the latter is used to encode modification texts. Since M-CIG is a vision-and-language (V&L) task, we take advantage of large-scale V&L pre-training by using CLIP, which is pretrained on 400M image-text pairs, as the backbone of both encoders. Specifically, we use the vision part of CLIP, which is a vision transformer (ViT) [11], as the backbone of the image encoder, and use the language part of CLIP, which is a GPT-like [46] auto-regressive language model, as the backbone of the text encoder. In each encoder, we append a linear projection layer after the backbone, which finally yields the representations. For both linear projection layers, we set their output dimensionality to the same value d, which is a hyper-parameter denoting the dimensionality of the joint representation space.

Additionally, as shown in Figure 2, we also include a noisy image encoder in CDD-ICM, which is used to encode the noisy target images obtained in the denoising diffusion mechanism. It has the same structure as the image encoder, but holds different trainable parameters.

2.3. Compositional Fusion

At each turn of M-CIG, to extract clues for the generation of the target image, we perform compositional fusion between the reference image and the modification text. As shown in Figure 2, we include a fusion module in CDD-ICM, which fuses the representation of the reference image with that of the modification text. Actually, we can interpret each turn of M-CIG from the perspective of QA. Specifically, we can regard the reference image as a context, the modification text as a relevant question, and the target image as the corresponding answer. In this way, to implement the fusion module, we borrow the following gated fusion mechanism from a QA method [67]:

$$f(u, v) = g \odot h + (1 - g) \odot u$$

$$g = \text{sigmoid}(W_g[u; v; u \odot v; u - v] + b_g)$$
(1)

$$h = \text{gelu}(W_h[u; v; u \odot v; u - v] + b_h)$$

where W_g and W_h are trainable weight matrices, b_g and b_h are trainable bias vectors, \odot denotes element-wise multiplication, and [;] denotes vector concatenation. In the fusion

module, we set u to the representation of the reference image, set v to that of the modification text, and thereby obtain f(u, v) as the compositional fusion result.

2.4. Conditional Denoising Diffusion

To generate the target image at each turn of M-CIG, we learn conditional denoising diffusion (CDD), which is to perform the generation using a denoising diffusion mechanism conditioned on the compositional fusion result between the reference image and the modification text. As proposed by [16], the denoising diffusion mechanism consists of a forward diffusion process, which gradually injects noises to the target image, and a reverse denoising process, which gradually erases the injected noises.

With the target image denoted as $x_0 \sim q(x_0)$, the forward diffusion process is a pre-defined Markov chain of T time steps $\langle 1, \ldots, T \rangle$, where the state transfers from x_0 all the way to x_T . Specifically, at each time step t, we obtain x_t by injecting a Gaussian noise $\epsilon_t \sim N(0, I)$ to x_{t-1} :

$$q(x_t|x_{t-1}) = N(\sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)I)$$

$$\Leftrightarrow \quad x_t = \sqrt{\alpha_t}x_{t-1} + \sqrt{1 - \alpha_t}\epsilon_t$$
(2)

where α_t is a hyper-parameter used to control the noise scale. It is easy to derive that with $\prod_{i=1}^{t} \alpha_i$ denoted as $\bar{\alpha}_t$, we can actually obtain x_t by directly injecting ϵ_t to x_0 :

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

$$\Leftrightarrow \quad x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_t$$
(3)

As shown in Figure 2, we include a noise injector in CDD-ICM, which executes Equation 3 at an arbitrary time step tto obtain x_t as a noisy target image. To implement the noise injector, we set each $\alpha_t \in \{\alpha_1, \ldots, \alpha_T\}$ in the following way proposed by [40]:

$$\alpha_t = \frac{s(t)}{s(t-1)}$$

$$s(t) = \cos^2\left(\frac{t/T + 0.008}{1.008} \cdot \frac{\pi}{2}\right)$$
(4)

With the compositional fusion result between the reference image and the modification text denoted as c, the reverse denoising process is a parameterized Markov chain of T time steps $\langle T, \ldots, 1 \rangle$, where the state transfers from x_T all the way back to x_0 conditioned on c. Specifically, at each time step t, given x_t and the condition c, we obtain x_{t-1} using a neural network θ :

$$p_{\theta}(x_{t-1}|x_t, c) = N\left(\mu_{\theta}(x_t, t, c), \Sigma_{\theta}(x_t, t, c)\right)$$

$$\Leftrightarrow \quad x_{t-1} = \mu_{\theta}(x_t, t, c) + \Sigma_{\theta}^{\frac{1}{2}}(x_t, t, c)\xi_t$$
(5)

where $\xi_t \sim N(0, I)$. According to [16], we can learn θ by minimizing the following variational lower-bound (VLB)

loss:

$$L^{vlb} = E_{x_0 \sim q(x_0), t \sim U\{1, T\}} [L_t^{vlb}]$$

$$L_t^{vlb} = \begin{cases} -\log p_\theta(x_{t-1} | x_t, c), & \text{if } t = 1\\ D_{KL} (q(x_{t-1} | x_t, x_0) | | p_\theta(x_{t-1} | x_t, c)), & \text{else} \end{cases}$$
(6)

Using Bayes' theorem, it can be derived that with $1 - \alpha_t$ denoted as β_t , $q(x_{t-1}|x_t, x_0)$ is the following Gaussian distribution:

$$q(x_{t-1}|x_t, x_0) = N\left(\tilde{\mu}(x_t, x_0), \beta_t I\right)$$

$$\tilde{\mu}(x_t, x_0) = \frac{\sqrt{\bar{\alpha}_{t-1}\beta_t}}{1 - \bar{\alpha}_t} x_0 + \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} x_t \qquad (7)$$

$$\tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t$$

Based on Equation 3 and Equation 7, we parameterize $\mu_{\theta}(x_t, t, c)$ in the following way proposed by [16]:

$$\mu_{\theta}(x_t, t, c) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(x_t, t, c) \right)$$
(8)

where $\epsilon_{\theta}(x_t, t, c)$ is a prediction to ϵ_t . Besides, we also parameterize $\Sigma_{\theta}(x_t, t, c)$ in the following way proposed by [40]:

$$\Sigma_{\theta}(x_t, t, c) = \exp\left(\log\frac{\beta_t}{\tilde{\beta}_t}\rho_{\theta}(x_t, t, c) + \log\tilde{\beta}_t\right)$$
(9)

where $\rho_{\theta}(x_t, t, c)$ is a fraction used to interpolate between $\log \beta_t$ and $\log \tilde{\beta}_t$. As shown in Figure 2, we include a denoising U-Net in CDD-ICM, which performs the above parameterizations and thus can be seen as θ . To implement the denoising U-Net, we make three changes to the U-Net structure used by [8]. First, we replace the class embedding with the condition c. Second, we concatenate x_t with the reference image along the channel dimension, and thereby use the result as the input. Third, we concatenate the patch representations of the reference image with the token representations of the modification text, and thereby use the result to augment each attention layer as suggested by [39]. The output of the denoising U-Net is divided into two parts along the channel dimension, which are separately used as $\epsilon_{\theta}(x_t, t, c)$ and $\rho_{\theta}(x_t, t, c)$.

According to [16], to learn the above reverse denoising process, instead of minimizing L^{vlb} , we can actually minimize the following mean squared error (MSE) loss:

$$L^{mse} = E_{x_0 \sim q(x_0), t \sim U\{1, T\}} [||\epsilon_t - \epsilon_\theta(x_t, t, c)||^2]$$
(10)

Compared with L^{vlb} , L^{mse} is not only simpler but also more effective. However, minimizing L^{mse} cannot bring any learning signal to $\Sigma_{\theta}(x_t, t, c)$. To benefit from learning $\Sigma_{\theta}(x_t, t, c)$, we combine L^{mse} with L^{vlb} as suggested by [40]. Specifically, we minimize a CDD loss L^{cdd} , which is calculated as follows:

$$L^{cdd} = L^{mse} + \gamma L^{vlb} \tag{11}$$

where γ is a hyper-parameter used to control the weight of L^{vlb} . Additionally, we also stop the gradients of L^{vlb} from flowing to $\epsilon_{\theta}(x_t, t, c)$.

2.5. Image Compositional Matching

At each turn of M-CIG, the semantic quality of the generated target image depends on the condition of the denoising diffusion mechanism, which is the compositional fusion result between the reference image and the modification text. Although learning CDD ensures that the condition is learned, this effect is implicit. To explicitly enhance the condition so that it embodies more clues about the target image, we learn image compositional matching (ICM) as an auxiliary objective of CDD, which is to align the compositional fusion result with the representation of the target image.

To learn ICM, we adopt the InfoNCE loss [41] used in the contrastive pre-training of CLIP, and apply it to individual turns constituting M-CIG samples. Specifically, given a mini-batch of n (reference image, modification text, target image) triples $\{(a_1, m_1, z_1), \dots, (a_n, m_n, z_n)\}$, each of which denotes an individual turn picked from an M-CIG sample, we treat them as positive samples, and generate $n^2 - n$ negative samples by replacing the target image z_i in each positive sample (a_i, m_i, z_i) separately with the other n-1 target images $\{z_1,\ldots,z_n\} - \{z_i\}$. For each of the positive and negative samples, suppose that we have already obtained the compositional fusion result between the reference image and the modification text and the representation of the target image, then we calculate the cosine similarity between them. As a result, we construct a similarity matrix $S \in \mathbb{R}^{n \times n}$, where the element at the *i*-th row and the *j*-th column corresponds to the sample (a_i, m_i, z_i) . It is easy to see that the diagonal elements in S correspond to the positive samples, while the other elements correspond to the negative samples. Based on S, we minimize an ICM loss L^{icm} , which is calculated as follows:

$$L^{icm} = \frac{1}{n} \operatorname{tr} \left(-\log\left(\operatorname{softmax}\left(\frac{S}{\tau}\right)\right) \right) + \frac{1}{n} \operatorname{tr} \left(-\log\left(\operatorname{softmax}\left(\frac{S^{\top}}{\tau}\right)\right) \right)$$
(12)

where τ is a trainable temperature scalar, $tr(\cdot)$ denotes calculating matrix trace, and $softmax(\cdot)$ is calculated along the row dimension. In this way, the compositional fusion result between a (reference image, modification text) pair will be close to the representation of the real target image, while apart from those of the fake ones.

Besides, to enable ICM guidance, which will be introduced later, we also learn noise-aware image compositional matching (N-ICM). Specifically, we replace the above target images with their noisy variants, which are obtained using the noise injector, and encode these noisy target images using the noisy image encoder. On this basis, we minimize an N-ICM loss L^{n-icm} , which is calculated in the same way as we calculate L^{icm} .

2.6. Training and Inference

For the training, we disassemble M-CIG samples into individual turns and thereby apply teacher forcing. On this basis, we divide the training into three stages, where in each stage, we minimize a different loss through mini-batch gradient descent to update the corresponding trainable components of CDD-ICM. Specifically, in the first stage, we minimize L^{icm} to update the image encoder, the text encoder, and the fusion module. In the second stage, we minimize the following joint loss, which is a combination of L^{cdd} and L^{icm} , to update the image encoder, the text encoder, the fusion module, and the denoising U-Net:

$$L^{joint} = L^{cdd} + \delta L^{icm} \tag{13}$$

where δ is a hyper-parameter used to control the weight of L^{icm} . In the third stage, we freeze the image encoder, the text encoder, and the fusion module, and minimize L^{n-icm} to update the noisy image encoder. To effectively fine-tune CLIP, we set the backbone learning rate as a product of the global learning rate and a backbone activity ratio η , which is a hyper-parameter. From the perspective of transfer learning, η controls the trade-off between the knowledge transferred from CLIP and that embodied in the training data.

For the inference, at each turn of M-CIG, we use the image encoder to encode the reference image, use the text encoder to encode the modification text, use the fusion module to perform compositional fusion based on the encoding results, and use the denoising U-Net to iteratively execute Equation 5 from the time step T until the time step 1, where the condition c is set to the compositional fusion result. From Equation 3 and Equation 4, it can be derived that if T is large enough, then $q(x_T) \approx N(0, I)$, thus we sample x_T from N(0, I) at the time step T. To accelerate the inference, we traverse only a part of the time steps, which are uniformly distributed among all of them, and make this process deterministic as suggested by [60]. Finally, we obtain x_0 as the generated target image at the time step 1.

Moreover, we also perform ICM guidance and classifierfree guidance to boost performance. Our ICM guidance is similar to the CLIP guidance of [39]. Specifically, suppose that we have minimized L^{n-icm} in the training, then in the inference, instead of using $\mu_{\theta}(x_t, t, c)$ in Equation 5, we use $\hat{\mu}_{\theta}(x_t, t, c)$, which is obtained by perturbing $\mu_{\theta}(x_t, t, c)$ using the gradient of L^{n-icm} with respect to x_t :

$$\hat{\mu}_{\theta}(x_t, t, c) = \mu_{\theta}(x_t, t, c) + \psi \Sigma_{\theta}(x_t, t, c) \nabla_{x_t} L^{n-\iota cm}$$
(14)

where ψ is a hyper-parameter used to control the perturbation scale. Our classifier-free guidance is similar to that of [39]. Specifically, in the training, when calculating $\epsilon_{\theta}(x_t, t, c)$ in Equation 10, we set the condition c to $\vec{0}$ with a probability of λ , which is a hyper-parameter. On this basis, in the inference, instead of using $\epsilon_{\theta}(x_t, t, c)$ in Equation 8, we use $\hat{\epsilon}_{\theta}(x_t, t, c)$, which is obtained by perturbing $\epsilon_{\theta}(x_t, t, c)$ using $\epsilon_{\theta}(x_t, t, \vec{0})$:

$$\hat{\epsilon}_{\theta}(x_t, t, c) = \phi \epsilon_{\theta}(x_t, t, c) + (1 - \phi) \epsilon_{\theta}(x_t, t, \vec{0})$$
(15)

where ϕ is a hyper-parameter used to control the perturbation scale.

3. Related Works

3.1. Image Manipulation

The goal of image manipulation is to modify specific attributes of an image while avoiding unintended changes or generating a completely new image. Existing works can be split into two main categories: image-to-image translation and text-conditioned image manipulation.

Image-to-Image Translation. The image-to-image translation aims to generate an output image only conditioning on an input image, *i.e.*, uni-modal condition. Image inpainting and image super-resolution are two typical image-to-image translation tasks. In recent years, deep learning has achieved great success in image inpainting. Context Encoders [44] first explores to utilize conditional GANs. Multiple variants [29, 69, 72, 73] of U-Net [50] have been proposed for image inpainting. Some works explore multi-stage generation by taking object edges [38], structures [48], or semantic segmentation maps [64] as intermediate clues. In terms of super-resolution, most early works are regression-based and trained with MSE loss [9, 10, 23, 70]. Auto-regressive models [7, 42] and GAN-based methods [19, 26, 32, 35, 68] have also shown high quality results.

Text-Conditioned Image Manipulation. The textconditioned image manipulation targets generating an output image conditioned on both the input image and text, *i.e.*, multi-modal condition. The input text can be a caption-like description of the target image, and the editing is usually single-turn. TAGAN [37] employs word-level local discriminators to preserve text-irrelevant content. ManiGAN [28] first selects image regions and then correlates the regions with semantic words. DiffusionCLIP [22] is a robust framework that utilizes the pre-trained diffusion models and CLIP loss for image manipulation.

The input text can also be user-provided text instructions that describe desired modifications, such as adding, changing, or removing the objects in images. Generating an image based on provided instructions and an input image is dubbed as the compositional image generation (CIG) task in this paper. [1] achieve great performance on the benchmarks CSS [65] and Fashion Synthesis [75] by designing an improved image & text composition layer and a multi-modal similarity module. [74] propose a GAN-based method to locally modify image features and show remarkable results on both CSS and Abstract Scene [76]. Afterward, the M-CIG task presents a more challenging setting compared to the above single-turn CIG task. [12] first propose the M-CIG task known as Generative Neural Visual Artist (GeNeVA) task, which requires iteratively generating an image according to ongoing linguistic input. [14] introduce the self-supervised counterfactual reasoning (SSCR) framework to tackle the data scarcity problem. LatteGAN [34] improves desired object generation by introducing a Latte module and a textconditioned U-Net discriminator. Our research work targets the M-CIG task, following [12], we conducted experiments on CoDraw [25] and i-CLEVR [12].

3.2. Diffusion Models

Diffusion models (DMs) [58], which formulate the data sampling process as an iterative denoising procedure, are closely related to a large family of methods for learning generative models as transition operators of Markov chains [3, 27, 56, 58, 61]. Many research works concentrate on improving the diffusion process of DMs. [62] propose to estimate the gradients of data distribution via score matching and produce samples via Langevin dynamics. Denoising diffusion probabilistic models (DDPMs) [16], which optimize a variational lower bound to the log-likelihood, can achieve comparable sample quality as GANs [5,20]. Denoising diffusion implicit models (DDIMs) [59] speed up the sampling process while enabling near-perfect inversion [8]. The improved DDPM [40] introduces several modifications to achieve competitive likelihoods without sacrificing sample quality. The latent diffusion [49] model is applied in latent space instead of pixel space to enable an efficient diffusion process. Although GANs have achieved plausible results in image synthesis, they are usually difficult to train and tend to limit the diversity of the generated images [57,66]. DMs are more stable during training and demonstrate comparable or even better performance for image synthesis [8, 17, 63].

Motivated by the progress in developing DMs, some research works explore text or image conditional diffusion mechanisms. While certain diffusion models solely utilize an input image for conditioning, such as PALETTE [51] and SR3 [53], those that condition on both an input image and text are more pertinent to our work. GLIDE [39] is a text-guided diffusion model, where classifier-free guidance yields higher-quality images than CLIP guidance. DALL-E 2 [47] initially generates a CLIP [45] image embedding given a text caption and then generates an image conditioned on the image embedding. ImageGen [52] exhibits a deep level of language understanding which enables high-fidelity image generation. Stable Diffusion [49] is an efficient latent diffusion model and has achieved superior image synthesis performance. [30] compose pre-trained text-guided diffusion models to improve structured generalization for image generation. However, in the context of this paper, it should

be noted that the textual input in the aforementioned models refers to a caption-like description, rather than iterative instructions on image manipulations. Designing diffusionbased methods for the M-CIG task remains a relatively underexplored area, presenting challenges in iteratively modifying images with instructions and conditioning the denoising diffusion mechanism on multi-modalities. To address this gap in the literature, we propose a diffusion-based approach coupled with auxiliary ICM objectives to enhance the visual and semantic fidelity of generated images in the M-CIG task.

4. Experiments

4.1. Datasets

To verify the effectiveness of CDD-ICM, we conduct experiments on the following two M-CIG benchmark datasets:

CoDraw. CoDraw contains 8K M-CIG samples for training, 1K for validation, and 1K for test. The number of turns per M-CIG sample varies between 1 and 14 with an average of 4.25. The reference images and the target images contain 58 classes of clip-art-style objects, such as boys, girls, and trees. The modification texts are conversations between a teller and a drawer.

i-CLEVER. i-CLEVER contains 6K M-CIG samples for training, 2K for validation, and 2K for test. The number of turns per M-CIG sample is always 5. The reference images and the target images contain 24 classes of colored geometric objects, such as red spheres, yellow cubes, and blue cylinders. The modification texts are sentences indicating the addition of a new object.

Both datasets adopt four metrics to evaluate the semantic quality of the generated target images: precision, recall, F1 score, and relational similarity (RSIM). Specifically, each dataset comes with an object detector, which is trained on the dataset by [12]. In the evaluation, the object detector is applied to both the generated target images and the ground-truth target images. On this basis, precision, recall, and F1 score are calculated for each turn of M-CIG by comparing the object presence in the generated target image with that in the ground-truth target image:

$$precision = \frac{|O_{gen} \cap O_{gt}|}{|O_{gen}|}$$
$$recall = \frac{|O_{gen} \cap O_{gt}|}{|O_{gt}|}$$
(16)
$$F1 = 2 \cdot \frac{precision \cdot recall}{precision + recall}$$

where O_{gen} and O_{gt} denote the objects in the generated target image and the ground-truth target image, respectively. RSIM is calculated on the last turn of M-CIG by comparing the object topology in the generated target image with that in the ground-truth target image:

$$RSIM = recall \cdot \frac{|E_{gen} \cap E_{gt}|}{|E_{gt}|}$$
(17)

Method	CoDraw				i-CLEVR			
ivicinou	Precision	Recall	F1	RSIM	Precision	i-CLEVR Recall I 84.72 88 46.39 56 92.96 92 96.93 97 99.94 99	F1	RSIM
GeNeVA-GAN [12]	66.64	52.66	58.83	35.41	92.39	84.72	88.39	74.02
SSCR [13]	58.17	56.61	57.38	39.11	73.75	46.39	56.96	34.54
TIRG [21]	76.56	73.40	72.40	46.64	94.30	92.96	93.71	77.55
LatteGAN [34]	81.50	78.37	77.51	54.16	97.72	96.93	97.26	83.21
CDD-ICM (ours)	90.61	87.55	89.05	57.39	99.99	99.94	99.96	85.66

Table 1. Performance comparison on CoDraw and i-CLEVR.

where E_{gen} and E_{gt} denote the edges interconnecting $O_{gen} \cap O_{gt}$ in the generated target image and the ground-truth target image, respectively.

4.2. Implementation Details

We use PyTorch [43] to implement CDD-ICM, and use HuggingFace's Transformers [71] to load CLIP. In the three CLIP-based encoders, we adopt the basic version of CLIP (i.e. CLIP-ViT-B/32) as the backbone, set the backbone activity ratio η to 0.001, and set the output dimensionality d of all the linear projection layers to 512. In the denoising diffusion mechanism, we use 1000 time steps for the training (*i.e.* T = 1000), and use the 250 time steps uniformly distributed among them for the inference. For ICM guidance, we set ψ to 2. For classifier-free guidance, we set the value of λ to 0.2 and ϕ to 3. To calculate L^{cdd} , we set γ to 1.5. To calculate L^{icm} and L^{n-icm} , we initialize τ to e^{-1} . To calculate L^{joint} in the second training stage, we set δ to 0.1. For the optimization in each training stage, we apply an AdamW optimizer [31] with an initial learning rate of 0.0001 and a weight decay factor of 0.01. We perform the optimization on 8 NVIDIA V100 16GB GPUs in parallel, and set the mini-batch size on each GPU to 32. We calculate the average loss on the validation subset after every 10 epochs. If the resulting number is reduced, then we save the current CDD-ICM model, otherwise we restore the CDD-ICM model to the previous saved version. We decay the learning rate by 50% after each restoration, and terminate the optimization after the 5th restoration.

4.3. Experimental Results

On each dataset, we train a CDD-ICM model by using the training set for optimization and using the validation set for model selection. To compare CDD-ICM with the existing M-CIG methods, we use the test set for evaluation, and finally report the precision, recall, F1 score, and RSIM over all the M-CIG samples in the test set. As shown in Table 1, we achieve SOTA performance on both datasets. Specifically, on CoDraw, CDD-ICM outperforms the existing M-CIG methods by a large margin, where the advantage in precision, recall, and F1 score is relatively larger than that in RSIM. On i-CLEVR, although the existing M-CIG methods did not leave too much room for improvement in precision, recall, and F1 score, CDD-ICM is still better than them in these metrics, reaching almost perfect numbers, and also outperforms them in RSIM.

From all the M-CIG methods in Table 1, we observe two regularities. On the one hand, the performance of these methods on CoDraw is generally worse than that on i-CLEVR. By comparing CoDraw with i-CLEVR, we speculate that this is mainly because the modification texts in CoDraw are commonly longer and more complicated than those in i-CLEVR, which makes it more difficult to generate the desired target images. On the other hand, the performance of these methods in object presence, which is reflected by precision, recall, and F1 score, is generally better than that in object topology, which is reflected by RSIM. By investigating both datasets, we speculate that this is mainly because good performance in object presence just requires correctly identifying the names and attributes of objects from the modification texts, while that in object topology usually requires comprehensively understanding the modification texts.

4.4. Case Study

To visually demonstrate the capability of CDD-ICM, we select several representative M-CIG samples from the test set of both datasets, and use the inference results of CDD-ICM on them as demo cases. A demo case from CoDraw and another one from i-CLEVR are shown in Figure 3, and more demo cases are available in the appendix. In the demo cases from CoDraw, the generated target images contain most of the desired objects, but there are still some extra and missing objects. Besides, it also shows that on CoDraw, CDD-ICM struggles with accurately positioning and orienting objects. In the demo cases from i-CLEVR, the generated target images are very similar to the ground-truth target images, where only a few objects are misaligned.

4.5. Ablation Study

To probe the performance contribution from each design point of CDD-ICM, we conduct the following five ablation experiments. As shown in Table 2, in each ablation experiment, we change a corresponding design point, and report the resulting F1 score and RSIM on each dataset.



Figure 3. Demo cases from CoDraw and i-CLEVR. For the convenience of display, we only include the utterances of the drawer in the modification texts of CoDraw.

Method	Col	Draw	i-CLEVR		
Withou	F1	RSIM	F1	RSIM	
CDD-ICM	89.05	57.39	99.96	85.66	
w/o ICM	75.86	50.58	89.92	74.49	
w/o ICM Guidance	87.69	56.94	96.27	81.32	
w/o Classifier-free Guidance	86.34	56.11	94.91	80.28	
Fine-tuning of CLIP: Frozen	80.63	54.89	87.24	72.44	
Fine-tuning of CLIP: Fully-Trainable	58.52	39.97	65.83	44.76	
w/o Iterative Setting	92.51	75.93	100.00	96.20	

Table 2. Results of ablation experiments.

ICM. We disable ICM by skipping the first training stage and setting δ to 0 when calculating L^{joint} in the second training stage. As a result, we observe a significant drop in F1 score and RSIM. This verifies the effectiveness of learning ICM as an auxiliary objective of CDD. Besides, we also observe a significant drop in the converging speed of the second training stage. This further verifies that learning ICM is beneficial for learning CDD.

ICM Guidance. We disable ICM guidance by skipping the third training stage and setting ψ to 0. As a result, we observe a drop in F1 score and RSIM. This verifies the effectiveness of ICM guidance.

Classifier-free Guidance. We disable classifier-free guidance by setting λ to 0 and setting ϕ to 1. As a result, we

observe a drop in F1 score and RSIM. This verifies the effectiveness of classifier-free guidance.

Fine-tuning of CLIP. For the fine-tuning of CLIP, which is controlled by the backbone activity ratio η , we examine two extreme cases. On the one hand, we make CLIP frozen by setting η to 0. On the other hand, we make CLIP fully-trainable by setting η to 1. As a result, we observe a significant drop in F1 score and RSIM in both cases. This verifies the necessity of applying η .

Iterative Setting. In the evaluation, we disable the iterative setting by using the ground-truth target image at each turn of M-CIG as the reference image of the next turn, which actually downgrades M-CIG to CIG. As a result, we observe a significant rise in F1 score and RSIM. This verifies that M-CIG is more challenging than CIG.

5. Conclusion and Limitation

In this paper, we focus on M-CIG, which is a challenging and practical image generation task, and propose a diffusionbased method named CDD-ICM, which achieves SOTA performance on CoDraw and i-CLEVR. The limitation of CDD-ICM mainly lies in its inference efficiency. Although we have accelerated the inference of CDD-ICM by traversing only a part of the time steps in a deterministic manner, it still takes 3 GPU seconds for CDD-ICM to generate a target image, which is much slower than the GAN-based methods. In the future, we plan to further accelerate the inference of CDD-ICM by applying latent diffusion models [49] and knowledge distillation methods [33, 55].

References

- Kenan E. Ak, Ying Sun, and Joo-Hwee Lim. Learning crossmodal representations for language-based image manipulation. 2020 IEEE International Conference on Image Processing (ICIP), pages 1601–1605, 2020. 5
- [2] Martin Arjovsky and Leon Bottou. Towards principled methods for training generative adversarial networks. In *ICLR*, 2017. 1
- [3] Yoshua Bengio, Eric Laufer, Guillaume Alain, and Jason Yosinski. Deep generative stochastic networks trainable by backprop. In *ICML*, pages 226–234. PMLR, 2014. 6
- [4] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. In *ICLR*, 2019. 1
- [5] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. *ICLR*, 2019. 6
- [6] Andrew Brock, Theodore Lim, James M Ritchie, and Nick Weston. Neural photo editing with introspective adversarial networks. In *ICLR*, 2017. 1
- [7] Ryan Dahl, Mohammad Norouzi, and Jonathon Shlens. Pixel recursive super resolution. In *ICCV*, pages 5439–5448, 2017.
- [8] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. In *NeurIPS*, 2021. 2, 4, 6
- [9] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Learning a deep convolutional network for image superresolution. In *ECCV*, pages 184–199. Springer, 2014. 5
- [10] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Image super-resolution using deep convolutional networks. *IEEE transactions on pattern analysis and machine intelligence*, 38(2):295–307, 2015. 5
- [11] 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. 3
- [12] Alaaeldin El-Nouby, Shikhar Sharma, Hannes Schulz, Devon Hjelm, Layla El Asri, Samira Ebrahimi Kahou, Yoshua Bengio, and Graham W Taylor. Tell, draw, and repeat: Generating and modifying images based on continual linguistic instruction. In *ICCV*, 2019. 1, 2, 6, 7
- [13] Tsu-Jui Fu, Xin Wang, Scott Grafton, Miguel Eckstein, and William Yang Wang. SSCR: Iterative language-based image editing via self-supervised counterfactual reasoning. In *EMNLP*, 2020. 1, 7
- [14] Tsu-Jui Fu, Xin Wang, Scott Grafton, Miguel Eckstein, and William Yang Wang. Sscr: Iterative language-based image editing via self-supervised counterfactual reasoning. In *EMNLP*, pages 4413–4422, 2020. 6
- [15] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NeurIPS*, 2014. 1
- [16] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020. 1, 2, 3, 4, 6

- [17] Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. *Journal of Machine Learning Research*, 2022. 2, 6
- [18] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. arXiv preprint arXiv:2207.12598, 2022. 2
- [19] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. In *ICLR*, 2018. 5
- [20] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *CVPR*, pages 4401–4410, 2019. 6
- [21] E Ak Kenan, Ying Sun, and Joo Hwee Lim. Learning crossmodal representations for language-based image manipulation. In 2020 IEEE International Conference on Image Processing (ICIP), 2020. 7
- [22] Gwanghyun Kim, Taesung Kwon, and Jong Chul Ye. Diffusionclip: Text-guided diffusion models for robust image manipulation. In *CVPR*, 2022. 2, 5
- [23] Jiwon Kim, Jung Kwon Lee, and Kyoung Mu Lee. Accurate image super-resolution using very deep convolutional networks. In *CVPR*, pages 1646–1654, 2016. 5
- [24] Jin-Hwa Kim, Nikita Kitaev, Xinlei Chen, Marcus Rohrbach, Byoung-Tak Zhang, Yuandong Tian, Dhruv Batra, and Devi Parikh. CoDraw: Collaborative drawing as a testbed for grounded goal-driven communication. In ACL, 2019. 2
- [25] Jin-Hwa Kim, Nikita Kitaev, Xinlei Chen, Marcus Rohrbach, Byoung-Tak Zhang, Yuandong Tian, Dhruv Batra, and Devi Parikh. Codraw: Collaborative drawing as a testbed for grounded goal-driven communication. In ACL, pages 6495– 6513, 2019. 6
- [26] Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al. Photorealistic single image super-resolution using a generative adversarial network. In *CVPR*, pages 4681–4690, 2017. 5
- [27] Daniel Levy, Matthew D. Hoffman, and Jascha Sohl-Dickstein. Generalizing hamiltonian monte carlo with neural networks. In *ICLR*, 2018. 6
- [28] Bowen Li, Xiaojuan Qi, Thomas Lukasiewicz, and Philip HS Torr. Manigan: Text-guided image manipulation. In *CVPR*, pages 7880–7889, 2020. 5
- [29] Hongyu Liu, Bin Jiang, Yibing Song, Wei Huang, and Chao Yang. Rethinking image inpainting via a mutual encoderdecoder with feature equalizations. In *ECCV*, pages 725–741. Springer, 2020. 5
- [30] Nan Liu, Shuang Li, Yilun Du, Antonio Torralba, and Joshua B Tenenbaum. Compositional visual generation with composable diffusion models. In *ECCV*, pages 423–439. Springer, 2022. 6
- [31] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2019. 7
- [32] Andreas Lugmayr, Martin Danelljan, Luc Van Gool, and Radu Timofte. Srflow: Learning the super-resolution space with normalizing flow. In *ECCV*, pages 715–732. Springer, 2020.

- [33] Eric Luhman and Troy Luhman. Knowledge distillation in iterative generative models for improved sampling speed. arXiv preprint arXiv:2101.02388, 2021. 8
- [34] Shoya Matsumori, Yuki Abe, Kosuke Shingyouchi, Komei Sugiura, and Michita Imai. Lattegan: Visually guided language attention for multi-turn text-conditioned image manipulation. *IEEE Access*, 2021. 1, 6, 7
- [35] Sachit Menon, Alexandru Damian, Shijia Hu, Nikhil Ravi, and Cynthia Rudin. Pulse: Self-supervised photo upsampling via latent space exploration of generative models. In *CVPR*, pages 2437–2445, 2020. 5
- [36] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. In *ICLR*, 2018. 1
- [37] Seonghyeon Nam, Yunji Kim, and Seon Joo Kim. Textadaptive generative adversarial networks: manipulating images with natural language. *NeurIPS*, 2018. 5
- [38] Kamyar Nazeri, Eric Ng, Tony Joseph, Faisal Z Qureshi, and Mehran Ebrahimi. Edgeconnect: Generative image inpainting with adversarial edge learning. *arXiv preprint arXiv:1901.00212*, 2019. 5
- [39] Alex 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. 2, 4, 5, 6
- [40] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *ICML*, 2021. 2, 3, 4, 6
- [41] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. arXiv preprint arXiv:1807.03748, 2018. 4
- [42] Niki Parmar, Ashish Vaswani, Jakob Uszkoreit, Lukasz Kaiser, Noam Shazeer, Alexander Ku, and Dustin Tran. Image transformer. In *ICML*, pages 4055–4064. PMLR, 2018.
 5
- [43] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *NeurIPS*, 2019. 7
- [44] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. Context encoders: Feature learning by inpainting. In *CVPR*, pages 2536–2544, 2016. 5
- [45] 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, 6
- [46] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. *OpenAI*, 2018. 3
- [47] 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. 2, 6

- [48] Yurui Ren, Xiaoming Yu, Ruonan Zhang, Thomas H Li, Shan Liu, and Ge Li. Structureflow: Image inpainting via structureaware appearance flow. In *ICCV*, pages 181–190, 2019. 5
- [49] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 2, 6, 8
- [50] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. Unet: Convolutional networks for biomedical image segmentation. In Medical Image Computing and Computer-Assisted Intervention-MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18, pages 234–241. Springer, 2015. 5
- [51] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. In ACM SIGGRAPH 2022 Conference Proceedings, 2022. 2, 6
- [52] 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. 2, 6
- [53] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi. Image superresolution via iterative refinement. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. 2, 6
- [54] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. In *NeurIPS*, 2016. 1
- [55] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. In *ICLR*, 2022. 8
- [56] Tim Salimans, Diederik Kingma, and Max Welling. Markov chain monte carlo and variational inference: Bridging the gap. In *ICML*, pages 1218–1226. PMLR, 2015. 6
- [57] Konstantin Shmelkov, Cordelia Schmid, and Karteek Alahari. How good is my gan? In ECCV, pages 213–229, 2018. 6
- [58] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *ICML*, 2015. 1, 6
- [59] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*, 2021. 1, 6
- [60] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*, 2021. 5
- [61] Jiaming Song, Shengjia Zhao, and Stefano Ermon. A-nice-mc: Adversarial training for mcmc. *NeurIPS*, 30, 2017. 6
- [62] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. In *NeurIPS*, 2019. 1, 6
- [63] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *ICLR*, 2021. 2, 6
- [64] Yuhang Song, Chao Yang, Yeji Shen, Peng Wang, Qin Huang, and C.-C. Jay Kuo. Spg-net: Segmentation prediction and guidance network for image inpainting. In *British Machine Vision Conference 2018, BMVC 2018, Newcastle, UK, September 3-6, 2018*, page 97. BMVA Press, 2018. 5

- [65] Nam Vo, Lu Jiang, Chen Sun, Kevin Murphy, Li-Jia Li, Li Fei-Fei, and James Hays. Composing text and image for image retrieval-an empirical odyssey. In *CVPR*, pages 6439–6448, 2019. 5
- [66] Hao Wang, Guosheng Lin, Steven CH Hoi, and Chunyan Miao. Cycle-consistent inverse gan for text-to-image synthesis. In *Proceedings of the 29th ACM International Conference* on Multimedia, pages 630–638, 2021. 6
- [67] Wei Wang, Ming Yan, and Chen Wu. Multi-granularity hierarchical attention fusion networks for reading comprehension and question answering. In *ACL*, 2018. 2, 3
- [68] Xintao Wang, Ke Yu, Shixiang Wu, Jinjin Gu, Yihao Liu, Chao Dong, Yu Qiao, and Chen Change Loy. Esrgan: Enhanced super-resolution generative adversarial networks. In Proceedings of the European conference on computer vision (ECCV) workshops, pages 0–0, 2018. 5
- [69] Yi Wang, Xin Tao, Xiaojuan Qi, Xiaoyong Shen, and Jiaya Jia. Image inpainting via generative multi-column convolutional neural networks. *NeurIPS*, 2018. 5
- [70] Zhaowen Wang, Ding Liu, Jianchao Yang, Wei Han, and Thomas Huang. Deep networks for image super-resolution with sparse prior. In *ICCV*, pages 370–378, 2015. 5
- [71] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Huggingface's transformers: State-of-the-art natural language processing. ArXiv, 2019. 7
- [72] Zhaoyi Yan, Xiaoming Li, Mu Li, Wangmeng Zuo, and Shiguang Shan. Shift-net: Image inpainting via deep feature rearrangement. In *ECCV*, pages 1–17, 2018. 5
- [73] Yanhong Zeng, Jianlong Fu, Hongyang Chao, and Baining Guo. Learning pyramid-context encoder network for highquality image inpainting. In *CVPR*, pages 1486–1494, 2019.
 5
- [74] Tianhao Zhang, Hung-Yu Tseng, Lu Jiang, Weilong Yang, Honglak Lee, and Irfan Essa. Text as neural operator: Image manipulation by text instruction. In *Proceedings of the 29th* ACM International Conference on Multimedia, pages 1893– 1902, 2021. 5
- [75] Shizhan Zhu, Raquel Urtasun, Sanja Fidler, Dahua Lin, and Chen Change Loy. Be your own prada: Fashion synthesis with structural coherence. In *ICCV*, pages 1680–1688, 2017.
 5
- [76] C Lawrence Zitnick and Devi Parikh. Bringing semantics into focus using visual abstraction. In *CVPR*, pages 3009–3016, 2013. 6

A. Demo Cases from CoDraw



Figure 5



Figure 6



Figure 7

B. Demo Cases from i-CLEVR



Figure 9





Figure 11