

DICE: DISCRETE INVERSION ENABLING CONTROLLABLE EDITING FOR MULTINOMIAL DIFFUSION AND MASKED GENERATIVE MODELS

Xiaoxiao He¹, Ligong Han^{1,2†}, Quan Dao¹, Song Wen¹, Minhao Bai¹, Di Liu¹, Han Zhang³, Martin Renqiang Min⁴, Felix Juefei-Xu⁵, Chaowei Tan¹, Bo Liu⁶, Kang Li¹, Hongdong Li⁷, Junzhou Huang⁸, Faez Ahmed⁹, Akash Srivastava² & Dimitris N. Metaxas¹

¹Rutgers University ²MIT-IBM Watson AI Lab ³Google DeepMind ⁴NEC Labs America
⁵NYU ⁶Walmart Global Tech ⁷ANU ⁸UT Arlington ⁹Massachusetts Institute of Technology

† Project Lead & Corresponding Author Project: [\[Website\]](#)

ABSTRACT

Discrete diffusion models have achieved success in tasks like image generation and masked language modeling but face limitations in controlled content editing. We introduce **DICE** (Discrete Inversion for Controllable EditIng), the first approach to enable precise inversion for discrete diffusion models, including multinomial diffusion and masked generative models. By recording noise sequences and masking patterns during the reverse diffusion process, **DICE** enables accurate reconstruction and flexible editing of discrete data without the need for pre-defined masks or attention manipulation. We demonstrate the effectiveness of **DICE** across both image and text domains, evaluating it on models such as VQ-Diffusion, Paella, and RoBERTa. Our results show that **DICE** preserves high data fidelity while enhancing editing capabilities, offering new opportunities for fine-grained content manipulation in discrete spaces. Code is available at [\[link\]](#).

1 INTRODUCTION

Diffusion models have emerged as a powerful class of generative models, achieving remarkable success in high-fidelity image and video synthesis (Ho et al., 2020; Dhariwal & Nichol, 2021; Rombach et al., 2022; Ramesh et al., 2022; Ho et al., 2022; OpenAI, 2024). These models generate data by iteratively denoising samples from a simple noise distribution, effectively reversing a diffusion process that gradually corrupts data. Broadly, diffusion models can be categorized into continuous and discrete types, each tailored to different data modalities and applications.

Continuous diffusion models operate in continuous spaces, leveraging stochastic differential equations (SDEs) or their deterministic counterparts, ordinary differential equations (ODEs), to model the forward and reverse diffusion processes (Song et al., 2020; 2021). Advances such as flow matching (Lipman et al., 2022; Liu et al., 2022) have enhanced their efficiency and flexibility. These models have been successfully applied in various domains, including image editing (Meng et al., 2021; Avrahami et al., 2022; Mokady et al., 2022; Han et al., 2023; 2024; Zhang et al., 2023), medical imaging (He et al., 2023), and solving inverse problems (Chung et al., 2022; Stathopoulos et al., 2024). In image editing, continuous diffusion models enable controlled manipulation of images while preserving consistency with the underlying data distribution. A key capability enabling this is *inversion*—the process of reversing the diffusion model to recover the original noise vector or latent representation that could have generated a given data sample. Two main inversion approaches exist: deterministic inversion using ODEs (e.g., DDIM Inversion (Song et al., 2021)) and stochastic inversion by recording noise sequences (e.g., CycleDiffusion (Wu & De la Torre, 2022), DDPM Inversion (Dhariwal & Nichol, 2021)).

Discrete diffusion models are designed for inherently discrete data such as text or image tokens (Esser et al., 2021b). They adapt the diffusion framework to discrete spaces by defining appropriate transition kernels that corrupt and restore discrete data (Hoogeboom et al., 2021; Austin et al., 2021; Gu et al., 2022). Prominent examples include multinomial diffu-

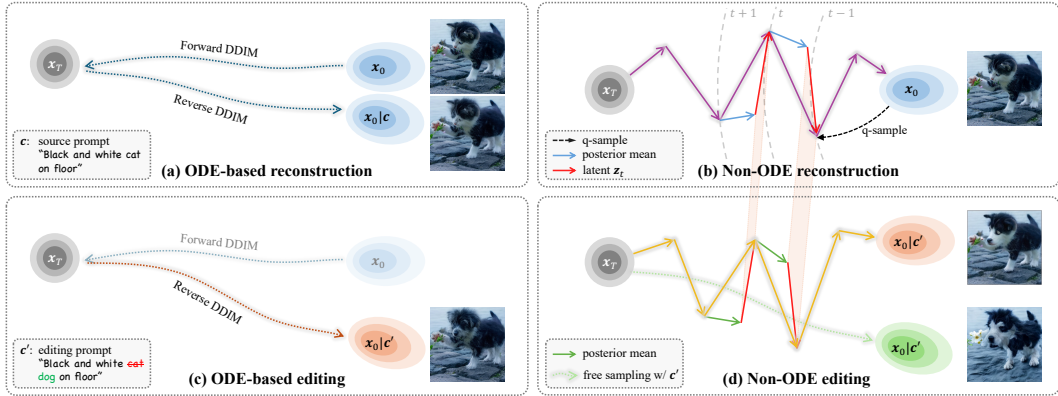


Figure 2: Here we demonstrate the two types of reconstruction and editing paradigms, namely ODE-based and Non-ODE based. (a,c) shows the ODE-based editing and reconstructions, while it provides accurate editing and reconstruction performances, it highly depends on the underlying ODE trajectory, which is not feasible in the discrete diffusion. However, the Non-ODE editing samples a trajectory by directly adding noise to x_0 and record the difference between the predicted x_{t-1} and the sampled x_{t-1} as indicated in the red arrow. In this way, we are able to reconstruct/edit the image without the strong condition of having an underlying ODE.

sion (Hoogeboom et al., 2021; Gu et al., 2022), D3PM (Austin et al., 2021), and masked generative models like MaskGIT (Chang et al., 2022), Muse (Chang et al., 2023). Despite their success in generation tasks, discrete diffusion models face limitations in controlled content editing. For instance, masked generative models achieve image editing through masked inpainting, where regions are masked and regenerated based on new conditions. However, this approach lacks the ability to inject information from the masked area into the inpainting process, limiting fine-grained control over the editing outcome (as illustrated in Figure 1).

Moreover, existing ODE-based inversion techniques developed for continuous diffusion models are not directly applicable to discrete diffusion models due to inherent differences in data representation and diffusion processes. This gap hinders the ability to perform precise inversion and controlled editing in discrete spaces. To address this challenge, we propose **DICE** (Discrete Inversion for Controllable Editing), the first inversion algorithm for discrete diffusion models to the best of our knowledge. Our method extends the stochastic inversion approach to discrete diffusion models, including both multinomial diffusion and masked generative models. The core idea is to record the noise sequence needed to recover a stochastic trajectory in the reverse diffusion process. Specifically, given an artificial trajectory where latent states have low correlation, we fit reverse sampling steps to this trajectory and save the residuals between targets and predictions. This process *imprints* the information of the original input data into the recorded residuals. During editing or inference, the recorded residuals are added back, allowing us to inject and control the amount of information introduced into the inference process.

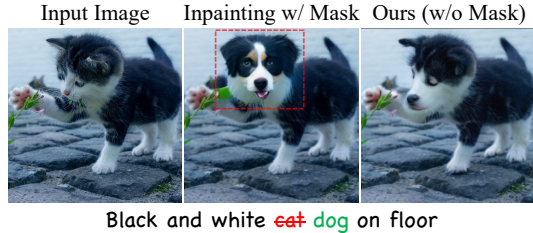


Figure 1: Illustration of the limitation of masked inpainting method. Here, we want to change the cat to a dog. Inpainting with masked generation inadvertently modifies the orientation of the head, resulting in a less favourable result. With our discrete inversion method, we are able to edit the image while preserving other properties of the object being edited. This is achieved by injecting the information from the input image into the logit space. Dotted red box indicates the mask.

Our approach enables accurate reconstruction of the original input data and facilitates controlled editing without the need for predefined masks or attention map manipulation. It provides a flexible framework for fine-grained content manipulation in discrete spaces, overcoming the limitations of

existing methods. We validate the effectiveness of DICE through extensive experiments on both image and text modalities. We evaluate our method on models such as VQ-Diffusion (Gu et al., 2022), Paella (Rampas et al., 2022), and RoBERTa (Liu et al., 2019), demonstrating its versatility across different types of discrete generative models. Additionally, we introduce a novel text-editing dataset to further showcase our method’s capabilities and to facilitate future research in this area. Our contributions can be summarized as follows:

- We introduce DICE, an inversion algorithm for discrete diffusion models, including multinomial diffusion and masked generative models. By recording and injecting noise sequences or masking patterns, DICE enables accurate reconstruction and controlled editing of discrete data without the need for predefined masks or attention manipulation.
- We validate the effectiveness of DICE through comprehensive experiments on both image and text modalities, demonstrating its versatility across different types of discrete generative models.
- We show that our approach can transform a model primarily trained for understanding tasks, such as RoBERTa, into a competitive generative model for text generation and editing, illustrating the potential for extending discrete diffusion models to new applications.

2 RELATED WORK

Discrete diffusion. D3PM (Austin et al., 2021) and Multinomial Diffusion (Hoogeboom et al., 2021) spearheaded the study of diffusion processes in discrete spaces by developing a corruption mechanism for categorical data. Following those works, Esser et al. (2021a) and Gu et al. (2022) introduced the VQ-GAN as a way to discretize the image into tokens. Also, extending to the natural language processing, Devlin et al. (2018) and Liu et al. (2019) proposed a bidirectional transformer for language understanding, which can be viewed as a discrete diffusion model (Wang & Cho, 2019). Additionally, Campbell et al. (2022) proposed discrete diffusion models with continuous time, while Lou et al. (2023) extended score matching (Song & Ermon, 2019) to discrete spaces by learning probability ratios. Gat et al. (2024) proposed discrete flow matching to extend the flow matching to discrete space. MaskGIT (Chang et al., 2022), Muse (Chang et al., 2023) and MMVID (Han et al., 2022) introduced efficient non-autoregressive methods for image generation by iteratively remasking and reprediction.

Diffusion inversion. Diffusion inversion aims to find an encoding or latent representation of the input signal that can be used to reconstruct the original data. Traditional approaches to diffusion inversion are based on neural ODEs (Chen et al., 2018), such as DDIM inversion (Song et al., 2021) and flow matching (Lipman et al., 2022; Liu et al., 2022), where deterministic trajectories are used for inversion. Another class of methods focuses on stochastic differential equations (SDEs) (Song et al., 2020), including models like CycleDiffusion (Wu & De la Torre, 2022) and DDPM Inversion (Huberman-Spiegelglas et al., 2024), which rely on tracking noise or residuals along a stochastic path to recover the input. Our approach generalizes the concept of DDPM Inversion by extending it to discrete diffusion models, enabling effective inversion in both continuous and discrete settings.

Inversion-based image editing. DDIM inversion (Song et al., 2021) has served as a foundational technique for various diffusion-based image editing approaches. In many image editing tasks, DDIM-type methods are often employed alongside guidance techniques like Prompt-to-Prompt (Hertz et al., 2022), which manipulate cross-attention maps. In contrast, DDPM inversion-based (Huberman-Spiegelglas et al., 2024) approaches are more user-friendly, as they do not require

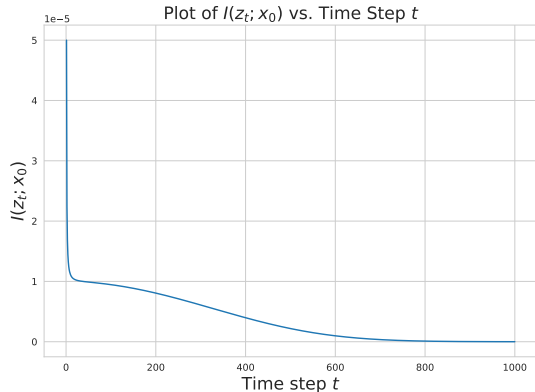


Figure 3: Mutual information between z_t and x_0 . Computed with a simple DDPM setting by assuming $x_0 \sim \mathcal{N}(\mathbf{0}, I)$.

cross-attention manipulations. To address issues such as inaccurate reconstruction and error accumulation, Null-text Inversion (Mokady et al., 2022) introduces test-time optimization of null embeddings, ensuring the reconstruction trajectory aligns more closely with the DDIM inversion path. Negative-prompt Inversion (Miyake et al., 2023; Han et al., 2024) further improves time efficiency by providing a closed-form solution to an approximate inversion problem, reducing computational costs while maintaining competitive reconstruction quality.

3 METHODS

3.1 PRELIMINARIES

Masked generative modeling. Masked generative modeling is widely used in representation learning for both natural language processing and computer vision. It works by masking parts of the input and training the model to reconstruct the missing data. In models like BERT (Devlin et al., 2018) and RoBERTa (Liu et al., 2019), masked tokens ([MASK]) are predicted based on the surrounding context, excelling in text completion and embedding representation learning. For image generation, Paella (Rampas et al., 2022) adapts this approach for text-conditional image generation by renoising tokens instead of masking. The inference process in masked generative models typically involves iterative renoise/remask and repredict steps.

Multinomial Diffusion. Denoting $\mathbf{x}_0 \in \{1, \dots, K\}^D$ as a data point of dimension D . We use $\mathbf{v}(x_t^{(i)})$ to denote the one-hot column vector representation of the i -th entry of \mathbf{x}_t . To simplify notation, in the following we drop index i and any function that operates on vector \mathbf{x}_t is populated along its dimension. Diffusion model defines a Markov chain $q(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})$ that gradually adds noise to the data \mathbf{x}_0 for T times so that \mathbf{x}_T contains little to no information. Discrete diffusion model (Hoogeboom et al., 2021; Austin et al., 2021; Gu et al., 2022) proposed an alternative likelihood-based model for categorical data, and defines the forward process following:

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \text{Cat}(\mathbf{v}(\mathbf{x}_t); \boldsymbol{\pi} = \mathbf{Q}_t \mathbf{v}(\mathbf{x}_{t-1})). \quad (1)$$

where \mathbf{Q}_t is the transition matrix between adjacent states following mask-and-replace strategy. The posterior distribution given \mathbf{x}_0 has a closed-form solution,

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \frac{(\mathbf{Q}_t^\top \mathbf{v}(\mathbf{x}_t)) \odot (\overline{\mathbf{Q}}_{t-1} \mathbf{v}(\mathbf{x}_0))}{\mathbf{v}(\mathbf{x}_t)^\top \mathbf{Q}_t \mathbf{v}(\mathbf{x}_0)}. \quad (2)$$

where $\overline{\mathbf{Q}}_t = \mathbf{Q}_t \cdots \mathbf{Q}_1$ is the cumulative transition matrix. The details of \mathbf{Q}_t and $\overline{\mathbf{Q}}_t$ are given in the supplementary materials. The inference process is as below:

$$\boldsymbol{\pi}_\theta(\mathbf{x}_t, t) = p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \sum_{\tilde{\mathbf{x}}_0=1}^K q(\mathbf{x}_{t-1}|\mathbf{x}_t, \tilde{\mathbf{x}}_0) p_\theta(\tilde{\mathbf{x}}_0|\mathbf{x}_t), \quad (3)$$

with $p_\theta(\tilde{\mathbf{x}}_0|\mathbf{x}_t)$ is parameterized by a neural network. We gradually denoise from \mathbf{x}_T to \mathbf{x}_0 using 3. For numerical stability, the implementation uses log space instead of probability space. Masked generative models can be viewed as a special case of multinomial diffusion models with an additional *absorbing* state (or the [MASK] state). Its training objective can be viewed as a reweighted ELBO (Bond-Taylor et al., 2022).

3.2 DISCRETE INVERSION FOR CONTROLLABLE EDITING

Non ODE-based inversion. ODE-based generative models, such as DDIM and flow matching, define an ODE trajectory. Due to the deterministic nature of ODEs, inversion can be achieved by solving the ODE using the Euler method in forward direction, ensuring reconstruction based on the inherent properties of the ODE. In contrast, another line of research focuses on SDE-based models, such as CycleDiffusion (Wu & De la Torre, 2022) and DDPM Inversion (Huberman-Spiegelglas et al., 2024). Broadly speaking, these approaches ensure reconstruction by recording the noises or residuals that are required to reproduce the stochastic trajectory. CycleDiffusion records the Gaussian noise \mathbf{z}_t during sampling from posterior $p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0 = \mathbf{x}_0)$ and injects information of the input signal by feeding the true \mathbf{x}_0 . DDPM Inversion, on the other hand, incorporates information

into \mathbf{z}_t by fitting the reverse process into an artificial stochastic trajectory obtained by independent `q-sample`. For both CycleDiffusion and DDPM Inversion, the key idea is to utilize the Gaussian reparameterization trick, $x = \mu + \sigma z \Leftrightarrow x \sim \mathcal{N}(x; \mu, \sigma^2)$, and keeping track of the “noise” that could have generated the sample from mean. For discrete diffusion models, we utilize the Gumbel-Max trick (Maddison et al., 2014; Jang et al., 2016), $x = \arg \max (\log(\boldsymbol{\pi}) + \mathbf{g}) \Leftrightarrow x \sim \text{Cat}(x; \boldsymbol{\pi})$. Figure 2 provides an intuition of the proposed method.

Inverting multinomial diffusion. Similar to Huberman-Spiegelglas et al. (2024), we start by sampling a stochastic trajectory, $\{\mathbf{x}_t\}$, a sequence of independent `q-sample`’s from $q(\mathbf{x}_t|\mathbf{x}_0)$ (we populate the following sampling operation along the dimension of \mathbf{x}_t),

$$x_t = \arg \max (\log(q(x_t|x_0)) + \mathbf{g}), \quad \text{with} \quad (4)$$

$$q(x_t|x_0) = \text{Cat}(x_t; \boldsymbol{\pi} = \overline{\mathbf{Q}}_t \mathbf{v}(x_0)) \quad \text{and} \quad \mathbf{g} \sim \text{Gumbel}(\mathbf{0}, \mathbf{I}).$$

Note that here we use the Gumbel softmax trick (Jang et al., 2016), which is equivalent to sampling from categorical distribution $q(x_t|x_0)$.

$$\begin{aligned} \mathbf{y}_{t-1} &= \log(\text{onehot}(\mathbf{x}_{t-1})), \quad \text{and} \\ \hat{\mathbf{y}}_{t-1} &= \log(\boldsymbol{\pi}_\theta(\mathbf{x}_t, t)), \\ \mathbf{z}_t &:= \mathbf{y}_{t-1} - \hat{\mathbf{y}}_{t-1} \end{aligned} \quad (5)$$

Note that here the latent $\mathbf{z}_t \in \mathbb{R}^{D \times K}$.

In this reverse process, the latent space $\{\mathbf{x}_T, \mathbf{z}_T, \mathbf{z}_{t-1}, \dots, \mathbf{z}_1\}$ together with the fixed discrete diffusion model $\boldsymbol{\pi}_\theta$ also uniquely define the same stochastic trajectory $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_T$. The detailed algorithm is given in Algorithm 2.

Algorithm 1 Discrete Inversion for Masked Generative Modeling

Inversion:

- 1: $\mathbf{y}_0 \leftarrow \mathcal{D}(\mathbf{x}_0, \mathbf{c}, t=0)$
- 2: Sample noise token map \mathbf{n}
- 3: **for** t from 1 to T **do**
- 4: $\mathbf{m}_t \leftarrow \text{GenerateMask}(t)$ ▷ Sampling masks according to inference algorithm
- 5: $\mathbf{x}_t \leftarrow \mathbf{x}_0 \odot (\mathbf{1} - \mathbf{m}_t) + \mathbf{n} \odot \mathbf{m}_t$
- 6: $\hat{\mathbf{y}}_{0|t} \leftarrow \mathcal{D}_\theta(\mathbf{x}_t, \mathbf{c}, t=t)$
- 7: $\mathbf{z}_t \leftarrow \mathbf{y}_0 - \hat{\mathbf{y}}_{0|t}$ ▷ Eq 6
- 8: **end for**

Editing/Sampling:

- 9: **for** t from τ to 1 **do**
 - 10: $\hat{\mathbf{y}}_{0|t} \leftarrow \mathcal{D}_\theta(\mathbf{x}_t, \mathbf{c}', t=t)$
 - 11: $\mathbf{g} \sim \text{Gumbel}(\mathbf{0}, \mathbf{I})$
 - 12: $\tilde{\mathbf{y}}_0 \leftarrow \hat{\mathbf{y}}_{0|t} + \lambda_1 \cdot \mathbf{z}_t + \lambda_2 \cdot \mathbf{g}$
 - 13: $\tilde{\mathbf{x}}_0 \leftarrow \arg \max \tilde{\mathbf{y}}_0$
 - 14: $\mathbf{x}_{t-1} \leftarrow \tilde{\mathbf{x}}_0 \odot (\mathbf{1} - \mathbf{m}_{t-1}) + \mathbf{n} \odot \mathbf{m}_{t-1}$
 - 15: **end for**
 - 16: Return \mathbf{x}_0 .
-

Algorithm 2 Discrete Inversion for Multinomial Diffusion

Inversion:

- 1: **for** t from 1 to T **do**
- 2: $\mathbf{x}_t \sim q(\mathbf{x}_t|\mathbf{x}_0)$ ▷ Independent q-sample using Eq 4
- 3: $\mathbf{y}_t \leftarrow \log(\text{onehot}(\mathbf{x}_t))$
- 4: **end for**
- 5: **for** t from T to 1 **do**
- 6: $\hat{\mathbf{y}}_{t-1} \leftarrow \log(\boldsymbol{\pi}_\theta(\mathbf{x}_t, \mathbf{c}, t))$ ▷ Log posterior using Eq 3
- 7: $\mathbf{z}_t \leftarrow \mathbf{y}_{t-1} - \hat{\mathbf{y}}_{t-1}$ ▷ Eq 5
- 8: **end for**

Editing/Sampling:

- 9: **for** t from τ to 1 **do**
 - 10: $\hat{\mathbf{x}}_0 \leftarrow p_\theta(\mathbf{x}_0|\mathbf{x}_t = \arg \max \mathbf{y}_t)$
 - 11: $\mathbf{g} \sim \text{Gumbel}(\mathbf{0}, \mathbf{I})$
 - 12: $\mathbf{y}_{t-1} \leftarrow \log(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \hat{\mathbf{x}}_0; \mathbf{c}')) + \lambda_1 \cdot \mathbf{z}_t + \lambda_2 \cdot \mathbf{g}$
 - 13: **end for**
 - 14: Return $\mathbf{x}_0 = \arg \max \mathbf{y}_0$.
-

Inverting masked generative models. In masked generative modeling, the stochastic trajectory \mathbf{x}_t is constructed according to the specific inference algorithm of the model in use. For example, in Paella Rampas et al. (2022), the masking is *inclusive*, meaning that as the time step t increases, the set of masked tokens grows. In contrast, the Unleashing Transformer Bond-Taylor et al. (2022) employs *random* masking at each step, where masks are generated independently using the `q-sample` function. Without loss of generality, we define a denoiser function \mathcal{D}_θ (parameterized by θ). This denoiser outputs the *logits* of the predicted unmasked data given the noisy tokens \mathbf{x}_t . Since in this case, the categorical sampling happens at sampling from the denoiser’s prediction, we therefore define an corresponding latent sequence:

$$\begin{aligned} \hat{\mathbf{y}}_{0|t} &= \log(p_\theta(\mathbf{x}_0|\mathbf{x}_t)) = \mathcal{D}_\theta(\mathbf{x}_t, t) \\ \mathbf{z}_t &:= \mathbf{y}_0 - \hat{\mathbf{y}}_{0|t}. \end{aligned} \quad (6)$$

With our proposed latent space, accurate reconstruction is guaranteed. However, for editing tasks, this level of precision may not be ideal if the latent variable z_t dominates the generation process. The detailed algorithm is given in Algorithm 1.

To provide more flexibility, we introduce the hyperparameters τ , λ_1 , and λ_2 , which allow for finer control over the editing process. Specifically, τ represents the starting (and largest) timestep at which the editing process begins, while λ_1 controls the amount of information injected from the original input, and λ_2 governs the introduction of random noise.

Analysis. We describe a simple yet prototypical example of DDPM and compute the mutual information between encoded latents and the input signal.

Remark 3.1. *Given a simple Gaussian DDPM with $x_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, latents $\{z_t\}$ are obtained with DDPM inversion (Huberman-Spiegelglas et al., 2024), then the mutual information between z_t and x_0 is:*

$$I(z_t; x_0) = \frac{D}{2} \log\left(\frac{\beta_t^2 \bar{\alpha}_{t-1} + 1 - \bar{\alpha}_{t-1} + \alpha_t(1 - \bar{\alpha}_t)}{1 - \bar{\alpha}_{t-1} + \alpha_t(1 - \bar{\alpha}_t)}\right). \quad (7)$$

The mutual information between z_t and x_0 is illustrated in Figure 3. We observe that the amount of information encoded from x_0 into z_t decreases as t increases, motivating us to explore different scheduling strategies for λ 's (see Figure 7).

4 EXPERIMENTS

In this section, we demonstrate the effectiveness of our proposed inversion methods on both image and language diffusion models. Our experiments show that the methods can preserve identity in both vision and language tasks while successfully making the intended changes. The implementation details can be reviewed in Supplementary Materials.

4.1 IMAGE DIFFUSION MODEL

For the image diffusion model, we mainly investigate the use of absorbing state discrete model (Austin et al., 2021) including a masked generative model, Paella, and a multinomial diffusion model, VQ-Diffusion. We demonstrate the inversion reconstruction ability and image editing performance in both categories with DICE.

Dataset. The Prompt-based Image Editing Benchmark (PIE-Bench) by (Ju et al., 2023) is a recently introduced dataset designed to evaluate text-to-image (T2I) editing methods. The dataset assesses language-guided image editing in 9 different scenarios with 700 images. The benchmark’s detailed annotations and variety of editing tasks were instrumental in thoroughly assessing our method’s capabilities, ensuring a fair and consistent comparison with existing approaches.

4.1.1 INVERSION RECONSTRUCTION

In this section, we evaluate the accuracy of inversion without editing. This is achieved by first inverting the image and then using the recorded latent code to reconstruct the original image.

Evaluation Metrics. Here, we evaluate the image similarity by PSNR, LPIPS, MSE and SSIM of the original and the generated image under the same prompt with DICE and masked generation.

Quantitative Analysis. The reconstruction performance of our method, as shown in Table 1, far surpasses the baseline Inpainting + Paella model across all metrics. In the case of masked inpainting, all image tokens are replaced with randomly sampled tokens, meaning the model lacks any prior information about the original image. As a result, the reconstructed image differs significantly from the one being inverted, leading to lower similarity scores. In contrast, our method demonstrates near-perfect reconstruction, as indicated by the metrics, and notably produces an identical image without the errors typically introduced by the VQ-VAE/GAN quantization process, as seen in the results marked with ([†]). This highlights the superior accuracy and consistency of our approach in generating high-fidelity reconstructions.

Method	Metric			
	PSNR \uparrow	LPIPS $_{\times 10^3}$ \downarrow	MSE $_{\times 10^4}$ \downarrow	SSIM $_{\times 10^2}$ \uparrow
Inpainting+Paella	10.50	565.11	1002.09	30.13
Ours+Paella	30.91	39.81	11.07	90.22
Ours[†]+Paella	Inf	0.07	0.01	99.99

Table 1: **Inversion Reconstruction performance** [†] The metric is calculated between the original image and its inverted counterpart. Due to the encoding and decoding steps in the VQ-VAE/GAN process, some inaccuracies are introduced by the quantization. The PSNR is `Inf` due to the reconstruction of our method yielding the same image after the VQ-VAE/GAN process.

4.1.2 EDITING PERFORMANCE

In this section, we discuss the editing performance of our proposed method. Since there is no discrete diffusion inversion exists, we compare our method with masked generation as indicated in the original paper. In addition to that, we also demonstrate the metric from continuous counterparts.

Evaluation Metrics. To demonstrate the effectiveness and efficiency of our proposed inversion method, we employ eight metrics covering three key aspects: structure distance, background preservation, and edit prompt-image consistency, as outlined in Ju et al. (2023). We utilize the structure distance metric proposed by Tumanyan et al. (2023) to measure the structural similarity between the original and generated images. To evaluate how well the background is preserved outside the annotated editing mask, we use Peak Signal-to-Noise Ratio (PSNR), Learned Perceptual Image Patch Similarity (LPIPS) (Zhang et al., 2018), Mean Squared Error (MSE), and Structural Similarity Index Measure (SSIM) (Wang et al., 2004). We also assess the consistency between the edit prompt and the generated image using CLIP (Radford et al., 2021) Similarity Score (Wu et al., 2021), which is calculated over the whole image and specifically within the regions defined by the editing mask.

Results. In Table 2, we demonstrate the quantitative result of `DICE` using Paella and VQ-Diffusion compared to continuous diffusion model and also inpainting. Notably, our approach with the Paella model achieves the lowest structure distance 11.34, outperforming all other methods, including the continuous diffusion models. Additionally, while the DDPM Inversion with Stable Diffusion v1.4 shows the highest CLIP similarity scores for both whole and edited regions, our method maintains competitive CLIP similarity with Paella. Given the significant reduction in structure distance, our method offers a superior balance between structural preservation and semantic alignment in edits. Furthermore, when combined with VQ-Diffusion, our method continues to show strong performance. The results in Table 3 clearly demonstrate the superior background preservation capabilities of our method compared to DDIM+SD1.4. All four metrics underscore the structural consistency of our approach in preserving the unedited regions of the image. These results show the effectiveness of our method in maintaining background integrity during editing and provide evidence that information about the original image is instilled into the latent space of `DICE`.

In Figure 4, we show the editing results for both Paella and VQ-Diffusion using `DICE`. Both models successfully modify real images according to the target prompts. In all cases, our results exhibit both high fidelity to the input image and adherence to the target prompt.

4.2 LANGUAGE DIFFUSION MODEL

In this section, we evaluate `DICE` on RoBERTa (Liu et al., 2019), a text discrete diffusion model, to generate sentences with opposing sentiments while preserving structural similarities. We begin with two prompts—one with a positive sentiment and another with a negative sentiment. Each prompt contains two sentences: the first sentence indicates the sentiment type and sets the contextual background, and the second sentence is the target for inversion and generation. Initially, we invert the second sentence of the negative sentiment prompt using the entire prompt as context, which produces a noised token representation of that sentence. Next, we condition the model on the positive sentiment by concatenating the first sentence of the positive sentiment prompt with the noised token of the inverted negative sentence. This setup guides the model to generate a new second sentence that mirrors the structure of the original negative sentence but expresses a positive sentiment instead.

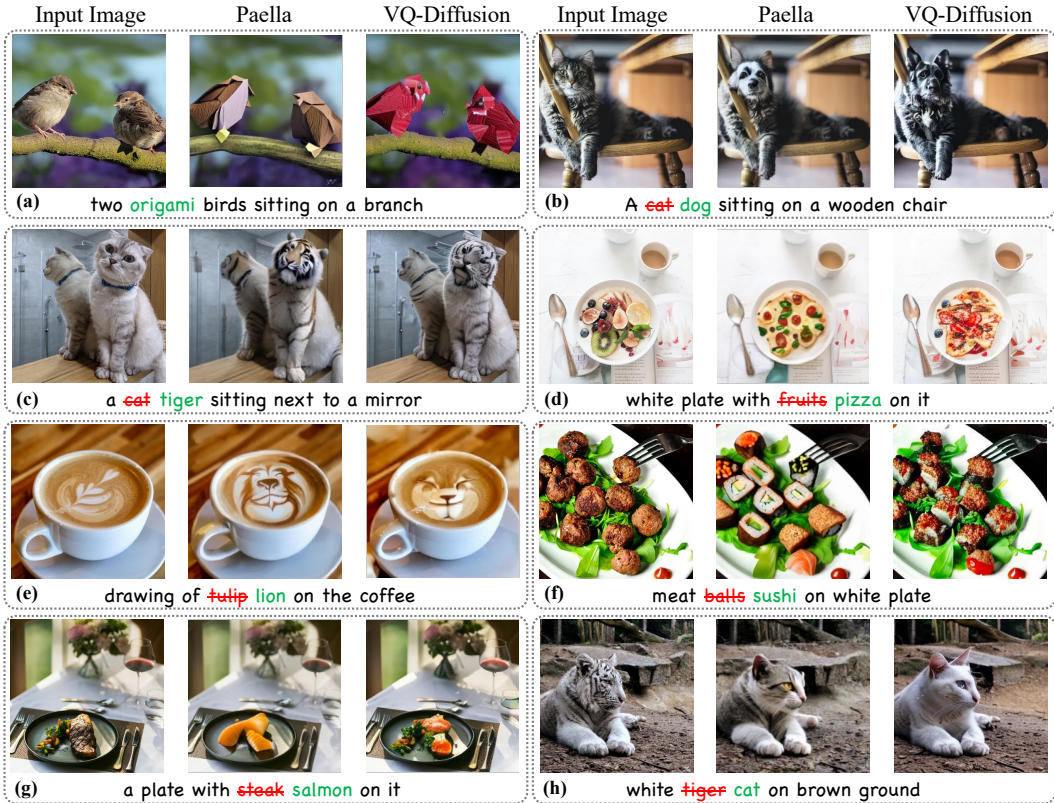


Figure 4: **Visualization of editing results.** Editing results for our method using Paella and VQ-Diffusion are presented, along with their corresponding prompts. The results demonstrate that our method can effectively modify the input image according to the target prompt while preserving the image structure. Editing with masked generative model (Paella (Rampas et al., 2022)) is more stable and easier than with multinomial diffusion models (VQ-Diffusion (Gu et al., 2022)).

Through this process, we assess the model’s capability to invert and generate text that aligns with a specified sentiment while retaining the original sentence’s structural elements.

Inversion Process. In our experiment, we specifically focus on inverting the second sentence, indicated as red in Table 6, while keeping the first sentence intact (black), as it usually contains essential context. During the reverse process, we aim to reconstruct/edit the second sentence by recovering it from the noised tokens acquired in the inversion phase.

Dataset	
1. Positive Sentiment:	Thanks to her efforts. The event was a huge success.
Negative Sentiment:	Despite her efforts. The event was a complete disaster.
2. Positive Sentiment:	This book is definitely interesting. I can’t put it down; it’s full of surprises.
Negative Sentiment:	This book is definitely interesting. I can’t wait to finish it; it’s so predictable.
3. ...	

Dataset Generation. In order to evaluate the editing performance, we designed and proposed a new dataset called Sentiment Editing. The objective is to edit the sentiment of the sentence while preserving the structure of the sentence and also sticking to the theme of the sentence. Here, we demonstrate two sets of sentences in our dataset. Please refer to supplementary materials for the process of generating the dataset and more examples.

	Method		Structure	CLIP Similarity	
	Inverse	Editing	Distance $\times 10^3$ ↓	Whole ↑	Edited ↑
Continuous	DDIM+SD1.4	P2P	69.43*	25.01*	22.44*
	Null-Text + SD1.4	P2P	13.44*	24.75*	21.86*
	Negative-Prompt + SD1.4	P2P	16.17*	24.61*	21.87*
	DDPM-Inversion + SD1.4	Prompt	22.12	26.22	<u>23.02</u>
Discrete	Inpainting + Paella	Prompt	91.10	25.36	23.42
	Ours + Paella	Prompt	11.34	23.79	21.23
	Ours + VQ-Diffusion [†]	Prompt	<u>12.70</u>	23.85	21.02

Table 2: **Quantitative results on image editing performance.** Comparison of our proposed method with the masked inpainting with the discrete diffusion model Paella, as well as continuous diffusion model (Stable Diffusion v1.4) using DDIM inversion. “P2P” refers to Prompt-to-Prompt (Hertz et al., 2022), and “Prompt” denotes editing performed solely through forward edit prompts. Entries marked with an asterisk (*) are cited from Ju et al. (2023). [†]: For VQ-Diffusion, the images are down-sampled to 256×256 . It is important to note that due to differences in base models and editing algorithms, the metrics across methods are not directly comparable. However, our method significantly outperforms both inpainting and strong baselines (e.g., Null-Text Inversion + SD1.4) in terms of structural preservation. As expected, inpainting achieves a high CLIP score since it directly generates image patches based on the target prompt.

Method		Background Preservation			
Inverse	Editing	PSNR ↑	LPIPS $\times 10^3$ ↓	MSE $\times 10^4$ ↓	SSIM $\times 10^2$ ↑
DDIM+SD1.4	P2P	17.87	208.80	219.88	71.14
Ours+Paella	Prompt	27.29	52.90	43.76	89.79

Table 3: **Background Preservation.** Quantitative comparison of background preservation between our proposed method and DDIM+SD 1.4, achieved by masking the edited region and calculating image similarity with the unedited masked image. The inpainting is served as upper bound since only the masked region are edited and background are not modified.

4.2.1 INVERSION RECONSTRUCTION

Similar to the image generation section, we first demonstrate the inversion and reconstruction capabilities of the proposed methods. This process involves inverting the sentences, followed by using the same prompt to generate the reconstructed version of the second sentence.

Evaluation Metric. For reconstruction, we use Hit Rate, which is defined as the proportion of cases where each method generates an identical sentence to the original. In addition, we compute the Semantic Textual Similarity (STS) score by measuring the cosine similarity between the sentence embeddings, using the model proposed by Reimers (2019) *et al.*

Quantitative Analysis. Table 4 compares DICE with Masked Generation using RoBERTa across two metrics: Accuracy and Semantic Textual Similarity. Our method significantly surpasses Masked Generation in both metrics, demonstrating that our z_t latent space effectively captures the information of the sentence being inverted and facilitates its subsequent reconstruction.

4.2.2 SENTENCE EDITING

In this section, we evaluate the editing performance of the proposed inversion method on RoBERTa. In Table 6, the sentence shown in black under the negative prompt column is input during the inversion process. The sentence that is being inverted is displayed in red. For editing, the prompt is then substituted with the black sentence on the right, and noise is added at the end for the forward process. The output of the forward process for the noise is presented in blue.

Method	Metric	
	Accuracy $\times 10^2$ \uparrow	Textual Similarity $\times 10^2$ \uparrow
Inverse+Model		
Masked Generation+RoBERTa	0.0	6.57
Ours+RoBERTa	99.74	99.90

Table 4: **Text Inversion Reconstruction Performance.** Quantitative comparisons of the text reconstruction performance by Masked Generation and DICE method using RoBERTa as the language model.

Method	Metric	
	Structure Preservation $\times 10^2$ \uparrow	Sentiment Correctness $\times 10^2$ \uparrow
Inverse+Model		
Masked Generation+RoBERTa	29.80	12.94
Ours+RoBERTa	94.76	72.51

Table 5: **Text Editing Performance.** Evaluation of the text editing performance between Masked Generation and DICE using ChatGPT as a classifier.

Evaluation Metric. For the sentence editing task, we evaluate the generated sentences based on two criteria: (1) structural preservation, which assesses whether the sentence structure is retained, and (2) sentiment correctness, which evaluates whether the sentiment of the edited sentence aligns with the sentiment of the original prompt. Both the structural preservation rate and sentiment correctness rate are calculated using ChatGPT-4 (Achiam et al., 2023) as a classifier. The details of using ChatGPT for evaluation can be reviewed in Supplementary Materials.

Results. Table 5 presents a comparative analysis of two text editing methods that both employ RoBERTa, focusing on the effectiveness in terms of Structure Preservation and Sentiment Correctness. Our method significantly outperforms masked generation in both metrics. This difference highlights the superior capability of our inversion method to encode the original structure of the text in the latent space and the flexibility to adjust its sentiment more accurately. In Table 6, we demonstrate both the initial prompt and the edited result. Our approach retains the sentence structure of the negative prompt while modifying its sentiment to a more positive one.

Negative Prompt	Our Edited Results
Negative Sentiment: This book is definitely interesting. I can't wait to finish it; it's so predictable.	Positive Sentiment: This book is definitely interesting. I can't wait to see it; it sounds so beautiful.
Negative Sentiment: The new office space is fantastic. It's cramped and lacks proper facilities.	Positive Sentiment: The new office space is fantastic. It's spacious and has great facilities.
Negative Sentiment: Despite her efforts. The event was a complete disaster.	Positive Sentiment: Thanks to her efforts. This event was a fantastic comedy game.
Negative Sentiment: Regarding the lecture. It was dull and confusing.	Positive Sentiment: Regarding the lecture. It was clear and surprising.
Negative Sentiment: Despite the initial problems. The project ended in failure.	Positive Sentiment: Despite the initial problems. New project still in progress.
Negative Sentiment: Regarding the new app. It's complicated and not useful.	Positive Sentiment: Regarding the new app. It's On and It's Epic.
Negative Sentiment: Reflecting on my environmental initiatives. It's challenging to maintain, and progress is slow.	Positive Sentiment: Reflecting on my environmental initiatives. It's easy to understand, and progress is undeniable.

Table 6: **Editing results of our method with RoBERTa.** The sentences in black are the prompts used for inversion and editing in their respective column. The sentence in red is the one being inverted, and the blue sentence represents the editing result.

5 CONCLUSION

In this paper, we introduced DICE (Discrete Inversion for Controllable Editing), an inversion algorithm for discrete diffusion models, including multinomial diffusion and masked generative models. By leveraging recorded noise sequences and masking patterns during the reverse diffusion process,

DICE enables accurate reconstruction and flexible editing of discrete data without the need for pre-defined masks or cross-attention manipulation. Our experiments across multiple models and modalities, such as images and text, demonstrate the effectiveness of DICE in preserving data fidelity while enhancing editing capabilities. Furthermore, we demonstrate the potential of DICE for converting RoBERTa, a model traditionally focused on data understanding, into a generative model for text generation and editing. We believe that DICE enhances the capabilities of discrete generative models, offering new opportunities for fine-grained content manipulation in discrete spaces.

REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarlow, and Rianne Van Den Berg. Structured denoising diffusion models in discrete state-spaces. *Advances in Neural Information Processing Systems*, 34:17981–17993, 2021.
- Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended diffusion for text-driven editing of natural images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 18208–18218, 2022.
- Sam Bond-Taylor, Peter Hessey, Hiroshi Sasaki, Toby P Breckon, and Chris G Willcocks. Unleashing transformers: Parallel token prediction with discrete absorbing diffusion for fast high-resolution image generation from vector-quantized codes. In *European Conference on Computer Vision*, pp. 170–188. Springer, 2022.
- Andrew Campbell, Joe Benton, Valentin De Bortoli, Thomas Rainforth, George Deligiannidis, and Arnaud Doucet. A continuous time framework for discrete denoising models. *Advances in Neural Information Processing Systems*, 35:28266–28279, 2022.
- Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman. Maskgit: Masked generative image transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11315–11325, 2022.
- Huiwen Chang, Han Zhang, Jarred Barber, AJ Maschinot, Jose Lezama, Lu Jiang, Ming-Hsuan Yang, Kevin Murphy, William T Freeman, Michael Rubinstein, et al. Muse: Text-to-image generation via masked generative transformers. *arXiv preprint arXiv:2301.00704*, 2023.
- Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. *Advances in neural information processing systems*, 31, 2018.
- Hyungjin Chung, Jeongsol Kim, Michael T Mccann, Marc L Klasky, and Jong Chul Ye. Diffusion posterior sampling for general noisy inverse problems. *arXiv preprint arXiv:2209.14687*, 2022.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34, 2021.
- Patrick Esser, Robin Rombach, Andreas Blattmann, and Bjorn Ommer. Imagebart: Bidirectional context with multinomial diffusion for autoregressive image synthesis. *Advances in neural information processing systems*, 34:3518–3532, 2021a.
- Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *CVPR*, pp. 12873–12883, 2021b.
- Itai Gat, Tal Remez, Neta Shaul, Felix Kreuk, Ricky TQ Chen, Gabriel Synnaeve, Yossi Adi, and Yaron Lipman. Discrete flow matching. *arXiv preprint arXiv:2407.15595*, 2024.
- 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 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10696–10706, 2022.

- Ligong Han, Jian Ren, Hsin-Ying Lee, Francesco Barbieri, Kyle Olszewski, Shervin Minaee, Dimitris Metaxas, and Sergey Tulyakov. Show me what and tell me how: Video synthesis via multimodal conditioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3615–3625, 2022.
- Ligong Han, Yinxiao Li, Han Zhang, Peyman Milanfar, Dimitris Metaxas, and Feng Yang. Svdiff: Compact parameter space for diffusion fine-tuning. *arXiv preprint arXiv:2303.11305*, 2023.
- Ligong Han, Song Wen, Qi Chen, Zhixing Zhang, Kunpeng Song, Mengwei Ren, Ruijiang Gao, Anastasis Stathopoulos, Xiaoxiao He, Yuxiao Chen, et al. Proxedit: Improving tuning-free real image editing with proximal guidance. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 4291–4301, 2024.
- Xiaoxiao He, Chaowei Tan, Ligong Han, Bo Liu, Leon Axel, Kang Li, and Dimitris N Metaxas. Dmcvr: Morphology-guided diffusion model for 3d cardiac volume reconstruction. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 132–142. Springer, 2023.
- Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*, 2022.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.
- Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022.
- Emiel Hoogeboom, Didrik Nielsen, Priyank Jaini, Patrick Forré, and Max Welling. Argmax flows and multinomial diffusion: Learning categorical distributions. *Advances in Neural Information Processing Systems*, 34:12454–12465, 2021.
- Inbar Huberman-Spiegelglas, Vladimir Kulikov, and Tomer Michaeli. An edit friendly ddpn noise space: Inversion and manipulations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12469–12478, 2024.
- Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*, 2016.
- Xuan Ju, Ailing Zeng, Yuxuan Bian, Shaoteng Liu, and Qiang Xu. Direct inversion: Boosting diffusion-based editing with 3 lines of code. *arXiv preprint arXiv:2310.01506*, 2023.
- Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
- Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- Aaron Lou, Chenlin Meng, and Stefano Ermon. Discrete diffusion language modeling by estimating the ratios of the data distribution. *arXiv preprint arXiv:2310.16834*, 2023.
- Chris J Maddison, Daniel Tarlow, and Tom Minka. A* sampling. *Advances in neural information processing systems*, 27, 2014.
- Chenlin Meng, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073*, 2021.

- Daiki Miyake, Akihiro Iohara, Yu Saito, and Toshiyuki Tanaka. Negative-prompt inversion: Fast image inversion for editing with text-guided diffusion models. *arXiv preprint arXiv:2305.16807*, 2023.
- Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for editing real images using guided diffusion models. *arXiv preprint arXiv:2211.09794*, 2022.
- OpenAI. Sora: Video generation model, 2024. URL <https://openai.com/index/video-generation-models-as-world-simulators>. Accessed: 2024-10-09.
- 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 *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- 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.
- Dominic Rampas, Pablo Pernias, and Marc Aubreville. A novel sampling scheme for text-and image-conditional image synthesis in quantized latent spaces. *arXiv preprint arXiv:2211.07292*, 2022.
- N Reimers. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*, 2019.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 10684–10695, June 2022.
- Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=StlgiaRCHLP>.
- Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in Neural Information Processing Systems*, 32, 2019.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.
- Anastasis Stathopoulos, Ligong Han, and Dimitris Metaxas. Score-guided diffusion for 3d human recovery. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 906–915, 2024.
- Narek Tumanyan, Michal Geyer, Shai Bagon, and Tali Dekel. Plug-and-play diffusion features for text-driven image-to-image translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1921–1930, 2023.
- Alex Wang and Kyunghyun Cho. Bert has a mouth, and it must speak: Bert as a markov random field language model. *arXiv preprint arXiv:1902.04094*, 2019.
- Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.
- Wikipedia contributors. Gumbel distribution — Wikipedia, The Free Encyclopedia. https://en.wikipedia.org/wiki/Gumbel_distribution, 2024. [Online; accessed 8-October-2024].
- Chen Henry Wu and Fernando De la Torre. Unifying diffusion models’ latent space, with applications to cyclediffusion and guidance. *arXiv preprint arXiv:2210.05559*, 2022.
- Chenfei Wu, Lun Huang, Qianxi Zhang, Binyang Li, Lei Ji, Fan Yang, Guillermo Sapiro, and Nan Duan. Godiva: Generating open-domain videos from natural descriptions. *arXiv preprint arXiv:2104.14806*, 2021.

Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 586–595, 2018.

Zhixing Zhang, Ligong Han, Arnab Ghosh, Dimitris N Metaxas, and Jian Ren. Sine: Single image editing with text-to-image diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6027–6037, 2023.

A DETAILS ON MULTINOMIAL DIFFUSION MODELS

Definition of Q_t with mask-and-replace strategy. Following mask-and-replace strategy as:

$$Q_t = \begin{bmatrix} \alpha_t + \beta_t & \beta_t & \beta_t & \cdots & 0 \\ \beta_t & \alpha_t + \beta_t & \beta_t & \cdots & 0 \\ \beta_t & \beta_t & \alpha_t + \beta_t & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \gamma_t & \gamma_t & \gamma_t & \cdots & 1 \end{bmatrix}, \quad (8)$$

given $\alpha_t \in [0, 1]$, $\beta_t = (1 - \alpha_t - \gamma_t)/K$ and γ_t the probability of a token to be replaced with a [MASK] token.

Cumulative transition matrix. The cumulative transition matrix \bar{Q}_t and $q(x_t|x_0)$ can be computed via closed form:

$$\bar{Q}_t \mathbf{v}(x_0) = \bar{\alpha}_t \mathbf{v}(x_0) + (\bar{\gamma}_t - \bar{\beta}_t) \mathbf{v}(K+1) + \bar{\beta}_t \mathbf{1}, \quad (9)$$

where $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$, $\bar{\gamma}_t = 1 - \prod_{i=1}^t (1 - \gamma_i)$, and $\bar{\beta}_t = (1 - \bar{\alpha}_t - \bar{\gamma}_t)/(K+1)$ can be calculated and stored in advance.

B ANALYSIS ON MUTUAL INFORMATION

Proof of Remark 3.1.

Proof. We assumed that \mathbf{x}_0 satisfies standard Gaussian distribution $\mathcal{N}(\mathbf{0}, I_D)$. Since

$$\mathbf{x}_t = \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_t$$

where both \mathbf{x}_{t-1} and $\boldsymbol{\epsilon}_t$ are independent standard Gaussian random variables, \mathbf{x}_t is also standard Gaussian, and in each dimension

$$\text{Cov}(\mathbf{x}_t, \mathbf{x}_{t-1}) = \sqrt{\alpha_t},$$

which leads to

$$\hat{\mu}_t(\mathbf{x}_t) = \mathbb{E}(\mathbf{x}_{t-1}|\mathbf{x}_t) = \sqrt{\alpha_t} \mathbf{x}_t.$$

Therefore,

$$\begin{aligned} \mathbf{z}_t &= \mathbf{x}'_{t-1} - \hat{\mu}_t(\mathbf{x}_t) \\ &= (\sqrt{\alpha_{t-1}} \mathbf{x}_0 + \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon}) - \sqrt{\alpha_t} (\sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}') \\ &= \beta_t \cdot \sqrt{\alpha_{t-1}} \mathbf{x}_0 + \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon} + \sqrt{\alpha_t (1 - \alpha_t)} \boldsymbol{\epsilon}'. \end{aligned}$$

Let

$$E = \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon} + \sqrt{\alpha_t (1 - \alpha_t)} \boldsymbol{\epsilon}'$$

which is a Gaussian error term independent to \mathbf{x}_0 with mean 0 and variance $1 - \alpha_{t-1} + \alpha_t (1 - \alpha_t)$. Thus we can calculate the mutual information

$$\begin{aligned} I(\mathbf{z}_t; \mathbf{x}_0) &= H(\mathbf{z}_t) - H(\mathbf{z}_t|\mathbf{x}_0) \\ &= H(\mathbf{z}_t) - H(E) \\ &= \frac{D}{2} \log(2\pi e (\beta_t^2 \alpha_{t-1} + 1 - \alpha_{t-1} + \alpha_t (1 - \alpha_t))) - \frac{D}{2} \log(2\pi e (1 - \alpha_{t-1} + \alpha_t (1 - \alpha_t))) \\ &= \frac{D}{2} \log\left(\frac{\beta_t^2 \alpha_{t-1} + 1 - \alpha_{t-1} + \alpha_t (1 - \alpha_t)}{1 - \alpha_{t-1} + \alpha_t (1 - \alpha_t)}\right). \end{aligned}$$

□

C IMPLEMENTATION DETAILS

For all reconstruction task, we employ a $\tau = 1.0$ and $\lambda_1 = 1.0$, $\lambda_2 = 0.0$ with 32 sampling steps and 26 renoising steps. For editing tasks, the hyper-parameters are summarized in Table 7.

All models are implemented in PyTorch 2.0 and inferenced on a single NVIDIA A100 40GB.

Editing Experiment		Hyper-parameters			
Method	Configuration	CFG	λ_1	λ_2	τ
Paella	Set 1	10.0	0.7	0.3	0.9
VQ-Diffusion	Set 1	5.0	0.2	0.8	1.0
RoBERTa Sentiment	Set 1	-	0.2	0.8	0.7
	Set 2	-	0.25	0.75	0.7

Table 7: Hyper-parameters for Paella, VQ-Diffusion, and RoBERTa sentiment editing experiments. For sentiment editing task with RoBERTa, we utilize two sets of hyper-parameters empirically due to the variance in the sentence length.

D ABLATION STUDIES

D.1 NOISE INJECTION FUNCTION

Addition. In the main text we have adopted the *addition* function as noise injection function,

$$\tilde{\mathbf{y}} = \log(\boldsymbol{\pi}) + \lambda_1 \cdot \mathbf{z} + \lambda_2 \cdot \mathbf{g}.$$

This is a natural form inspired by the Gumbel-Max trick: thinking of $\lambda_1 \cdot \mathbf{z}$ as a correction term, then $\log(\boldsymbol{\pi}) + \lambda_1 \cdot \mathbf{z}$ is the corrected logit and λ_2 is the inverse of temperature of the logit to control the sharpness of the resulting categorical distribution, as

$$\arg \max (\log(\boldsymbol{\pi}) + \lambda_1 \cdot \mathbf{z} + \lambda_2 \cdot \mathbf{g}) = \arg \max \left(\frac{1}{\lambda_2} (\log(\boldsymbol{\pi}) + \lambda_1 \cdot \mathbf{z}) + \mathbf{g} \right), \quad \lambda_2 > 0.$$

λ_1 then controls how much correction we would like to introduce in the original logit.

Variance preserving. From another perspective, \mathbf{z} is the artificial ‘‘Gumbel’’ noise that could have been sampled to realize the target tokens. Then, if we treat \mathbf{z} as Gumbel noise and want to perturb it with random Gumbel noise, addition does not result in a Gumbel distribution. One way is to approximate this sum with another Gumbel distribution. If $G_1 \sim \text{Gumbel}(\mu_1, \beta_1)$, $G_2 \sim \text{Gumbel}(\mu_2, \beta_2)$ and $G = \lambda_1 G_1 + \lambda_2 G_2$, then the moment matching *Gumbel approximation* for G is

$$\begin{aligned} & \text{Gumbel}(\mu_G, \beta_G), \quad \text{with} \\ & \beta_G = \sqrt{\lambda_1^2 \beta_1^2 + \lambda_2^2 \beta_2^2}, \\ & \mu_G = \lambda_1 \mu_1 + \lambda_2 \mu_2 + \gamma(\lambda_1 \beta_1 + \lambda_2 \beta_2 - \beta_G), \end{aligned}$$

where $\gamma \approx 0.5772$ is the Euler-Mascheroni constant. We consider the *variance preserving* form:

$$\tilde{\mathbf{y}} = \log(\boldsymbol{\pi}) + \sqrt{\lambda_1} \cdot \mathbf{z} + \sqrt{\lambda_2} \cdot \mathbf{g}, \quad \lambda_1 + \lambda_2 = 1.$$

Max. The third way is inspired by the property of Gumbel distribution (Wikipedia contributors, 2024), that if G_1, G_2 are iid random variables following $\text{Gumbel}(\mu, \beta)$ then $\max\{G_1, G_2\} - \beta \log 2$ follows the same distribution. We also consider the *max* function for noise injection:

$$\tilde{\mathbf{y}} = \log(\boldsymbol{\pi}) + \max\{\lambda_1 \cdot \mathbf{z}, \lambda_2 \cdot \mathbf{g}\}.$$

D.2 HYPERPARAMETER SEARCH

In this section, we analyze the impact of varying hyperparameters $\lambda_1, \lambda_2, \tau$, and CFG scale on the quality of image generation and adherence to textual descriptions, quantified through Structure Distance and CLIP similarity. The hyperparameters play specific roles: λ controls the amount of noise introduced in each reverse step, τ governs the percentage of tokens replaced with random tokens during inversion, and Classifier-Free Guidance (CFG) scales the influence of the text prompt during image synthesis. To limit the search space and simplify the ablation, we choose $\lambda_1 = \lambda$ and $\lambda_2 = 1 - \lambda$ and vary the value of λ . Evaluation metrics are given in Figure 5.

Effect of λ_1 and λ_2 : With a fixed CFG of 10.0, the graphs indicate that increasing λ results in a rise in Structure Distance, suggesting a decline in structural integrity of the images. This increase in

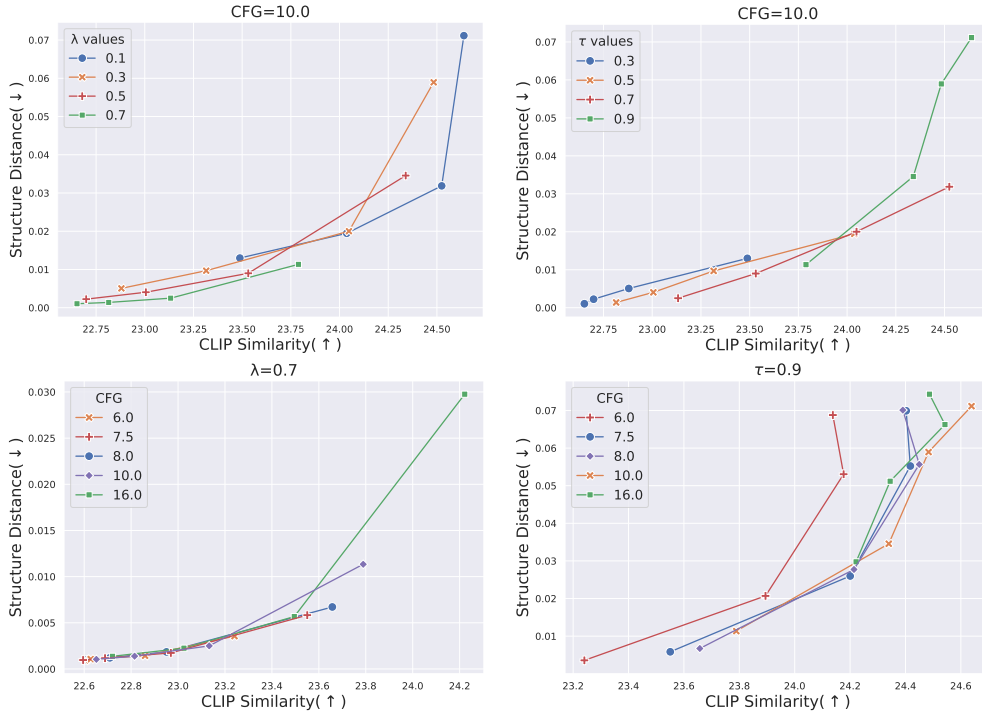


Figure 5: **The effect of hyperparameters $\lambda_1, \lambda_2, \tau, \text{CFG}$ on the Structure Distance (\downarrow) and CLIP similarity (\uparrow) with addition function as noise inject function.** In our implementation, to limit the search space, we choose $\lambda_1 = \lambda$ and $\lambda_2 = 1 - \lambda$ for simplicity.

noise appears to allow for greater exploration of the generative space at the expense of some loss in image clarity.

Effect of τ : Higher τ values, particularly at 0.9, show a notable rise in Structure Distance as CLIP similarity increases. This implies that more token replacement can lead to images that align better with the text prompts but may suffer in maintaining structural fidelity, likely due to x_T contains less information of the original image while λ injects additional noise during editing phase.

Effect of CFG Scale: Varying CFG at a fixed λ of 0.7 and τ of 0.9 reveals that higher CFG values substantially improve Structure Distance, but to an extent (CFG of 10). Beyond this point, further increases in CFG do not yield significant improvements in structural quality, indicating a diminishing return on higher guidance levels. This plateau suggests that while increasing CFG helps in aligning the generated images more closely with the text prompts initially, the benefits in structural integrity and clarity become less visible as CFG values exceed a certain threshold. This finding underscores the need for a balanced approach in setting CFG, where too much guidance may not necessarily lead to better outcomes in terms of image quality and fidelity to the textual description.

Effect of noise injection function: We also conducted evaluations using a variance-preserving noise injection function by setting $\lambda_1 = \sqrt{\lambda}$ and $\lambda_2 = \sqrt{1 - \lambda}$. The results of these experiments are presented in Figure 6. As for the `max` function, we performed a manual inspection of the visual examples generated with this function. The quality of these examples was noticeably inferior, we therefore omit the corresponding evaluation curves from our analysis.

In conclusion, this ablation study demonstrates that increasing λ and τ can enhance adherence to text prompts through broader explorations in generative spaces, yet this benefit is offset by a decrease in the structural quality of the images. On the other hand, raising CFG values enhances the structural integrity of images to a certain threshold, after which the improvements plateau, indicating a ceiling to the effectiveness of higher CFG settings. This analysis offers empirical guidance for selecting hyperparameters, balancing the trade-offs between text alignment and image quality to optimize image synthesis outcomes.

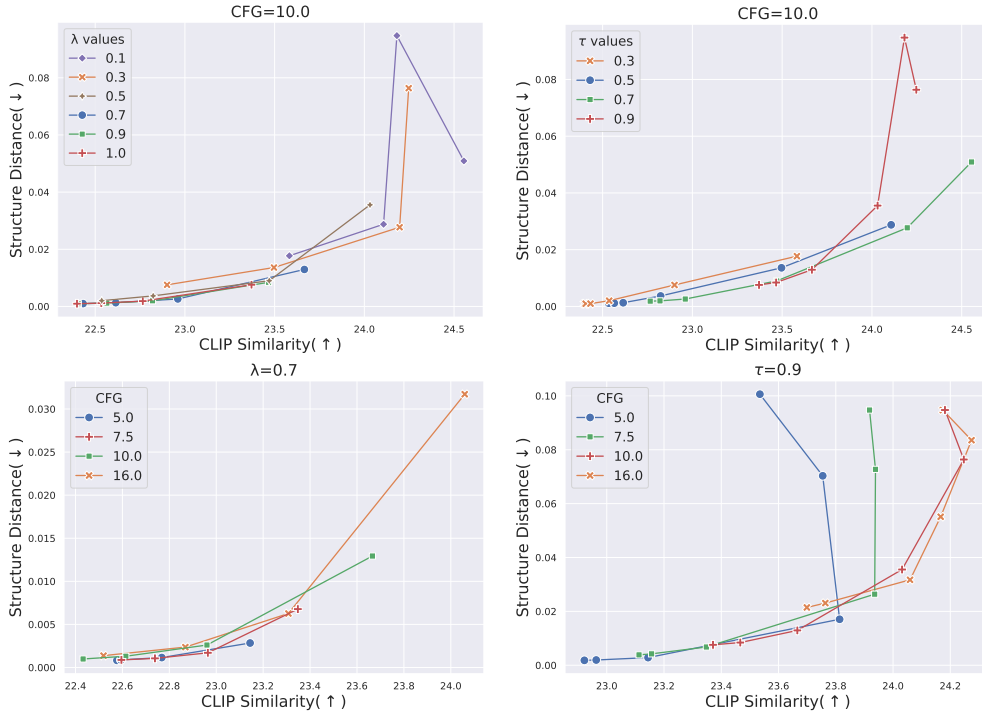


Figure 6: **The effect of hyperparameters λ_1, λ_2 with variance preserving scheme.** We set $\lambda_1 = \sqrt{\lambda}$ and $\lambda_2 = \sqrt{1 - \lambda}$.

E ADDITIONAL RESULTS ON IMAGE EDITING

Reconstruction result with Paella. In Figure 8 we demonstrates the inversion reconstruction result with Paella using our proposed method.

Image editing with diversity. As shown in Figure 10, our method enables diverse image editing results through stochastic variation. The first three rows demonstrate the impact of varying both the inversion masks and the injected Gumbel noise, while the last two rows focus on variations produced by changing only the inversion masks.

F ADDITIONAL RESULTS ON TEXT EDITING

Dataset generation. To generate the dataset, we utilize ChatGPT-4o with the following prompt:

User

Generate 200 pairs of sentences that contains the same meaning, but one with positive sentiment and one with negative sentiment. For both positive sentiment and negative sentiment, you need to write two sentences with the first part being a hint of the sentiment and the second part being the actual content. The first part for both sentences should be same. write in the format like:
 hint. positive.
 hint. negative.
 Make sure that there are two lines for each pairs. Also, the hint should provide enough context and both positive and negative sentiment should be related to the hint. Do not repeat the hint, also make sure that there is only two sentences in each of the line, one is the hint and the other is about the sentiment.

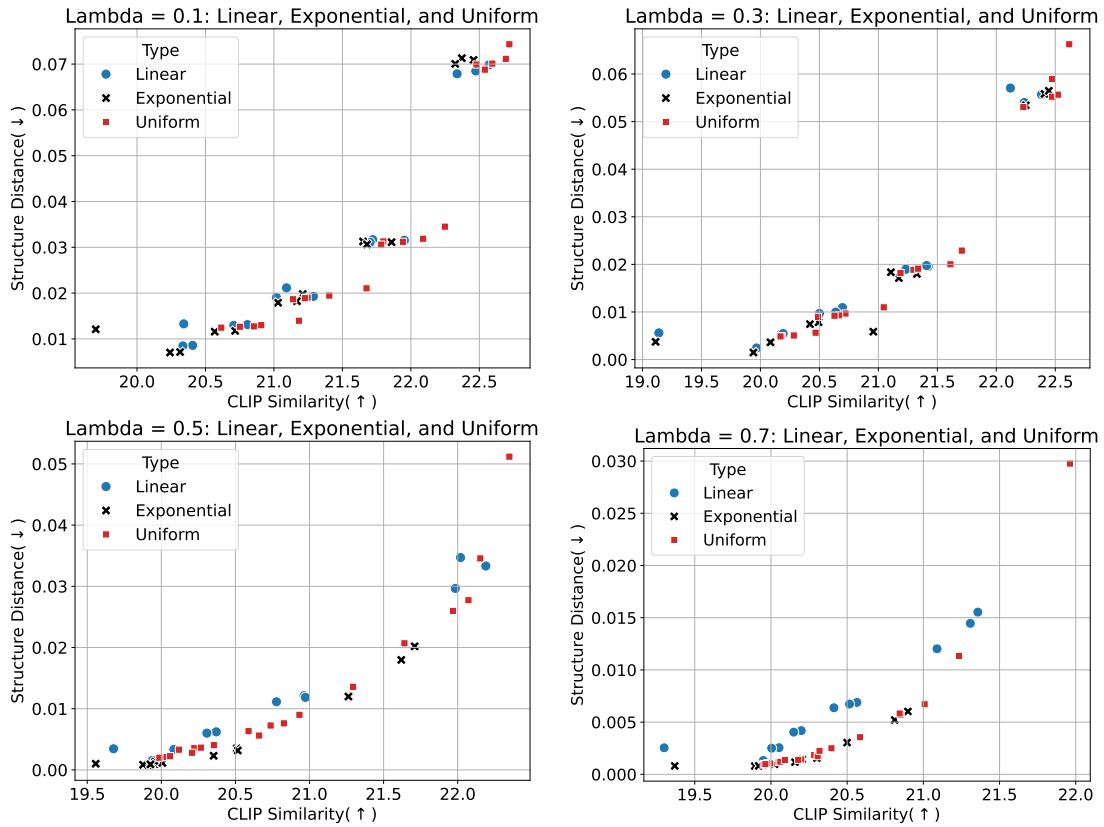


Figure 7: **The effect of different λ schedule on the Structure Distance (\downarrow) and CLIP similarity (\uparrow).** In our implementation, to limit the search space, we choose $\lambda_1 = \lambda$ and $\lambda_2 = 1 - \lambda$ for simplicity.

ChatGPT

1. Thanks to her efforts. The event was a huge success.
Despite her efforts. The event was a complete disaster.
2. ...

The sentences is then added with a prefix to indicates the sentiment of the context. Here we demonstrates a subset of our generated dataset:

1. Positive Sentiment: Thanks to her efforts. The event was a huge success.
Negative Sentiment: Despite her efforts. The event was a complete disaster.
2. Positive Sentiment: This book is definitely interesting. I can't put it down; it's full of surprises.
Negative Sentiment: This book is definitely interesting. I can't wait to finish it; it's so predictable.
3. Positive Sentiment: The new office space is fantastic. It's spacious and perfect for productivity.
Negative Sentiment: The new office space is fantastic. It's cramped and lacks proper facilities.
4. Positive Sentiment: Thanks to her efforts. The event was a huge success.
Negative Sentiment: Despite her efforts. The event was a complete disaster.
5. Positive Sentiment: Regarding the lecture. It was insightful and engaging.
Negative Sentiment: Regarding the lecture. It was dull and confusing.
6. Positive Sentiment: Despite the initial problems. The project was a success.
Negative Sentiment: Despite the initial problems. The project ended in failure.
7. Positive Sentiment: Regarding the new app. It's user-friendly and very helpful.
Negative Sentiment: Regarding the new app. It's complicated and not useful.
8. Positive Sentiment: Reflecting on my environmental initiatives. Implementing changes has reduced my carbon footprint.
Negative Sentiment: Reflecting on my environmental initiatives. It's challenging to maintain, and progress is slow.
9. Positive Sentiment: The business proposal was well-received. The ideas were innovative, and the presentation was convincing.
Negative Sentiment: The business proposal was rejected. The ideas were impractical, and the presentation was unconvincing.
10. Positive Sentiment: The training program was highly effective. It boosted skills and confidence, and everyone left motivated.
Negative Sentiment: The training program was ineffective. It didn't teach much, and most people left feeling unmotivated.
11. ...

Evaluation. Below, we demonstrate the prompt used for evaluating the editing results:

User

Given three sentences, confirm that the second sentence is roughly the same sentence structure as the first sentence, then confirm that the second sentence has positive sentiment. Output only two numbers with each number indicating whether the corresponding criteria is satisfied. Use 1 for satisfied and 0 for not satisfied. The sentences are given below:
The event was a complete disaster.
This event was a fantastic comedy game.

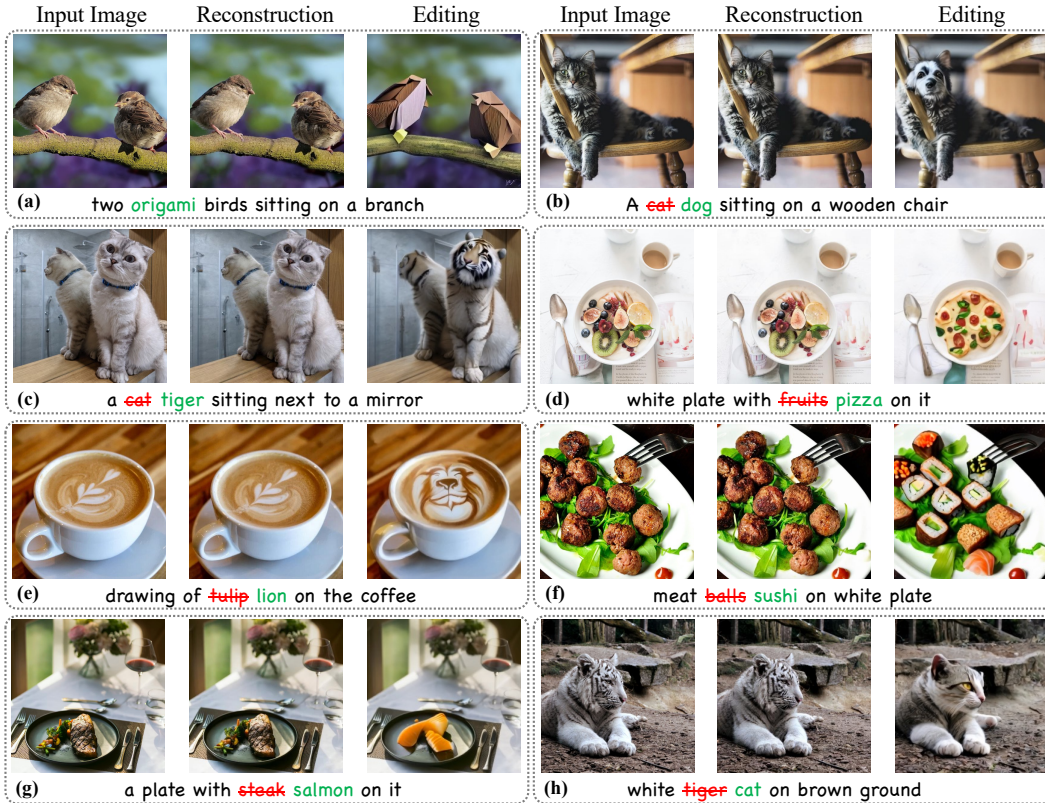


Figure 8: Reconstruction and editing result with DICE+Paella.

ChatGPT

1 1

Comparison between masked inpainting and DICE. In Figure 9 we demonstrates the reconstruction and editing results with our DICE and Masked Inpainting.

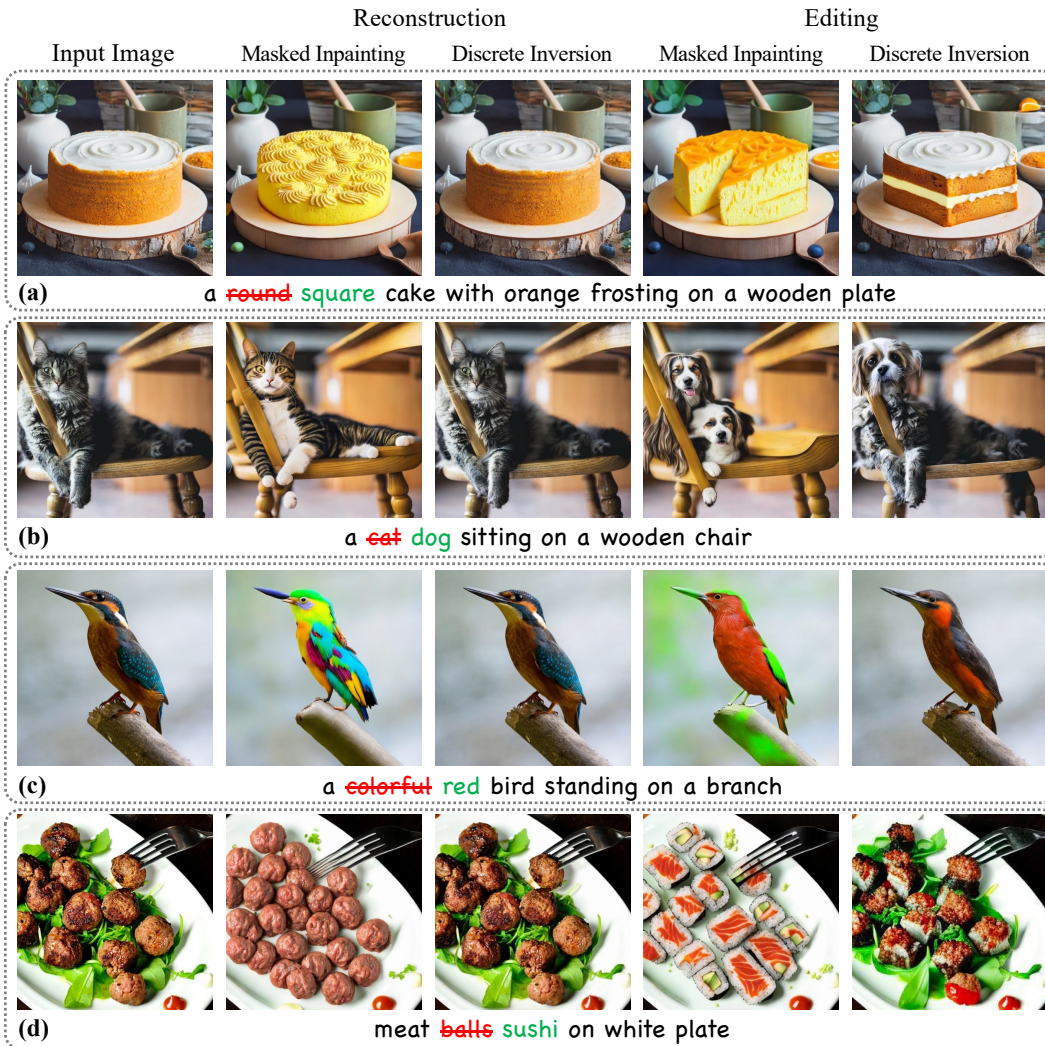


Figure 9: **Reconstruction and editing result with DICE and masked inpainting.** Notice that for reconstruction, we use the **red** prompt, but for editing we use the **green** prompt.

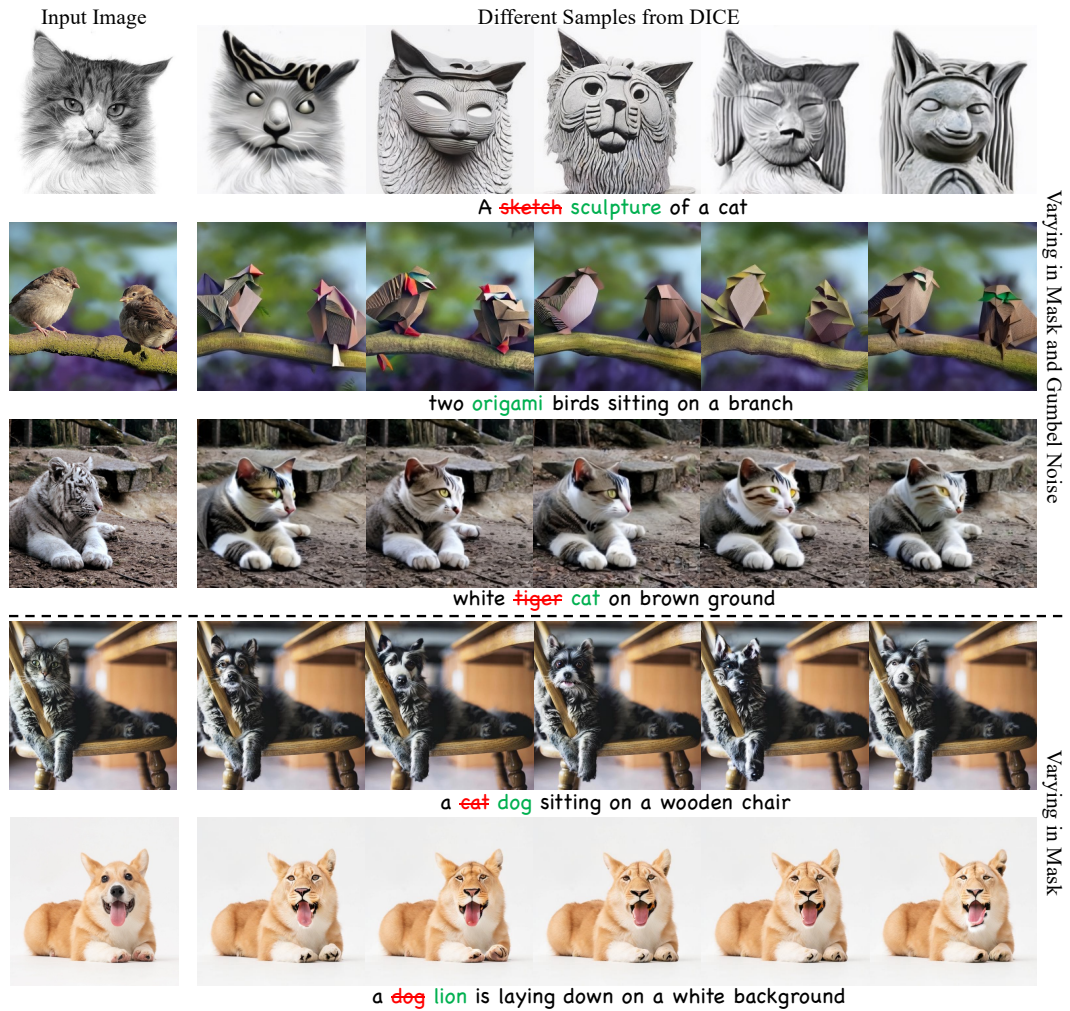


Figure 10: **Image Editing with Diversity.** Due to the stochastic nature of our method, we can generate diverse outputs. The first three rows illustrate variations in both inversion masks and injected Gumbel noise ($\lambda_1 = 0.7$, $\lambda_2 = 0.3$). The last two rows demonstrate variations using only inversion masks ($\lambda_1 = 1$, $\lambda_2 = 0$).