# FreeCond: Free Lunch in the Input Conditions of Text-Guided Inpainting

Teng-Fang Hsiao, Bo-Kai Ruan, Sung-Lin Tsai, Yi-Lun Wu, Hong-Han Shuai
National Yang Ming Chiao Tung University

tfhsiao.ee13@nycu.edu.tw, bkruan.ee11@@nycu.edu.tw, tsai412504004.ee12@nycu.edu.tw

yilun.ee08@nycu.edu.tw, hhshuai@nycu.edu.tw

Figure 1. Comparison of T2I inpainting methods with FreeCond, applied across various mask types: "multi-masks" (column 1 and 4), "precise masks" (columns 2, 5, and 6), and "rough masks" (columns 3 and 7) with complex prompts and unrelated image contexts. By integrating FreeCond, existing inpainting baselines obtain better "instruction-following" performance.

## Abstract

*In this study, we aim to determine and solve the deficiency of Stable Diffusion Inpainting (SDI) in following the instruction of both prompt and mask. Due to the training bias from masking, the inpainting quality is hindered when the prompt instruction and image condition are not related. Therefore, we conduct a detailed analysis of the internal representations learned by SDI, focusing on how the mask input influences the cross-attention layer. We observe that adapting text key tokens toward the input mask enables the model to selectively paint within the given area. Leveraging these insights, we propose FreeCond, which adjusts only the input mask condition and image condition. By increasing the latent mask value and modifying the frequency of image condition, we align the cross-attention features with the model's training bias to improve generation quality without additional computation, particularly when user inputs are complicated and deviate from the training setup. Extensive experiments demonstrate that FreeCond can enhance any SDI-based model, e.g., yielding up to a 60% and 58% improvement of SDI and SDXLI in the CLIP score. The code and appendix are available in our repository at https://github.com/basiclab/FreeCond [1].*

## 1. Introduction

Text-to-image (T2I) inpainting seeks to fill specified masked areas based on user-provided text prompts. Stable Diffusion Inpainting (SDI), a tailored adaptation of

---

[1]Due to arXiv file size limitations, we provide an abbreviated version of the paper here; the full version can be accessed in the repository.

"*A fluffy panda juggling teacups*"   "*Three penguins playing music*..."

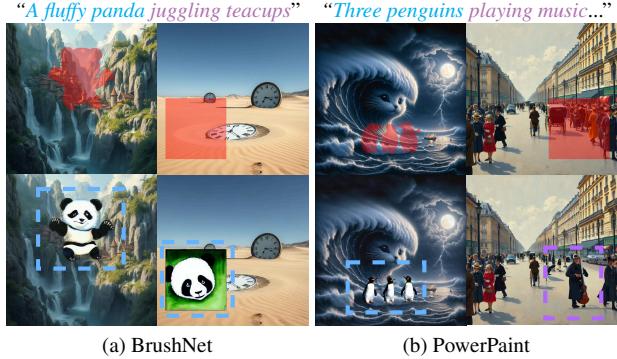(a) BrushNet                           (b) PowerPaint

Figure 2. Comparison of existing SOTA methods. BrushNet rigidly follows the mask instructions but only partially adheres to the prompt. PowerPaint produces outputs that are harmonious with the image context but at the cost of reduced prompt-adherence. FreeCond addresses these limitations, as shown in Fig. 1.

Stable Diffusion [26], is widely used for its effectiveness in achieving high-quality inpainting aligned with text prompts. However, the SDI training process employs a random masking strategy, which often hinders the model's ability to follow prompts accurately and fit masks precisely, especially when the prompt lacks contextual relevance to the surrounding image. We refer to these dual issues of "prompt-adherence" and "mask-fitting" collectively as the "instruction-following" problem: the model prioritizes contextual coherence over generating content that strictly adheres to both the complex prompt details and the user-specified mask, as illustrated in the second row of Fig. 1.

To address these limitations, methods such as HD-Painter [20], BrushNet [13], and PowerPaint [42] have been developed. BrushNet, for example, leverages segmentation-based training data and a ControlNet-like structure, allowing it to learn a direct link between the input mask and prompt. PowerPaint incorporates training with dilated segmentation masks and task-specific tokens, enabling flexible object inpainting that better conforms to varied shapes. While these approaches effectively reduce the "mask-fitting" issue, they primarily optimize for simple prompts and often lack the generalization capability required for complex prompt adherence, as shown in Fig. 2. Observing the limitations of training-only modifications, **we propose to directly modify the model's behavior by adjusting its learned mechanism.**

In this paper, we contend that effective instruction-following relies on the differential noise predictions: *conditional versus unconditional*, modulated via classifier-free guidance (CFG) [9]. This differential is notably manifested in the outputs of the cross-attention layer, where prompt tokens receive significantly higher attention within the masked area than in surrounding regions. Consequently, the query and key features within the cross-attention layer must dynamically adapt to the input mask. This adapta-

tion focuses on generating coherent new content within the masked regions while concurrently preserving the integrity of the surrounding context by selectively enhancing features related to the mask in the corresponding channels.

To address these challenges, we propose FreeCond, a training-free method that requires no extra computation. Specifically, FreeCond filters the high-frequency components of the image condition in the early diffusion steps, reducing the contextual information of the image condition. Furthermore, FreeCond induces a stronger feature shift in the cross-attention layer by scaling the mask condition, enabling stronger activation of the masked area. By adjusting the input conditions, FreeCond significantly enhances both prompt-adherence and mask-fitting while preserving overall harmony. Notably, FreeCond, as a more general form of noise prediction function, can be seamlessly integrated with other SDI-based methods [13, 20, 23, 39, 42], as it directly enhances the original SDI backbone

Finally, we propose FCIBench (FreeCond Inpainting Benchmark), a new benchmark with 600 inpainting pairs, to handle complex inpainting scenarios. Compared with existing inpainting benchmarks [13, 16], FCIBench includes precise masks, rough masks, and multi-masks, along with complex prompts that are unrelated to image condition, as shown in the first row of Fig. 1. This variety enables a more comprehensive evaluation of SDI across diverse inpainting conditions. Our expanded benchmark thus helps a thorough assessment of the performance of different models across varied prompts and mask configurations. Experimental results demonstrate that FreeCond consistently improves performance across models and benchmarks, achieving a 60% increase in CLIP score [24] of SDI backbone and a 1% increase of existing SOTA. The contributions can be summarized as follows.

- We conduct an in-depth analysis of the SDI model's mechanism, enhancing the interpretability by examining its learned bias of relying on image context and its capability to selectively inpaint within the masked area.
- We introduce FreeCond, a novel noise prediction function, that addresses the instruction-following limitations of SDI-based models without adding computational overhead, especially in scenarios where the complex prompt instruction is unrelated to the image condition.
- We provide FCIBench to evaluate inpainting methods in scenarios across precise mask, rough mask, and multi-mask settings, along with complex prompts that are unrelated to image conditions, extending the evaluation to more diverse scenarios.

2

## 2. Related Works

### 2.1. Image Inpainting

Image inpainting focuses on repainting specified regions while ensuring coherence with the surrounding image. Various non-text-guided inpainting methods have been developed to achieve this [4, 12, 18, 22, 36–38, 41], alongside the emergence of text-to-image inpainting methods [1, 21, 25, 29, 33, 39]. SDI [26] pioneered the integration of a random masking strategy into its training objective, producing harmonized outputs. However, this approach often prioritizes image conditioning over adherence to instructions. To improve this, recent methods [13, 20, 31, 42] have introduced solutions focused primarily on enhancing mask-fitting. Despite these advancements, these methods often lack prompt-adherence when handling complex instructions. In contrast, FreeCond leverages insights into the inpainting mechanism to improve instruction-following across any SDI-based inpainting model, achieving a balanced performance between mask-fitting and prompt-adherence.

### 2.2. Unveiling the Mechanism of T2I Models

T2I diffusion models possess powerful image-generation capabilities. To fully harness this potential, recent works have explored various training-free modifications based on in-depth analyses of different components [3, 5–8, 10, 11, 17, 19, 34, 35]. Notably, FreeU [28] reveals that the skip-connection primarily retains texture details, while the backbone captures more semantic information. By adjusting the balance between them, FreeU enhances semantic accuracy with minimal impact on detail and computational costs.

Our analysis reveals that the image and mask conditions play similar roles: the image condition provides contextual details, while the mask condition regulates prompt influence. By modulating both conditions, we achieve improved instruction-following without additional costs, providing a "free lunch" in performance enhancement.

## 3. Analysis of SDI model

In this section, we provide an in-depth analysis of the SDI model. Firstly, we identify the training bias of the SDI model in Sec. 3.2 that the model heavily relies on the image condition to generate the content. This mechanism can cause the generated output unrelated to the input prompt. Secondly, we demonstrate in Sec. 3.3 that increasing the size of the mask provides SDI with more potential solutions to integrating prompt instructions into the image context, resulting in outputs that more closely follow the instruction. Finally, in Sec. 3.4, we test the hypotheses that interpret the internal mechanics of the SDI model, with a particular focus on the relationship between the inpainting conditioning and the cross-attention layer. The analysis of the SDI model helps us understand the factors leading to successful instruction-following inpainting.

### 3.1. Stable Diffusion Inpainting Model (SDI)

The SDI model receives a prompt $p$, an image $I \in \mathbb{R}^{H \times W \times 3}$, and a mask $M \in \mathbb{R}^{H \times W}$ that specifies the inpainting area. To prevent the model from copying content directly from $I$, the masked area of $I$ is set to zero, yielding $I^c = (\mathbf{1} - M) \odot I$. Since the UNet operates in VAE [14] latent space, $I^c$ is encoded as the image condition $z^c = \mathcal{E}(I^c) \in \mathbb{R}^{(H/4) \times (W/4) \times 4}$ and $\mathcal{E}$ denotes the pretrained VAE encoder. To match the input size, the SDI model uses an interpolated mask condition $M^c \in \mathbb{R}^{(H/4) \times (W/4)}$, created by downsampling $M$ with nearest-neighbor interpolation. The final inpainting noise prediction from the diffusion model is $\epsilon_\theta(z_t, z^c, M^c, t, p)$, where $\epsilon_\theta$ represents the SDI UNet model, $z_t$ is the noise latent at timestep $t \in [0, T]$, and initial noise $z_T$ is sampled from $\mathcal{N}(\mathbf{0}, \mathbf{I})$. To control the influence of prompt $p$, we follow classifier-free guidance (CFG) [9], modifying the noise prediction with a scaling parameter $w \in \mathbb{R}$:

$$\hat{\epsilon}_\theta(z_t, z^c, M^c, t, p) = \epsilon_\theta(z_t, z^c, M^c, t, \varnothing)$$
$$+ w\big(\epsilon_\theta(z_t, z^c, M^c, t, p) - \epsilon_\theta(z_t, z^c, M^c, t, \varnothing)\big) \quad (1)$$

To evaluate the inpainting results in our study, we use a set of six metrics adopted from BrushBench [13]. These metrics cover three key areas: **Image Quality**, measured by Image Reward (IR) [32], HPS [30], and Aesthetic Score (AS) [27]; **Background Preservation**, assessed using PSNR and LPIPS[40]; and **Instruction Following**, evaluated through CLIP [24]. Additionally, we introduce a novel Intersection-over-Union (IoU) score to specifically capture the mask-fitting quality, complementing the CLIP Score by distinguishing mask accuracy from prompt-adherence. This score is computed by the IoU between input mask and auto-labeled mask via SAM [15], as detailed in Appendix.

### 3.2. Influence of Image Condition $z^c$

The random masking strategy of SDI is via creating masked data by randomly masking 25% of image areas in LAION-5B [27], aiming to enhance generalizability across various inputs. To investigate the mask distribution under this random strategy, we define three mask placements: "not masked," "partially masked," and "fully masked." Although the exact SDI training mask distribution is not accessible, our analysis on the COCO dataset as a surrogate reveals that, with a 25% mask coverage, over 80% of training data falls under the "not masked" or "partially masked" categories (see the Appendix for details). Thus, we hypothesize that **the random masking design optimizes SDI for maintaing overall image harmony rather than strict prompt-adherence**. For instance, in the second row of Fig. 3a, when

*"zebra"*

*"airplane"*

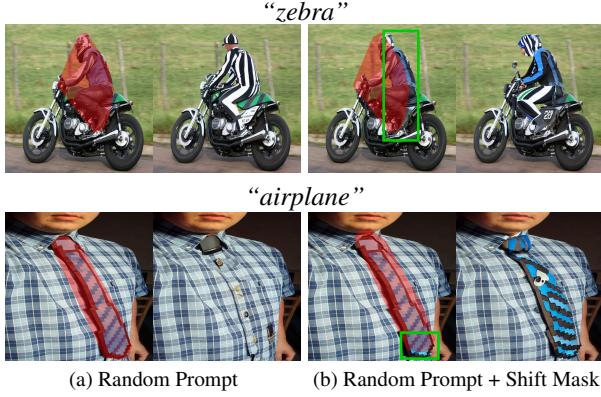(a) Random Prompt      (b) Random Prompt + Shift Mask

Figure 3. Visualization of contextual influence: A random prompt, unrelated to the image condition, is assigned. The input mask is shown in columns 1 and 3, along with the corresponding prompt, while shifted areas are highlighted with a green frame. The resulting outputs are displayed in columns 2 and 4.

|  | AS | LPIPS | $\text{CLIP}_{\text{IN}}$ | $\text{CLIP}_{\text{GT}}$ |
|---|---|---|---|---|
| Original | 5.89 | 0.04 | 18.66 | 18.66 |
| Shift Mask | 5.89 | 0.04 | 18.61 | 18.61 |
| Random Prompt | 5.79 | 0.04 | 15.57 | 15.62 |
| Random + Shift | 5.66 | 0.03 | 15.09 | 16.38 |

Table 1. Table of SDI model on different settings toward inpainting COCO dataset. $\text{CLIP}_{\text{IN}}$ denotes the CLIP similarity toward the input prompt, while $\text{CLIP}_{\text{GT}}$ for the ground truth prompt.

prompted with *"airplane"* instead of the actual ground-truth *"tie"*, SDI ignores the prompt and generates contextually consistent but unrelated content.

To validate this, we select 600 mask and ground-truth pairs from COCO [16]. In the "random prompt" setting, ground truth prompts are replaced with random ones. Results in Tab. 1 show that the CLIP score for ground truth prompts ($\text{CLIP}_{\text{GT}}$) closely matches that for input prompts ($\text{CLIP}_{\text{IN}}$), indicating that SDI favors contextual coherence over strict prompt-adherence. For example, in Fig. 3a, SDI interprets contextual hints by generating an object riding on a motorcycle rather than a zebra as prompted, producing a person with zebra-patterned clothing instead. In the "random prompt + shifted mask" setting, we shift the mask by 25 pixels to add more ground-truth information into the image condition $z^c$. This adjustment decreases instruction-following accuracy, reflected by a 3% drop in $\text{CLIP}_{\text{IN}}$ and a 5% increase in $\text{CLIP}_{\text{GT}}$ compared to the "random prompt" case. In Fig. 3b, when an object, like a person or tie, is visible in $z^c$, SDI can revert the whole object. This analysis confirms that **the context provided by $z^c$ significantly limits SDI's instruction following**.



(a) Ground Truth Prompt + Mask Dilation
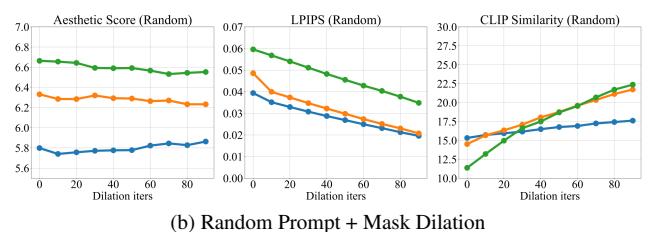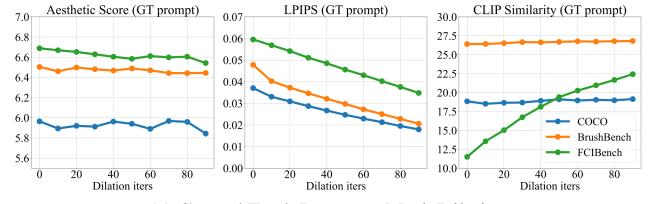


(b) Random Prompt + Mask Dilation

Figure 4. Illustration of mask size impact on inpainting metrics.

## 3.3. Influence of Input Mask $M$

In the preceding analysis, we observe that the inpainting result can be hugely guided by the image condition $z^c$, leading to its deficiency in the instruction following. Here, we explore how to adjust the input mask $M$ to promote instruction-following outputs across varied inputs. Studying the SDI model under complex scenarios—such as multiple or rough masks, unrelated prompt instruction $p$ for reference image $I$—requires a more comprehensive benchmark. However, since the COCO dataset includes only the precise masks, and its prompts for generating are simple and highly related to the image context, we propose FCIBench, which incorporates rough masks, multi-mask, and complex prompts with unrelated contexts of image condition. FCIBench compensates the shortage of existing benchmarks [13, 16], as shown in Appendix.

Building on our observation in the first row of Fig. 3a that the prompt *"zebra"* and ground-truth *"person"* coexist, we explore which modifications to the input mask $M$ can facilitate this coexistence across different scenarios. Intuitively, we hypothesize that **increasing the mask size may provide SDI with more potential solutions to integrate prompt instructions into the image context, rather than simply disregarding the prompt**. The results, illustrated in Fig. 4, reveal that in both scenarios, AS remains nearly constant, indicating that image quality is nearly invariant to mask size. Additionally, LPIPS decrease as the non-masked area became smaller. Finally, as mask size increases, CLIP consistently improves, especially in Fig. 4b where the prompt was unrelated to the background context. This finding supports our hypothesis that **increasing the mask size can enhance prompt-adherence in the model's output**.

## 3.4. The Mechanism Behind Inpainting

In Sec. 3.2, we identify that SDI's deficiency in instruction-following stems from its preference for maintaining har-

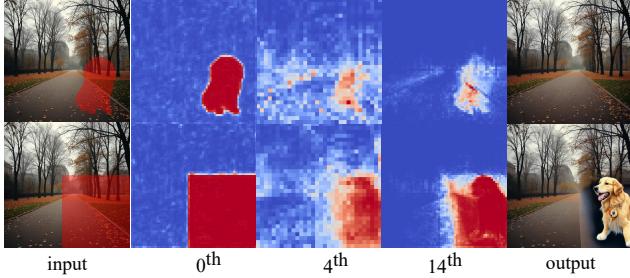"A golden retriever wearing astronaut gear, in cyberpunk style"

input     0th     4th     14th     output

Figure 5. A self-attention visualization in different layers. The attention from $M$ is colored.
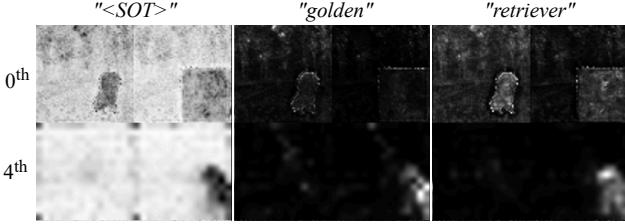


"<SOT>"     "golden"     "retriever"

0th

4th

Figure 6. A cross-attention visualization of Fig. 5 in different cross-attention layers. The attention follows the input mask shape in the first layer, adapting to the output shape in the deeper layer.

mony within the image context. Although in Sec. 3.3 we show that simply enlarging the mask $M$ can better balance instruction-following and contextual harmony, this solution is impractical, as we seek to modify only specific regions within the given mask. This discussion raises two key questions: **(1) How does the image condition influence features within the masked area?** and **(2) How does prompt information selectively affect only the masked area?** By uncovering the underlying mechanisms behind these questions, we can explore ways to refine SDI's learned behavior.

Our initial hypothesis to the first question is that **features within the masked area become progressively diluted by background elements during down-sampling and self-attention operations**. This is illustrated in Fig. 5, where we compared two cases, "precise mask" (row 1) and "large rough mask" (row 2). For both cases, the mask shape is clearly visible in early layers. However, in the subsequent layers (4th, 14th) of the first case, attention within the masked area becomes further diluted by background elements. This results in the final image output with background-like elements in the masked area that are totally unrelated to the prompt. By contrast, in case 2, the attention of the generated object within the masked area successfully deviates from the background elements, aligning closely with the generated object. This results in a more instruction-following outcome, supporting our hypothesis.

To address the second question, we hypothesize that, given the architecture of SDI, **certain channels within the cross-attention key are highly adapted to the mask in-**
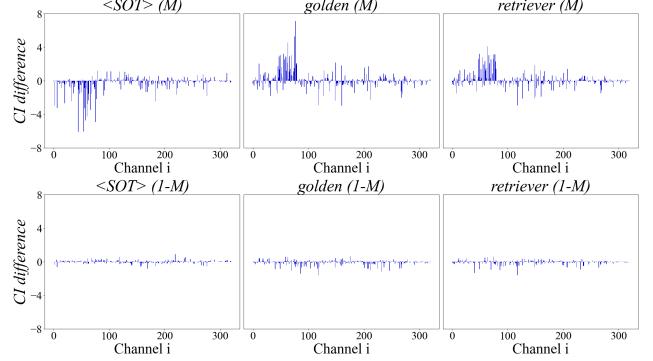


Figure 7. An illustration of the difference of channel influence indicator ($\Delta CI$) across different channels, in the region $M$ (first row) and the region $(1 - M)$ (second row).

**put, selectively enhance prompt response toward the masked region**. This adaptation is illustrated in Fig. 6 and is supported by the classifier-free guidance differences observed in Appendix. To strengthen this hypothesis, we measure the numerical influence made by input $M$. As we mentioned in Sec. 3.1, the mask input $M$ is processed into $M^c$ and $z^c$. In the initial input layer of UNet, these inputs are Concatenated and forwarded into $h_0 = \Phi_0(\text{Concat}([z_t, M^c, z^c])) \in \mathbb{R}^{(H/4) \times (W/4) \times 320}$, where $\Phi_0$ is the first convolutional layer of the denoising UNet. Following this, the cross-attention layer projects $h_0$ feature into the query representation $Q = W_Q h_0$, where $Q \in \mathbb{R}^{(H/4 \times W/4) \times 320}$. Simultaneously, the prompt embedding $p$ is projected into $K = W_K \Psi(p), V = W_V \Psi(p)$, where $\Psi$ is the CLIP text encoder and $K, V \in \mathbb{R}^{77 \times 320}$. The cross-attention output is then computed as $\text{Attention}(Q, K, V) = \text{Softmax}(QK^T/\sqrt{d})V$.

To test our assumption that specific channels of query $Q$ are highly adapted to the given key token $k \in \mathbb{R}^{320}$ in certain feature channels, we define a Channel Influence Indicator ($CI$) here:

$$CI(Q, M, k, i) = \frac{1}{\sum_j \bar{M}_j} \sum_{j=1}^{H \times W/16} (Q_j \odot k)_i \cdot \bar{M}_j, \quad (2)$$

where $\odot$ denotes the Hadamard product, the subscript $i$ refers to the $i$-th element of a vector, and $\bar{M}$ represents the flattened version of $M$. Since the sum of $CI$ across different channels is positively correlated with the $QK^T$ computation in cross-attention, the $CI$ indicator offers a means to visualize the influence introduced by $M$ within specific feature channels.

To measure the difference leading by mask input $M$, we choose the input mask of the second row of Fig. 5, denoted as $M^l$, compared with the zero matrix $M^n$. We then define $Q^l = W_Q \Phi_0(\text{Concat}([z_t, M^l, z^l])$ and $Q^n = W_Q \Phi_0(\text{Concat}([z_t, M^n, z^n])$, where $z^l$ is the image condition given $M^l$, similarly the $z^n$ and $M^n$.
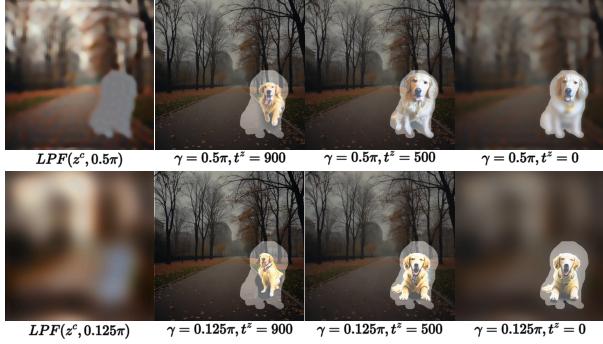
5

Figure 8. A illustration of the $z^{fc}$ (column 1) and the output leading by different values of $t^z$ (columns 2,3,and 4). The input mask is highlighted by overlaying it onto the output images.

To demonstrate that mask input ($M^l$ in this case) significantly influences $QK^T$ computations in certain channels of given masked area $M^l$, we plot $\Delta CI = CI(Q^l, M^l, k, i) - CI(Q^n, M^l, k, i)$, at initial timestep t=T, where both $Q^l$ and $Q^n$ share the same noise latent $z_t$. In Fig. 7, we observe that the $\Delta CI$ shifts more markedly within $M^l$ area than in $(1 - M^l)$ area. For non-informative tokens (*i.e.* "<SOT>"), $\Delta CI$ decreases significantly, leading to a relative increase in attention toward other informative tokens. For meaningful tokens such as "*golden*" and "*retriever*", $\Delta CI$ increases, especially within the first 80 channels. This finding supports our hypothesis that **cross-attention key features adapt specifically to mask input, particularly in the first 80 channels, enabling selective prompt influence within** $M$. More evaluations can be found in Appendix.

## 4. Method

In Sec. 3.2, we identify that the image context provided by $z^c$ can impede instruction-following in the SDI model. In Sec. 3.3 and Sec. 3.4, we observe that the $M^c$ plays an important role in the cross-attention layer, the inclusion of $M^c$ leading to prompt-adherence by shifting the cross-attention features. Building on these insights, we propose FreeCond to directly reduce the heavy reliance on the image context provided by $z^c$ and increase the feature shift led by $M^c$. With FreeCond, we can achieve improved instruction-following for the SDI-based approach in a post-hoc manner without additional fine-tuning or computational costs.

### 4.1. FreeCond Image Condition

In light of the phenomenon observed in Sec. 3.2, where inpainting outputs can be significantly influenced or dominated by the image condition $z^c$, it appears reasonable to reduce the influence of $z^c$ to improve insturction following. However, since the inpainting model relies on $z^c$ to preserve the background, any adjustments to $z^c$ can compromise mask preservation. Nonetheless, based on the nature of the
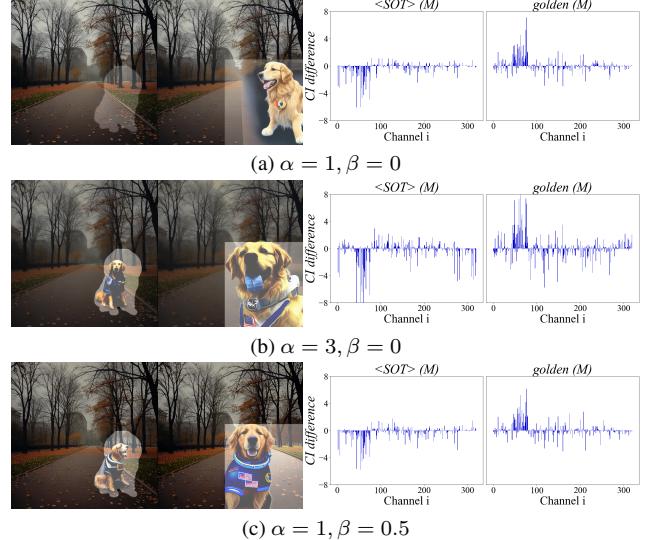


Figure 9. A illustration of the effect of $M^{fc}$ and corresponding CI plot for "large mask" (columns 3 and 4). The input mask is highlighted by overlaying it onto the output images.

T2I diffusion process, as described in [2, 28, 35], we note that low-frequency components are formed in early steps while high-frequency details emerge in later steps. In other words, we can still largely preserve the background in the final output by inputting only the low-frequency portion of $z^c$ in the early step and then transitioning to the original $z^c$. We define the FreeCond image condition as:

$$z^{fc} = \begin{cases} LPF(z^c, \gamma), & \text{if } t \geq t^{fc} \\ z^c, & \text{if } t < t^{fc} \end{cases} \tag{3}$$

where $LPF(z^c, \gamma)$ is a low-pass filter that excludes high-frequency components above the threshold $\gamma$, and $t^{fc}$ is the timestep control parameter, with lower values of $t^{fc}$ producing a blurrier output. The effect of $z^{fc}$ is demonstrated in Fig. 8. By modifying $z^{fc}$ in the early step, such as setting $t^{fc} \in [0.5T, 0.9T]$, we can effectively improve instruction-following with minimal impact on background preservation. Since the $LPF$ filters out high-frequency image information, the overall image context is disrupted, thus enhancing instruction-following by reducing interference from the original image context.

### 4.2. FreeCond Mask Condition

Building on our observations in Sec. 3.4, we find that the T2I inpainting effect of the SDI model is raised by the shifting in cross-attention features with non-zero mask input $M$. Further analysis in Appendix reveals that while both $z^c$ and $M^c$ affect the inpainting outcome, shifts in cross-attention features are primarily driven by $M^c$ values. Based on this, we explore the potential to enhance cross-attention feature shifts by manipulating the mask condition

| Methods | Image Quality | | | Background Preservation | | Instruction Following | |
|---|---|---|---|---|---|---|---|
| | ImageReward↑ | HPS ↑ | Aesthetic↑ | PSNR↑ | LPIPS ↓ | CLIP Score ↑ | IoU Score↑ |
| SDI [26] | -1.21/1.22/-1.95 | 0.24/0.27/0.17 | **5.90**/6.50/**6.69** | 25.95/27.26/25.54 | **0.04/0.04/0.06** | 18.67/26.41/11.45 | 0.49/0.62/0.07 |
| SDI$^{fc}$ | -1.23/1.14/-1.29 | 0.24/0.27/0.20 | 5.86/**6.54**/6.44 | 24.79/26.73/24.58 | **0.04**/0.05/0.07 | 18.82/26.46/18.27 | 0.65/0.70/0.54 |
| CNI [39] | -1.24/1.15/-1.95 | 0.24/0.27/0.17 | 5.78/6.44/6.63 | **26.23/27.69/25.81** | **0.04/0.04/0.06** | 18.83/26.04/11.38 | 0.57/0.60/0.12 |
| CNI$^{fc}$ | -1.24/1.11/-1.52 | 0.24/0.27/0.19 | 5.76/6.46/6.49 | 25.58/27.06/25.23 | **0.04**/0.05/0.07 | 19.10/25.97/16.01 | 0.67/0.67/0.45 |
| HDP [20] | -1.19/1.18/-1.40 | **0.25**/0.27/0.20 | 5.79/6.46/6.57 | 24.95/27.02/25.23 | **0.04**/0.05/**0.06** | 19.14/26.45/17.08 | 0.60/0.63/0.36 |
| HDP$^{fc}$ | -1.27/1.20/-1.15 | 0.24/**0.28**/0.20 | 5.75/6.52/6.43 | 23.63/26.00/23.48 | 0.05/0.05/0.08 | 19.17/26.41/19.37 | 0.77/0.68/0.67 |
| PP [42] | -1.19/1.21/-1.13 | 0.22/0.27/0.19 | 5.79/6.30/6.40 | 25.63/27.62/25.29 | 0.05/0.05/0.07 | 18.74/27.02/19.05 | 0.59/0.55/0.43 |
| PP$^{fc}$ | **-1.15**/1.20/-1.12 | 0.22/0.27/0.19 | 5.73/6.33/6.39 | 25.41/27.37/24.32 | 0.05/0.05/0.08 | 19.12/**27.05**/19.43 | 0.67/0.59/0.52 |
| BN [13] | -1.24/**1.24**/-1.08 | 0.24/0.27/**0.21** | 5.77/6.53/6.38 | 24.89/26.37/24.35 | 0.05/0.06/0.07 | 19.22/26.50/19.96 | 0.83/**0.75**/0.77 |
| BN$^{fc}$ | -1.31/1.21/**-1.05** | 0.23/0.27/**0.21** | 5.77/6.53/6.43 | 24.21/25.38/23.49 | 0.05/0.06/0.08 | **19.27**/26.50/**20.18** | **0.85/0.75/0.78** |
| SDXL [23] | -1.06/1.32/-1.72 | 0.25/0.29/0.19 | 5.74/6.40/6.55 | 24.61/25.78/25.00 | 0.03/0.03/0.04 | 19.09/26.96/14.16 | 0.53/0.68/0.09 |
| SDXL$^{fc}$ | -0.94/1.30/-0.78 | 0.25/0.29/0.22 | 5.69/6.34/6.56 | 24.15/26.12/24.59 | 0.04/0.04/0.05 | 19.77/27.16/22.36 | 0.60/0.64/0.44 |

Table 2. Quantitative results showing improvements achieved by FreeCond (denoted with $^{fc}$) across three benchmarks: **COCO, Brush-Bench, and FCIBench, separated by "/" respectively**. Note: we discovered that the BrushBench results reported in [13] were calculated with an NSFW detector; as NSFW detection may vary, we disable it here to ensure a more precise evaluation.

$M^c$. Accordingly, we introduce a FreeCond mask condition, $M^{fc}$, which scales up the value of $M^c$ to induce a stronger cross-attention feature shift, thereby improving prompt alignment. Another observation, shown in Fig. 7, is that the CI indicator in the $(1 - M)$ region is subtly impacted by the mask $M$. Thus, increasing mask values within the $(1 - M)$ region can amplify feature shifts within $M$. To test this, we define the FreeCond mask condition:

$$M^{fc} = \alpha \cdot M^c + \beta \cdot (1 - M^c) \qquad (4)$$

where $\alpha$ and $\beta$ are scaling factors to control the influence of $M^c$ and $(1 - M^c)$. The impact of $M^{fc}$ is illustrated in Fig. 9. Compared to the baseline results in Fig. 9a, the output with a $M^{fc}$ exhibits greater prompt-adherence. For instance, in the "precise mask" condition, instead of filling the background element, the *"golden retriever wearing astronaut gear"* appears. Additionally, in the "large mask" setting, the golden retriever now includes the *"astronaut gear"*. We also provide the response of the CI indicator for $M^{fc}$ in the right side of Fig. 9 to show that modifying the latent mask $M^c$ can effectively enhance feature shifts within the cross-attention layer.

With our proposed alteration, the noise prediction function from Eq. (1) can be generalized as $\hat{\epsilon}_\theta(z_t, z^{fc}, M^{fc}, t, p)$. As FreeCond only modifies the input, it is compatible with other SDI-based models, detailed in Appendix.

# 5. Experiment

## 5.1. Experiment Setting

We conduct experiments on three datasets: COCO [16], BrushBench [13], and our proposed FCIBench, each comprising 600 instruction pairs (details in Appendix). To account for the inherent differences between the baseline methods, we provide the specific hyperparameter settings and further discussion in Appendix. The computational metrics used for evaluation are detailed in Sec. 3.1.

## 5.2. Experiment Results

Tab. 2 presents the quantitative improvements achieved with the inclusion of FreeCond. We compare FreeCond with the original SDI [26] and its variants, including ControlNet Inpainting (CNI) [39], HD-Painter (HDP) [20], PowerPaint (PP) [42], and BrushNet (BN) [13]. Additionally, we assess SDXL [23], a much larger model, as a reference to showcase the zero-shot improvements by FreeCond, it is not directly compared with other baselines.

In our proposed FCIBench, shown in the right section of each block. Designed as a more challenging benchmark, FreeCond demonstrates substantial gains, achieving a 60% increase over the original SDI model and a 1% improvement over BrushNet, the existing SOTA. Additionally, FreeCond improves metrics such as IR, HPS, and IoU across all baselines. Notably, both BrushNet and PowerPaint benefit from FreeCond with modest increases in IoU yet larger gains in CLIP score, highlighting FreeCond's capability to further enhance prompt-adherence in SOTA methods that are optimized for mask-fitting. For the two widely used datasets, COCO and BrushBench—represented in the left and middle sections of each block—current inpainting methods demonstrate similar performance levels. FreeCond advances this upper limit, increasing BrushNet's CLIP score from 19.22 to 19.27 on COCO and PowerPaint's CLIP score from 27.02 to 27.05 on BrushBench. Overall, FreeCond enhances performance across both instruction following and image quality, with improvements in IR, HPS, and CLIP scores.

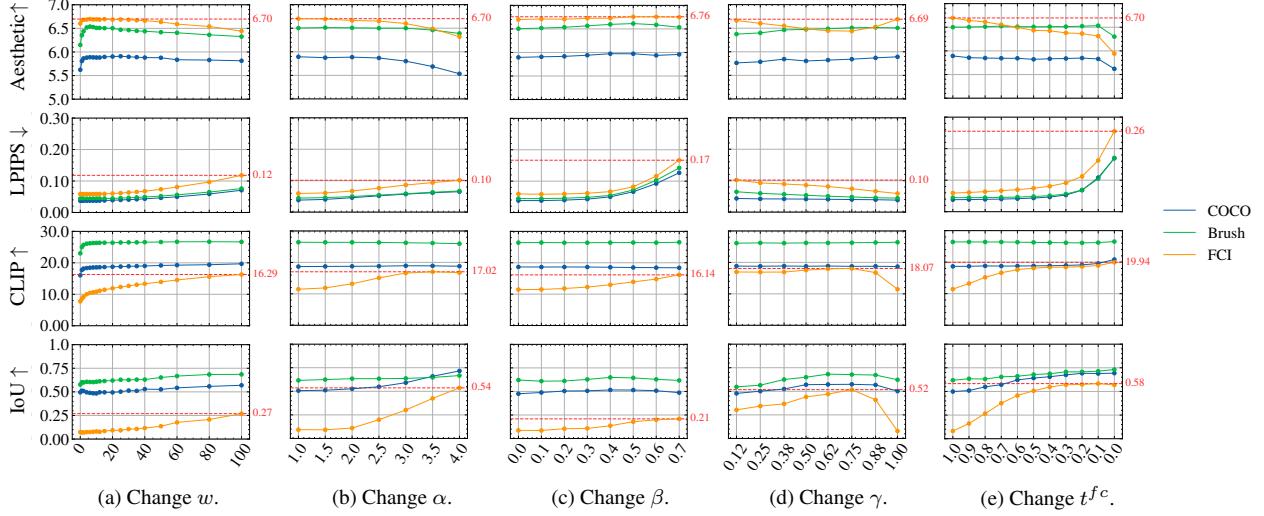Nevertheless, modifying learned mechanism with

Figure 10. Illustration of the influence of CFG ($w$) and each hyperparameter of FreeCond($\alpha, \beta, \gamma, t^{fc}$), highest values are denoted.



Figure 11. The qualitative illustration of Fig. 10, the change compared to normal SDI is colored.

FreeCond can lead to minor degradations in detail-oriented metrics, such as AS, PSNR, and LPIPS. These minor distortions, are generally imperceptible to human preference, as evidenced by the increases in the human-preference-based metrics IR and HPS. These results are more apparent in the qualitative results, discussed further in Appendix.

### 5.3. Ablation Study

In Fig. 10 and Fig. 11, we examine the impact of adjusting five components: (a) the classifier-free guidance (CFG) scale $w$ [9], (b) the inner-mask scale $\alpha$, (c) the outer-mask scale $\beta$, (d) the LPF threshold $\gamma$ with a fixed $t^{fc} = 25$, and (e) the LPF timestep $t^{fc}$ with $\gamma = 0.75\pi$. For each test, we fix the parameters at $(w, \alpha, \beta, \gamma, t^{fc}) = (15, 1, 0, \pi, T)$ (the default configuration of original SDI) and vary only one parameter at a time. Based on quantitative and qualitative outcomes, we summarize our findings below.

**Effect of $w$.** As discussed in Sec. 3.2, SDI's random masking strategy prioritizes generating objects within the mask rather than strict mask conformity. Therefore, increasing $w$ in Fig. 11a primarily enhances prompt-related details without substantially increasing object size. This outcome is further reflected in the **lesser improvement of the IoU score compared to the CLIP score** in Fig. 10a.

**Effect of $\alpha$.** As explained in Sec. 4, increasing $\alpha$ intensifies the cross-attention response within the masked area, **enhancing both prompt-adherence and mask-fitting**, as illustrated in Fig. 11b and Fig. 9b. However, excessively high $\alpha$ disrupts the learned feature distribution, leading to over-saturated results and a drop in AS.

**Effect of $\beta$.** Unlike other parameters, increasing $\beta$ **enhances both the CLIP score and AS**, indicating a stronger self-attention interaction between $M$ and $1 - M$, which results in a more harmonious output. However, as shown in Fig. 11c and Fig. 9c, higher $\beta$ values also increase background distortion, reflected by the LPIPS in Fig. 10c.

**Effect of $z^{fc}$ (controlled by $\gamma$ and $t^{fc}$).** These parameters control the frequency components of $z^{fc}$, which play a key role in reducing contextual influence and establishing prompt-related structures at early timesteps. As shown in Fig. 10e, increasing $t^{fc}$ significantly **improves both CLIP and IoU scores**. This is reflected in Fig. 11d, where mask-fitting is improved while prompt-adherence is lacking (*e.g.*, the *"moss and flowers"* are not fully generated).

## 6. Conclusion

In this work, we identify an instruction-following deficiency that persists across current SDI-based inpainting methods, particularly when complex prompts are provided alongside unrelated image conditions. Through an in-depth investigation of the SDI mechanism, we discover that its selective inpainting capability within masked areas stems from a feature shift in the cross-attention layer. Based on this insight, we propose FreeCond—a training-free plug-in that introduces no additional computation overhead. Unlike classifier-free guidance, FreeCond enhances not only prompt-adherence but also mask-fitting and image quality. However, excessive parameter adjustments can degrade image quality, highlighting the need for careful tuning.

## References

[1] Omri Avrahami, Ohad Fried, and Dani Lischinski. Blended latent diffusion. *ACM Transactions on Graphics*, 42(4):1–11, 2023. 3

[2] Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Qinsheng Zhang, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, et al. ediff-i: Text-to-image diffusion models with an ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*, 2022. 6

[3] Mingdeng Cao, Xintao Wang, Zhongang Qi, Ying Shan, Xiaohu Qie, and Yinqiang Zheng. Masactrl: Tuning-free mutual self-attention control for consistent image synthesis and editing. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22560–22570, 2023. 3

[4] 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*, pages 11315–11325, 2022. 3

[5] Chieh-Yun Chen, Li-Wu Tsao, Chiang Tseng, and Hong-Han Shuai. A cat is a cat (not a dog!): Unraveling information mix-ups in text-to-image encoders through causal analysis and embedding optimization. *arXiv preprint arXiv:2410.00321*, 2024. 3

[6] Jiwoo Chung, Sangeek Hyun, and Jae-Pil Heo. Style injection in diffusion: A training-free approach for adapting large-scale diffusion models for style transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8795–8805, 2024.

[7] Yingying Deng, Xiangyu He, Fan Tang, and Weiming Dong. Z*: Zero-shot style transfer via attention reweighting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6934–6944, 2024.

[8] 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. 3

[9] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 2, 3, 8

[10] Teng-Fang Hsiao, Bo-Kai Ruan, and Hong-Han Shuai. Training-and-prompt-free general painterly harmonization using image-wise attention sharing. *arXiv preprint arXiv:2404.12900*, 2024. 3

[11] Inbar Huberman-Spiegelglas, Vladimir Kulikov, and Tomer Michaeli. An edit friendly ddpm noise space: Inversion and manipulations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12469–12478, 2024. 3

[12] Changho Jo, Woobin Im, and Sung-Eui Yoon. In-n-out: Towards good initialization for inpainting and outpainting. *arXiv preprint arXiv:2106.13953*, 2021. 3

[13] Xuan Ju, Xian Liu, Xintao Wang, Yuxuan Bian, Ying Shan, and Qiang Xu. Brushnet: A plug-and-play image inpainting model with decomposed dual-branch diffusion. *arXiv preprint arXiv:2403.06976*, 2024. 2, 3, 4, 7

[14] Diederik P Kingma. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 3

[15] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023. 3

[16] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, pages 740–755. Springer, 2014. 2, 4, 7

[17] Nan Liu, Shuang Li, Yilun Du, Antonio Torralba, and Joshua B Tenenbaum. Compositional visual generation with composable diffusion models. In *European Conference on Computer Vision*, pages 423–439. Springer, 2022. 3

[18] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11461–11471, 2022. 3

[19] Yang Luo, Yiheng Zhang, Zhaofan Qiu, Ting Yao, Zhineng Chen, Yu-Gang Jiang, and Tao Mei. Freeenhance: Tuning-free image enhancement via content-consistent noising-and-denoising process. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 7075–7084, 2024. 3

[20] Hayk Manukyan, Andranik Sargsyan, Barsegh Atanyan, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. Hd-painter: high-resolution and prompt-faithful text-guided image inpainting with diffusion models. *arXiv preprint arXiv:2312.14091*, 2023. 2, 3, 7

[21] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073*, 2021. 3

[22] Jialun Peng, Dong Liu, Songcen Xu, and Houqiang Li. Generating diverse structure for image inpainting with hierarchical vq-vae. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10775–10784, 2021. 3

[23] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023. 2, 7

[24] 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*, pages 8748–8763. PMLR, 2021. 2, 3

[25] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. Pmlr, 2021. 3

[26] 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*, pages 10684–10695, 2022. 2, 3, 7

[27] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems*, 35:25278–25294, 2022. 3

[28] Chenyang Si, Ziqi Huang, Yuming Jiang, and Ziwei Liu. Freeu: Free lunch in diffusion u-net. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4733–4743, 2024. 3, 6

[29] Su Wang, Chitwan Saharia, Ceslee Montgomery, Jordi Pont-Tuset, Shai Noy, Stefano Pellegrini, Yasumasa Onoe, Sarah Laszlo, David J Fleet, Radu Soricut, et al. Imagen editor and editbench: Advancing and evaluating text-guided image inpainting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18359–18369, 2023. 3

[30] Xiaoshi Wu, Keqiang Sun, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score: Better aligning text-to-image models with human preference. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2096–2105, 2023. 3

[31] Shaoan Xie, Zhifei Zhang, Zhe Lin, Tobias Hinz, and Kun Zhang. Smartbrush: Text and shape guided object inpainting with diffusion model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22428–22437, 2023. 3

[32] Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation. *Advances in Neural Information Processing Systems*, 36, 2024. 3

[33] Shiyuan Yang, Xiaodong Chen, and Jing Liao. Uni-paint: A unified framework for multimodal image inpainting with pretrained diffusion model. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 3190–3199, 2023. 3

[34] Mingyang Yi, Aoxue Li, Yi Xin, and Zhenguo Li. Towards understanding the working mechanism of text-to-image diffusion model. *arXiv preprint arXiv:2405.15330*, 2024. 3

[35] Mingyang Yi, Aoxue Li, Yi Xin, and Zhenguo Li. Towards understanding the working mechanism of text-to-image diffusion model. *arXiv preprint arXiv:2405.15330*, 2024. 3, 6

[36] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S Huang. Generative image inpainting with contextual attention. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5505–5514, 2018. 3

[37] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S Huang. Free-form image inpainting with gated convolution. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4471–4480, 2019.

[38] Yanhong Zeng, Jianlong Fu, Hongyang Chao, and Baining Guo. Aggregated contextual transformations for high-resolution image inpainting. *IEEE Transactions on Visualization and Computer Graphics*, 29(7):3266–3280, 2022. 3

[39] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847, 2023. 2, 3, 7

[40] 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/CVF Conference on Computer Vision and Pattern Recognition*, pages 586–595, 2018. 3

[41] Chuanxia Zheng, Tat-Jen Cham, and Jianfei Cai. Pluralistic image completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1438–1447, 2019. 3

[42] Junhao Zhuang, Yanhong Zeng, Wenran Liu, Chun Yuan, and Kai Chen. A task is worth one word: Learning with task prompts for high-quality versatile image inpainting. *arXiv preprint arXiv:2312.03594*, 2023. 2, 3, 7