000 MINORITYPROMPT: TEXT TO MINORITY IMAGE 001 GENERATION VIA PROMPT OPTIMIZATION 002 003

Anonymous authors

Paper under double-blind review

ABSTRACT

We investigate the generation of minority samples using pretrained text-to-image (T2I) latent diffusion models. Minority instances, in the context of T2I generation, can be defined as ones living on low-density regions of *text-conditional* data distributions. They are valuable for various applications of modern T2I generators, such as data augmentation and creative AI. Unfortunately, existing pretrained T2I diffusion models primarily focus on high-density regions, largely due to the influence of guided samplers (like CFG) that are essential for producing high-quality generations. To address this, we present a novel framework to counter the high-densityfocus of T2I diffusion models. Specifically, we first develop an online prompt optimization framework that can encourage the emergence of desired properties during inference while preserving semantic contents of user-provided prompts. We subsequently tailor this generic prompt optimizer into a specialized solver that promotes the generation of minority features by incorporating a carefully-crafted likelihood objective. Our comprehensive experiments, conducted across various types of T2I models, demonstrate that our approach significantly enhances the capability to produce high-quality minority instances compared to existing samplers.



Figure 1: Example results from our minority generation approach using SDXL-Lightning. Our 048 framework is designed to produce unique *minority* samples w.r.t. user-provided prompts, which are 049 rarely generated by standard samplers like DDIM (Song et al., 2020a). Due to its low-likelihood encouraging nature, our sampler often demonstrates counteracting results against demographic bi-051 ases in text-to-image models (Friedrich et al., 2023). See the samples in the last row for instance, 052 where our sampler mitigates prevalent age and racial biases (e.g., associating "man" with "young" and "woman" with "white") by modifying the demographic traits of the subjects.

1

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

055 056

Text-to-image (T2I) generative models (Xu et al., 2018; Ramesh et al., 2021; Nichol et al., 2021) have recently attracted substantial interest for their capability to convert textual descriptions into visually striking images. At the forefront of the surge are diffusion models (Song & Ermon, 2019; Ho et al., 2020), augmented by guidance techniques (Dhariwal & Nichol, 2021; Ho & Salimans, 2022) such as classifier-free guidance (CFG) (Ho & Salimans, 2022). The guided T2I samplers encourage generations from high-density regions of a data manifold (Dhariwal & Nichol, 2021), producing realistic images that faithfully respect the provided prompts.

A key challenge is that the inherent high density focus of modern T2I samplers makes it difficult to generate *minority* samples – instances that reside in low-density regions of the manifold. This limitation is particularly significant as T2I-generated data is increasingly incorporated in downstream applications (Tian et al., 2024a;b; Afkanpour et al., 2024) where the majority-focused bias within the data may be perpetuated. Furthermore, the unique attributes found in minority instances are crucial for applications like creative AI (Rombach et al., 2022; Han et al., 2022), where generating novel and highly creative outputs is essential.

070 In this work, we present a novel approach dubbed as *MinorityPrompt* that counteracts the high-071 density bias of T2I samplers to improve their capability of minority generation. Our framework is 072 built upon the concept of prompt optimization, an intuitive technique that exhibits strong perfor-073 mance in enhancing T2I diffusion models for various tasks (Gal et al., 2022; Chung et al., 2023b; 074 Park et al., 2024). Unlike existing T2I-based online prompt-tuning methods that modify the entire input prompts (e.g., by optimizing their text-embeddings during inference), our approach updates 075 the prompts in a *selective* fashion to preserve the intended semantics while encouraging generations 076 of unique low-density features. Specifically during inference, we incorporate learnable tokens into 077 the input prompts, e.g., by appending them to the end of the text. We then adjust the embeddings of these tokens across sampling timesteps, targetting the minimization of a likelihood metric de-079 signed to capture the uniqueness of noisy intermediate samples. See Figure 2 for an overview. An additional benefit of our token-based approach is that it offers enhanced semantic controllability, 081 enabling users to express specific desired semantics in generated samples by selecting appropriate 082 initialization words for the learnable token embeddings. To further improve the performance of our 083 sampler, we provide new design choices that can be synergistically employed with our approach 084 for T2I minority generation. Comprehensive experiments validate that our method can significantly 085 improve the ability of creating minority instances of modern widely-adopted T2I models (including Stable Diffusion (SD) (Rombach et al., 2022)) with minimal compromise in sample quality and text-086 image alignment. In addition, we emphasize that our framework can work on distilled backbones 087 like SDXL-Lightning (Lin et al., 2024), which demonstrates its robustness and practical relevance. 088 As an additional application, we explore the potential of our prompt optimization framework to 089 improve the diversity of T2I models, further exhibiting its versatility as a general-purpose solver 090 applicable across various tasks. 091

Given that our prompt optimization is performed in an online manner, does not require expensive fine-tuning of T2I models, and is entirely *self-contained*, *i.e.*, implementable solely with a pretrained T2I model, we believe our approach open a new avenue for creative AI, emphasizing the practical relevance of our framework.

096 097

2 RELATED WORK

098

099 The generation of minority samples has been explored in a range of different scenarios and gener-100 ative frameworks (Yu et al., 2020; Lin et al., 2022; Sehwag et al., 2022; Qin et al., 2023; Huang & 101 Jafari, 2023; Um & Ye, 2023; 2024). However, significant progress has been recently made with the 102 introduction of diffusion models, due to their ability to faithfully capture data distributions (Sehwag 103 et al., 2022; Um & Ye, 2023; 2024). As an initial effort, Sehwag et al. (2022) incorporate separately-104 trained classifiers into the sampling process of diffusion models to yield guidance for low-density 105 regions. The approach by Um & Ye (2023) shares similar intuition of integrating an additional classifier into the reverse process for low-density guidance. A limitation is that their methods rely upon 106 external classifiers that are often difficult to obtain, especially for large-scale datasets such as T2I 107 benchmarks (Schuhmann et al., 2022). The challenge was recently addressed by Um & Ye (2024)



Figure 2: **Overview of MinorityPrompt.** Unlike existing online prompt tuning approaches that adjust the entire text-embedding (*e.g.*, the output of the text-encoder) during inference, our framework focuses on optimizing a dedicated *token-embedding* to better preserve the semantics within the prompt. Specifically given a user-prompt (*e.g.*, "A portrait of a dog"), we integrate a placeholder string (*e.g.*, S in the figure) into the prompt, marking the position of the learnable token embedding v. With the text-embedding C_v that incorporates the contents of v, we update v on-the-fly during the inference process to maximize the reconstruction loss of the denoised version of z_t (i.e., \hat{z}_0^1 in the figure). The optimized token v^* is subsequently used to progress the inference at the corresponding timestep; see Section 3 for details.

123

124

125

126

127

128

129

- 132
- 133

134

where the authors develop a self-contained minority sampler that works without expensive extra
components (such as classifiers). However, their method is tailored for canonical image benchmarks
(like LSUN (Yu et al., 2015) and ImageNet (Deng et al., 2009)) and exhibits limited performance
gain in more challenging scenarios like T2I generation. Moreover, none of these approaches have
explored the dimension of prompt optimization specifically for minority generation, which is the
central focus of our framework.

141 A related yet distinct objective is enhancing the diversity of diffusion models, an area that has been 142 relatively overlooked compared to improving their quality. Significant progress was recently made in Sadat et al. (2023), where the authors demonstrated that adding noise perturbations, if gradually 143 annealed over time, to conditional embeddings could greatly enhance the diversity of generated sam-144 ples. However, unlike our approach, their method focuses on producing diverse samples that remain 145 consistent with the ground-truth data distribution, rather than targeting the low-density regions of 146 the distribution. Another notable contribution was done by Corso et al. (2023). Their idea is to repel 147 intermediate latent samples that share the same condition, thereby encouraging the final generated 148 samples to exhibit distinct features. A disadvantage is that it requires generating multiple instances 149 for each prompt, which can be redundant in many practical scenarios. 150

Prompt optimization has been widely explored in the context of T2I diffusion models due to their 151 strong dependence on language models. This approach has exhibited significant performance across 152 various tasks, including inverse problems (Chung et al., 2023b) and image editing (Park et al., 2024; 153 Mokady et al., 2023). A key difference is that most existing methods in these lines tune the en-154 tire prompts to find the ones that best perform the focused tasks (e.g., minimizing data consistency 155 loss (Chung et al., 2023b)). In contrast, our framework updates only the attached learnable tokens, 156 thereby preserving the original prompt's semantics while encouraging the emergence of low-density 157 features. Additional use cases of prompt tuning include personalization (Gal et al., 2022; Han et al., 158 2023) and object counting (Zafar et al., 2024). Similar to ours, their frameworks introduce vari-159 able tokens and tune their embeddings. However, their optimizations aim to learn visual concepts captured in user-provided images, whereas our focus is to invoke low-density features through opti-160 mized prompts. Also, their methods are not online, requiring separate training procedure which can 161 be potentially expensive.

162 3 METHOD

163 164

Our focus is to generate high-quality minority instances using text-to-image (T2I) diffusion models, which faithfully reflect user-provided prompts while featuring unique visual attributes rarely produced via standard generation techniques¹. To this end, we start with providing a brief overview on T2I diffusion frameworks and the essential background necessary to understand the core of our work. We subsequently present our proposed framework for minority generation based on the idea of prompt optimization.

170 171

179

3.1 BACKGROUND AND PRELIMINARIES

The task of T2I diffusion models is to generate an output image $x_0 \in \mathbb{R}^d$ from a random noise vector $z_T \in \mathbb{R}^k$ (where typically k < d), given a user-defined text prompt \mathcal{P} . Similar to standard (non-T2I) diffusion frameworks, the core of T2I diffusion sampling lies in an iterative denoising process that progressively removes noise from z_T until a clean version z_0 is obtained. This denoising capability is learned through noise-prediction training (Ho et al., 2020; Song & Ermon, 2019), mathematically written as:

$$\max_{\boldsymbol{\rho}} \mathbb{E}_{\boldsymbol{z}_{0}, y, \boldsymbol{\epsilon} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I}), t \sim \text{Unif}\{1, \dots, T\}} [\|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\boldsymbol{\theta}}(\boldsymbol{z}_{t}, \mathcal{C})\|_{2}^{2}],$$

where $z_0 := \mathcal{E}(x_0)$, yielded by passing a training image x_0 through a compressive model \mathcal{E} (*e.g.*, the encoder of VQ-VAE (Esser et al., 2021; Rombach et al., 2022)). Here, z_t represents a noiseperturbed version of z_0 , given by $z_t := \sqrt{\alpha_t} z_0 + \sqrt{1 - \alpha_t} \epsilon$, where $\{\alpha_t\}_{t=1}^T$ defines the noiseschedule. ϵ_{θ} refers to a T2I diffusion model parameterized to predict the noise ϵ , and \mathcal{C} represents the embedding of the text prompt \mathcal{P} . See below for details on how to obtain \mathcal{C} from \mathcal{P} .

Once trained, T2I generation can be done by starting from $z_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and implementing an iterative noise removal process guided by the text embedding C. A common approach is to follow the deterministic DDIM sampling (Song et al., 2020a; Chung et al., 2023a):

$$\boldsymbol{z}_{t-1} = \sqrt{\alpha_{t-1}} \hat{\boldsymbol{z}}_0(\boldsymbol{z}_t, \mathcal{C}) + \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon}_{\boldsymbol{\theta}}(\boldsymbol{z}_t, \mathcal{C})$$

where $\hat{\boldsymbol{z}}_0(\boldsymbol{z}_t, \mathcal{C}) \coloneqq \frac{1}{\sqrt{\alpha_t}} \left(\boldsymbol{z}_t - \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{\boldsymbol{\theta}}(\boldsymbol{z}_t, \mathcal{C}) \right).$ (1)

190 191

189

198

Here $\hat{z}_0(z_t, C)$ indicates a denoised estimate of z_t conditioned on the text embedding C, implemented via Tweedie's formula (Chung et al., 2022).

To further strengthen the impact of text conditioning, classifier-free guidance (CFG) (Ho & Salimans, 2022) is commonly integrated into the sampling process. In particular, one can obtain a high-density-focused noise estimation through extrapolation using an unconditional prediction:

 $\tilde{\boldsymbol{\epsilon}}_{\boldsymbol{\theta}}^{w}(\boldsymbol{z}_{t}, \mathcal{C}) \coloneqq w \boldsymbol{\epsilon}_{\boldsymbol{\theta}}(\boldsymbol{z}_{t}, \mathcal{C}) + (1 - w) \boldsymbol{\epsilon}_{\boldsymbol{\theta}}(\boldsymbol{z}_{t}),$

where $\epsilon_{\theta}(z_t)$ indicates an unconditional noise prediction, often implemented via null-text conditioning (Ho & Salimans, 2022). CFG refers to the technique that employs $\tilde{\epsilon}_{\theta}^{w}(z_t, C)$ in place of $\epsilon_{\theta}(z_t, C)$ (in Eq. (1)), which has been shown in various scenarios to significantly improve both sample quality and text alignment yet at the expense of diversity (Sadat et al., 2023).

Text processing. A key distinction from non-T2I diffusion models is the incorporation of the text embedding C, a continuous vector yielded by a text encoder T (such as BERT (Devlin, 2018)) based on the user prompt P. To obtain this embedding, each word (or sub-word) in P is first converted into a token – an index in a pre-defined vocabulary. Each token is then mapped to a unique embedding vector through an index-based lookup. These token-wise embedding vectors, often referred to as *token* embeddings, are typically learned as part of the text encoder. The token embeddings are then passed through a transformer model, yielding the final text embedding C. For simplicity, we denote this text processing operation as the forward pass of the text encoder T; thus, C = T(P).

Prompt optimization. In the context of T2I diffusion models, prompt tuning is performed by intervening in the text-processing stage. A common approach is to adjust the text embedding C

¹More formally, this can be expressed as drawing instances from $S_c := \{z \in \mathcal{M}_c : p_{\theta}(z|\mathcal{C}) < \epsilon\}$, where ²¹⁵ \mathcal{C} is the prompt, \mathcal{M}_c represents the (latent) data manifold associated with \mathcal{C} , and p_{θ} denotes the probability density captured by the T2I diffusion model. Here ϵ is a small positive constant.

216 over inference time, which is widely adopted in existing online prompt optimizers (Chung et al., 217 2023b; Park et al., 2024). Specifically at sampling timestep t, existing online prompt tuners can be 218 formulated as the following optimization problem: 219

$$\mathcal{C}_t^* \coloneqq \arg\max_{\mathcal{C}} \mathcal{J}(\boldsymbol{z}_t, \mathcal{C}), \tag{2}$$

where z_t is a noisy latent at step t and \mathcal{J} represents a task-specific objective function, such as data 222 consistency in inverse problems (Chung et al., 2023b). Once C_t^* is obtained, it is used as a drop-223 in replacement for C at time t (e.g., in Eq. (1)), encouraging the desired property to manifest in 224 subsequent timesteps.

225 A problem is that the optimization in Eq. (2) may lead to a loss of user-intended semantics in \mathcal{P} , 226 due to the comprehensive updating of the entire text-embedding \mathcal{C} . This is critical, especially in 227 the context of our focused T2I minority generation where preserving prompt semantics is essen-228 tial; see Table 2a for our empirical results that support this. One can resort to tuning the null-text 229 embedding while keeping C intact (as suggested by Mokady et al. (2023)). However, this method 230 requires reserving the null-text dimension for this specific purpose, limiting its potential use for im-231 proving sample quality or serving other functions. In the following sections, we present an online 232 prompt optimization framework designed to better preserve semantics. Building on this foundation, we develop our T2I minority sampler, which promotes the generation of minority features while 233 maintaining both sample quality and text-alignment performance. 234

235

246

247

220

221

236 3.2 SEMANTIC-PRESERVING PROMPT OPTIMIZATION

237 The key idea of our optimization approach is to incorporate learnable tokens into a given prompt \mathcal{P} 238 and update its embedding *on-the-fly* during inference. Specifically, we append a placeholder string² 239 S to the prompt \mathcal{P} , which acts as a mark for the learnable tokens. For instance, the augmented 240 prompt could be $\mathcal{P}_{\mathcal{S}} :=$ "A portrait of a dog \mathcal{S} ". This additional string is treated as a new vocabulary 241 item for the text-encoder \mathcal{T} . We assign a token embedding v to \mathcal{S} , and denote the text encoder 242 incorporating it as $\mathcal{T}(\cdot; \boldsymbol{v})$.

243 We propose optimizing this embedding v rather than \mathcal{C} . The proposed online prompt optimization 244 at sampling step t can then be formalized as follows: 245

$$\boldsymbol{v}_t^* \coloneqq \arg \max_{\boldsymbol{v}} \mathcal{J}(\boldsymbol{z}_t, \mathcal{C}_{\boldsymbol{v}}), \tag{3}$$

where $C_{v} := \mathcal{T}(\mathcal{P}_{\mathcal{S}}; v)$. Afterward, the optimized text-embedding $C_{v_{t}^{*}}$ is obtained by text-processing 248 $\mathcal{P}_{\mathcal{S}}$ with the updated token-embedding of \mathcal{S} , therefore $\mathcal{C}_{\boldsymbol{v}_{t}^{*}} := \mathcal{T}(\mathcal{P}_{\mathcal{S}}; \boldsymbol{v}_{t}^{*})$. 249

250 Note that our optimization does not affect the embeddings of the tokens w.r.t. the original prompt 251 \mathcal{P} . This is inherently more advantageous for preserving semantics compared to existing methods, 252 which alter the entire text-embedding C and thereby effectively impact all token embeddings. We 253 also highlight that unlike existing learnable-token-based approaches that share the same embedding throughout inference (Gal et al., 2022; Han et al., 2023; Zafar et al., 2024), our framework allows 254 the token embedding v to change over timesteps t. This adaptive feature offers potential advan-255 tages, since the role of v in maximizing \mathcal{J} can vary with the changing nature of z_t across different 256 timesteps. This point is also implied in previous works that employ adaptive text-embeddings over 257 time (Chung et al., 2023b; Park et al., 2024). 258

Intuitively, our optimization can be understood as capturing a specific concept relevant to noisy 259 latent z_t within the token v_t^* , guided by the objective function \mathcal{J} . Thanks to its general design that 260 accommodates any arbitrary objective function \mathcal{J} , this framework is versatile and can be employed 261 in various contexts beyond minority generation. For instance, it can be used to diversify the outputs 262 of T2I models. See details in Table 3b. 263

264 265

269

3.3 MINORITY PROMPT: MINORITY-FOCUSED PROMPT TUNING

266 We now specialize the generic solver in Eq. (3) for the task of minority generation. The key ques-267 tion is how to formulate an appropriate objective function \mathcal{J} for this purpose. To address this, we 268

²The placeholder string can be placed at any position in the prompt, but we empirically found that inserting it at the end of the prompt yields the best performance; see Table 6b for details.

Al	gorithm 1 MinorityPrompt	Algorithm 2 Prompt optimization				
Re	quire: $\epsilon_{\theta}, \mathcal{T}, \mathcal{D}, \boldsymbol{v}_{T}^{(0)}, \mathcal{P}_{\mathcal{S}}, \mathcal{C}, N, K, w, T, s, \lambda.$	1:	function OptimizeEmb $(\boldsymbol{z}_t, \boldsymbol{v}_t^{(0)}, \boldsymbol{\epsilon}_{\boldsymbol{\theta}}, \mathcal{T}, K, s, \lambda)$			
1:	$oldsymbol{z}_T \sim \mathcal{N}(oldsymbol{0},oldsymbol{I})$	2:	for $k \leftarrow 1$ to K do			
2:	for $t \leftarrow T$ to 1 do	3:	$\mathcal{C}_{oldsymbol{v}} \leftarrow \mathcal{T}(\mathcal{P}_{\mathcal{S}}; oldsymbol{v}_t^{(k-1)})$			
3:	${\mathcal C}_{oldsymbol{v}_{*}^{*}} \leftarrow {\mathcal C}$	4:	$oldsymbol{\epsilon}_{oldsymbol{ heta}}^1 \leftarrow oldsymbol{\epsilon}_{oldsymbol{ heta}}(oldsymbol{z}_t, oldsymbol{\mathcal{C}}_{oldsymbol{v}})$			
4:	if $t \mod N = 0$ then	5:	$\hat{oldsymbol{z}}_0^1 \leftarrow (oldsymbol{z}_t - \sqrt{1 - lpha_t}oldsymbol{\epsilon}_{oldsymbol{ heta}}^1)/\sqrt{lpha_t}$			
5:	$\boldsymbol{v}_t^* \leftarrow \text{OptimizeEmb}(\boldsymbol{z}_t, \boldsymbol{v}_t^{(0)}, \boldsymbol{\epsilon}_{\boldsymbol{\theta}}, \mathcal{T}, K, s, \lambda)$	6:	$oldsymbol{\epsilon} \sim \mathcal{N}(oldsymbol{0}, oldsymbol{I})$			
6:	$\mathcal{C}_{oldsymbol{v}^*_{oldsymbol{ au}}} \leftarrow \mathcal{T}(\mathcal{P}_{\mathcal{S}};oldsymbol{v}^*_t)$	7:	$oldsymbol{z}_{s t,0} \leftarrow \sqrt{lpha_s} \hat{oldsymbol{z}}_0^1 + \sqrt{1-lpha_s} oldsymbol{\epsilon}$			
7:	end if	8:	$\boldsymbol{\epsilon_{\theta}^2} \leftarrow \boldsymbol{\epsilon_{\theta}}(\boldsymbol{z}_{s t,0}, \mathcal{C})$			
8:	$\tilde{\boldsymbol{\epsilon}}_{\boldsymbol{\theta}}^{w} \leftarrow w \boldsymbol{\epsilon}_{\boldsymbol{\theta}}(\boldsymbol{z}_{t}, \mathcal{C}_{\boldsymbol{v}_{t}^{*}}) + (1-w) \boldsymbol{\epsilon}_{\boldsymbol{\theta}}(\boldsymbol{z}_{t})$	9:	$\hat{\boldsymbol{z}}_{0}^{2} \leftarrow (\boldsymbol{z}_{s t,0} - \sqrt{1 - \alpha_{s}} \boldsymbol{\epsilon}_{\boldsymbol{\theta}}^{2}) / \sqrt{\alpha_{s}}$			
9:	$\hat{oldsymbol{z}}_0^w \leftarrow (oldsymbol{z}_t - \sqrt{1 - lpha_t} \widetilde{\epsilon}_{oldsymbol{ heta}}^w) / \sqrt{lpha_t}$	10.	$\mathcal{I} \leftarrow \ \hat{z}_1^1 - \operatorname{sq}(\hat{z}_2^2)\ _2^2 + \lambda\ \operatorname{sq}(\hat{z}_1^1) - \hat{z}_2^2\ _2^2$			
10	$oldsymbol{z}_{t-1} \leftarrow \sqrt{lpha_{t-1}} \hat{oldsymbol{z}}_0^w + \sqrt{1-lpha_{t-1}} ilde{eta}_{oldsymbol{ heta}}^w$	11.	$u^{(k)} = u^{(k-1)} + b \dim C = d(\mathcal{J})$			
11	$\boldsymbol{n}^{(0)} \leftarrow \boldsymbol{n}^*$	11:	$v_t \leftarrow v_t + \text{AdamGrad}(J_t)$			
12	end for	12:				
13	return $\mathbf{r}_{0} \leftarrow \mathcal{D}(\mathbf{r}_{0})$	13:	return $v_t \leftarrow v_t$			
15	$\mathcal{L}(\mathbf{z}_0) \leftarrow \mathcal{L}(\mathbf{z}_0)$	14:	ena function			

draw inspiration from Um & Ye (2024), employing their likelihood metric as the starting point for developing our objective function.

Since the metric was originally defined in the pixel domain using non-T2I diffusion models (see Sec-289 tion B.1 for details), we initially perform a naive adaptation to accommodate the latent space of 290 interest, $z_t \in \mathbb{R}^k$, and integrate text conditioning using CFG as is typical in the T2I context (Kim et al., 2023). The adapted version of the metric reads: 292

$$\mathcal{J}(\boldsymbol{z}_t, \mathcal{C}) \coloneqq \mathbb{E}_{\boldsymbol{\epsilon}} \left[\| \hat{\boldsymbol{z}}_0^w(\boldsymbol{z}_t, \mathcal{C}) - \operatorname{sg}(\hat{\boldsymbol{z}}_0^w(\boldsymbol{z}_{s|t, 0}^w, \mathcal{C})) \|_2^2 \right],$$
(4)

where $\hat{z}_0^w(z_t, C)$ represents a clean estimate of z_t using the CFG noise term $\tilde{\epsilon}_{\theta}^w(z_t, C)$. Here $z_{s|t,0}^w$ 295 indicates a noised version of $\hat{z}_0^w(z_t, \mathcal{C})$ w.r.t. timestep $s: z_{s|t,0}^w \coloneqq \sqrt{\alpha_s} \hat{z}_0^w(z_t, \mathcal{C}) + \sqrt{1 - \alpha_s} \epsilon$, and 296 297 $\hat{z}_0^w(z_{s|t,0}^w, C)$ is a clean version of $z_{s|t,0}^w$ conditioned on C. $sg(\cdot)$ denotes the stop-gradient operator 298 for reducing computational cost when used in guided sampling (Um & Ye, 2024). Notice that the 299 squared L2 error is used as the discrepancy loss, rather than the originally used LPIPS (Zhang et al., 300 2018), due to its incompatibility with our latent space. The quantity in Eq. (4) is interpretable as a 301 reconstruction loss of $\hat{z}_0^w(z_t, \mathcal{C})$. As exhibited in Um & Ye (2024), the loss may become large if 302 z_t (represented by $\hat{z}_0^w(z_t, C)$) contains highly-unique minority features that often vanish during the reconstruction process. The comprehensive details regarding the original metric due to Um & Ye 303 (2024) are provided in Section B.1. 304

305 Considering Eq. (4) as the objective function, a natural approach for minority-focused prompt tuning 306 would be to incorporate C_v and optimize for the best v: 307

$$\boldsymbol{v}_{t}^{*} \coloneqq \arg \max_{\boldsymbol{v}} \mathcal{J}(\boldsymbol{z}_{t}, \mathcal{C}_{\boldsymbol{v}})$$

where $\mathcal{J}(\boldsymbol{z}_{t}, \mathcal{C}_{\boldsymbol{v}}) \coloneqq \mathbb{E}_{\boldsymbol{\epsilon}} \left[\| \hat{\boldsymbol{z}}_{0}^{w}(\boldsymbol{z}_{t}, \mathcal{C}_{\boldsymbol{v}}) - \operatorname{sg}(\hat{\boldsymbol{z}}_{0}^{w}(\boldsymbol{z}_{s|t, 0}^{w}, \mathcal{C}_{\boldsymbol{v}})) \|_{2}^{2} \right].$ (5)

However, we argue that this naively extended framework has theoretical issues that lead to limited 311 performance gain over standard samplers. Specifically, three aspects of this objective weaken the 312 desired connection to the target log-likelihood log $p_{\theta}(z_0 \mid C)$ that we aim to capture: (i) the reliance 313 on the CFG-based clean predictions; (ii) obstructed gradient flow through the second term in the 314 squared-L2 loss; and (iii) the incorporation of C_v within the second term in the loss. See Section A.2 315 on a detailed analysis on these points. 316

317 Hence, we propose the following optimization to address the theoretical issues:

$$\boldsymbol{v}_{t}^{*} \coloneqq \arg \max_{\boldsymbol{v}} \mathcal{J}_{\mathcal{C}}(\boldsymbol{z}_{t}, \mathcal{C}_{\boldsymbol{v}})$$

where $\mathcal{J}_{\mathcal{C}}(\boldsymbol{z}_{t}, \mathcal{C}_{\boldsymbol{v}}) \coloneqq \mathbb{E}_{\boldsymbol{\epsilon}} \left[\| \hat{\boldsymbol{z}}_{0}(\boldsymbol{z}_{t}, \mathcal{C}_{\boldsymbol{v}}) - \hat{\boldsymbol{z}}_{0}(\boldsymbol{z}_{s|t,0}, \mathcal{C}) \|_{2}^{2} \right]$ (6)

319 320 321

318

308

310

270

286

287

291

293 294

where $\hat{z}_0(z_t, C_v) \coloneqq (z_t - \sqrt{1 - \alpha_t} \epsilon_{\theta}(z_t, C_v)) / \sqrt{\alpha_t}$, indicating a non-CFG clean estimate. We 322 found that the proposed optimization maintains a close connection to the focused log-likelihood. 323 Below we provide a formal statement of our finding. See Section A.1 for the proof.



Figure 3: Enhanced semantic controllability by MinorityPrompt. The samples in the first column are generations due to DDIM using the two base prompts (*e.g.*, "A portrait of a dog" for the first row). The third and fifth columns exhibit generated samples from our framework, where we selected the corresponding word embeddings as the starting points of the prompt optimizations. For comparison, we also present DDIM samples produced using attached prompts with the corresponding words in the the second and fourth columns. For instance in the first row, the image in the second column corresponds to the generation due to "A portrait of a dog *dirty*". All samples were obtained using SDXL-Lightning (Lin et al., 2024)

Proposition 1. The objective function in Eq. (6) is equivalent (upto a constant factor) to the negative ELBO w.r.t. $\log p_{\theta}(\hat{z}_0(z_t, C_v) \mid C)$ when integrated over timesteps with $\bar{w}_s \coloneqq \alpha_s/(1 - \alpha_s)$:

$$\sum_{s=1}^{T} \bar{w}_s \mathcal{J}_{\mathcal{C}}(\boldsymbol{z}_t, \mathcal{C}_{\boldsymbol{v}}) = \sum_{s=1}^{T} \mathbb{E}_{\boldsymbol{\epsilon}} [\|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\boldsymbol{\theta}}(\sqrt{\alpha_s} \hat{\boldsymbol{z}}_0(\boldsymbol{z}_t, \mathcal{C}_{\boldsymbol{v}}) + \sqrt{1 - \alpha_s} \boldsymbol{\epsilon}, \mathcal{C})\|_2^2$$
$$\gtrsim -\log p_{\boldsymbol{\theta}}(\hat{\boldsymbol{z}}_0(\boldsymbol{z}_t, \mathcal{C}_{\boldsymbol{v}}) \mid \mathcal{C}).$$

Intuitively, our optimization seeks to make the text-conditioned clean view $\hat{z}_0(z_t, C_v)$ of the current sample z_t as unique as possible, from the perspective of the log-likelihood $\log p_{\theta}(\hat{z}_0(z_t, C_v)|C)$.

Techniques for improvement. In practice, we found that our optimization could be further stabilized by introducing a sg-related trick into the objective function:

$$\tilde{\mathcal{J}}_{\mathcal{C}} \coloneqq \mathcal{J}_{\mathcal{C}}^{1} + \lambda \mathcal{J}_{\mathcal{C}}^{2}, \quad \lambda > 0$$
where
$$\mathcal{J}_{\mathcal{C}}^{1} \coloneqq \mathbb{E}_{\epsilon} \left[\left\| \hat{\boldsymbol{z}}_{0}(\boldsymbol{z}_{t}, \mathcal{C}_{\boldsymbol{v}}) - \operatorname{sg}\left(\hat{\boldsymbol{z}}_{0}(\boldsymbol{z}_{s|t,0}, \mathcal{C}) \right) \right\|_{2}^{2} \right] \qquad (7)$$

$$\mathcal{J}_{\mathcal{C}}^{2} \coloneqq \mathbb{E}_{\epsilon} \left[\left\| \operatorname{sg}\left(\hat{\boldsymbol{z}}_{0}(\boldsymbol{z}_{t}, \mathcal{C}_{\boldsymbol{v}}) \right) - \hat{\boldsymbol{z}}_{0}(\boldsymbol{z}_{s|t,0}, \mathcal{C}) \right\|_{2}^{2} \right].$$

In our empirical results, setting $\lambda = 1$ consistently produces the best performance across all con-sidered T2I models. We note that this technique allows the gradient flow through the second term (contrary to the case of Eq. (5)), thereby sidestepping the gradient blocking issue that we mentioned earlier. Another significant improvement comes from the use of an annealed timestep s, which was originally adhered to a fixed value in Um & Ye (2024). We empirically found that employing an an-nealing schedule based on the inverse of the sampling step (e.g., s = T - t) outperforms other fixed choices of s. Similar to Um & Ye (2024), we conduct our prompt optimization intermittently (*i.e.*, once every N sampling steps) to reduce computational costs. We found that during non-optimizing steps, employing the base prompt C instead of C_v (with the most recently updated token embed-ding) yields improvements in text-alignment and sample quality. See Algorithms 1 and 2 for the pseudocode of our approach.

Enhanced semantic controllability. A key benefit of our prompt optimization approach is its ability to provide an additional dimension of semantic control over the generated samples. Specifically, by selecting an appropriate initial point for v (*i.e.*, $v_T^{(0)}$ in Algorithm 1), such as a word embedding



"Three women in purple dresses on their cellphones..."

"A group of young children standing next to each other"

Figure 4: Sample comparison on SDXL-Lightning. Generated samples from three different approaches: (i) DDIM (Song et al., 2020a); (ii) SGMS (Um & Ye, 2024); (iii) MinorityPrompt (ours). Six distinct prompts were used for this comparison, and random seeds were shared across all three methods.

with relevant semantics, one can impart the desired semantics to the generated output; see Figure 3 for instance. Note that the controllability is not achievable with existing minority samplers that rely upon latent-space optimizations (Sehwag et al., 2022; Um & Ye, 2023; 2024). In addition, we found that properly choosing a initial word can yield improved minority generation performance compared to approaches that rely upon random starting points; see Table 2c for details.

4 EXPERIMENTS

4.1 Setup

T2I backbones and dataset. Our experiments were conducted using three distinct versions of Stable Diffusion (SD) (Rombach et al., 2022), encompassing both standard and distilled versions to demonstrate the robustness of our approach. Specifically, we consider: (i) SDv1.5; (ii) SDv2.0; (iii) SDXL-Lightning (SDXL-LT) (Lin et al., 2024). For all pretrained models, we employed the widely-adopted HuggingFace checkpoints trained on LAION (Schuhmann et al., 2022) without any further modifications. As convention, we randomly selected 10K captions from the validation set of MS-COCO (Lin et al., 2014) for the generations with SDv1.5 and SDv2.0 while using 5K captions for the SDXL-Lightning results.

Baselines. The same four baselines were considered over all SD versions: (i) the standard DDIM (Song et al., 2020a); (ii) a null-prompted DDIM; (iii) CADS (Sadat et al., 2023); (iv) SGMS (Um & Ye, 2024). The null-prompted DDIM serves as a naive baseline that attempts to encourage unique sampling by incorporating a proper null-text prompt, such as "commonly-looking". CADS (Sadat et al., 2023) is the state-of-the-art diversity-focused sampler that may rival our approach in minority generation, while SGMS (Um & Ye, 2024) is the state-of-the-art of mi-nority generation outside the T2I domain. We adhered to standard sampling setups for all methods. Specifically, 50 DDIM steps (*i.e.*, T = 50) with w = 7.5 were used for SDv1.5 and SDv2.0, while w = 1.0 was employed for the 4-step SDXL-Lightning model.

Evaluations. For evaluating text-alignment and user-preference, we consider three distinct quantities: (i) ClipScore (Hessel et al., 2021); (ii) PickScore (Kirstain et al., 2023); (iii) Image-Reward (Xu et al., 2023). We additionally employ two metrics for quality and diversity: Precision and Re-

Model	Method	CLIPScore \uparrow	PickScore \uparrow	ImageReward ↑	Precision \uparrow	Recall \uparrow	Likelihood↓
	DDIM	31.4801	21.4830	0.2106	0.5907	0.6328	1.0367
	DDIM + null	31.1007	21.5391	0.2422	0.5660	0.6236	1.0339
SDv1.5	CADS	31.4178	21.2836	0.1012	0.5696	0.6346	1.0127
	SGMS	31.1665	21.2126	0.1230	0.4943	0.5960	<u>0.9540</u>
	MinorityPrompt	31.5376	21.3111	0.2352	0.5671	0.6228	0.8971
	DDIM	31.8490	21.6801	0.3821	0.5930	0.6292	1.1100
	DDIM + null	31.7223	21.7190	0.4024	0.5861	0.6308	1.0769
SDv2.0	CADS	31.7687	21.5225	0.2981	0.5811	0.6194	1.0851
	SGMS	31.4750	21.4457	0.2981	0.5166	0.6130	<u>0.9898</u>
	MinorityPrompt	31.9586	21.5958	0.4249	0.6047	0.6100	0.9143
	DDIM	31.5714	22.6822	0.7317	0.5306	0.6648	0.6102
	DDIM + null	31.5754	22.7124	0.7397	0.5194	0.6602	0.6093
SDXL-LT	CADS	31.0837	22.3690	0.4946	0.5244	0.6594	0.6038
	SGMS	31.3589	22.5866	0.6759	0.4868	0.6968	0.5470
	MinorityPrompt	31.3838	22.6157	0.7042	0.4758	<u>0.6928</u>	0.5463

Table 1: **Quantitative comparisons.** "SDXL-LT" denotes SDXL-Lightning (4-step version) (Lin et al., 2024). "DDIM + null" indicates a baseline that leverages a properly-chosen null-prompt to encourage minority generations, where we used "commonly-looking" for the results herein. "CADS (Sadat et al., 2023)" is the state-of-the-art in diverse sampling, while SGMS (Um & Ye, 2024) denotes a minority sampler similar to ours, representing the state-of-the-art outside the T2I context. "Likelihood" represents log-likelihood values measured in bpd (bits per dimension).

call (Kynkäänniemi et al., 2019). For the likelihood of generated samples, we rely upon the exact
likelihood computation method based on PF-ODE as proposed by Song et al. (2020b). Notably, we
do not include Fréchet Inception Distance (FID) (Heusel et al., 2017) as an evaluator, since FID
measures closeness to baseline real data (*e.g.*, the MS-COCO validation set), which diverges from
our focus on promoting generations in low-density regions.

4.2 RESULTS

446

447

448

449

450

451 452

458 459

460

468

Qualitative comparisons. Figure 4 presents a comparison of generated samples of our approach with two baselines. Notice that our MinorityPrompt tends to yield highly more distinct and complex features (*e.g.*, intricate visual elements (Arvinte et al., 2023; Serrà et al., 2019)) compared to the baseline samplers, demonstrating its effectiveness even with distilled pretrained models. A significant observation, also reflected in Figure 1, is that MinorityPrompt often counters the inherent demographic biases of T2I models, *e.g.*, by adjusting age or skin color. See the samples in the second and third rows of the figure. A more extensive set of generated samples, including those from SDv1.5 and v2.0, can be found in Section D.2.

Quantitative evaluations. Table 1 exhibits performance comparisons across three distinct T2I models. Observe that our sampler outperforms all baselines in generating low-likelihood samples, while maintaining reasonable performance in text-to-alignment and user preference; also see Figure 5 where we present the distributions of log-likelihood. However, the performance trends for SDXL-LT results differ slightly from those of SDv1.5 and SDv2.0 across all tailored samplers, with particularly degraded results. We attribute this to the small-step nature of distilled models, which offer fewer opportunities to intervene in the sampling process, thereby limiting the potential for quantitative improvements.

476 Ablation studies. Table 2 investigates the impact of key design choices in our framework. Specif-477 ically, Table 2a highlights the benefits of optimizing small sets of token embeddings, which outper-478 form alternatives targeting text or null-text embeddings in both text alignment and log-likelihood. 479 The advantage of using the proposed objective function Eq. (6) is exhibited in Table 2b, where the 480 naively-extended framework based on Eq. (5) demonstrates significant performance gap compared 481 to our carefully-crafted approach. Table 2c explores various initialization techniques for v. While 482 all methods yield substantial improvements over the unoptimized sampler (see "unoptimized" in 483 Table 2b for comparison), we observe that further gains can be achieved with properly chosen initial words. A more comprehensive analysis and ablation study, encompassing additional design choices 484 and applications to trending sampling techniques such as CFG++ (Chung et al., 2024), is presented 485 in Section C.1.

Target	$CS\uparrow$	$LL\downarrow$	Method	$\mathbf{CS}\uparrow$	$LL\downarrow$	Туре	$\mathbf{CS}\uparrow$	$LL\downarrow$
Text	31.3503	0.9263	Unoptimized	31.4395	1.0465	Default	31.5154	0.935
Null-text	31.1089	1.0175	Naive (Eq. (5))	30.2994	0.9245	Gaussian	31.5054	0.942
ken (ours)	31.6465	0.9006	Ours (Eq. (6))	31.7369	0.9230	Word init	31.7369	0.923

Table 2: Ablation study results. "CS" denotes ClipScore (Hessel et al., 2021), while 'LL' indicates 492 log-likelihood. "Text" is the optimization framework focused on updating the text-embedding C, 493 and "Null-text" refers to the one that adjusts the null-text embedding (as in (Mokady et al., 2023)). 494 "Unoptimized" corresponds to the standard DDIM sampler. "Default" denotes the case that simply 495 employs the default embedding assigned with an added learnable token, while "Gaussian" initializes 496 v from a multivariate Gaussian distribution constructed using the mean and variance of the token 497 embeddings from the text-encoder \mathcal{T} . "Word init" indicates initializing with a specific word embed-498 ding. We used SDv1.5 for the results herein. 499

500 **Further applications.** Beyond our primary focus on minority generation, we additionally inves-501 tigate the potential of our framework, specifically in the perspectives of fairness and diversity. Ta-502 ble 3a presents one such instance. Although not explicitly designed to address demographic biases, our minority sampler demonstrates the ability to counteract gender bias and produce more neutral 504 generation results. This corroborates with our qualitative observations made in Figures 1 and 4.

Another area of investigation involves diversity, where we validate the versatility of our prompt optimizer in Eq. (3) for fostering diverse generation. To achieve this, we develop a new objective function aimed at encouraging diversity within a sampling batch that shares the same prompt (i.e., similar to the goal in Corso et al. (2023)) by enforcing repulsion between generated instances:

$$\bar{\mathcal{J}} \coloneqq \sum_{i=1}^{B} \sum_{k \neq i} \|\hat{\boldsymbol{z}}_0(\boldsymbol{z}_t^{(i)}, \mathcal{C}_{\boldsymbol{v}}) - \hat{\boldsymbol{z}}_0(\boldsymbol{z}_t^{(j)}, \mathcal{C}_{\boldsymbol{v}})\|_2^2,$$
(8)

where \mathcal{B} is the batch size, and $\{z_t^{(i)}\}_{i=1}^{\mathcal{B}}$ denotes the noisy instances in the batch. We found that incorporating this objective into Eq. (3) yields impressive results, even rivaling the state-of-the-art diverse sampler (Sadat et al., 2023); see Table 3b for details.

Method	$ p_{\text{female}} - p_{\text{male}} \downarrow$	Method	ClipScore \uparrow	$PickScore \uparrow$	ImReward \uparrow	$\text{Prec} \uparrow$	$\operatorname{Rec}\uparrow$	$\mathrm{IBS}\downarrow$
DDIM	0.3600	DDIM	31.4393	21.2478	0.0121	0.5860	0.6390	0.6164
CADS	0.2686	CADS	31.2692	21.0262	-0.0976	0.5620	0.5980	0.5494
Ours	0.2342	Ours	31.2724	21.0404	-0.1204	0.5480	0.6316	0.5439

(a) Debiasing effect of ours

(b) Effectiveness of our diversity-focused framework in Eq. (8)

Table 3: (a) Bias-mitigating impact of MinorityPrompt. p_{female} indicates the proportion of females in generated samples via gender-neutral prompts that include "a person" (e.g., "A person doing karate in a field at night"). On the other hand, p_{male} is the proportion of males in the samples. We employed SDXL-LT for these results. (b) Effectiveness of our diversity-focused prompt optimization framework. "ImReward" denotes Image-Reward metric. "IBS" represents In-Batch 526 Similarity, a diversity metric (Corso et al., 2023) that evaluates the cosine similarity in the DINO feature space (Caron et al., 2021). The results were obtained on SDv1.5.

528 529 530

491

505

506

507

508

513

514

521 522

523

524

525

527

5 CONCLUSION

531

532 We developed a novel framework for generating minority samples in the context of T2I generation. 533 Built upon our prompt optimization framework that updates the embeddings of additional learnable 534 tokens, our minority sampler offers significant performance improvements both in text-alignment and low-likelihood generation compared to existing approaches. To accomplish this, we meticu-536 lously tailor the objective function with theoretical justifications and implement several techniques for further enhancements. Beyond our main interest of minority generation, we further demonstrated the potential of our framework in promoting fairness and diversity. During this process, we 538 also showed that the proposed optimization framework can serve as a general solution, with potential applicability to various optimization tasks associated with T2I generation.

540 ETHICS STATEMENT

541 542

One potential concern associated with our approach is the possibility of its malicious use to inhibit
the generation of minority-featured samples. For instance, this could occur by flipping the sign of
the objective function Eq. (6), yielding a focus on high-density generations. It is crucial to recognize
this risk and to ensure that our proposed framework is employed responsibly to foster fairness and
inclusivity in generative modeling.

547 548

549

556

565

566

567

579

586

Reproducibility

To ensure the reproducibility of our experiments, we provide a comprehensive description regarding the employed pretrained models for our experiments. All experimental settings, including hyperparameter choices, are detailed in Section B.2. Additionally, we include the average running time of our algorithm along with specific details about the computer configuration in the same section. Finally, to assist with replication efforts, we have made our code available in a public repository: https://github.com/anonymous-6898/MinorityPrompt.

- 557 REFERENCES
- Arash Afkanpour, Vahid Reza Khazaie, Sana Ayromlou, and Fereshteh Forghani. Can generative
 models improve self-supervised representation learning? *arXiv preprint arXiv:2403.05966*, 2024.
- Marius Arvinte, Cory Cornelius, Jason Martin, and Nageen Himayat. Investigating the adversarial robustness of density estimation using the probability flow ode. *arXiv preprint arXiv:2310.07084*, 2023.
 - Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2021.
- 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.
- Hyungjin Chung, Suhyeon Lee, and Jong Chul Ye. Decomposed diffusion sampler for accelerating
 large-scale inverse problems. *arXiv preprint arXiv:2303.05754*, 2023a.
- 573
 574 Hyungjin Chung, Jong Chul Ye, Peyman Milanfar, and Mauricio Delbracio. Prompt-tuning latent diffusion models for inverse problems. *arXiv preprint arXiv:2310.01110*, 2023b.
- Hyungjin Chung, Jeongsol Kim, Geon Yeong Park, Hyelin Nam, and Jong Chul Ye.
 Cfg++: Manifold-constrained classifier free guidance for diffusion models. *arXiv preprint arXiv:2406.08070*, 2024.
- Gabriele Corso, Yilun Xu, Valentin De Bortoli, Regina Barzilay, and Tommi Jaakkola. Particle guidance: non-iid diverse sampling with diffusion models. *arXiv preprint arXiv:2310.13102*, 2023.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hi erarchical image database. In 2009 IEEE conference on computer vision and pattern recognition,
 pp. 248–255. Ieee, 2009.
- Jacob Devlin. 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, 34:8780–8794, 2021.
- 592 Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image
 593 synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recogni tion*, pp. 12873–12883, 2021.

618

632

- Felix Friedrich, Manuel Brack, Lukas Struppek, Dominik Hintersdorf, Patrick Schramowski, Sasha Luccioni, and Kristian Kersting. Fair diffusion: Instructing text-to-image generation models on fairness. *arXiv preprint arXiv:2302.10893*, 2023.
- Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel
 Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual
 inversion. *arXiv preprint arXiv:2208.01618*, 2022.
- Inhwa Han, Serin Yang, Taesung Kwon, and Jong Chul Ye. Highly personalized text embedding for
 image manipulation by stable diffusion. *arXiv preprint arXiv:2303.08767*, 2023.
- Jiyeon Han, Hwanil Choi, Yunjey Choi, Junho Kim, Jung-Woo Ha, and Jaesik Choi. Rarity
 score: A new metric to evaluate the uncommonness of synthesized images. *arXiv preprint arXiv:2206.08549*, 2022.
- Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter.
 Gans trained by a two time-scale update rule converge to a local nash equilibrium. Advances in neural information processing systems, 30, 2017.
- Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- Chunsan Hong, ByungHee Cha, and Tae-Hyun Oh. Cas: A probability-based approach for universal condition alignment score. In *International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=E780aH2s3f.
- Gaofeng Huang and Amir Hossein Jafari. Enhanced balancing gan: Minority-class image generation. *Neural computing and applications*, 35(7):5145–5154, 2023.
- Jeongsol Kim, Geon Yeong Park, Hyungjin Chung, and Jong Chul Ye. Regularization by texts for
 latent diffusion inverse solvers. *arXiv preprint arXiv:2311.15658*, 2023.
- Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Pick a-pic: An open dataset of user preferences for text-to-image generation. 2023.
- Tuomas Kynkäänniemi, Tero Karras, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Improved
 precision and recall metric for assessing generative models. *Advances in Neural Information Processing Systems*, 32, 2019.
- Tuomas Kynkäänniemi, Miika Aittala, Tero Karras, Samuli Laine, Timo Aila, and Jaakko Lehtinen.
 Applying guidance in a limited interval improves sample and distribution quality in diffusion models, 2024. URL https://arxiv.org/abs/2404.07724.
- Alexander C Li, Mihir Prabhudesai, Shivam Duggal, Ellis Brown, and Deepak Pathak. Your diffusion model is secretly a zero-shot classifier. *arXiv preprint arXiv:2303.16203*, 2023.
- 642
 643
 644
 644
 644
 644
 644
 644
 644
 644
 644
- 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 *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pp. 740–755. Springer, 2014.

658

665

667

672

679

681

648 649	Zinan Lin, Hao Liang, Giulia Fanti, Vyas Sekar, Rahul Anand Sharma, Elahe Soltanaghaei, Anthony Rowe, Hun Namkung, Zaoxing Liu, Daehyeok Kim, et al. Raregan: Generating samples for rare
650 651	classes. <i>arXiv preprint arXiv:2203.10674</i> , 2022.
652 653	Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for editing real images using guided diffusion models. In <i>Proceedings of the IEEE/CVF Conference</i>
654	on Computer Vision and Pattern Recognition, pp. 6038–6047, 2023.

- 655 Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, 656 Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with 657 text-guided diffusion models. arXiv preprint arXiv:2112.10741, 2021.
- Geon Yeong Park, Jeongsol Kim, Beomsu Kim, Sang Wan Lee, and Jong Chul Ye. Energy-based 659 cross attention for bayesian context update in text-to-image diffusion models. Advances in Neural 660 Information Processing Systems, 36, 2024. 661
- 662 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor 663 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-664 performance deep learning library. Advances in neural information processing systems, 32, 2019.
- Yiming Qin, Huangjie Zheng, Jiangchao Yao, Mingyuan Zhou, and Ya Zhang. Class-balancing 666 diffusion models. arXiv preprint arXiv:2305.00562, 2023.
- 668 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, 669 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual 670 models from natural language supervision. In International conference on machine learning, pp. 671 8748-8763. PMLR, 2021.
- Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, 673 and Ilya Sutskever. Zero-shot text-to-image generation. In International conference on machine 674 *learning*, pp. 8821–8831. Pmlr, 2021. 675
- 676 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-677 resolution image synthesis with latent diffusion models. In Proceedings of the IEEE/CVF confer-678 ence on computer vision and pattern recognition, pp. 10684–10695, 2022.
- Seyedmorteza Sadat, Jakob Buhmann, Derek Bradely, Otmar Hilliges, and Romann M Weber. Cads: 680 Unleashing the diversity of diffusion models through condition-annealed sampling. arXiv preprint arXiv:2310.17347, 2023. 682
- 683 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi 684 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An 685 open large-scale dataset for training next generation image-text models. Advances in Neural Information Processing Systems, 35:25278–25294, 2022. 686
- 687 Vikash Sehwag, Caner Hazirbas, Albert Gordo, Firat Ozgenel, and Cristian Canton. Generating high 688 fidelity data from low-density regions using diffusion models. In *Proceedings of the IEEE/CVF* 689 Conference on Computer Vision and Pattern Recognition, pp. 11492–11501, 2022. 690
- 691 Joan Serrà, David Álvarez, Vicenç Gómez, Olga Slizovskaia, José F Núñez, and Jordi Luque. In-692 put complexity and out-of-distribution detection with likelihood-based generative models. arXiv 693 preprint arXiv:1909.11480, 2019.
- Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. arXiv 695 preprint arXiv:2010.02502, 2020a. 696
- 697 Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. 698 Advances in Neural Information Processing Systems, 32, 2019. 699
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben 700 Poole. Score-based generative modeling through stochastic differential equations. arXiv preprint 701 arXiv:2011.13456, 2020b.

702 703 704	Yonglong Tian, Lijie Fan, Kaifeng Chen, Dina Katabi, Dilip Krishnan, and Phillip Isola. Learning vision from models rivals learning vision from data. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 15887–15898, 2024a.
705 706 707 708	Yonglong Tian, Lijie Fan, Phillip Isola, Huiwen Chang, and Dilip Krishnan. Stablerep: Synthetic images from text-to-image models make strong visual representation learners. <i>Advances in Neural Information Processing Systems</i> , 36, 2024b.
709 710	Soobin Um and Jong Chul Ye. Don't play favorites: Minority guidance for diffusion models. <i>arXiv</i> preprint arXiv:2301.12334, 2023.
711 712 713	Soobin Um and Jong Chul Ye. Self-guided generation of minority samples using diffusion models. <i>arXiv preprint arXiv:2407.11555</i> , 2024.
714 715 716	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, 2023.
717 718 719 720 721	Tao Xu, Pengchuan Zhang, Qiuyuan Huang, Han Zhang, Zhe Gan, Xiaolei Huang, and Xiaodong He. Attngan: Fine-grained text to image generation with attentional generative adversarial networks. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 1316–1324, 2018.
722 723 724	Fisher Yu, Ari Seff, Yinda Zhang, Shuran Song, Thomas Funkhouser, and Jianxiong Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. <i>arXiv</i> preprint arXiv:1506.03365, 2015.
725 726 727 728	Ning Yu, Ke Li, Peng Zhou, Jitendra Malik, Larry Davis, and Mario Fritz. Inclusive gan: Improving data and minority coverage in generative models. In <i>European Conference on Computer Vision</i> , pp. 377–393. Springer, 2020.
729	Oz Zafar, Lior Wolf, and Idan Schwartz. Iterative object count optimization for text-to-image diffu-
130	sion models. <i>arxiv preprud arxiv.</i> 2400.11721, 2024.
731 732 733	Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In <i>CVPR</i>, 2018.
731 732 733 734	Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In <i>CVPR</i>, 2018.
731 732 733 734 735	Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In <i>CVPR</i>, 2018.
730 731 732 733 734 735 736 727	Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In <i>CVPR</i>, 2018.
730 731 732 733 734 735 736 737 738	Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In <i>CVPR</i>, 2018.
730 731 732 733 734 735 736 737 738 739	 Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In <i>CVPR</i>, 2018.
730 731 732 733 734 735 736 737 738 739 740	 Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In <i>CVPR</i>, 2018.
730 731 732 733 734 735 736 737 738 739 740 741	 Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In <i>CVPR</i>, 2018.
730 731 732 733 734 735 736 737 738 739 740 741 742	 Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In <i>CVPR</i>, 2018.
730 731 732 733 734 735 736 737 738 739 740 741 742 743	 Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In <i>CVPR</i>, 2018.
730 731 732 733 734 735 736 737 738 739 740 741 742 743 744	 Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In <i>CVPR</i>, 2018.
730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745	 Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In <i>CVPR</i>, 2018.
730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746	 Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In <i>CVPR</i>, 2018.
730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747	Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In <i>CVPR</i> , 2018.
730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748	 Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In <i>CVPR</i>, 2018.
730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 744 745 746 747 748 749	Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In <i>CVPR</i> , 2018.
730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 744 745 746 747 748 749 750	Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In <i>CVPR</i> , 2018.
730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 744 745 746 747 748 749 750 751	Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In <i>CVPR</i> , 2018.
730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 744 745 746 747 748 749 750 751 752	Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In <i>CVPR</i> , 2018.
730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 744 745 746 747 748 749 750 751 752 753	Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In <i>CVPR</i> , 2018.
 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 754 	Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In <i>CVPR</i> , 2018.

A THEORETICAL RESULTS

A.1 PROOF OF PROPOSITION 1

Proposition 1. The objective function in Eq. (6) is equivalent (upto a constant factor) to the negative ELBO w.r.t. $\log p_{\theta}(\hat{z}_0(z_t, C_v) \mid C)$ when integrated over timesteps with $\bar{w}_s := \alpha_s/(1 - \alpha_s)$:

$$\sum_{s=1}^{T} \bar{w}_{s} \mathcal{J}_{\mathcal{C}}(\boldsymbol{z}_{t}, \mathcal{C}_{\boldsymbol{v}}) = \sum_{s=1}^{T} \mathbb{E}_{\boldsymbol{\epsilon}}[\|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\boldsymbol{\theta}}(\sqrt{\alpha_{s}}\hat{\boldsymbol{z}}_{0}(\boldsymbol{z}_{t}, \mathcal{C}_{\boldsymbol{v}}) + \sqrt{1 - \alpha_{s}}\boldsymbol{\epsilon}, \mathcal{C})\|_{2}^{2}] \qquad (9)$$
$$\gtrsim -\log p_{\boldsymbol{\theta}}(\hat{\boldsymbol{z}}_{0}(\boldsymbol{z}_{t}, \mathcal{C}_{\boldsymbol{v}}) \mid \mathcal{C}).$$

Proof. Remember the definition of the objective function in Eq. (6):

$$\mathcal{J}_{\mathcal{C}}(oldsymbol{z}_t,\mathcal{C}_{oldsymbol{v}})\coloneqq \mathbb{E}_{oldsymbol{\epsilon}}ig[\|\hat{oldsymbol{z}}_0(oldsymbol{z}_t,\mathcal{C}_{oldsymbol{v}}) - \hat{oldsymbol{z}}_0(oldsymbol{z}_{s|t,0},\mathcal{C})\|_2^2ig].$$

Plugging this into the LHS of Eq. (9) yields:

$$\sum_{s=1}^{T} \bar{w}_{s} \mathcal{J}_{\mathcal{C}}(\boldsymbol{z}_{t}, \mathcal{C}_{\boldsymbol{v}}) = \sum_{s=1}^{T} \frac{\alpha_{s}}{1 - \alpha_{s}} \mathbb{E}_{\boldsymbol{\epsilon}} \left[\| \hat{\boldsymbol{z}}_{0}(\boldsymbol{z}_{t}, \mathcal{C}_{\boldsymbol{v}}) - \hat{\boldsymbol{z}}_{0}(\boldsymbol{z}_{s|t,0}, \mathcal{C}) \|_{2}^{2} \right]$$
$$= \sum_{s=1}^{T} \frac{\alpha_{s}}{1 - \alpha_{s}} \mathbb{E}_{\boldsymbol{\epsilon}} \left[\left\| \frac{1}{\sqrt{\alpha_{s}}} (\boldsymbol{z}_{s|t,0} - \sqrt{1 - \alpha_{s}} \boldsymbol{\epsilon}) - \frac{1}{\sqrt{\alpha_{s}}} (\boldsymbol{z}_{s|t,0} - \sqrt{1 - \alpha_{s}} \boldsymbol{\epsilon}_{\boldsymbol{\theta}}(\boldsymbol{z}_{s|t,0}, \mathcal{C})) \right\|_{2}^{2} \right]$$
$$= \sum_{s=1}^{T} \frac{\alpha_{s}}{1 - \alpha_{s}} \mathbb{E}_{\boldsymbol{\epsilon}} \left[\left\| \frac{\sqrt{1 - \alpha_{s}}}{\sqrt{\alpha_{s}}} (\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\boldsymbol{\theta}}(\boldsymbol{z}_{s|t,0}, \mathcal{C})) \right\|_{2}^{2} \right]$$
$$= \sum_{s=1}^{T} \mathbb{E}_{\boldsymbol{\epsilon}} [\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\boldsymbol{\theta}}(\sqrt{\alpha_{s}} \hat{\boldsymbol{z}}_{0}(\boldsymbol{z}_{t}, \mathcal{C}_{\boldsymbol{v}}) + \sqrt{1 - \alpha_{s}} \boldsymbol{\epsilon}, \mathcal{C}) \|_{2}^{2}],$$
(10)

where the second equality is from the definitions of $z_{s|t,0}$ and $\hat{z}_0(z_{s|t,0}, C)$:

$$oldsymbol{z}_{s|t,0} \coloneqq \sqrt{lpha_s} \hat{oldsymbol{z}}_0(oldsymbol{z}_t, oldsymbol{\mathcal{C}}_{oldsymbol{v}}) + \sqrt{1 - lpha_s} oldsymbol{\epsilon} \ \hat{oldsymbol{z}}_0(oldsymbol{z}_{s|t,0}, oldsymbol{\mathcal{C}}) \coloneqq rac{1}{\sqrt{lpha_s}} (oldsymbol{z}_{s|t,0} - \sqrt{1 - lpha_s} oldsymbol{\epsilon}_{oldsymbol{ heta}}(oldsymbol{z}_{s|t,0}, oldsymbol{\mathcal{C}})).$$

Note that the last expression in Eq. (10), which is the same as the RHS of Eq. (9), is equivalent (up to a constant) to the expression of the negative ELBO w.r.t. $\hat{z}_0(z_t, C_v)$ (Ho et al., 2020; Li et al., 2023). The distinction here is that now we use a text-conditional diffusion model $\epsilon_{\theta}(\cdot, C)$ that approximates $\log p_{\theta}(\cdot|C)$. This completes the proof.

A.2 THEORETICAL ISSUES ON EQ. (5)

We continue from Section 3.3 to scrutinize the theoretical challenges that arise in the naivelyextended optimization framework in Eq. (5). To proceed, we first restate the objective function in Eq. (5):

$$\mathcal{J}(\boldsymbol{z}_t, \mathcal{C}_{\boldsymbol{v}}) \coloneqq \mathbb{E}_{\boldsymbol{\epsilon}} \big[\| \hat{\boldsymbol{z}}_0^w(\boldsymbol{z}_t, \mathcal{C}_{\boldsymbol{v}}) - \operatorname{sg}(\hat{\boldsymbol{z}}_0^w(\boldsymbol{z}_{s|t,0}^w, \mathcal{C}_{\boldsymbol{v}})) \|_2^2 \big].$$

Remember that we identified three theoretical issues that impair the connection to the target loglikelihood log $p_{\theta}(z_0 | C)$: (i) the reliance on the CFG-based clean predictions; (ii) obstructed gradient flow through the second term in the squared-L2 loss; and (iii) the incorporation of C_v within the second term in the loss.

CFG-based clean prediction. We start by examining the first point, the pathology due to the CFG-based clean predictions. Suppose we incorporate the CFG-based clean predictions \hat{z}_0^w in our framework Eq. (6), in place of the non-CFG terms \hat{z}_0 . The objective function then becomes:

$$\mathcal{J}^w_\mathcal{C}(oldsymbol{z}_t,\mathcal{C}_oldsymbol{v})\coloneqq \mathbb{E}_{oldsymbol{\epsilon}}[\|\hat{oldsymbol{z}}^w_0(oldsymbol{z}_t,\mathcal{C}_oldsymbol{v})-\hat{oldsymbol{z}}^w_0(oldsymbol{z}_{s|t,0}^w,\mathcal{C})\|_2^2]$$

To see its connection to log-likelihood, let us consider the weighted sum of this objective with $\bar{w}_s := \alpha_s/(1-\alpha_s)$ (as in Proposition 1). Manipulating the averaged objective similarly as in Section A.1 then yields:

813 814 815

821

$$\sum_{s=1}^{T} \bar{w}_{s} \mathcal{J}_{\mathcal{C}}^{w}(\boldsymbol{z}_{t}, \mathcal{C}_{\boldsymbol{v}}) = \sum_{s=1}^{T} \bar{w}_{s} \mathbb{E}_{\boldsymbol{\epsilon}} \left[\| \hat{\boldsymbol{z}}_{0}^{w}(\boldsymbol{z}_{t}, \mathcal{C}_{\boldsymbol{v}}) - \hat{\boldsymbol{z}}_{0}^{w}(\boldsymbol{z}_{s|t,0}^{w}, \mathcal{C}) \|_{2}^{2} \right]$$
$$= \sum_{s=1}^{T} \frac{\alpha_{s}}{1 - \alpha_{s}} \mathbb{E}_{\boldsymbol{\epsilon}} \left[\left\| \frac{1}{\sqrt{\alpha_{s}}} (\boldsymbol{z}_{s|t,0}^{w} - \sqrt{1 - \alpha_{s}} \boldsymbol{\epsilon}) - \frac{1}{\sqrt{\alpha_{s}}} (\boldsymbol{z}_{s|t,0}^{w} - \sqrt{1 - \alpha_{s}} \boldsymbol{\epsilon}_{\boldsymbol{\theta}}^{w}(\boldsymbol{z}_{s|t,0}^{w}, \mathcal{C})) \right\|_{2}^{2} \right]$$
$$= \sum_{s=1}^{T} \mathbb{E}_{\boldsymbol{\epsilon}} [\| \boldsymbol{\epsilon} - \tilde{\boldsymbol{\epsilon}}_{\boldsymbol{\theta}}^{w}(\sqrt{\alpha_{s}} \hat{\boldsymbol{z}}_{0}^{w}(\boldsymbol{z}_{t}, \mathcal{C}_{\boldsymbol{v}}) + \sqrt{1 - \alpha_{s}} \boldsymbol{\epsilon}, \mathcal{C}) \|_{2}^{2}]. \tag{11}$$

B22 Observe that in the RHS of Eq. (11), we see the CFG noise estimation term $\tilde{\epsilon}_{\theta}^{w}$, instead of ϵ_{θ} as in Eq. (10). This comes from the use of $\hat{z}_{0}^{w}(\boldsymbol{z}_{s|t,0}^{w}, \mathcal{C})$ in the second term of the squared-L2 loss. Since $\tilde{\epsilon}_{\theta}^{w}$ represents a distinct probability density, say $\tilde{p}_{\theta}(\cdot | \mathcal{C})$, the averaged objective in Eq. (11) is no longer connected to our focused conditional log-likelihood log $p_{\theta}(\cdot | \mathcal{C})$.

One may wonder whether the use of CFG for the first term in the squared-L2 loss of Eq. (6) is safe. However, we claim that it is also problematic. To show this, we derive the associated log-likelihood, which is immediate with the algebra used for Eq. (11):

830 831

832

848

854

855 856

857 858

859

$$\sum_{s=1}^{I} \bar{w}_s \mathbb{E}_{\boldsymbol{\epsilon}} \left[\| \hat{\boldsymbol{z}}_0^w(\boldsymbol{z}_t, \mathcal{C}_{\boldsymbol{v}}) - \hat{\boldsymbol{z}}_0(\boldsymbol{z}_{s|t,0}^w, \mathcal{C}) \|_2^2 \right] \gtrsim -\log p_{\boldsymbol{\theta}}(\hat{\boldsymbol{z}}_0^w(\boldsymbol{z}_t, \mathcal{C}_{\boldsymbol{v}}) \mid \mathcal{C}).$$

833 We see that now the diffusion model (represented by p_{θ}) should estimate the conditional log-density 834 w.r.t. the CFG clean prediction $\hat{z}_0^w(z_t, \mathcal{C}_v)$. We argue that this estimation may be inaccurate, since the CFG clean sample in the T2I context is potentially off-manifold. As analyzed in Chung et al. 835 (2024), the CFG clean prediction $\hat{z}_0^w(z_t, C_v)$ is in fact an extrapolation between $\hat{z}_0(z_t, C_v)$ and 836 $\hat{z}_0(z_t)$ (controlled by w). As a result, it may deviate from the data manifold, particularly for high w 837 values commonly used in standard T2I scenarios; see Figure 3 in Chung et al. (2024) for details. This 838 off-manifold issue is especially pronounced during the initial phase of inference, as also reported in 839 other studies (Kynkäänniemi et al., 2024). See Table 4 for experimental results that support this 840 claim. 841

Obstructed gradient. Now we move onto the second issue. From the above analysis, we saw that the noise prediction in the second term is crucial for relating the objective function to the log-likelihood, meaning that allowing gradient flow through the second term is essential for accurate likelihood optimization. However, blocking the gradient via the stop-gradient on the second term contradicts this theoretical intuition. We found that the use of stop-gradient actually degrades performance; see Table 4 for instance.

 C_v in the second term. The reasoning behind the third challenge follows naturally from the previous analyses. In this case, the corresponding log-likelihood term can be derived as:

$$\sum_{s=1}^{T} \bar{w}_{s} \mathbb{E}_{\boldsymbol{\epsilon}} \big[\| \hat{\boldsymbol{z}}(\boldsymbol{z}_{t}, \mathcal{C}_{\boldsymbol{v}}) - \hat{\boldsymbol{z}}_{0}(\boldsymbol{z}_{s|t,0}, \mathcal{C}_{\boldsymbol{v}}) \|_{2}^{2} \big] \gtrsim -\log p_{\boldsymbol{\theta}}(\hat{\boldsymbol{z}}_{0}(\boldsymbol{z}_{t}, \mathcal{C}_{\boldsymbol{v}}) \mid \mathcal{C}_{\boldsymbol{v}}).$$

We see that C_{v} appears in conditioning variable, which diverges from our interest of approximating $\log p_{\theta}(\cdot | C)$. See Table 4 for experimental results that corroborate this.

B SUPPLEMENTARY DETAILS

 \mathbf{T}

B.1 DETAILS ON THE METRIC IN UM & YE (2024)

We continue from Section 3.3 to provide additional details on the likelihood metric developed by Um & Ye (2024). This original version is defined on pixel space $x_0 \in \mathbb{R}^d$ (rather than latent domain $z_0 \in \mathbb{R}^k$ as ours), formally written as (Um & Ye, 2024):

$$\mathcal{J}(\boldsymbol{x}_t;s) \coloneqq \mathbb{E}_{\boldsymbol{\epsilon}} \left[d(\hat{\boldsymbol{x}}_0(\boldsymbol{x}_t), \hat{\boldsymbol{x}}_0(\boldsymbol{x}_{s|t,0})) \right]_{\boldsymbol{\epsilon}}$$

where x_t is a noisy pixel-domain image, and $\hat{x}_0(x_t)$ represents a clean estimate of x_t : $\hat{x}_0(x_t) := (x_t - \sqrt{1 - \alpha_t} \epsilon'_{\theta}(x_t))/\sqrt{\alpha_t}$, where ϵ'_{θ} denotes a pixel diffusion model (different from our ϵ_{θ}). Here $x_{s|t,0}$ indicates a noised version of $\hat{x}_0(x_t)$ according to timestep s: $x_{s|t,0} := \sqrt{\alpha_s} \hat{x}_0(x_t) + \sqrt{1 - \alpha_s} \epsilon$, and $\hat{x}_0(x_{s|t,0})$ is a denoised version of $x_{s|t,0}$. d indicates a discrepancy metric (e.g., LPIPS (Zhang et al., 2018)). This quantity is interpretable as a reconstruction loss of $\hat{x}_0(x_t)$, and theoretically, it is an estimator of the negative log-likelihood of $\hat{x}_0(x_t)$ (Um & Ye, 2024).

Similar to ours, the authors in Um & Ye (2024) employs this metric as a guidance function for
minority sampling, sharing similar spirit as ours. In doing so, they propose several techniques such as stop-gradient, learning-rate scheduling, and the incorporation of LPIPS as *d*. Their proposed metric for the guidance function is expressible as:

 $\mathcal{J}(\boldsymbol{x}_t; s) \coloneqq \eta_t \mathbb{E}_{\boldsymbol{\epsilon}} [\text{LPIPS}(\hat{\boldsymbol{x}}_0(\boldsymbol{x}_t), \text{sg}(\hat{\boldsymbol{x}}_0(\boldsymbol{x}_{s|t,0})))],$

where η_t indicates learning rate at time t designed to decrease over time, and LPIPS is the perceptual metric proposed by Zhang et al. (2018). Although this approach offers considerable advantages in traditional image generation tasks (such as unconditional generation), it is not optimized for T2I generation, which presents unique challenges and requires more specialized techniques. This is confirmed by our experimental results, where a straightforward extension of their framework yields only modest performance improvements. See Table 2b for details.

882 883 884

B.2 IMPLEMENTATION DETAILS

885 Pretrained models and baselines. We employed the official checkpoints provided in HuggingFace 886 for all three pretrained models. For the null-prompted DDIM baselines, we employed "commonly-887 looking" as the null-text prompt for all three pretrained models. The CADS baselines were primarily obtained using the recommended settings in the paper (Sadat et al., 2023), while we adjusted the hy-889 perparameters on SDXL-Lightning for adaptation to distilled models. Specifically, we set $\tau_1 = 0.8$, 890 $\tau_2 = 1.0$, and s = 0.1, while keeping other settings unchanged. For SGMS, we respected the origi-891 nal design choices (like the use of sg) and tuned the remaining hyperparameters to attain the optimal performance in the T2I context. In particular, we used the squared-L2 loss as the discrepancy metric 892 and employed s = 0.75T. For their latent optimizations, we employed Adam optimizer (Kingma, 893 2014) (as ours) with learning rates between 0.005 and 0.01. Similar to ours, latent updates were 894 performed intermittently, with N = 3 (*i.e.*, one update per three sampling steps). Each latent opti-895 mization consisted of three distinct update steps: K = 3. 896

Evaluations. The ClipScore values reported in our paper were due to torchmetrics³. For PickScore and Image-Reward, we employed the implementations provided in the official code repositories⁴⁵. Precision and Recall were computed with k = 5 using the official codebase of Han et al. (2022)⁶. The log-likelihood values were evaluated based on the implementation of Hong et al. (2024)⁷. In-Batch Similarity that we used in the diversity optimization (in Table 3b) were computed with the repository of Corso et al. (2023)⁸.

Hyperparameters. Our results were obtained using s = T - t, and we used Adam optimizer with 903 K = 3, similar to SGMS. Learning rates were set between 0.001 and 0.002 across all experiments. 904 We shared the same intermittent update rate of N = 3 with SGMS. For initializing v, we shared 905 the same word embedding for "cool" for the main results (presented in Table 1). The number of 906 learnable tokens for our approaches was set to 1. As described in Section 3.3, we globally used 907 $\lambda = 1$ across all experiments. For the experiments on SDXL-Lightning that involves two distinct 908 text-encoders, we employed a single Adam optimizer to jointly update both embedding spaces to 909 minimize parameter complexity. We also synchronized other design choices for the two encoders, 910 e.g., sharing the same initial token embedding.

⁹¹¹ 912

^{912 &}lt;sup>3</sup>https://lightning.ai/docs/torchmetrics/stable/multimodal/clip_score. 913 html

^{914 &}lt;sup>4</sup>https://github.com/yuvalkirstain/PickScore

^{915 &}lt;sup>5</sup>https://github.com/THUDM/ImageReward

^{916 &}lt;sup>6</sup>https://github.com/hichoe95/Rarity-Score

^{917 &}lt;sup>7</sup>https://github.com/unified-metric/unified_metric

⁸https://github.com/gcorso/particle-guidance

8	Method	CLIPScore ↑	PickScore \uparrow	ImageReward ↑	Precision \uparrow	Recall ↑	Likelihood ↓
	DDIM	31.4395	21.4570	0.1845	0.6070	0.7094	1.0465
	Eq. (6) (proposed)	31.7369	21.3522	0.2839	0.5420	0.7340	0.9230
	Eq. (6) + \hat{z}_0^w	30.5193	20.7307	-0.1468	0.4890	0.7182	0.9399
	Eq. (6) + sg	31.6597	21.3114	0.2738	0.5230	0.7284	0.9290
	Eq. (6) + C_v	<u>31.6676</u>	21.3652	0.2808	0.5550	0.7262	0.9281
	Eq. (6) + all $(i.e., Eq. (5))$	30.2994	20.4840	-0.1944	0.4760	0.6864	0.9245

925 Table 4: Impacts of theoretical flaws in Eq. (5). "+ \hat{z}_0^w " indicates the case that further incorporates 926 the CFG clean predictions into Eq. (6). "+ sg" refers to the one employing the stop-gradient on $\hat{z}_0(z_{s|t,0}, C)$. "+ C_v " represents the setting of feeding C_v in the computations of $\hat{z}_0(z_{s|t,0}, C)$ in 927 place of C. "+ all" is the case that employs all the above three flawed choices, *i.e.*, Eq. (5). We 928 observe clear performance benefits of our theory-driven design choices over the naive framework 929 in Eq. (5). The results were obtained on SDv1.5. 930

931

932 933

934 935

936 937

938

939

940

941

Target	$CS\uparrow$	$LL\downarrow$	Туре	$\mathbf{CS}\uparrow$	$LL\downarrow$		Туре	$CS\uparrow$	LI
Eq. (6)	31.3658	0.5449	s = 0.75T	31.4534	0.9469		C	31.7548	0.97
Eq. (7)	31.4194	0.5449	s = T - t	31./369	0.9230		$\mathcal{C}_{\boldsymbol{v}^*}$	31.7871	0.95
(a) Infl	luence of s	g-trick	(b) Impact of adaptive s				(c)	Effect of us	sing $\mathcal C$

Table 5: Effectiveness of our new techniques. "CS" denotes ClipScore (Hessel et al., 2021), while 'LL' indicates log-likelihood. " \mathcal{C} " refers to the use of \mathcal{C} during sampling steps without prompt optimization (when incorporating an intermittent prompt update, *i.e.*, N > 1). On the other hand, " \mathcal{C}_{v^*} " refers to the use of optimized token embeddings in the latest steps. Our results show that the proposed design choices consistently outperform naive approaches. The results in (a) were obtained using SDXL-Lightning, while SDv1.5 was employed for (b) and (c).

942 943 944

945

946

947

948

949

950

951

952

953

954 955

Computational complexity. The inference time for DDIM is approximately 1.136 seconds per sample, with CADS requiring a similar amount of time. The complexities of SGMS and our approach are rather higher due to the inclusion of backpropagation and iterative updates of latents or prompts. Specifically, SGMS takes 5.756 seconds per sample, while our sampler requires slightly more time -6.205 seconds per sample – which we attribute to the additional backpropagation pass introduced by our removal of gradient-blocking. All computations herein were performed on SDv2.0 using a single NVIDIA A100 GPU.

Other details. Our implementation is based on PyTorch Paszke et al. (2019), and experiments were performed on twin NVIDIA A100 GPUs. Code is available at https://github.com/ anonymous-6898/MinorityPrompt. For

С ADDITIONAL ABLATIONS AND DISCUSSIONS

956 957 958

959

FURTHER ABLATION STUDIES C.1

Table 4 exhibits the individual impacts of the three theoretical flaws in the naive framework 960 in Eq. (5). We highlight that our new design choices motivated by a set of careful theoretical analy-961 ses yields significant advantages specifically in preserving text-alignment and user-preference. This 962 further validates our framework as a powerful minority sampler that achieves high-quality genera-963 tion. 964

Table 5 explores the impact of our techniques developed for further improvements in Section 3.3. We 965 see consistent enhancements over naive design choices. A key insight from Table 5c is that reusing 966 token embeddings optimized at earlier timesteps, denoted as " C_{v^*} " in the table, offers limited benefit 967 compared to simply using the base prompts C. This finding highlights the evolutionary nature of 968 our prompt-tuning framework, which supports continual updates to embeddings across sampling 969 timesteps. 970

Table 6 investigates the design choices related to learnable tokens in our framework. Observe in Ta-971 ble 6a that our framework consistently delivers significant performance gains across different initial

Ini	t word	$CS\uparrow$	$LL\downarrow$	Position	$CS\uparrow$	$LL\downarrow$	# of tokens	$\mathbf{CS}\uparrow$	$LL\downarrow$
"unc	ommon"	31.6971	0.8868	_	31.4395	1.0465	1	31.6465	0.9006
"s	pecial"	31.6178	0.9342	Prefix	31.5519	0.9249	2	31.5866	0.9163
	cool"	31.7369	0.9230	Postfix	31.7369	0.9230	4	31.4989	0.9419

Table 6: Exploring the design space of learnable tokens. "Init word" indicates the word embed-978 ding used for initializing v. "-" refers to standard DDIM sampling without prompt optimization. "Prefix" denotes prepending the placeholder string S to P, while "Postfix" indicates appending it to the end of \mathcal{P} . "# of tokens" represents the number of tokens assigned to the string \mathcal{S} . We observe that the proposed approach is not highly sensitive to the choice of initial word, and as suggested, 982 attaching S at the end of the prompts yields the best performance. Additionally, using a single token is sufficient to achieve performance gains. We used SDv1.5 for the results herein. 984

Method	CLIPScore \uparrow	PickScore ↑	ImageReward \uparrow	Precision \uparrow	Recall \uparrow	Likelihood ↓
DDIM	31.4395	21.4570	0.1845	0.6070	0.7094	1.0465
DDIM-CFG++	31.4755	21.4490	0.1938	<u>0.5710</u>	0.7100	1.0452
Ours	<u>31.7369</u>	21.3522	0.2839	0.5420	0.7340	0.9230
Ours-CFG++	31.7627	21.3399	0.3062	0.5540	0.7284	0.9183

Table 7: Compatibility with other sampling techniques. "DDIM-CFG++" represents the DDIM 993 sampler integrated with CFG++ (Chung et al., 2024), while "Ours-CFG++" is our MinorityPrompt 994 framework implemented with CFG++. We highlight that MinorityPrompt demonstrates strong com-995 patibility, delivering significant performance gains even when combined with CFG++. The guidance 996 weights for the CFG++ cases were set to 0.6, *i.e.*, the recommended setting in the paper (Chung et al., 997 2024). We used SDv1.5 for the results herein. 998

999

977

979

980

981

983

985 986 987

1000

1001

1002 word embeddings. Regarding the position of S, appending it to the end of the prompts yields better 1003 results. We speculate that prepending may have a greater impact on the semantics of the text embeddings due to the front-weighted nature of the training process for the CLIP text encoders (Radford 1004 et al., 2021) employed in our T2I models. As exhibited in Table 6c, a single token is sufficient to 1005 realize the performance benefits of our approach. The performance degradation observed with increasing tokens is likely due to their heightened influence on semantics, similar to the effect of S's 1007 position. 1008

1009 Table 7 investigates the performance of MinorityPrompt when integrated with CFG++ (Chung et al., 2024) across different sampling techniques. We see that MinorityPrompt consistently outperforms 1010 standard DDIM-CFG++ by a notable margin, demonstrating the robustness and adaptability of our 1011 approach. This compatibility with CFG++ highlights the flexibility of MinorityPrompt, enabling 1012 substantial gains even when leveraging advanced conditioning strategies. 1013

- 1014
- 1015
- 1016

C.2 LIMITATIONS AND DISCUSSION 1017

1018

A disadvantage is that our framework introduces additional computational costs (similar to Um & 1020 Ye (2024)), particularly when compared to standard samplers like DDIM. As noted in Section B.2, 1021 this is mainly due to the incorporation of backpropagation and iterative updates of prompts. Additionally, the removal of gradient-blocking, aimed at restoring the theoretical connection to the target conditional density, further contributes to the overhead. Future work could focus on optimizing 1023 these processes to reduce computational demands. One potential approach is to develop an approx-1024 imation of our objective that mitigates the need for extensive backpropagation while maintaining its 1025 alignment with the target log-likelihood.



Figure 5: **Comparison of log-likelihood distributions.** The likelihood values were measured using the PF-ODE-based computation proposed by Song et al. (2020b). We observe that MinorityPrompt better produces low-likelihood instances compared to the considered baselines across all three pretrained models.

1042 D ADDITIONAL EXPERIMENTAL RESULTS

1044 D.1 LOG-LIKELIHOOD DISTRIBUTIONS

Figure 5 exhibits the log-likelihood distributions for MinorityPrompt and the baseline models across all three pretrained architectures. We see that MinorityPrompt consistently produces lower loglikelihood instances, further demonstrating its improved capability of generating minority samples. The distributions for SDXL-Lightning are more dispersed than in other scenarios, which may be attributed to the larger latent space upon which SDXL-Lightning is based. The competitive results compared to SGMS observed in SDXL-Lightning may arise from the limited optimization opportunities available in distilled models (as discussed in the manuscript).

1052 1053 1054

1041

D.2 ADDITIONAL GENERATED SAMPLES

To facilitate a more comprehensive qualitative comparison among the samplers, we provide an extensive showcase of generated samples for all the focused T2I pretrained models. See Figures 6–8 for details.

1058	
1059	
1060	
1061	
1062	
1063	
1064	
1065	
1066	
1067	
1068	
1069	
1070	
1071	
1072	
1073	
1074	
1075	
1076	
1077	
1078	



Figure 6: Generated samples on SDv1.5. Generated samples from three distinct samplers: (i)
DDIM (Song et al., 2020a); (ii) SGMS (Um & Ye, 2024); (iii) MinorityPrompt (ours). Random seeds were shared across all three methods.



Figure 7: **Generated samples on SDv2.0.** Generated instances from three different techniques: (i) DDIM (Song et al., 2020a); (ii) SGMS (Um & Ye, 2024); (iii) MinorityPrompt (ours). We shared the same random seeds across all three approaches.



Figure 8: Additional generated samples on SDXL-Lightning. Generated samples from three different approaches: (i) DDIM (Song et al., 2020a); (ii) SGMS (Um & Ye, 2024); (iii) MinorityPrompt (ours). We employed the same initial noises across all three samplers.