

CONTRASTIVE CFG: GUIDING DIFFUSION SAMPLING BY CONTRASTING POSITIVE AND NEGATIVE CONCEPTS

Anonymous authors

Paper under double-blind review

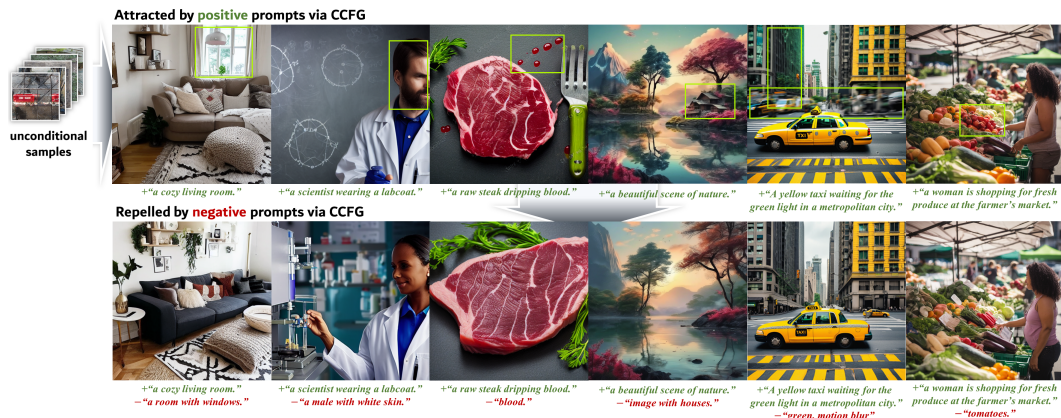


Figure 1: Image samples generated from StableDiffusion 1.5 and SDXL by our proposed Contrastive CFG (CCFG) guidance term, using *positive* prompts (written in green) and *negative* prompts (written in red). The samples in the same column were from the same initial noise.

ABSTRACT

As Classifier-Free Guidance (CFG) has proven effective in conditional diffusion model sampling for improved condition alignment, many applications use a negated CFG term as a Negative Prompting (NP) to filter out unwanted features from samples. However, simply negating CFG guidance creates an inverted probability distribution, often distorting samples away from the marginal distribution. Inspired by recent advances in conditional diffusion models for inverse problems, here we present a novel method to achieve guidance toward the given condition using contrastive loss. Specifically, our guidance term aligns or repels the denoising direction based on the given condition through contrastive loss, achieving a similar guiding effect to traditional CFG for positive conditions while overcoming the limitations of existing negative guidance methods. Experimental results demonstrate that our approach effectively injects or removes the given concepts while maintaining sample quality across diverse scenarios, from simple class conditions to complex and overlapping text prompts.

1 INTRODUCTION

Classifier-Free Guidance (CFG) (Ho & Salimans, 2022) forms the key basis of modern text-guided generation with diffusion models. From Bayes rule, CFG constructs a Bayesian classifier $\nabla_{\mathbf{x}} \log p(\mathbf{c}|\mathbf{x}) = \nabla_{\mathbf{x}} \log p(\mathbf{x}|\mathbf{c}) - \nabla_{\mathbf{x}} \log p(\mathbf{x})$ without training additional external classifiers (Dhariwal & Nichol, 2021). In practice, it is common to emphasize the classifier vector direction with some constant γ , which corresponds to sharpening the posterior, i.e. $p(\mathbf{x})p(\mathbf{c}|\mathbf{x})^\gamma$. While this may potentially distort the sampling distribution (Du et al., 2023), it is known that the sample quality along with the text alignment increases.

While CFG was originally devised for *adhering* to the target condition, more often than not, there are cases where it is desirable to *avoid* sampling from some conditions. Canonical examples include

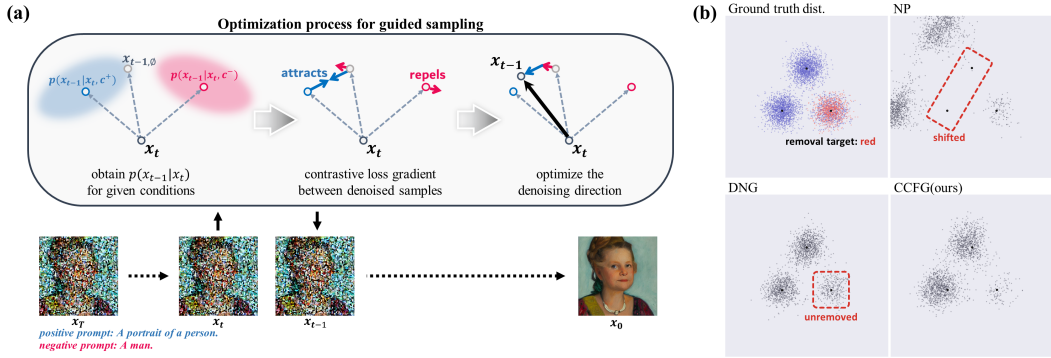


Figure 2: **Overview of the proposed guided sampling of CCFG.** (a) We pose guided sampling as an optimization problem that minimizes the contrastive loss of the positive and negative prompts, which has no computational overhead, yet avoids pitfalls of previous strategies such as NP. (b) Output distribution with different negative sampling methods on a toy dataset with two classes.

conditions that describe the poor quality of the image, or conditions that are related to harmful content (Gandikota et al., 2023; Wu et al., 2024). Often referred to as *Negative Prompting (NP)*, most implementations simply negate the vector direction of CFG, corresponding to inversely weighting with the classifier probability, i.e. $p(\mathbf{x})/p(c|\mathbf{x})^\gamma$. Nonetheless, such naïve implementation often leads to a decrease in the sample quality (Kim et al., 2025a; Armandpour et al., 2023). Specifically, due to the unboundness of an inverted probability distribution, sampling from $p(\mathbf{x})/p(c|\mathbf{x})^\gamma$ leads to sampling from the low-density region or completely outside of the original density support.

To mitigate these drawbacks, by leveraging the recent advances in guided sampling methods (Chung et al., 2023; Kim et al., 2025b; Chung et al., 2024b; Yu et al., 2023; Ye et al., 2024), we reformulate the conditional guidance as an optimization problem with a positive or negative prompt-conditioned contrastive loss (Gutmann & Hyvärinen, 2010; Chen et al., 2020). Then, we derive a reverse diffusion sampling strategy with the optimized denoising direction. This results in a simple modification of CFG to the sampling process with little computational overhead. The resulting process, termed *Contrastive CFG (CCFG)*, optimizes the denoising direction by attracting or repelling the denoising direction for the given condition. Furthermore, the attracting and repelling forces are automatically controlled along the sampling process. Through extensive experiments, we verify that CCFG not only guides the samples to satisfy the wanted conditions, but also successfully avoid undesirable concepts while preserving the sample quality.

2 RELATED WORKS

2.1 CLASSIFIER-FREE GUIDANCE FOR DIFFUSION MODELS.

Diffusion models (Song et al., 2021b; Ho et al., 2020) are a class of generative models that learn the score function (Hyvärinen & Dayan, 2005) of the data distribution, and use this score function to reverse the forward noising process. The forward process, denoted with the time index t , is governed by a Gaussian kernel that the underlying data distribution $p(\mathbf{x}_0) \equiv p(\mathbf{x})$ eventually approximates the standard normal distribution at time $t = T$, i.e. $p(\mathbf{x}_T) \approx \mathcal{N}(0, \mathbf{I})$. The variance preserving forward transition kernel (Ho et al., 2020) is given as $p(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})$. The reverse generative process follows a stochastic differential equation (SDE) (Song et al., 2021b) governed by the score function $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$. To estimate this score function, one typically uses epsilon matching (Ho et al., 2020)

$$\theta^* = \arg \min_{\theta} \mathbb{E} [\|\epsilon_{\theta}(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon) - \epsilon\|_2^2], \quad (1)$$

which can be shown to be equivalent to denoising score matching (Vincent, 2011), $s_{\theta}(\mathbf{x}_t) = \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) = -\frac{1}{\sqrt{1 - \bar{\alpha}_t}}\epsilon_{\theta}(\mathbf{x}_t)$. By Tweedie’s formula (Efron, 2011), one can recover the posterior mean $\hat{\mathbf{x}}(\mathbf{x}_t) = \mathbb{E}[\mathbf{x}_0|\mathbf{x}_t] = \frac{1}{\sqrt{\bar{\alpha}_t}}(\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t}\epsilon_{\theta}(\mathbf{x}_t))$. Moreover, it is common practice to train a *conditional* score function conditioned on the text prompt (Ho & Salimans, 2022) with random dropping to use it flexibly, either as $\epsilon_{\theta}(\mathbf{x}_t, c)$ or $\epsilon_{\theta}(\mathbf{x}_t) := \epsilon_{\theta}(\mathbf{x}_t, \emptyset)$, where \emptyset refers to the

108 null condition. In practice, a popular way of generating images through reverse sampling is through
 109 DDIM sampling (Song et al., 2021a), where a single iteration can be written as
 110

$$111 \hat{\mathbf{x}}_{\emptyset}(\mathbf{x}_t) = (\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_{\theta}(\mathbf{x}_t, \emptyset)) / \sqrt{\bar{\alpha}_t} \quad (2)$$

$$112 \mathbf{x}_{t-1} = \sqrt{\bar{\alpha}_{t-1}} \hat{\mathbf{x}}_{\emptyset}(\mathbf{x}_t) + \sqrt{1 - \bar{\alpha}_{t-1}} \epsilon_{\theta}(\mathbf{x}_t, \emptyset), \quad (3)$$

113 where we defined $\hat{\mathbf{x}}_{\emptyset}(\mathbf{x}_t)$ as the posterior mean without the conditioning. Iterating (2) and (3)
 114 amounts to sampling from $p(\mathbf{x})$.
 115

116 Plugging the wanted condition \mathbf{c}^+ into the conditional epsilon $\epsilon_{\theta}(\mathbf{x}_t, \mathbf{c}^+)$ to sample from the con-
 117 ditional distribution $p(\mathbf{x}|\mathbf{c}^+)$ does not work well in practice, due to the guidance effect being too
 118 weak. To mitigate this downside, it is standard to use classifier-free guidance (CFG) (Ho & Salimans,
 119 2022) at sampling time. The key idea is to use

$$120 \hat{\epsilon}_{\mathbf{c}^+}^{\gamma}(\mathbf{x}_t) := \hat{\epsilon}_{\emptyset}(\mathbf{x}_t) + \gamma(\hat{\epsilon}_{\mathbf{c}^+}(\mathbf{x}_t) - \hat{\epsilon}_{\emptyset}(\mathbf{x}_t)), \quad (4)$$

121 where we defined $\hat{\epsilon}_{\mathbf{c}}(\mathbf{x}_t) := \epsilon_{\theta}(\mathbf{x}_t, \mathbf{c})^1$. Running DDIM sampling with $\hat{\epsilon}_{\mathbf{c}^+}^{\gamma}(\mathbf{x}_t)$ in the place of
 122 $\hat{\epsilon}_{\emptyset}(\mathbf{x}_t)$ leads to sampling from the gamma-powered distribution $p^{\gamma}(\mathbf{x}|\mathbf{c}^+) \propto p(\mathbf{x})p(\mathbf{c}^+|\mathbf{x})^{\gamma}$, a
 123 sharpened posterior. In this way, adherence to the condition \mathbf{c}^+ is emphasized.
 124
 125

126 2.2 GUIDED SAMPLING METHODS FOR DIFFUSION MODELS

127
 128 Another popular way of posterior sampling with diffusion models is to use guided sampling, where
 129 the noised sample \mathbf{x}_t (Chung et al., 2023; Yu et al., 2023; Song et al., 2023; Ye et al., 2024) or
 130 its posterior mean $\hat{\mathbf{x}}_{\emptyset}(\mathbf{x}_t)$ (Chung et al., 2024b; He et al., 2023) is added by appropriate guidance
 131 terms according to the gradient of a pre-defined energy function $\ell(\cdot)$ to be minimized. For instance,
 132 Decomposed Diffusion Sampling (Chung et al., 2024b) has the form of

$$133 \mathbf{x}_{t-1} = \sqrt{\bar{\alpha}_{t-1}} (\hat{\mathbf{x}}_{\emptyset} - \omega_t \nabla_{\hat{\mathbf{x}}_{\emptyset}} \ell(\hat{\mathbf{x}}_{\emptyset})) + \sqrt{1 - \bar{\alpha}_{t-1}} \hat{\epsilon}_{\emptyset}, \quad (5)$$

134 where ω_t is the step size, to generate sample \mathbf{x}_0 with small $\ell(\mathbf{x}_0)$. Chung et al. (2024a) presented
 135 that CFG can also be seen as such guided sampling, when we replace $\ell(\hat{\mathbf{x}}_{\emptyset})$ in (5) with the following
 136 loss function similar to Score Distillation Sampling (SDS) (Poole et al., 2023).
 137

$$138 \ell_{SDS}(\mathbf{x}) := \|\epsilon_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x} + \sqrt{1 - \bar{\alpha}_t} \epsilon, \mathbf{c}) - \epsilon\|_2^2 \\
 139 = \frac{\bar{\alpha}_t}{1 - \bar{\alpha}_t} \|\hat{\mathbf{x}}_{\mathbf{c}} - \mathbf{x}\|_2^2 \quad (6)$$

140 This leads to the alternative version of CFG called CFG++ (Chung et al., 2024a) with $\gamma := \frac{2\bar{\alpha}_t}{1 - \bar{\alpha}_t} \omega_t$,
 141

$$142 \mathbf{x}_{t-1} = \sqrt{\bar{\alpha}_{t-1}} (\hat{\mathbf{x}}_{\emptyset} + \gamma(\hat{\mathbf{x}}_{\mathbf{c}} - \hat{\mathbf{x}}_{\emptyset})) + \sqrt{1 - \bar{\alpha}_{t-1}} \hat{\epsilon}_{\emptyset}, \quad (7)$$

143 which is shown as a weight-rescaled version of CFG where the improvement is especially dominat-
 144 ing at the early stage of reverse sampling.
 145

146 2.3 NEGATIVE GUIDANCE IN DIFFUSION MODEL SAMPLING

147 In contrast to CFG that samples toward the \mathbf{c}^+ , *negative* guidance aims to avoid samples that meet
 148 the unwanted condition \mathbf{c}^- . Negative guidance can be utilized as itself when samples from a broad
 149 unconditional distribution, yet excluding certain unwanted features, are desired. When the uncon-
 150 ditional distribution is too broad to be meaningful (e.g., text-to-image models), negative guidance is
 151 often combined with the positive guidance of CFG.

152 **Inversely proportional distribution-based methods.** The simplest case (Gandikota et al., 2023;
 153 Koulischer et al., 2025; Armandpour et al., 2023; Liu et al., 2022), often called *negative prompt-*
 154 *ing*(NP), negates the guidance direction of CFG in (4) such as
 155

$$156 \hat{\epsilon}_{\mathbf{c}^-}^{\gamma}(\mathbf{x}_t) := \hat{\epsilon}_{\emptyset}(\mathbf{x}_t) - \gamma(\hat{\epsilon}_{\mathbf{c}^-}(\mathbf{x}_t) - \hat{\epsilon}_{\emptyset}(\mathbf{x}_t)), \quad (8)$$

157
 158
 159
 160
 161 ¹With a slight abuse of notation, we omit the dependence on \mathbf{x}_t when it is clear from context.

where the goal is to avoid sampling from c^- . Note that this corresponds to sampling from $p^{-\gamma}(\mathbf{x}|c^-) := p(\mathbf{x})/p(c^-|\mathbf{x})^\gamma$, a joint distribution inversely proportional to the posterior likelihood. When the goal is to sample from c^+ while avoiding c^- , one uses (Ban et al., 2024)

$$\begin{aligned}\hat{\epsilon}_{c^+,c^-}^\gamma &:= \hat{\epsilon}_\emptyset + \gamma(\hat{\epsilon}_{c^+} - \hat{\epsilon}_\emptyset) - \gamma(\hat{\epsilon}_{c^-} - \hat{\epsilon}_\emptyset) \\ &= \hat{\epsilon}_\emptyset + \gamma(\hat{\epsilon}_{c^+} - \hat{\epsilon}_{c^-})\end{aligned}\tag{9}$$

In both cases, pushing away from c^- is governed by the *negation* of the vector direction $(\hat{\epsilon}_{c^-} - \hat{\epsilon}_\emptyset)$.

Several works proposed empirical improvements for aggregating positive and negative concepts in (9) for the situation where the negative guidance term conflicts with the pre-given positive conditions. This includes the on-the-fly control of the negative guidance scale (Schramowski et al., 2023), or canceling the negative guidance component parallel to the CFG term (Armandpour et al., 2023).

Utilizing complementary conditions. Koulischer et al. (2025) proposed dynamic negative guidance (DNG) and proposed sampling from $p(\mathbf{x})p(\neg c^-|\mathbf{x})^\gamma$, where $\neg c^-$ is defined as a hypothetical condition such that $p(\neg c^-|\mathbf{x}) = 1 - p(c^-|\mathbf{x})$. This $\neg c^-$ can also be seen as a union of all possible input conditions except c^- , as a complementary set of c^- . Applying the Bayes rule, one can show that

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\neg c^-) = \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) - \gamma(\mathbf{x}_t, c^-)(\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|c^-) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)),\tag{10}$$

where $\gamma(\mathbf{x}_t, c^-) := \frac{p(c^-|\mathbf{x}_t)}{1-p(c^-|\mathbf{x}_t)}$, which can be approximated during the reverse diffusion process to adjust the guidance scale dynamically.

3 CONTRASTIVE CLASSIFIER-FREE GUIDANCE

3.1 MOTIVATION

The downside of the NP term in (8) can be explained in two different aspects: the sampling distribution involves the inverted probability $p(c^-|\mathbf{x})^{-\gamma}$, which quickly dominates over $p(\mathbf{x})$ as γ grows (Koulischer et al., 2025) and heavily distorts the sampling distribution off the supports of the marginal data distribution. This can also be seen from the NP term $(\hat{\epsilon}_{c^-} - \hat{\epsilon}_\emptyset)$ becoming bigger in regions far from c^- , unnecessarily pushing further away. Figure 2-(b) visualizes this issue with a toy dataset consisting of two classes, where two modes contain the blue class and the other mode is mixed with blues and reds. When NP sampling was done with the red class, the output distribution avoided the red mode but the other two modes were entirely shifted from the original location. The ideal behavior of negative guidance is that the strength of the guidance term decreases to zero as \mathbf{x}_t becomes irrelevant to the given condition, being equivalent to unconditional sampling.

Meanwhile, sampling from the complementary condition $\neg c^-$ (Koulischer et al., 2025) resolved this issue by weighting the NP guidance term. However, the biggest limitation is that its sample faithfully avoids c^- only when $p(c^-|\mathbf{x}_t) \approx 1$ for all \mathbf{x}_t that satisfies c^- . If the conditions are not mutually exclusive, the output distribution still includes the data that agrees with c^- , as the samples are still retained in the third node as a blue class in Figure 2-(b). This limitation is especially crucial for text-to-image (T2I) models (e.g. StableDiffusion (Rombach et al., 2022)) which take a wide range of overlapping prompts, and therefore each prompt has low $p(c|\mathbf{x})$.

3.2 DERIVATION OF CCFG

Similar to redefining positive CFG as a guided sampling that assimilates the denoising direction with the conditioned noise (Chung et al., 2024a; Kim et al., 2025b), we implement sampling that negates certain conditions with an appropriate objective function. Considering this task as optimizing the data in the sampling process closer to or further away from a given condition, we propose using contrastive loss (Hadsell et al., 2006; Chen et al., 2020) as the objective function for positive/negative prompt guidance. Contrastive loss assimilates features with the same context while distancing semantically unrelated features. Among various contrastive loss concepts, we found that Noise Contrastive Estimation (NCE) (Gutmann & Hyvärinen, 2010) best suits our situation.

Specifically, NCE parameterizes the data by performing logistic regression to discriminate the observed data from the noise data distribution. With a modeled data distribution $p_\theta(\mathbf{x})$ parameterized

by θ and pre-defined noise distribution $q(\mathbf{x})$, the logit of \mathbf{x} for the logistic regression is defined as

$$l_\theta(\mathbf{x}) := \log p_\theta(\mathbf{x}) - \log q(\mathbf{x}), \quad (11)$$

which formulates the NCE loss as

$$\begin{aligned} \mathcal{L}_{NCE} &:= -y \log \sigma(l_\theta(\mathbf{x})) - (1-y) \log(1 - \sigma(l_\theta(\mathbf{x}))) \\ &= -y \log \frac{p_\theta(\mathbf{x})}{p_\theta(\mathbf{x}) + q(\mathbf{x})} - (1-y) \log \frac{q(\mathbf{x})}{p_\theta(\mathbf{x}) + q(\mathbf{x})} \end{aligned} \quad (12)$$

where y is 1 if \mathbf{x} is the real data and 0 if \mathbf{x} is sampled from $q(\mathbf{x})$. Although NCE mostly updates the model parameter with a given positive or noise data, we optimize a datapoint with a pre-trained model that nicely parameterizes $p_\theta(\mathbf{x})$.

When we guide the denoising of the noised data \mathbf{x}_t to satisfy \mathbf{c} , we set $q(\mathbf{x}) := p(\mathbf{x}_{t-1}|\mathbf{x}_t, \emptyset)$ and parameterize $p_\theta(\mathbf{x})$ with the pre-trained diffusion model $p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{c})$ to optimize ϵ with the corresponding NCE loss:

$$\ell^+(\epsilon) := -\log \frac{p_\theta(\hat{\mathbf{x}}_{t-1}(\epsilon))}{p_\theta(\hat{\mathbf{x}}_{t-1}(\epsilon)) + q(\hat{\mathbf{x}}_{t-1}(\epsilon))} \quad (13)$$

Here, $\hat{\mathbf{x}}_{t-1}(\epsilon)$ is a mean prediction of \mathbf{x}_{t-1} from \mathbf{x}_t with ϵ . DDIM (Song et al., 2021a) states a closed form representation of $p(\mathbf{x}_{t-1}|\mathbf{x}_t)$ with a given noise prediction ϵ as a Gaussian distribution $\mathcal{N}(\boldsymbol{\mu}(\mathbf{x}_t, \epsilon), \sigma_t^2)$ with a certain choice of variance σ_t^2 and

$$\boldsymbol{\mu}(\mathbf{x}_t, \epsilon) = \frac{\sqrt{\bar{\alpha}_{t-1}}}{\sqrt{\bar{\alpha}_t}} \mathbf{x}_t - \left(\frac{\sqrt{\bar{\alpha}_{t-1}}\sqrt{1-\bar{\alpha}_t}}{\sqrt{\bar{\alpha}_t}} - \sqrt{1-\bar{\alpha}_{t-1}-\sigma_t^2} \right) \epsilon \quad (14)$$

Therefore, $p(\mathbf{x}_{t-1}|\mathbf{x}_t, \emptyset) = \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}_t, \hat{\epsilon}_\emptyset), \sigma_t^2)$ and $p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{c}) = \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}_t, \hat{\epsilon}_\mathbf{c}), \sigma_t^2)$. By replacing $\hat{\mathbf{x}}_{t-1}(\epsilon)$ with $\boldsymbol{\mu}(\mathbf{x}_t, \epsilon)$, the loss function in (13) can be simplified as

$$\ell^+(\epsilon) = -\log \frac{e^{-\|\boldsymbol{\mu}(\mathbf{x}_t, \epsilon) - \boldsymbol{\mu}(\mathbf{x}_t, \hat{\epsilon}_\mathbf{c})\|_2^2 / 2\sigma_t^2}}{e^{-\|\boldsymbol{\mu}(\mathbf{x}_t, \epsilon) - \boldsymbol{\mu}(\mathbf{x}_t, \hat{\epsilon}_\mathbf{c})\|_2^2 / 2\sigma_t^2} + e^{-\|\boldsymbol{\mu}(\mathbf{x}_t, \epsilon) - \boldsymbol{\mu}(\mathbf{x}_t, \hat{\epsilon}_\emptyset)\|_2^2 / 2\sigma_t^2}} \quad (15)$$

$$= -\log \frac{e^{-\tau_t \|\epsilon - \hat{\epsilon}_\mathbf{c}\|_2^2}}{e^{-\tau_t \|\epsilon - \hat{\epsilon}_\mathbf{c}\|_2^2} + e^{-\tau_t \|\epsilon - \hat{\epsilon}_\emptyset\|_2^2}} \quad (16)$$

Here, τ_t is a coefficient that depends on the noise scheduling and the selection of σ_t , working as the ‘‘temperature’’ on the logit (Chen et al., 2020). Although τ_t can be defined as timestep-dependent based on the diffusion reverse process, we found that setting τ_t as a constant is sufficient to yield stable and preferable results, while being simple to implement. A more detailed analysis of the selection of τ_t can be found in Appendix B.3.

If we take the derivative of $\ell^+(\epsilon)$ with respect to ϵ at $\epsilon = \hat{\epsilon}_\emptyset$ and update ϵ with step size γ_t , we can obtain the positive guidance term:

$$\hat{\epsilon}_{\mathbf{c}, \ell^+}^{\gamma_t, \tau_t} := \hat{\epsilon}_\emptyset - \gamma_t \nabla_\epsilon \ell^+(\hat{\epsilon}_\emptyset) \quad (17)$$

$$= \hat{\epsilon}_\emptyset + \frac{2\omega_t}{1 + e^{-\tau_t \|\hat{\epsilon}_\emptyset - \hat{\epsilon}_\mathbf{c}\|_2^2}} (\hat{\epsilon}_\mathbf{c} - \hat{\epsilon}_\emptyset) \quad (18)$$

Here, $\omega_t := \tau_t \gamma_t$ governs the overall guidance strength, which works similarly to the guidance scale in conventional CFG.

Similarly, to utilize the same objective function to avoid \mathbf{c} , we set $q(\mathbf{x})$ and $p_\theta(\mathbf{x})$ the same as above but treat \mathbf{x}_t as a sample from $q(\mathbf{x})$. When we perform gradient descent on $\hat{\epsilon}_\emptyset$ to minimize the following loss term, we get

$$\ell^-(\epsilon) := -\log \frac{q(\hat{\mathbf{x}}_{t-1}(\epsilon))}{p_\theta(\hat{\mathbf{x}}_{t-1}(\epsilon)) + q(\hat{\mathbf{x}}_{t-1}(\epsilon))} \quad (19)$$

$$= -\log \frac{e^{-\tau_t \|\epsilon - \hat{\epsilon}_\emptyset\|_2^2}}{e^{-\tau_t \|\epsilon - \hat{\epsilon}_\mathbf{c}\|_2^2} + e^{-\tau_t \|\epsilon - \hat{\epsilon}_\emptyset\|_2^2}} \quad (20)$$

Algorithm 1 DDIM deterministic sampling with CCFG

Require: positive prompt c^+ , negative prompt c^- , $\{\omega_t\}_{t=1}^T > 0$, $\{\tau_t\}_{t=1}^T > 0$

- 1: Initialize $x_t \sim \mathcal{N}(0, \mathbf{I})$
- 2: **for** $i = T$ **to** 1 **do**
- 3: $\lambda^+ = \frac{2}{1 + e^{-\tau_t \|\hat{\epsilon}_\emptyset(x_t) - \hat{\epsilon}_{c^+}(x_t)\|^2}}$ **if** $c^+ \neq \text{None}$ **else** 0
- 4: $\lambda^- = \frac{2e^{-\tau_t \|\hat{\epsilon}_\emptyset(x_t) - \hat{\epsilon}_{c^-}(x_t)\|^2}}{1 + e^{-\tau_t \|\hat{\epsilon}_\emptyset(x_t) - \hat{\epsilon}_{c^-}(x_t)\|^2}}$ **if** $c^- \neq \text{None}$ **else** 0
- 5: $\hat{\epsilon}_c^{\omega_t}(x_t) = \hat{\epsilon}_\emptyset(x_t) + \omega_t \lambda^+ (\hat{\epsilon}_{c^+}(x_t) - \hat{\epsilon}_\emptyset(x_t)) - \omega_t \lambda^- (\hat{\epsilon}_{c^-}(x_t) - \hat{\epsilon}_\emptyset(x_t))$
- 6: $\hat{x}_c^{\omega_t}(x_t) = (x_t - \sqrt{1 - \alpha_t} \hat{\epsilon}_c^{\omega_t}(x_t)) / \sqrt{\alpha_t}$
- 7: $x_t = \sqrt{\alpha_{t-1}} \hat{x}_c^{\omega_t}(x_t) + \sqrt{1 - \alpha_{t-1}} \hat{\epsilon}_c^{\omega_t}(x_t)$
- 8: **end for**
- 9: **return** x_t

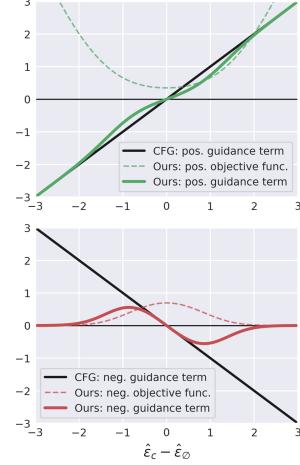


Figure 3: Left: A detailed algorithm for diffusion model sampling with CCFG. Right: The plot of the utilized objective function and its guidance term, in positive (up) and negative guidance (down). For visual simplicity, we considered one-dimensional data and fixed τ_t to 1.0.

which results in the following update:

$$\hat{\epsilon}_{c, \ell^-}^{\gamma_t, \tau_t} := \hat{\epsilon}_\emptyset - \frac{2\omega_t e^{-\tau_t \|\hat{\epsilon}_\emptyset - \hat{\epsilon}_c\|_2^2}}{1 + e^{-\tau_t \|\hat{\epsilon}_\emptyset - \hat{\epsilon}_c\|_2^2}} (\hat{\epsilon}_c - \hat{\epsilon}_\emptyset). \quad (21)$$

Explicit derivation of (18) and (21) can be found in Appendix B.2.

Algorithm 1 congregates the two scenarios and performs deterministic² DDIM sampling to satisfy or remove the given condition c . Compared to the conventional CFG, an additional weighting coefficient was added to the guidance term of each step. We plot the objective functions and their guidance term of CCFG in the right side of Figure 3 as the difference between $\hat{\epsilon}_\emptyset$ and $\hat{\epsilon}_c$ changes. As $\|\hat{\epsilon}_c - \hat{\epsilon}_\emptyset\|$ increases, the positive guidance term approximates a linear function as the original CFG. On the other hand, the negative contrastive loss and its guidance term flatten out to 0 as $\|\hat{\epsilon}_c - \hat{\epsilon}_\emptyset\|$ gets bigger, while the negative guidance of NP diverges. This demonstrates the instability of the NP term, and how CCFG can resolve this issue by canceling the gradient guidance when $\hat{\epsilon}_\emptyset$ from the current sample x_t is sufficiently unrelated to the condition. The behavior of the objective function is also changed by the choice of τ_t , whose effect on the changes of the guidance term is described in Appendix B.3.

Overall, we list the advantages of the proposed CCFG as follows: First, it faithfully removes unwanted concepts compared to complementary condition-based methods such as DNG. Concurrently, contrary to NP, it minimizes the influence on the sampling odds of unrelated data. Finally, the negative guidance of CCFG is a stand-alone method, while many empirical adjustments of NP require the presence of positive prompts to be applicable Armandpour et al. (2023); Schramowski et al. (2023).

4 EXPERIMENTAL RESULTS

We tested how CCFG steers the sample to satisfy or exclude certain conditions, along with its effect on the sample quality. We then thoroughly compared its performance to the following baselines: For scenarios where only one of the positive and negative conditions is used, we tested the conventional CFG, CFG++ (Chung et al., 2024a), NP, and DNG. In more complex scenarios involving both positive and negative conditions, we also tested Perp-Neg (Armandpour et al., 2023), SLD (Schramowski et al., 2023) and SAFREE (Yoon et al., 2024) as empirical improvements of NP in the presence of positive conditions. We conducted experiments across a range of scenarios where CFG is applicable, from simple class-conditioned image generation to text-to-image (T2I) generation.

²The selected variance σ_t for τ_t is only used to obtain the guidance term, and it's not required to perform DDIM sampling in a non-deterministic manner with σ_t .

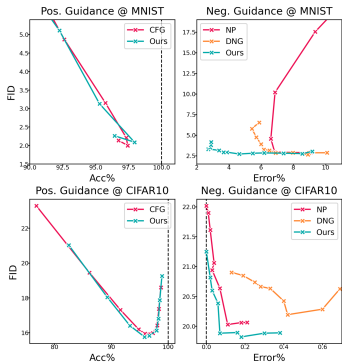


Figure 4: Left: The plot between the classification accuracy and FID with class-conditional and class-removal sampling, on MNIST and CIFAR10. Right: Quantitative results for positive and negative guidance on ImageNet-1k.

4.1 CLASS-CONDITIONED MODELS

Positive-only and negative-only guidance. To show a clear effect of concept negation and its trade-off between sample quality, we first employed MNIST (LeCun et al., 1998) and CIFAR10 (Krizhevsky et al., 2009), which comprise only 10 classes (e.g., 10% chance of unwanted class in unconditional sampling) where each class is clearly distinctive from one another. We sampled 10,000 images with a positive or negative guidance for a randomly selected class, and calculated their classification accuracy to quantify the alignment to the given positive or negative class. The classification was done by a separate external classifier. Fréchet Inception Distance (FID) (Heusel et al., 2017) between the sampled images and the real data was measured to quantify the overall sampled image quality. See Appendix D.2 for more details of the experimental settings.

Figure 4 shows the curve of the classification rate against FID with varying guidance scales. For the positive guidance, the curve of CCFG almost lies on the curve of the original CFG, indicating that positive CCFG behaves similarly to the widely-used guidance term. The curve of the negative CCFG is located at the bottom left in both MNIST and CIFAR10, outperforming the NP and DNG. In contrast to NP, which corrupts the image quality and DNG, which fails to thoroughly remove the unwanted classes, CCFG achieve similar or better precision with higher sample quality.

Guidance with positive and negative prompts. Beyond datasets with categorically distinct classes (e.g., MNIST, CIFAR10), large-scale class-conditional settings such as ImageNet-1k (Deng et al., 2009) contain many fine-grained pairs whose features are strongly entangled (e.g., *malalmute* vs. *siberian husky*). This motivates evaluating guidance schemes for a task that (i) pull samples toward a desired class while (ii) repelling them from a closely related, potentially confounding class.

We curated 100 fine-grained class pairs with strong visual proximity (c^+ , c^-) from ImageNet-1k (See Appendix D.2 for the construction) and generated 100 images for each pair, applying positive guidance toward c^+ while applying negative guidance to c^- . We report (i) Positive accuracy: fraction of samples classified as c^+ by an external classifier; (ii) Negation ratio: fraction not classified as c^- ; and (iii) Overall accuracy: the product of the two, simultaneously capturing the alignment and exclusion. To quantify perceptual quality independent of class labels, we report the mean Aesthetic Score (Schuhmann, 2022) over all generated samples. For reference against standard positive-only guidance, we also measure positive accuracy by sampling 100 random classes and generating 100 images per class with positive guidance only.

The table in Figure 4 summarizes the trade-off between class alignment, successful negation, and perceptual quality. CCFG attains the highest *overall accuracy* across pairs, indicating the most faithful target-class alignment while reliably excluding the forbidden class. NP, whose inverted weighting can over-repel, shows degradation of the positive accuracy and the perceptual quality. DNG, whose usage of complementary-conditions can under-repel when classes overlap, shows the lowest negation ratio. Perp-Neg, while demonstrating the highest negation ratio, suffers minimum positive accuracy compared to SLD, DNG and CCFG. Meanwhile, CCFG obtains the highest positive accuracy and near-perfect negation without sacrificing perceptual quality, demonstrating that CCFG’s



Figure 5: Image samples from StableDiffusion 1.5 using positive prompts (written in green) and negative prompts (written in red), across various negative sampling methods. The samples in the same column were sampled from the same initial noise. The regions where the negative prompt hasn't fully removed are highlighted in blue, and the regions where they failed to follow the positive prompt faithfully are highlighted in yellow.

StableDiffusion 1.5													
Method	GenEval[↑]							DPGBench[↑]					
	single_obj.	two_obj.	count	colors	position	color_attr.	overall	attr.	entity	global	relation	other	overall
CFG	0.959	0.301	0.388	0.723	0.028	0.045	0.407	0.734	0.731	0.839	0.794	0.556	0.635
CFG++	0.991	0.434	0.290	0.739	0.047	0.055	0.426	0.722	0.720	0.769	0.810	0.552	0.630
CCFG	0.986	0.422	0.284	0.782	0.058	0.058	0.432	0.737	0.730	0.830	0.800	0.560	0.642

SDXL													
Method	GenEval[↑]							DPGBench[↑]					
	single_obj.	two_obj.	count	colors	position	color_attr.	overall	attr.	entity	global	relation	other	overall
CFG	0.980	0.682	0.459	0.859	0.128	0.225	0.555	0.807	0.827	0.854	0.871	0.668	0.751
CFG++	0.972	0.691	0.453	0.872	0.128	0.225	0.556	0.797	0.821	0.872	0.875	0.644	0.744
CCFG	0.984	0.717	0.441	0.886	0.123	0.240	0.565	0.810	0.822	0.863	0.873	0.600	0.752

Table 1: Quantitative results on GenEval and DPGBench. **Bold**: best performance.

contrastive weighting intrinsically attenuates negative guidance when the sample is already far from c^- , preventing unnecessary distributional shift.

4.2 TEXT-CONDITIONED MODELS

Now we test the performance of CCFG in the context of T2I generation, where the conditioning signal is much complex and intertwined. The main results were obtained from StableDiffusion 1.5(SD1.5) and SDXL (Podell et al., 2024), where additional experiments on a flow-matching model of StableDiffusion 3 (Esser et al., 2024) are in Appendix B.4. See Figure 1 and Figure 11 for the examples of T2I generation scenarios with CCFG using positive and negative prompts.

Guidance with positive prompts only. Tab. 1 lists the performance of CFG and CCFG on GenEval and DPGBench, the text-to-image generation benchmarks with given text prompts to follow. CCFG showed similar or better performance to CFG in most of the subtasks, and achieved a slightly higher overall score in both benchmarks and a larger model scale on SDXL.

Guidance with positive and negative prompts. Figure 5 illustrates the result of different negative sampling algorithms from SD1.5 with several hand-crafted positive and negative prompt pairs. We observed that due to the characteristic of the overlapping text condition, the calculated guidance scale of DNG was insufficient to faithfully remove the concept in the negative prompt(e.g., cane in case (e)). The exaggerated negative guidance of NP induces sample bias(e.g., negative prompt “strawberries” removed *any* red fruit in case (b)) and even hampers the affinity to the positive prompt(e.g.,

StableDiffusion 1.5											
Method	stand-alone	positive prompt[↑]				negative prompt[↓]				overall[↑]	
		CLIP	IR	HPSv2	VLM	CLIP	IR	HPSv2	VLM	VLM _{all}	Aesthetic
NP	✓	0.2612	0.2967	0.2572	0.7981	0.1257	-2.0176	0.1637	0.0849	0.7303	5.5309
DNG	✓	0.2702	0.5794	0.2610	0.8454	0.1408	-1.8813	0.1753	0.1117	0.7484	5.5947
Perp-Neg		0.2724	0.6219	<u>0.2642</u>	0.8462	0.1364	-1.9037	0.1798	0.0943	0.7664	5.6017
SLD		0.2687	0.5640	0.2622	<u>0.8485</u>	0.1346	-1.9184	0.1715	0.0931	<u>0.7695</u>	<u>5.6053</u>
SAFREE		<u>0.2658</u>	<u>0.4329</u>	<u>0.2595</u>	<u>0.8272</u>	<u>0.1341</u>	<u>-1.9551</u>	<u>0.1709</u>	<u>0.0909</u>	<u>0.7470</u>	<u>5.5969</u>
CCFG	✓	<u>0.2703</u>	<u>0.5791</u>	0.2695	0.8540	<u>0.1339</u>	-1.9356	<u>0.1709</u>	<u>0.0894</u>	0.7777	5.6211

SDXL											
Method	stand-alone	positive prompt[↑]				negative prompt[↓]				overall[↑]	
		CLIP	IR	HPSv2	VLM	CLIP	IR	HPSv2	VLM	VLM _{all}	Aesthetic
NP	✓	0.2665	0.9658	0.2851	0.8902	0.1313	-1.9254	0.1704	0.1093	0.7929	5.9190
DNG	✓	0.2745	1.1652	0.2936	0.9040	0.1456	-1.7628	0.1842	0.1386	0.7787	5.9891
Perp-Neg		0.2735	1.1919	<u>0.2976</u>	0.9158	0.1384	-1.8324	0.1792	0.1192	0.8066	6.0155
SLD		0.2733	1.1743	0.2957	0.9185	0.1392	-1.8179	0.1802	0.1150	0.8128	6.0125
SAFREE		<u>0.2729</u>	<u>1.1831</u>	<u>0.2962</u>	<u>0.9160</u>	<u>0.1389</u>	<u>-1.8529</u>	<u>0.1781</u>	<u>0.1182</u>	<u>0.8077</u>	<u>5.9987</u>
CCFG	✓	<u>0.2739</u>	<u>1.1840</u>	0.2978	<u>0.9164</u>	<u>0.1381</u>	-1.8234	0.1790	0.1139	<u>0.8120</u>	6.0158

Table 2: Quantitative results on PIE-Bench-Neg. **Bold**: best performance, underline: second best.

“tree” and “traveler” in cases (d) and (e)). We also observed that the quality of the images can be heavily degraded, with unrecognizable contents and artifacts for cases (a) and (c). While Perp-Neg and SLD can integrate negative and positive guidance to reduce conflicts with the positive prompt, they may still fail to suppress unwanted features without careful hyperparameter search. CCFG has successfully eliminated the concept of negative prompts while maintaining positive prompts and image quality. See Appendix B.4 for a quantitative analysis of these case studies.

For a more thorough evaluation beyond case studies, we tested the behavior of CCFG and other baselines on carefully designed benchmarks. Nonetheless, very few benchmarks provide or evaluate the output image with both positive and negative prompts; most focus on how positive prompts are suppressed by the opposing negative prompts (Schramowski et al., 2023). Therefore, we designed to build a negative prompt that “*overlaps but doesn’t conflict*” with its corresponding positive prompt, based on source and target prompt pairs from known image editing benchmarks. Specifically, for given source and target prompt pairs in PIE-Bench (Ju et al., 2024), we compare the two prompts and utilize their mismatching objects, features, or style. This disagreement is utilized as a negative prompt, and its corresponding positive prompt is set to one of the two original prompts that doesn’t contain the disagreement. We refer to this dataset as *PIE-Bench-Neg*, where a more detailed construction process can be found in Appendix D.3.

We used various text-alignment metrics of CLIP cosine similarity (Radford et al., 2021), ImageReward (Xu et al., 2024), and HPS-v2 (Wu et al., 2023) to evaluate the faithfulness to the given positive or negative prompts. However, these metrics have the limitation of conflating image quality, prompt alignment, and human preference in a single value (e.g., the score with negative prompts can decrease by not the removal of unwanted features, but simply image corruption). To address this limitation, we adopt the VLM-as-a-judge framework (Lee et al., 2024; Chen et al., 2024) to perceptually evaluate the prompt alignment on a scale of 0 to 1 (See Appendix D.3 for details). To solely quantify the quality of the image, we measured the Aesthetic Score of the samples.

Tab. 2 shows the performance of CCFG and other baselines on our PIE-Bench-Neg benchmark. DNG shows the highest alignment scores with the negative prompt due to weak negative guidance, and NP shows degradation in both alignment with the positive prompt and the Aesthetic score. The drawback of these two methods led to lower VLM_{all}, which is the accuracy for both positive prompt alignment and negative prompt removal according to the VLM. Meanwhile, CCFG had a better trade-off between two different guidance, and achieves more output images that perceptually satisfy the given conditions with a better image quality; overall, it showed comparable or better performance with empirical baselines of Perp-Neg and SLD. We conclude that CCFG stands as a flexible and competitive method in the T2I generation task, on top of its broader applicability for negative-only guidance as discussed in Section 4.1.

5 CONCLUSION

In this work, we proposed CCFG, a sampling method to guide diffusion samples to satisfy or repel the given condition by optimizing the denoising direction with contrastive loss. Through the experiments on various datasets, we showed that CCFG resembles the positive guidance of the tra-

ditional CFG and resolves the downsides of the previous negative guidance methods, resulting in high-quality samples while faithfully aligning or avoiding specific attributes. Despite its effective computational overload and performance, one limitation could be the absence of an analytic closed-form sampling distribution corresponding to the proposed CCFG sampling. It would be a possible future work to propose negative guidance that allows a probabilistic interpretation while preserving the suggested advantages of CCFG. We believe that CCFG can be easily integrated and benefit various downstream applications and modalities that require negative sampling.

REFERENCES

- Mohammadreza Armandpour, Ali Sadeghian, Huangjie Zheng, Amir Sadeghian, and Mingyuan Zhou. Re-imagine the negative prompt algorithm: Transform 2d diffusion into 3d, alleviate janus problem and beyond. *arXiv preprint arXiv:2304.04968*, 2023.
- Yuanhao Ban, Ruochen Wang, Tianyi Zhou, Minhao Cheng, Boqing Gong, and Cho-Jui Hsieh. Understanding the impact of negative prompts: When and how do they take effect? In *European Conference on Computer Vision*, pp. 190–206. Springer, 2024.
- Michael Bostock. Imagenet hierarchy, 2019. URL <https://observablehq.com/@mbostock/imagenet-hierarchy>.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pp. 1597–1607. PMLR, 2020.
- Zhaorun Chen, Yichao Du, Zichen Wen, Yiyang Zhou, Chenhang Cui, Zhenzhen Weng, Haoqin Tu, Chaoqi Wang, Zhengwei Tong, Qinglan Huang, et al. Mj-bench: Is your multimodal reward model really a good judge for text-to-image generation? *arXiv preprint arXiv:2407.04842*, 2024.
- Hyungjin Chung, Jeongsol Kim, Michael Thompson Mccann, Marc Louis Klasky, and Jong Chul Ye. Diffusion posterior sampling for general noisy inverse problems. In *International Conference on Learning Representations*, 2023.
- 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*, 2024a.
- Hyungjin Chung, Suhyeon Lee, and Jong Chul Ye. Decomposed diffusion sampler for accelerating large-scale inverse problems. In *The Twelfth International Conference on Learning Representations*, 2024b.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Prafulla Dhariwal and Alexander Quinn Nichol. Diffusion models beat GANs on image synthesis. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, 2021.
- Yilun Du, Conor Durkan, Robin Strudel, Joshua B Tenenbaum, Sander Dieleman, Rob Fergus, Jascha Sohl-Dickstein, Arnaud Doucet, and Will Sussman Grathwohl. Reduce, reuse, recycle: Compositional generation with energy-based diffusion models and mcmc. In *International conference on machine learning*, pp. 8489–8510. PMLR, 2023.
- Bradley Efron. Tweedie’s formula and selection bias. *Journal of the American Statistical Association*, 106(496):1602–1614, 2011.
- Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first international conference on machine learning*, 2024.

- 540 Weichen Fan, Amber Yijia Zheng, Raymond A Yeh, and Ziwei Liu. Cfg-zero*: Improved classifier-
541 free guidance for flow matching models. *arXiv preprint arXiv:2503.18886*, 2025.
- 542
- 543 Rohit Gandikota, Joanna Materzynska, Jaden Fiotto-Kaufman, and David Bau. Erasing concepts
544 from diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer*
545 *Vision*, pp. 2426–2436, 2023.
- 546 Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle
547 for unnormalized statistical models. In *Proceedings of the thirteenth international conference on*
548 *artificial intelligence and statistics*, pp. 297–304. JMLR Workshop and Conference Proceedings,
549 2010.
- 550
- 551 Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant
552 mapping. In *2006 IEEE computer society conference on computer vision and pattern recognition*
553 *(CVPR’06)*, volume 2, pp. 1735–1742. IEEE, 2006.
- 554 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
555 nition. *CoRR*, abs/1512.03385, 2015. URL <http://arxiv.org/abs/1512.03385>.
- 556
- 557 Yutong He, Naoki Murata, Chieh-Hsin Lai, Yuhta Takida, Toshimitsu Uesaka, Dongjun Kim, Wei-
558 Hsiang Liao, Yuki Mitsufuji, J Zico Kolter, Ruslan Salakhutdinov, et al. Manifold preserving
559 guided diffusion. *arXiv preprint arXiv:2311.16424*, 2023.
- 560 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter.
561 Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances*
562 *in neural information processing systems*, 30, 2017.
- 563
- 564 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint*
565 *arXiv:2207.12598*, 2022.
- 566
- 567 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*
568 *Neural Information Processing Systems*, 33:6840–6851, 2020.
- 569 Aapo Hyvärinen and Peter Dayan. Estimation of non-normalized statistical models by score match-
570 ing. *Journal of Machine Learning Research*, 6(4), 2005.
- 571
- 572 Xuan Ju, Ailing Zeng, Yuxuan Bian, Shaoteng Liu, and Qiang Xu. Pnp inversion: Boosting
573 diffusion-based editing with 3 lines of code. *ICLR*, 2024.
- 574 Jeongsol Kim, Geon Yeong Park, Hyungjin Chung, and Jong Chul Ye. Regularization by texts for
575 latent diffusion inverse solvers. In *The Thirteenth International Conference on Learning Repre-*
576 *sentations*, 2025a.
- 577
- 578 Jeongsol Kim, Geon Yeong Park, and Jong Chul Ye. Dreamsampler: Unifying diffusion sampling
579 and score distillation for image manipulation. In *European Conference on Computer Vision*, pp.
580 398–414. Springer, 2025b.
- 581 Felix Koulischer, Johannes Deleu, Gabriel Raya, Thomas Demeester, and Luca Ambrogioni. Dy-
582 namic negative guidance of diffusion models. In *The Thirteenth International Conference on*
583 *Learning Representations*, 2025.
- 584
- 585 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
586 2009.
- 587
- 588 Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to
589 document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 590 Seongyun Lee, Seungone Kim, Sue Park, Geewook Kim, and Minjoon Seo. Prometheus-vision:
591 Vision-language model as a judge for fine-grained evaluation. In Lun-Wei Ku, Andre Martins,
592 and Vivek Srikumar (eds.), *Findings of the Association for Computational Linguistics: ACL 2024*,
593 pp. 11286–11315, Bangkok, Thailand, August 2024. Association for Computational Linguistics.
doi: 10.18653/v1/2024.findings-acl.672.

- 594 Nan Liu, Shuang Li, Yilun Du, Antonio Torralba, and Joshua B Tenenbaum. Compositional visual
595 generation with composable diffusion models. In *European Conference on Computer Vision*, pp.
596 423–439. Springer, 2022.
- 597 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of
598 the IEEE/CVF international conference on computer vision*, pp. 4195–4205, 2023.
- 600 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe
601 Penna, and Robin Rombach. SDXL: Improving latent diffusion models for high-resolution image
602 synthesis. In *The Twelfth International Conference on Learning Representations*, 2024.
- 603 Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d
604 diffusion. In *The Eleventh International Conference on Learning Representations*, 2023.
- 606 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
607 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
608 models from natural language supervision. In *International conference on machine learning*, pp.
609 8748–8763. PMLR, 2021.
- 610 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
611 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Con-
612 ference on Computer Vision and Pattern Recognition*, pp. 10684–10695, 2022.
- 613 Patrick Schramowski, Manuel Brack, Björn Deiseroth, and Kristian Kersting. Safe latent diffusion:
614 Mitigating inappropriate degeneration in diffusion models. In *Proceedings of the IEEE/CVF
615 Conference on Computer Vision and Pattern Recognition*, pp. 22522–22531, 2023.
- 617 Christoph Schuhmann. Laion-aesthetics. [https://laion.ai/blog/
618 laion-aesthetics/](https://laion.ai/blog/laion-aesthetics/), 2022. Accessed: 2023-11-10.
- 619 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *Interna-
620 tional Conference on Learning Representations*, 2021a.
- 622 Jiaming Song, Qinsheng Zhang, Hongxu Yin, Morteza Mardani, Ming-Yu Liu, Jan Kautz, Yongxin
623 Chen, and Arash Vahdat. Loss-guided diffusion models for plug-and-play controllable generation.
624 In *International Conference on Machine Learning*, pp. 32483–32498. PMLR, 2023.
- 625 Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben
626 Poole. Score-based generative modeling through stochastic differential equations. In *9th Interna-
627 tional Conference on Learning Representations, ICLR, 2021b*.
- 628 Gemma Team. Gemma 3 technical report. *arXiv:2503.19786*, 2025.
- 629 Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural compu-
630 tation*, 23(7):1661–1674, 2011.
- 631 Jing Wu, Trung Le, Munawar Hayat, and Mehrtash Harandi. Erasediff: Erasing data influence in
632 diffusion models, 2024.
- 633 Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li.
634 Human preference score v2: A solid benchmark for evaluating human preferences of text-to-
635 image synthesis. *arXiv preprint arXiv:2306.09341*, 2023.
- 636 Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao
637 Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation.
638 *Advances in Neural Information Processing Systems*, 36, 2024.
- 642 Haotian Ye, Haowei Lin, Jiaqi Han, Minkai Xu, Sheng Liu, Yitao Liang, Jianzhu Ma, James Y Zou,
643 and Stefano Ermon. Tfg: Unified training-free guidance for diffusion models. *Advances in Neural
644 Information Processing Systems*, 37:22370–22417, 2024.
- 645 Jaehong Yoon, Shoubin Yu, Vaidehi Patil, Huaxiu Yao, and Mohit Bansal. Safree: Training-free and
646 adaptive guard for safe text-to-image and video generation. *arXiv preprint arXiv:2410.12761*,
647 2024.

648 Jiwen Yu, Yinhuai Wang, Chen Zhao, Bernard Ghanem, and Jian Zhang. Freedom: Training-free
649 energy-guided conditional diffusion model. In *Proceedings of the IEEE/CVF International Con-*
650 *ference on Computer Vision*, pp. 23174–23184, 2023.
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

Algorithm 2 Deterministic sampling on flow-matching models with CCFG**Require:** positive prompt c^+ , negative prompt c^- , $\{\omega_t\}_{t=0}^1 > 0$, $\{\tau_t\}_{t=0}^1 > 0$

- 1: Initialize $\mathbf{x}_t \sim \mathcal{N}(0, \mathbf{I})$
- 2: **for** $i = 0$ **to** 1 **do**
- 3: $\lambda^+ = \frac{2}{1+e^{-\tau_t \|\hat{\mathbf{v}}_{\varnothing}(\mathbf{x}_t) - \hat{\mathbf{v}}_{c^+}(\mathbf{x}_t)\|^2}}$ **if** $c^+ \neq \text{None}$ **else** 0
- 4: $\lambda^- = \frac{2e^{-\tau_t \|\hat{\mathbf{v}}_{\varnothing}(\mathbf{x}_t) - \hat{\mathbf{v}}_{c^-}(\mathbf{x}_t)\|^2}}{1+e^{-\tau_t \|\hat{\mathbf{v}}_{\varnothing}(\mathbf{x}_t) - \hat{\mathbf{v}}_{c^-}(\mathbf{x}_t)\|^2}}$ **if** $c^- \neq \text{None}$ **else** 0
- 5: $\hat{\mathbf{v}}_c^{\omega_t}(\mathbf{x}_t) = \hat{\mathbf{v}}_{\varnothing}(\mathbf{x}_t) + \omega_t \lambda^+ (\hat{\mathbf{v}}_{c^+}(\mathbf{x}_t) - \hat{\mathbf{v}}_{\varnothing}(\mathbf{x}_t)) - \omega_t \lambda^- (\hat{\mathbf{v}}_{c^-}(\mathbf{x}_t) - \hat{\mathbf{v}}_{\varnothing}(\mathbf{x}_t))$
- 6: $\mathbf{x}_t = \mathbf{x}_t + \hat{\mathbf{v}}_c^{\omega_t}(\mathbf{x}_t) dt$
- 7: **end for**
- 8: **return** \mathbf{x}_t

Algorithm 3 DDIM deterministic sampling with CCFG guidance on the posterior mean**Require:** $\{\rho_t\}_{t=1}^T > 0$, $\{\tau_t\}_{t=1}^T > 0$, positive prompt c^+ , negative prompt c^-

- 1: Initialize $\mathbf{x}_t \sim \mathcal{N}(0, \mathbf{I})$
- 2: **for** $i = T$ **to** 1 **do**
- 3: $\lambda^+ = \frac{2}{1+e^{-\tau_t \|\hat{\epsilon}_{\varnothing}(\mathbf{x}_t) - \hat{\epsilon}_{c^+}(\mathbf{x}_t)\|^2}}$ **if** $c^+ \neq \text{None}$ **else** 0
- 4: $\lambda^- = \frac{2e^{-\tau_t \|\hat{\epsilon}_{\varnothing}(\mathbf{x}_t) - \hat{\epsilon}_{c^-}(\mathbf{x}_t)\|^2}}{1+e^{-\tau_t \|\hat{\epsilon}_{\varnothing}(\mathbf{x}_t) - \hat{\epsilon}_{c^-}(\mathbf{x}_t)\|^2}}$ **if** $c^- \neq \text{None}$ **else** 0
- 5: $\hat{\epsilon}_c^{\rho_t}(\mathbf{x}_t) = \hat{\epsilon}_{\varnothing}(\mathbf{x}_t) + \rho_t \lambda^+ (\hat{\epsilon}_{c^+}(\mathbf{x}_t) - \hat{\epsilon}_{\varnothing}(\mathbf{x}_t)) - \rho_t \lambda^- (\hat{\epsilon}_{c^-}(\mathbf{x}_t) - \hat{\epsilon}_{\varnothing}(\mathbf{x}_t))$
- 6: $\hat{\mathbf{x}}_c^{\rho_t}(\mathbf{x}_t) = (\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \hat{\epsilon}_c^{\rho_t}(\mathbf{x}_t)) / \sqrt{\bar{\alpha}_t}$
- 7: $\mathbf{x}_{t-1} = \sqrt{\bar{\alpha}_{t-1}} \hat{\mathbf{x}}_c^{\rho_t}(\mathbf{x}_t) + \sqrt{1 - \bar{\alpha}_{t-1}} \hat{\epsilon}_{\varnothing}(\mathbf{x}_t)$
- 8: **end for**
- 9: **return** \mathbf{x}_t

A CCFG ON FLOW-MATCHING MODELS

Our derivation of CCFG leverages the distribution of the next denoised sample with different conditions, which came from the stochastic sampling process. On the other hand, recent text-to-image models such as StableDiffusion 3 (Esser et al., 2024) are constructed within the flow-matching framework, where generation is modeled as the solution of an Ordinary Differential Equation (ODE) driven by a learned velocity field. Unlike diffusion models based on stochastic forward–reverse processes, flow-matching models employ deterministic dynamics of (22) over the time interval $[0, 1]$ that transport Gaussian noise toward the data distribution.

$$d\mathbf{x}_t = \mathbf{v}_{\theta}(\mathbf{x}_t, t) dt, \quad \mathbf{x}_0 \sim \mathcal{N}(0, \mathbf{I}). \quad (22)$$

That said, according to the Fokker-Planck equation, there exists a set of stochastic differential equations that has the same marginal distribution $p(\mathbf{x}_t)$ as (22),

$$d\mathbf{x}_t = \left(\mathbf{v}_{\theta}(\mathbf{x}_t, t) + \frac{t\sigma_t^2}{2(1-t)} \left(\mathbf{v}_{\theta}(\mathbf{x}_t, t) - \frac{1}{t} \mathbf{x}_t \right) \right) dt + \sigma_t d\mathbf{B}_t, \quad \mathbf{x}_0 \sim \mathcal{N}(0, \mathbf{I}), \quad (23)$$

where σ_t is an arbitrary choice of time-dependent diffusion coefficient. The stochasticity in (23) enables the same analogy that we derived CCFG in diffusion models in Section 3.2. Algorithm 2 describes the deterministic sampling process on flow-matching models with positive and negative prompt guidances by CCFG.

We’ve conducted the same experiments as those done in Section 4.2, where we tested CCFG on GenEval, DPGBenchmark, and PIE-Bench-Neg. Note that since other baselines in the main manuscript were not proposed for the flow-matching model, we compared CCFG with a naive CFG, NP, and CFG-zero (Fan et al., 2025) as an alternative for CFG on flow-matching models. Table 3 and Table 4 show that CCFG achieved both positive prompt alignment and negative prompt removal with a

756 similar trend in diffusion models, demonstrating the generalizability of CCFG across different model
757 types. See Figure 11 for the example of CCFG sampling on StableDiffusion 3.
758

759 B ADDITIONAL RESULTS

760 B.1 CCFG IN PERSPECTIVE OF GUIDED SAMPLING ON THE POSTERIOR MEAN

761 In the suggested algorithm of CCFG in Alg. 1, we applied the suggested contrastive loss to optimize
762 the denoising direction $\hat{\epsilon}_\emptyset(\mathbf{x}_t)$, with the same spirit of the conventional CFG that modifies $\hat{\epsilon}_\emptyset(\mathbf{x}_t)$
763 directly. This enabled us to conduct a fair comparison with the conventional CFG using an equal
764 guidance scale.
765

766 Meanwhile, recall that the motivation for constructing CCFG was to pose positive and negative
767 prompt sampling as guided sampling similar to CFG++ (Chung et al., 2024b), which optimizes the
768 posterior mean $\hat{\mathbf{x}}_\emptyset(\mathbf{x}_t)$ with the iteration of (7). In order to follow this scheme, one would have to
769 use the null noise $\hat{\epsilon}_\emptyset(\mathbf{x}_t)$ in the renoising process, as shown in lines 9~10 of Alg. 3.
770

771 Here, we show that these two different implementations can be equivalent with appropriate guidance
772 scale scheduling. Considering a single DDIM reverse sampling step, the obtained \mathbf{x}_{t-1} from \mathbf{x}_t by
773 Alg. 3 can be written as follows:
774

$$775 \mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t + \left(\sqrt{1 - \bar{\alpha}_{t-1}} - \frac{\sqrt{1 - \bar{\alpha}_t}}{\sqrt{\alpha_t}} \right) \hat{\epsilon}_\emptyset - \rho_t \lambda \frac{\sqrt{1 - \bar{\alpha}_t}}{\sqrt{\alpha_t}} (\hat{\epsilon}_c - \hat{\epsilon}_\emptyset). \quad (24)$$

776 In contrast, in Alg. 1, the iteration reads
777

$$778 \mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t + \left(\sqrt{1 - \bar{\alpha}_{t-1}} - \frac{\sqrt{1 - \bar{\alpha}_t}}{\sqrt{\alpha_t}} \right) \hat{\epsilon}_\emptyset - \omega_t \lambda \left(\sqrt{1 - \bar{\alpha}_{t-1}} - \frac{\sqrt{1 - \bar{\alpha}_t}}{\sqrt{\alpha_t}} \right) (\hat{\epsilon}_c - \hat{\epsilon}_\emptyset). \quad (25)$$

781 Therefore, by setting
782

$$783 \rho_t = \omega_t \left(1 - \frac{\sqrt{\alpha_t} \sqrt{1 - \bar{\alpha}_{t-1}}}{\sqrt{1 - \bar{\alpha}_t}} \right), \quad (26)$$

784 we can show the equivalence between Alg. 1 and Alg. 3. Hence, for a direct and fair comparison of
785 CCFG against conventional CFG, we chose CCFG with Alg. 1 as our implementation.
786

787 B.2 DERIVATION OF THE GUIDANCE TERM FROM CCFG OBJECTIVE FUNCTION

788 Note that
789

$$790 f(\mathbf{x}; \mathbf{a}, \mathbf{b}, \tau_t) := -\log \frac{e^{-\tau_t \|\mathbf{x} - \mathbf{a}\|_2^2}}{e^{-\tau_t \|\mathbf{x} - \mathbf{a}\|_2^2} + e^{-\tau_t \|\mathbf{x} - \mathbf{b}\|_2^2}} \quad (27)$$

$$791 \Leftrightarrow \nabla_{\mathbf{x}} f = -\frac{2\tau_t(\mathbf{a} - \mathbf{b})e^{-\tau_t \|\mathbf{x} - \mathbf{b}\|_2^2}}{e^{-\tau_t \|\mathbf{x} - \mathbf{a}\|_2^2} + e^{-\tau_t \|\mathbf{x} - \mathbf{b}\|_2^2}} \quad (28)$$

792 When we put ($\mathbf{x} = \hat{\epsilon}_\emptyset$, $\mathbf{a} = \hat{\epsilon}_c$, $\mathbf{b} = \hat{\epsilon}_\emptyset$) and ($\mathbf{x} = \hat{\epsilon}_\emptyset$, $\mathbf{a} = \hat{\epsilon}_\emptyset$, $\mathbf{b} = \hat{\epsilon}_c$) into (28), we get the
793 positive and negative CCFG guidance term for (18) and (21), respectively.
794

795 B.3 ANALYSIS ON THE SELECTION OF CCFG HYPERPARAMETERS

796 In the process of simplifying (15) into (16), the explicit expression of τ_t in terms of σ_t can be
797 obtained as
798

$$800 \frac{\|\boldsymbol{\mu}(\mathbf{x}_t, \boldsymbol{\epsilon}) - \boldsymbol{\mu}(\mathbf{x}_t, \hat{\epsilon}_c)\|_2^2}{2\sigma_t^2} = \frac{\left(\frac{\sqrt{\alpha_{t-1}}\sqrt{1-\bar{\alpha}_t}}{\sqrt{\alpha_t}} - \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2} \right)^2}{2\sigma_t^2} \|\boldsymbol{\epsilon} - \hat{\epsilon}_c\|_2^2 \quad (29)$$

$$801 := \tau_t \|\boldsymbol{\epsilon} - \hat{\epsilon}_c\|_2^2 \quad (30)$$

$$802 \Leftrightarrow \tau_t := \frac{\left(\frac{\sqrt{\alpha_{t-1}}\sqrt{1-\bar{\alpha}_t}}{\sqrt{\alpha_t}} - \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2} \right)^2}{2\sigma_t^2} \quad (31)$$

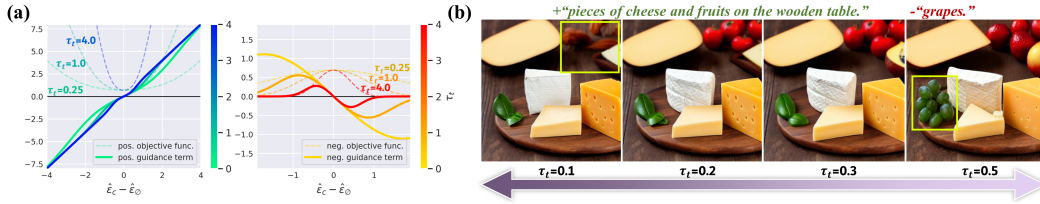


Figure 6: (a) The objective function and their guidance term of CCFG with different selection of τ_t , in positive guidance (left) and negative guidance (right). The guidance scale $\hat{\epsilon}_c^{\omega_t}$ is fixed to 1.0. (b) The examples of CCFG sampling with different τ_t , using SD1.5 and the same initial noise.

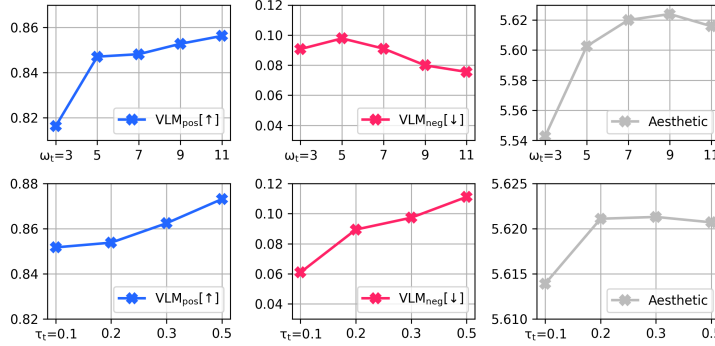


Figure 7: Performance of CCFG on PIE-Bench-Neg with different selection of ω_t (top row) and τ_t (bottom row).

where $0 \leq \sigma_t^2 \leq 1 - \bar{\alpha}_{t-1}$.

Figure 6-(a) shows the objective functions and the corresponding guidance terms of CCFG as τ_t changes. In any τ_t , the guidance term for positive CCFG approximates a linear function for the original CFG, when $\|\hat{\epsilon}_c - \hat{\epsilon}_\emptyset\|$ is sufficiently large or near 0. The guidance coefficient λ for the positive prompt becomes constant when τ_t is 0 or diverges to infinity, making positive CCFG equivalent to the original CFG. In the case of the negative guidance, λ becomes 0 when τ_t diverges to infinity, and no negative guidance is applied. As τ_t goes to 0, λ again becomes constant, making negative CCFG equivalent to NP. In general, the choice of τ_t affects more on the negative CCFG since it governs the range of $\|\hat{\epsilon}_c - \hat{\epsilon}_\emptyset\|$ in which the negative guidance is applied sufficiently.

Figure 6-(b) illustrates this behavior in T2I generation with SD1.5, where too large τ_t disables proper negative guidance and too small τ_t introduces the drawback of NP(*e.g.* applying too much negative guidance disrupts following the positive prompt).

Figure 7 shows the performance of CCFG on PIE-Bench-Neg with different selection of ω_t and τ_t . The positive prompt alignment(VLM_{pos}) and negative prompt removal(VLM_{neg}) improve as the guidance scale ω_t increases, while the overall quality of the image(Aesthetic Score) has a sweet spot at a suitable ω_t . This trend agrees with a well-known observation of traditional CFG. Meanwhile, as illustrated in Figure 6-(b), a smaller τ_t can more thoroughly remove the negative concept, but also hinders the positive prompt alignment.

B.4 TEXT-TO-IMAGE SAMPLING WITH POSITIVE AND NEGATIVE PROMPTS

Figure 8 shows the quantitative results on five T2I case studies in Figure 5 with CLIP cosine similarity score, ImageReward, HPS-v2, and perceptual accuracy judged by VLM. Here, we show the average score over the five cases for each method. We sampled 200 images for each method and case. The ideal behavior as an effective and safe negative sampling method is to maintain alignment with positive prompts while decreasing alignment with negative prompts, compared to when no negative guidance is applied.

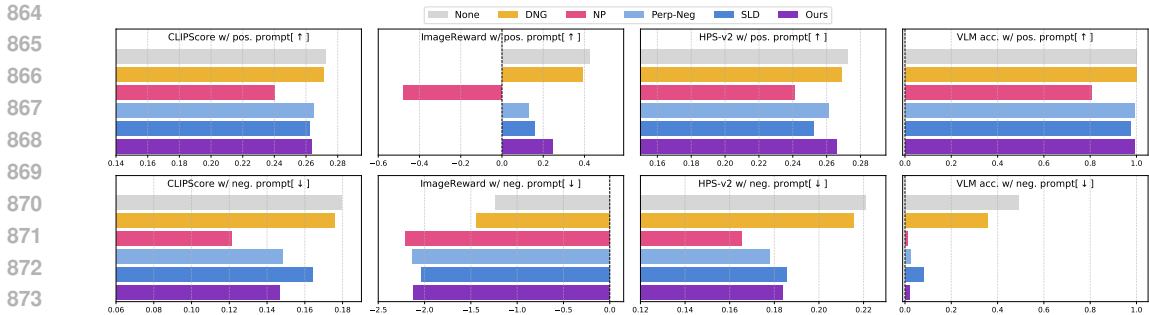


Figure 8: Quantitative evaluation results over the case scenarios in Figure 5. “None” denotes cases where only the positive guidance is applied.

StableDiffusion 3													
Method	GenEval[↑]							DPGBench[↑]					
	single_obj.	two_obj.	count	colors	position	color_attr.	overall	attr.	entity	global	relation	other	overall
CFG	0.984	0.740	0.559	0.823	0.237	0.449	0.632	0.861	0.878	0.844	0.910	0.780	0.825
CCFG	0.988	0.748	0.575	0.840	0.250	0.463	0.644	0.865	0.883	0.839	0.906	0.816	0.827

Table 3: Quantitative results with StableDiffusion 3 on GenEval and DPGBench. **bold**: best performance.

DNG showed the smallest alignment metric decrease between its samples and the given positive prompt, but failed to thoroughly rule out the features in the negative prompt; the alignment with the negative prompt was marginally decreased, and VLM often detected the attributes described in the negative prompt. Meanwhile, NP showed the lowest alignment score not only with the negative prompts, but also for the positive prompts. Also, almost 20% of its output was judged to not satisfy the positive prompt by VLM. While Perp-Neg and SLD can orchestrate the negative guidance with the positive guidance to prevent the misalignment between the positive prompt, they also occasionally fail to remove the unwanted features without careful hyperparameter search. Our proposed CCFG removes the undesired properties enough to be undetected by VLM, and still maintains the agreement with the positive prompts.

In Figure 11, we display additional image generation examples by CCFG and given positive-negative prompt pairs from StableDiffusion 1.5, SDXL, and StableDiffusion 3. Leveraging the ability of CCFG for an effective and stable concept negation, we can resolve certain undesirable features such as objects, motions, and multiple composed features. Furthermore, CCFG can help the model generate minority samples and remove potentially harmful objects by negating potential internal bias or inappropriate contents.

B.5 POSSIBLE FAILURE CASES

Figure 9 presents several failure cases of CCFG on text-to-image generation tasks; when the negative condition involves multiple attributes simultaneously, CCFG might suppresses only a subset of the undesired aspects. And when the undesired negative prompt describes subtle continuous attributes, the negative guidance strength of CCFG may be weakened in the process of concept removal, resulting in outputs on the borderline of unwanted features.

C THE BEHAVIOR OF NP ON DIFFERENT DIFFUSION SAMPLING SCENARIOS

So far, we’ve shown the theoretical pitfall of the NP and demonstrated its negative effect on the sample distribution and its quality. Meanwhile, we observed that the criticality of NP’s drawback has been mitigated as the learned data distribution and the nature of the used condition get more complex: the NP term completely shifts the sample distribution out from the marginal support in the toy example of Figure 2, but often produced reasonable images for T2I generation tasks. Indeed, many off-the-shelf implementations for the T2I generation model have used this guidance to this day for concept negation (Gandikota et al., 2023; Koulischer et al., 2025; Armandpour et al., 2023; Liu et al., 2022).

StableDiffusion 3										
Method	positive prompt[↑]				negative prompt[↓]				overall[↑]	
	CLIP	IR	HPSv2	VLM	CLIP	IR	HPSv2	VLM	VLM _{all}	Aesthetic
NP	0.2616	1.0302	0.2813	0.9415	0.1335	-1.9418	0.1722	0.1135	0.8346	5.7214
CCFG	0.2722	1.2976	0.2891	0.9630	0.1316	-1.7148	0.1958	0.1185	0.8489	5.7751

Table 4: Quantitative results with StableDiffusion 3 on PIE-Bench-Neg. **bold**: best performance.



Figure 9: Examples of failure cases on text-to-image generation with positive and negative prompts via CCFG.

We insist that this phenomenon occurs for the following reasons: First, as the possible conditions become more diverse and their relationship with the data distribution becomes more complex, it becomes harder for the model to accurately parameterize the entire unconditional and conditional distribution. It is well known that many other T2I models perform relatively poorly for pure unconditional generation, and extra distribution sharpening with CFG is almost essential for acceptable image quality. When this happens, even an unbounded guidance term can enhance the sample, until it overshoots and introduces artifacts. For more simplified scenarios like MNIST, where the model can accurately learn the target data distribution, the downside of NP becomes more evident. The second reason comes from $p(c^-|\mathbf{x})^{-\gamma}$, the factor that NP introduced to the sampling distribution. In a class-conditioned dataset where each class doesn't overlap too much, $p(c^-|\mathbf{x})$ can be close to 1 if the given image is sampled from class c^- . This contrasts $p(c^-|\mathbf{x})^{-\gamma}$ with different \mathbf{x} significantly and assigns a high probability for the region with low $p(c^-|\mathbf{x})$, even for \mathbf{x} with low marginal probability. On the other hand, a single image can be described by numerous different texts, therefore the magnitude of $p(c|\mathbf{x})$ for an image-text pair is small no matter how well c describes the given image.

Nonetheless, we believe that the fundamental flaws of NP have been already demonstrated with various experiments, and a risk of failure still remains for T2I models.

D EXPERIMENTAL DETAILS

D.1 SAMPLING HYPERPARAMETERS

Tab. 5 lists the hyperparameter values used for DNG and CCFG for the results in the main manuscript. Apart from the guidance scale ω that is shared by all negative sampling methods we've tested, DNG requires a prior for the given condition $p(c)$, an affine transformation weight τ' and bias δ for its guidance term calculation. We first took the values from the authors of DNG (Koulisher et al., 2025) as a basis, then additionally searched and tuned for the best performance on each dataset. We note that DNG only showed acceptable performances with the heavily exaggerated condition prior $p(c)$ which are far from the reasonable value (e.g. 0.1 for MNIST and CIFAR10).

D.2 CLASS-CONDITIONED MODELS

For the experiments about CCFG on MNIST and CIFAR10, we trained DDPM (Ho et al., 2020) for each dataset from scratch which enables both conditional and unconditional generation to perform CFG, CCFG, and other sampling methods. The external classifier was imported from the following checkpoints, where the ViT-base is trained with the classification of MNIST and CIFAR10, respec-

Dataset	DNG	SLD	CCFG
MNIST	$p(c)=0.25, \tau'=0.25, \delta=0$	(N/A)	$\tau=1.0$
CIFAR10	$p(c)=0.5, \tau'=0.1, \delta=0.0002$	(N/A)	$\tau=0.25$
ImageNet-1k	$p(c)=0.01, \tau'=0.2, \delta=0.0002$	$\delta=15, s_S=200, \lambda=0.00, s_m=0.00, \beta_m=0.4$	$\tau=0.5$
StableDiffusion 1.5	$p(c)=0.01, \tau'=0.2, \delta=0.003$	$\delta=10, s_S=1000, \lambda=0.01, s_m=0.3, \beta_m=0.4$	$\tau=0.2$
SDXL	$p(c)=0.01, \tau'=0.2, \delta=0.003$	$\delta=10, s_S=1000, \lambda=0.01, s_m=0.3, \beta_m=0.4$	$\tau=0.1$
StableDiffusion 3	(N/A)	(N/A)	$\tau=0.005$

Table 5: Sampling hyperparameters used for DNG, SLD, and CCFG.

tively. The diffusion models were trained with 500 noise timesteps. For the inference, we used deterministic DDIM sampling with NFE=50.

For the experiment on ImageNet-1k, We prepared 100 class pairs with strong visual proximity (c^+, c^-) (e.g., *Maltese_dog - Shih-Tzu*, *Italian_greyhound - whippet*, *Samoyed - Pomeranian*, *racer - sports_car*, *beach_wagon - cab*, *bathing_cap - shower_cap*, *garden_spider - wolf_spider*) according to ImageNet hierarchy Bostock (2019) and the following pipeline:

- For each ImageNet-1k class, we retrieve the corresponding WordNet ID (wnid) and its hypernyms.
- We construct a pool of candidate pairs by including all pairs of classes that share an identical hypernym.
- From this pool, we uniformly sample 100 class pairs to obtain the final benchmark set, while ensuring that every class in the 100 class pairs is unique.

The sampling was done at 256×256 using an ImageNet-1k-pretrained DiT-XL/2 (Peebles & Xie, 2023) by deterministic DDIM sampling with NFE=50. For the external classifier, we used a frozen ResNet50 (He et al., 2015) pretrained on ImageNet-1k using torchvision’s default IMAGENET1K_V2 weights.

D.3 TEXT-CONDITIONED MODELS

For the SD1.5 and SDXL image sampling, we used DDIM deterministic sampling with 50 NFE and a guidance scale of 7.5. In the flow-matching model sampling with StableDiffusion 3, we used a deterministic sampling with 28 NFE and a guidance scale of 3.0. In the case of DNG and SLD, we used one of the hyperparameter configurations provided by the authors.

To build the PIE-Bench-Neg benchmark, we processed each pair of source and target prompts in PIE-Bench (Ju et al., 2024) according to the given affiliation that describes the editing task:

- `add_object`: Compared to the source prompt (e.g., “*a bee illustration on a paper*”), the newly added objects in the target prompt are listed in the attribute of `aspect_mapping` (e.g., “*leaves, flower*”). These objects become the negative prompt, and the source prompt is the positive prompt.
- `delete_object`: The objects from the source prompt (e.g., “*a wathet bird sitting on a branch of yellow flowers*”) removed in the target prompt (e.g., “*a wathet bird sitting on a branch*”) are listed in the attribute of `aspect_mapping` (e.g., “*flowers*”). These objects become the negative prompt, and the target prompt is the positive prompt.
- `change_object`, `change_attribute_content`, `change_background`: The changed appearances from the source prompt (e.g., “*fishes and kelp in the ocean*”) are marked in the target prompt (e.g., “*[sharks] and [flowers] in the ocean*”), which becomes the negative prompt (e.g., “*sharks, flowers*”). The source prompt is used for the positive prompt.
- `change_attribute_pose`, `change_attribute_color`, `change_attribute_material`: These editing tasks require certain features of the objects in the source prompt (e.g., “*a brown bear is sleeping with a sign that says no sleep*”) to be changed, and the mapping between them can be found in key-value pairs of `aspect_mapping` (e.g., “*bear*”: “*black*”, “*sign*”: “*white*”). We combine the objects and

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

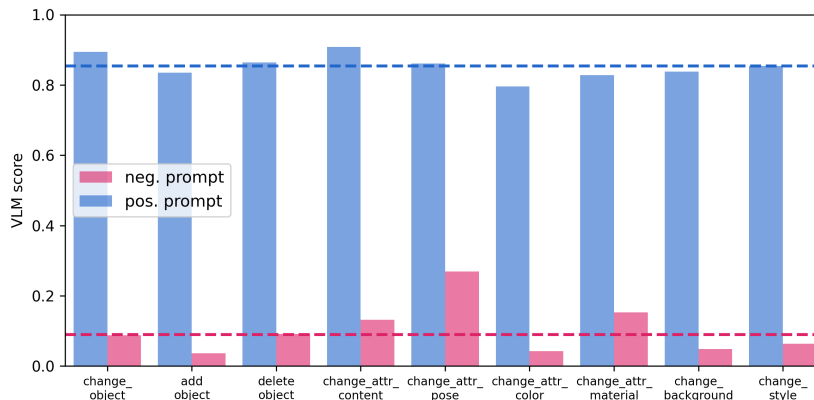


Figure 10: The performance of CCFG in PIE-Bench-Neg prompts from different PIE-Bench categories. The dotted blue and pink line represents the overall mean VLM score for positive and negative prompts, respectively.

their target editing features to construct a negative prompt(e.g., “black bear, white sign”), and the original source prompt is used as the positive prompt.

- `change_style`: The target style is described in the attribute of `blended_words`(e.g., “baroque, watercolor”). These styles become the negative prompt, and the source prompt(e.g., “a cartoon giraffe sitting down on a white background”) is the positive prompt.

In total, we prepared 560 positive prompts, paired with the corresponding negative prompt that is correlated but not entirely opposed to the positive prompt.

The performance of CCFG in PIE-Bench-Neg prompts from each of these category is shown in Figure 10. Here, we observed that CCFG is relatively less effective on the series of “change.attribute_XX” categories, which tend to yield relatively lower positive prompt VLM scores(e.g., `change_attr_color`) or higher negative prompt VLM scores(e.g., `change_attr_pose`). This can be because the negative prompts from these categories require removing specific attributes of something explicitly contained in the positive prompt, therefore conveying overlapping concepts and making the task more difficult. We also found that having nouns as negative prompts(e.g., `change_object`) or adjectives(e.g., `change_style`) shows similar performances.

To evaluate a perceptual agreement of the image with the given prompt, separated from the image quality, we employed an open-source VLM Gemma3-27B (Team, 2025). We utilize Gemma3 to decompose the given prompt and check the given image’s agreement with each component to eventually rate the image. Prompt 1 was used to evaluate the perceptual accuracy of the given image to follow the text prompt using VLM. The last line of the output is interpreted as a fraction between 0 and 1, where 1 is the highest possible agreement score. Similar to how we calculated the overall accuracy in Section 4.1, we multiply the VLM score for the positive and the negative prompt as VLM_{ou} to quantify how the model output matches for both prompt conditions.

E THE USE OF LARGE LANGUAGE MODELS (LLMs)

LLMs were not involved in research ideation, methodological design, or the writing of the manuscript. The authors retain full responsibility for all scientific content.

1080
 1081
 1082
 1083
 1084
 1085
 1086
 1087
 1088
 1089
 1090
 1091
 1092
 1093
 1094
 1095
 1096
 1097
 1098
 1099
 1100
 1101
 1102
 1103
 1104
 1105
 1106
 1107
 1108
 1109
 1110
 1111
 1112
 1113
 1114
 1115
 1116
 1117
 1118
 1119
 1120
 1121
 1122
 1123
 1124
 1125
 1126
 1127
 1128
 1129
 1130
 1131
 1132
 1133

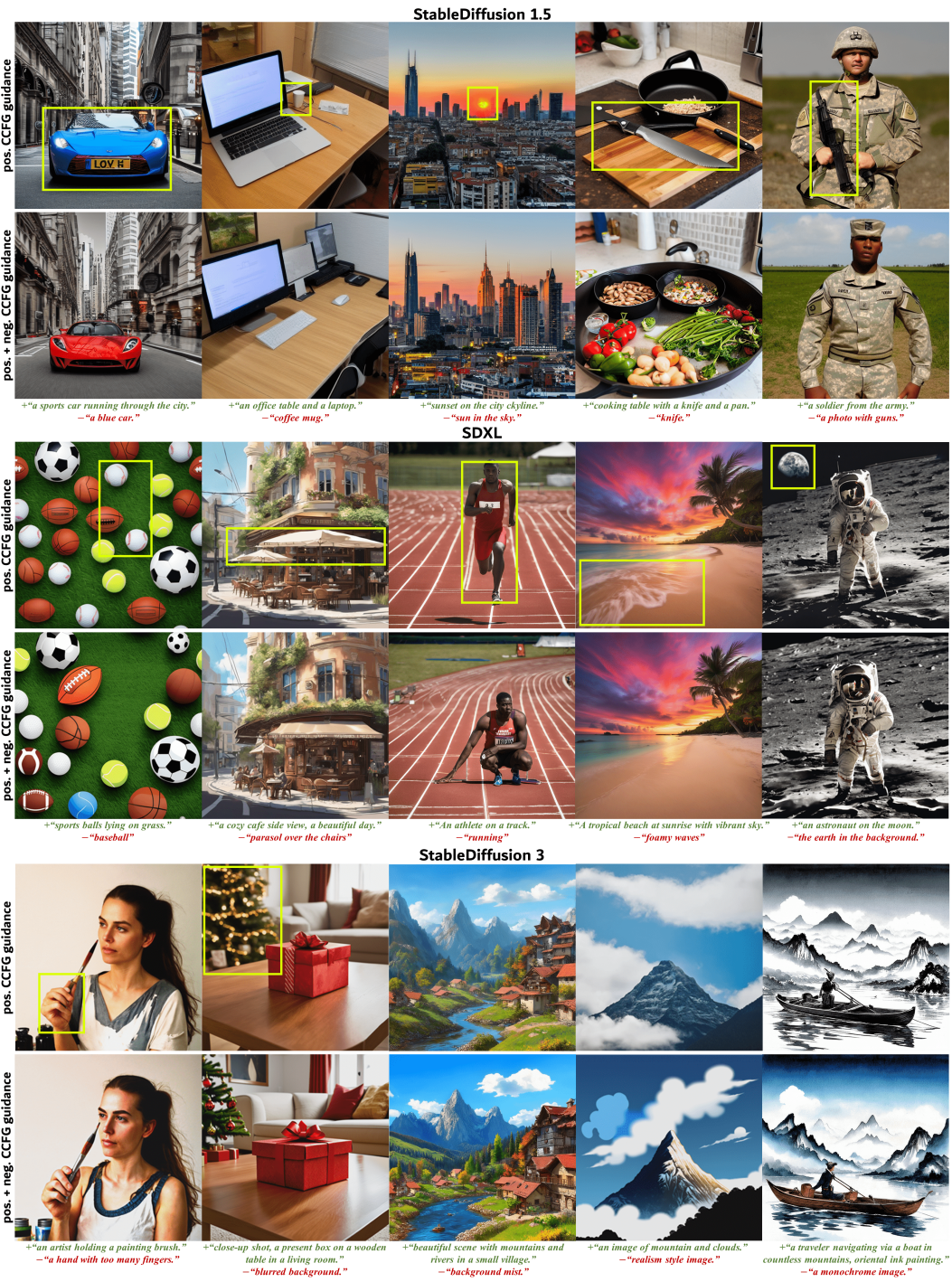


Figure 11: Examples of text-to-image generation results with positive and negative prompts via CCFG, sampled from StableDiffusion 1.5, SDXL, and StableDiffusion 3. In each model, the samples in the same column were sampled from the same initial noise.

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

prompt 1

I'll give you an image and a text prompt. You have to evaluate how well the image satisfies the prompt. In the process, you first list all the objects and appearances (with its corresponding features) described in the prompt. Don't count the word "image" in the prompt as an object. Then, you observe the image and decide whether each detected object/appearance in the prompt is contained in the image.

The last line of your response should be a format of "Final answer: a over b". Here, "a" is the number of objects that are both contained in the prompt and the image, and "b" is the total number of objects in the prompt. Here is the example:

If the prompt is "an image with a blue cat and a red dog" and the image contains a blue cat and a yellow dog, then the prompt contains "a blue cat" and "a red dog". "A blue cat" is contained in the image, and "a red dog" is not contained in the image. Therefore, the answer should be "1 over 2".

Now, the input sentence is "<caption>". Observe the given image and answer.