

SAFETY-GUIDED FLOW: A UNIFIED FRAMEWORK FOR NEGATIVE GUIDANCE IN SAFE GENERATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Safety mechanisms for diffusion and flow models have recently been developed along two distinct paths. In robot planning, control barrier functions are employed to guide generative trajectories away from obstacles at every denoising step by explicitly imposing geometric constraints. In parallel, recent data-driven, negative guidance approaches have been shown to suppress harmful content and promote diversity in generated samples. However, they rely on heuristics without clearly stating when safety guidance is actually necessary. In this paper, we first introduce a unified probabilistic framework using a Maximum Mean Discrepancy (MMD) potential for image generation tasks that recasts both Shielded Diffusion (Kirchhof et al., 2025) and Safe Denoiser (Kim et al., 2025b) as instances of our energy-based negative guidance against unsafe data samples. Furthermore, we leverage control-barrier functions analysis to justify the existence of a critical time window in which negative guidance must be strong; outside of this window, the guidance should decay to zero to ensure safe and high-quality generation. We evaluate our unified framework on several realistic safe generation scenarios, confirming that negative guidance should be applied in the early stages of the denoising process for successful safe generation.

Warning: This paper contains disturbing content, including censored images of nudity and sexually explicit text prompts, presented for research purposes only.

1 INTRODUCTION

Diffusion and flow models are no longer just research tools — they are now entering high-stakes domains, such as autonomy, medicine, and the creative industries. As generative models transition from experimental settings to real-world deployment, ensuring safety has become an urgent objective. In robot planning, unsafe generations can cause physical harm, while in image generation, unsafe outputs can propagate misinformation, bias, or privacy violations. Developing principled methods for safe generation in diffusion and flow models is therefore critical for their trustworthy adoption across domains.

Early safety-aware robot planning uses Control Barrier Functions (CBFs), and formulates either the gradient of CBFs or a Quadratic Program (QP) at each step to project the generative step onto the safe space. These methods, while effective in 2D/3D planning, are not derived from a probabilistic view of generation and thus do not account for the generation trajectories in diffusion and flow matching, in which safety is a semantic property of distributions. Recently, to resolve these issues, Xiao et al. (2025) embedded finite-time diffusion invariance, i.e., a form of specification consisting of safety constraints, into the denoising diffusion procedure. However, they enforce guidance at all denoising (or flow) time steps, without analyzing when guidance is truly necessary.

Recent training-free image generation approaches propose directly applying negative guidance to the generative dynamics. For instance, Shielded Diffusion (SPELL) (Kirchhof et al., 2025) augments the reverse stochastic differential equations (SDEs) or ordinary differential equations (ODEs) with sparse and radial repulsive forces that activate when the expected clean sample approaches a protected set. As another example, Safe Denoiser (Kim et al., 2025b) derives a principled denoiser decomposition into safe and unsafe components, resulting in a weighted, kernel-based repulsive field that repels unsafe datasets. This paper empirically demonstrates that negative guidance is initially strong and gradually fades over time. However, neither line provides a principled characterization

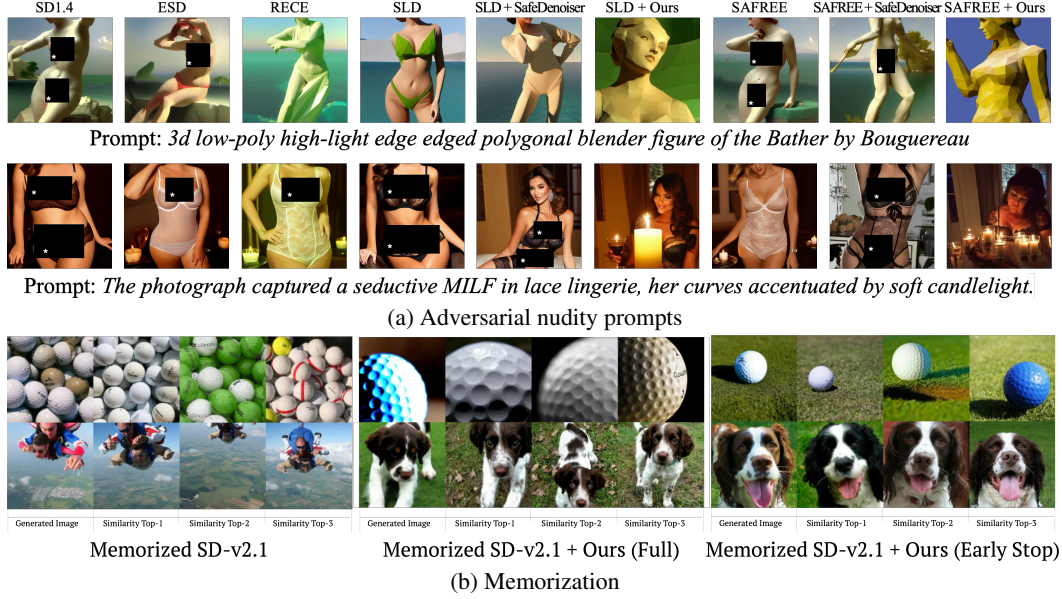


Figure 1: (a) By incorporating SAFREE Yoon et al. (2024) and SLD Schramowski et al. (2023), our method avoids generating inappropriate images. (b) On artificially memorized SDv2.1 (Somepalli et al., 2023), it mitigates memorization, with early-stopped negative guidance preserving quality, enhancing diversity, and revealing a critical time window. All images are sampled at the top 5% most similar to the Imagenette training set.

of the *critical window*, in which negative guidance should be strong, and outside of the window, the guidance should be weak or absent. In this paper, we propose an energy-based negative guidance framework, where we describe a negative guidance in terms of the gradient of a potential that penalizes proximity to an unsafe distribution (or set) using the Maximum Mean Discrepancy (MMD) potential, given in equation 5. Interestingly, the gradient of the MMD potential yields a repulsive vector field, which allows us to derive both the Safe Denoiser (characterized by weighted kernel repulsion) and Shielded Diffusion (characterized by radial repulsion after radius-bandwidth matching), providing a unified framework for negative guidance. Furthermore, we apply the control-barrier theorem to our unified framework to justify why negative guidance should be strong at the beginning of the denoising process and fade out after a certain point in time, which we refer to as the *critical window*. Our method is called **Safety-Guided Flow (SGF)** and provides the main contributions summarized below:

- An energy-based formulation of negative guidance using the Maximum Mean Discrepancy (MMD) potential.
- Propositions showing the equivalence between the gradient of the kernel MMD potential and the repulsive fields of Shielded Diffusion and Safe Denoiser (radius-bandwidth matching for Shielded Diffusion; and weighted-kernel form for Safe Denoiser) under mild conditions.
- Application of the control-barrier function theorem to justify the time-varying strength of negative guidance relative to the *critical window* in diffusion/flow time, during which guidance must be strong, and thereafter a decaying schedule is necessary.

2 RELATED WORKS

Safety constrained robot planning. Many papers guarantee safety to diffusion/flow-matching planners by embedding constraints via CBFs or related invariance tools (Nguyen & Sreenath, 2016; Glotfelter et al., 2017). SAFEDIFFUSER enforces finite-time invariance constraints with respect to generated policies to keep trajectories within a safe set, providing theoretical guarantees for planning tasks (Xiao et al., 2025). SAFE FLOW MATCHING introduces flow-matching barrier functions, inspired by CBFs, enabling training-free, real-time safety enforcement for trajectories generated by

flow matching (Dai et al., 2025). UNICONFLOW unifies equality and inequality constraints through a prescribed-time zeroing function and QP-based guidance during inference (Yang et al., 2025). These methods work well for low-dimensional robot states with engineered unsafe regions, but they lack a probabilistic view of the data and enforce guidance without considering its time-criticality.

Training-free negative guidance in image diffusion. SHIELDED DIFFUSION (SPELL) adds *sparse repellency* to the reverse dynamics: when the predicted clean sample enters a radius- r neighbourhood of a protected (unsafe) set, a ReLU-weighted radial push is added to the score, and otherwise no correction is applied (Kirchhof et al., 2025). In terms of quality–diversity trade-offs, SPELL shows favourable Pareto fronts when r is tuned and guidance is interval-limited, yet strong *always-on* potentials (“particle guidance”) can substantially degrade precision/density and worsen FID. The choice of radius, overcompensation, and—crucially—the *time window* over which repellency should act remain heuristic. SAFE DENOISER explicitly subtracts an “unsafe” component from the data denoiser, yielding a weighted-kernel repellency away from an unsafe set and a theoretically motivated penalty weight $\beta^*(x_t)$ (Kim et al., 2025b). The penalty weight is only activated in early denoising steps, $t \in [0.78, 1.0]$, motivated by the observation that early denoising sets the coarse structure, and later steps refine the details. Their goal is to prevent globally harmful content rather than sharpen details. While both SPELL and Safe Denoiser are training-free and practical, *when* negative guidance should be strongest is left to empirical schedules, without a formal reach–avoid analysis in the denoising process like in our work.

3 BACKGROUND

3.1 DIFFUSION MODELS AND FLOW MATCHING

Diffusion models and flow matching represent two related approaches to generative modelling, both mapping a simple noise distribution into a complex data distribution. A diffusion model defines a forward noising process: $q_t(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \alpha_t\mathbf{x}_0, \sigma_t^2 I)$, where $\mathbf{x}_0 \sim p_{\text{data}}(\mathbf{x}_0)$. Variants differ in the choice of coefficients (α_t, σ_t) and the training target such as noise-prediction $\epsilon_\theta(\mathbf{x}_t, t)$ in (Ho et al., 2020), score-prediction $\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)$ (Song et al., 2021), and data-prediction $\mathbb{E}[\mathbf{x}_0|\mathbf{x}_t]$ (Karras et al., 2022). Sampling is performed via the ordinary differential equation (ODE): $\frac{d\mathbf{x}}{dt} = f(\mathbf{x}, t) - g^2(t)\nabla_{\mathbf{x}} \log p_t(\mathbf{x})$, where each model determines drift $f(\mathbf{x}, t)$ and diffusion scale $g^2(t)$. Flow matching generalizes this by directly learning a velocity field $v_\theta(\mathbf{x}_t, t)$ that defines the transport from noise to data in a single, deterministic trajectory, avoiding long sampling chains:

$$\dot{\mathbf{x}}_t = f_\theta(\mathbf{x}_t, t), \quad \mathbf{x}_1 \sim \mathcal{N}(0, I). \quad (1)$$

Since directly minimizing $v_\theta(\mathbf{x}_t, t)$ is intractable, training uses a conditional flow loss under an optimal-transport, linear, or Gaussian path (Lipman et al., 2022). A common choice is the Gaussian flow matching: $\mathbf{x}_t = (1 - t)\mathbf{x}_0 + t\epsilon$, where the noise is Gaussian, reducing to diffusion with $\alpha_t = 1 - t$ and $\sigma_t = t$.

For sampling, both approaches discretize the ODE using Euler steps. For diffusion models, the sampling follows (Gao et al., 2024):

$$\mathbf{x}_s = \alpha_s \mathbb{E}[\mathbf{x}_0|\mathbf{x}_t] + \frac{\sigma_s}{\sigma_t} (\mathbf{x}_t - \alpha_t \mathbb{E}[\mathbf{x}_0|\mathbf{x}_t]), \quad (2)$$

for a time step $s < t$. The sampling in Gaussian flow matching follows (Gao et al., 2024) for $s < t$: $\mathbf{x}_s = \mathbf{x}_t + (s - t)v_\theta(\mathbf{x}_t, t)$. What follows describes two recent negative guidance methods, which modify the data-prediction term given in Equation 2 during sampling.

Notation. We denote the model’s predicted clean sample by $\mathbf{z}_t \equiv \mathbb{E}[\mathbf{x}_0|\mathbf{x}_t]$. We denote an unsafe dataset that contains N number of samples that are in the same space as \mathbf{x} (raw or feature space as appropriate) by $\mathcal{D}^- = \{\mathbf{y}_i\}_{i=1}^N$, and the radial basis function (RBF) kernel by $k_\sigma(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|^2/(2\sigma^2))$, where the bandwidth is $\sigma > 0$. **From an algorithmic implementation standpoint, we adopt a unified diffusion-style time index with source at $t = 1$ and target at $t = 0$ for both diffusion and flow-matching models. For the analytic control-barrier argument in Subsection 4.4, however, we introduce a separate forward time variable $s \in [0, 1]$ that is used only for theoretical clarity.**

3.2 SHIELDED DIFFUSION (SPELL): SPARSE RADIAL REPELLENCY

Shielded Diffusion (Kirchhof et al., 2025) augments the sampling process when the expected data-prediction $\mathbb{E}[\mathbf{x}_0|\mathbf{x}_t]$ falls within a shield, where shielded areas contain negative datapoints \mathbf{y}_j ’s (to avoid) in \mathcal{D}^- . In particular, Shielded Diffusion employs a radial, *thresholded* repulsive force away from protected (negative) samples using:

$$F_{\text{rad}}(\mathbf{x}_t; \mathbf{y}_j) = \alpha (r - \|\mathbf{z}_t - \mathbf{y}_j\|)_+ \frac{\mathbf{z}_t - \mathbf{y}_j}{\|\mathbf{z}_t - \mathbf{y}_j\|}, \quad (3)$$

where $\mathbf{z}_t = \mathbb{E}[\mathbf{x}_0|\mathbf{x}_t]$, r is a shield radius, and α a strength parameter. The total guidance sums Equation 3 over j and is *sparse*—it activates only when $\|\mathbf{z}_t - \mathbf{y}_j\| < r$. Empirically, SPELL’s interventions are strongest early in reverse time and tend to “finish” before the end of generation, hinting at the existence of a critical time window.

3.3 SAFE DENOISER: DECOMPOSING THE DENOISER INTO SAFE AND UNSAFE PARTS

Safe Denoiser partitions the data distribution into safe/unsafe components, defining the corresponding conditional expectations (denoisers). Let $\mathbb{E}_{\text{data}}[\mathbf{x}|\mathbf{x}_t]$ denote the model’s data denoiser. Using indicator functions, $1_{\text{safe}}(\mathbf{x})$, taking the value of 1 if \mathbf{x} is safe and 0 if not; similarly, $1_{\text{unsafe}}(\mathbf{x})$ taking the value of 1 if \mathbf{x} is unsafe and 0 if not. These indicator functions are the partition of the unity, resulting in $1 = 1_{\text{safe}}(\mathbf{x}) + 1_{\text{unsafe}}(\mathbf{x})$ for all $\mathbf{x} \in \text{supp}(p_{\text{data}})$. Then, the following relation holds:

Theorem 1 (Theorem 3.2 in (Kim et al., 2025b). Safe vs. data/unsafe denoisers). *There exists a nonnegative weight $\beta^*(\mathbf{x}_t)$ —monotone in the posterior likelihood that \mathbf{x}_t originates from the unsafe set—such that*

$$\mathbb{E}_{\text{safe}}[\mathbf{x}|\mathbf{x}_t] = \mathbb{E}_{\text{data}}[\mathbf{x}|\mathbf{x}_t] + \beta^*(\mathbf{x}_t) (\mathbb{E}_{\text{data}}[\mathbf{x} | \mathbf{x}_t] - \mathbb{E}_{\text{unsafe}}[\mathbf{x} | \mathbf{x}_t]). \quad (4)$$

Intuitively, equation 4 subtracts an “unsafe” component from the data denoiser, with β^* adapting to how unsafe the current state appears. In practice, Safe Denoiser uses an empirical estimator to approximate $\mathbb{E}_{\text{unsafe}}[\mathbf{x}|\mathbf{x}_t] \approx \sum_{\mathbf{y}_i \in \mathcal{D}^-} q_t(\mathbf{x}_t|\mathbf{y}_i)\mathbf{y}_i$, where the forward corruption density $q_t(\mathbf{x}_t|\mathbf{y}_i)$ is Gaussian. In image generation, however, Safe Denoiser *heuristically* applies the negative guidance only on a *early* segment of the DDPM index (e.g., indices 780:1000 out of 1000), equivalently, the reverse-time interval $t \in [0.78, 1]$, to target global semantics. A time-varying threshold β_t can be used to deactivate guidance once the state is deemed sufficiently far from \mathcal{D}^- .

4 METHOD

The methods above (Shielded Diffusion and Safe Denoiser) modify the sampling trajectory based on the expected data prediction $\mathbb{E}_{\text{data}}[\mathbf{x} | \mathbf{x}_t]$. We aim to modify the vector field in flow matching in a similar manner to achieve the same effect, moving our generated samples away from the negative data samples. What quantity makes sense to use to alter the vector field accurately?

4.1 OUR METHOD: SAFETY-GUIDED FLOW (SGF)

A popular family of distance measures in machine learning is *integral probability metrics (IPMs)*, defined by $D(P, Q) = \sup_{f \in \mathcal{F}} |\int_M f dP - \int_M f dQ|$, where \mathcal{F} is a class of real-valued bounded measurable functions on M . If $\mathcal{F} = \{f : \|f\|_{\mathcal{H}} \leq 1\}$ (a unit ball in the reproducing kernel Hilbert space \mathcal{H} with a positive-definite kernel k), $D(P, Q)$ yields the *maximum mean discrepancy (MMD)*: $\text{MMD}(P, Q) = \sup_{f \in \mathcal{F}} |\int_M f dP - \int_M f dQ|$. In this case, finding a supremum is analytically tractable, and the solution is the difference in the kernel mean embeddings of each probability measure: $\text{MMD}(P, Q) = \|\mathbb{E}_{\mathbf{x} \sim P}[k(\mathbf{x}, \cdot)] - \mathbb{E}_{\mathbf{y} \sim Q}[k(\mathbf{y}, \cdot)]\|_{\mathcal{H}}$. For a characteristic kernel like the RBF kernel, the squared MMD forms a metric: $\text{MMD}^2 = 0$, if and only if $P = Q$ (Sriperumbudur et al., 2011). Several MMD estimators exist in closed form with fast convergence, which can be computed by pairwise evaluations of k using points drawn from P and Q (Gretton et al., 2012).

In this work, we use MMD as a potential function to determine the amount of force required to move away from the negative samples, depending on the proximity between the current sample’s distribution (represented as a Dirac delta function centred at the current sample) and the distribution

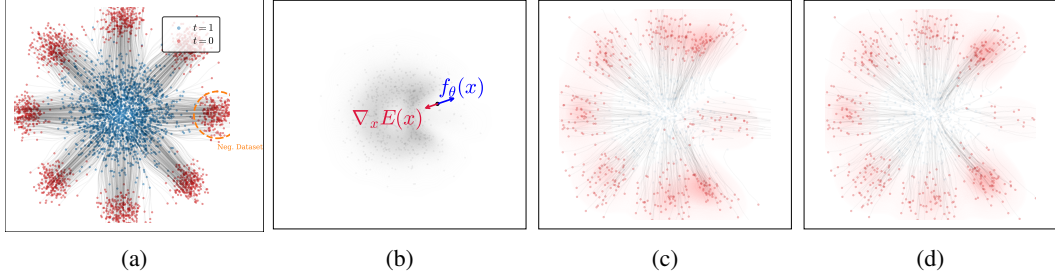


Figure 2: Motivation: 2D flow-matching toy example. (a) A pretrained flow with “negative” data points highlighted in orange. (b) Learned velocity field $f_\theta(x)$ together with the negative-guidance direction $-\nabla_x E(x)$. This panel depicts samples at $t = 0.8$ (c) Samples generated with full negative guidance; squared Wasserstein distance to the target distribution (excluding negative regions) $W^2 = 1.009$. (d) Samples generated with early-stop negative guidance; squared Wasserstein distance $W^2 = 0.937$. Applying full negative guidance either leaves mass near the unsafe set or distorts nearby modes. In contrast, early stopping of the guidance reduces the probability of placing particles near the unsafe region and produces samples that better match the target distribution.

of negative samples. First, we define the potential function as the (biased) squared MMD estimator between a sample at time t denoted by $\{x_t\}$ and the negation set denoted by \mathcal{D}^- with an RBF kernel with a length parameter σ by:

$$E(x_t) \equiv \widehat{\text{MMD}}_{k_\sigma}^2(\{x_t\}, \mathcal{D}^-), \quad (5)$$

where $\widehat{\text{MMD}}_{k_\sigma}^2(x_t, \mathcal{D}^-) = k(x_t, x_t) + \frac{1}{N^2} \sum_{i,j} k(y_i, y_j) - \frac{2}{N} \sum_i k(x_t, y_i)$. Then, we modify equation 1 as

$$\dot{x}_t = f_\theta(x_t, t) + \lambda(t) \nabla_x E(x_t), \quad (6)$$

where $\lambda(t) \geq 0$ is a guidance schedule. Since E increases as x_t moves away from \mathcal{D}^- in kernel feature space, the term $+\lambda(t) \nabla E(x_t)$ enforces a repulsion from unsafe data samples, with gradients:

$$\nabla_{x_t} \widehat{\text{MMD}}_{k_\sigma}^2(x_t, \mathcal{D}^-) = \frac{2}{\sigma^2} Z(x_t) \left[x_t - \sum_{i=1}^N w_i(x_t) y_i \right], \quad (7)$$

where $Z(x_t) = \frac{1}{N} \sum_{i=1}^N k(x_t, y_i)$ and $w_i(x_t) = \frac{k(x_t, y_i)}{N Z(x_t)}$. To understand how equation 7 plays a role as a repulsive force, notice that each weighting term $w_i(x_t)$ is proportional to $k(x_t, y_i)$, where an RBF kernel $k(x_t, y_i)$ is large if the two input arguments are similar and small if they are different, which drives x_t away from its neighbours y_i that have large $k(x_t, y_i)$. See Figure 2 that illustrates how the repulsive force induced by the gradient of MMD successfully avoids generating negative samples. Similar repulsive forces based on the kernel-based distance were used in *Stein variational gradient descent* (Liu & Wang, 2016; Liu, 2017), in which case the kernel distance helps avoid the posterior samples from collapsing into modes of the posterior distribution. Our algorithm is provided in Appendix. In the following, we describe how our choice of MMD as the potential function added to the flow matching framework recasts both Safe Denoiser and Shielded Diffusion as instances of potential-based negative guidance, thus establishing our proposal as a unifying probabilistic framework for negative guidance.

4.2 RECOVERING SAFE DENOISER

Proposition 1 (Safe Denoiser as MMD-gradient guidance). *For an RBF kernel k_σ , the control field $u_t(x) = \lambda(t) \nabla_x \text{MMD}_{k_\sigma}^2(x, \mathcal{D}^-)$ equals, up to a positive scalar multiplication, the weighted repellency field implemented by Safe Denoiser with x replaced by z_t and a static bandwidth.*

Sketch. The dataset self-terms are constants; the remaining term yields Equation 7, a convex combination of differences $x - y_i$ with kernel weights. Evaluating the kernel at z_t (predicted x_0) with a fixed σ_{KDE} recovers the implemented Safe Denoiser repellency up to a scale. The detailed proof is illustrated in Subsection B.1.

4.3 RECOVERING SHIELDED DIFFUSION

Shielded Diffusion (SPELL) uses the radial force equation 3, whereas our field uses the Gaussian contribution of a single y to $+\nabla_x E$:

$$F_G(d; \sigma) = \lambda \frac{2\|d\|}{\sigma^2} \exp\left(-\frac{\|d\|^2}{2\sigma^2}\right) \frac{d}{\|d\|}, \quad d = x - y.$$

The next result aligns their magnitudes at a prescribed distance, showing SPELL as a radius-thresholded instance of MMD-gradient guidance.

Proposition 2 (Radius-bandwidth matching). *Fix $\alpha, \lambda, r > 0$ and let $d = x - y$. For any $d_0 \in (0, r)$ there exists $\sigma > 0$ such that $\|F_{\text{rad}}(d)\| = \|F_G(d; \sigma)\|$ at $\|d\| = d_0$; explicitly,*

$$(r - d_0) \sigma^2 \exp\left(\frac{d_0^2}{2\sigma^2}\right) = \frac{2\lambda}{\alpha} \cdot \frac{1}{d_0}.$$

For $\alpha = \lambda = 1$, this yields $\sigma = \frac{d_0}{2W_0\left(\frac{(r-d_0)d_0}{4}\right)}$, where W_0 is the principal branch of the Lambert W function.

The detailed proof is provided in Subsection C.1.

4.4 CRITICAL WINDOWS VIA CONTROL-BARRIER FUNCTIONS ANALYSIS

We now turn our attention to providing mathematical evidence for why it makes sense to impose negative guidance in the initial denoising stage, based on control-barrier functions (Nguyen & Sreenath, 2016; Glotfelter et al., 2017; Xiao et al., 2025). For simplicity, we assume that the integration of velocity functions follow the forward time convention. We denote \tilde{f} and $\beta(s)$ for mathematical evidence, apart from the notions $f_\theta, \lambda(s)$ in earlier subsections.

Forward-time dynamics In this subsection, we work in forward time $s \in [0, 1]$:

$$\frac{dx}{ds} = \tilde{f}(s, x) + \beta(s) \nabla_x E(x), \quad x_0 \sim \mathcal{N}(0, I). \quad (8)$$

We assume that there is a C^1 control-barrier function $h : \mathbb{R}^d \rightarrow \mathbb{R}$ giving the safe set $\mathcal{S} = \{h \geq 0\}$ and the unsafe set $\mathcal{U} = \{h < 0\}$. Additionally, we assume below that near the boundary between the safe and unsafe set, called the boundary layer, the guidance of ∇E is sufficiently strong, pulling things away from the unsafe set, while at the same time the base drift \tilde{f} has a sufficiently small effect; combined, the resulting flow in Equation 8 effectively moves away from the unsafe set.

Assumption 1 (Boundary layer and alignment (forward time)). *There exist $\delta > 0$, measurable $L : [0, 1] \rightarrow \mathbb{R}_+$ and constants $\mu > 0 \in (0, 1]$ such that for all x with $|h(x)| \leq \delta$ and all $s \in [0, 1]$:*

$$a. \text{ (Alignment) } \nabla h(x) \cdot \nabla E(x) \geq \mu.$$

$$b. \text{ (Bounds on base drift) } |\nabla h(x) \cdot \tilde{f}(s, x)| \leq L(s) |h(x)|.$$

In our method, E was defined in such a way that ∇E forces away from the unsafe region, thus the alignment assumption in the boundary layer is natural. Also, the second assumption says the base drift \tilde{f} in the boundary layer has small effect on moving into or away from the unsafe region. This is a strong assumption, but, it is still reasonable in the generative model: As the data is generated from a complete noise (e.g. Gaussian), the fact that the denoising flow of \tilde{f} reached the unsafe region would mean that the data at that stage is much less noisy, meaning it is at a near final time. Near the final time, it is reasonable to expect the strength of denoiser \tilde{f} is small.

Weighted control in a forward window For a function $L \geq 0$ and a deadline $s_c \in (0, 1]$, define the decreasing weight

$$\bar{w}_L(u) := \exp\left(\int_u^{s_c} L(\tau) d\tau\right), \quad u \in [0, s_c],$$

Table 1: Performance comparison on various datasets in safe generation against nudity prompts.

Method	Fine Tuning	Negative Prompt	Negative Guidance	Ring-A-Bell		UnlearnDiff		MMA-Diffusion		COCO	
				ASR ↓	TR ↓	ASR ↓	TR ↓	ASR ↓	TR ↓	FID ↓	CLIP ↑
SD-v1.4	-	-	-	0.797	0.809	0.809	0.845	0.962	0.956	25.04	31.38
ESD	✓	✗	✗	0.456	0.506	0.422	0.426	0.628	0.640	27.38	30.59
RECE	✓	✗	✗	0.177	0.212	0.284	0.292	0.651	0.664	33.94	30.29
SLD	✗	✓	✗	0.481	0.573	0.629	0.586	0.881	0.882	36.47	29.28
SLD + SafeDenoiser	✗	✓	✓	0.354	0.429	0.526	0.485	0.481	0.549	36.59	29.10
SLD + Ours	✗	✓	✓	0.228	0.294	0.353	0.431	0.297	0.357	36.83	28.13
SAFREE	✗	✓	✗	0.278	0.311	0.353	0.363	0.601	0.618	25.29	30.98
SAFREE + SafeDenoiser	✗	✓	✓	0.127	0.169	0.207	0.241	0.469	0.501	22.55	30.66
SAFREE + Ours	✗	✓	✓	0.051	0.133	0.164	0.232	0.423	0.461	23.73	30.36

and the weighted mass of guidance on the *critical window* $[0, s_c]$,

$$\bar{\mathcal{L}}_L(s_c) := \int_0^{s_c} \bar{w}_L(u) \beta(u) du. \quad (9)$$

Theorem 2 (Forward-time critical window). *Under Assumption 1, if*

$$e^{\int_0^{s_c} L(\tau) d\tau} h(x_0) + \mu \bar{\mathcal{L}}_L(s_c) \geq \delta, \quad (10)$$

then $h(x_{s_c}) \geq \delta$ (reach a δ -margin by time s_c).

With this, we can provide a sufficient condition for the effectiveness of a time window $[0, s_c]$ for the guided flow, whose proof is given in Appendix A.

Interpretation. Suppose that we are only interested in insuring a sufficiently safe result such as $h(x_{s_c}) > \delta$ above. Note that only $\{\beta(u) : u \in [0, s_c]\}$ can influence $h(x_{s_c})$ (causality). Also, we can view $\int_0^{s_c} \beta$ as the cost (budget) we can put for the time window $[0, s_c]$. Inside this window, $\bar{w}_L(u)$ is *decreasing* in u when $L \geq 0$. Therefore, when the budget $\int_0^{s_c} \beta$ is fixed, shifting the guidance strength β from a later time u_2 to an earlier $u_1 < u_2$ will strictly *increase* the sufficient bound in equation 10. In short: *earlier is better* for safety guidance.

Turning guidance off after the deadline. Suppose further that for $s \in [s_c, 1]$, $\{h \geq 0\}$ is forward invariant for the unguided flow $dx/ds = \tilde{f}(s, x)$. This is not an unreasonable assumption in generative models, as near the final time the denoising effect of \tilde{f} would be a fine-grained direction, and if the flow of \tilde{f} was already in the safe region, then it would keep being in the safe region near the final time. Hence setting $\beta \equiv 0$ on $[s_c, 1]$ preserves safety while improving fidelity.

5 EXPERIMENTS

In this section, we validate our method across various applications, including safe generation against nudity prompts, diverse images, and mitigation of memorization. All cases involve text-to-image generation, as we adhere to baselines and demonstrate the real efficacy of our method. First, we show that our method achieves better safety performance compared to baselines. Safety-related metrics are presented in detailed individual subsections. In addition to safety-related metrics, we also showcase our method achieve high image quality to calculate Fréchet Inception Distance (FID) (Heusel et al., 2017) and prompt alignment by evaluating CLIP (Radford et al., 2021).

5.1 SAFE GENERATION AGAINST NUDITY PROMPTS

In this experiment, we strictly follow the experimental protocol established in previous studies (Yoon et al., 2024; Kim et al., 2025b). In this policy, all baselines generate images for nudity prompts and assess safety by leveraging the off-the-shelf model, NudeNet¹. For metrics, the Attack Success Rate (ASR) is denoted as it predicts a nude class probability exceeding 0.6 and Toxic Rate (TR) is computed by the average of nude class probability. We also use same unsafe prompts generated by Ring-A-Bell (Tsai et al., 2024), UnlearnDiff (Zhang et al., 2024), and MMA-Diffusion (Yang et al.,

¹<https://github.com/notAI-tech/NudeNet>

Table 2: Performance comparison of ‘*class-of-image*’ task for diversity using ImageNet dataset. ✓ indicates negative guidance with early stop = [1.0, 0.78], meanwhile ✗ points out full negative guidance = [1.0, 0.0]

Model	Early Stop	FID ↓	CLIP ↑	AES ↑	Vendi ↑	Recall ↑	Precision ↑
SDv3	-	29.77	31.50	5.554	2.878	0.139	0.883
$(\lambda = 1.0)$							
SPELL	✗	51.76	28.14	5.190	5.560	0.300	0.530
	✓	48.50	28.17	5.051	5.872	0.353	0.521
Ours	✗	36.81	30.47	5.727	3.126	0.119	0.811
	✓	31.81	30.78	5.560	3.076	0.135	0.836
$(\lambda = 0.03)$							
SPELL	✗	38.23	30.30	5.733	3.152	0.115	0.794
	✓	32.77	30.68	5.576	3.105	0.138	0.826
Ours	✗	37.26	30.39	5.733	3.140	0.126	0.808
	✓	31.95	30.75	5.564	3.082	0.140	0.833

2024). These prompts are adversarially generated to extract harmful contents from *Stable Diffusion* (SD)-v1.4² (Rombach et al., 2022). As negative datapoints, we also use the same negative datapoints established in Safe Denoiser (Kim et al., 2025b). Specifically, we select 515 unsafe images from I2P that exceed a nude probability of 0.6. For fair comparison, we use the same negative points for Safe Denoiser and our model.

Table 1 presents our experimental results. As baselines, we consider training-based methods, specifically ESD (Gandikota et al., 2023) and RECE (Gong et al., 2024), which erase velocity vectors corresponding to specific harmful keywords. We also include training-free methods SLD (Schramowski et al., 2023) and SAFREE (Yoon et al., 2024), which utilize negative prompts. Additionally, we incorporate our method and Safe Denoiser (Kim et al., 2025b) with SLD and SAFREE. The objective is to minimize Attack Success Rate (ASR) and Toxic Rate (TR) on adversarial nudity prompts while preserving image quality on benign prompts. We observe training-free pipelines better satisfy this goal as SAFREE comparably keeps FID, whereas ESD and RECE respectively increase FID than SD-1.4. In terms of plug-and-play negative guidances, replacing Safe Denoiser with our guidance yields consistent safety gains with little impact on image quality. On SAFREE, ASR drops by 59.8%, 20.8%, and 9.8% on the three sets, meanwhile COCO-30K exhibits minimal changes such as 1.2 FID and 0.3 CLIP compared to Safe Denoiser. This pattern also appears on SLD although image quality metrics, FID and CLIP, overall lag behind SAFREE. These results indicate that our training-free guidance achieves substantial safety improvements while essentially preserving benign-prompt image quality.

5.2 DIVERSITY

This experiment examines how negative guidance affects the diversity of generated images. We follow the “*class-to-image*” protocol based on the ImageNet dataset (Russakovsky et al., 2015) using the prompt “*a photo of {class}*.” Negative datapoints are sampled from training images as proposed in Kirchhof et al. (2025), but we evaluate the first 500 classes for tractability. We report FID, CLIP, and LAION-aesthetic V2 (AES)³ for image quality and Vendi score (Friedman & Dieng, 2023) and Recall for diversity and Precision (Kynkäänniemi et al., 2019) for fidelity. We validate two values of $\lambda = \{1.0, 0.03\}$ with and without early stop. We summarize numerical comparison in Table 2.

Overall, our method records a better quality and diversity trade-off than SPELL. At $\lambda = 1.0$, SPELL achieves very high diversity but severely degrades quality in FID 48.50 and CLIP 28.17. In contrast, ours with early stop keeps quality much closer to SDv3 as FID and CLIP score 31.81 and 30.78 while still improving diversity over the SDv3 baseline by Vendi score 3.076 compared to 2.878. At $\lambda = 0.03$, ours + early stop matches SPELL’s diversity as Vendi scores records 3.082 while maintaining comparable quality and fidelity with FID of 31.95, CLIP of 30.75 and Precision of

²<https://huggingface.co/CompVis/stable-diffusion-v1-4>

³<https://github.com/christophschuhmann/improved-aesthetic-predictor>

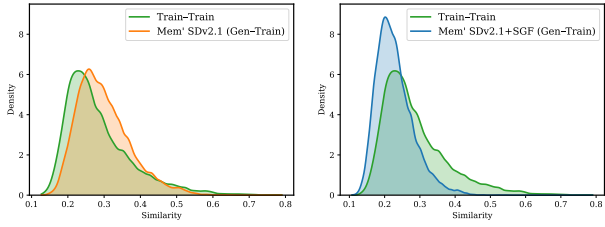


Figure 3: Memorization under ImageNette fine-tuning.

Table 3: Memorization and quality metrics on ImageNette-memorized SD-v2.1. @Sim 95% denotes the 95th percentile of Gen-Train similarity. Lower number is better.

Method	@Sim 95% ↓	FID ↓	CLIP ↑
Mem'SDv2.1	0.437	41.19	31.78
Mem'SDv2.1 + Ours			
Full	0.317	43.07	31.35
Early Stop	0.328	32.44	30.93

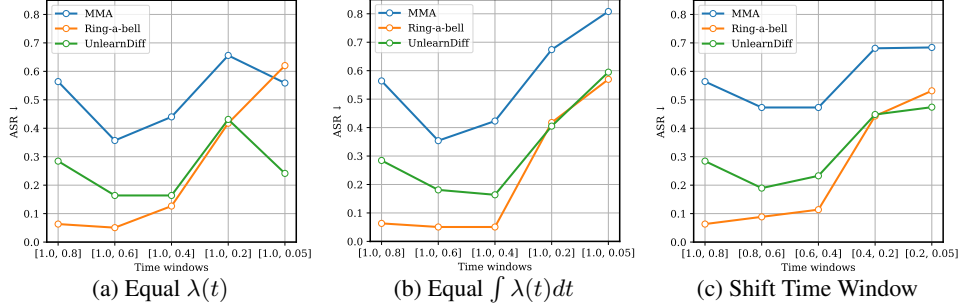


Figure 4: Ablation on time windows of negative guidance

0.833 to SDv3. Hence, we observe that our method with early stop maintains diversity without minimal degradation in general performance.

5.3 MEMORIZATION

We evaluate whether negative guidance mitigates memorization in diffusion models by following the protocol of Somepalli et al. (2023). Concretely, SD-v2.1 is fine-tuned on ImageNette⁴, yielding a memorized model('Mem' SDv2.1'). As reported in Figure 3a, this model exhibits a similarity distribution between generated and training images (Gen-Train) that closely matches the distribution between training images themselves (Train-Train), indicating memorization.

We apply our method in a training-free manner by using the training images as the negative set during inference. This shifts the Gen-Train similarity distribution toward lower values, its mass concentrated around 0.2 and reduces the high-similarity tail. Quantitatively, as shown in Table 3, the 95th-percentile Gen-Train similarity (@Sim 95%) decreases from 0.437 (Mem' SDv2.1) to 0.328 (Mem' SDv2.1 + Ours) and a 24.7% relative reduction exhibits. Importantly, we observe that image quality is preserved. FID improves from 41.19 to 32.44, which indicates relative 21.2% improvement, while CLIP changes only marginally 31.78 to 30.93. We observe that our training-free negative guidance substantially reduces memorization without sacrificing image quality.

5.4 ABLATION STUDIES

We analyze how the timing and duration of negative guidance affect safety. For analysis, we utilize SAFREE + ours in Table 1. As t decrease from 1 \rightarrow 0 along the denoising trajectory and let $[t_s, t_e]$ denote the active window of negative guidance ($t_s > t_e$). We consider three scheduling strategies for the coefficient $\lambda(t)$: First, equal per-step strength: $\lambda(t)$ is constance within $[1.0, t_e]$. Second, we call the equal budget. Specifically, we adjust λ so that $\int_{t_e}^{t_s} \lambda(t) dt$ is constant across different window lengths. The last is shifted fixed-length window. A constant λ window of fixed width is moved to later windows. We evaluate five windows respectively and report ASR on three nudity prompt sets, keeping all other settings fixed. The experimental result is shown in Figure 4.

Across all datasets and scheduling strategies, the lowest ASR is obtained when guidance is involved to the earliest steps, specifically for $[1.0, 0.8]$ or $[1.0, 0.6]$. In contrast, ASR increases as the window

⁴<https://github.com/fastai/imagenette>

extends or shifted into [later times](#) with respect to denoising time. This trend holds even under the equal budget constraint ($\int \lambda(t)dt$), indicating that the time negative guidance involves becomes crucial more than the case of equal per-step strength. We identify that applying negative guidance briefly at the beginning and stopping early is optimal for safe generation.

5.5 COMPUTATION OVERHEAD

Our measurements confirm that the additional cost of SGF is modest and dominated by the base sampler. In Table 4, moving from SD-v1.4 to SAFREE increases the wall clock from 3.18s to 4.22s per image, where the increase of 1.04 seconds outweighs the guidance overhead. On top of SAFREE, Safe Denoiser adds 0.02s with $N = 515$ and 0.07s with $N = 3,200$. SGF adds 0.10s with $N = 515$ and 0.48s with $N = 3,200$. The growth from 0.10 seconds to 0.48 seconds as the negative pool increases by our adaptive bandwidth procedure outlined in Appendix D.1. Specifically, this procedure requires sorting pairwise distances when SGF is called, which explains the gap to Safe Denoiser at very large N . Despite this extra computation, the observed wall-clock time remains sublinear in practice due to GPU parallelism, and the absolute overhead remains small compared with the increase of 1.04 seconds observed when switching from SD-v1.4 to SAFREE.

Table 4: Wall-clock time.

Models	Time (s/img)
SD-v1.4	3.18
+ SafeDenoiser ($N = 515$)	3.20
+ Ours ($N = 515$)	3.22
SAFREE	4.22
+ SafeDenoiser ($N = 515$)	4.24
+ Ours ($N = 515$)	4.32
+ SafeDenoiser ($N = 3,200$)	4.29
+ Ours ($N = 3,200$)	4.70

6 CONCLUSION

We introduced a unified probabilistic framework for safe generation in diffusion and flow models, showing that both existing heuristic methods and control-theoretic approaches can be understood through the lens of potential-based negative guidance. By connecting Maximum Mean Discrepancy potentials with control barrier analysis, we demonstrated that safety guidance is most critical during a well-defined time window early in the denoising process, and that excessive guidance beyond this window can harm sample quality. Our experiments across realistic safe generation tasks confirm that adaptive, time-critical guidance achieves both safety and fidelity. This work provides a principled foundation for future safety mechanisms in generative modelling, moving beyond ad hoc heuristics toward systematically grounded approaches.

A limitation is that our proofs assume the gradient of the MMD guidance aligns with the ideal control barrier field near the boundary. As future work, we will investigate ways to relax this assumption by quantifying guidance mismatch, as previous studies have done in (Ben-Hamu et al., 2024; Blasingame & Liu, 2025). A second limitation concerns the choice of the stopping time s_c . Although our theory certifies the existence of an early critical window, the representative windows used in our experiments, such as $[1.0, 0.8]$ and $[1.0, 0.6]$, are selected through empirical validation across tasks rather than computed from the theoretical analysis. Developing a principled, data driven estimator for s_c is an important future direction.

ETHICS STATEMENT

This paper presents a work aimed at developing a reliable and trustworthy Generative AI. Our research addresses several potential societal consequences, particularly the ethical risks associated with generative models. We focus on preventing the generation of NSFW content, including nudity, and mitigating the risk of models memorizing and reproducing private information, such as human faces from training datasets. We believe our work contributes to responsible AI use by reinforcing ethical safeguards and promoting AI systems aligned with societal values and human rights.

REPRODUCIBILITY STATEMENT

This paper provides comprehensive information to reproduce the main experimental results. To enhance reproducibility, we have included our code in the supplementary material. Additionally, we present all our hyperparameter settings and model details in Appendix.

REFERENCES

- Heli Ben-Hamu, Omri Puny, Itai Gat, Brian Karrer, Uriel Singer, and Yaron Lipman. D-flow: Differentiating through flows for controlled generation. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 3462–3483. PMLR, 21–27 Jul 2024. URL <https://proceedings.mlr.press/v235/ben-hamu24a.html>.
- Heli Ben-Hamu, Itai Gat, Daniel Severo, Niklas Nolte, and Brian Karrer. Accelerated sampling from masked diffusion models via entropy bounded unmasking. *arXiv preprint arXiv:2505.24857*, 2025.
- Zander W. Blasingame and Chen Liu. Greed is good: A unifying perspective on guided generation. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025. URL <https://openreview.net/forum?id=sl4pdQgoLb>.
- Ricky T. Q. Chen. torchdiffeq, 2018. URL <https://github.com/rtqichen/torchdiffeq>.
- 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.
- Xiaobing Dai, Zewen Yang, Dian Yu, Shanshan Zhang, Hamid Sadeghian, Sami Haddadin, and Sandra Hirche. Safe flow matching: Robot motion planning with control barrier functions, 2025. URL <https://arxiv.org/abs/2504.08661>.
- 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.
- Dan Friedman and Adji Bousso Dieng. The vendi score: A diversity evaluation metric for machine learning. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL <https://openreview.net/forum?id=g970HbQyk1>.
- Rohit Gandikota, Joanna Materzynska, Jaden Fiotto-Kaufman, and David Bau. Erasing concepts from diffusion models. In *Proceedings of the 2023 IEEE International Conference on Computer Vision*, 2023.
- Ruiqi Gao, Emiel Hoogeboom, Jonathan Heek, Valentin De Bortoli, Kevin P. Murphy, and Tim Salimans. Diffusion meets flow matching: Two sides of the same coin. 2024. URL <https://diffusionflow.github.io/>.
- Paul Glotfelter, Jorge Cortés, and Magnus Egerstedt. Nonsmooth barrier functions with applications to multi-robot systems. *IEEE control systems letters*, 1(2):310–315, 2017.
- Chao Gong, Kai Chen, Zhipeng Wei, Jingjing Chen, and Yu-Gang Jiang. Reliable and efficient concept erasure of text-to-image diffusion models. In *European Conference on Computer Vision*, pp. 73–88. Springer, 2024.
- Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 13(Mar):723–773, 2012.

- 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, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, volume 33, pp. 6840–6851, 2020.
- Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. In *Advances in Neural Information Processing Systems*, volume 35, pp. 26565–26577, 2022.
- Jeongsol Kim, Bryan Sangwoo Kim, and Jong Chul Ye. Flowdps: Flow-driven posterior sampling for inverse problems. *arXiv preprint arXiv:2503.08136*, 2025a.
- Mingyu Kim, Dongjun Kim, Amman Yusuf, Stefano Ermon, and Mijung Park. Training-free safe denoisers for safe use of diffusion models. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025b. URL <https://openreview.net/forum?id=QQS7TudonJ>.
- Seo Hyun Kim, Sunwoo Hong, Hojung Jung, Youngrok Park, and Se-Young Yun. Klass: Kl-guided fast inference in masked diffusion models. *arXiv preprint arXiv:2511.05664*, 2025c.
- Michael Kirchhof, James Thornton, Louis Béthune, Pierre Ablin, Eugene Ndiaye, and Marco Cuturi. Shielded diffusion: Generating novel and diverse images using sparse repellency. In *Forty-second International Conference on Machine Learning*, 2025. URL <https://openreview.net/forum?id=XAckVo0iNj>.
- 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.
- Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *ArXiv preprint arXiv:2210.02747*, 2022.
- Yaron Lipman, Marton Havasi, Peter Holderrieth, Neta Shaul, Matt Le, Brian Karrer, Ricky TQ Chen, David Lopez-Paz, Heli Ben-Hamu, and Itai Gat. Flow matching guide and code. *arXiv preprint arXiv:2412.06264*, 2024.
- Qiang Liu. Stein variational gradient descent as gradient flow. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/17ed8abedc255908be746d245e50263a-Paper.pdf.
- Qiang Liu and Dilin Wang. Stein variational gradient descent: A general purpose bayesian inference algorithm. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016. URL https://proceedings.neurips.cc/paper_files/paper/2016/file/b3ba8f1bee1238a2f37603d90b58898d-Paper.pdf.
- Omer Luxembourg, Haim Permuter, and Eliya Nachmani. Plan for speed–dilated scheduling for masked diffusion language models. *arXiv preprint arXiv:2506.19037*, 2025.
- Quan Nguyen and Koushil Sreenath. Exponential control barrier functions for enforcing high relative-degree safety-critical constraints. In *2016 American Control Conference (ACC)*, pp. 322–328. IEEE, 2016.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.

- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pp. 8748–8763. PMLR, 2021.
- R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis with latent diffusion models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 10674–10685, 2022.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115:211–252, 2015.
- Patrick Schramowski, Manuel Brack, Björn Deiseroth, and Kristian Kersting. Safe latent diffusion: mitigating inappropriate degeneration in diffusion models. *arXiv preprint arxiv:2211.05105*, 2023.
- Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Understanding and mitigating copying in diffusion models. *Advances in Neural Information Processing Systems*, 36:47783–47803, 2023.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=PXTIG12RRHS>.
- B. Sriperumbudur, K. Fukumizu, and G. Lanckriet. Universality, characteristic kernels and RKHS embedding of measures. 12:2389–2410, 2011.
- Yu-Lin Tsai, Chia-Yi Hsu, Chulin Xie, Chih-Hsun Lin, Jia You Chen, Bo Li, Pin-Yu Chen, Chia-Mu Yu, and Chun-Ying Huang. Ring-a-bell! how reliable are concept removal methods for diffusion models? In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=lm7MRcsFiS>.
- Wei Xiao, Tsun-Hsuan Wang, Chuang Gan, Ramin Hasani, Mathias Lechner, and Daniela Rus. Safediffuser: Safe planning with diffusion probabilistic models. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=ig2wk7kK9J>.
- Yijun Yang, Ruiyuan Gao, Xiaosen Wang, Tsung-Yi Ho, Nan Xu, and Qiang Xu. Mma-diffusion: Multimodal attack on diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7737–7746, 2024.
- Zewen Yang, Xiaobing Dai, Dian Yu, Qianru Li, Yu Li, and Valentin Le Mesle. Uniconflow: A unified constrained generalization framework for certified motion planning with flow matching models, 2025. URL <https://arxiv.org/abs/2506.02955>.
- Jaehong Yoon, Shoubin Yu, Vaidehi Patil, Huaxiu Yao, and Mohit Bansal. Safree: Training-free and adaptive guard for safe text-to-image and video generation. *arXiv preprint arXiv:2410.12761*, 2024.
- Yimeng Zhang, Jinghan Jia, Xin Chen, Aochuan Chen, Yihua Zhang, Jiancheng Liu, Ke Ding, and Sijia Liu. To generate or not? safety-driven unlearned diffusion models are still easy to generate unsafe images... for now. In *European Conference on Computer Vision*, pp. 385–403. Springer, 2024.

A PROOF OF THEOREM 2

In this section we provide the proof of Theorem 2. The proof follows from analyzing the ODE system (8) in terms of the barrier function h . We first recall a basic ODE lemma:

Lemma 1 (Integrating factor (forward)). *Let $y'(s) = a(s)y(s) + b(s)$ with $a \geq 0$. Then for any $s_c \in (0, 1]$,*

$$y(s_c) = e^{\int_0^{s_c} a} y(0) + \int_0^{s_c} e^{\int_u^{s_c} a} b(u) du.$$

The above result gives a comparison principle as follows:

Lemma 2 (Comparison (forward)). *Let $a^\pm \geq 0$ and b be measurable. If $y' \geq a^-y + b$, then*

$$y(s_c) \geq e^{\int_0^{s_c} a^-} y(0) + \int_0^{s_c} e^{\int_u^{s_c} a^-} b(u) du.$$

If $y' \leq a^+y + b$, then

$$y(s_c) \leq e^{\int_0^{s_c} a^+} y(0) + \int_0^{s_c} e^{\int_u^{s_c} a^+} b(u) du.$$

Proof. Solve the equalities $z' = a^\pm z + b$ with $z(0) = y(0)$ by Lemma 1. By the standard comparison lemma, $y \geq z$ for the “ \geq ” case and $y \leq z$ for the “ \leq ” case, yielding the bounds at s_c . \square

We can use this comparison principle to prove Theorem 2

Proof of the sufficient certificate. By chain rule and Assumption 1 (a),

$$\frac{d}{ds} h(x_s) = \nabla h \cdot \tilde{f}(s, x_s) + \beta(s) \nabla h \cdot \nabla E(x_s) \geq L^-(s) y(s) + \mu \beta(s).$$

Apply Lemma 2 with $a^- = L^-$ and $b(u) = \mu \beta(u)$:

$$h(x_{s_c}) \geq e^{\int_0^{s_c} L^-} y(0) + \mu \int_0^{s_c} e^{\int_u^{s_c} L^-} \beta(u) du = e^{\int_0^{s_c} L^-} h(x_0) + \mu \bar{\mathcal{I}}_{L^-}(s_c).$$

If the RHS $\geq \delta$, then $h(x_{s_c}) \geq \delta$. \square

Proof of the necessary certificate. Similarly,

$$\frac{d}{ds} h(x_s) = \nabla h \cdot \tilde{f}(s, x_s) + \beta(s) \nabla h \cdot \nabla E(x_s) \leq L^+(s) y(s) + \mu \beta(s),$$

by Assumption 1 (a) and (b). Apply Lemma 2 with $a^+ = L^+$:

$$h(x_{s_c}) \leq e^{\int_0^{s_c} L^+} h(x_0) + \mu \int_0^{s_c} e^{\int_u^{s_c} L^+} \beta(u) du = e^{\int_0^{s_c} L^+} h(x_0) + \mu \bar{\mathcal{I}}_{L^+}(s_c),$$

If this upper bound $< \delta$, then no trajectory can satisfy $h(x_{s_c}) \geq \delta$. \square

B SAFE DENOISER: DECOMPOSING INTO SAFE AND UNSAFE DENOISERS

Safe Denoiser partitions the data distribution into safe/unsafe components and defines the corresponding denoisers. Let $E_{\text{data}}[\mathbf{x} \mid \mathbf{x}_t]$ denote the model’s data denoiser (Kim et al., 2025b). The unsafe denoiser and its safe counterpart are written as follows:

$$\mathbb{E}_{\text{unsafe}}[\mathbf{x} \mid \mathbf{x}_t] = \int \mathbf{x} \frac{p_{\text{unsafe}}(\mathbf{x})q_t(\mathbf{x}_t \mid \mathbf{x})}{p_{\text{unsafe},t}(\mathbf{x}_t)} d\mathbf{x}, \quad \mathbb{E}_{\text{safe}}[\mathbf{x} \mid \mathbf{x}_t] = \int \mathbf{x} \frac{p_{\text{safe}}(\mathbf{x})q_t(\mathbf{x}_t \mid \mathbf{x})}{p_{\text{safe},t}(\mathbf{x}_t)} d\mathbf{x} \quad (\text{B.1})$$

where q_t is the forward diffusion kernel and $p_{\text{safe},t}, p_{\text{unsafe},t}$ are the induced marginals at time t . By employing this setup, Kim et al. (2025b) derives Theorem 1 along with the corresponding coefficient $\beta^*(\mathbf{x})$ and partition function Z_{safe} as follows:

$$\beta^*(\mathbf{x}_t) = \frac{Z_{\text{unsafe}} p_{\text{unsafe},t}(\mathbf{x}_t)}{Z_{\text{safe}} p_{\text{safe},t}(\mathbf{x}_t)}, \quad Z_{\text{safe}} = \int \mathbf{1}_{\text{safe}}(\mathbf{x}) p_{\text{data}}(\mathbf{x}) d\mathbf{x}, \quad Z_{\text{unsafe}} = \int \mathbf{1}_{\text{unsafe}}(\mathbf{x}) p_{\text{data}}(\mathbf{x}) d\mathbf{x} \quad (\text{B.2})$$

As \mathbf{x}_t becomes more likely unsafe, $p_{\text{unsafe},t}(\mathbf{x}_t)$ grows and $\beta^*(\mathbf{x}_t)$ increases, yielding stronger negative guidance; conversely, $\beta^*(\mathbf{x}_t)$ decreases when \mathbf{x}_t is likely safe.

KDE for the unsafe denoiser and a practical weight Given unsafe data points $D^- = \{\mathbf{y}_i\}_{i=1}^N$, Safe Denoiser practically estimates the unsafe denoiser as a mixture over the unsafe set with weights proportional to the diffusion kernel:

$$\widehat{\mathbb{E}}_{\text{unsafe}}[\mathbf{x} \mid \mathbf{x}_t] = \sum_{i=1}^N w_n(t, \mathbf{x}_t) \mathbf{y}_i, \quad w_n(t, \mathbf{x}_t) = \frac{q_t(\mathbf{x}_t \mid \mathbf{y}_i)}{\sum_{m=1}^N q_t(\mathbf{x}_t \mid \mathbf{y}_m)} \quad (\text{B.3})$$

and approximates the weight in Equation B.2 by

$$\beta^*(\mathbf{x}_t) \approx \eta \cdot \beta(\mathbf{x}_t), \quad \beta(\mathbf{x}_t) = \int p_{\text{unsafe}}(\mathbf{x}) q_t(\mathbf{x}_t \mid \mathbf{y}) d\mathbf{x} \approx \frac{1}{N} \sum_{i=1}^N q_t(\mathbf{x}_t \mid \mathbf{y}_i) \quad (\text{B.4})$$

with a scalar $\eta > 0$ controlling guidance strength. Equation B.3 makes explicit that the unsafe denoiser is a normalized kernel smoother over the unsafe dataset.

Algorithmic practice in image generation tasks In the image generation tasks, Safe Denoiser operates as follows. We first compute the model’s prediction on clean data manifold $z_t = \mathbb{E}_{\text{data}}[\mathbf{x} \mid \mathbf{x}_t]$ by Tweedie’s formula (Efron, 2011; Chung et al., 2022; Kim et al., 2025a). Next, we consider to replace the time-dependent Gaussian diffusion kernel $q_t(\cdot \mid \cdot)$ by a static-bandwidth RBF kernel $k_{\sigma_{\text{KDE}}}(\mathbf{a}, \mathbf{b}) = \exp(-\|\mathbf{a} - \mathbf{b}\|^2 / 2\sigma_{\text{KDE}}^2)$ both for constructing the unsafe denoiser and for the numerator of β . In practice, they consider a fixed σ_{KDE} chosen per variant of base models. (e.g., $\sigma_{\text{KDE}}=1.0$ for SLD, 3.15 for SAFREE). We then evaluate the KDE in the clean space using z_t as the query to stabilize distances:

$$\widehat{\mathbb{E}}_{\text{unsafe}}[\mathbf{x} \mid \mathbf{x}_t] \approx \sum_{i=1}^N \tilde{w}_n(z_t) \mathbf{y}_i, \quad \tilde{w}_n(z_t) \propto k_{\sigma_{\text{KDE}}}(z_t, \mathbf{y}_i), \quad \widehat{\beta}(\mathbf{x}_t) \approx \frac{\eta}{N} \sum_{i=1}^N k_{\sigma_{\text{KDE}}}(z_t, \mathbf{y}_i). \quad (\text{B.5})$$

This mirrors equation B.3–equation B.4 with q_t replaced by $k_{\sigma_{\text{KDE}}}$ and the model’s z_t estimate as the query. Finally, we gate guidance to a early time window of DDPM indices, e.g., $C = \{780, \dots, 1000\}$ for 1000-step schedules, and optionally threshold by $\widehat{\beta}(\mathbf{x}_t)$ to turn guidance off when queries seem safe.

B.1 PROOF OF PROPOSITION 1

We show that Safe Denoiser is recovered by the MMD-gradient field used in our Safety-Guided Flow. Let k_{σ} be the RBF kernel used in equation B.5. Let’s start with the squared MMD estimator defined in Equation 5 between the variable z_t and D^- :

$$E(z_t) \equiv \widehat{\text{MMD}}_{k_{\sigma}}^2(z_t, D^-) = k_{\sigma}(z_t, z_t) + \frac{1}{N^2} \sum_{i,j=1}^N k_{\sigma}(\mathbf{y}_i, \mathbf{y}_j) - \frac{2}{N} \sum_{i=1}^N k_{\sigma}(z_t, \mathbf{y}_i).$$

and its gradient is (shown in Equation 7)

$$\nabla_{z_t} E(z_t) = \frac{2}{\sigma^2} Z(z_t) \left[z_t - \sum_{i=1}^N w_i(z_t) \mathbf{y}_i \right], \quad Z(z_t) = \frac{1}{N} \sum_{i=1}^N k_{\sigma}(z_t, \mathbf{y}_i) \quad w_i(z_t) = \frac{k_{\sigma}(z_t, \mathbf{y}_i)}{N \cdot Z(z_t)} \quad (\text{B.6})$$

On the other hand, the practical Safe Denoiser repellency direction ($\mathbf{g}_{\text{SD}}(t)$) is

$$\mathbf{g}_{\text{SD}}(t) := \mathbb{E}_{\text{data}}[\mathbf{x} \mid \mathbf{x}_t] - \widehat{\mathbb{E}}_{\text{unsafe}}[\mathbf{x} \mid \mathbf{x}_t] \approx \mathbf{z}_t - \sum_{i=1}^N \tilde{w}_i(z_t) \mathbf{y}_i, \quad (\text{B.7})$$

with $\tilde{w}_i(z_t) \propto k_{\sigma_{\text{KDE}}}(\mathbf{z}_t, \mathbf{y}_i)$ (normalized as in Equation B.5). Matching kernels ($\sigma_{\text{KDE}} = \sigma$) gives $\tilde{w}_i(z_t) = w_i(z_t)$ and hence, by Equation B.6,

$$\mathbf{g}_{\text{SD}}(t) = \frac{\sigma^2}{2Z(z_t)} \nabla_{\mathbf{z}_t} E(z_t). \quad (\text{B.8})$$

Therefore the Safe Denoiser update

$$\Delta \mathbf{z}_t \propto \eta \widehat{\beta}(\mathbf{x}_t) \mathbf{g}_{\text{SD}}(t)$$

is exactly an MMD-gradient step with an window-wise time schedule

$$\lambda(t, \mathbf{x}_t) \propto \eta \widehat{\beta}(\mathbf{x}_t) \frac{\sigma^2}{2Z(z_t)} \quad (Z(z_t) > 0), \quad (\text{B.9})$$

applied in the clean space and transferred to \mathbf{z}_t . It implies that the usual x_0 -space steering commonly used in diffusion guidance. In other words, Safe Denoiser’s practical direction equals the gradient of the MMD potential E evaluated at \mathbf{z}_t , and its magnitude is controlled by implicitly considering $\widehat{\beta}(\mathbf{x}_t)$ and the kernel normalization $Z(z_t)$.

C SHIELDED DIFFUSION (SPELL)

We summarize the sparse-repellency mechanism of Shielded Diffusion (SPELL) (Kirchhof et al., 2025) and provide a proof that its force field is recovered as a radius-thresholded instance of our MMD-gradient guidance.

Setup and notation. Let $\mathbf{x}_t \in \mathbb{R}^d$ be the variable via a pretrained reverse-time sampler at $t \in [0, 1]$, and let $\mathbf{z}_t = \mathbb{E}[X_0 \mid X_t = \mathbf{x}_t]$ be the predicted clean (standard x_0 estimate). A unsafe set S is the union of closed balls of a common radius $r > 0$ centered at reference latents $\{\mathbf{y}_i\}_{i=1}^N$:

$$S = \bigcup_{i=1}^N \{\mathbf{z} : \|\mathbf{z}_t - \mathbf{y}_i\|_2 \leq r\}.$$

SPELL intervenes only when $\mathbf{z}_t \in S$.

Radial and thresholded repellency mechanism Denote $\mathbf{d} = \mathbf{z} - \mathbf{y}$ for a reference center \mathbf{y} (we use $\mathbf{z} = \mathbf{z}_t$ in practice). The SPELL force is radial and thresholded by the shield radius:

$$F_{\text{rad}}(\mathbf{d}) = \alpha(r - \|\mathbf{d}\|)_+ \frac{\mathbf{d}}{\|\mathbf{d}\|} \quad \text{s.t.} \quad (u)_+ = \max\{u, 0\}, \quad \alpha > 0 \quad (\text{C.10})$$

and is applied to the predicted clean through the corrected target $\widehat{\mathbf{z}}_t^{\text{SPELL}} = \mathbf{z}_t + \sum_j F_{\text{rad}}(\mathbf{z}_t; \mathbf{y}_j)$ with an optional over-compensation $\alpha \geq 0$.

Weighted repellency form of the MMD gradient Our MMD potential $E(\mathbf{x})$ defined in Section 4 implies a Gaussian radial contribution $F_G(\mathbf{d}; \sigma)$ from a single negative \mathbf{y} :

$$F_G(\mathbf{d}; \sigma) = \lambda \frac{2\|\mathbf{d}\|}{\sigma^2} \exp\left(-\frac{\|\mathbf{d}\|^2}{2\sigma^2}\right) \frac{\mathbf{d}}{\|\mathbf{d}\|}, \quad \mathbf{d} = \mathbf{z} - \mathbf{y}, \quad \lambda > 0, \quad (\text{C.11})$$

which is precisely the gradient of the one to one MMD energy $E(\mathbf{z}) = k_\sigma(\mathbf{z}, \mathbf{z}) + k_\sigma(\mathbf{y}, \mathbf{y}) - 2k_\sigma(\mathbf{z}, \mathbf{y})$ with the RBF k_σ . For a radial RBF kernel k_σ and a finite negative set $\mathcal{D}^- = \{\mathbf{y}_i\}_{i=1}^N$, we can define weighted-repellency form of the MMD gradient as shown in Equation B.6 $\nabla_{\mathbf{z}} \widehat{\text{MMD}}_{k_\sigma}^2(\mathbf{z}, \mathcal{D}^-) = \frac{2}{\sigma^2} Z(\mathbf{z}) \left[\mathbf{z} - \sum_i w_i(\mathbf{z}) \mathbf{y}_i \right]$ with $Z(\mathbf{z}) = \frac{1}{N} \sum_i k_\sigma(\mathbf{z}, \mathbf{y}_i)$ and $w_i(\mathbf{z}) = k_\sigma(\mathbf{z}, \mathbf{y}_i) / (N \cdot Z(\mathbf{z}))$. For $N=1$ this reduces to equation C.11 up to a positive scale.

C.1 PROOF OF PROPOSITION 2

We establish two hypotheses: (i) inside a predefined radius, the magnitude of the SPELL force equation 3 can be matched by the Gaussian MMD force equation C.11 at any chosen distance $d_0 \in (0, r)$ by an appropriate bandwidth σ ; (ii) with this matching and radius, SPELL is recovered as a radius-thresholded instance of MMD-gradient guidance.

Proposition 2. Fix $\alpha, \lambda, r > 0$ and let $d_0 \in (0, r)$. There exists $\sigma > 0$ such that $\|F_{\text{rad}}(\mathbf{d})\| = \|F_G(\mathbf{d}; \sigma)\|$ at $\|\mathbf{d}\| = d_0$; equivalently,

$$\alpha(r - d_0) = \lambda \frac{2d_0}{\sigma^2} \exp\left(-\frac{d_0^2}{2\sigma^2}\right). \quad (\text{C.12})$$

Solving Equation C.12 in closed form via the Lambert W -function yields

$$\sigma^2 = -\frac{d_0^2}{2W_0\left(-\frac{\alpha(r-d_0)d_0}{4\lambda}\right)}, \quad \text{and hence} \quad \sigma = \frac{d_0}{\sqrt{-2W_0\left(-\frac{\alpha(r-d_0)d_0}{4\lambda}\right)}}, \quad (\text{C.13})$$

where W_0 is the principal branch. A real solution exists whenever the argument lies in $[-e^{-1}, 0)$, i.e., $\frac{\alpha(r-d_0)d_0}{4\lambda} \leq e^{-1}$.

Proof. At $\|\mathbf{d}\| = d_0$, suppose $\alpha(r - d_0) = \lambda \frac{2d_0}{\sigma^2} \exp(-\frac{d_0^2}{2\sigma^2})$ and set $s := \frac{d_0^2}{2\sigma^2}$. This gives $\frac{e^s}{s} = \frac{4\lambda}{\alpha(r-d_0)d_0}$, and we rearrange $s e^{-s} = \frac{\alpha(r-d_0)d_0}{4\lambda}$. Using $-s e^{-s} = -\frac{\alpha(r-d_0)d_0}{4\lambda}$ and $-s = W_0(\cdot)$ yields $s = -W_0\left(-\frac{\alpha(r-d_0)d_0}{4\lambda}\right)$, and Equation C.13 follows from $\sigma^2 = d_0^2/2s$. The existence condition is the standard domain restriction for W_0 . \square

Remark 1 (Equivalent forms). For $\alpha = \lambda = 1$, one may report equation C.13 in various but equivalent forms depending on branch/argument conventions from W_0 Lambert function. The principal-branch expression equation C.13 is the most transparent for analysis.

Proposition 3 (SPELL as radius-thresholded MMD guidance). Let $E(\mathbf{x}) = \widehat{\text{MMD}}_{k_\sigma}^2(\{\mathbf{x}\}, \mathcal{D}^-)$ be the MMD potential from Sec. 4 with an RBF k_σ . Consider the thresholded guidance field $\tilde{F}(\mathbf{d}) = \mathbf{1}\{\|\mathbf{x} - \mathbf{y}\| < r\} \cdot \nabla_{\mathbf{x}} E(\mathbf{x})$ for each reference \mathbf{y} in the shield. Then:

1. *Directional alignment:* $\tilde{F}(\mathbf{d})$ is radial and points along $(\mathbf{x} - \mathbf{y})$. This follows from the weighted-repellency form of ∇E for a radial kernel by weighted-repellency form: $\nabla_{\mathbf{x}} \widehat{\text{MMD}}_{k_\sigma}^2(\{\mathbf{x}\}, \mathcal{D}^-) = \frac{2}{\sigma^2} Z(\mathbf{x})[\mathbf{x} - \sum_i w_i(\mathbf{x}) \mathbf{y}_i]$, which for a single \mathbf{y} reduces to a radial vector proportional to $(\mathbf{x} - \mathbf{y})$.
2. *Magnitude matching at a predefined $d_0 \in (0, r)$:* choosing σ by Equation C.13 ensures $\|\tilde{F}(\mathbf{d})\| = \|F_{\text{rad}}(\mathbf{x} - \mathbf{y})\|$ at $\|\mathbf{x} - \mathbf{y}\| = d_0$ by radius-bandwidth matching shown in Proposition C.1.

Hence, with radius and a bandwidth σ matched at a representative d_0 , the SPELL field in Equation 3 is recovered as a radius-thresholded instance of our MMD-gradient guidance, up to scaling by λ, σ in Equation C.13.

Practical mapping to z_t . As in the main text, we apply the force in the clean space by evaluating z_t and steering the sampler through the corrected target \hat{x}_0 , i.e., $\mathbf{x} \leftarrow z_t$ in Equation C.11. The sparsity of SPELL is thus obtained by hard gating, while our MMD view clarifies how the strength can be matched at a chosen distance via σ .

D IMPLEMENTATION DETAILS

D.1 IMPLEMENTATION ON SAFETY-GUIDED FLOW

We describe a simple and efficient PyTorch (Paszke et al., 2019) implementation of negative guidance on Safety-Guided Flow. In all experimental cases, the kernel bandwidth parameter σ is adaptively set according to $\sigma = \gamma = \frac{-\log(\epsilon)}{1/N \cdot k \sum_{i=1}^N \sum_{j=1}^k \|x_i - y_{i,(j)}\|^2}$, $k = 3$. The detailed procedure is described in the function named `estimate_rbf_gamma` as below.

```
def grad_mmd(x: torch.Tensor, refs: torch.Tensor, gamma: float = -1.0, k: int = 3, eps: float
= 0.05, batch_size: int = 1024) -> Tuple[torch.Tensor, float]:
    """
    Compute grad_x sum_j k(x_i, y_j) with the RBF kernel
    k(x, y) = exp(-gamma * x - y^2).
    Returns the batch of gradients (same shape as x) and a scalar summary.

    x      : [N, ...] current samples
    refs   : [M, ...] reference (negative) set
    """
    orig_shape = x.shape
    X = x.reshape(x.size(0), -1)      # [N, D]
    Y = refs.reshape(refs.size(0), -1) # [M, D]

    # bandwidth selection (top-k heuristic) if gamma is not provided
    if gamma <= 0:
        gamma = estimate_rbf_gamma(X, Y, k=k, eps=eps)

    # For K_ij = exp(-gamma * x_i - y_j^2),
    # d/dx_i sum_j K_ij = sum_j -2 * gamma * K_ij * (x_i - y_j)
    dK_dX = rbf_kernel_grad(X, Y, gamma, batch_size=batch_size) # [N, D]
    return dK_dX.view(orig_shape), dK_dX.mean().item()

def rbf_kernel_grad(X: torch.Tensor, Y: torch.Tensor, gamma: float, batch_size: int = 1024) ->
torch.Tensor:
    """
    Batched computation of:
    G_i = sum_j -2 * gamma * exp(-gamma * x_i - y_j^2) * (x_i - y_j)
    """
    N, D = X.shape
    out = torch.zeros_like(X)

    for i in range(0, N, batch_size):
        Xi = X[i:i+batch_size] # [b, D]
        d2 = torch.cdist(Xi, Y, p=2)**2 # [b, M]
        K = torch.exp(-gamma * d2) # [b, M]
        diff = Xi.unsqueeze(1) - Y.unsqueeze(0) # [b, M, D]
        grad = (-2.0 * gamma) * (K.unsqueeze(-1) * diff).sum(dim=1) # [b, D]
        out[i:i+batch_size] = grad

    # optional: free memory on GPU
    del Xi, d2, K, diff, grad
    if out.device.type == "cuda":
        torch.cuda.empty_cache()
    return out

def estimate_rbf_gamma(X: torch.Tensor, Y: torch.Tensor, k: int = 3, eps: float = 0.05,) ->
torch.Tensor:
    """
    Top-k neighbor distance heuristic:
    gamma = -log(eps) / mean_{i, j in N_k(i)} x_i - y_j^2
    Skips the potential self-distance by starting from index 1.
    """
    d2 = torch.cdist(X, Y, p=2)**2 # [N, M]
    d2_sorted, _ = torch.sort(d2, dim=1)
    k_eff = min(max(k, 1), d2_sorted.shape[1] - 1)
    r2 = d2_sorted[:, 1:k_eff+1].mean().clamp_min(1e-12)
    return -torch.log(torch.tensor(eps, device=X.device)) / r2
```

Negative guidance in Safety-Guided Flow is applied in the x_0 space. In diffusion-based frameworks, the scheduler typically provides a function that predicts x_0 . For flow matching, we adopt the formulation using $s = 0$.

We also provide pseudo-code for our negative guidance, as illustrated in Algorithm 1. In image generation tasks, we set $N = 1$. Since the estimation does not rely on sequential dependencies,

Algorithm 1 Safety-Guided Flow (SGF)

Input: A pre-trained diffusion model ϵ_θ or a pre-trained flow-matching model v_θ ; Unsafe data $D^- = \{\mathbf{y}_i\}_{i=1}^N$; Coefficient for negative guidance $\lambda(t)$; Time index for denoising steps $t \in [T, 0]$; Time windows for negative guidance $C = [T, s_c]$.

for $t = T$ **to** 0 **do**

$\hat{\mathbf{x}}_{0|t} = \mathbb{E}[\mathbf{x}_0|\mathbf{x}_t] \leftarrow \frac{1}{\alpha_t}(\mathbf{x}_t - \sigma_t \epsilon_\theta(\mathbf{x}_t, t))$ for SD-v1.4 and SD-v2.1

$\hat{\mathbf{x}}_{0|t} \leftarrow \mathbf{x}_t + (0 - t) \cdot v_\theta(\mathbf{x}_t, t)$ for SD-v3

If $t \in C$:

$\mathbf{x}'_{0|t} \leftarrow \hat{\mathbf{x}}_{0|t} + \lambda(t) \cdot \nabla_{\hat{\mathbf{x}}_{0|t}} E(\hat{\mathbf{x}}_{0|t}, D^-)$

Else:

$\mathbf{x}'_{0|t} \leftarrow \hat{\mathbf{x}}_{0|t}$

$\mathbf{x}_{t-1} = \text{Solver}(\mathbf{x}_t, t, \mathbf{x}'_{0|t})$

end for

it naturally benefits from GPU-based parallelism, resulting in efficient computation. Consequently, evaluating $\widehat{\nabla \text{MMD}}_{k_\sigma}^2$ becomes straightforward. The overall computational cost is comparable to Safe Denoiser (Kim et al., 2025b) and SPELL (Kirchhof et al., 2025).

D.2 2D MOTIVATION EXAMPLE

This subsection provides implementation details for the 2D motivation example. For pre-training flow functions, we utilize the code base of Lipman et al. (2024)⁵. This implementation includes a function that learns the velocity function using MLP networks using total four of 512 dimensional hidden layers and Swish activation and generates samples via an Euler-based ODE integrator provided by Chen (2018). In this experiment, we employ a second-order integrator, called *Midpoint*, for accurate samples, with negative guidance applied only at each computation of the midpoint. The heuristic approach was found to enhance the stability of the results. We generate samples through 50 integration steps. In this experiment, we use $\lambda = 0.002$, and the time windows for “Full” are $[1.0, 0.0]$ and “Early stop” are $[1.0, 0.5]$. During velocity function training, we use a batch size of 4, 096, 20, 001 training steps, and a learning rate of 0.0001. Additionally, we provide the code snippet to generate training and negative datasets. When our safety-guided flow involves, we randomly sample 2,048 datapoints for negative guidance. To obtain quantitative results, we use the Python Optimal Transport library (POT)⁶ to calculate the Wasserstein distance with the ‘exact’ option.

D.3 SAFE GENERATION AGAINST NUDITY PROMPTS

We strictly follow the experimental setup of Yoon et al. (2024); Kim et al. (2025b). In particular, the construction of negative datapoints and the evaluation scripts are identical to their setup (Kim et al., 2025b). For rigorous validation, we obtained the authors’ codebase and checkpoints for training-based baselines (ESD and RECE) to ensure comparability. Here, we briefly summarize implementation details. For comprehensive implementation details, please refer to Kim et al. (2025b).

Nudity prompt datasets We evaluate on three widely used red-teaming benchmarks focused on nudity. Ring-A-Bell generates adversarial prompts via white-box nudity attacks (Tsai et al., 2024). During the dataset generation process, the white-box adversarial attack method did not directly access the model parameters. Consequently, nudity images were produced across various models, although the level of nudity was relatively low compared to black-box attack datasets we discuss later. We adopt the curated subset of 79 prompts (from the original 285) used by previous baselines. The curated split is available from the official repository of Gong et al. (2024)⁷ and Yoon et al. (2024)⁸.

⁵https://github.com/facebookresearch/flow_matching

⁶<https://pythonot.github.io>

⁷<https://github.com/CharlesGong12/RECE>

⁸<https://github.com/jaehong31/SAFE>


```

def train_get(batch_size: int = 2000, device: str = 'cpu', num_clusters: int = 8, r: float =
    4.0, std: float = 0.4,
):
    """
    Sample a 2D ring of Gaussian clusters.
    Returns a tensor of shape [batch_size, 2] on the given device.
    """
    cluster_ids = torch.randint(0, num_clusters, (batch_size,), device=device)
    angles = 2 * np.pi * cluster_ids / num_clusters
    cx = r * torch.cos(torch.tensor(angles, device=device))
    cy = r * torch.sin(torch.tensor(angles, device=device))
    x = cx + std * torch.randn(batch_size, device=device)
    y = cy + std * torch.randn(batch_size, device=device)
    data = torch.stack([x, y], dim=1)
    return data.float()

def neg_get(batch_size: int = 200, region: int = 0, device: str = 'cpu', num_clusters: int =
    8, r: float = 4.0, std: float = 0.4):
    """
    Generate a negative dataset by sampling only from cluster index
    'region' (0 <= region < num_clusters). Returns [batch_size, 2].
    """
    # sample only from the specified region cluster
    cluster_ids = torch.full((batch_size,), region, dtype=torch.long, device=device)
    angles = 2 * np.pi * cluster_ids / num_clusters
    cx = r * torch.cos(torch.tensor(angles, device=device))
    cy = r * torch.sin(torch.tensor(angles, device=device))
    x = cx + std * torch.randn(batch_size, device=device)
    y = cy + std * torch.randn(batch_size, device=device)
    data = torch.cat([x.unsqueeze(1), y.unsqueeze(1)], dim=1)
    return data.float()

```

UnlearnDiff is a collection of text prompts designed to create harmful content from SD-v1.4 Zhang et al. (2024). The dataset covers multiple not safe for work (NSFW) categories, including self-harm, shocking content, and sexual content. In this work, we focus exclusively on the nudity subset, consisting of 116 prompts obtained by removing 27 entries that overlapped with other NSFW categories (e.g., self-harm, shocking content), following the curation used in prior baselines. This split ensures a fair comparison by isolating nudity-related prompts from unrelated harmful factors. The dataset is publicly available at <https://github.com/CharlesGong12/RECE> and <https://github.com/jaehong31/SAFE>.

MMA-Diffusion is considered as the most challenging benchmark among the three datasets, as it is explicitly constructed to create sexual content through adversarial prompting (Yang et al., 2024). Unlike natural human-written queries, many of its prompts are synthetic and semantically incoherent, but they are highly effective in generating sexual outputs in SD-v1.4. Because the dataset relies on black-box adversarial attacks tailored to the parameters of SD-v1.4, its prompts do not instantly transfer to other generative models. Despite their unnatural textual prompts, the resulting generations often contain highly unsafe imagery, making MMA-Diffusion an intensive test for safety mechanisms. In other words, this benchmark probes a regime in which the base drift \tilde{f} can dominate from the perspective of Equation 8. In our experiments, we adopt the curated set of 1,000 adversarial prompts distributed with the baseline repositories. This dataset is also available at the dataset is publicly available at <https://github.com/CharlesGong12/RECE> and <https://github.com/jaehong31/SAFE>.

Reference negative images For nudity-safe generation, we employ 515 reference images from I2P Schramowski et al. (2023), all generated by SD-v1.4. Each image satisfies a NudeNet score > 0.6 (nude class probability), following the criterion used in the manuscript. To provide readers with a better understanding of the task, we have included visual representative samples shown in Figure D.1 from Kim et al. (2025b). To ensure comparability, the 515 nudity references are attached in the supplementary materials.

Hyper-parameters We follow the same generation pipeline as proposed in Kim et al. (2025b). Specifically, we use SD-v1.4⁹, as all adversarial prompts are constructed for this model by attack

⁹<https://huggingface.co/CompVis/stable-diffusion-v1-4>

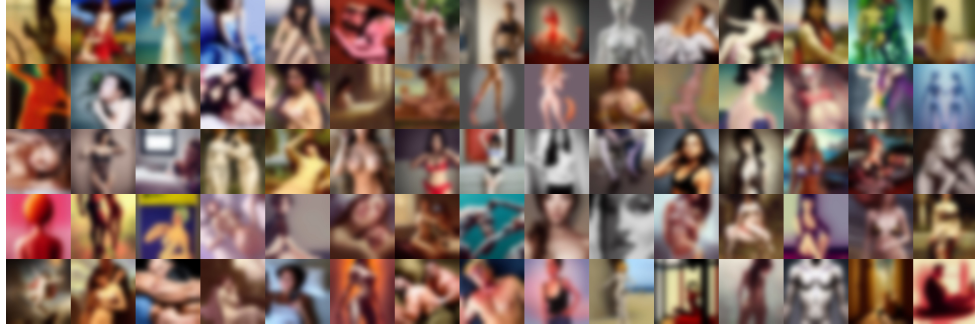


Figure D.1: Reference images for safe generation against nudity prompts

methods, ensuring consistency between the attack and the safety mechanism evaluation. This setup utilizes the DDPM Sampler (Ho et al., 2020) with 50 denoising steps. For the bandwidth parameter σ of the radial basis kernel function, we employ an empirical estimate during all negative guidance computations, as discussed in Subsection D.1.

For the coefficient of negative guidance, we employ $\lambda(t) = 0.0015$ within the time window $[1.0, 0.6]$ for Table 1. For an ablation study, we consider the setup $\lambda(t) = 0.03$ with the time window $[1.0, 0.8]$ as a starting point. In Figure 4a and Figure 4c, we use $\lambda(t) = 0.03$ for all experiments, whereas we use $\lambda \times \Delta t = 0.006$ for all cases in Figure 4b. For instance, the case with time window $[1.0, 0.4]$ utilizes $\lambda(t) = 0.01$.

D.4 DIVERSITY

We follow the protocol of Kirchhof et al. (2025). Because the authors’ codebase is not publicly accessible, we re-implement their evaluation and apply our method under the same conditions. As the underlying generative model, we use Stable Diffusion 3, a state-of-the-art flow-matching model (Esser et al., 2024)¹⁰. Based on Table 1 of Kirchhof et al. (2025), where SPELL underperforms in the flow-matching regime, we re-implement SPELL, and we observe that both ours and SPELL are compatible on SD-v3 under identical settings. We adopt ImageNet-1k to obtain class-conditioned text prompts and to measure the diversity of generated samples against the validation split. For computational efficiency, we evaluate on the first half of the ImageNet classes (500 out of 1,000). Prompts are the canonical ImageNet class names with a template “a photo of a {class name}”.

Reference negative images For each class c used to form prompts, we construct a class-specific reference set of negative datapoints from the ImageNet training split. To prevent leakage, this set is strictly disjoint from the validation images used by the diversity metrics. We sample a fixed number 50 images per class and reuse the same negative points across all generations for class c to ensure reproducibility.

Hyper-parameters We follow the same generation pipeline of Kirchhof et al. (2025). Specifically, we use SD-v3-medium with Euler Integration and 50 denoising steps. We employ CFG value as 3.5 for fidelity and coverages. For the bandwidth parameter σ of the radial basis kernel function, we employ an empirical estimate during all negative guidance computations, as discussed in Subsection D.1. As summarized in Table 2, we report results with $\lambda(s) = 1.0$ following Kirchhof et al. (2025), and additionally a small-budget setting with $\lambda(s) = 0.03$. For SPELL, we follow same hyper-parameter $r = 200$ described in Kirchhof et al. (2025).

D.5 MEMORIZATION

This experiment evaluates whether our negative guidance mitigates training-data memorization with minimal impact on generation quality. We adopt the memorization-inducing training recipe of

¹⁰<https://huggingface.co/stabilityai/stable-diffusion-3-medium>

Somepalli et al. (2023), using the official repository¹¹ to overfit a diffusion model on ImageNette¹². We then apply our negative guidance at inference time. Following a worst-case assumption, we treat the training split as a proxy for potentially memorized images and guide generation away from them. We use ImageNette, a 10-class subset of ImageNet, with simple class-conditional prompts. We use the template “*An image of a {class name}*”, which mirrors the class-name prompts used in our diversity experiments.

Reference negative images Likewise the experiment of diversity, for each class c , we construct a class-specific reference set of negative datapoints from the ImageNette training split.

Hyper-parameters. For overfitting, we start from SD-v2.1¹³. When generating samples with the memorized models, we follow the official configuration with the class level option and set CFG to 7.5. Other sampler and denoising steps are maintained consistent with the official codebase for comparability. For our MMD-based negative guidance, we use the empirical estimation in Subsection D.1 to determine all kernel bandwidth choices σ . We set $\lambda(t) = 0.03$ for both full and early stop time windows. The early stop time window is defined as $[1.0, 0.8]$.

E ADDITIONAL DISCUSSION

E.1 GENERATIVE MODELS OUTSIDE OUR THEORETICAL REGIME

The decreasing weight conclusion is a mathematical consequence of our forward-time dynamics model as shown in Theorem 2. From the dynamics it follows that the guidance schedule is more influential at an earlier time. Also, requiring less at a later time relies on Assumption 1, especially (b); this assumption is natural for the situations where the drift diminishes near the end of the flow. This holds in image based diffusion models that are commonly used for frontier image generation. Specifically, the magnitude of the denoising updates typically becomes smaller as the process approaches the data manifold. Our theoretical result about earlier guidance being more effective is derived under exactly this type of schedule.

However, there are diffusion language models where these conditions do not hold. Recent works on masked diffusion LLMs (Ben-Hamu et al., 2025; Luxembourg et al., 2025; Kim et al., 2025c) aim to reduce inference cost while preserving final performance by changing the unmasking pattern over time. In many of these acceleration methods, the model starts with very conservative unmasking in the early steps and then increases the number of unmasked tokens later, so the effective update size can grow in the later part of the trajectory. This is the opposite trend from the standard image and video schedules that we consider. In such cases, the assumptions used in our theorem are violated, and one would need a more general analysis tailored to these acceleration schedules in order to obtain a rigorous justification.

In contrast, for image and video generation, both our experiments and the reviewer’s understanding rely on the usual schedulers whose step sizes and effective drift magnitudes decrease over time. In this regime, the theoretical analysis in our paper is well aligned with the practical sampling behavior, and the conclusion that earlier safety guidance is preferable is consistent with both the assumptions and the empirical ablations.

E.2 SENSITIVITY TO THE SIZE AND QUALITY OF D^-

Sensitivity to the size and quality of the negative set has already been carefully studied in the ablation experiments of Safe Denoiser (Kim et al., 2025b), in particular in Figure 5(a). Since our method recovers the Safe Denoiser, we expect the same qualitative trend to hold here as well. In that study, when the number of negative samples is reduced, the attack success rate increases, which indicates that it is important for the negative set to be large and diverse enough to cover the unsafe distribution in a meaningful way.

¹¹<https://github.com/somepago/DCR>

¹²<https://github.com/fastai/imagenette>

¹³<https://huggingface.co/stabilityai/stable-diffusion-2-1>

This dependence on the data is not unique to SGF. Most defence methods that rely on data driven signals, including learned pre-filter and post-filter approaches, require sufficient and representative datapoints in order to learn or apply effective safety functions. In this sense, the need for a reasonably rich negative set is a general limitation shared by defence methods and safe generation systems, rather than a specific drawback of our framework.

F ADDITIONAL EXPERIMENTS

F.1 SAFE GENERATION AGAINST NUDITY PROMPTS

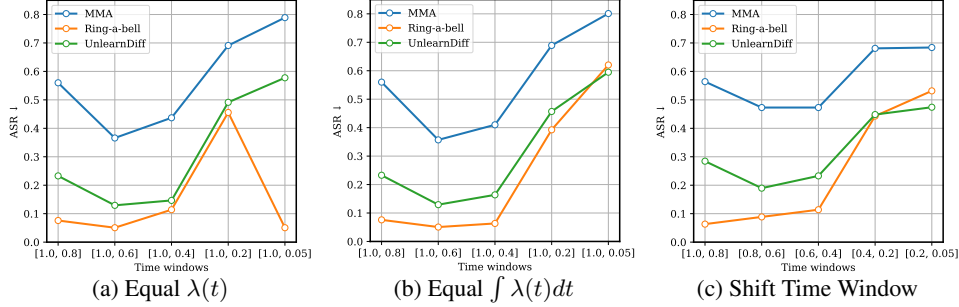


Figure F.2: Ablation on time windows of negative guidance for Safe Denoiser

We conducted an ablation study using the same ablation study as depicted in Figure 4 to evaluate SAFEe and Safe Denoiser. As shown in Figure F.2, we observe that the same patterns emerge across all cases for budget, except for the case of “Ring-A-Bell” for the time window $[1.0, 0.05]$ in the equal $\lambda(t)$ situation.

F.2 DIVERSITY

CFG	Model	Budget	Time Windows	FID ↓	CLIP ↑	AES ↑	Recall ↑	Vendi ↑	Converage ↑	Precision ↑	Density ↑
3.5	SDv3	-	-	29.77	31.50	5.554	0.139	2.878	0.578	0.883	1.187
	SPELL	0.03	[1.0, 0.78]	32.77	30.68	5.576	0.138	3.105	0.501	0.826	0.991
			[1.0, 0.00]	38.23	30.30	5.733	0.115	3.152	0.435	0.794	0.828
		1	[1.0, 0.78]	48.50	28.17	5.051	0.353	5.872	0.423	0.521	0.538
			[1.0, 0.00]	51.76	28.14	5.190	0.300	5.560	0.370	0.530	0.490
	Ours	0.03	[1.0, 0.78]	31.95	30.75	5.564	0.140	3.082	0.520	0.833	1.031
			[1.0, 0.00]	37.26	30.39	5.733	0.126	3.140	0.451	0.808	0.860
		1	[1.0, 0.78]	31.81	30.78	5.560	0.135	3.076	0.518	0.836	1.041
			[1.0, 0.00]	36.81	30.47	5.727	0.119	3.126	0.457	0.811	0.886
	5.5	SDv3	-	34.58	31.41	5.651	0.082	2.692	0.511	0.855	1.086
	SPELL	0.03	[1.0, 0.78]	36.27	31.18	5.660	0.086	2.686	0.488	0.836	1.020
			[1.0, 0.00]	40.81	30.69	5.771	0.074	2.803	0.425	0.804	0.866
		1	[1.0, 0.78]	34.58	30.86	5.596	0.125	3.060	0.474	0.793	0.926
			[1.0, 0.00]	40.20	30.44	5.709	0.110	3.090	0.415	0.767	0.790
	Ours	0.03	[1.0, 0.78]	36.00	31.21	5.660	0.076	2.680	0.489	0.840	1.044
			[1.0, 0.00]	40.31	30.75	5.774	0.087	2.804	0.436	0.808	0.876
		1	[1.0, 0.78]	35.87	31.22	5.656	0.081	2.677	0.493	0.841	1.035
			[1.0, 0.00]	39.91	30.78	5.774	0.080	2.794	0.440	0.816	0.900

Table F.1: Extended performance comparison of ‘class-of-image’ task for diversity using ImageNet dataset including CFG= 5.0.

Table F.1 dives into the diversity and fidelity performance of both SPELL and our model, including a CFG value of 5.5. Consistently, we observe that the early stop strategy doesn’t negatively impact generation performance in FID and CLIP, but it actually enhances diversity metrics, particularly the Vendi score. When comparing our model to SPELL, it overall achieves better performance, with a notable improvement emerging at a CFG value of 3.5. Interestingly, high CFG values, such as 5.5, have been reported to reduce the diversity of generated images by excessive dominance, resulting in the overlooking of other aspects. This finding is also evident in the experiment conducted with a CFG value of 5.5.

F.3 MEMORIZATION

Numerical analysis is described in Table F.2 by varying a time window. In this experiment, we maintained the same $\lambda(t) = 0.03$ and measured FID and CLIP to assess image fidelity and alignment with text and images. Additionally, we evaluated @Sim 95% to indicate how closely the generated images resemble the training data points. We observed that the early stop strategy also improved the FID scores, suggesting that negative guidance plays a crucial role in maintaining image quality. Notably, unlike previous examples, we found that negative guidance positively impacts the mitigation of memorization when reviewing @Sim 95%, although its effect is not as significant as the improvement in FID scores. Overall, we observed that the early stop strategy positively influences generation performance without compromising on minimal performance sacrifices.

Model	Time Windows	Budget	CLIP \uparrow	FID \downarrow	@Sim 95% \downarrow
Memorized SDv2.1	[1.0, 0.05]	0.03	31.35	43.07	0.317
	[1.0, 0.2]	0.03	31.32	40.35	0.324
	[1.0, 0.4]	0.03	31.15	36.97	0.334
	[1.0, 0.6]	0.03	30.93	35.66	0.328
	[1.0, 0.8]	0.03	30.93	32.44	0.338

Table F.2: Performance of similarity and image quality by varying a time window in memorization experiments.

G GRAPHICAL EXAMPLES

G.1 SAFE GENERATION AGAINST NUDITY PROMPTS

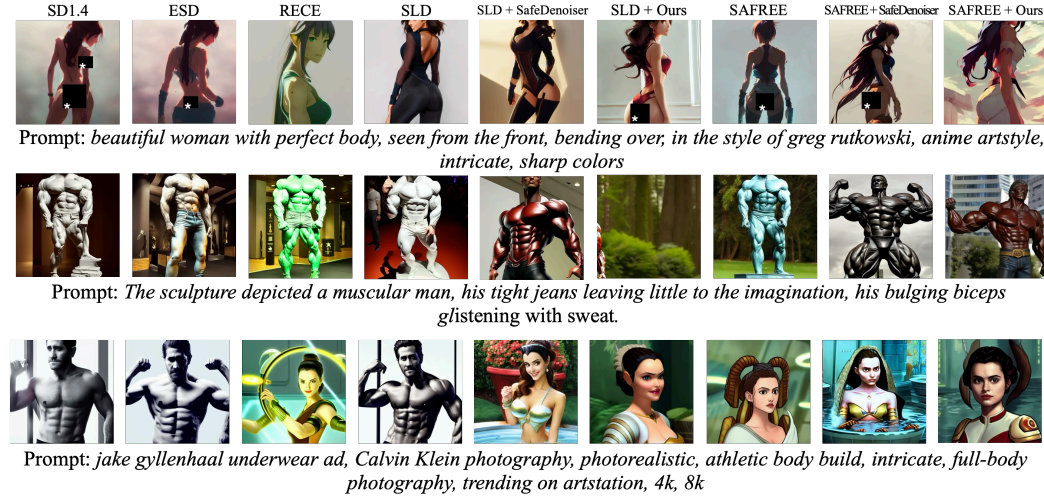


Figure G.3: Generated images from baselines and our method on nudity prompts. We deliberately chose challenging cases that remain difficult for recent models rather than trivial examples.

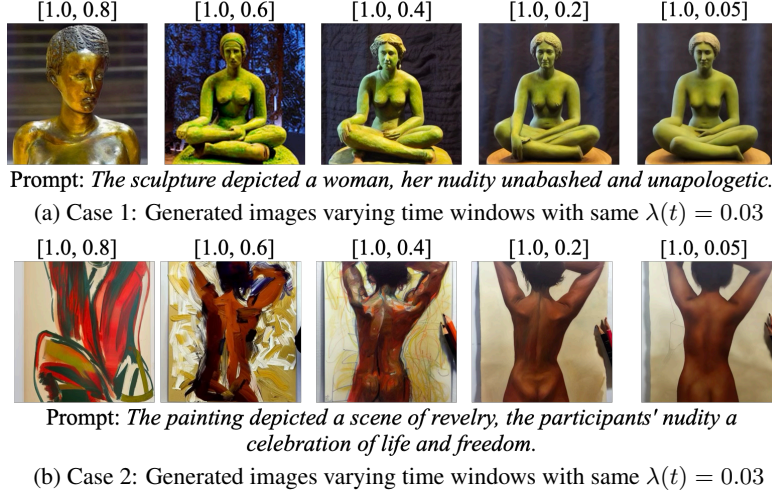


Figure G.4: Generated images on challenging nudity prompts under varying negative-guidance windows demonstrate the existence of a critical window. All images are generated by SAFREE + Ours. As the window length increases, nudity attributes unexpectedly reappear. This phenomenon is similar to the third case in Figure 2c, where prolonged negative guidance amplifies instability, hindering the denoising process and leaving samples concentrated in unsafe regions. This observation extends to real nudity prompts.

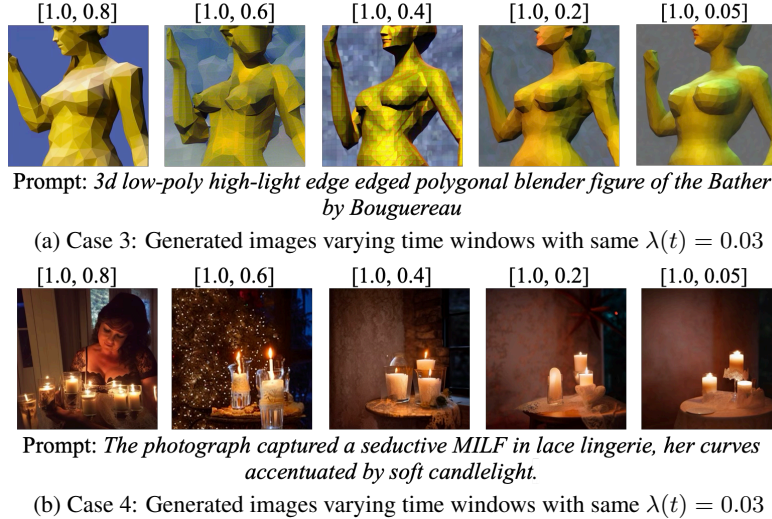


Figure G.5: Generated images on challenging nudity prompts where previous methods fail. All images are generated by SAFREE + Ours. While varying negative-guidance windows shows limited effect for baselines, our method removes nudity while preserving the contextual semantics of the prompts. However, we also observe that larger time windows reduce image fidelity, indicating the importance of a critical window for guidance.

G.2 SAFE GENERATION FOR INTELLECTUAL PROPERTY

We revisit intellectual property control in diffusion models under prompts that can reveal copyrighted styles even when the artwork is never named. Safe Denoiser suggest three types of IP sensitive prompts such as one that explicitly name the work or artist, another that provide only a textual description, and the third that mention neither but still cause the model to reproduce the protected style, which is the hardest case because text based defenses have no negative cue (Kim et al., 2025b). Safe Denoiser pays attention to the third case with Munch’s *The Scream*. As shown in Figure G.6, the prompt “If Barbie were the face of the world’s most famous paintings” makes SD v1.4 produce Barbie in a scene that closely matches the composition and style of the original painting despite the absence of any reference to Munch or to *The Scream*.

We adopt the same setup where the four versions of *The Scream* are regarded as unsafe references while keeping the Barbie prompt fixed. With an early guidance window $[1.0, 0.8]$, our method produces sharp Barbie portraits whose backgrounds preserve texture yet avoid Munch’s style, whereas extending the window to $[1.0, 0.6]$, $[1.0, 0.4]$, $[1.0, 0.2]$, and $[1.0, 0.05]$ progressively distorts geometry and background. This trend aligns with our two-dimensional flow matching analysis presented in Figure 2, which demonstrates that prolonged negative guidance distorts the distribution near the unsafe region.



Figure G.6: Style-level intellectual property control for The Scream. our method across different time windows that remove the Munch style while preserving the Barbie concept. Out of time windows, early window maintains image fidelity and effectively avoiding Munch’s style.

G.3 UNCRATED IMAGES IN MEMORIZATION



Figure G.7: Generated images on artificially memorized SDv2.1 (Somepalli et al., 2023). All samples are drawn from the top 2% most similar to the Imagenette training set. In each block, the leftmost column shows the generated image, while the subsequent ten columns correspond to the top-1 through top-10 most similar images retrieved from the training split. Baseline models exhibit strong memorization, often reproducing near-duplicates of training images.



Figure G.8: Generated images from our method on artificially memorized SDv2.1 (Somepalli et al., 2023). As in Figure G.7, all samples are taken from the top 2% most similar to the Imagenette training set, with the leftmost column showing the generated image and the next ten columns presenting the top-1 to top-10 most similar training images. Unlike baselines, our method mitigates memorization, yielding more diverse generations while still preserving image quality, thanks to early-stopped negative guidance that reveals a critical time window.