

# Dynamic Negative Guidance of Diffusion Models: Towards Immediate Content Removal

Félix Koulischer, Johannes Deleu, Gabriel Raya, Thomas Demeester\*, Luca Ambrogioni\*

Radboud Universiteit

Contact: felix.koulischer@ugent.be

TL;DR: A novel theoretically grounded *Dynamic Negative Guidance* scheme is introduced as alternative for the Negative Prompting algorithm.

## Motivation

We derive a new *Dynamic Negative Guidance* scheme from first principles that better preserves image diversity when compared to the theoretically ungrounded Negative Prompting (NP) algorithm. Preserving image diversity is key to both improve FIDs and reduce potential ethical concerns.

## Theoretical derivation

Negative Prompting originated as a simple change of sign from the Classifier-Free Guidance (CFG) equation. This implies sampling from an ill-defined distribution:

$$p(\mathbf{x}, \mathbf{c}) \propto p(\mathbf{x}) / p(\mathbf{c} | \mathbf{x})^\lambda \quad (1)$$

Instead, we remove the asymptote by reconstructing the desired conditional from the undesired conditional:

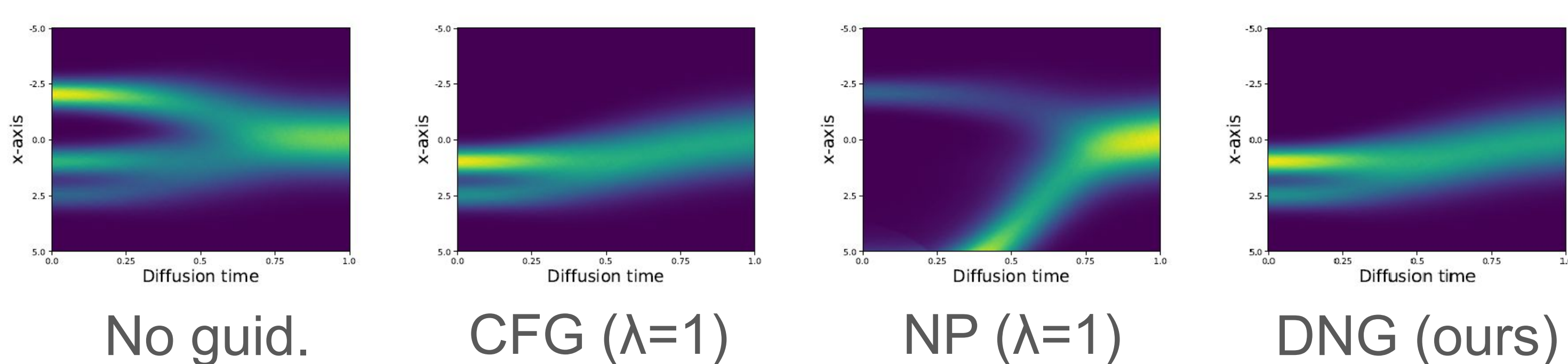
$$p_t(\mathbf{x} | \mathbf{c}_+) \propto p_t(\mathbf{x}) (1 - p_t(\mathbf{c} | \mathbf{x})) \quad (2)$$

From a score based perspective, this results in:

$$\nabla_{\mathbf{x}} \log p_t(\mathbf{x} | \mathbf{c}_+) = \nabla_{\mathbf{x}} \log p_t(\mathbf{x}) + \nabla_{\mathbf{x}} \log (1 - p_t(\mathbf{c} | \mathbf{x})) \quad (3)$$

$$\lambda(\mathbf{x}, t) = \left[ -\frac{p_t(\mathbf{c} | \mathbf{x})}{1 - p_t(\mathbf{c} | \mathbf{x})} \right] \nabla_{\mathbf{x}} \log p_t(\mathbf{c} | \mathbf{x})$$

We find that a *dynamic* state-time modulated guidance scale is required. This guidance scale is large when the posterior probability is close to one, and small otherwise. Therefore, only regions related to the negative prompt are affected, reducing the invasiveness of the method.

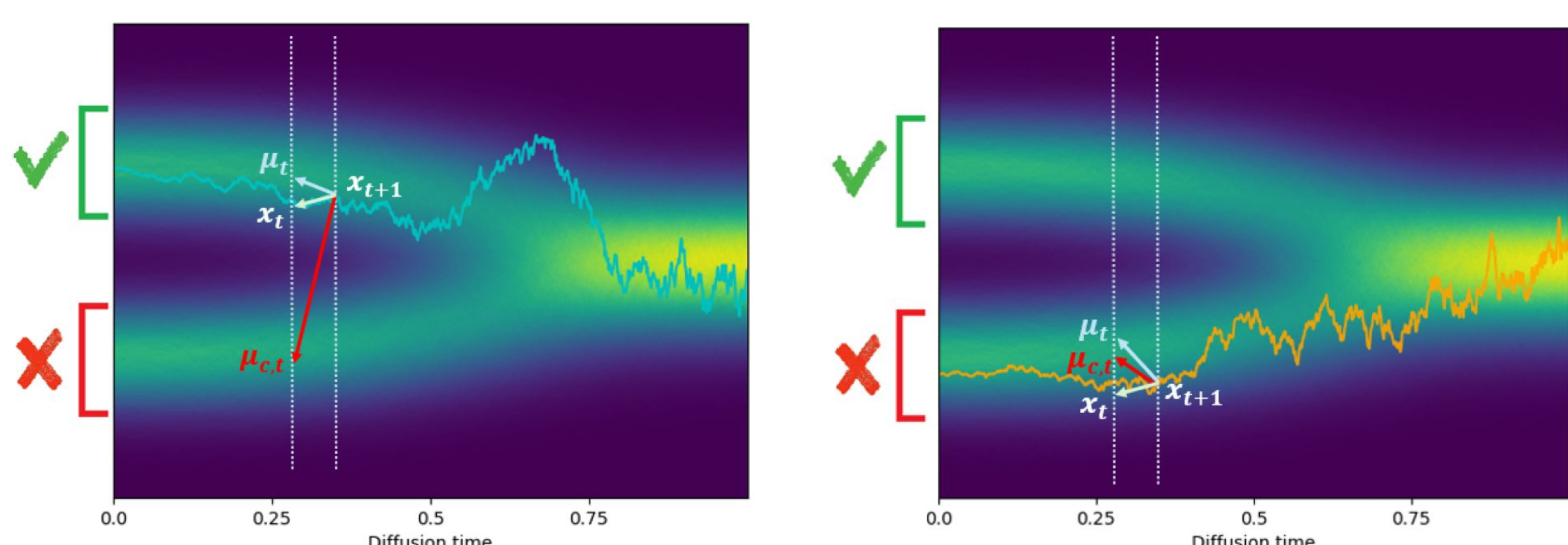


## Posterior approximation

Unfortunately, the posterior likelihood is most often unavailable at inference. We approximate it by tracking the Markov Chain along the diffusion trajectory. We find a recursive formula that solely depends on already computed quantities:

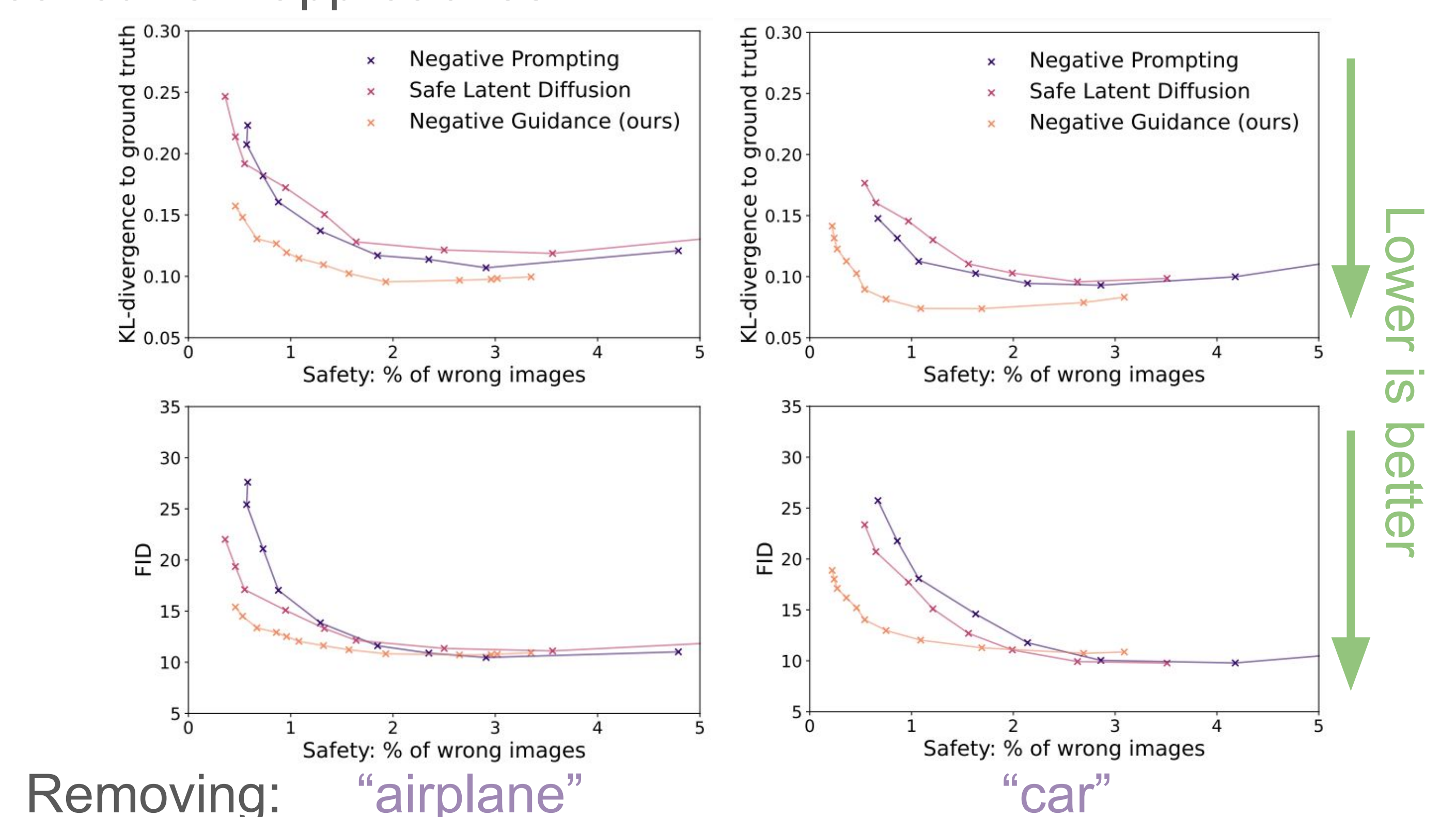
$$\log p_t(\mathbf{c} | \mathbf{x}_t) = \log p_{t+1}(\mathbf{c} | \mathbf{x}_{t+1}) - \frac{1}{2\sigma_t^2} (\|\mathbf{x}_t - \mu_{t,\theta}(\mathbf{x}_{t+1} | \mathbf{c})\|^2 - \|\mathbf{x}_t - \mu_{t,\theta}(\mathbf{x}_{t+1})\|^2) \quad (4)$$

neg. prompted      unconditional



## Class-removal Experiments

To evaluate how well the diversity of a model is affected by various negative guidance scheme, we propose a class removal setting. The goal is to avoid the generation of a specific class using solely an unconditional model and a model trained on that single class. DNG outperforms concurrent approaches

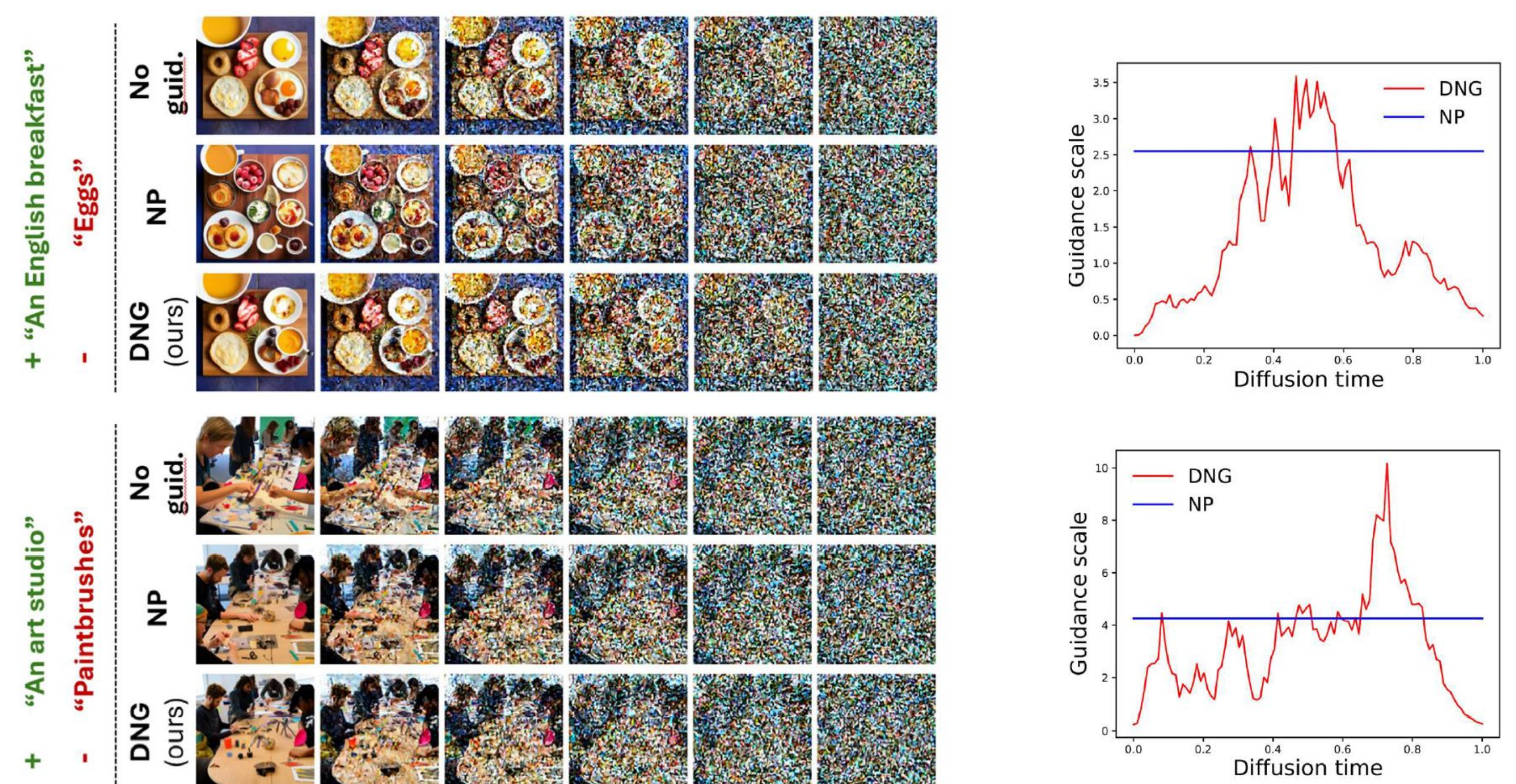


Removing: "airplane"

"car"

## T2I Experiments

In the context of T2I, our *dynamic guidance scale* preserves image diversity by deactivating itself for unrelated negative prompts. Here are some illustrative trajectories:



## Conclusion

The theoretically optimal *dynamic* negative guidance scale can be approximated by tracking the posterior using the discrete markov chain to improve sample diversity when using negative prompting-like techniques. Our self-regulated guidance scale is further trajectory and prompt specific, providing further insights into when guidance is required

Further questions?  
Scan this or ask us!

