

Towards Understanding the Mechanism of CFG

Xiang Li, Rongrong Wang, Qing Qu

I. INTRODUCTION

Problem: Conditional diffusion models generate images by progressively denoising a random noise $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \sigma(T)^2 \mathbf{I})$ to a clean image with the probabilistic ODE:

$$d\mathbf{x} = -\sigma(t)\nabla_{\mathbf{x}} \log p(\mathbf{x}|\mathbf{c}; \sigma(t))dt, \quad (1)$$

where $\sigma(t)$ is a predefined schedule. However, this standard (naive) sampling approach does not lead to high quality samples, as shown in Figure 1 (left). In contrast, state of the art diffusion models generate high quality images with Classifier-Free-Guidance (CFG):

$$d\mathbf{x} = -\sigma(t)\nabla_{\mathbf{x}} (\log p(\mathbf{x}|\mathbf{c}; \sigma(t))dt + \gamma g(\mathbf{x}, t)dt), \quad (2)$$

$g(\mathbf{x}, t) = \nabla_{\mathbf{x}} \log p(\mathbf{x}|\mathbf{c}; \sigma(t)) - \nabla_{\mathbf{x}} \log p(\mathbf{x}|\sigma(t))$ is the difference between conditional and unconditional scores. With CFG, the quality of conditional generation greatly improves, as shown in Figure 1 (right). This work aims to understand the underlying mechanisms of CFG and we choose linear diffusion models as a prototype since CFG has similar effects on them, as shown in the right part of Figure 1.

II. CFG IN LINEAR MODELS

In linear diffusion models, CFG guidance $g(\mathbf{x}; \sigma(t))$ can be decomposed as:

$$(\tilde{\Sigma}_{c,t} - \tilde{\Sigma}_{uc,t})(\mathbf{x}_t - \boldsymbol{\mu}_c) + \gamma(\mathbf{I} - \tilde{\Sigma}_{uc,t})(\mathbf{u}_c - \mathbf{u}_{uc}). \quad (3)$$

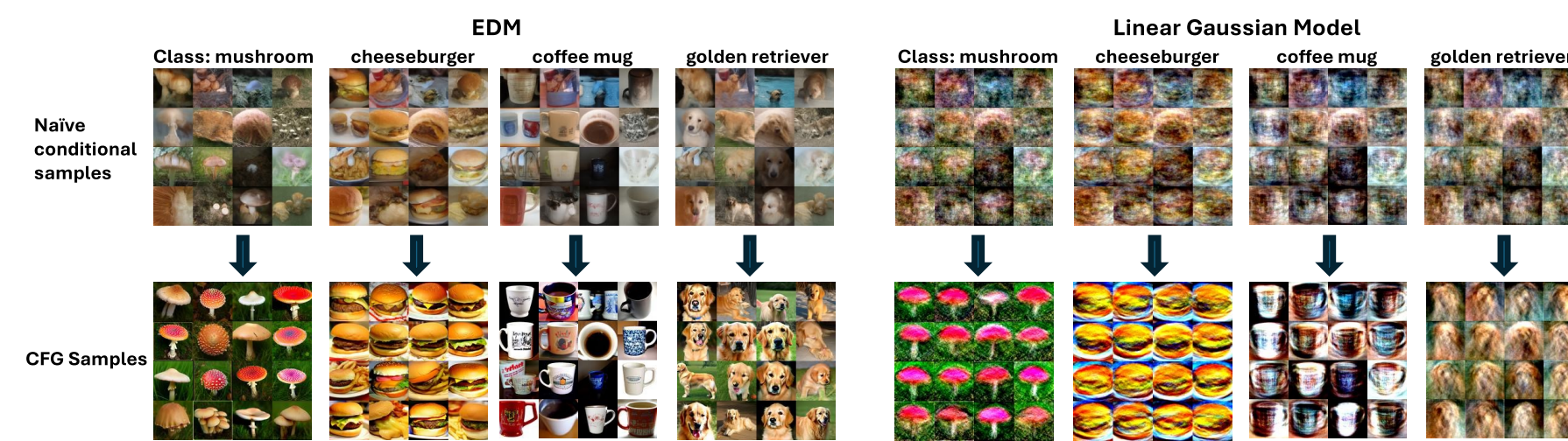


Figure 1: Left and right parts demonstrate CFG's Effect on nonlinear and linear diffusion models respectively.

Here:

$$\tilde{\Sigma}_{c,t} - \tilde{\Sigma}_{uc,t} = V_+ \Sigma_+ V_+ + V_- \Sigma_- V_- \quad (4)$$

, which is the difference between conditional and unconditional data covariance, and it can be decomposed with eigendecomposition, separating the positive eigenvector (positive CPC) and negative eigenvector (negative CPC). On the other hand, $\gamma(\mathbf{I} - \tilde{\Sigma}_{uc,t})(\mathbf{u}_c - \mathbf{u}_{uc})$ approximately shifts the sample towards the difference between the conditional mean $\boldsymbol{\mu}_c$ and the unconditional mean $\boldsymbol{\mu}_{uc}$.

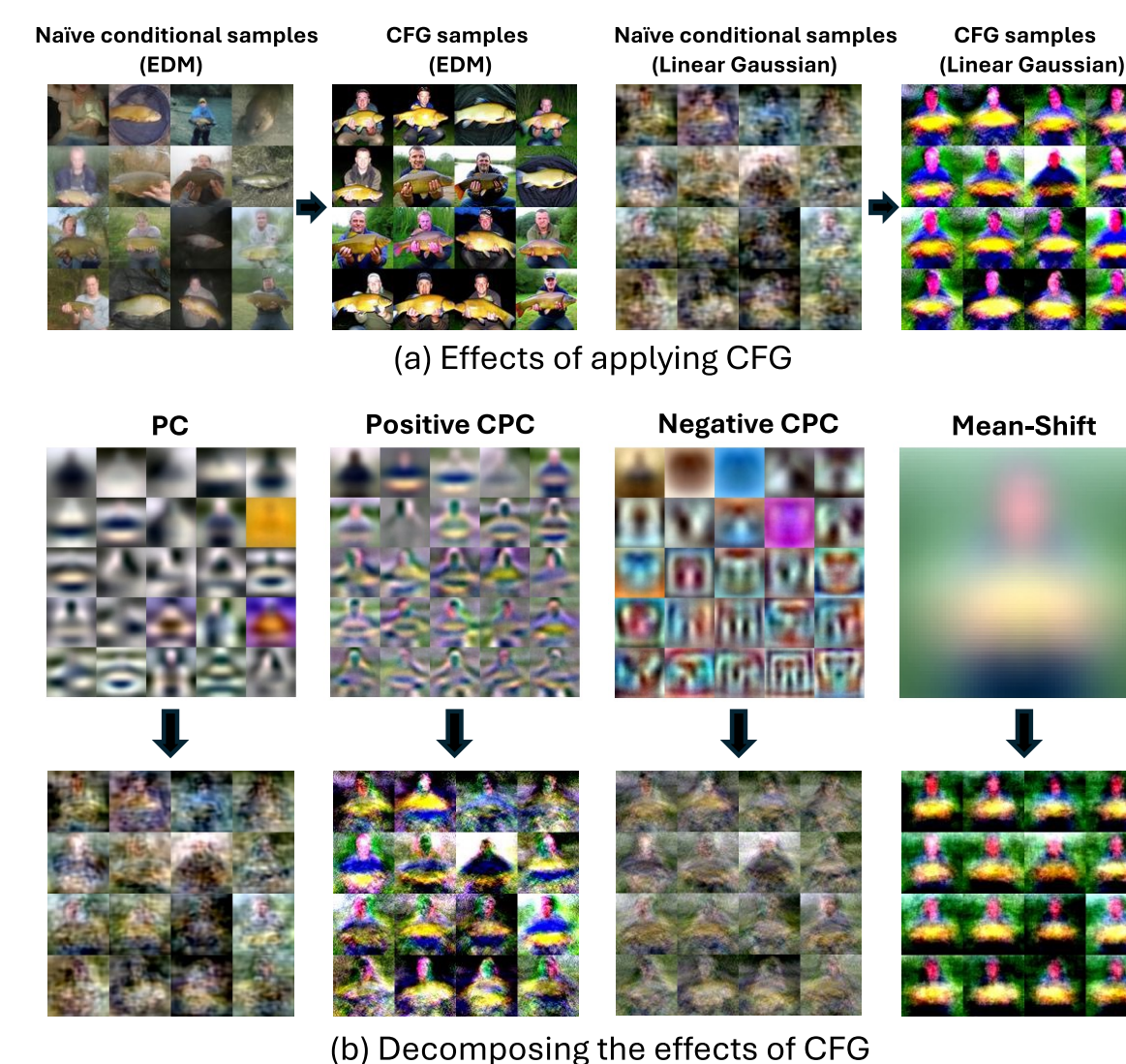


Figure 2: Distinct effects of different CFG components.

By decomposing CFG in this way, we can separately examine their distinctive effects, which are shown in Figure 2.

In short, the positive CPC guidance enhances class specific features of the dataset, the negative CPC guidance suppresses conditional-unrelated features in the unconditional dataset and the mean-shift guidance approximately shifts the samples towards the class mean.

III. CFG IN NONLINEAR MODELS

For a wide range of noise levels (high noise), diffusion models exhibit linearity and applying linear CFG leads to similar effects as the actual CFG as shown in Figure 3. On the other hand, for low noise levels, diffusion models are highly non-linear. In this regime, we can construct CPC guidance from the denoiser's Jacobians, which represents the posterior covariances. Such guidance could lead to similar effects as the actual nonlinear CFG (see the paper for more details).

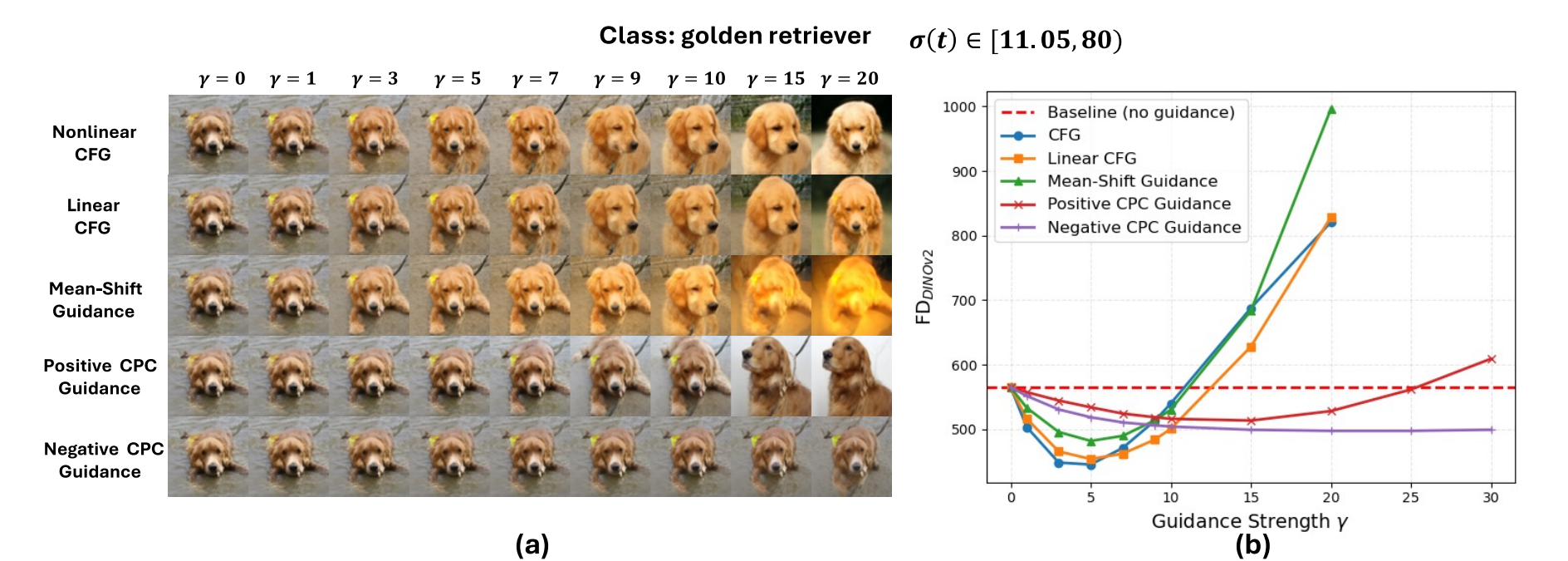


Figure 3: Distinct effects of different CFG components.

The main takeaway is that, CFG contrasts between the conditional and unconditional data, highlighting class-specific information while suppressing unrelated ones.