

Tackling the Generative Learning Trilemma with DENOISING DIFFUSION GANS

Vikram Voleti
Christopher Beckham
Jae Hyun Lim
PhD students, Mila

TACKLING THE GENERATIVE LEARNING TRILEMMA WITH DENOISING DIFFUSION GANS

Zhisheng Xiao*

The University of Chicago
zxiao@uchicago.edu

Karsten Kreis

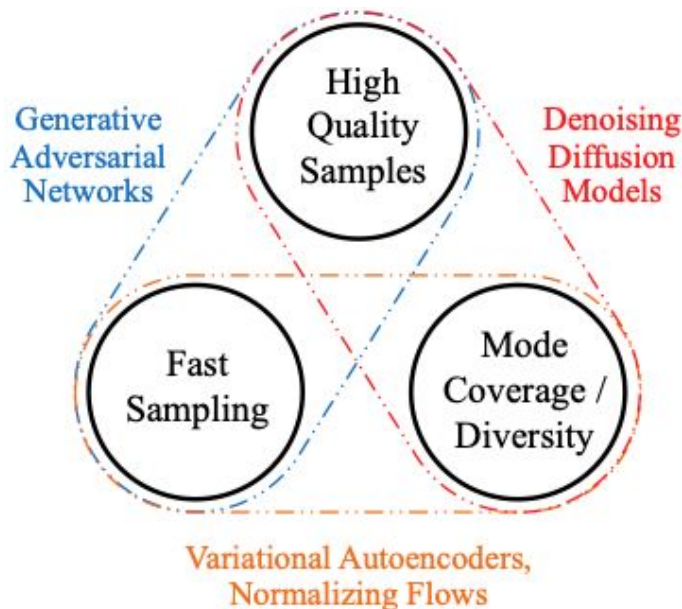
NVIDIA
kkreis@nvidia.com

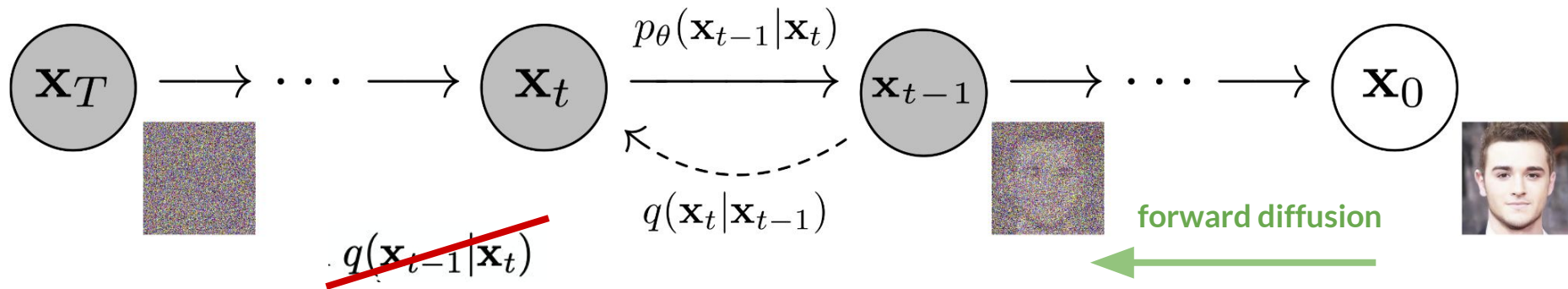
Arash Vahdat

NVIDIA
avahdat@nvidia.com

ABSTRACT

A wide variety of deep generative models has been developed in the past decade. Yet, these models often struggle with simultaneously addressing three key requirements including: high sample quality, mode coverage, and fast sampling. We call the challenge imposed by these requirements *the generative learning trilemma*, as the existing models often trade some of them for others. Particularly, denoising diffusion models have shown impressive sample quality and diversity, but their expensive sampling does not yet allow them to be applied in many real-world applications. In this paper, we argue that slow sampling in these models is fundamentally attributed to the Gaussian assumption in the denoising step which is justified only for small step sizes. To enable denoising with large steps, and hence, to reduce the total number of denoising steps, we propose to model the denoising distribution using a complex multimodal distribution. We introduce *denoising diffusion generative adversarial networks* (*denoising diffusion GANs*) that model each denoising step using a multimodal conditional GAN. Through extensive evaluations, we show that denoising diffusion GANs obtain sample quality and diversity competitive with original diffusion models while being 2000 \times faster on the CIFAR-10 dataset. Compared to traditional GANs, our model exhibits better mode coverage and sample diversity. To the best of our knowledge, denoising diffusion GAN is the first model that reduces sampling cost in diffusion models to an extent that allows them to be applied to real-world applications inexpensively. Project page and code: <https://nvlabs.github.io/denoising-diffusion-gan>.



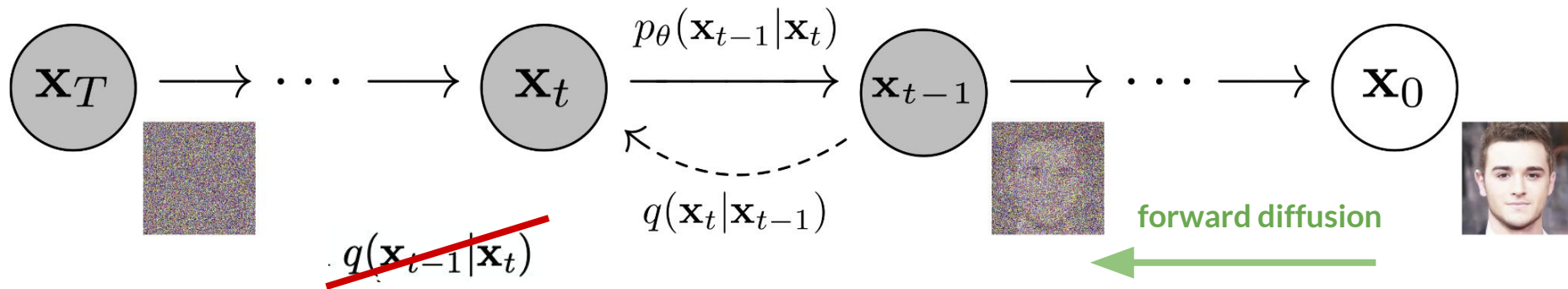


In diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020), there is a forward process that gradually adds noise to the data $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ in T steps with pre-defined variance schedule β_t :

$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t \geq 1} q(\mathbf{x}_t|\mathbf{x}_{t-1}), \quad q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}), \quad (1)$$

where $q(\mathbf{x}_0)$ is a data-generating distribution. The reverse denoising process is defined by:

$$p_\theta(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t \geq 1} p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t), \quad p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I}), \quad (2)$$



$$\mathcal{L} = - \sum_{t \geq 1} \mathbb{E}_{q(\mathbf{x}_t)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t) \| p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))] + C,$$

Denoising distribution $p(\cdot)$ is assumed to be Gaussian, and with infinitesimal step sizes, it can also be shown that the reversal of q (the denoising form of q) is as well

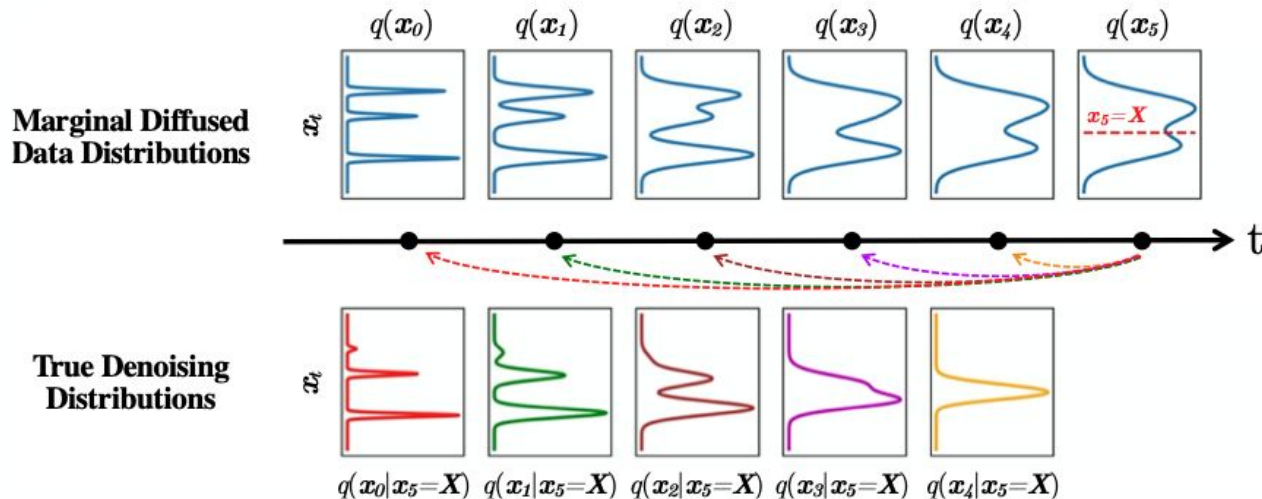
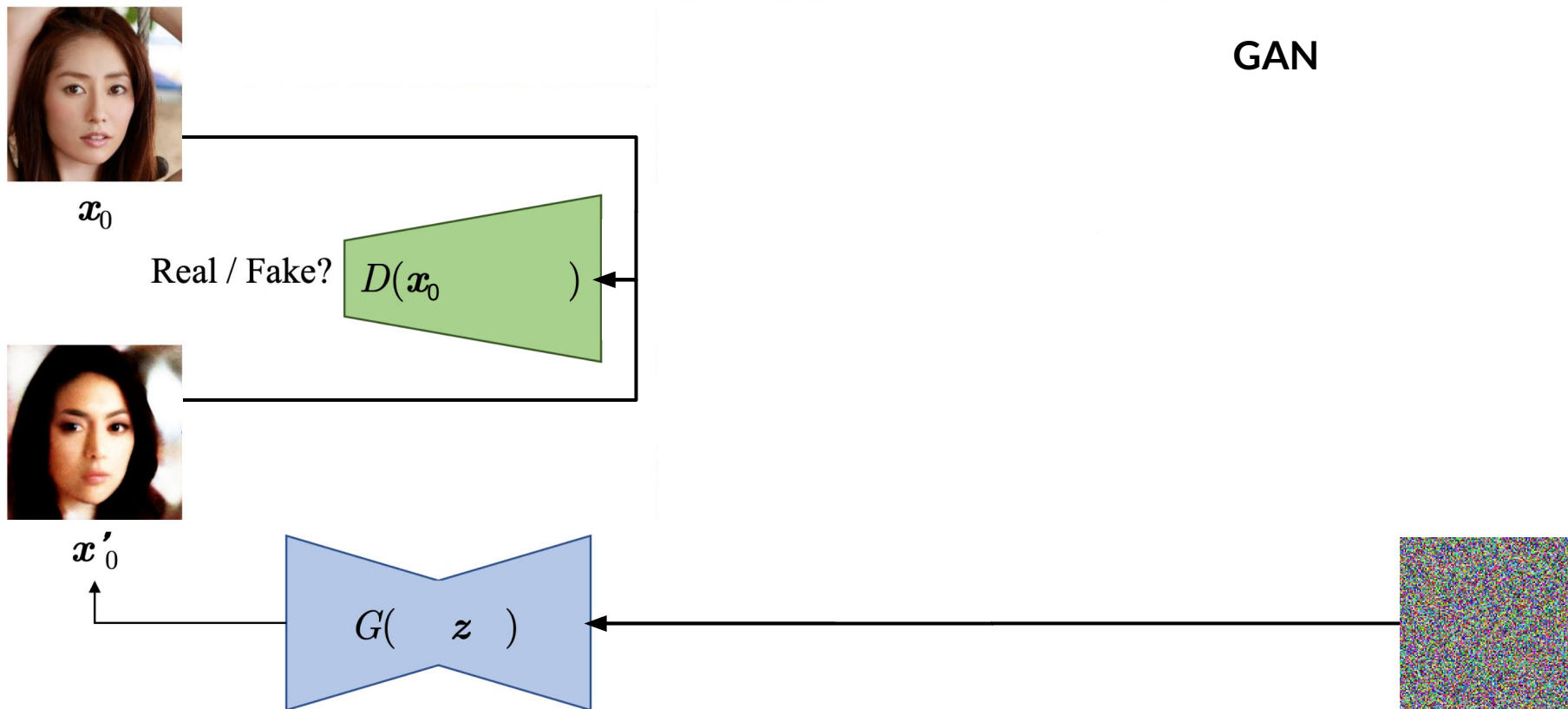
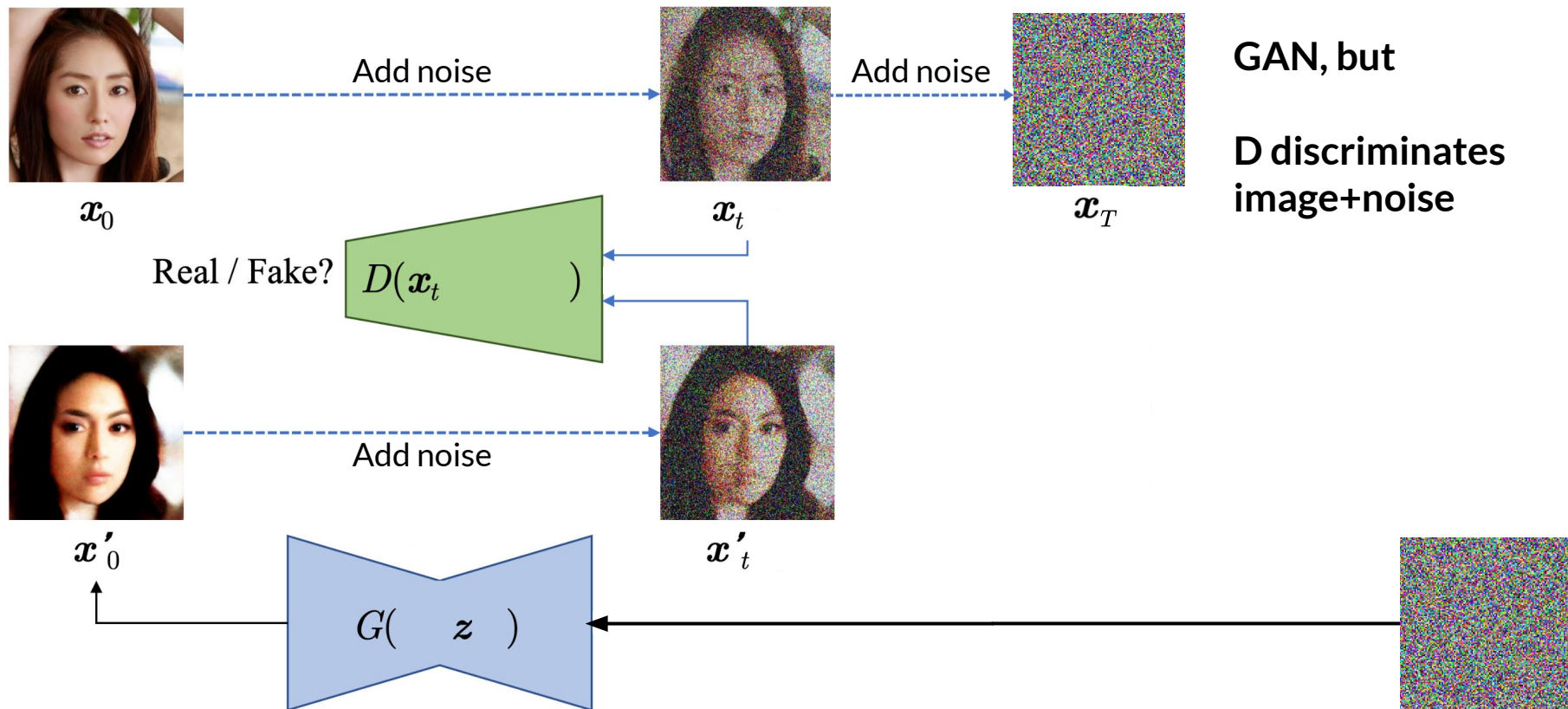
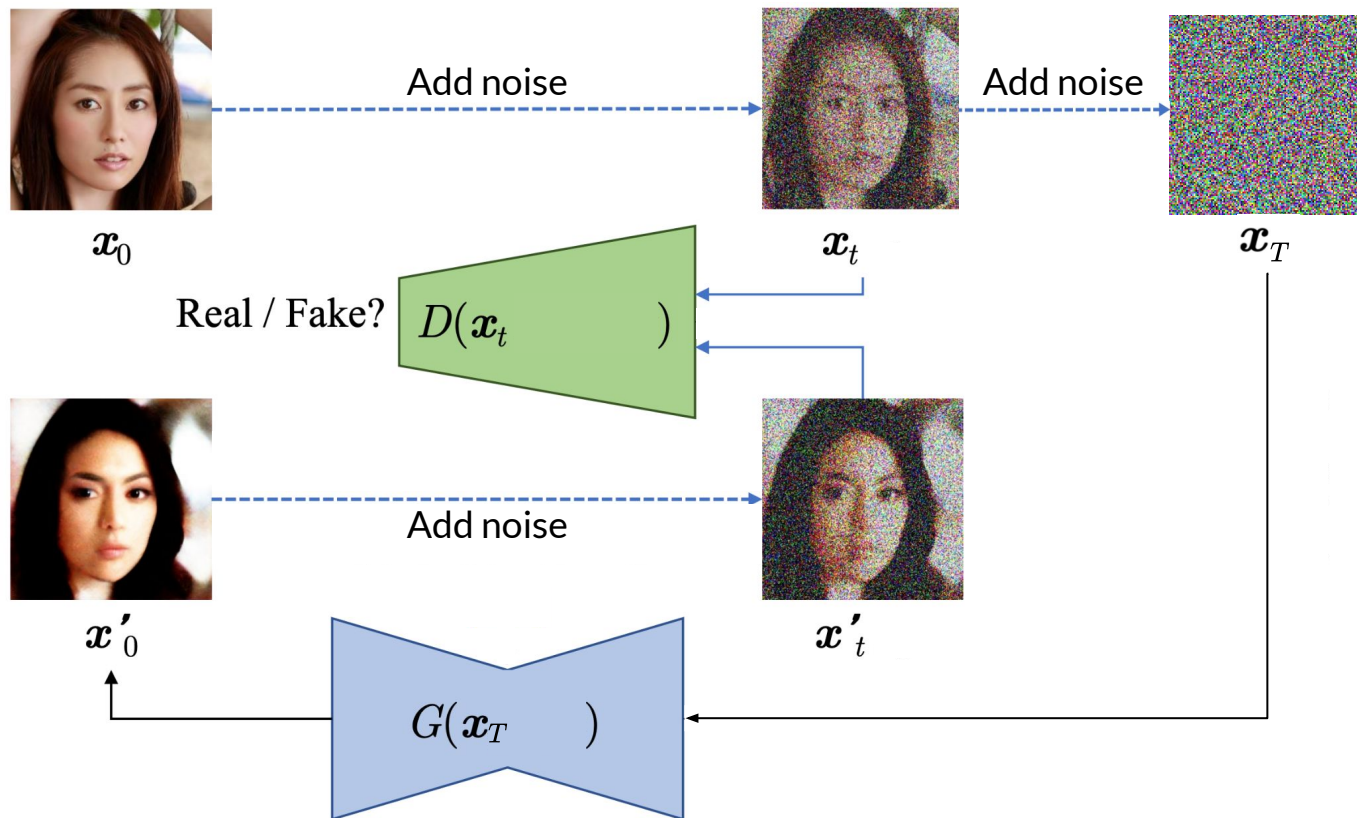


Figure 2: **Top:** The evolution of 1D data distribution $q(\mathbf{x}_0)$ through the diffusion process. **Bottom:**, The visualization of the true denoising distribution for varying step sizes conditioned on a fixed \mathbf{x}_5 . The true denoising distribution for a small step size (i.e., $q(\mathbf{x}_4 | \mathbf{x}_5 = \mathbf{X})$) is close to a Gaussian distribution. However, it becomes more complex and multimodal as the step size increases.

GAN



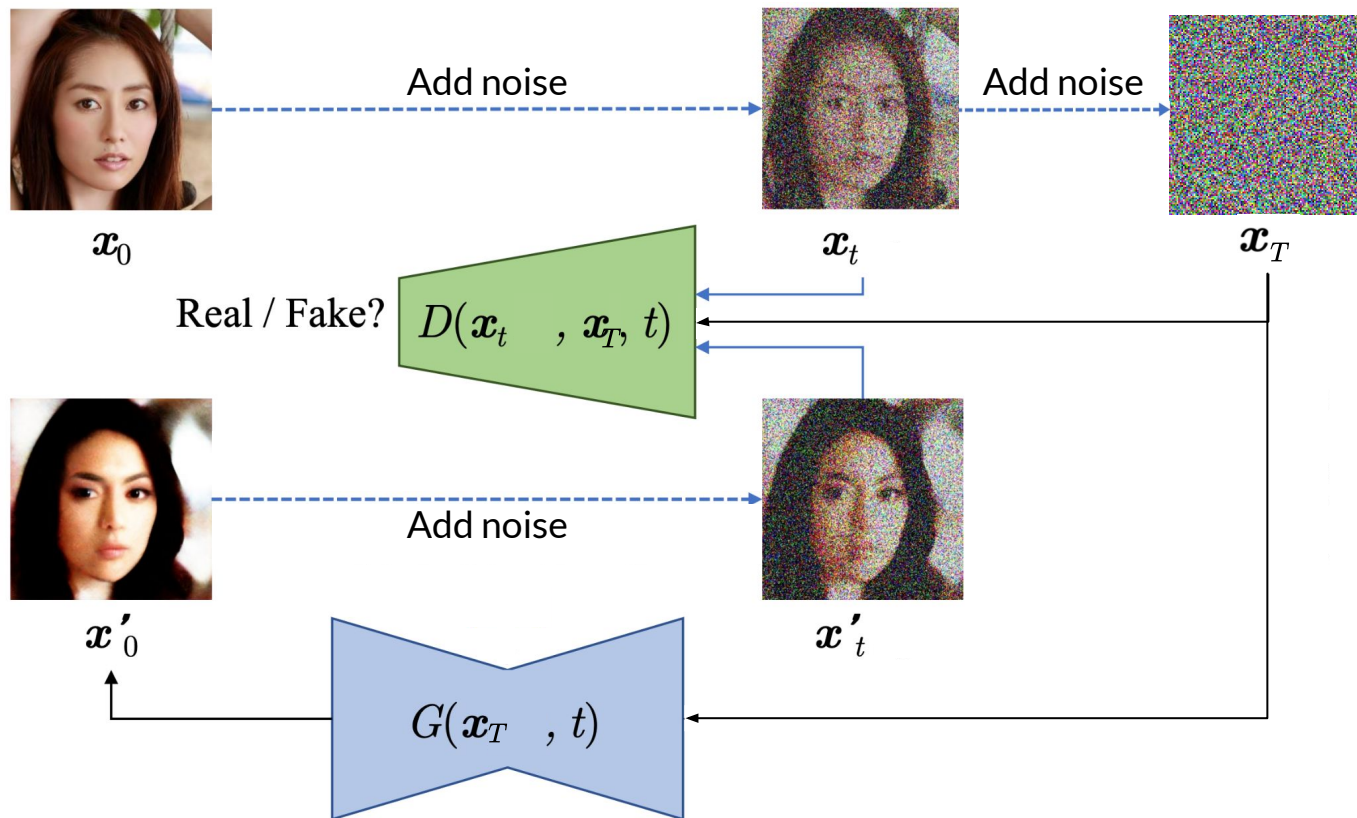




GAN, but

D discriminates
image+noise,

input to G is
real image + noise

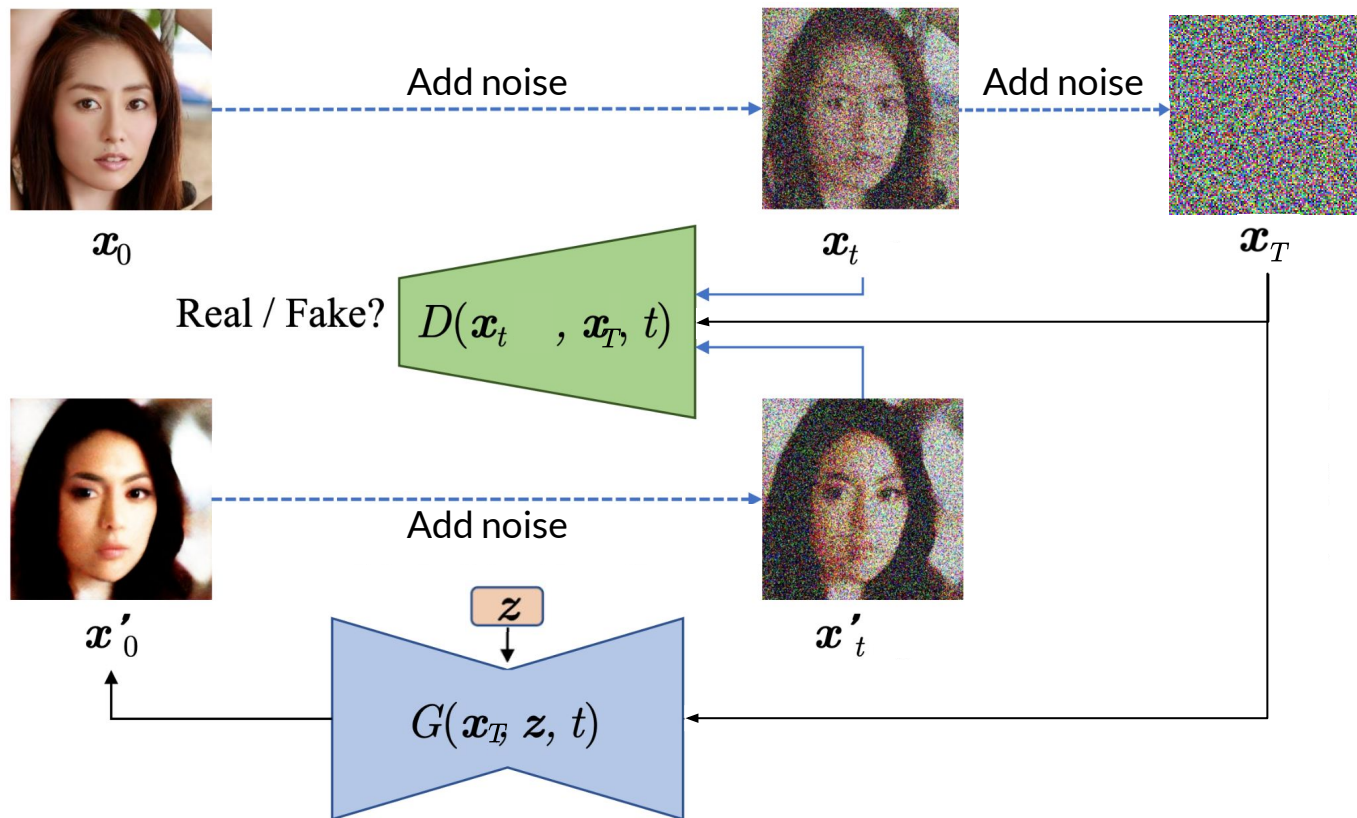


GAN, but

**D discriminates
image+noise,**

**input to G is
real image + noise,**

**D is conditioned on
real image + noise**



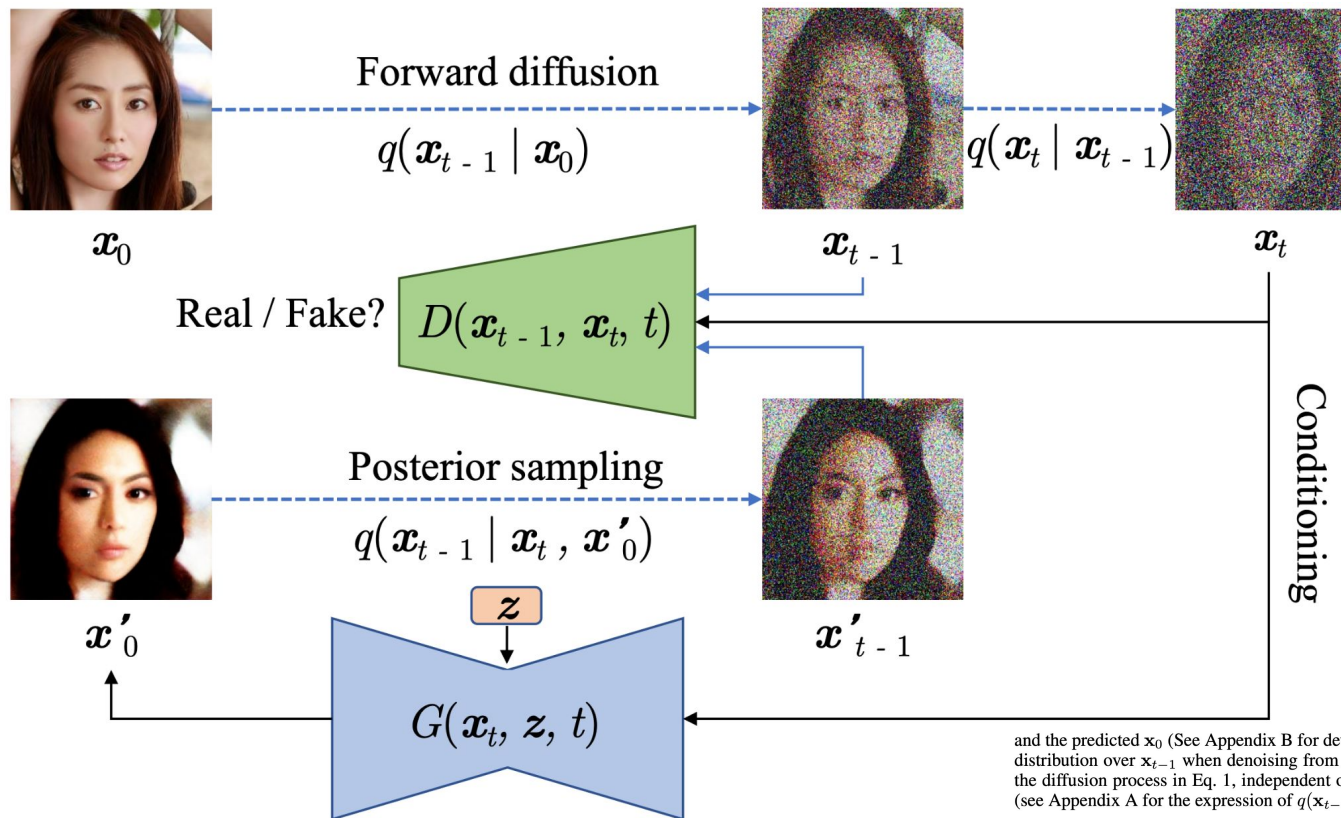
GAN, but

D discriminates
image+noise,

input to G is
real image + noise,

D is conditioned on
real image + noise,

G is conditioned on
random noise



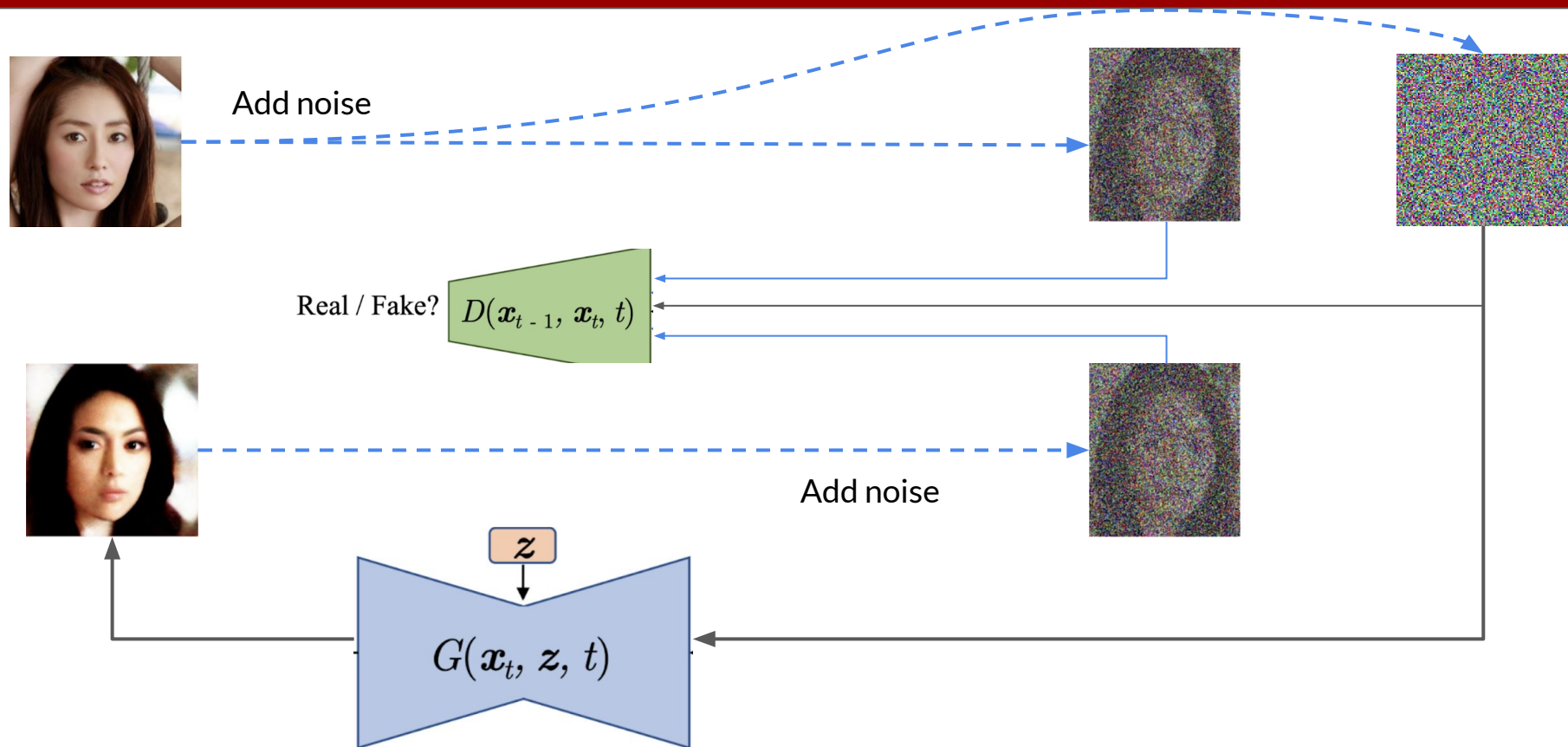
Iterate, et voila!

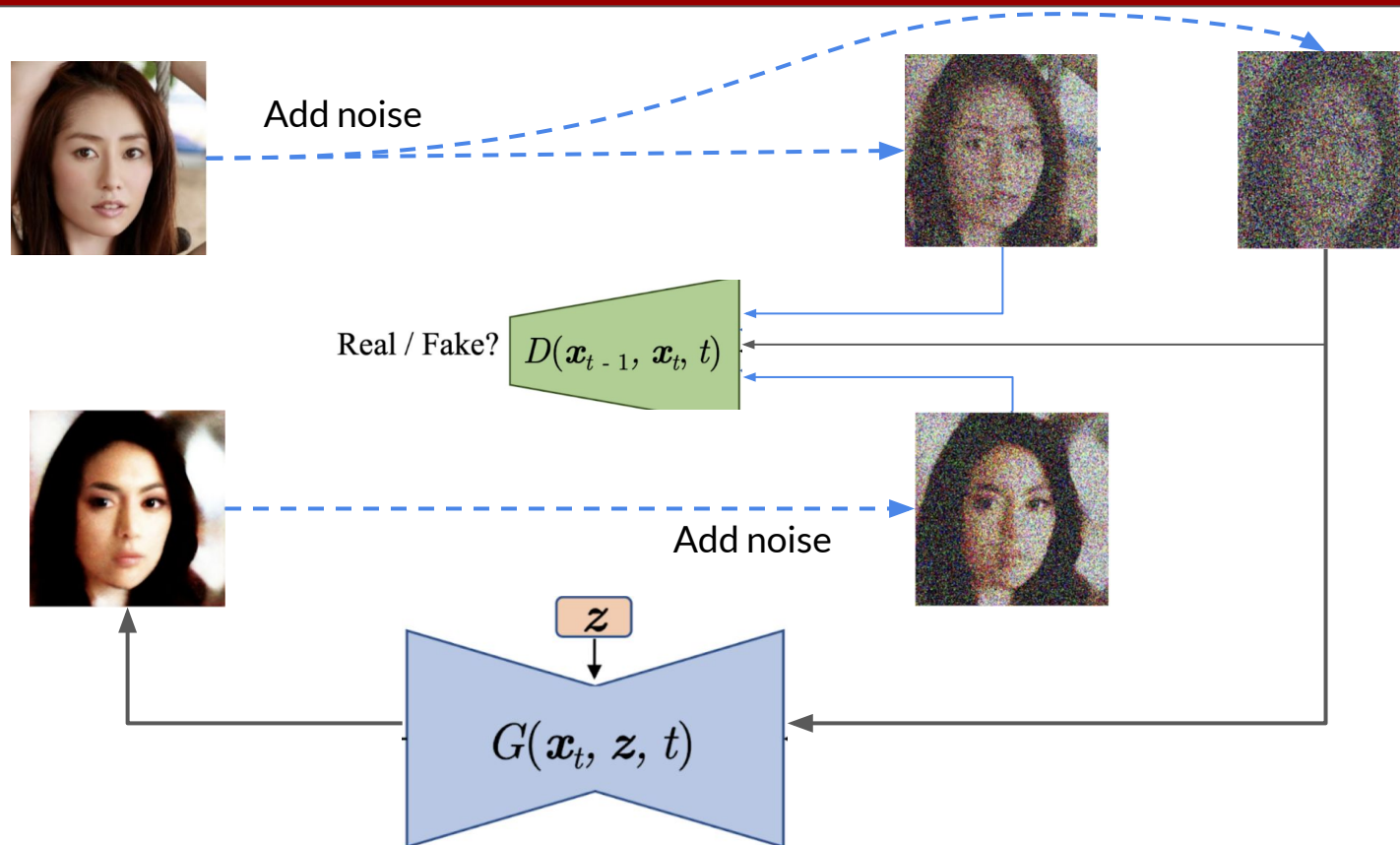
**Denoising
Diffusion GAN!**

and the predicted x_0 (See Appendix B for details). The distribution $q(x_{t-1} | x_0, x_t)$ is intuitively the distribution over x_{t-1} when denoising from x_t towards x_0 , and it always has a Gaussian form for the diffusion process in Eq. 1, independent of the step size and complexity of the data distribution (see Appendix A for the expression of $q(x_{t-1} | x_0, x_t)$). Similarly, we define $p_\theta(x_{t-1} | x_t)$ by:

$$p_\theta(x_{t-1} | x_t) := \int p_\theta(x_0 | x_t) q(x_{t-1} | x_0, x_t) dx_0 = \int p(z) q(x_{t-1} | x_t, x_0 = G_\theta(x_t, z, t)) dz, \quad (6)$$

where $p_\theta(x_0 | x_t)$ is the implicit distribution imposed by the GAN generator $G_\theta(x_t, z, t) : \mathbb{R}^N \times \mathbb{R}^L \times \mathbb{R} \rightarrow \mathbb{R}^N$ that outputs x_0 given x_t and an L -dimensional latent variable $z \sim p(z) := \mathcal{N}(z; 0, \mathbf{I})$.





— ground-truth pdf

— perturbed

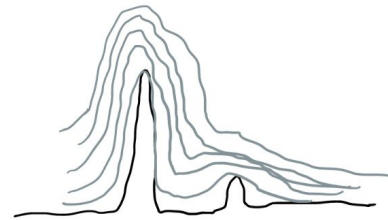
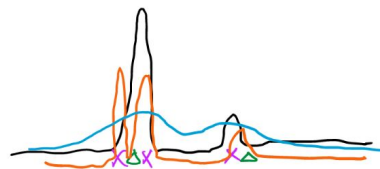
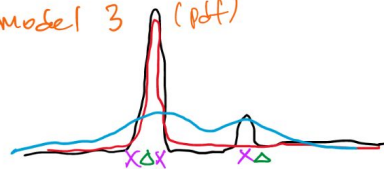
X training data

△ test data

— model 1 (pdf)

— model 2 (pdf)

— model 3 (pdf)



ex) GAN w/o reg

$$D_{KL}(\text{model 1} \parallel g_t) < D_{KL}(\text{model 2} \parallel g_t)$$

ex) VAE

$$D_{KL}(g_t \parallel \text{model 1}) > D_{KL}(g_t \parallel \text{model 2})$$

ex) SM

$$SM(g_t \parallel \text{model 3}) > SM(g_t \parallel \text{model 2})$$
 (or model 1)
 = matching the gradient.

Thank you!