

Analysing

Adversarial Robustness of VAE and

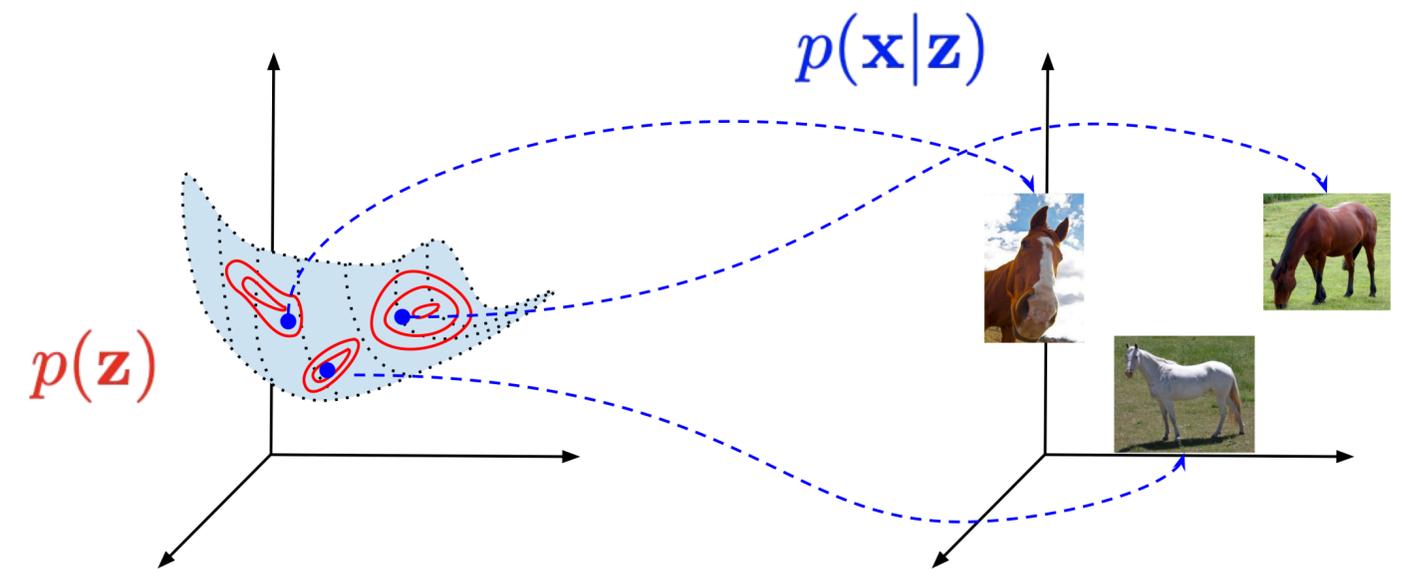
Denoising Abilities of DGM

Anna Kuzina

Vrije Universiteit Amsterdam

Latent Variable Models

Observing a finite sample x_1, \dots, x_N we are interested in the underlying distribution $p(x)$



Latent Variable Models

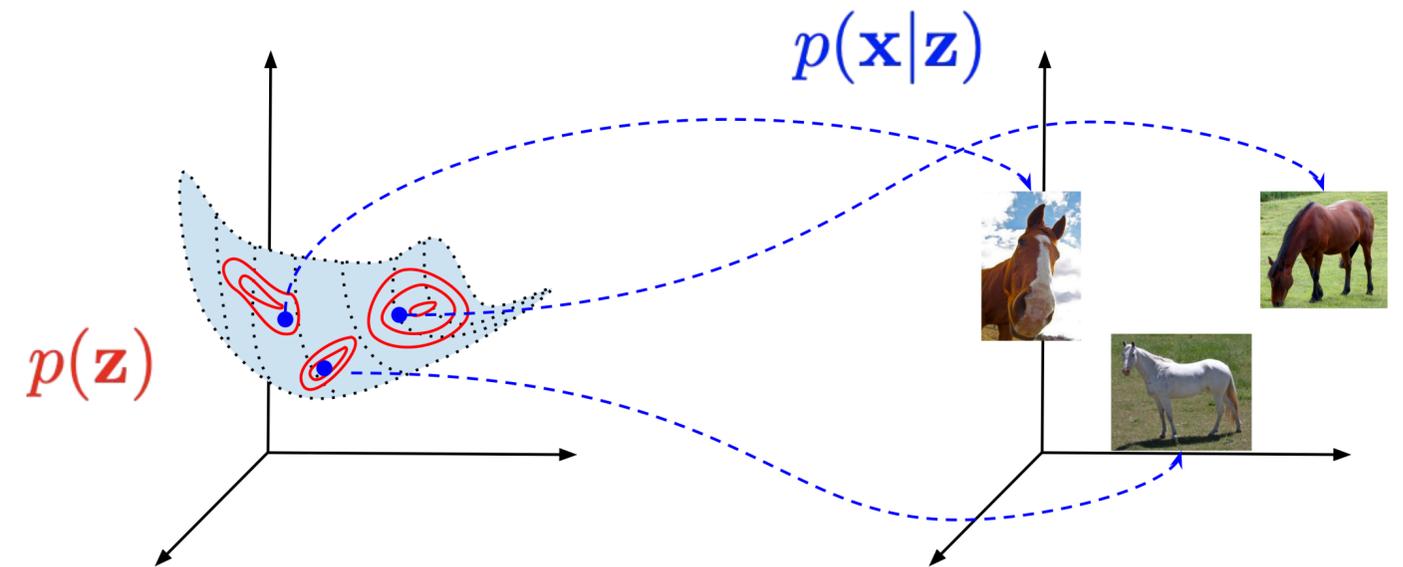
Observing a finite sample x_1, \dots, x_N we are interested in the underlying distribution $p(x)$

Assume the model

$$p_{\theta}(x) = \int p_{\theta}(x|z)p_{\theta}(z)dz$$

aka generative process:

$$z \sim p_{\theta}(z)$$
$$x \sim p_{\theta}(x|z)$$



Latent Variable Models

Observing a finite sample x_1, \dots, x_N we are interested in the underlying distribution $p(x)$

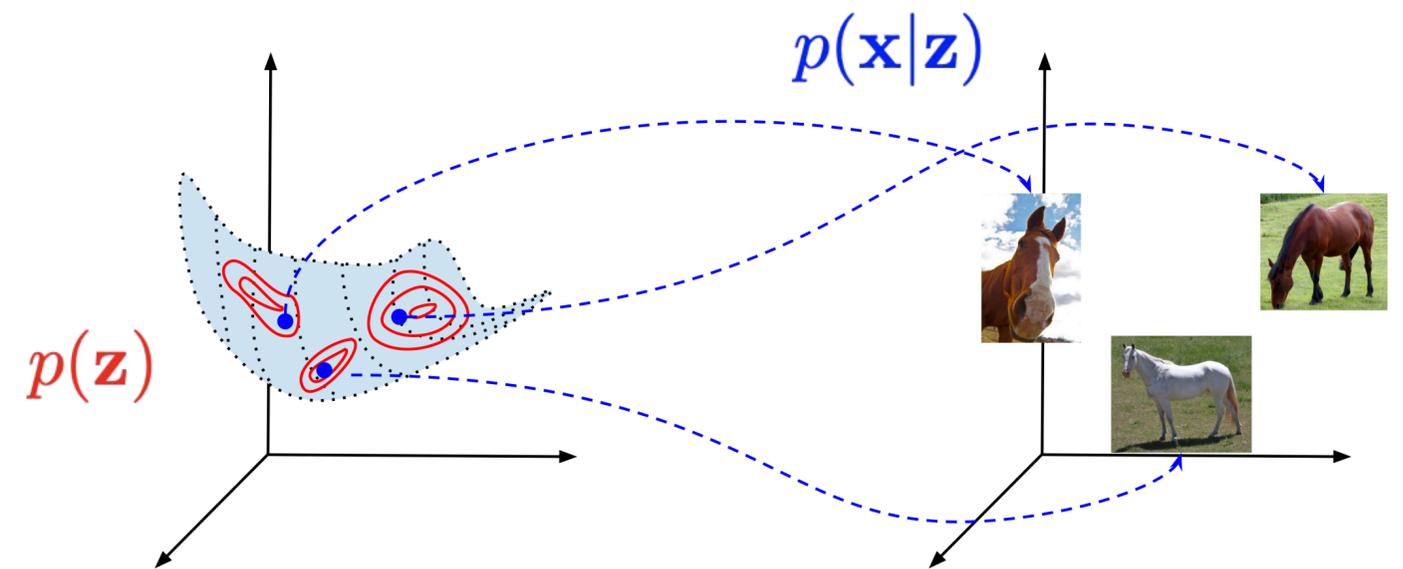
Assume the model

$$p_{\theta}(x) = \int p_{\theta}(x|z)p_{\theta}(z)dz$$

aka generative process:

$$z \sim p_{\theta}(z)$$
$$x \sim p_{\theta}(x|z)$$

Unknown parameters



Latent Variable Models

Observing a finite sample x_1, \dots, x_N we are interested in the underlying distribution $p(x)$

Assume the model

$$p_{\theta}(x) = \int p_{\theta}(x|z)p_{\theta}(z)dz$$

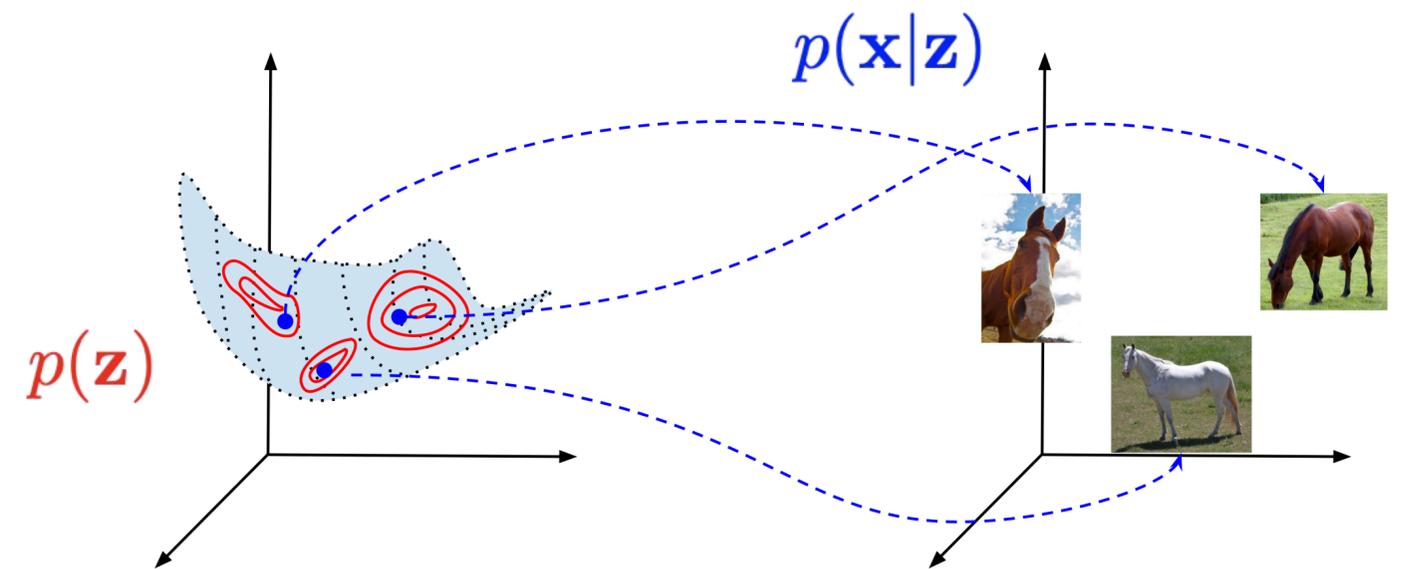
aka generative process:

$$z \sim p_{\theta}(z)$$

$$x \sim p_{\theta}(x|z)$$

MLE objective:

$$\max_{\theta} \sum_n \ln p_{\theta}(x_n)$$



Latent Variable Models

Observing a finite sample x_1, \dots, x_N we are interested in the underlying distribution $p(x)$

Assume the model

$$p_{\theta}(x) = \int p_{\theta}(x|z)p_{\theta}(z)dz$$

aka generative process:

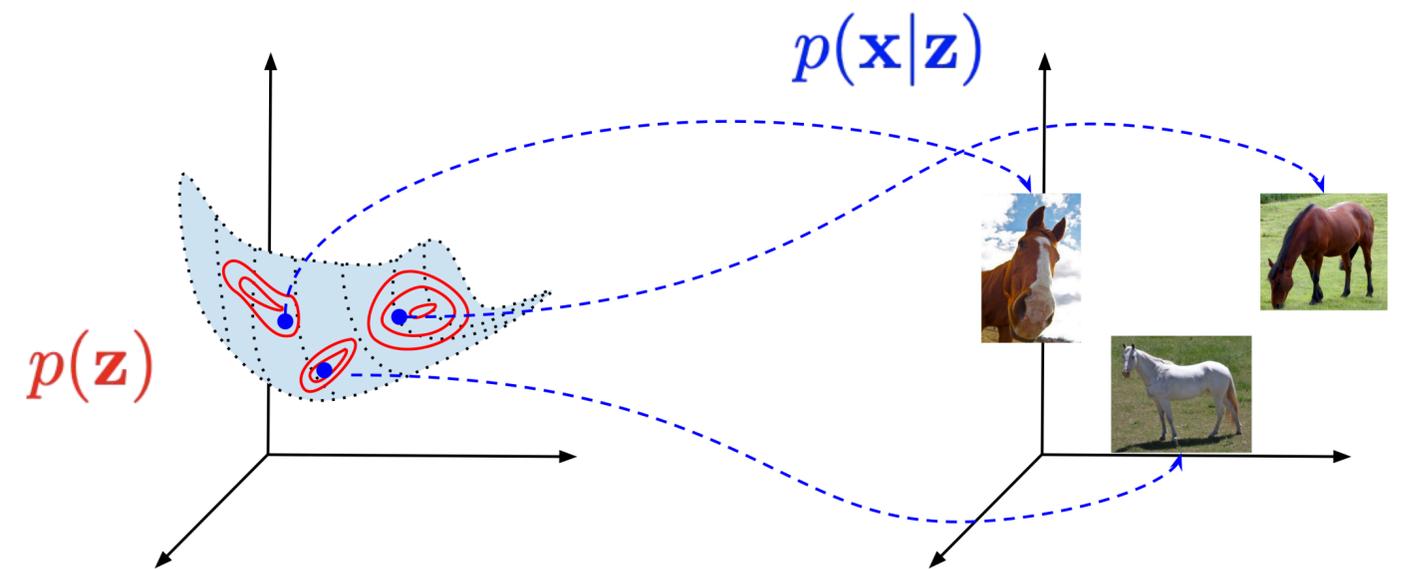
$$z \sim p_{\theta}(z)$$

$$x \sim p_{\theta}(x|z)$$

MLE objective:

$$\max_{\theta} \sum_n \ln p_{\theta}(x_n)$$

intractable



Variational Autoencoder

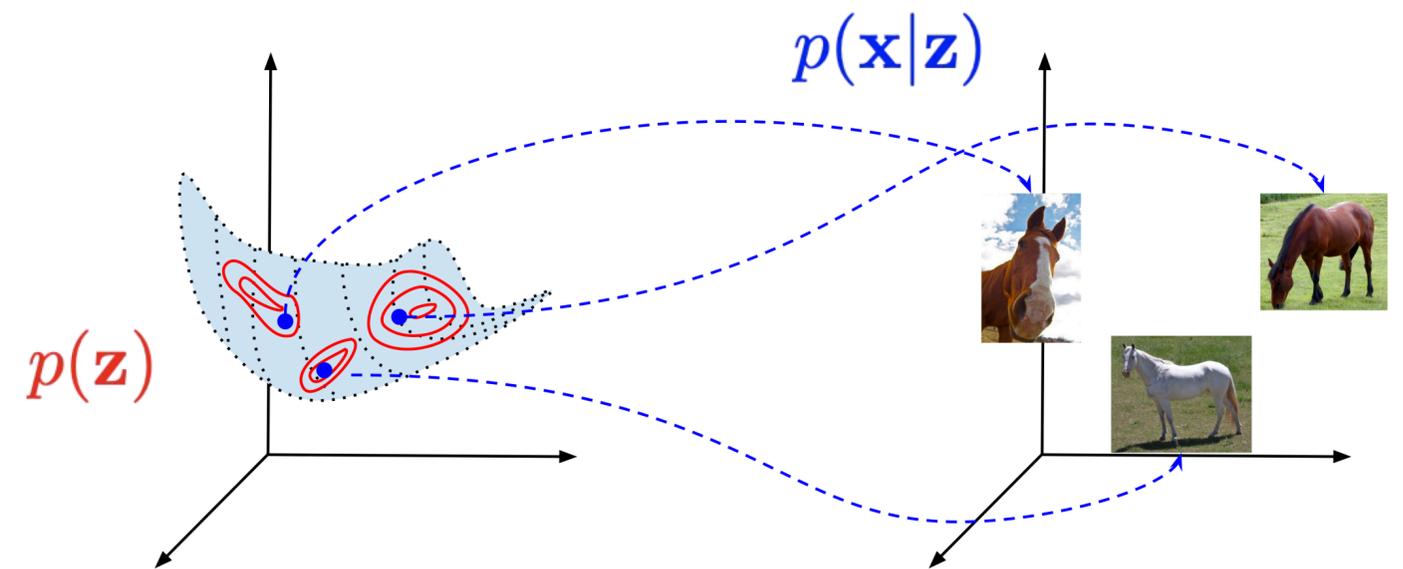
Generative Model:

$$p_{\theta}(x|z)p_{\theta}(z)$$

Let us define one more model:

$$q_{\phi}(z|x)$$

More
unknown
parameters



Variational Autoencoder

Generative Model:

$$p_{\theta}(x|z)p_{\theta}(z)$$

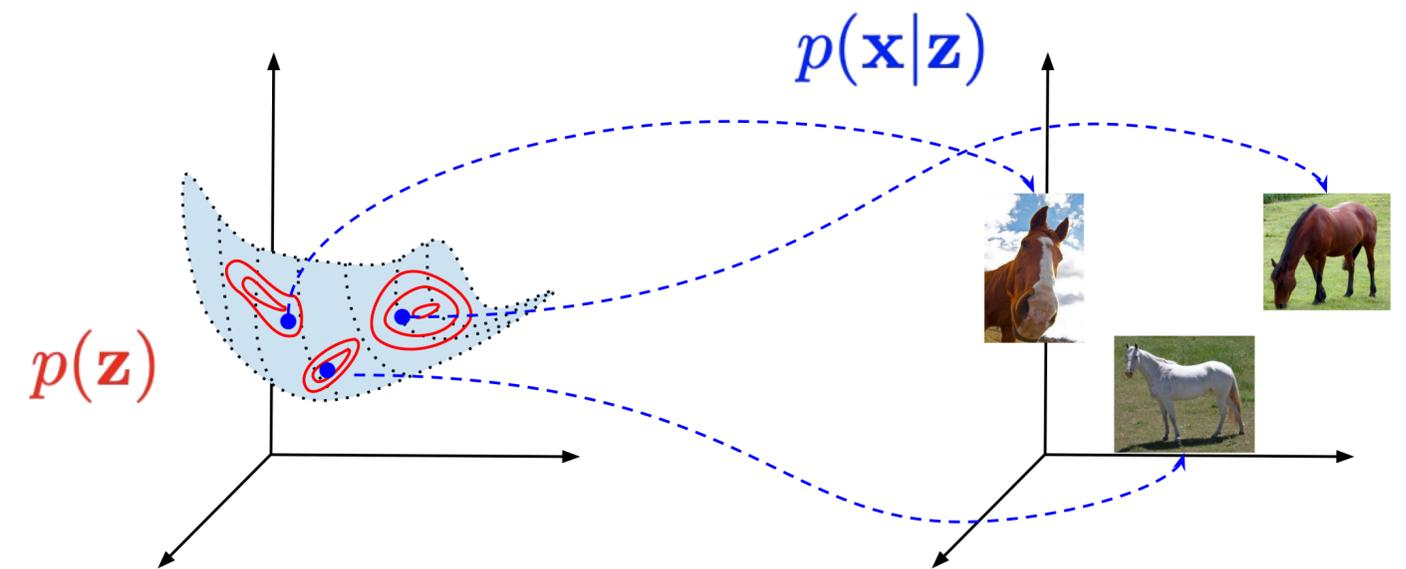
Let us define one more model:

$$q_{\phi}(z|x)$$

Tractable objective:

$$\ln p_{\theta}(x) \geq \underbrace{\mathbb{E}_{q_{\phi}(z|x)} \ln p_{\theta}(x|z) - \text{KL}[q_{\phi}(z|x) || p_{\theta}(z)]}_{\text{ELBO}}$$

ELBO



Variational Autoencoder

Generative Model:

$$p_{\theta}(x|z)p_{\theta}(z)$$

Let us define one more model:

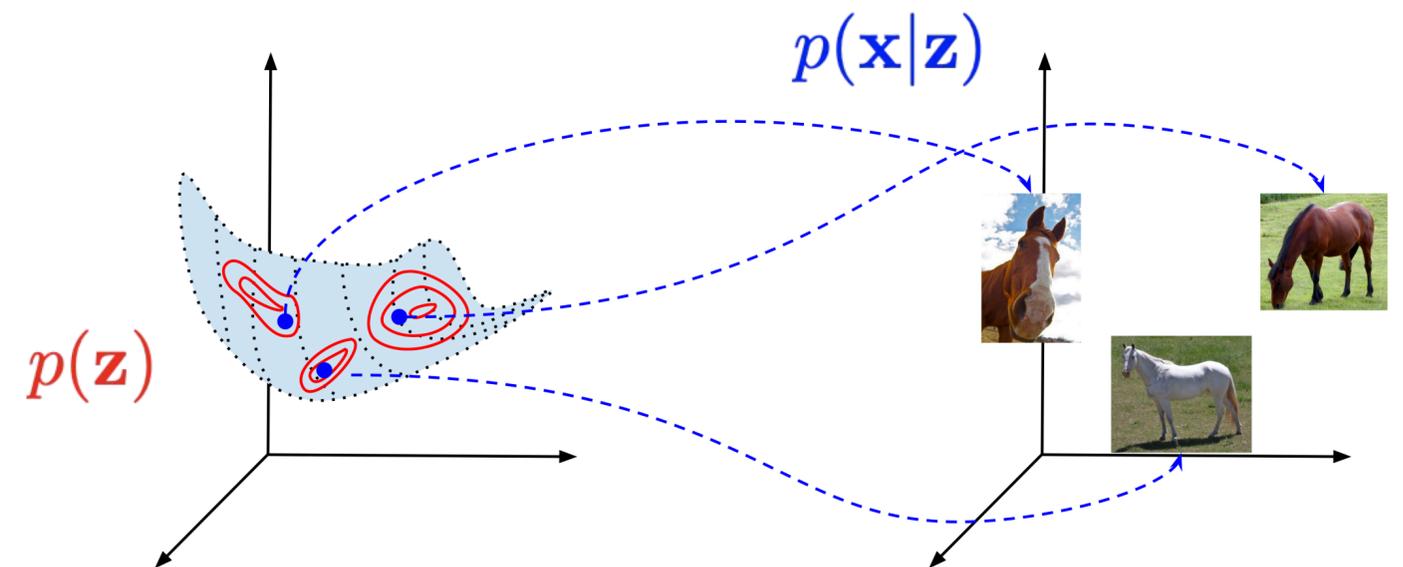
$$q_{\phi}(z|x)$$

- Inference Model
- Variational Posterior
- Encoder

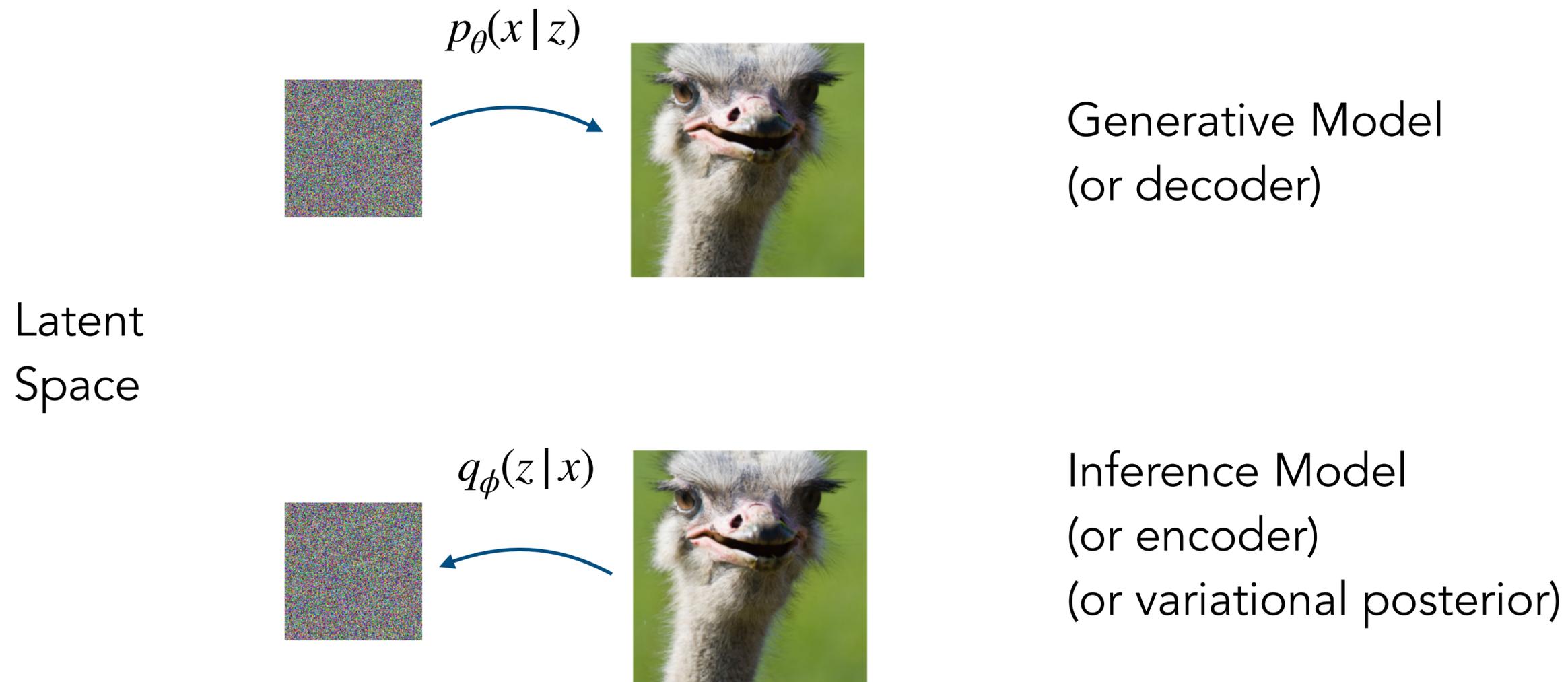
Turns out it can help us to get a tractable objective:

$$\ln p_{\theta}(x) \geq \underbrace{\mathbb{E}_{q_{\phi}(z|x)} \ln p_{\theta}(x|z) - \text{KL}[q_{\phi}(z|x) || p_{\theta}(z)]}_{\text{ELBO}}$$

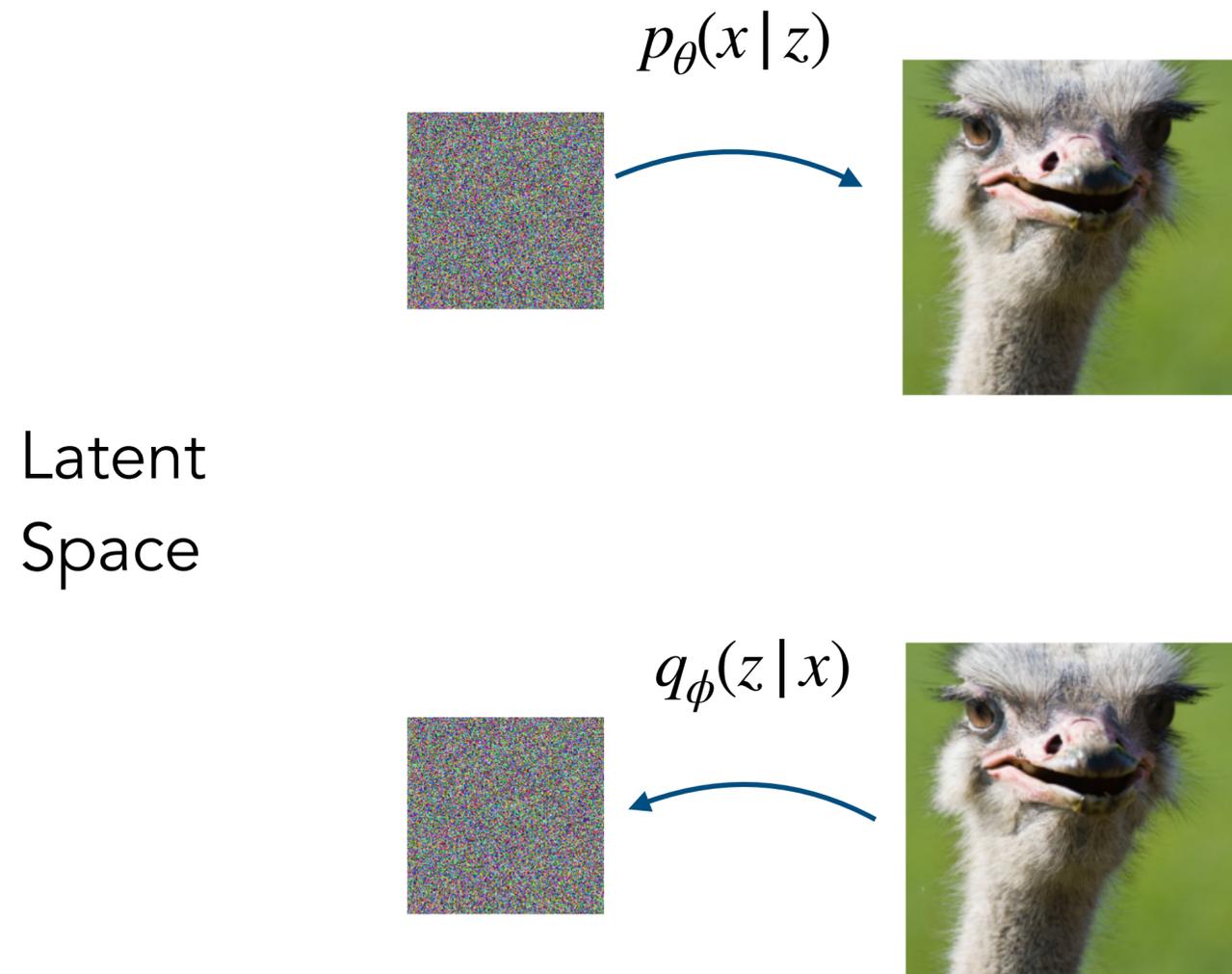
ELBO



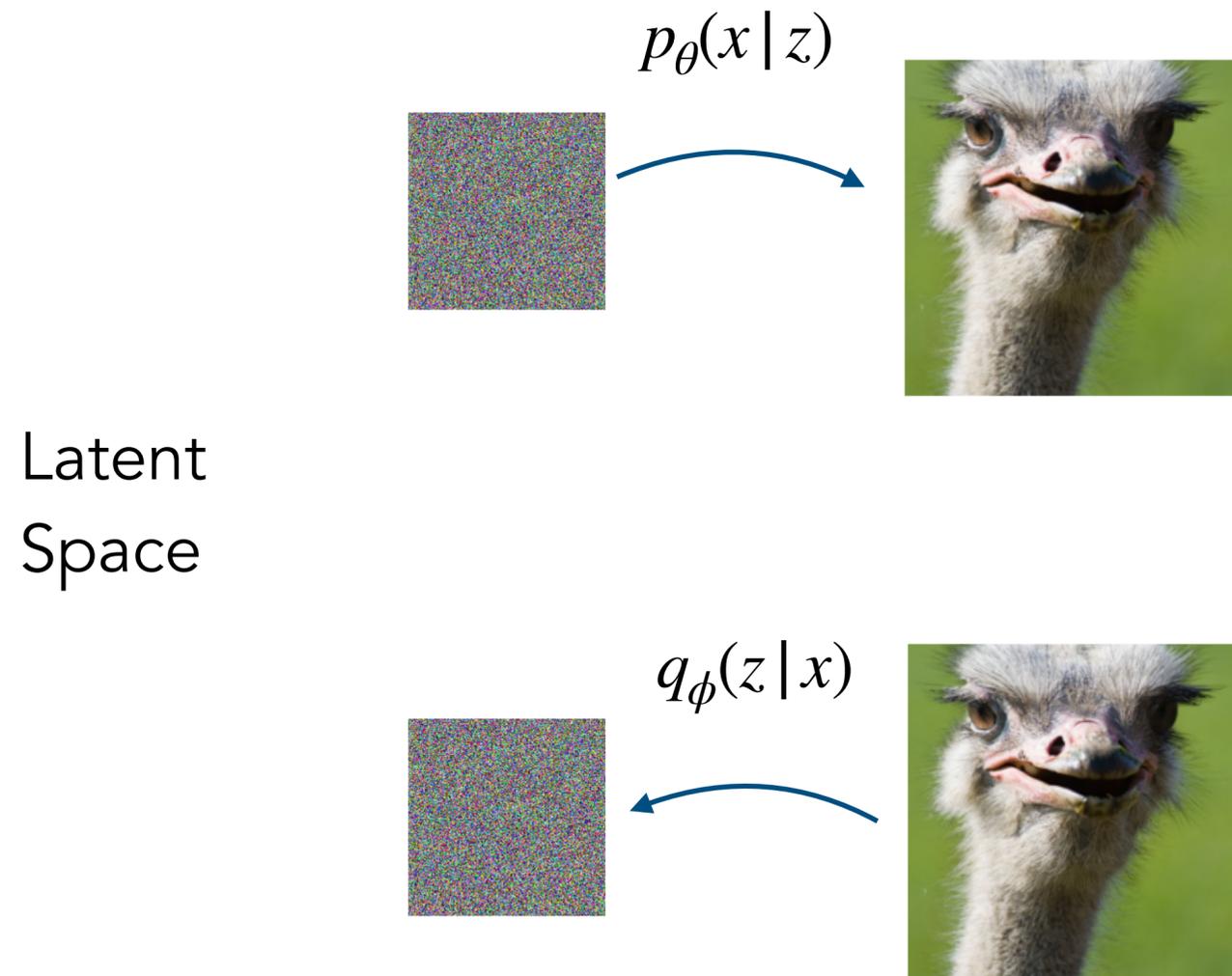
Variational Autoencoder



Are VAEs robust to Adversarial Attacks?



Are VAEs robust to Adversarial Attacks?



No :(

But there is something we can do with that

Adversarial Attack

$$x^a = x^r + \varepsilon, \quad \|\varepsilon\| < \delta$$

x^a "looks" like reference x^r , but is "perceived" differently

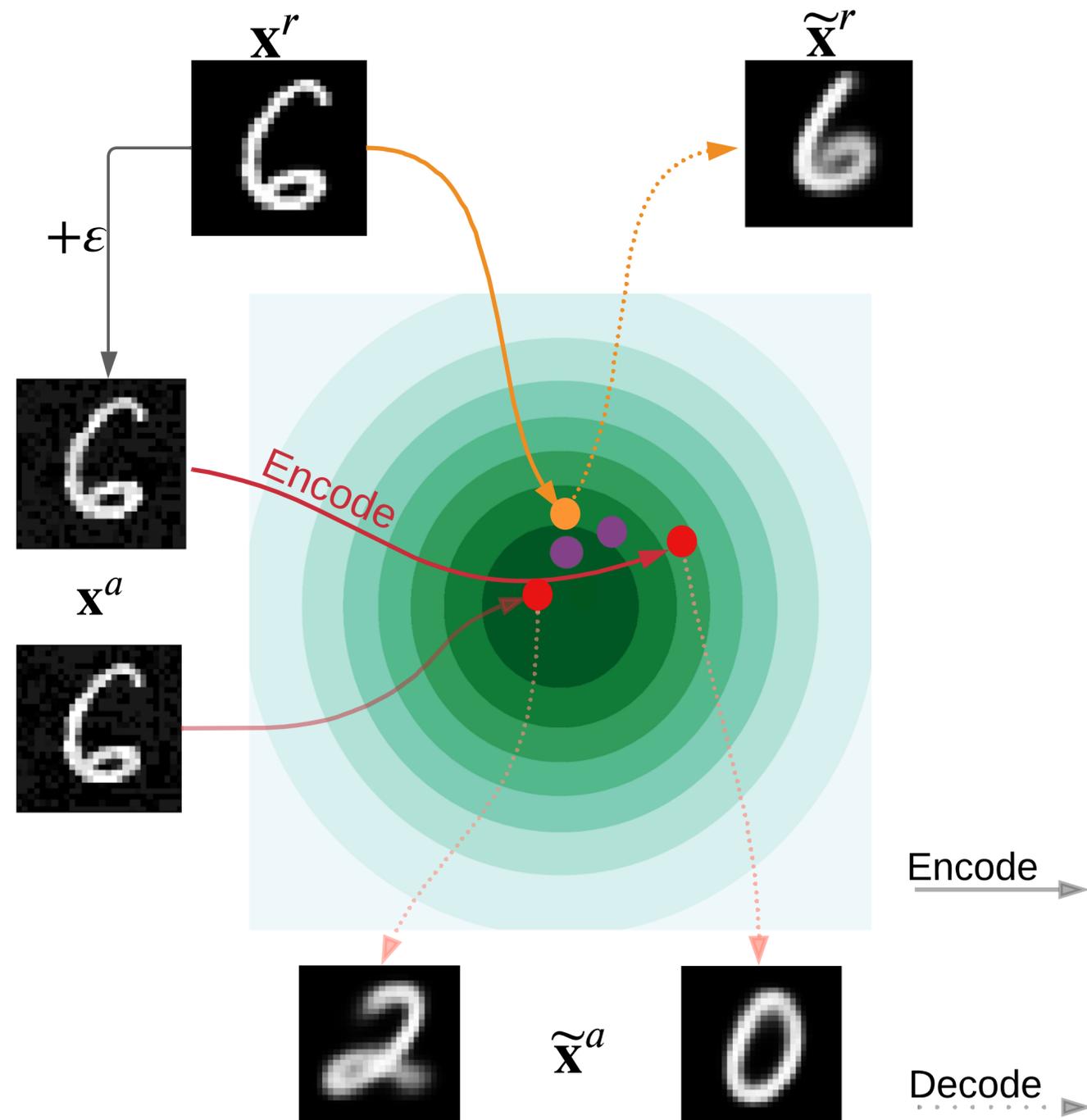
Attacker solves optimisation problem:

$$\varepsilon = \arg \max_{\|\varepsilon\| < \delta} \Delta [f(x^r + \varepsilon), f(x^r)]$$

Distance metric

What is being attacked

Adversarial Attack on VAEs



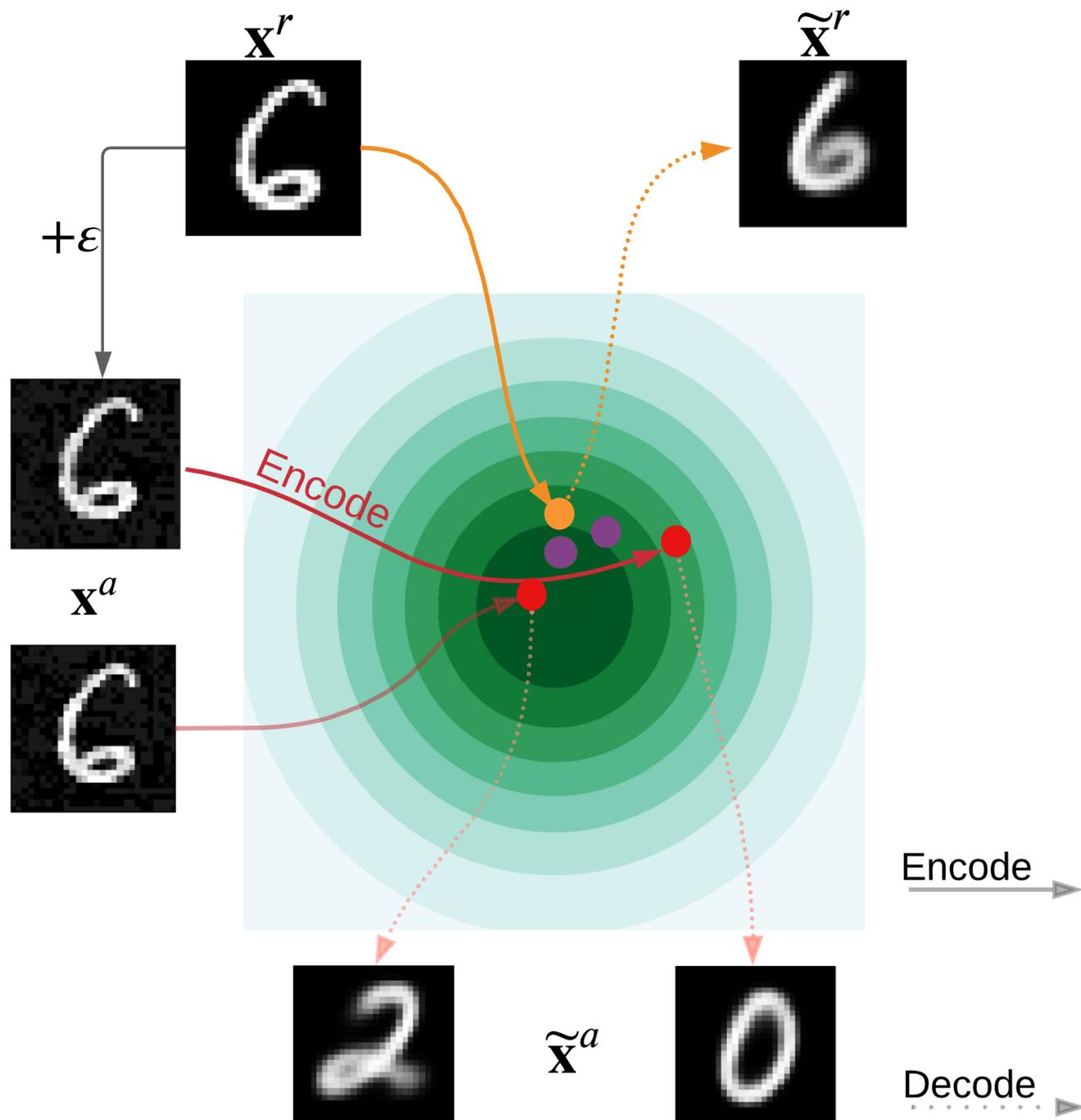
$$x^a = x^r + \epsilon, \quad \|\epsilon\| < \delta$$

x^a "looks" like reference x^r , but is "perceived" differently

Attacker solves optimisation problem:

$$\epsilon = \arg \max_{\|\epsilon\| < \delta} \text{SKL} \left[q_{\phi}(z | x^r + \epsilon), q_{\phi}(z | x^r) \right]$$

Defence Strategy



$$z^a \sim q_\phi(z|x^a) \quad \text{vs} \quad z^r \sim q_\phi(z|x^r)$$

Let's use true posterior instead:

Target density:

$$p_\theta(z|x^a) \propto p(z)p_\theta(x^a|z)$$

Make T steps of MCMC (starting from z^a)

$$z^{(T)} \sim q^{(T)}(z|x^a)$$

Final Algorithm

1. (Defender)

Train a VAE:

$$q_\phi(z|x), p(z), p_\theta(x|z)$$

2. (Attacker)

For a given x^r , construct the attack x^a

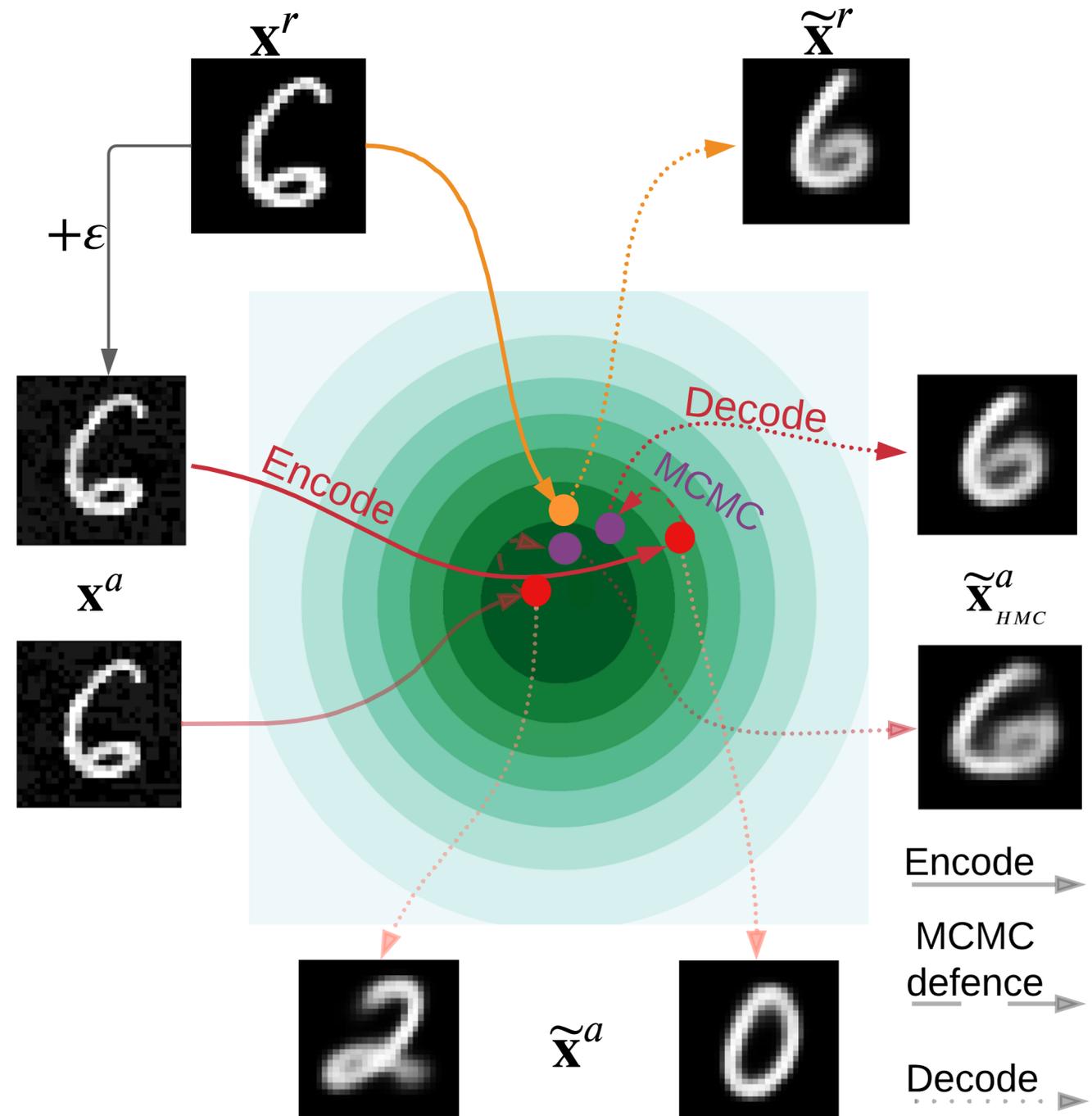
$$x^a = x^r + \varepsilon, \quad \|\varepsilon\| < \delta$$

s.t $q_\phi(z|x^a)$ is "far enough" from $q_\phi(z|x^r)$

3. (Defender)

Run T steps of HMC with the target $\propto p(z)p_\theta(x^a|z)$

Use $z := z^{(T)}$ to decode / in downstream task



Why it works?

Theoretical justification

t steps of MCMC

"desired" distribution

$$\text{TV}[q^{(t)}(z | x^a) \| q_\phi(z | x^r)] \leq \sqrt{\frac{1}{2} \text{KL} [q^{(t)}(z | x^a) \| p_\theta(z | x^a)]} + \sqrt{\frac{1}{2} \text{KL} [q_\phi(z | x^r) \| p_\theta(z | x^r)]} + o(\sqrt{\|\varepsilon\|})$$

How good
is defence

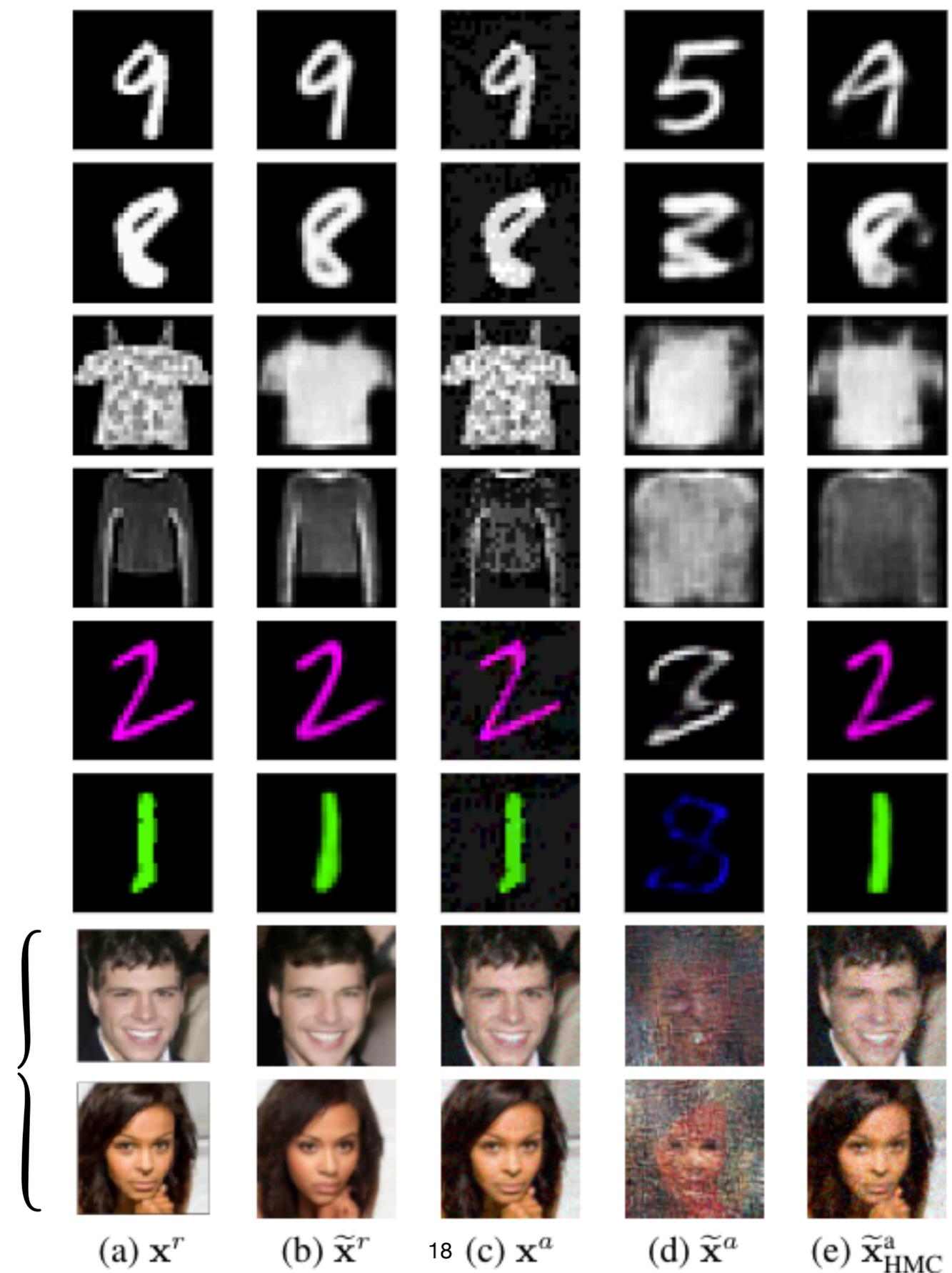
Goes to 0 with $t \rightarrow \infty$

How good VAE is
(approximation gap)

Attack
radius

Empirical Results

NVAE:
deep hierarchical VAE

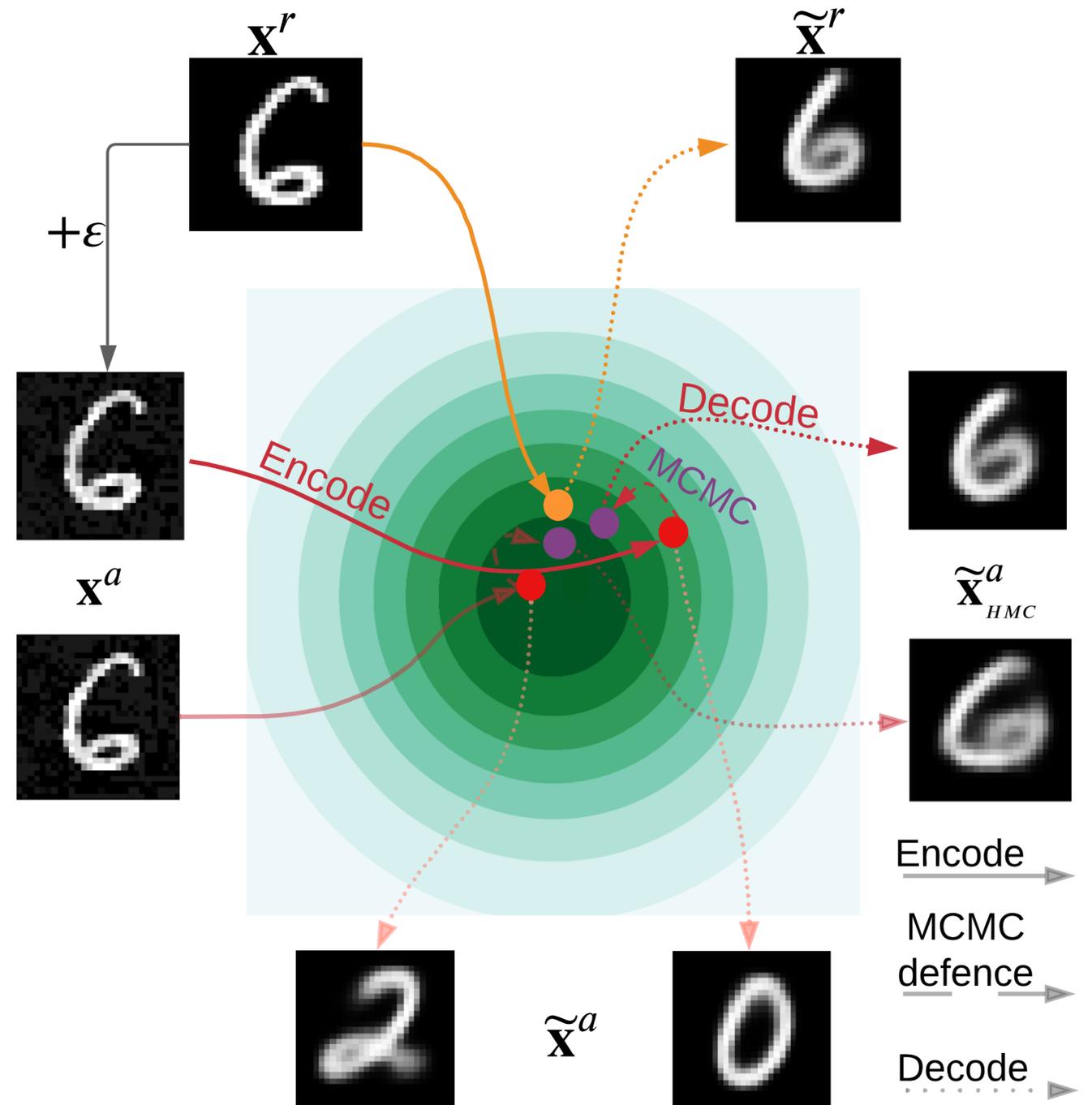


Take Home Message

Latent representations of the data learned by VAE are vulnerable to adversarial attacks

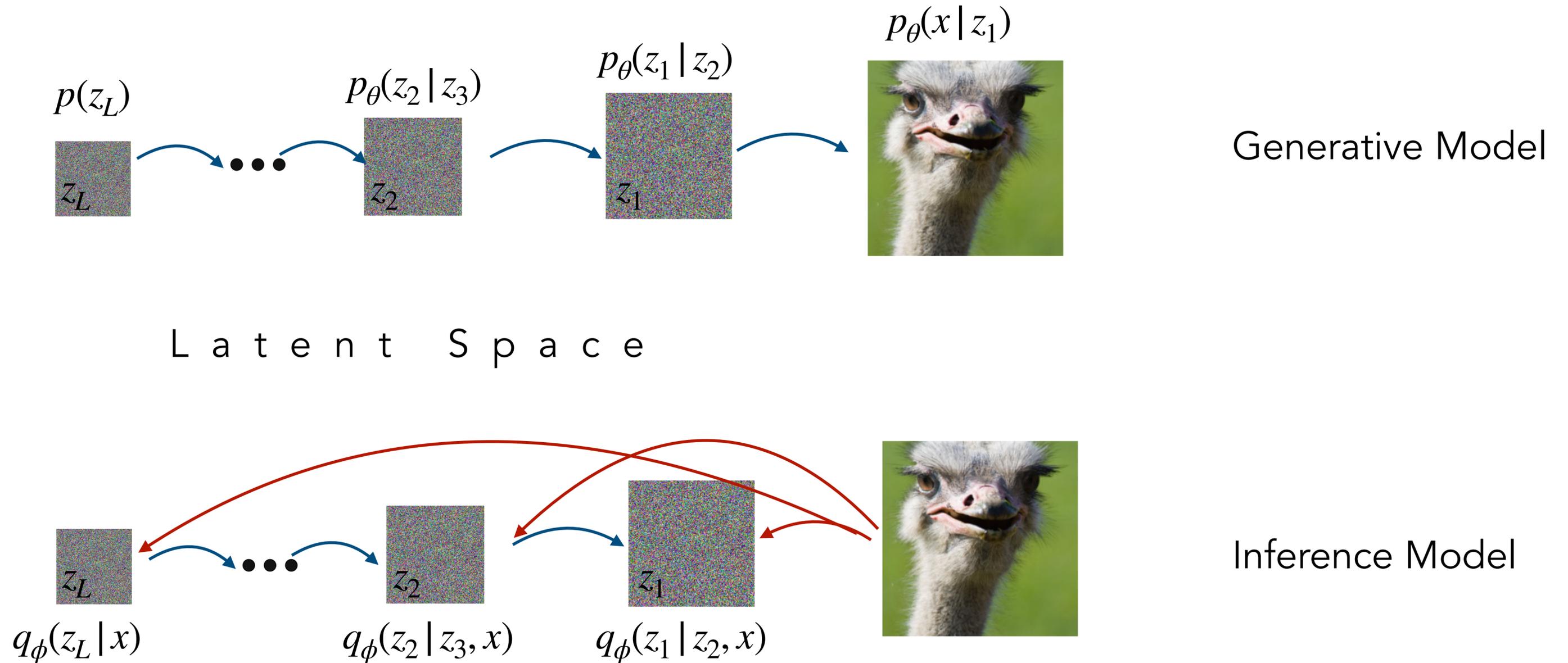
We can use decoder to alleviate the effect of the attack

We "pay" with the inference time for the increased robustness



Hierarchical Variational Autoencoder

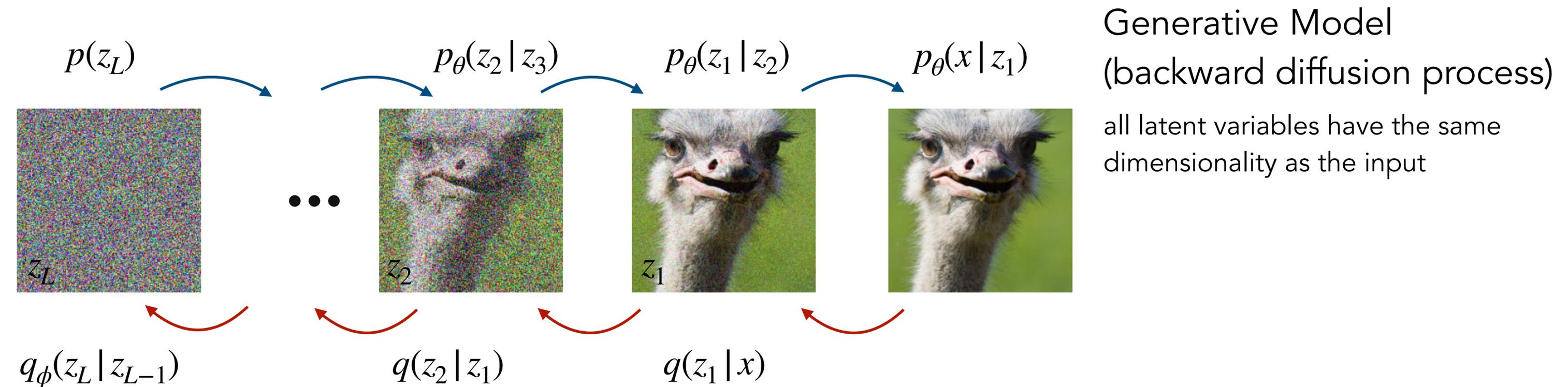
What if $z = (z_1, \dots, z_L)$?



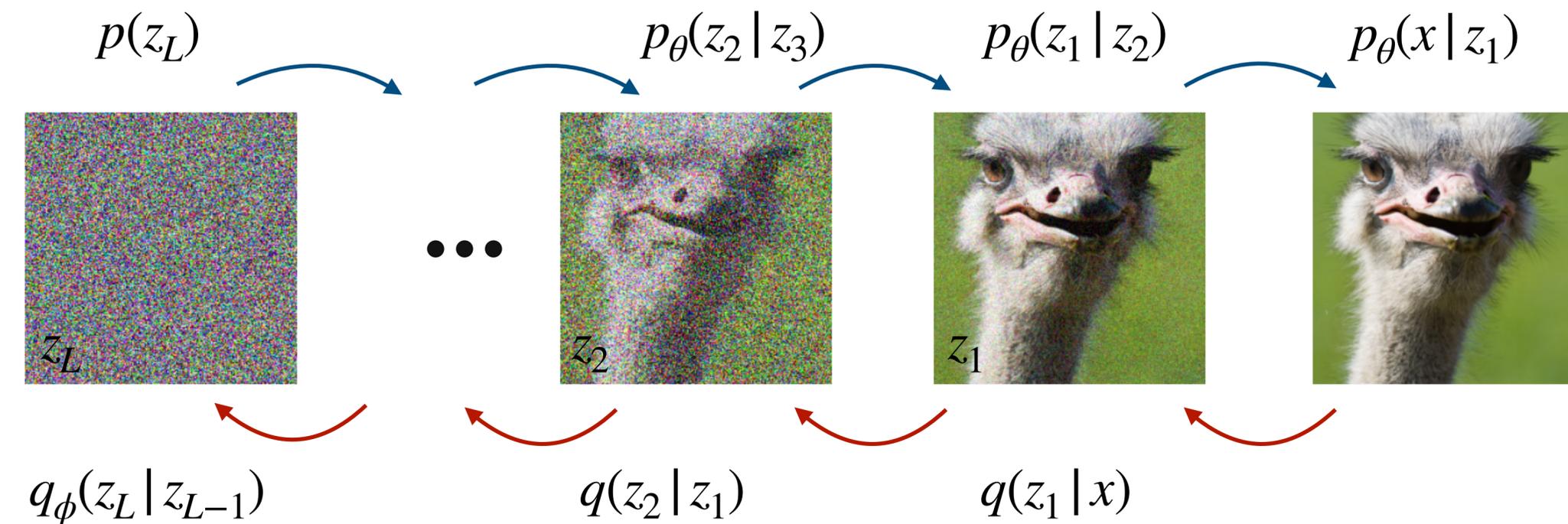
Hierarchical Variational Autoencoder



Diffusion-based Generative Models



Diffusion-based Generative Models



Generative Model
(backward diffusion process)

all latent variables have the same dimensionality as the input

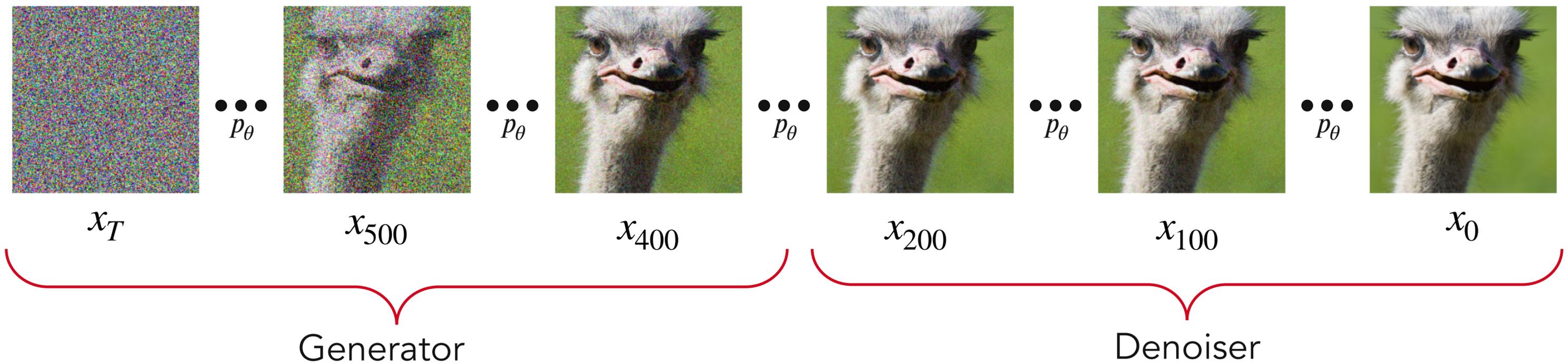
Fixed Inference Model
(or diffusion process)

$$q(z_l | z_{l-1}) = \mathcal{N}(z_l | \sqrt{1 - \beta_l} z_{l-1}, \beta_l I)$$

add gaussian noise according to pre-defined schedule β_1, \dots, β_L

Diffusion-base Generative Models

Can we split the DGM into “generator” and “denoiser”?

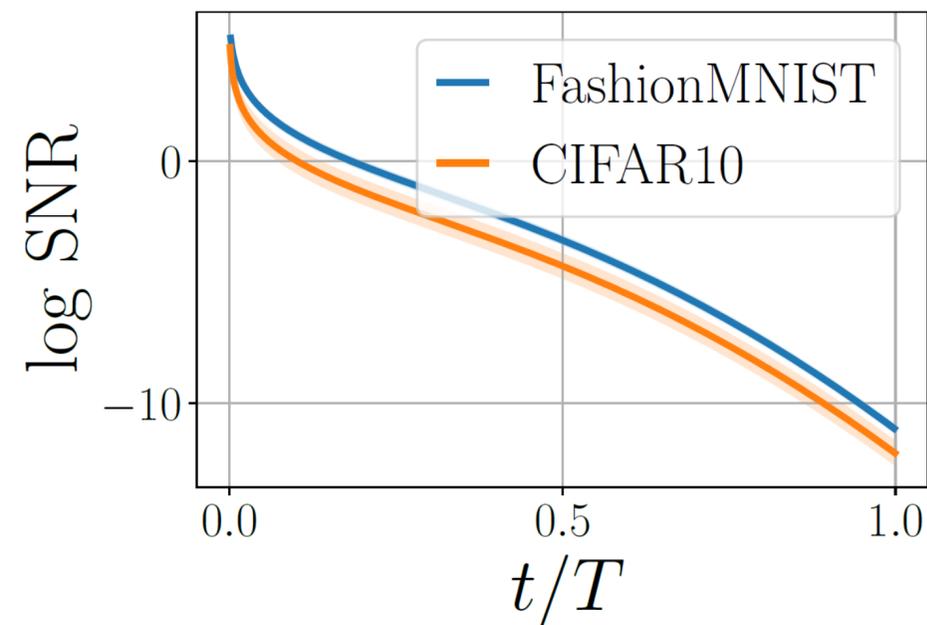


Signal-to-noise ratio

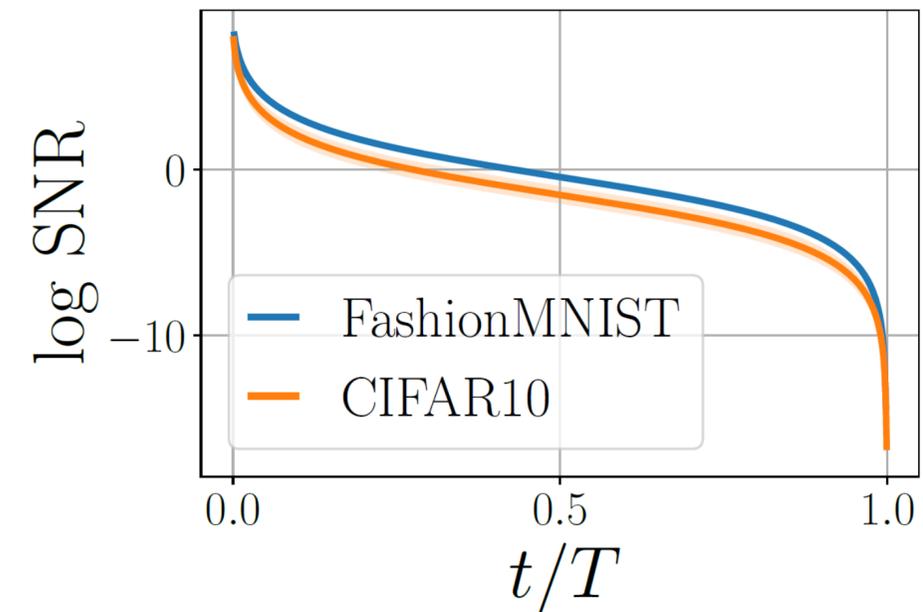
Observation 1: We have strong signal within first 10-20% steps of the forward process

$$q(z_t | z_{t-1}) = \mathcal{N}(z_t | \sqrt{1 - \beta_t} z_{t-1}, \beta_t I)$$

Linear β schedule

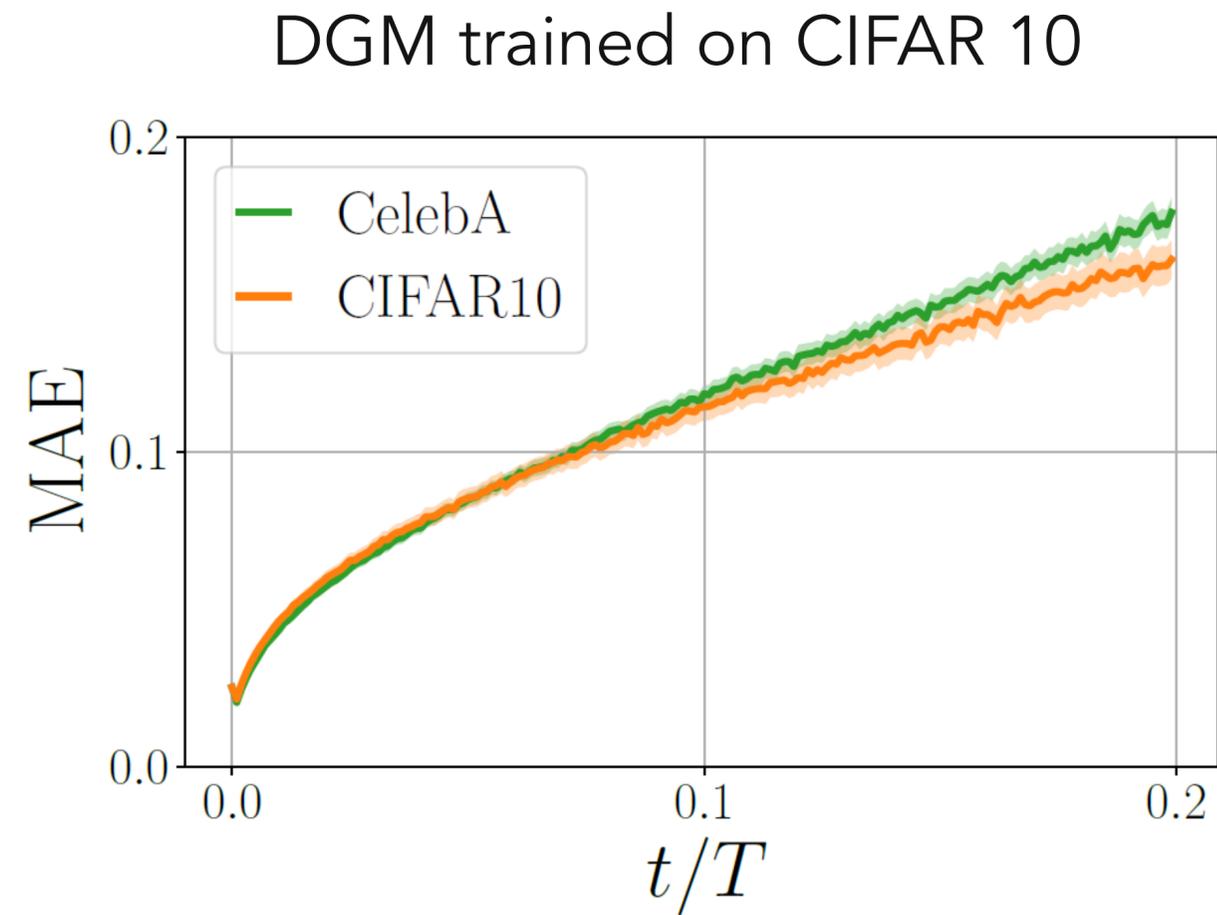


Cosine β schedule



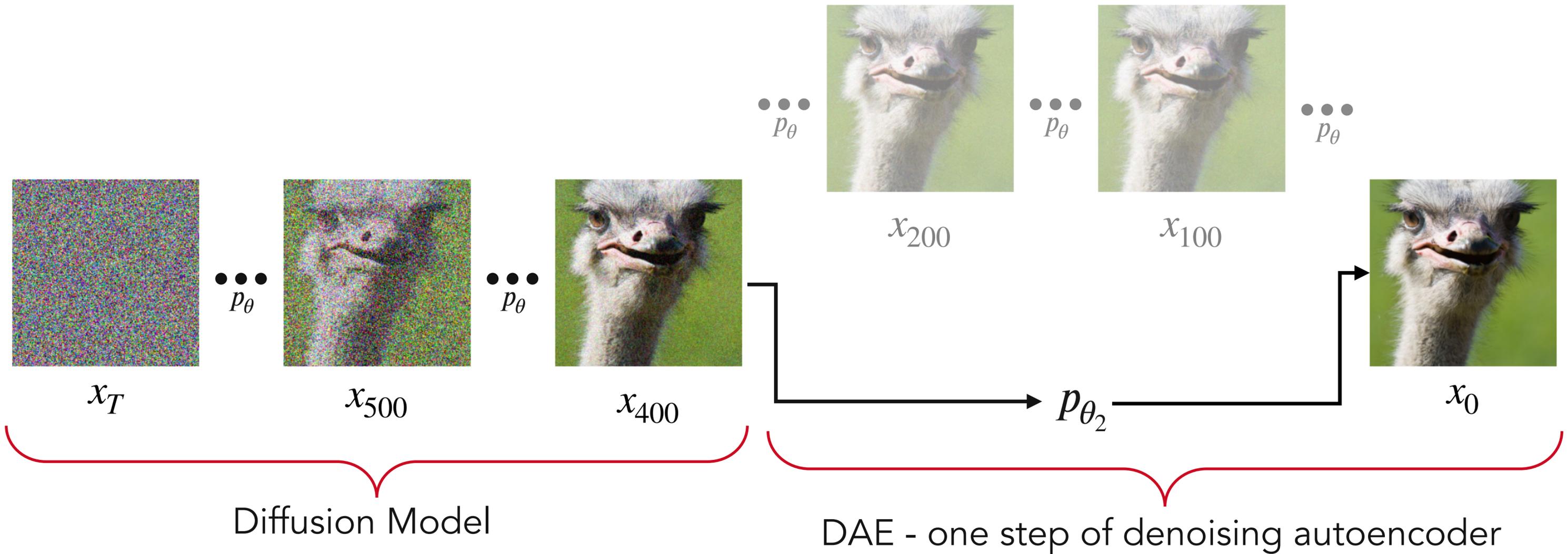
How good is denoising?

Observation 2: Backward process is capable of denoising the out-of-distribution data



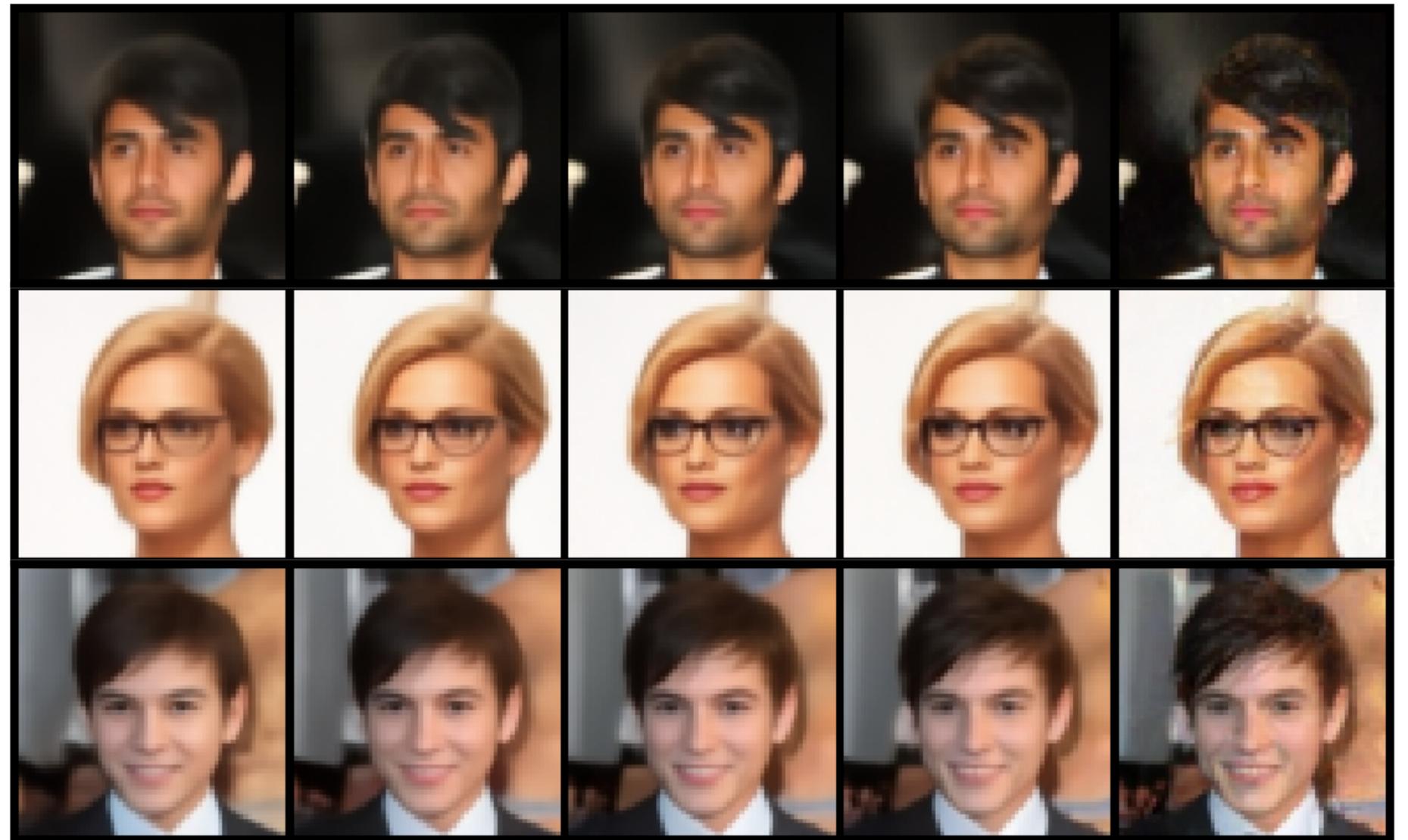
DAED: DGM + Denoising Autoencoder

Let's make denoising step **explicit**



DAED: DGM + Denoising Autoencoder

We replace up to 10% of DGM steps with a single DAE without significant drop in model's performance



$\beta_1 = 0.2$

$\beta_1 = 0.1$

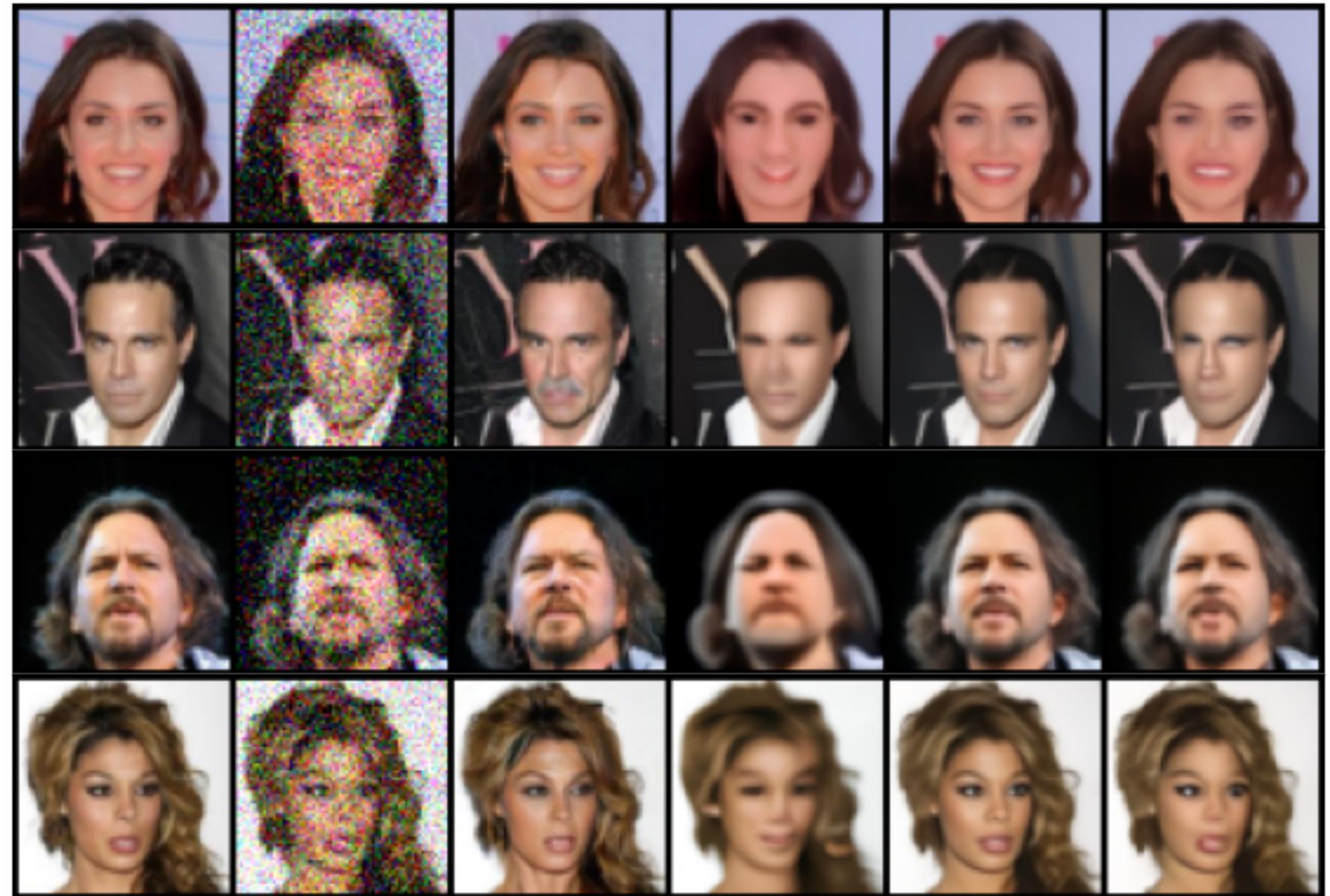
$\beta_1 = 0.05$

$\beta_1 = 0.025$

$\beta_1 = 0.001$

Transferability of noise removal between data distributions

DAED is better in removing noise from unseen data distribution



x_0

x_1
 $\beta_1 = 0.1$

DDGM
CelebA

DDGM
ImageNet

DAED
CelebA

DAED
ImageNet

Alleviating Adversarial Attacks on Variational Autoencoders with MCMC

NeurIPS 2022



the best co-authors:



Max
Welling



Jakub M.
Tomczak

On Analyzing Generative and Denoising Capabilities of Diffusion-based Generative Models

NeurIPS 2022



the best co-authors:



Kamil
Deja



Tomasz
Trzciński



Jakub M.
Tomczak

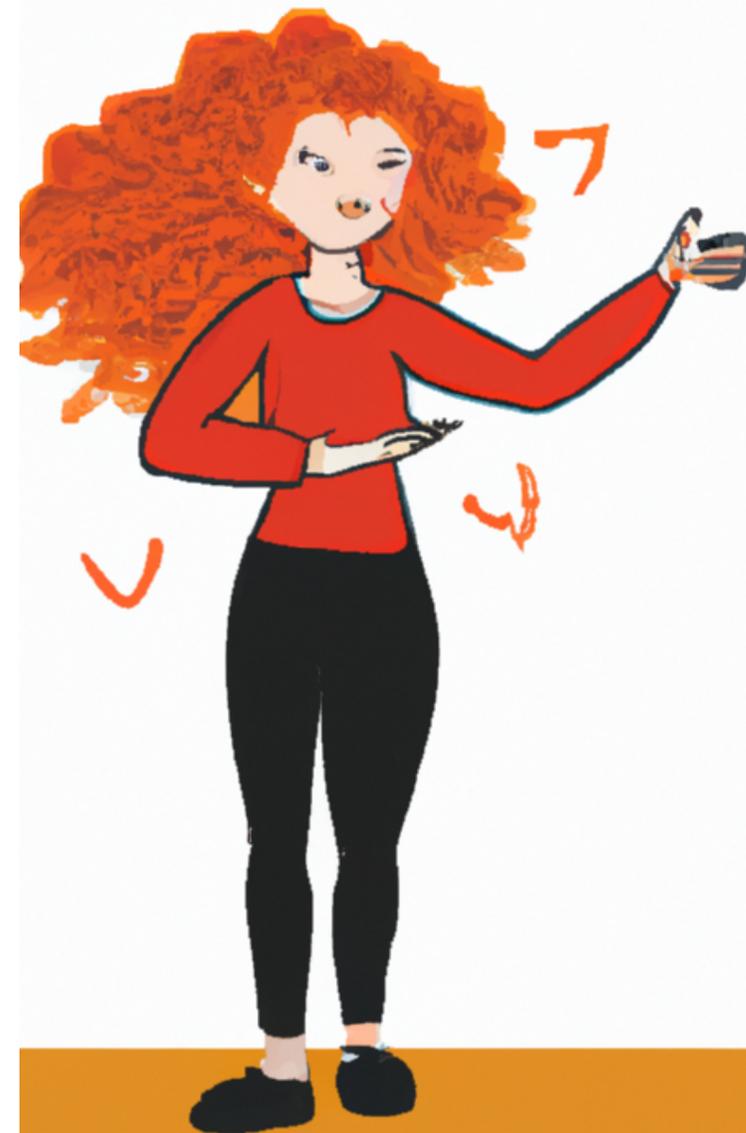
Thank you

Anna Kuzina

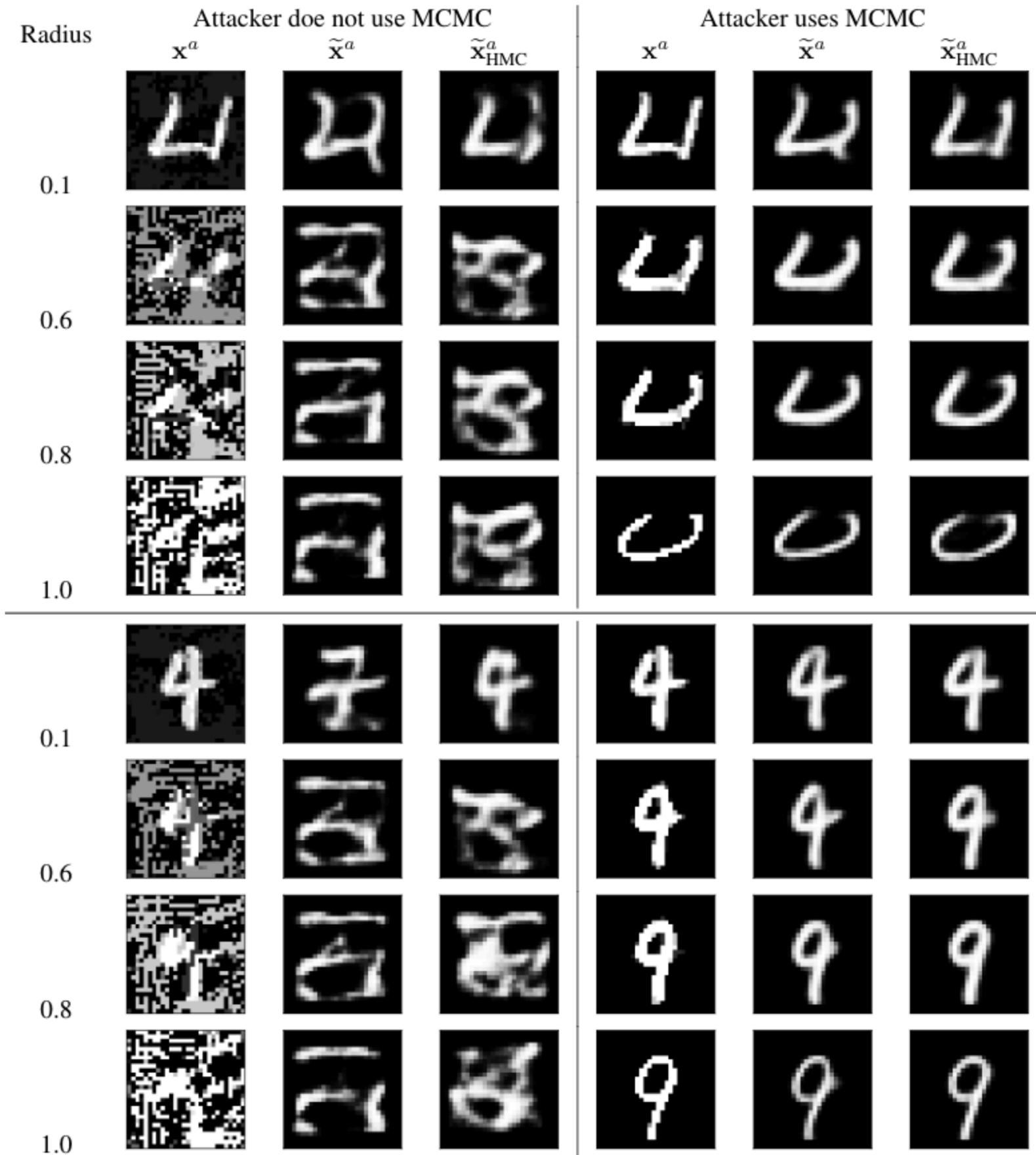
PhD student
Vrije Universiteit Amsterdam

 a.kuzina@vu.nl

 akuzina.github.io/



What if attacker knows the defence strategy?



Why it works?

Empirical Evidence

Given a reference point, one can evaluate posterior ratio for two latent codes:

$$\text{PR}(z_1, z_2) = \frac{p_{\theta}(z_1 | x^r)}{p_{\theta}(z_2 | x^r)}$$

Blue: reference latent code VS adversarial latent code

Orange: reference latent code VS adversarial latent code after HMC

