# Alleviating Adversarial Attacks on Variational Autoencoders with MCMC

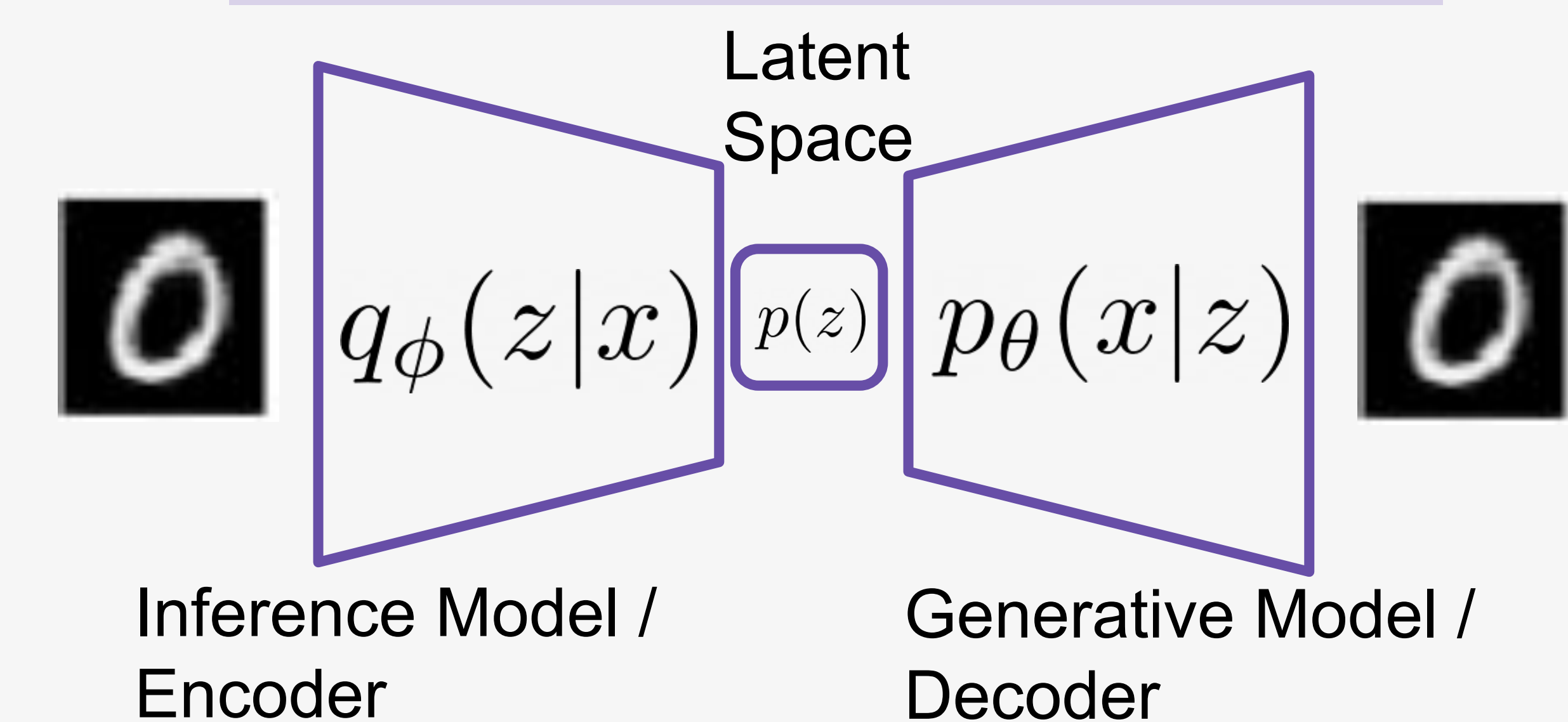Anna Kuzina,    Max Welling,    Jakub M. Tomczak

## Summary

VAE are not robust to adversarial attacks

We propose the way to alleviate the effect

Method does not require changing the training procedure

We have theoretical justification why it works

## Variational Auto-Encoder



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

Latent Space

Inference Model / Encoder

Generative Model / Decoder

## Hierarchical VAE

L latent variables $\mathbf{z} = (z_1, \ldots, z_L)$

## Adversarial Input

$$x^a = x^r + \varepsilon$$



$\tilde{x}^r$

$\tilde{x}^a$

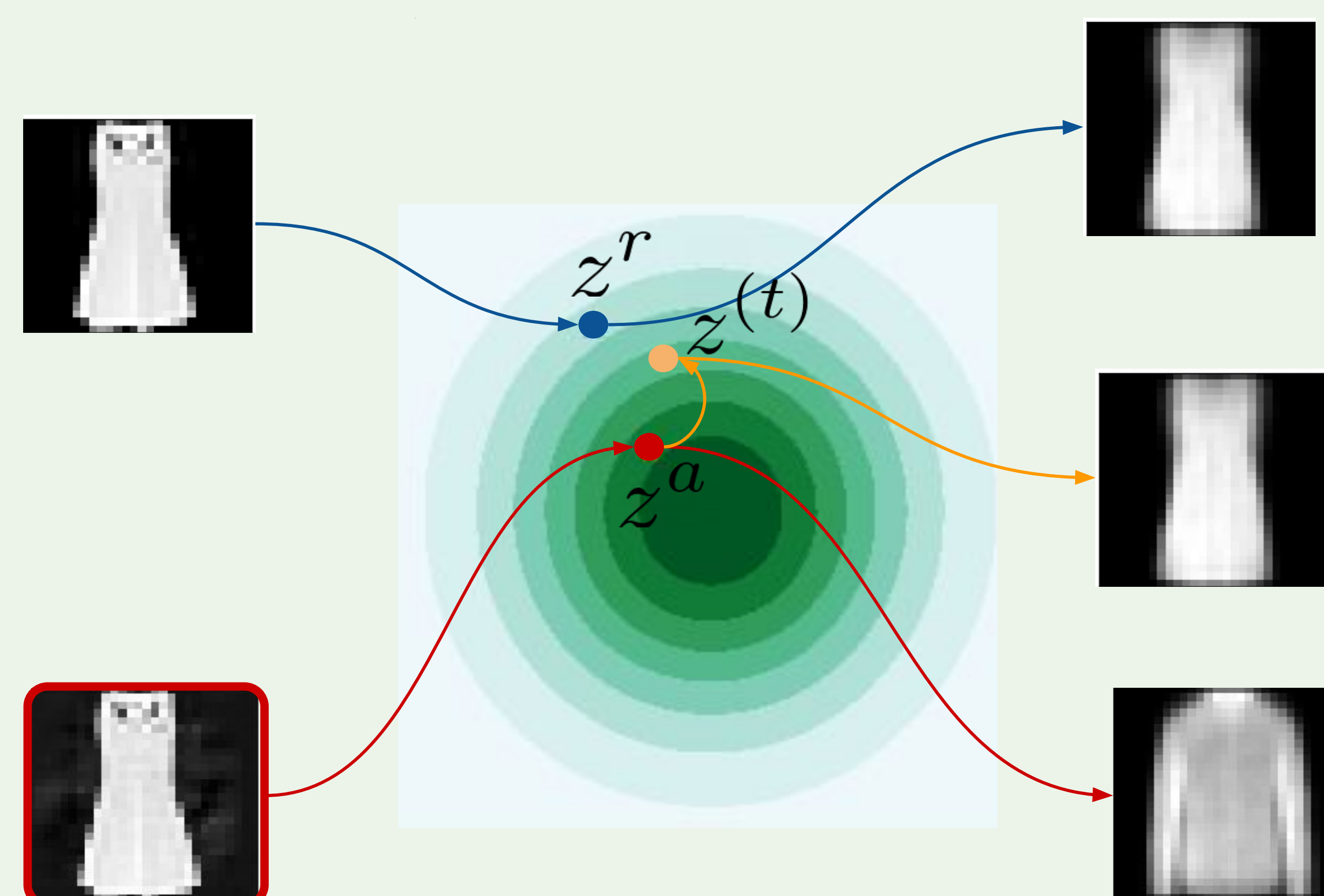## Can we reduce the effect of attack?

$z^r \sim q_\phi(z|x^r)$    is what we want

$z^a \sim q_\phi(z|x^a)$    is what we get instead

Let's sample from the true posterior $p_\theta(z|x^a) \propto p(z)p_\theta(x^a|z)$

$$z^{(t)} \sim q^{(t)}(z|x^a) = \int q_\phi(z_0|x^a)Q^{(t)}(z|z_0)dz_0$$



$z^r$    $z^{(t)}$    $z^a$

## Why does it work?

Gets smaller with each MCMC step

$$\mathrm{TV}[q^{(t)}(z|x^a)\|q_\phi(z|x^r)] \leq \sqrt{\tfrac{1}{2}\mathrm{KL}\left[q^{(t)}(z|x^a)\|p_\theta(z|x^a)\right]}$$
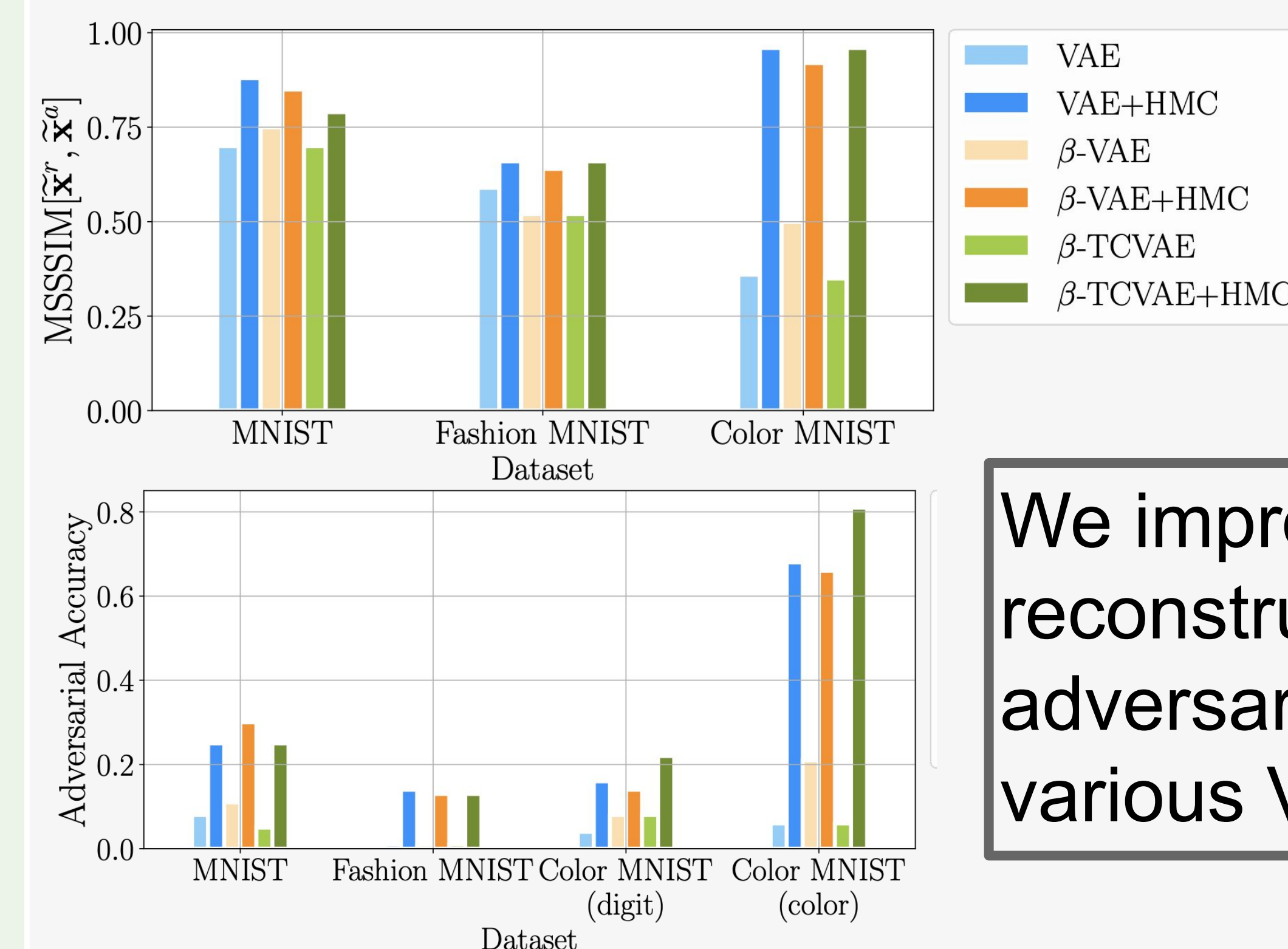
VAE amortization gap

$$+ \sqrt{\tfrac{1}{2}\mathrm{KL}\left[q_\phi(z|x^r)\|p_\theta(z|x^r)\right]}$$
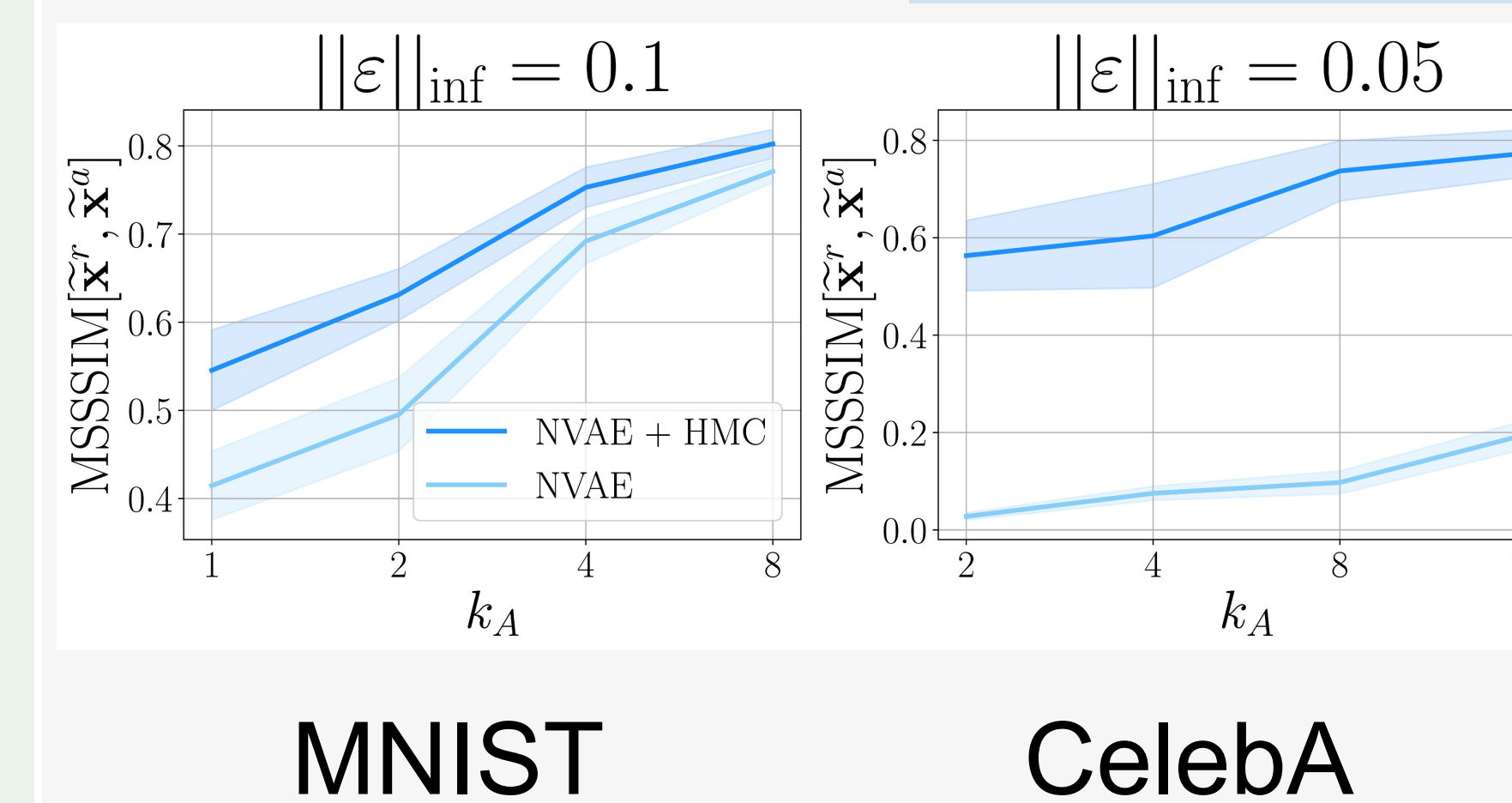
Attack radius

$$+ o(\sqrt{\|\varepsilon\|})$$

## Results: VAE



We improve both reconstruction quality and adversarial accuracy for various VAE modificationd

## Results: NVAE



$\|\varepsilon\|_{\mathrm{inf}} = 0.1$    $\|\varepsilon\|_{\mathrm{inf}} = 0.05$

MNIST    CelebA

Consider only last $k_A$ latent variables for attack construction

We observe that reconstructions and more similar to the reference when we use the proposed method



(a) $\mathbf{x}^r$    (b) $\tilde{\mathbf{x}}^r$    (c) $\mathbf{x}^a$    (d) $\tilde{\mathbf{x}}^a$    (e) $\tilde{\mathbf{x}}^a_{\mathrm{HMC}}$