

# Adversarial Attacks on Variational Autoencoders

**George Gondim-Ribeiro, *Pedro Tabacof* & Eduardo Valle**



# WHO AM I?

- ❖ PhD student at the University of Campinas in Brazil
- ❖ Data scientist at Nubank, credit card fintech
- ❖ Co-authored a few other papers on adversarial attacks, mostly during the Masters:
  - Exploring the space of adversarial images, 2016 IJCNN, with Eduardo Valle (55 citations)
  - Adversarial images for variational autoencoders, 2016 NIPS Adversarial Learning workshop, with Julia Tavares and Eduardo Valle (11 citations)
- ❖ Also interested in Bayesian deep learning and uncertainty in machine learning (current research topic)



# ADVERSARIAL IMAGES



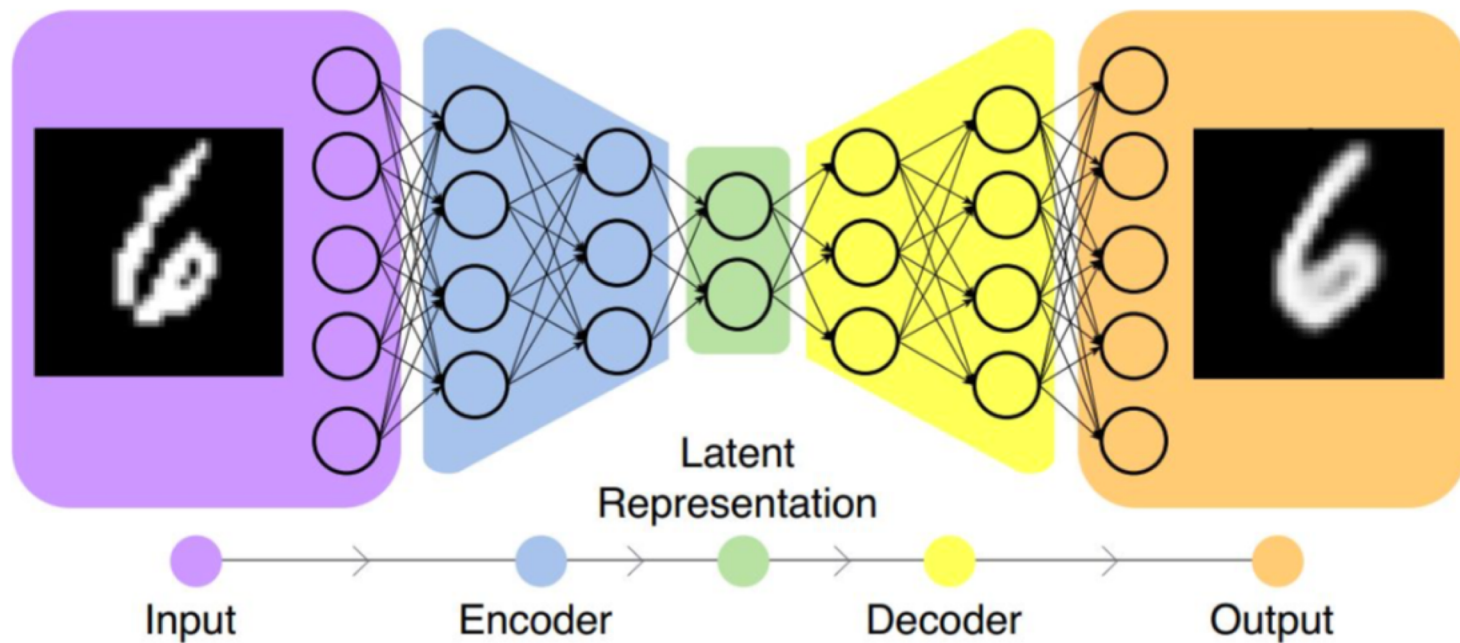
# ADVERSARIAL IMAGES

$$\begin{aligned} & \underset{\mathbf{d}}{\text{minimize}} && \|\mathbf{d}\| \\ & \text{subject to} && L \leq \mathbf{x} + \mathbf{d} \leq U \\ & && \mathbf{p} = f(\mathbf{x} + \mathbf{d}) \\ & && \max(p_1 - p_c, \dots, p_n - p_c) > 0 \end{aligned}$$

*$\mathbf{x}$  is the original image,  $\mathbf{d}$  is the distortion,  $\mathbf{x} + \mathbf{d}$  is the adversarial input,  $f$  is the classifier,  $p_i$  are the scores for each class (where  $c$  is the correct class), and  $L$  and  $U$  are the bounds for the input space.*



# VARIATIONAL AUTOENCODERS



# VARIATIONAL AUTOENCODERS

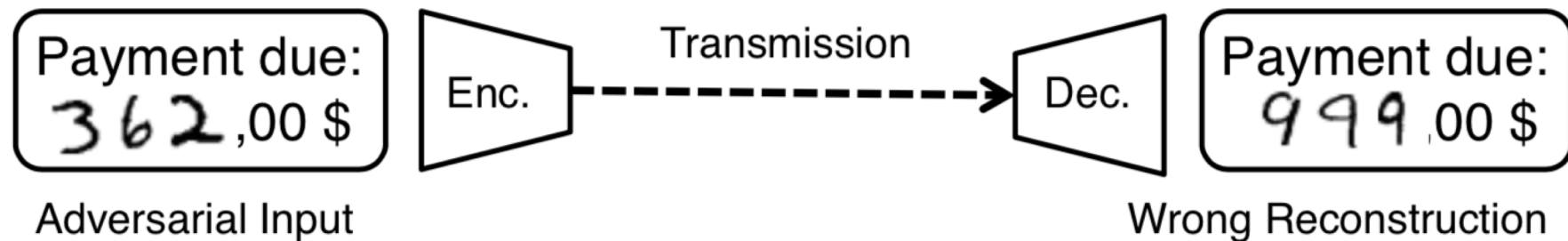
Maximize the Evidence Lower Bound (ELBO)

$$\log p_{\theta}(\mathbf{x}) \geq \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\psi}(\mathbf{x}|\mathbf{z})] - D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}))$$

Reconstruction

Regularization

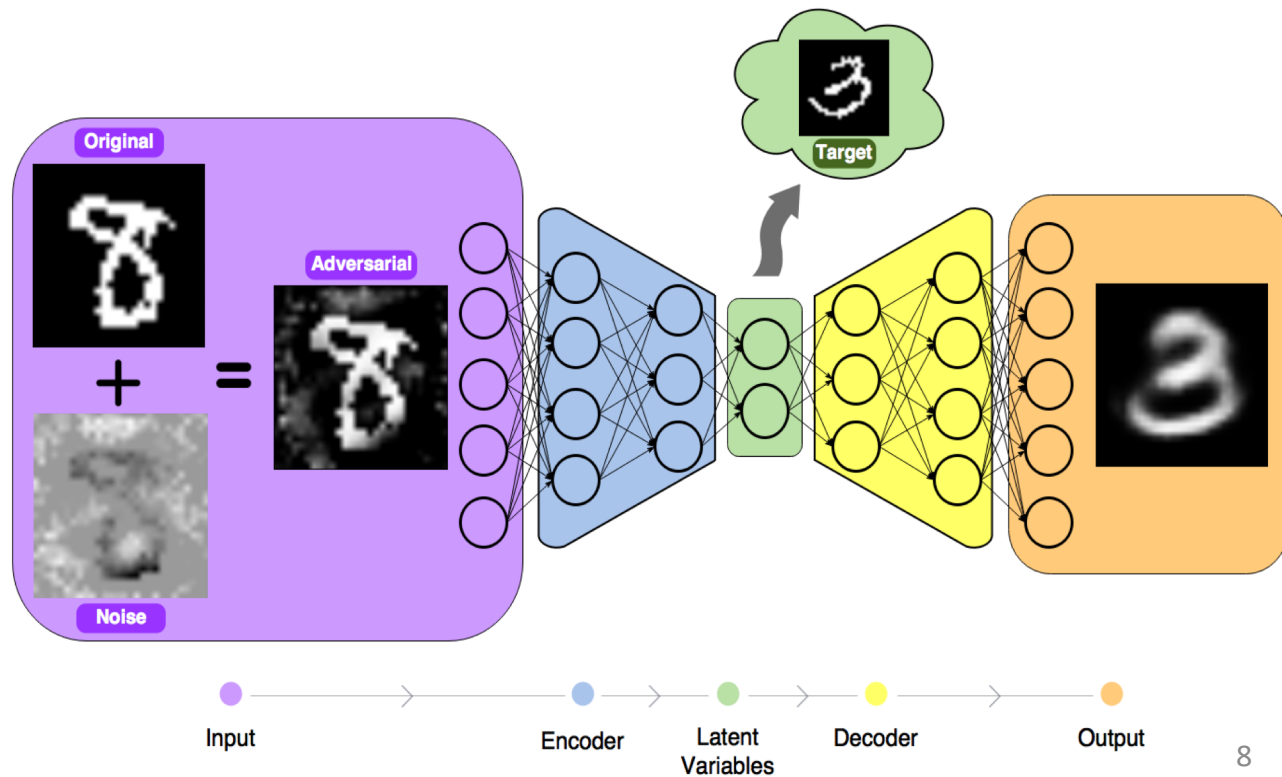
# MOTIVATION



# MAIN IDEA

We attack variational autoencoders with adversarial images. We aim not only to disturb the reconstruction, but also to fool the autoencoder into reconstructing a completely different target image.

We attack the latent representation, attempting to match it to the target image's, while keeping the input distortion as small as possible.



# THE ATTACK

We attack the *latent layer* — which is the information bottleneck of the autoencoder— with the optimization at the right.

The  $\Delta$  function we used was the KL divergence.

$$\begin{aligned} \min_{\mathbf{d}} \quad & \Delta(\mathbf{z}_a, \mathbf{z}_t) + C \|\mathbf{d}\| \\ \text{s.t.} \quad & L \leq \mathbf{x} + \mathbf{d} \leq U \\ & \mathbf{z}_a = \text{encoder}(\mathbf{x} + \mathbf{d}) \end{aligned}$$



# THE ATTACK

We also attack the *output reconstruction* — with the optimization at the right.

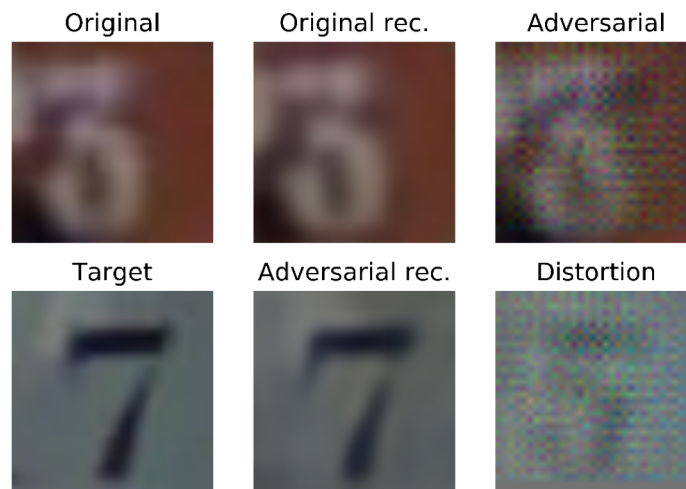
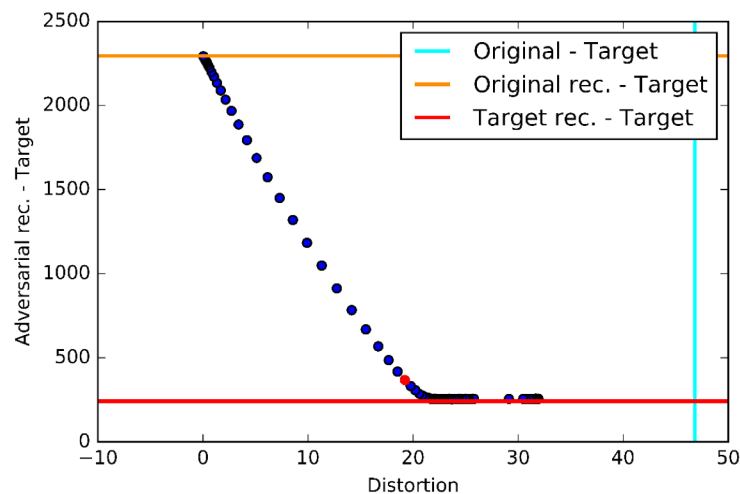
The  $\Delta$  function is the  $\ell_2$ -norm.

$$\begin{aligned} \min_{\mathbf{d}} \quad & \Delta(\mathbf{r}_a, \mathbf{I}_t) + C\|\mathbf{d}\| \\ \text{s.t.} \quad & L \leq \mathbf{x} + \mathbf{d} \leq U, \\ & \mathbf{z}_a = \text{encoder}(\mathbf{x} + \mathbf{d}), \\ & \mathbf{r}_a = \text{decoder}(\mathbf{z}_a) \end{aligned}$$

# THE ATTACK

Decreasing the regularizer  $C$  allows for bigger distortions...

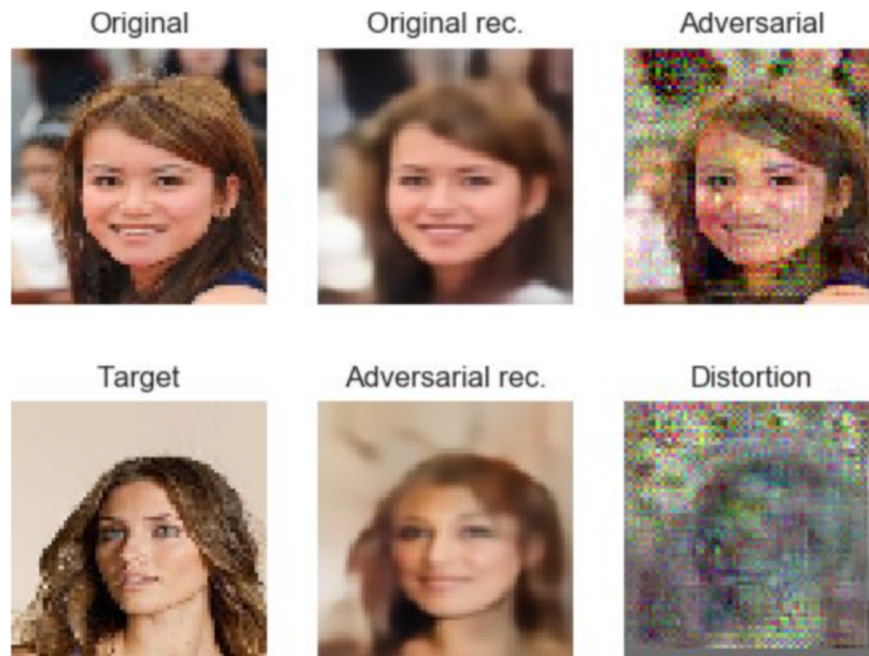
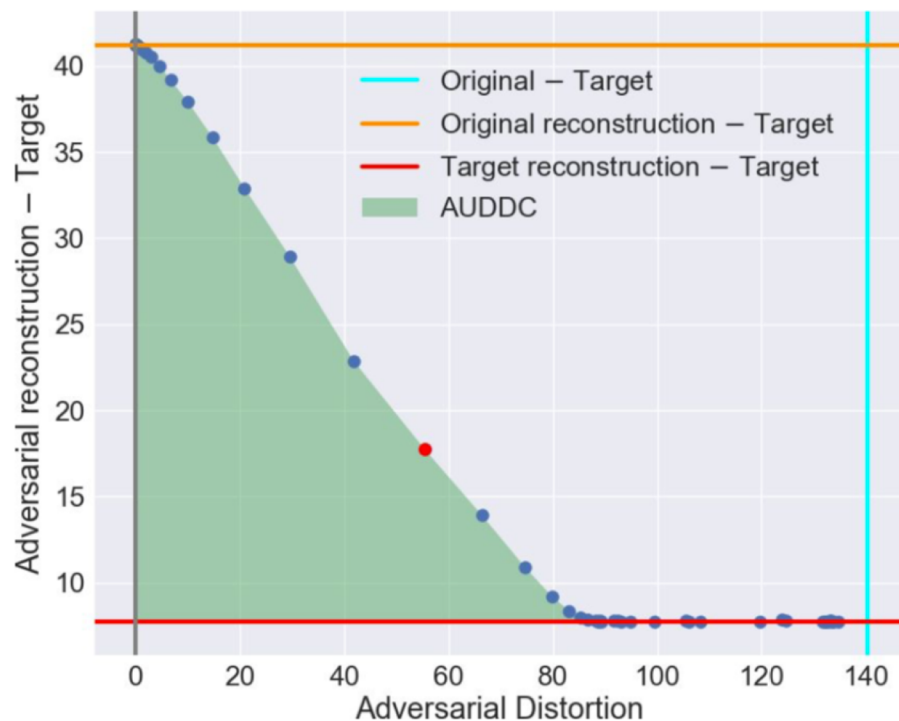
...bringing the adversarial reconstruction closer to the target



# THE METRIC

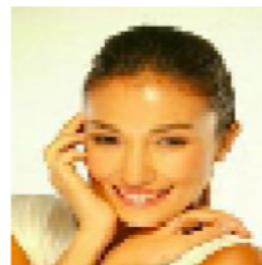
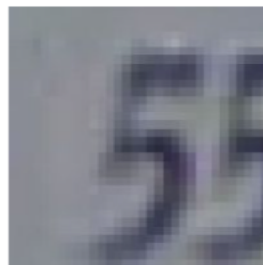
AUDDC: Area Under the Distortion-Distortion Curve

*From 0 (easiest attack possible) to 100 (hardest attack possible)*



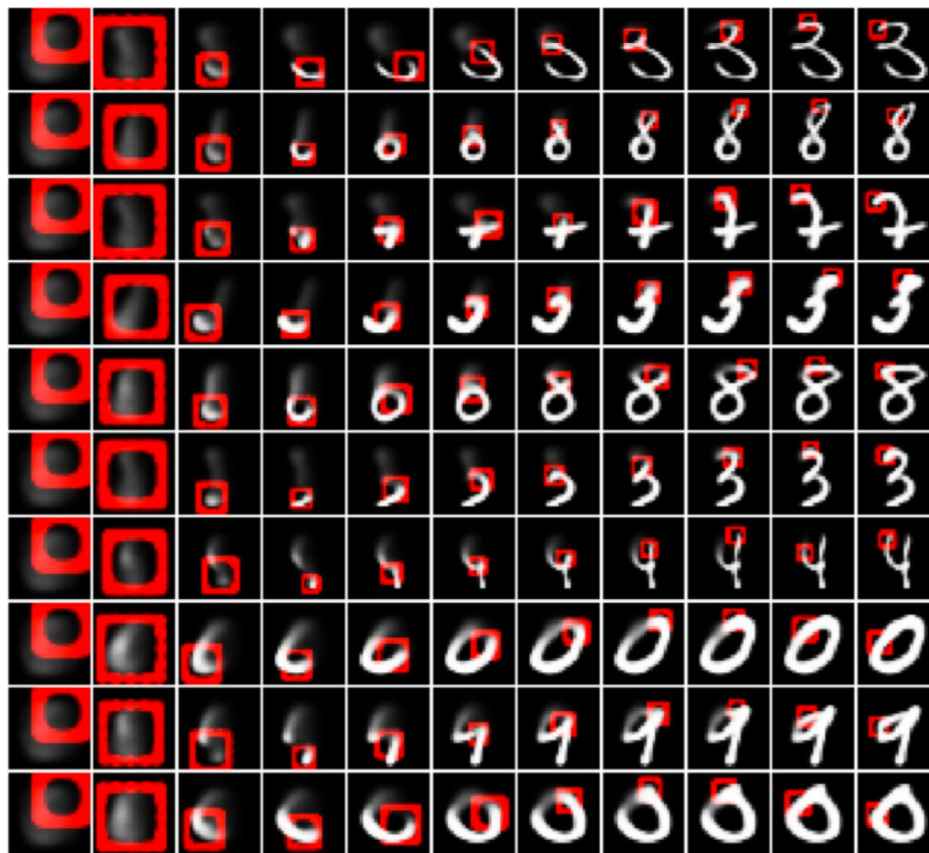
# METHODOLOGY

- Three datasets: MNIST, SVHN, and CelebA



- Models: fully-connected VAEs, convolutional VAEs (CVAE), and DRAW
- A point in the Distortion-Distortion Curve is the average of 128 attacks

# DRAW

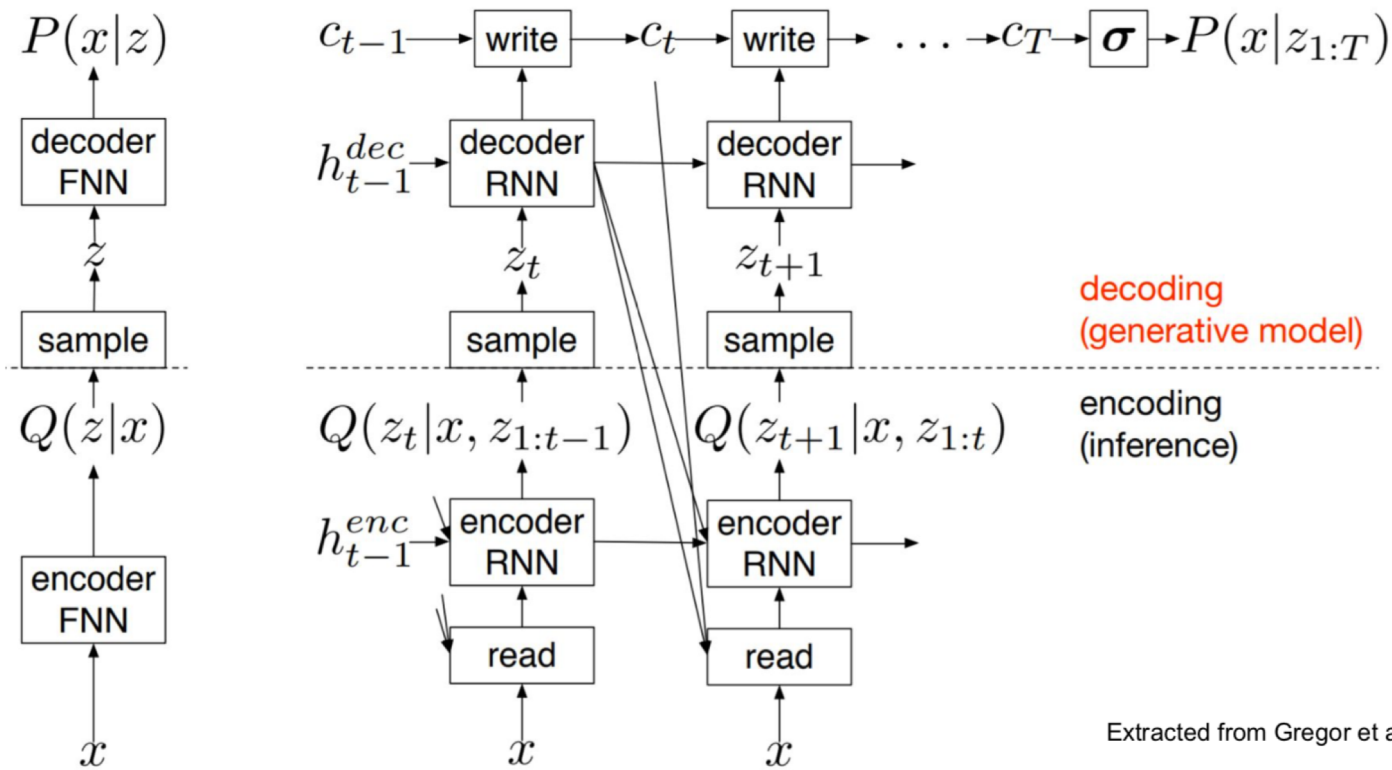


Time →

Extracted from Gregor et al., 2015

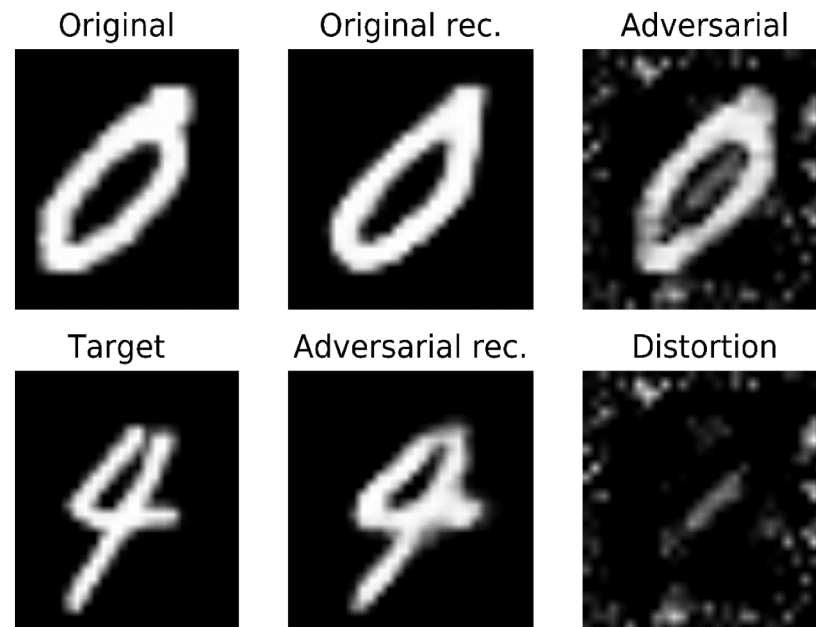
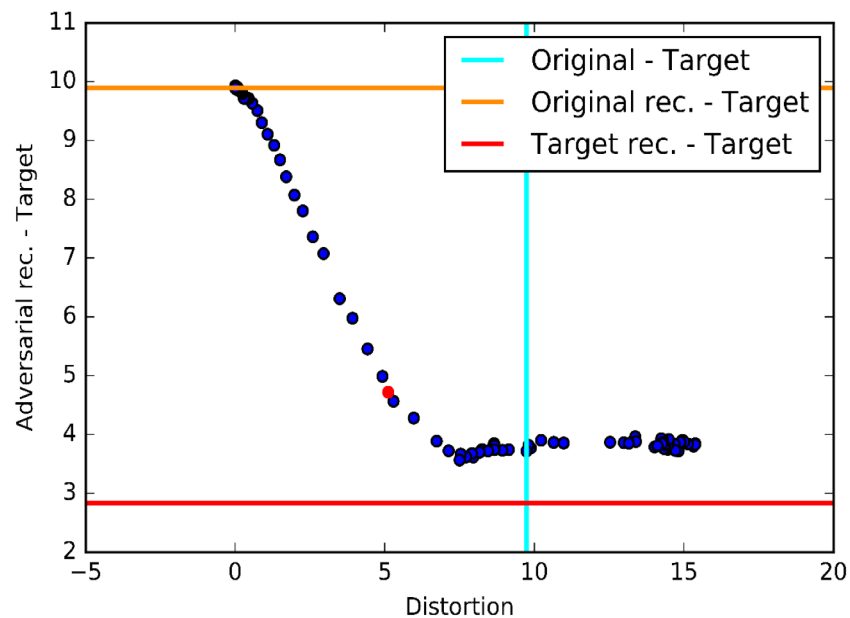


**DRAW**

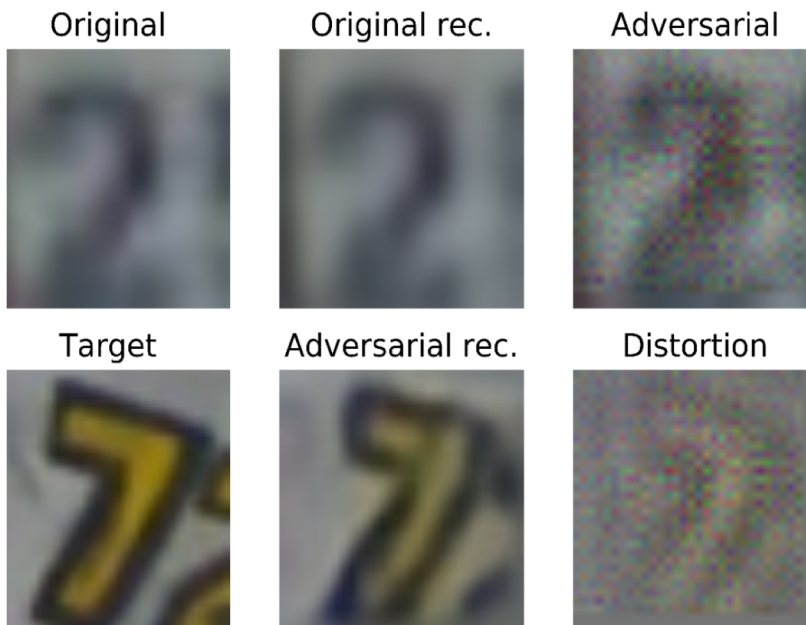
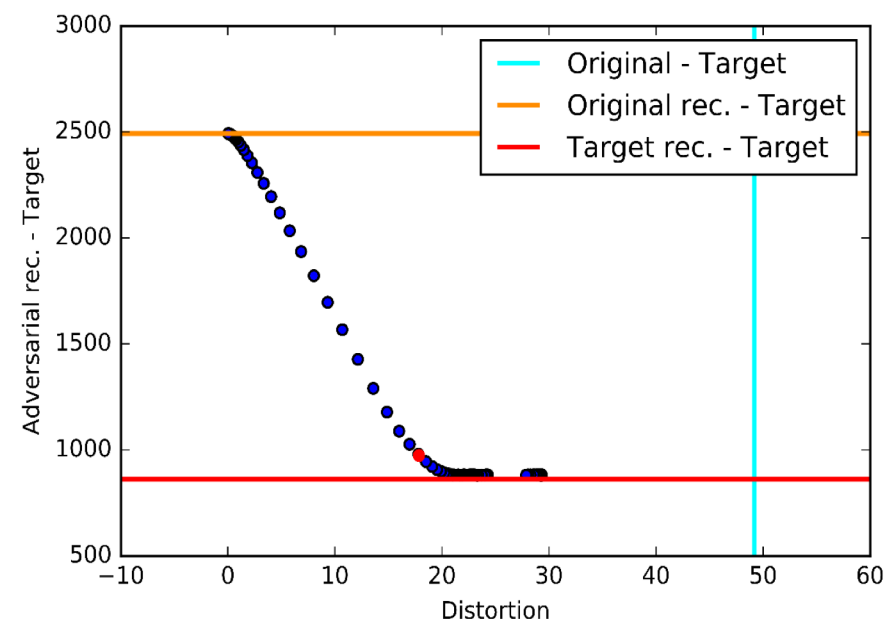


Extracted from Gregor et al., 2015

# MAIN FINDINGS



# MAIN FINDINGS

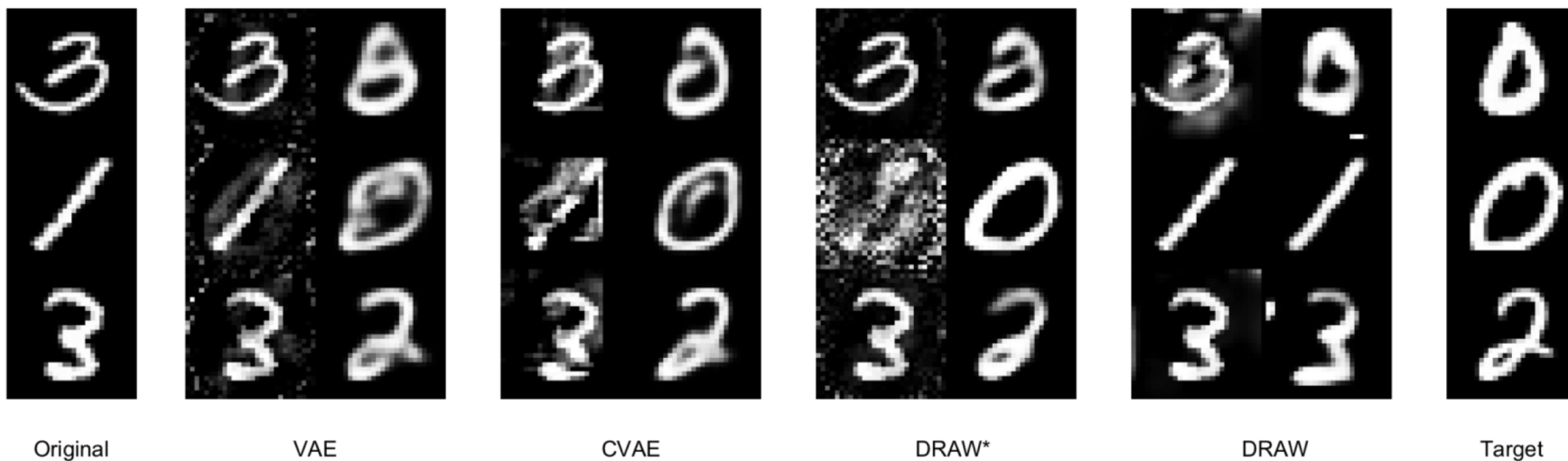


# MAIN FINDINGS

Steps	VAE —	CVAE —	DRAW* 1	DRAW 1	DRAW* 16	DRAW 16		
Attacks on latent representation								
MNIST	27 ± 2	35 ± 3	27 ± 1	35 ± 3	71 ± 5	<b>91 ± 3</b>	47 ± 3	
SVHN	19 ± 1	18 ± 1	09 ± 1	27 ± 2	74 ± 6	<b>96 ± 2</b>	41 ± 4	
CelebA	31 ± 1	28 ± 1	21 ± 2	36 ± 1	81 ± 4	<b>97 ± 1</b>	49 ± 4	
	25 ± 1	27 ± 2	19 ± 2	33 ± 1	75 ± 3	<b>95 ± 1</b>	46 ± 2	
Attacks on output								
MNIST	35 ± 2	56 ± 3	38 ± 2	48 ± 4	29 ± 3	<b>69 ± 4</b>	46 ± 2	
SVHN	19 ± 1	19 ± 2	13 ± 1	27 ± 2	21 ± 2	<b>34 ± 2</b>	22 ± 1	
CelebA	27 ± 1	24 ± 1	31 ± 3	35 ± 1	29 ± 2	<b>40 ± 1</b>	31 ± 1	
	27 ± 1	33 ± 3	27 ± 2	37 ± 2	26 ± 1	<b>47 ± 3</b>	33 ± 1	
All attacks								
MNIST	31 ± 2	45 ± 3	32 ± 2	42 ± 3	50 ± 5	<b>80 ± 3</b>	47 ± 2	
SVHN	19 ± 1	19 ± 1	11 ± 1	27 ± 1	47 ± 7	<b>65 ± 7</b>	31 ± 2	
CelebA	29 ± 1	26 ± 1	26 ± 2	36 ± 1	55 ± 6	<b>68 ± 7</b>	40 ± 2	
	26 ± 1	30 ± 2	23 ± 1	35 ± 1	51 ± 4	<b>71 ± 3</b>	39 ± 1	

\* Attention mechanism disabled.

# MAIN FINDINGS



\* Attention mechanism disabled.



# MAIN FINDINGS



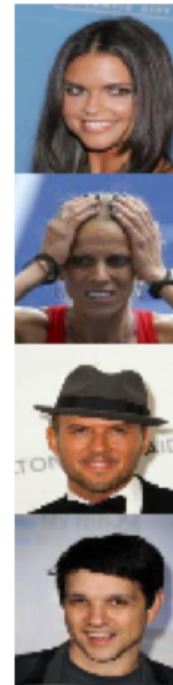
Original



Fully-connected VAE



DRAW



Target

# CONCLUSIONS

- ✓ We can attack autoencoders with adversarial images, by targeting their internal representations;
- ✓ The attack forces the autoencoder to reconstruct a different image;
- ✓ Autoencoders are, however, robust: success cases are hard to find and must be regularized “by hand”;
- ✓ The attack has a linear “give-and-take”: success in approaching the target output is proportional to the distortion of the input;
- ✓ The proposed metric (AUDDC) correlates well with qualitative results and provides a measure of robustness;
- ✓ DRAW is the most resistant architecture: attention and recurrence hinders the attack.

# Adversarial Attacks on Variational Autoencoders

[tabacof@dca.fee.unicamp.br](mailto:tabacof@dca.fee.unicamp.br)

