## **AVAE: Adversarial Variational Auto Encoder**

Antoine Plumerault<sup>\*,†</sup>, Hervé Le Borgne<sup>\*</sup>, Céline Hudelot<sup>†</sup>

\*: CEA List, †: Centrale Supelec

#### Two main types of generative models

٠	VAEs	have	several	advantages	over	GANs
---	------	------	---------	------------	------	------

GAN	VAE	
+ realistic images	+ disentangled latent space	
<ul> <li>mode collapse</li> </ul>	+ encoder model	
- difficult to invert	+ easy to train	
	<ul> <li>blurry images</li> </ul>	





Problematic: VAE fail to produce realistic images (w.r.t GANs)

- How can we explain this lack of realism ?
- Can we combine the best of VAEs and GANs ?

# Which problem for VAEs to produce realistic images ?

1. Information bottleneck:

$$\mathcal{L}_{\text{VAE}} = \underbrace{\mathbb{E}\left[\mathbb{E}_{q_{\theta_{e}}(z|x)}\left[-\log p_{\theta_{d}}(x|z)\right]\right]}_{\text{reconstruction error}} + \underbrace{l_{\theta}(x;z)}_{\text{mutual information}}$$
(1)  
+ 
$$\underbrace{\mathsf{KL}(p_{\theta_{e}}(z)||p(z))}_{\text{Figure 1}}$$

prior on z

- $\rightarrow$  incomplete information
- $\rightarrow$  mean value of all possible images
- $\rightarrow$  blurry results
- 2. Underestimation of natural image manifold dimensionality:
  - $\rightarrow$  approximation of the manifold with a simpler one
  - $\rightarrow$  uncertainty on other dimensions responsible of smaller variations (e.g. textures)
  - $\rightarrow$  mean value of all possible images
  - $\rightarrow$  blurry results (no texture in images)

## How GANs are able to produce realistic images ?

GANs also underestimate the dimension of the natural image manifold.

- $\rightarrow$  Question: How are they able to produce realistic images ?
- $\rightarrow$  Answer: Mode collapse !  $\rightarrow$  only a few but plausible texture configurations are generated.

Illustration on a toy example



#### dots: data dotted line: dashed line: points VAE manifold GAN manifold

#### How to solve the VAE problem ?

Objective: Create a reconstruction error  $\mathcal{L}_{\mathcal{Z}}$ : convert -density 400 input.pdf

picture.pngthat is powerful enough to favor accurate reconstructions.that does not favor blurry reconstruction to allow realistic reconstructions.



- cylinders: real data high-dimensional manifold
- black line: low-dimensional manifold of VAEs reconstructions
- arrows: gradient of different losses

### What properties such a reconstruction loss should satisfy ?

With reconstruction errors of the form  $\mathcal{L}_{\mathcal{Z}}(\hat{x}, x) = \frac{1}{2} ||f(\hat{x}) - g(x)||^2$  where:

- f is an arbitrary differentiable function
- g is a stochastic function

Optimal solutions  $\hat{x}^*(z)$  verifies:

 $f(\hat{x}^*(z)) = \mathbb{E}_{g(x) \sim p_{\theta_e}(g(x)|z)}[g(x)]$ (2)

- f(x̂) should carry the maximum of information about x̂ and g(x) should be close to f(x).
- Common optimum with the GAN objective  $\iff p(f(\hat{x}^*(z))) = p(f(x))$  for  $z \sim p(z)$  and  $x \sim p_D(x)$ .

#### A simple example: the MSE

 $\mathcal{L}_{\mathcal{Z}}(\hat{x}, x) = MSE(\hat{x}, x) = \frac{1}{2} ||\hat{x} - x||^2 \rightarrow \text{optimal solution:}$ 

 $\hat{x}^*(z) = \mathbb{E}_{x \sim p_{\theta_e}(x|z)}[x]$ (3)

- f(x̂) carry all the information about x̂ as it is the identity, and g(x) = f(x).
- Optimal solution = mix of likely solutions
   → blurry / unrealistic image.
   p(f(x̂\*(z))) = p(x̂\*(z)) ≠ p(x) =

### p(f(x)) for $z \sim p(z)$ and $x \sim p_{\mathcal{D}}(x)$ .

#### The AVAE framework

With:  

$$f(\hat{x}) = \frac{\mu_{\theta_e}(\hat{x})}{\sigma_{\theta_e}}$$

$$g(x) = \frac{\sqrt{1 - \sigma_{\theta_e}^2}}{\sigma_{\theta_e}} z$$

$$\mathcal{L}_{\mathcal{Z}}(\hat{x}) = \frac{1}{2} \left\| \frac{\mu_{\theta_e}(x) - \sqrt{1 - \sigma_{\theta_e}^2} z}{\sigma_{\theta_e}} \right\|^2$$
(4)

- $f(\hat{x})$  carry the information about  $\hat{x}$ contained in z, and  $g(x) = \frac{\sqrt{1-\sigma_{\theta_e}^2}}{\sigma_{\theta_e}} z \approx \frac{\mu_{\theta_e}(x)}{\sigma_{\theta_e}} = f(x)$
- $\mu_{\theta_{e}}(\hat{x}^{*}(z)) = \sqrt{1 \sigma_{\theta_{e}}^{2}} z \rightarrow p(\mu_{\theta_{e}}(\hat{x}^{*}(z))) = \mathcal{N}(\mu_{\theta_{e}}(x); 0, I \Sigma) = p(\mu_{\theta_{e}}(x)).$



#### Results

#### Quantitative results on CelebA:

		$MSE\downarrow$	LPIPS $\downarrow$	FID ↓	
	VAE	$0.03\pm0.00$	$0.18\pm0.00$	$60.04 \pm 0.47$	
	GAN	_	_	$14.54\pm0.41$	
ŀ	VAE/GAN	$0.07\pm0.00$	$0.09\pm0.00$	$26.45 \pm 4.66$	
	BiGAN	$0.18\pm0.01$	$0.16\pm0.00$	$18.49 \pm 5.06$	
	Ours	$0.05\pm0.00$	$0.11\pm0.00$	$15.01\pm0.82$	

- MSE: favorable to the VAE a priori.
- MSE: favorable to our approach a priori.
- LPIPS & FID: favorable to VAE/GAN a priori.

#### Qualitative results:







