

# Anomaly detection with generative methods

Coloma Ballester

Universitat Pompeu Fabra

Machine Learning in Astrophysics, March 5th, 2021

Joint work with



Pierrick Chatillon



Joan Sintes

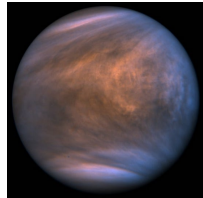
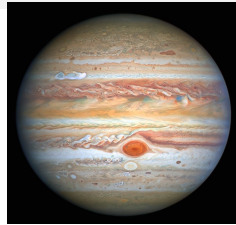
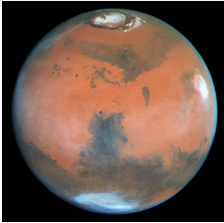


Patricia Vitoria

1. HistoryAD: An adversarial approach to anomaly detection
  - Introduction
  - Generative approaches
  - Generative adversarial-based anomaly detection
  - Proposed method
  - Experimental results
2. Semantic image inpainting through improved Wasserstein generative adversarial networks

Anomaly detection.

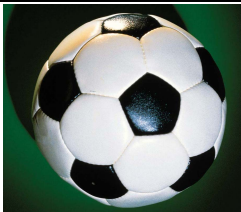
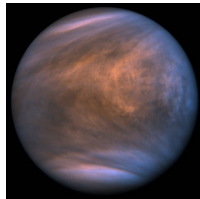
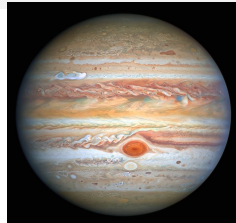
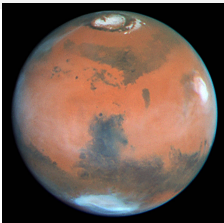
Anomalous?





Anomaly detection.

Anomalous?



Out-of-distribution

Anomaly detection.

Anomalous?



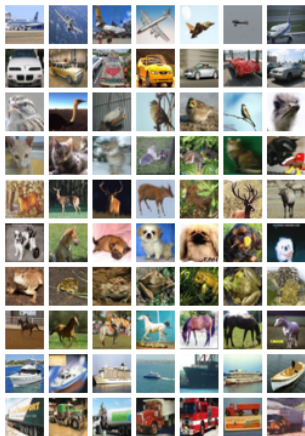
Out-of-distribution

# Anomalous? Out-of-distribution

MNIST



KMNIST



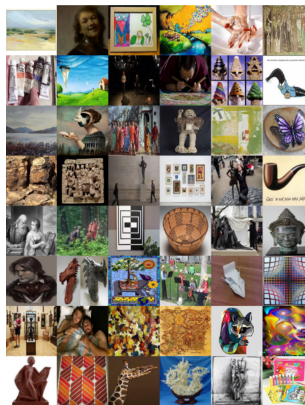
CIFAR-10



CelebA



SVHN



Tiny ImageNet



Generative methods that produce novel samples from high-dimensional data distributions, such as images, are currently very used.

Some of the most prominent approaches are

- autoregressive models (e.g., Van den Oord, et al., 2016)
- variational autoencoders (VAE) (e.g., Kingma, a Welling, 2013)
- generative adversarial networks (GAN) (e.g., Goodfellow et al., 2014)

- The goal is to learn the probability distribution  $\mathbb{P}_{\text{real}}$  of the given data, usually by learning its probability density.
- Classically, it is tackled by defining a parametric family of densities  $(P_{\theta})_{\theta \in \mathbb{R}^d}$  and finding the one that maximizes the likelihood on our data.

Problems:

- Does the real data distribution  $\mathbb{P}_{\text{real}}$  admit a density?
- Does the model density  $P_{\theta}$  exist?

- But, although natural images belong to high dimensional spaces, they contain geometric and semantic structure.
- Thus, following <sup>1</sup>, rather than estimating the density of  $\mathbb{P}_{\text{real}}$ , which may not exist, we can define a random variable  $Z$  with a fixed distribution  $\mathbb{P}_Z$  and pass it through a **parametric function**  $G_\theta : \mathcal{Z} \rightarrow \mathcal{X}$  (typically a neural network) that directly generates samples following a certain distribution  $\mathbb{P}_{\text{real}}$ .
- $\mathbb{P}_G = G_\# \mathbb{P}_Z$ , the **pushforward measure** of  $\mathbb{P}_Z$  through  $G$  (parametric density  $G_\# \mathbb{P}_Z$  through the neural network  $G$ ).
- By varying  $\theta$ , we can change this distribution  $\mathbb{P}_{G_\theta}$  and make it close to (converge, if possible) the real data distribution  $\mathbb{P}_{\text{real}}$ .

---

<sup>1</sup> Arjovsky, Chintala, and Bottou. Wasserstein GAN. 2017.

<sup>2</sup> Many authors: Peyré, Genevay, Cuturi, Brenier, Dieng, Lunz, Delon, Willett, Dumoulin, Schönlieb, Berthelot, Bengio, and many more.

First of the many GAN's papers (2014):

## Generative Adversarial Nets

---

**Ian J. Goodfellow<sup>\*</sup>, Jean Pouget-Abadie<sup>†</sup>, Mehdi Mirza, Bing Xu, David Warde-Farley,  
Sherjil Ozair<sup>‡</sup>, Aaron Courville, Yoshua Bengio<sup>§</sup>**

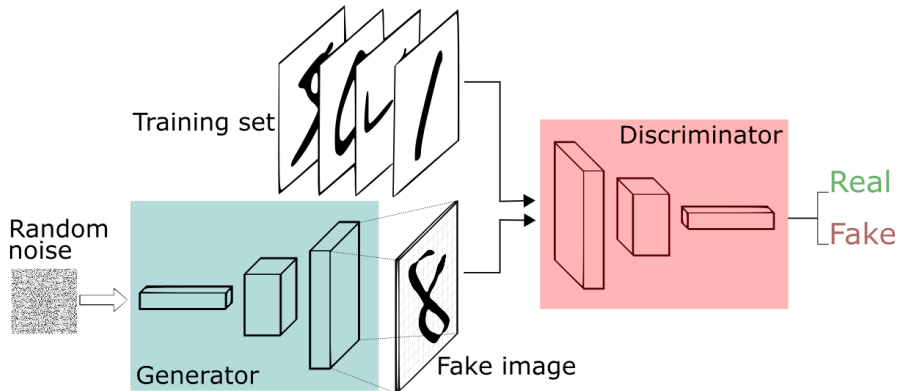
Département d'informatique et de recherche opérationnelle  
Université de Montréal  
Montréal, QC H3C 3J7

### Abstract

We propose a new framework for estimating generative models via an adversarial process, in which we simultaneously train two models: a generative model  $G$  that captures the data distribution, and a discriminative model  $D$  that estimates the probability that a sample came from the training data rather than  $G$ . The training procedure for  $G$  is to maximize the probability of  $D$  making a mistake. This framework corresponds to a minimax two-player game. In the space of arbitrary



# GAN Framework



# Training GAN

Original GAN objective:

$$\min_G \max_D V(G, D) = \mathbb{E}_{x \sim \mathbb{P}_{\text{real}}} [\log D(x)] + \mathbb{E}_{z \sim \mathbb{P}_z} [\log(1 - D(G(z)))]$$

Min max iterations: iterate the "two steps" until convergence (which may not happen)

- Updating the discriminator should make it better at discriminating between real images and generated ones (discriminator improves).
- Updating the generator makes it better at fooling the current discriminator (generator improves).

Eventually (we hope) that the generator gets so good that it is impossible for the discriminator to tell the difference between real and generated images. Discriminator guess = 0.5.

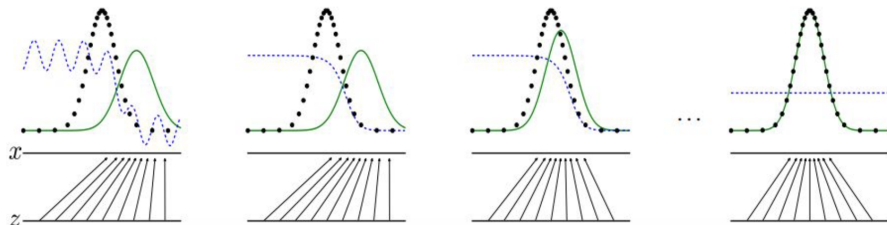


Image credits: Santiago Pascual, 2018

# Distances and divergences between probability distributions

- Vanila-GAN training objective:

$$\min_G \max_D V(G, D) = \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\text{real}}} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim \mathbb{P}_{\mathbf{z}}} [\log(1 - D(G(\mathbf{z})))]$$

- Under optimal discriminator  $D_G^*(\mathbf{x}) = \frac{p_{\text{real}}(\mathbf{x})}{p_{\text{real}}(\mathbf{x}) + p_G(\mathbf{x})}$ ,

$$\min_G V(G, D_G^*) = -\log(4) + 2 \cdot \delta_{\text{JS}}(\mathbb{P}_{\text{real}}, \mathbb{P}_G)$$

(where  $p_{\text{real}}, p_G$  densities)

- Jensen-Shannon Divergence:

$$\delta_{\text{JS}}(\mathbb{P}_1, \mathbb{P}_2) \triangleq \frac{1}{2} \left[ D_{\text{KL}} \left( \mathbb{P}_1 \parallel \frac{\mathbb{P}_1 + \mathbb{P}_2}{2} \right) + D_{\text{KL}} \left( \mathbb{P}_2 \parallel \frac{\mathbb{P}_1 + \mathbb{P}_2}{2} \right) \right]$$

where  $\mathbb{P}_1, \mathbb{P}_2 \in \text{Prob}(\mathcal{X})$ , space of probability distributions defined on  $\mathcal{X}$ ,  $\mathcal{X}$  a compact metric set (e.g., the space of images)

Wasserstein 1 Distance:

$$W_1(\mathbb{P}_1, \mathbb{P}_2) = \inf_{\pi \in \Pi(\mathbb{P}_1, \mathbb{P}_2)} \mathbb{E}_{x, y \sim \pi} (\|x - y\|),$$

By Kantorovitch-Rubenstein duality:

$$W_1(\mathbb{P}_1, \mathbb{P}_2) = \sup_{D \in \mathcal{D}} (\mathbb{E}_{x \sim \mathbb{P}_1} [D(x)] - \mathbb{E}_{y \sim \mathbb{P}_2} [D(y)]).$$

where  $\mathcal{D}$  denotes the set of 1-Lipschitz functions (i.e., <sup>1</sup>, the set of c-convex functions for the cost function  $c(x, y) = \|x - y\|$ ).

In these articles <sup>2,3</sup>, the training objectives are adapted to minimize  $W_1(\mathbb{P}_{\text{real}}, \mathbb{P}_G)$

---

<sup>1</sup> Villani. Optimal transport: old and new. 2008

<sup>2</sup> Arjovsky, et al. Wasserstein GAN. 2017.

<sup>3</sup> Gulrajani, et al. Improved training of Wasserstein GANs. 2017.

$$\delta(\mathbb{P}_1, \mathbb{P}_2) = \sup_{A \in \mathcal{F}} |\mathbb{P}_1(A) - \mathbb{P}_2(A)|$$

which represents the choice  $c(x, y) = \mathbb{1}_{x \neq y}$  in the optimal transport problem <sup>1</sup>.

$$\delta(\mathbb{P}_1, \mathbb{P}_2) = \frac{1}{2} \|\mathbb{P}_1 - \mathbb{P}_2\|_{TV}.$$

Kantorovitch-Rubenstein duality:

$$\delta(\mathbb{P}_1, \mathbb{P}_2) = \sup_{-1 \leq D \leq 1} (\mathbb{E}_{x \sim \mathbb{P}_1}[D(x)] - \mathbb{E}_{y \sim \mathbb{P}_2}[D(y)])$$

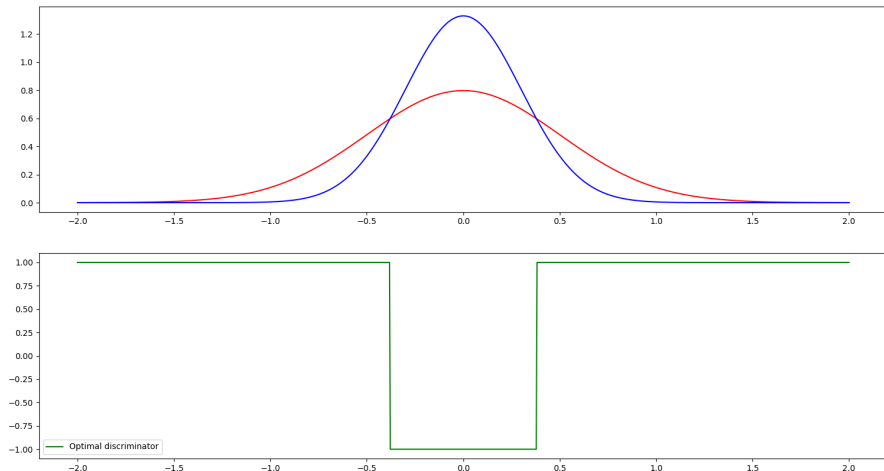
Taking  $\mu = \mathbb{P}_1 - \mathbb{P}_2$ , a signed measure, and  $(P, N)$  its Hahn decomposition ( $P = \{\mathbb{P}_1 > \mathbb{P}_2\}$ ), we can define the **optimal discriminator**  $D^* := \mathbb{1}_P - \mathbb{1}_N$

---

<sup>1</sup> Villani. Optimal transport: old and new. 2008

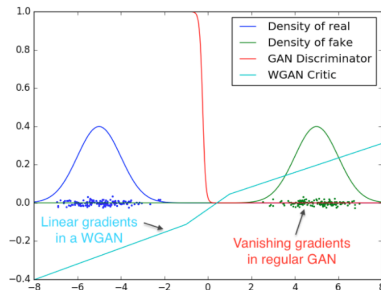
# Distances. Optimal dual variable

$\mathbb{P}_1$ ,  $\mathbb{P}_2$ , and the optimal  $D^* = \mathbb{1}_P - \mathbb{1}_N$



# Wasserstein Generative Adversarial Network

- Main vanilla-GANs problems: vanishing gradients, mode collapse <sup>1</sup>, non-continuity.
- In these articles <sup>2,3</sup>, the training objectives are adapted to minimize  $\mathbb{W}_1(\mathbb{P}_{\text{real}}, \mathbb{P}_G)$ .
- Wasserstein GAN (WGAN) uses an approximation of the Wasserstein distance. It is continuous everywhere and differentiable almost everywhere.



<sup>1</sup> Arjovsky, and Bottou. Towards principled methods for training generative adversarial networks. 2017

<sup>2</sup> Arjovsky, et al. Wasserstein GAN. 2017.

<sup>3</sup> Gulrajani, et al. Improved training of Wasserstein GANs. 2017.

**Theorem.** Let  $\mathbb{P}_{\text{real}}$  a fixed distribution over  $\mathcal{X}$ . Let  $Z$  be a random variable (e.g Gaussian) over another space  $\mathcal{Z}$ . Let  $G : \mathcal{Z} \times \mathbb{R}^d \rightarrow \mathcal{X}$  be a function, that will be denoted  $G_\theta(z)$  with  $z$  the first coordinate and  $\theta$  the second. Let  $\mathbb{P}_\theta$  denote the distribution of  $G_\theta(Z)$ . Then,

- 1 If  $G$  is continuous in  $\theta$ , so is  $W(\mathbb{P}_{\text{real}}, \mathbb{P}_\theta)$ .
- 2 If  $G$  is locally Lipschitz and satisfies the regularity assumption  $\mathbb{E}_{z \sim p}[L(\theta, z)] < +\infty$  on the local Lipschitz constants  $L(\theta, z)$ , then  $W(\mathbb{P}_{\text{real}}, \mathbb{P}_\theta)$  is continuous everywhere, and differentiable almost everywhere.
- 3 Statements 1-2 are false for the Jensen-Shannon divergence  $JS(\mathbb{P}_{\text{real}}, \mathbb{P}_\theta)$  and all the KLs.

The authors show that

- The assumption in 2 is satisfied for any feedforward neural network  $G_\theta$ , and thus  $W(\mathbb{P}_{\text{real}}, \mathbb{P}_\theta)$  is continuous everywhere and differentiable almost everywhere.
- $\nabla_\theta W(P_r, P_\theta) = -\mathbb{E}_{z \sim p(z)}[\nabla_\theta f_w(g_\theta(z))]$ , when both terms are well defined.

---

<sup>1</sup> Arjovsky, et al. Wasserstein GAN. 2017.



How to ensure to have a 1-Lipschitz discriminator?

- **WGAN: Weight clipping:**<sup>1</sup> clipping the parameters of the discriminators
  - Problems, including that it reduces the capacity of the discriminator.
- **WGAN-GP: Gradient penalty:**<sup>2</sup> penalizing the norm of discriminator gradients with respect to data samples during training to be less than 1.

$$\min_G \max_{D \in \mathcal{D}} \mathbb{E}_{\tilde{x} \sim \mathbb{P}_{real}} [D(\tilde{x})] - \mathbb{E}_{x \sim \mathbb{P}_G} [D(x)] - \lambda \mathbb{E}_{\tilde{x} \sim \mathbb{P}_{\tilde{x}}} [(\|\nabla_{\tilde{x}} D(\tilde{x})\|_2 - 1)^2]$$

where  $\mathbb{P}_{\tilde{x}}$  is implicitly defined sampling uniformly along straight lines between pairs of point sampled from the data distribution  $\mathbb{P}_{real}$  and the generator distribution  $\mathbb{P}_G$ .

- The dual variable  $D$  is expected to be positive on real data samples and negative on generated ones.

---

<sup>1</sup> Arjovsky, et al. Wasserstein GAN 2017.

<sup>2</sup> Gulrajani, et al. Improved Training of Wasserstein GANs. 2017.

## How can this be used for anomaly detection?

# GANs in anomaly detection, main strategies

- If we learned to generate normal data, only normal data can be reconstructed with such a generator <sup>1</sup>
- Use or create an auxiliary dataset of corrupted data (out of distribution) as negative data for a classifier (outlier exposure) <sup>2,3</sup>
- Corrupt the generator to provide anomalies <sup>4</sup>

---

<sup>1</sup> Schlegl, et al. AnoGAN. Unsupervised Anomaly Detection with Generative Adversarial Networks to Guide Marker Discovery. 2017.

<sup>2</sup> Hendrycks, et al. Deep Anomaly Detection with Outlier Exposure. 2019.

<sup>3</sup> Meinke& Hein: Towards neural networks which provably know when they don't know. 2020.

<sup>4</sup> Ngo, et al. Fence GAN: Towards Better Anomaly Detection. 2019.

History-based anomaly detector: an adversarial approach to anomaly detection

Joint work with Pierrick Chatillon

Method's idea: Oscillation during training provides anomalies

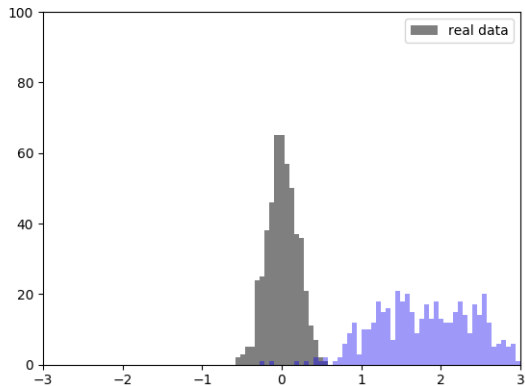


Figure: Generated distribution  $p_{G_t}$  oscillating around  $p_{\text{real}}$

Method's idea: Oscillation during training provides anomalies

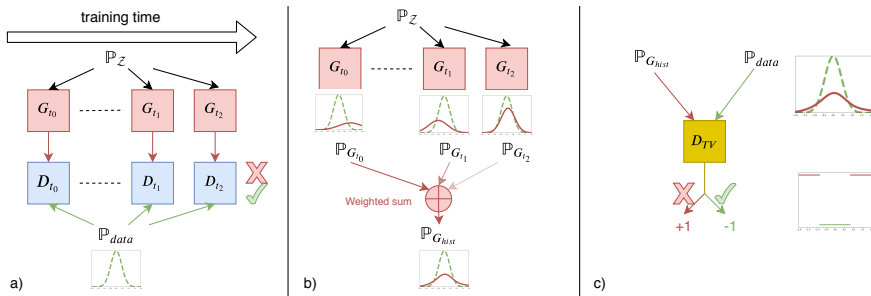
Figure: Generated distribution  $p_{G_t}$  oscillating around  $p_{\text{real}}$

Method's idea: Oscillation during training provides anomalies

Figure: Generated distribution  $p_{G_t}$  oscillating around  $p_{\text{real}}$

## Method: 'HistoryAD'

- Train a W-GAN on normal data while saving states
- Craft an 'anomalous' distribution
- Classify normal and anomalous data with the total variation framework
  - We use the obtained classifier  $D_{TV}$  as anomaly detector





## Which anomalous distribution?

Figure: Generated distribution  $p_{G_t}$  oscillating around  $p_{\text{real}}$

**Hypothesis:**  $\text{supp}(\mathbb{P}_{\text{real}}) \subset \text{supp}(\mathbb{P}_{G_{\text{hist}}})$ .

$\mathbb{P}_{G_{\text{hist}}}$  is our anomalous distribution: it is a **weighted average** of  $\mathbb{P}_{G_t}$  for the different states of  $G$  during training.

$$\mathbb{P}_{G_{\text{hist}}} \triangleq \int_1^{n_{\text{epoch}}} c \cdot G_t(\mathbb{P}_Z) \cdot e^{-\beta t} dt$$

To sample  $\mathbb{P}_{G_{\text{hist}}}$  :

- During W-GAN training, save the Generator's state at regular time steps
- When training  $D_{TV}$ , sample  $t$  along the exponential distribution, then sample  $z$  from  $\mathbb{P}_Z$ , and finally compute  $G_t(z)$ .  
i.e., in practice, we approximate  $\mathbb{P}_{G_{\text{hist}}}$  by sampling data from  $\mathbb{P}_{G_t}$  where  $t$  is a random variable of density of probability  $c \cdot \mathbb{1}_{[\alpha, n_{\text{epochs}}]} \cdot e^{-\beta t}$

$D_{TV}$  training objective:

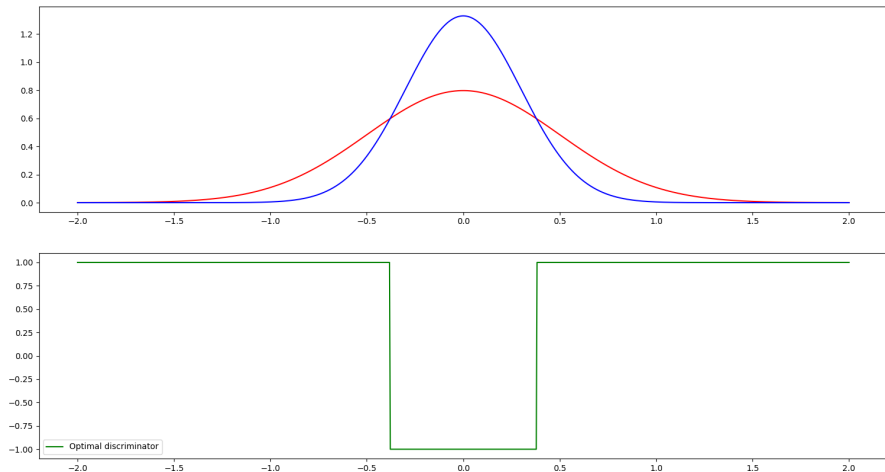
$$\sup_{-1 \leq D \leq 1} \left( \mathbb{E}_{x \sim \mathbb{P}_{\text{real}}} [D(x)] - \mathbb{E}_{y \sim \mathbb{P}_{G_{\text{hist}}}} [D(y)] \right)$$

For the optimal  $D_{TV}$ , the expression above is equal to  $\delta(\mathbb{P}_{\text{real}}, \mathbb{P}_{G_{\text{hist}}})$ , where  $\delta$  is the **total variation distance**:

$$\delta(\mathbb{P}_1, \mathbb{P}_2) = \sup_{A \text{ measurable}} |\mathbb{P}_1(A) - \mathbb{P}_2(A)|$$

# Optimal discriminator

$P_{\text{real}}$ ,  $P_{G_{\text{hist}}}$ , and the optimal  $D_{TV}$



Our method does **not depend** on any specific **perturbation of the GAN objective**, but rather on an **intrinsic property** of GAN training: the **oscillation** of the generated distribution around real data.

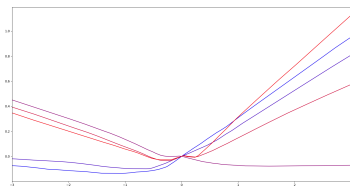
Our method can be seen as an extension of the 'early stopping in GANs' <sup>1</sup>

---

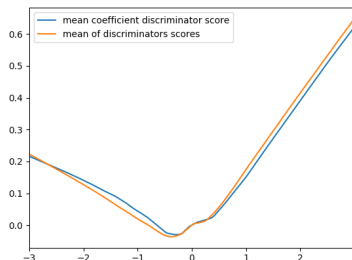
<sup>1</sup> Gu, et al. Semi-Supervised Outlier Detection Using a Generative and Adversary Framework. 2018.

## Which initialization for $D_{TV}$ ? Merging discrimination information

- The mean of the different discriminator states is a good candidate
- The network with mean parameters ( $W_{D_{TV}}^0 = \int_1^{n_{epoch}} c \cdot W_{D_t} \cdot e^{-\beta t} dt$ ) is a good approximation <sup>1</sup>



(b) Discriminator score output during training (from blue to red)



(c) Comparison of the average outputs of saved discriminators during training and the output of a discriminator with average coefficients.

<sup>1</sup> Wang, et al. Deep Network Interpolation for Continuous Imagery Effect Transition. 2019.

During  $D_{TV}$  training, we want  $-1 \leq D_{TV} \leq 1$ :

- ✗ non-linearity, bad gradient behavior, as the solution tends to  $-1$  or  $1$  almost everywhere.
- ✓  $\lambda_{bounded} \cdot d(x, [-1, 1])^2$ , smooth constraint loss term



One sample from each class of MNIST.



## TV-GAN Optimization process

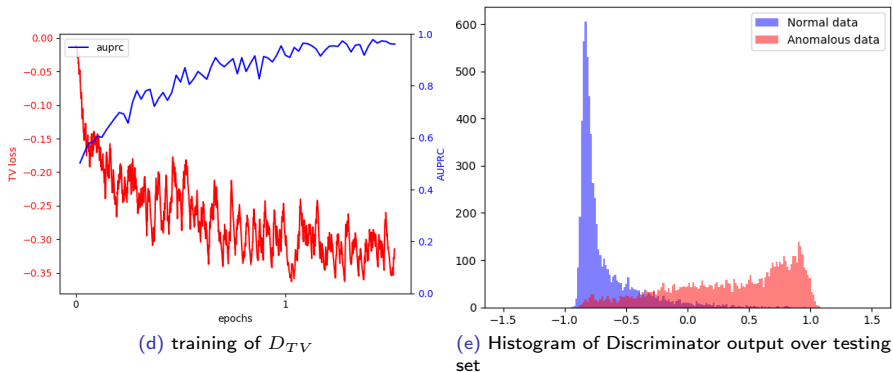
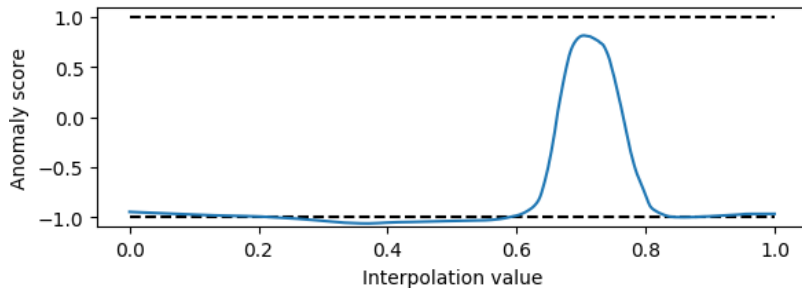


Figure: Successful training on digit 1, AUPCR around 0.95

## Experiments. Latent space linear interpolation

- anomaly score of  $G((1 - t)\mathbf{z}_1 + t\mathbf{z}_2)$



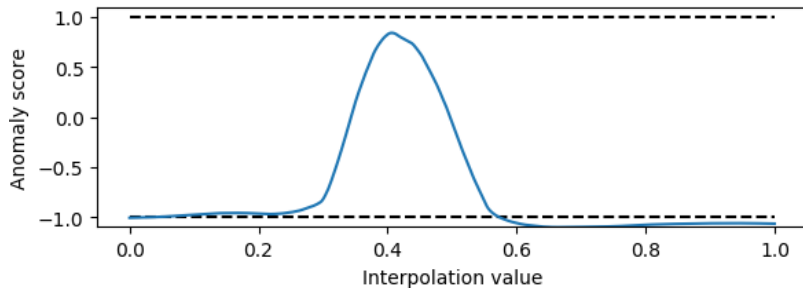
Corresponding interpolated images



Figure: History-GAN trained on MNIST

## Experiments. Latent space linear interpolation

- anomaly score of  $G((1 - t)\mathbf{z}_1 + t\mathbf{z}_2)$



Corresponding interpolated images

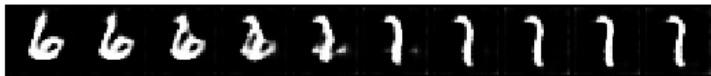


Figure: History-GAN trained on MNIST

## Experiments. Gaussian Noise on normal data

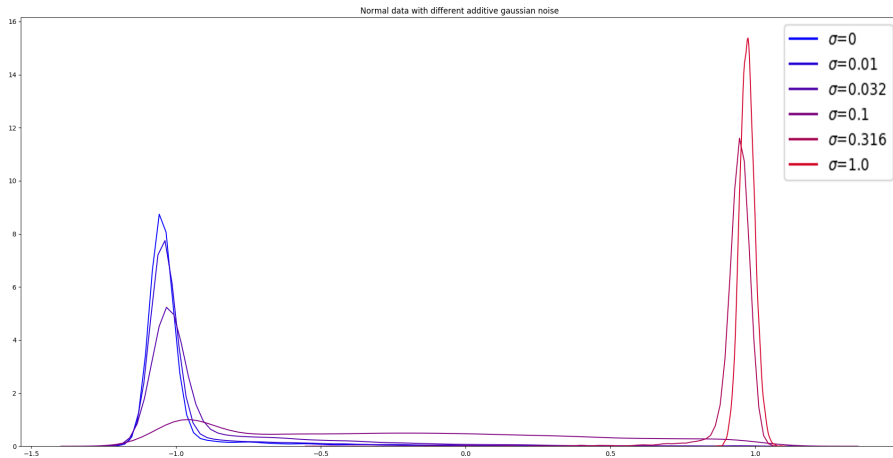
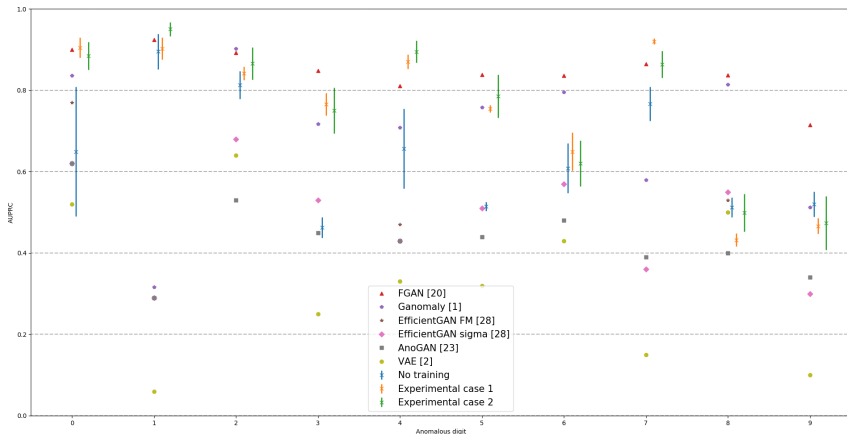


Figure: Density of distribution of anomaly scores

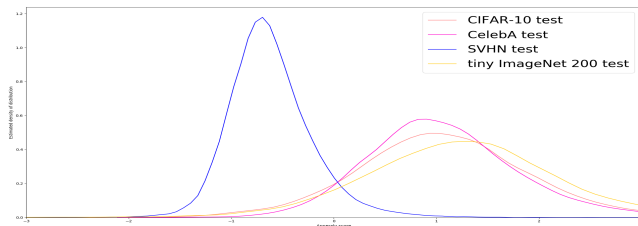
## Results. Comparison to others on MNIST experiment



**Figure:** mean AUPRC for each digit of MNIST for experimental case 2, compared with other methods (performances of methods provided by the authors of Fence-GAN)

# Results. Multiple datasets evaluation

## Method trained on SVHN and evaluated on several datasets



Approximate density of anomaly score distribution

test split	CIFAR-10	CelebA	Tiny ImageNet
AUPRC	0.941	0.976	0.949

Table AUPRC for SVHN compared to other datasets.

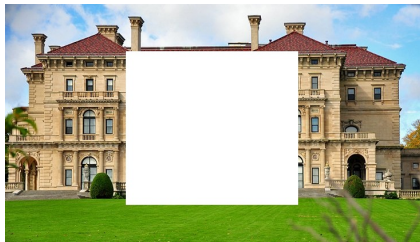
### 3. Image inpainting through an adversarial strategy

Joint work with Patricia Vitoria and Joan Sintes

# Image inpainting

- Image inpainting is also known as image completion, disocclusion or object removal. It aims to obtain a visually plausible completion of the image in a region in which data is missing due to damage or occlusion.

- Problem: When missing regions are large and moreover the missing information is unique in the sense that the information and redundancy available in the image is not useful to guide the completion, the task becomes even more challenging.





# Image inpainting

- Image inpainting is also known as image completion, disocclusion or object removal. It aims to obtain a visually plausible completion of the image in a region in which data is missing due to damage or occlusion.

- Problem: When missing regions are large and moreover the missing information is unique in the sense that the information and redundancy available in the image is not useful to guide the completion, the task becomes even more challenging.

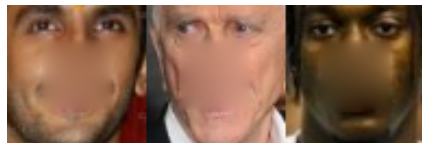
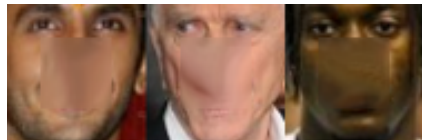
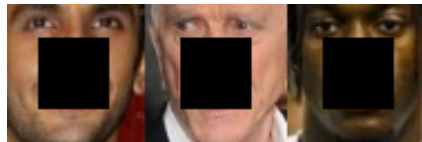


# Approaches for image inpainting

**Classical Methods** use redundancy of the incomplete input image.

- **Local Methods** (Masnou and Morel, 1998; Chan and Shen, 2001; Ballester et al. 2001; Esedoglu et al. 2003; Deriche et al. 2005; Figueiredo et al. 2007; Bertozzi et al. 2007; Weickert et al. 2012; Schönlieb et al. 14;...)

- **Non-local Methods** (Efros and Leung, 1999; Demanet et al, 2003; Wexler et al, 2004; Criminisi; Aujol; Bugeau; and many more).



# Approaches for image inpainting

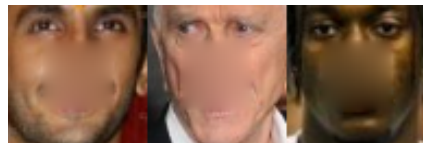
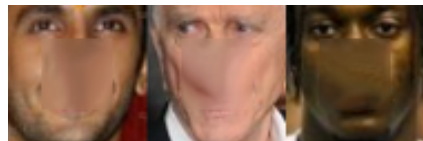
**Classical Methods** use redundancy of the incomplete input image.

- **Local Methods** (Masnou and Morel, 1998; Chan and Shen, 2001; Ballester et al. 2001; Esedoglu et al. 2003; Deriche et al. 2005; Figueiredo et al. 2007; Bertozzi et al. 2007; Weickert et al. 2012; Schönlieb et al. 14;...)

- **Non-local Methods** (Efros and Leung, 1999; Demanet et al, 2003; Wexler et al, 2004; Criminisi; Aujol; Bugeau; and many more).

**Semantic inpainting** infers large missing regions based on perception and image semantics (Nguyen et al. 2016; Yeh et al. 2017; Burlin et al. 2017; Pathak; Efros; Vedaldi; Zeng; and many more).

We use the understanding of more abstract and high level information provided by Generative Adversarial Models



# Semantic image inpainting

Method, based on an self-supervised adversarial strategy followed by an energy-based completion algorithm:

- **1st Step:** given a **dataset of (non-corrupted) images**, the data latent space is learned via an improved version of the Wasserstein GAN

$$\min_G \max_{D \in \mathcal{D}} \mathbb{E}_{\tilde{x} \sim \mathbb{P}_{real}} [D(\tilde{x})] - \mathbb{E}_{x \sim \mathbb{P}_G} [D(x)] - \lambda \mathbb{E}_{\hat{x} \sim \mathbb{P}_{\hat{x}}} [(\|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1)^2]$$

- **2nd Step:** given an **incomplete image**  $y$  and the converged generative adversarial  $G$  and  $D$ , a minimization procedure is performed to infer the missing content of the incomplete image by conditioning on the known regions

$$\hat{z} = \arg \min_z \{ \mathcal{L}_c(z|y, M) + \eta \mathcal{L}_p(z) \}$$

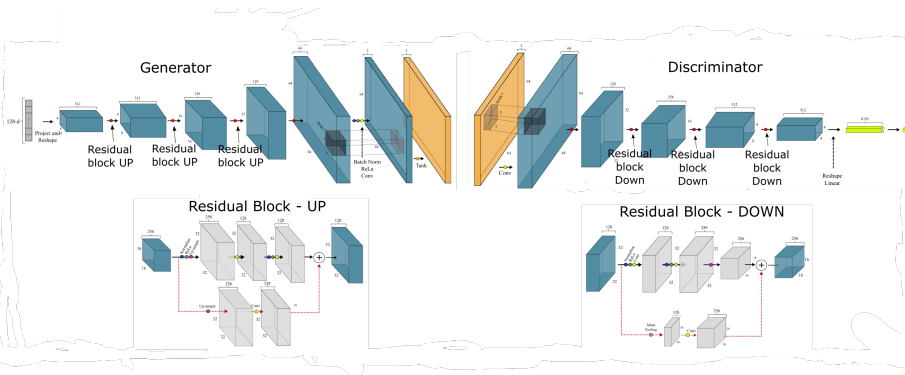
where  $\mathcal{L}_p$  is the **prior loss** and  $\mathcal{L}_c$  the **contextual loss** defined as

$$\mathcal{L}_c(z|y, M) = \alpha W \|M(G(z) - y)\| + \beta W \|M(\nabla G(z) - \nabla y)\|$$

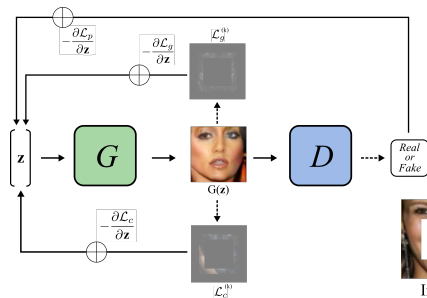
with  $M$  a binary mask equal to 1 on the known pixels,  $W$  a weight mask,  $\alpha, \beta, \eta > 0$ , and  $\mathcal{L}_p(z) = -D(G(z))$  favours realistic images, similar to the samples that are used to train the generative model.

# 1st Step: Train our generative model

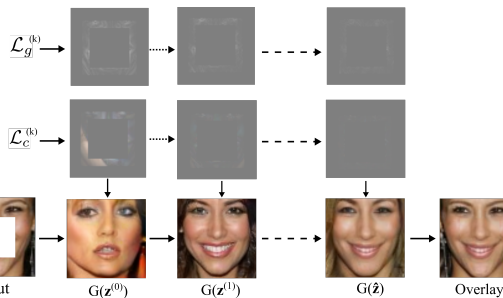
Architecture (based on the one of WGAN-GP plus some improvements)



## 2nd Step: Perform inpainting using an optimization method



Given a GAN model trained on real images, we iteratively update  $z$  to find the closest mapping on the latent image manifold, based on the designed loss function.



Manifold traversing when iteratively update  $z$  using back-propagation.  $z^{(0)}$  is random initialized;  $z^{(k)}$  denotes the result in  $k$ -th iteration; and  $\hat{z}$  denotes the final solution before the Poisson editing step is applied.

## 2st Step: Perform semantic inpainting. Losses

$$\hat{z} = \arg \min_z \{ \alpha L_c(z|y, M) + \beta L_g(z|y, M) + \eta L_p(z) \}$$

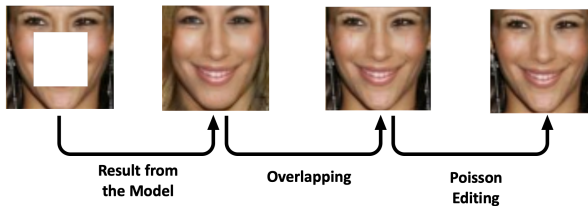
$$L_c(z|y, M) = W ||M(G(z) - y)||$$

$$L_g(z|y, M) = W ||M(\nabla G(z) - \nabla y)||$$

$$L_p(z) = -D_w(G_\theta(z))$$

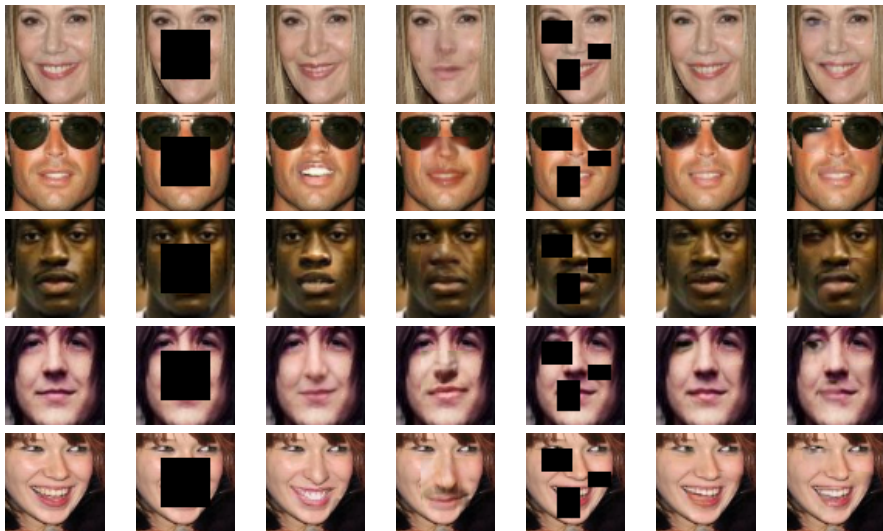
$$W(i) = \begin{cases} \sum_{j \in N_i} \frac{(1-M(j))}{|N_i|} & \text{if } M(i) \neq 0 \\ 0 & \text{if } M(i) = 0 \end{cases}$$

## 3rd Step: Poisson Editing





## Experimental results



Original

Masked

Ours

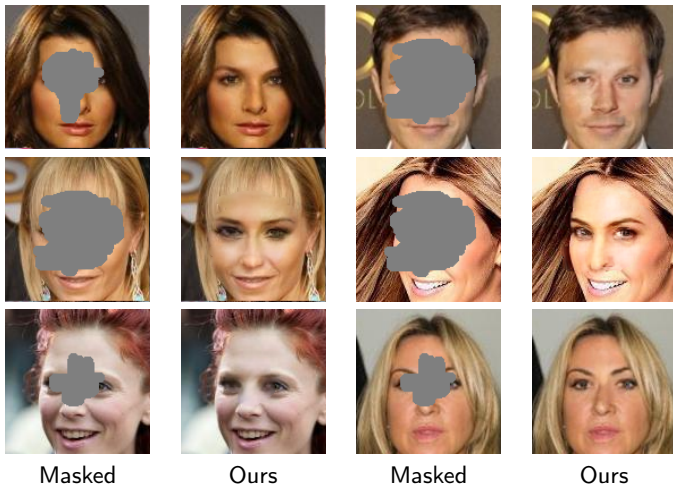
Yeh 2017

Masked

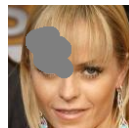
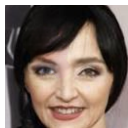
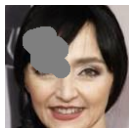
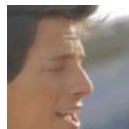
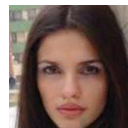
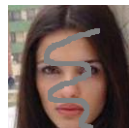
Ours

Yeh 2017

## Experimental results



## Experimental results



Masked

Ours

Masked

Ours

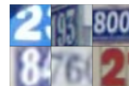
Masked

Ours

# Street View House Numbers (SVHN)



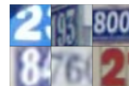
# Quantitative results



Loss formulation	CelebA dataset			SVHN dataset		
	MSE	PSNR	SSIM	MSE	PSNR	SSIM
(Yeh et al., 2017)	872.8672	18.7213	0.9071	1535.8693	16.2673	0.4925
(Yeh et al., 2017) adding gradient loss with $\alpha = 0.1$ , $\beta = 0.9$ and $\eta = 1.0$	832.9295	18.9247	0.9087	1566.8592	16.1805	0.4775
(Yeh et al., 2017) adding gradient loss with $\alpha = 0.5$ , $\beta = 0.5$ and $\eta = 1.0$	862.9393	18.7710	0.9117	1635.2378	15.9950	0.4931
(Yeh et al., 2017) adding gradient loss with $\alpha = 0.1$ , $\beta = 0.9$ and $\eta = 0.5$	<u>794.3374</u>	<u>19.1308</u>	<u>0.9130</u>	<u>1472.6770</u>	<u>16.4438</u>	<u>0.5041</u>
(Yeh et al., 2017) adding gradient loss with $\alpha = 0.5$ , $\beta = 0.5$ and $\eta = 0.5$	876.9104	18.7013	0.9063	1587.2998	16.1242	0.4818
Our proposed loss with $\alpha = 0.1$ , $\beta = 0.9$ and $\eta = 1.0$	855.3476	18.8094	0.9158	631.0078	20.1305	<b>0.8169</b>
Our proposed loss with $\alpha = 0.5$ , $\beta = 0.5$ and $\eta = 1.0$	<b>785.2562</b>	<b>19.1807</b>	<b>0.9196</b>	743.8718	19.4158	0.8030
Our proposed loss with $\alpha = 0.1$ , $\beta = 0.9$ and $\eta = 0.5$	862.4890	18.7733	0.9135	<b>622.9391</b>	<b>20.1863</b>	0.8005
Our proposed loss with $\alpha = 0.5$ , $\beta = 0.5$ and $\eta = 0.5$	833.9951	18.9192	0.9146	703.8026	19.6563	0.8000

Method	CelebA dataset			SVHN dataset		
	MSE	PSNR	SSIM	MSE	PSNR	SSIM
(Yeh et al., 2017)	622.1092	20.1921	0.9087	1531.4601	16.2797	0.4791
(Yeh et al., 2017) adding gradient loss with $\alpha = 0.1$ , $\beta = 0.9$ and $\eta = 1.0$	584.3051	20.4644	0.9067	1413.7107	16.6272	0.4875
(Yeh et al., 2017) adding gradient loss with $\alpha = 0.5$ , $\beta = 0.5$ and $\eta = 1.0$	600.9579	20.3424	0.9080	1427.5251	16.5850	0.4889
(Yeh et al., 2017) adding gradient loss with $\alpha = 0.1$ , $\beta = 0.9$ and $\eta = 0.5$	580.8126	20.4904	0.9115	1446.3560	16.5281	<u>0.5120</u>
(Yeh et al., 2017) adding gradient loss with $\alpha = 0.5$ , $\beta = 0.5$ and $\eta = 0.5$	<u>563.4620</u>	<u>20.6222</u>	0.9103	<u>1329.8546</u>	<u>16.8928</u>	0.4974
Our proposed loss with $\alpha = 0.1$ , $\beta = 0.9$ and $\eta = 1.0$	424.7942	21.8490	0.9281	168.9121	25.8542	0.8960
Our proposed loss with $\alpha = 0.5$ , $\beta = 0.5$ and $\eta = 1.0$	380.4035	22.3284	0.9314	221.7906	24.6714	<b>0.9018</b>
Our proposed loss with $\alpha = 0.1$ , $\beta = 0.9$ and $\eta = 0.5$	<b>321.3023</b>	<b>23.0617</b>	<b>0.9341</b>	<b>154.5582</b>	<b>26.2399</b>	0.8969
Our proposed loss with $\alpha = 0.5$ , $\beta = 0.5$ and $\eta = 0.5$	411.8664	21.9832	0.9292	171.7974	25.7806	0.8939

# Quantitative results



Loss formulation	CelebA dataset			SVHN dataset		
	MSE	PSNR	SSIM	MSE	PSNR	SSIM
(Yeh et al., 2017)	872.8672	18.7213	0.9071	1535.8693	16.2673	0.4925
(Yeh et al., 2017) adding gradient loss with $\alpha = 0.1$ , $\beta = 0.9$ and $\eta = 1.0$	832.9295	18.9247	0.9087	1566.8592	16.1805	0.4775
(Yeh et al., 2017) adding gradient loss with $\alpha = 0.5$ , $\beta = 0.5$ and $\eta = 1.0$	862.9393	18.7710	0.9117	1635.2378	15.9950	0.4931
(Yeh et al., 2017) adding gradient loss with $\alpha = 0.1$ , $\beta = 0.9$ and $\eta = 0.5$	<u>794.3374</u>	<u>19.1308</u>	<u>0.9130</u>	<u>1472.6770</u>	<u>16.4438</u>	<u>0.5041</u>
(Yeh et al., 2017) adding gradient loss with $\alpha = 0.5$ , $\beta = 0.5$ and $\eta = 0.5$	876.9104	18.7013	0.9063	1587.2998	16.1242	0.4818
Our proposed loss with $\alpha = 0.1$ , $\beta = 0.9$ and $\eta = 1.0$	855.3476	18.8094	0.9158	631.0078	20.1305	<b>0.8169</b>
Our proposed loss with $\alpha = 0.5$ , $\beta = 0.5$ and $\eta = 1.0$	<b>785.2562</b>	<b>19.1807</b>	<b>0.9196</b>	743.8718	19.4158	0.8030
Our proposed loss with $\alpha = 0.1$ , $\beta = 0.9$ and $\eta = 0.5$	862.4890	18.7733	0.9135	<b>622.9391</b>	<b>20.1863</b>	0.8005
Our proposed loss with $\alpha = 0.5$ , $\beta = 0.5$ and $\eta = 0.5$	833.9951	18.9192	0.9146	703.8026	19.6563	0.8000

Method	CelebA dataset			SVHN dataset		
	MSE	PSNR	SSIM	MSE	PSNR	SSIM
(Yeh et al., 2017)	622.1092	20.1921	0.9087	1531.4601	16.2797	0.4791
(Yeh et al., 2017) adding gradient loss with $\alpha = 0.1$ , $\beta = 0.9$ and $\eta = 1.0$	584.3051	20.4644	0.9067	1413.7107	16.6272	0.4875
(Yeh et al., 2017) adding gradient loss with $\alpha = 0.5$ , $\beta = 0.5$ and $\eta = 1.0$	600.9579	20.3424	0.9080	1427.5251	16.5850	0.4889
(Yeh et al., 2017) adding gradient loss with $\alpha = 0.1$ , $\beta = 0.9$ and $\eta = 0.5$	580.8126	20.4904	0.9115	1446.3560	16.5281	<u>0.5120</u>
(Yeh et al., 2017) adding gradient loss with $\alpha = 0.5$ , $\beta = 0.5$ and $\eta = 0.5$	<u>563.4620</u>	<u>20.6222</u>	0.9103	<u>1329.8546</u>	<u>16.8928</u>	0.4974
Our proposed loss with $\alpha = 0.1$ , $\beta = 0.9$ and $\eta = 1.0$	424.7942	21.8490	0.9281	168.9121	25.8542	0.8960
Our proposed loss with $\alpha = 0.5$ , $\beta = 0.5$ and $\eta = 1.0$	380.4035	22.3284	0.9314	221.7906	24.6714	<b>0.9018</b>
Our proposed loss with $\alpha = 0.1$ , $\beta = 0.9$ and $\eta = 0.5$	<b>321.3023</b>	<b>23.0617</b>	<b>0.9341</b>	<b>154.5582</b>	<b>26.2399</b>	0.8969
Our proposed loss with $\alpha = 0.5$ , $\beta = 0.5$ and $\eta = 0.5$	411.8664	21.9832	0.9292	171.7974	25.7806	0.8939

thank you!