# AE-OT-GAN: Training GANs from data specific latent distribution

Dongsheng An[1], Yang Guo[1], Min Zhang[2], Xin Qi[1], Na Lei[3], Shing-Tung Yau[4], and Xianfeng Gu[1]

[1]Department of Computer Science, Stony Brook University
[2]Brigham and Women's Hospital, Harvard Medical School
[3]DLUT-RU, Dalian University of Technology
[4]Department of Mathematics, Harvard University

## Abstract

*Though generative adversarial networks (GANs) are prominent models to generate realistic and crisp images, they are unstable to train and suffer from the mode collapse/mixture. The problems of GANs come from approximating the intrinsic discontinuous distribution transform map with continuous DNNs. The recently proposed AE-OT model addresses the discontinuity problem by explicitly computing the discontinuous optimal transform map in the latent space of the autoencoder. Though have no mode collapse/mixture, the generated images by AE-OT are blurry. In this paper, we propose the AE-OT-GAN model to utilize the advantages of the both models: generate high quality images and at the same time overcome the mode collapse/mixture problems. Specifically, we firstly embed the low dimensional image manifold into the latent space by training an autoencoder (AE). Then the extended semi-discrete optimal transport (SDOT) map from the uniform distribution to the empirical latent distribution is used to generate new latent codes. Finally, our GAN model is trained to generate high quality images from the latent distribution induced by the extended SDOT map. The distribution transform map from this dataset related latent distribution to the data distribution will be continuous, and thus can be well approximated by the continuous DNNs. Additionally, the paired data between the latent codes and the real images gives us further restriction about the generator and stabilizes the training process. Experiments on simple MNIST dataset and complex datasets like CIFAR10 and CelebA show the advantages of the proposed method.*

## 1. Introduction

Image generation has been one of the core topics in the area of computer vision for a long time. Thanks to the quick development of deep learning, numerous generative models are proposed, including encoder-decoder based models [22,
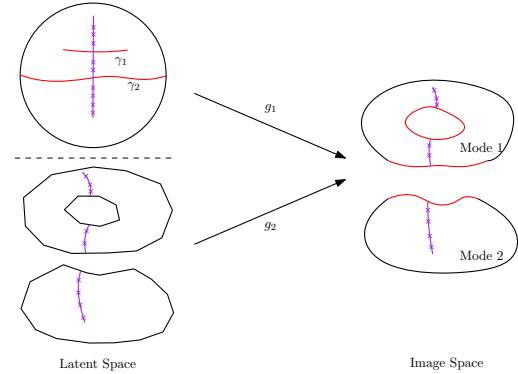


Figure 1. Distribution transport maps from two different latent distributions to the same data distribution. $g_1$ maps a *unimodal* latent distribution on the left to the data distribution on the right. Difference in the topology of their supporting manifolds will cause discontinuities of the map, which is hard to approximate by continuous neural networks. The singular set of $g_1$ consists of $\gamma_1$ and $\gamma_2$ (shown in red): continuous samplings of the source distribution are mapped to three disjoint segments (shown in purple). On the other hand, $g_2$ samples from a *suitably supported* latent distribution and is less likely to suffer from the discontinuity problem. Thus it can be well approximated by neural networks.

17, 2], generative adversarial networks (GANs) [16, 6, 42, 33, 5, 14], density estimator based models [39, 25, 24, 10] and energy based models [26, 47, 43, 11]. The encoder-decoder based models and GANs are the most prominent ones due to their capability to generate high quality images.

Intrinsically, the generator in a generative model aims to learn the real data distribution supported on the data manifold [37]. Suppose the distribution of a specific class of natural data $\nu_{gt}$ is concentrated on a low dimensional manifold $\chi$ embedded in the high dimensional data space. The encoder-decoder methods first attempt to embed the data into the latent space $\Omega$ through the encoder $f_\theta$, then samples from the latent distribution are mapped back to the manifold to generate new data by decoder $g_\xi$. While GANs, which have

no encoder, directly learn a map (generator) that transports a given prior low dimensional distribution to $\nu_{gt}$.

Usually, GANs are unstable to train and suffer from mode collapse [13, 30]. The difficulties come from the fact that the generator of a GAN model is trained to approximate the discontinuous distribution transport map from the *unimodal Gaussian distribution* to the *real data distribution* by the continuous neural networks [42, 2, 21]. In fact, when the supporting manifolds of the source and target distributions differ in topology or convexity, the OT map between them will be discontinuous [41], as illustrated in the map $g_1$ of Fig. 1. In practice, distribution transport maps can have complicated singularities, even when the ambient dimension is low (see e.g. [12]). This poses a great challenge for the generator training in standard GAN models.

To tackle the mode collapse and mode mixture problems caused by discontinuous transport maps, the authors of [2] proposed the AE-OT model. In this model, an autoencoder is used to map the images manifold $\chi$ into the latent manifold $\Omega$. Then, the semi-discrete optimal transport (SDOT) map $T$ from the uniform distribution $Uni([0, 1]^d)$ to the latent empirical distribution is explicitly computed via convex optimization approach. Then a piece-wise linear extension map of the SDOT, denoted by $\tilde{T}$, pushes forward the uniform distribution to a continuous latent distribution $\mu$, which in turn gives a good approximation of the latent distribution $\mu_{gt} = f_{\theta\#}\nu_{gt}$ ($f_{\theta\#}$ means the push forward map induced by $f_\theta$). Composing the continuous decoder $g_\xi$ and discontinuous $\tilde{T}$ together, i.e. $g_\xi \circ \tilde{T}(w)$, where $w$ is sampled from uniform distribution, this model can generate new images. Though have no mode collapse/mixture, the generated images look blurry. The framework of AE-OT is shown as follows:

$$(\nu_{gt}, \chi) \xrightarrow{f_\theta} (\mu_{gt}/\mu, \Omega) \xrightarrow{g_\xi} (\nu_{gt}, \chi)$$
$$\tilde{T} \uparrow$$
$$(Uni([0, 1]^d), \ [0, 1]^d)$$

In this work we propose the AE-OT-GAN framework to combine the advantages of the both models and generate high quality images without mode collapse/mixture. Specifically, after the training of the autoencoder and the computation of the extended SDOT map, we can directly sample from the latent distribution $\mu$ by applying $\tilde{T}(w)$ on the uniform distribution to train the GAN model. In contrast to the conventional GAN models, whose generators are trained to transport the latent Gaussian distribution to the data manifold distributions, our GAN model sample from the data inferred latent distribution $\mu$. The distribution transport map from $\mu$ to the data distribution $\nu_{gt}$ is continuous and thus can be well approximated by the generator (parameterized by CNNs), as shown in $g_2$ of Fig. 1. Moreover, the decoder

of the pre-trained autoencoder gives a warm start of the generator, so that the Kullback–Leibler divergence between real and fake batches of images have non-vanishing overlap in their supports during the training phase. Furthermore, the content loss and feature loss between paired latent codes and real input images regularize the adversarial loss and stabilize the GAN training. Experiments have shown efficacy and efficiency of our proposed model.

The contributions of the current work can be summarized as follows: **(1)** This paper proposes a novel AE-OT-GAN model that combines the strengths of AE-OT model and GAN model. It eliminates the mode collapse/mixture of GAN and removes the blurriness of the images generated by AE-OT. **(2)** The decoder of the autoencoder provides a good initialization of the generator of GAN. The number of iterations required to reach the equilibrium has been reduced by more than 100 times compared to typical GANs. **(3)** In addition to the adversarial loss, the explicit correspondence between the latent codes and the real images provide auxiliary constraints, namely the content loss, to the generator. **(4)** Our experiments demonstrate that our model can generate images consistently better than or comparable to the results of state-of-the-art methods.

## 2. Related Work

The proposed method in this paper is highly related to encoder-decoder based generation models, the generative adversarial networks (GANs), conditional GANs and the hybrid models that take the advantages of above.

**Encoder-decoder architecture** A breakthrough for image generating comes from the scheme of Variational Autoencoders (VAEs) (e.g. [22]), where the decoders approximate real data distributions from a Gaussian distribution in a variational approach (e.g [22] and [35]). Latter Yuri Burda et al. [45] lower the requirement of latent distribution and propose the importance weighted autoencoder (IWAE) model through a different lower bound. Bin and David [8] propose that the latent distribution of VAE may not be Gaussian and improve it by firstly training the original model and then generating new latent code through the extended ancestral process. Another improvement of the VAE is the VQ-VAE model [1], which requires the encoder to output discrete latent codes by vector quantisation, then the posterior collapse of VAEs can be overcome. By multi-scale hierarchical organization, this idea is further used to generate high quality images in VQ-VAE-2 [34]. In [17], the authors adopt the Wasserstein distance in the latent space to measure the distance between the distribution of the latent code and the given one and generate images with better quality. Different from the the VAEs, the AE-OT model [2] firstly embed the images into the latent space by autoencoder, then an extended semi-discrete OT map is computed to generate new latent code based on the fixed ones. Decoded by the decoder,
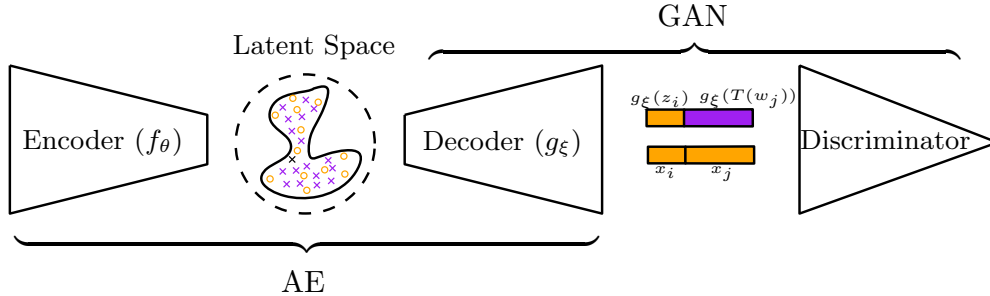
Figure 2. The framework of the proposed method. Firstly, the autoencoder is trained to embed the images into the latent space, the real latent codes are shown as the orange circles. Then we compute extended semi-discrete OT map $\tilde{T}$ to generate new latent codes in the latent space (the purple crosses). Finally, our GAN model is trained from the latent distribution induced by $\tilde{T}$ to the image distribution. Here the generator is just the decoder of the autoencoder. The fake batch (the bar with orange and purple colors) to train the discriminator is composed of two parts: the reconstructed images $g_\xi(z_i)$ of the real latent codes and the generated images $g_\xi(\tilde{T}(w))$ from the randomly generated latent codes with $w$ sampled from uniform distribution. The real batch (the bar with only orange color) is also composed of two parts: the real images $x_i$ corresponding to $z_i$, and the randomly selected images $x_j$.

new images can be generated. Although the encoder-decoder based methods are relatively simple to train, the generated images tend to be blurry.

**Generative adversarial networks** The GAN model [16] tries to alternatively update the generator, which maps the noise sampled from a given distribution to real images, and the discriminator differentiates the difference between the generated images and the real ones. If the generated images successfully fool the discriminator, we say the model is well trained. Later, [33] proposes a deep convolutions neural network (DCGAN) to generate images with better quality. While being a powerful tool in generating realistic samples, GANs can be hard to train and suffer from mode collapse problem [13]. After delicate analysis, [5] points out that it is the KL divergence the original GAN used causes these problems. Then the authors introduce the celebrated WGAN, which makes the whole framework easy to converge. To satisfy the lipschitz continuity required by WGAN, a lot of methods are proposed, including clipping [5], gradient penalty [14], spectral normalization [32] and so on. Later, Wu et al. [20] use the wasserstein divergence objective, which get rid of the lipschitz approximation problem and get a better result. Instead $L_1$ cost adopted by WGAN, Liu et.al [29] propose the WGAN-QC by taking the $L_2$ cost into consideration. Though various GANs can generate sharp images, they will theoretically encounter the mode collapse or mode mixture problem [13, 2].

**Hybrid models** To solve the blurry image problem of encoder-decoder architecture and the mode collapse/mixture problems of GANs, a natural idea is to compose them together. Larsen et al. [3] propose to combine the variational autoencoder with a generative adversarial network, and thus generate images better than VAEs. [31] matches the aggregated posterior of the hidden code vector of the autoencoder with an arbitrary prior distribution by a discriminator and then applies the model into tasks like semi-supervised classification and dimensionality reduction. BiGAN [19], with

the same architecture with ours, uses the discriminator to differentiate both the generated images and the generated latent code. Further, by utilizing the BigGAN generator [4], the BigBiGAN [9] extends this method to generate much better results. Here we also treat the BourGAN [42] as a hybrid model, because it firstly embeds the images into latent space by Bourgain theorem, then trains the GAN model by sampling from the latent space using the GMM model.

Conditional GANs are another kind of hybrid models that can also be treated as image-to-image transformation. For example, using an encoder-decoder architecture to build the connection between paired images and then differentiating the decoded images with the real ones by a discriminator, [18] is able to transform images of different styles. Further, SRGAN [7] uses similar architecture to get super resolution images from their low resolution versions. The SRGAN model is the most similar work to ours, as it also utilizes the content loss and adversarial loss. The main differences between this model and ours including: (i) SRGAN just uses the paired data, while the proposed method use both the paired data and generated new latent code to train the model; (ii) the visually meaningful features used by SRGAN are extracted from the pre-trained VGG19 network [36], while in our model, they come from the encoder itself. This makes them more reasonable especially under the scenes where the datasets are not included in those used to train the VGG.

## 3. The Proposed Method

In this section, we explain our proposed AE-OT-GAN model in detail. There are mainly three modules, an autoencoder (AE), an optimal transport mapper (OT) and a GAN model. Firstly, an AE model is trained to embed the data manifold $\chi$ into the latent space. At the same time, the encoder $f_\theta$ pushes forward the ground-truth data distribution $\nu_{gt}$ supported on $\chi$ to the ground-truth latent distribution $\mu_{gt}$ supported on $\Omega$ in the latent space. Secondly, we compute the semi-discrete OT map from the uniform distribution to

the empirical latent distribution. By extending the SDOT map, we can construct the continuous distribution $\mu$ that approximates the ground-truth latent distribution $\mu_{gt}$ well. Finally, starting from $\mu$ as the latent distribution, our GAN model is trained to generate both realistic and crisp images. The pipeline of our proposed model is illustrated in Fig. 2. In the following, we will explain the three modules one by one.

### 3.1. Data Embedding with Autoencoder

We model the real data distribution as a probability measure $\nu_{gt}$ supported on an $r$ dimensional manifold $\chi$ embedded in the $D$ dimensional Euclidean space $\mathbb{R}^D$ (ambient space) with $r \ll D$.

In the first step of our AE-OT-GAN model, we train an autoencoder (AE) to embed the real data manifold $\chi$ to be the latent manifold $\Omega$. In particular, training the AE model is equivalent to compute the encoding map $f_\theta$ and decoding map $g_\xi$

$$(\nu_{gt}, \chi) \xrightarrow{f_\theta} (\mu_{gt}, \Omega) \xrightarrow{g_\xi} (\nu_{gt}, \chi)$$

by minimizing the loss function:

$$\mathcal{L}(\theta, \xi) := \sum_{i=1}^{n} \|x_i - g_\xi \circ f_\theta(x_i)\|^2,$$

with $f_\theta$ and $g_\xi$ parameterized by standard CNNs ($\theta$ and $\xi$ are the parameters of the networks, respectively). Given densely sampling from the image manifold (detailed explanation is included in the supplementary) and ideal optimization (namely the loss function goes to 0), $f_\theta \circ g_\xi$ coincides with the identity map. After training, $f_\theta$ is a continuous, convertible map, namely a *homeomorphism*, and $g_\xi$ is the inverse homeomorphism. This means $f_\theta : \chi \to \Omega$ is an embedding, and pushes forward $\nu_{gt}$ to the latent data distribution $\mu_{gt} := f_{\theta\#}\nu_{gt}$. In practice, we only have the empirical data distribution given by $\hat{\nu}_{gt} = \frac{1}{n}\sum_{i=1}^{n}\delta(x - x_i)$, which is push forward to be the empirical latent distribution $\hat{\mu}_{gt} = \frac{1}{n}\sum_{i=1}^{n}\delta(z - z_i)$, where $n$ is the number of samples.

### 3.2. Constructing $\mu$ with Semi-Discrete OT Map

In this section, from the empirical latent distribution $\hat{\mu}_{gt}$, we construct a continuous latent distribution $\mu$ following [2] such that (i) it generalizes $\hat{\mu}_{gt}$ well, so that all of the modes are covered by the support of $\mu$ (ii) the support of $\mu$ has similar topology to that of $\mu_{gt}$, which ensures that the transport map from $\mu$ to $\nu_{gt}$ is continuous and (iii) it is efficient to sample from $\mu$.

To obtain $\mu$, the semi-discrete OT map $T$ from the uniform distribution $Uni([0, 1]^d)$ to $\hat{\mu}_{gt}$ is firstly computed. Here $d$ is the dimension of the latent space. By extending $T$ to be a piece-wise linear map $\tilde{T}$, we can construct $\mu$ as the
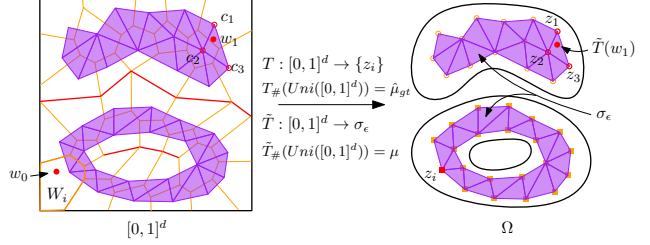


Figure 3. OT map $T$ and the extended OT map $\tilde{T}$ in 2D case. Here $T$ maps points (e.g. $w_0$) in each polyhedral cell $W_i$ (orange cells on the left) to the corresponding latent code $z_i$ (circles and squares on the right). The piece-wise linear $\tilde{T}$ maps triangulated regions in $[0, 1]^d$ to the simplicial complex $\sigma_\varepsilon$ in the latent space (shown in purple). Given the barycenters $c_i$'s of each $W_i$'s, each triangle $\Delta c_i c_j c_k$ is mapped to the corresponding simplex $[z_i, z_j, z_k]$. For example, $w_1$ in the triangle $\Delta c_0 c_1 c_2$ is mapped to $\tilde{T}(w_1)$ in the simplex $[z_0, z_1, z_2]$. Red lines in $[0, 1]^d$ illustrate the singular set of $T$, which corresponds to the pre-image of gaps or holes in $\Omega$.

push forward distribution of $Uni([0, 1]^d)$ under $\tilde{T}$:

$$(Uni([0, 1]^d), [0, 1]^d) \xrightarrow{\tilde{T}} (\mu, \Omega)$$

In the first step, we compute the semi-discrete OT map $T : [0, 1]^d \to \Omega$, with $T_\# Uni([0, 1]^d) = \hat{\mu}_{gt}$. Under $T$, the continuous domain of $[0, 1]^d$ is decomposed into cells $\{W_i\}$ with $T(w) = z_i$, $\forall w \in W_i$, with the Lebesgue measure of each $W_i$ to be $\frac{1}{n}$. The cell structure is shown in the left frame of Fig. 3 (the orange cells). Computational details of $T$ can be found in the supplementary material and [2].

Secondly, we extend the image domain of $T$ from the discrete latent codes $\{z_i\}$ to a continuous neighborhood $\sigma_\varepsilon$, which serves as the supporting manifold of $\mu$. Specifically, we construct a simplicial complex $\sigma_\varepsilon$ from the latent codes $\{z_i\}$. Here $\varepsilon > 0$ is a constant. The 0-skeleton of $\sigma_\varepsilon$, represented by $\sigma_\varepsilon^{(0)}$, is the set of all latent codes $\{z_i\}$. The we define its k-skeletons $\sigma_\varepsilon^{(k)}$ by $\sigma_\varepsilon^{(k)} = \{[z_1, z_2, \ldots, z_k] \mid \|z_i - z_j\|_2 \le \varepsilon, \forall 1 \le i < j \le k\}$ for $0 < k \le d$. The right frame of Fig. 3 shows an example of $\sigma_\varepsilon$. By assuming that the latent code is densely sampled from the latent manifold $\Omega$ and with an appropriate $\varepsilon$, $\sigma_\varepsilon$ will have consistent "hole" and "gap" structure with $\Omega$, in the sense of homology equivalence. Details are described in the supplementary material.

Finally, we define the piece-wise linear extended OT map $\tilde{T} : [0, 1]^d \to \sigma_\varepsilon$. Given a random sample $w$ sampled from $Uni([0, 1]^d)$, we can find the cell $W_i$ containing it. By computing the barycentric parameters $\lambda_j$'s with respect to the nearby mass centers $c_j$'s of the cells $W_j$'s, i.e. compute $\lambda_j$'s such that $w = \Sigma \lambda_j c_j$ with $0 \le \lambda_j \le 1$ and $\Sigma \lambda_j = 1$. Here $W_j$ represents the neighbour of $W_i$. Then $w$ is mapped to $\tilde{T}(w) := \Sigma \lambda_j T(c_j) = \Sigma \lambda_j z_j$ if the corresponding $z_j$'s form a simplex of $\sigma_\varepsilon$. Otherwise we map $w$ to $z_i$, i.e. $\tilde{T}(w) := T(c_i) = z_i$. As illustrated in Fig. 3, compared to the many-to-one semi-discrete OT map $T$, $\tilde{T}$

maps samples within the triangular areas (the purple triangles on the left frame) in $[0,1]^d$ *linearly* to the corresponding simplices in $\sigma_\varepsilon$ (the purple triangles on the right frame) in a bijective manner. We denote the pushed forward distribution under $\tilde{T}$ as $\mu_\varepsilon := \tilde{T}_\# Uni([0,1]^d)$.

**Theorem 1.** *The 2-Wasserstein distance between $\mu_\varepsilon$ and $\hat{\mu}_{gt}$ satisfies $W_2(\mu_\varepsilon, \hat{\mu}_{gt}) \leq \varepsilon$. Moreover, if the latent codes are densely sampled from the latent manifold $\Omega$, we have $W_2(\mu, \mu_{gt}) \leq 2\varepsilon$, $\mu$-almost surely.*

To avoid confusion, we omit the subscript $\varepsilon$ and denote $\mu_\varepsilon$ as $\mu$. With proof included in the supplementary material, this theorem tells us that as a continuous generalization of $\hat{\mu}_{gt}$, $\mu$ is a good approximation of $\mu_{gt}$. Also, we want to mention that $\tilde{T}$ is a piece-wise linear map that pushes forward $Uni([0,1]^d)$ to $\mu$, which makes the sampling from $\mu$ efficient and accurate.

### 3.3. GAN Training from $\mu$

The GAN model computes the transport map from the continuous latent distribution $\mu$ to the data distribution on the manifold.

$$(\mu, \Omega) \xrightarrow{g_\xi} (\nu_{gt}, \chi).$$

Our GAN model is based on the vanilla GAN model proposed by Ian Goodfellow et.al [16]. The generator $g_\xi$ is used to generate new images by sampling from the latent distributin $\mu$, while the discriminator $d_\eta$ is used to discriminate if the distribution of the generated images are the same with that of the real images. The training process is formalized to be a min-max optimization problem:

$$\min_\xi \max_\eta \mathcal{L}(\xi, \eta),$$

where the loss function is given by

$$\mathcal{L}(\xi, \eta) = \mathcal{L}_{adv} + \mathcal{L}_{feat} + \beta \mathcal{L}_{img} \qquad (1)$$

In our model, the loss function consists of three terms, the image content loss $\mathcal{L}_{img}$, the feature loss $\mathcal{L}_{feat}$ and the adversarial loss $\mathcal{L}_{adv}$. Here $\beta > 0$ is the weight of the content loss.

**Adversarial Loss** We adopt the vanilla GAN model [16] based on the Kullback–Leibler (KL) divergence. The key difference between our model and the original GAN is that our latent samples are drawn from the data related latent distribution $\mu$, instead of a Gaussian distribution. The adversarial loss is given by:

$$\mathcal{L}_{adv} = \min_\xi \max_\zeta E_{x \sim \nu_{gt}}[log\ d_\zeta(x)]$$
$$+ E_{z \sim \mu}[log(1 - d_\zeta(g_\xi(z)))]$$

According to [5], vanilla GAN is hard to converge because the supports of the distributions of real images and fake images may not intersect each other, which makes the KL divergence between them infinity. This issue is solved in our case, because (1) the training of AE gives a warm start to the generator, so at the beginning of the training, the generated distribution $g_{\xi\#}(\mu)$ is close to the real data distribution $\nu_{gt}$. (2) by delicate settings of the fake and real batches used to train the discriminator, we can keep the KL divergence between them converge well. In detail, as shown in Fig. 2, the fake batch is composed of both the reconstructed images from the real latent code (the orange circles) and the generated images from the generated latent code (the purple crosses), and the real batch includes both the real images corresponding to the real latent code and some randomly selected images.

**Content Loss** Recall that the generator can produce two types of images: images reconstructed by real latent codes and images from generated latent codes. Given a real sample $x_i$, its latent code is $z_i = f_\theta(x_i)$, the reconstructed image is $g_\xi(z_i)$. Each reconstructed image is represented as a triple $(x_i, z_i, g_\xi(z_i))$. Suppose there are $n$ reconstructed images in total, the content loss is given by

$$\mathcal{L}_{img} = \frac{1}{n} \sum_{i=1}^n \|g_\xi(z_i) - x_i\|_2^2 \qquad (2)$$

Where $g_\xi$ is the generator parameterized by $\xi$.

**Feature Loss** We adopt the feature loss similar to that in [7]. Given a reconstructed image triple $(x_i, z_i, g_\xi(z_i))$, we encode $g_\xi(z_i)$ by the encoder of AE. Ideally, the real image $x_i$ and the generated image $g_\xi(z_i)$ should be same, therefore their latent codes should be similar. We measure the difference between their latent codes by the feature loss. Furthermore, we can measure the difference between their intermediate features from different layers of the encoder.

Suppose the encoder is a network with $L$ layers, the output of the $l$th layer is denoted as $f_\theta^{(l)}$. The feature loss is given by

$$\mathcal{L}_{feat} := \frac{1}{n} \sum_{i=1}^n \sum_{l=1}^L \alpha^{(l)} \|f_\theta^{(l)}(x_i) - f_\theta^{(l)} \circ g_\xi(z_i)\|_2^2,$$

Where $\alpha^{(l)}$ is the weight of the feature loss of the $l$-th layer.

For reconstructed images $(x_i, z_i, g_\xi(z_i))$, the content loss and the feature loss force the generated image $g_\xi(z_i)$ to be the same with the real image $x_i$, therefore the manifold $g_\xi(\Omega)$ align well with the real data manifold $\chi$.

## 4. Expriments

To evaluate the proposed method, several experiments are conducted on simple dataset MNIST [27] and complex datasets including Cifar10 [23], CelebA [46] and CelebA-HQ [28].
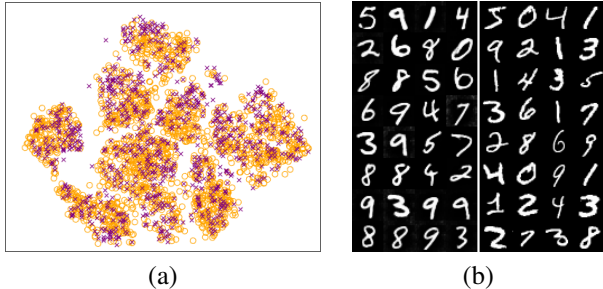
(a)                          (b)

Figure 4. (a) Latent code distribution. The orange circles represent the fixed latent code and the purple crosses are the generated ones. (b) Comparison between the generated digits (left) and the real digits (right).



(a)                          (b)

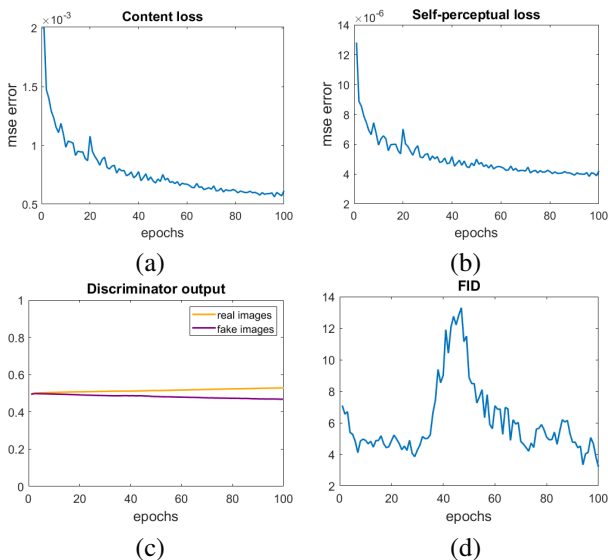(c)                          (d)

Figure 5. The curves for training on MNIST dataset [27] of each epoch, including the results of content loss (a) and self-perceptual loss (b), the discriminator output (c) and FIDs (d).

**Architecture** We adopt the InfoGAN [6] architecture as our GAN model to train the MNIST dataset. The standard and ResNet models used to train the Cifar10 dataset are the same with those used by SNGAN [32], and the architectures of WGAN-div [20] are used to train the CelebA dataset. The framework of encoder is set to be the mirror of the generators/decoders.

**Evaluation metrics** To illustrate the performance of the proposed method, we adopt the commonly used Frechet Inception distance (FID) [15] as our evaluation metrics. FID takes both the generated images and the real images into consideration. When the images are embedded into the feature space by inception network, two high dimensional Gaussian distributions are used to approximate the empirical distributions of the generated and real features, respectively. Finally, the FID is given by the difference between the two Gaussian distributions. Lower FID means better quality of the generated dataset. This metric has been proven to be effective in judging the performance of the generated models, and it serves as a standard for comparison with other works.

**Training details** To get rid of the vanishing gradient problem and make the model converge better, we use the following three strategies:

*(i) Train the discriminator using Batch Composition* There are two types of latent codes in our method: *the real latent codes* coming from encoding the real images by the encoder, and generated latent codes coming from the extended OT map. Correspondingly, there are two types of generated images, *the reconstructed images* from the real latent codes and *the generated images* from the generated latent codes.

To train the discriminator, both the fake batch and real batch are used. *The fake batch* consists of both randomly selected reconstructed images and generated images, and *the real batch* only includes real images, in which the first part has a one-to-one correspondence with the reconstructed images in the fake batch, as shown in Fig. 2. In all the experiments, the ratio between the number of generated images and reconstructed images in the fake batch is 3.

This strategy ensures that there is an overlap between the supports of the fake and real batches, so that the KL divergence is not infinity.

*(ii) Different learning rate* For better training, we use different learning rates for the generator and the discriminator as suggested by Heusel et al. in [15]. Specifically, we set the learning rate of the generator to be $lr_G = 2e - 5$ and that of the discriminator to be $lr_D = lr_G/R$, where $R > 1$. This improves the stability of the training process.

*(iii) Different inner steps* Another way to improve the training consistency of the whole framework is to set different update steps for the generator and discriminator. Namely, When the discriminator updated once, the generator updated $T$ times correspondingly. This strategy is the opposite of training vanilla GANs, which typically require multiple discriminator update steps per generator update step.

By setting $R$ and $T$, we can keep the discriminator output of the real images is slightly large than that of the generated images, which can better guide the training of the generator. For the MNIST dataset, $R = 15$ and $T = 3$; for the Cifar10 dataset, $R = 25$ and $T = 10$; and for the CelebA dataset, $R = 15$ and $T = 5$. In Eq. 1, $\beta = 2000$ and $\alpha^{(l)} = 0.06$ with $l < L$, where $L$ denotes the last layer of the encoder. $\alpha^L = 2.0/\|Z\|_2$ is used to regularize the loss of the latent codes.

With the above settings and the warm initialization of the generator from the pre-trained decoder, for each dataset, the total epochs for training is set to be 500, which is far less than the training of GANs (usually 10k~50k).

## 4.1. Convergence Analysis in MNIST

In this experiment, we evaluate the performance of our proposed model on MNIST dataset [27], which can be well embedded into the 64 dimensional latent space with the architecture of InfoGAN [6]. In Fig. 4(a), we visualize

(a) Epoch 0 (AE-OT)   (b) Epoch 80   (c) Epoch 160   (d) Epoch 240   (e) Ground-truth
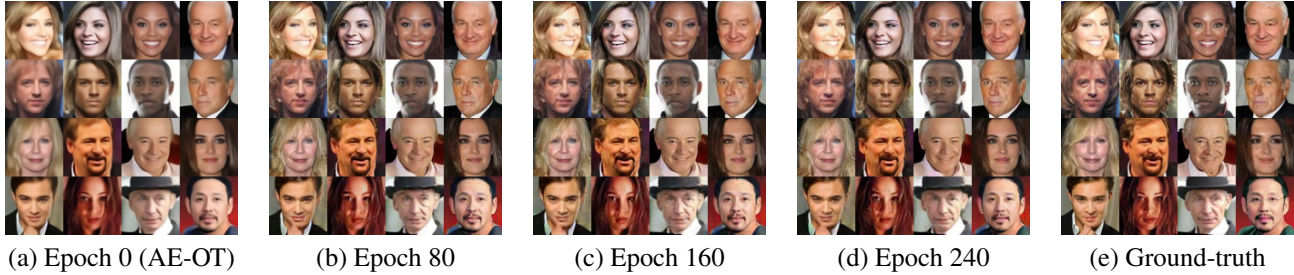
Figure 6. Evolution of the generator during training on the CelebA dataset [46]. Reconstructed images from real latent codes at different epochs are shown.



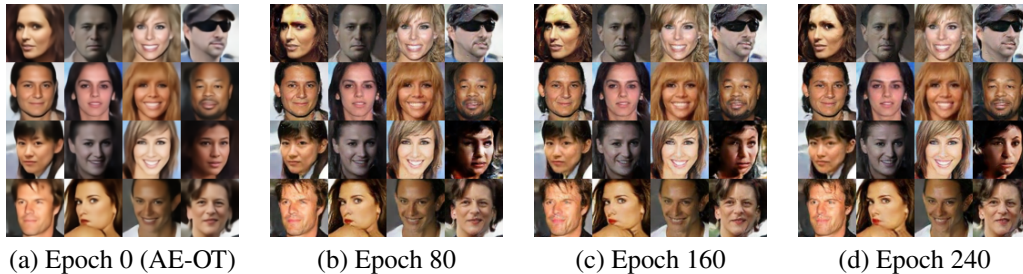(a) Epoch 0 (AE-OT)   (b) Epoch 80   (c) Epoch 160   (d) Epoch 240

Figure 7. Evolution of the generator during training on the CelebA dataset [46]. Generated images from generated latent codes at different epochs are shown.

| | CIFAR10 | | CelebA | |
|---|---|---|---|---|
| | Standard | ResNet | Standard | ResNet |
| WGAN-GP [14] | 40.2 | 19.6 | 21.2 | 18.4 |
| PGGAN [38] | - | 18.8 | - | 16.3 |
| SNGAN [32] | 25.5 | 21.7 | - | - |
| WGAN-div [20] | - | 18.1 | 17.5 | 15.2 |
| WGAN-QC [29] | - | - | - | 12.9 |
| AE-OT [2] | 34.2 | 28.5 | 24.3 | 28.6 |
| AE-OT-GAN | 25.2 | 17.1 | 11.2 | 7.8 |

Table 1. The comparison of FID between the proposed method and the state of the arts on Cifar10 and CelebA.

the real latent code (brown circles) and the generated latent codes (purple crosses) by t-SNE [40]. It is obvious that the support of the real latent distribution and that of the generated distribution align well. Frame (b) of Fig. 4 shows the comparison between the generated handwritten digits (left) and the real digits (right), which is very difficult for humans to distinguish.

To show the convergent property of the proposed method, we plot the related curves in Fig. 5. The frame (a) and (b) show the changes of the content loss about the images and latent codes, and both of them decrease monotonously. The frame (c) shows that the output of the discriminator for real images is only slightly larger than that for the fake images during the training process, which can help the generator generate more realistic digits. The frame (d) shows the evolution of FID and the final value is 3.2. For MNIST dataset, the best known FIDs with the same InfoGAN architecture are 6.7 and 6.4, reported in [30] and [2] respectively. This shows our model outperforms state-of-the-art.

## 4.2. Quality Evaluation on Complex Dataset

In this section, we compare with the state-of-the-art methods quantitatively and qualitatively.

**Progressive Quality Improvement** Firstly, we show the evolution results of the proposed method in Fig. 6 and Fig. 7 during GAN's training process. Quality of the generated images increases monotonously during the process. Images in first four frames of Fig. 6 illustrates the results reconstructed from the real latent codes by the decoder, with the last frame showing the corresponding ground-truth input images. By examining the frames carefully, it is obvious that as the increase of the epochs, the generated images become sharper and sharper, and eventually they are very close to the ground truth. Fig. 7 shows the generated images from some generated latent codes (therefore, no corresponding real images). Similarly. the images become sharper as the increase of epochs. Here we need to state that the 0 epoch stage means the images are generated by the original decoder, which are equivalent to the outputs of an AE-OT model [2]. Thus we can conclude that the proposed AE-OT-GAN improves the performance of AE-OT prominently.

| CT-GAN [44] | WGAN-GP [14] | WGAN-div [20] | WGAN-QC [29] | Proposed method |

Figure 8. The visual comparison between the proposed method and the state-of-the-arts on CelebA dataset [46] with ResNet architecture.



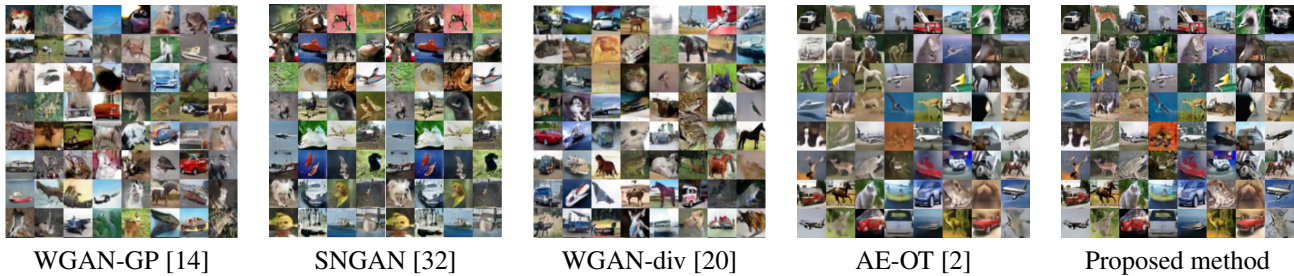| WGAN-GP [14] | SNGAN [32] | WGAN-div [20] | AE-OT [2] | Proposed method |

Figure 9. The visual comparison between the proposed method and the state-of-the-arts on Cifar10 dataset [23] with ResNet architecture.



Figure 10. The generation results of CelebA-HQ by the proposed method.

| PGGAN | WGAN-div | WGAN-QC | AE-OT-GAN |
|-------|----------|---------|-----------|
| 14.7  | 13.5     | 7.7     | 7.2       |

Table 2. The FIDs of the proposed method and the state-of-the-arts.

**Comparison on CelebA and CIFAR 10** Secondly, we compared with the state-of-the-arts including WGAN-GP [14], PGGAN [38], SNGAN [32], CTGAN [44], WGAN-div [20], WGAN-QC [29] and the recently proposed AE-OT model [2] on Cifar10 [23] and CelebA [46]. Tab. 1 shows the FIDs of the our method and the comparisons trained under both the standard and ResNet architectures. The FID of other methods come from the listed papers except those of the AE-OT, which are directly computed by our model (the results of epoch 0). From the table we can see that our method gets much better results than others on the CelebA dataset, both under the standard and the ResNet architecture. Also, the generated faces of the proposed method have less flaws compared to other GANs, as shown on Fig. 8. On Cifar10, the FIDs of our model are also comparable to the state-of-the-arts. And we also show some generated images on Fig. 9. The convergence curves for the both datasets can be found in the supplementary.

**Experiment on CelebA-HQ** Furthermore, We also test the proposed method on images with high resolution, namely the CelebA-HQ dataset with image size to be 256x256. The architecture used to train the model is illustrated in the supplementary. The parameters in our model is far less than that of [29, 20, 38], while the performance is better than theirs, as shown in Tab. 2. We also display several images generated in Fig. 10, which are crisp and visually realistic.

## 5. Conclusion and Future Work

In this paper, we propose the AE-OT-GAN model which composes the AE-OT model and vanilla GAN together. By utilizing the merits of the both models, our method can generate high quality images without mode collapse nor mode mixture. Firstly, the images are embedded into the latent space by autoencoder, then the SDOT map from uniform dis-

8

tribution to the empirical distribution supported on the latent code is computed. Sampling from the latent distribution by applying the extended SDOT map, we can train our GAN model. Moreover, the paired latent code and images give us additional constraints about the generator. Using the FID as metric, we show that the proposed model is able to generate images comparable or better than the state of the arts.

# References

[1] Koray Kavukcuoglu Aaron van den Oord, Oriol Vinyals. Neural discrete representation learning. In *NeurIPS*, 2017.

[2] Dongsheng An, Yang Guo, Na Lei, Zhongxuan Luo, Shing-Tung Yau, and Xianfeng Gu. Ae-ot: A new generative model based on extended semi-discrete optimal transport. In *International Conference on Learning Representations*, 2020.

[3] Hugo Larochelle Ole Winther Anders Boesen Lindbo Larsen, Søren Kaae Sønderby. Autoencoding beyond pixels using a learned similarity metric. 2016.

[4] Karen Simonyan Andrew Brock, Jeff Donahue. Large scale gan training for high fidelity natural image synthesis. In *International Conference on Learning Representations*, 2019.

[5] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *ICML*, pages 214–223, 2017.

[6] Xi Chen, Yan Duan, Rein Houthooft, John Schulman, Ilya Sutskever, and Pieter Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. In *Advances in Neural Information Processing Systems*, 2016.

[7] Ferenc Huszar Jose Caballero Andrew Cunningham Alejandro Acosta-Andrew Aitken Alykhan Tejani Johannes Totz Zehan Wang Wenzhe Shi Christian Ledig, Lucas Theis. Photo-realistic single image super-resolution using a generative adversarial network. 2017.

[8] Bin Dai and David Wipf. Diagnosing and enhancing VAE models. In *International Conference on Learning Representations*, 2019.

[9] Jeff Donahue and Karen Simonyan. Large scale adversarial representation learning. In *https://arxiv.org/abs/1907.02544*, 2019.

[10] Prafulla Dhariwal Durk P Kingma. Glow: Generative flow with invertible 1x1 convolutions. In *NeurIPS*, 2018.

[11] Song-Chun Zhu Ying Nian Wu Erik Nijkamp, Mitch Hill. On learning non-convergent non-persistent short-run mcmc toward energy-based model. *arXiv preprint arXiv:1904.09770*, 2019.

[12] Alessio Figalli. Regularity properties of optimal maps between nonconvex domains in the plane. *Communications in Partial Differential Equations*, 35(3):465–479, 2010.

[13] Ian Goodfellow. Nips 2016 tutorial: Generative adversarial networks. *arXiv preprint arXiv:1701.00160*, 2016.

[14] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of wasserstein gans. In *NIPS*, pages 5769–5779, 2017.

[15] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, Günter Klambauer, and Sepp Hochreiter. Gans

trained by a two time-scale update rule converge to a nash equilibrium. 2017.

[16] Mehdi Mirza Bing Xu David Warde-Farley Sherjil Ozair Aaron Courville Yoshua Bengio Ian J. Goodfellow, Jean Pouget-Abadie. Generative adversarial nets. 2014.

[17] Sylvain Gelly Bernhard Schoelkopf Ilya Tolstikhin, Olivier Bousquet. Wasserstein auto-encoders. In *ICLR*, 2018.

[18] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2017.

[19] Trevor Darrell Jeff Donahue, Philipp Krähenbühl. Adversarial feature learning. In *International Conference on Learning Representations*, 2017.

[20] Janine Thoma Dinesh Acharya Luc Van Gool Jiqing Wu, Zhiwu Huang. Wasserstein divergence for gans. In *ECCV*, 2018.

[21] Mahyar Khayatkhoei, Maneesh K. Singh, and Ahmed Elgammal. Disconnected manifold learning for generative adversarial networks. In *Advances in Neural Information Processing Systems*, 2018.

[22] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

[23] Alex Krizhevsky. Learning multiple layers of features from tiny images. *Tech report*, 2009.

[24] Jascha Sohl-Dickstein Laurent Dinh and Samy Bengio. Density estimation using real nvp. In *ICLR*, 2017.

[25] Yoshua Bengio Laurent Dinh, David Krueger. Nice: Nonlinear independent components estimation. *arXiv preprint arXiv:1410.8516*, 2014.

[26] Yann Lecun, Sumit Chopra, and Raia Hadsell. *A tutorial on energy-based learning*. 01 2006.

[27] Yann LeCun and Corinna Cortes. MNIST handwritten digit database. 2010.

[28] Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. Maskgan: Towards diverse and interactive facial image manipulation. *arXiv preprint arXiv:1907.11922*, 2019.

[29] Huidong Liu, Xianfeng Gu, and Dimitris Samaras. Wasserstein gan with quadratic transport cost. In *ICCV*, 2019.

[30] Mario Lucic, Karol Kurach, Marcin Michalski, Sylvain Gelly, and Olivier Bousquet. Are gans created equal? a large-scale study. In *Advances in neural information processing systems*, pages 698–707, 2018.

[31] Alireza Makhzani, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. Adversarial autoencoders. *arXiv preprint arXiv:1511.05644*, 2015.

[32] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. In *ICLR*, 2018.

[33] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. In *ICLR*, 2016.

[34] Ali Razavi, Aaron Oord, and Oriol Vinyals. Generating diverse high-fidelity images with vq-vae-2, 2019.

[35] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. *arXiv preprint arXiv:1401.4082*, 2014.

[36] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. 2014.

[37] J B Tenenbaum, V Silva, and J C Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290(5500):2391–232, 2000.

[38] Samuli Laine Jaakko Lehtinen Tero Karras, Timo Aila. Progressive growing of gans for improved quality, stability, and variation. In *ICLR*, 2018.

[39] Aaron van den Oord, Nal Kalchbrenner, Lasse Espeholt, koray kavukcuoglu, Oriol Vinyals, and Alex Graves. Conditional image generation with pixelcnn decoders. In *Advances in Neural Information Processing Systems*. 2016.

[40] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-SNE. *Journal of Machine Learning Research*, 2008.

[41] Cédric Villani. *Optimal transport: old and new*, volume 338. Springer Science & Business Media, 2008.

[42] Chang Xiao, Peilin Zhong, and Changxi Zheng. Bourgan: Generative networks with metric embeddings. In *NeurIPS*, 2018.

[43] Jianwen Xie, Yang Lu, Song Zhu, and Yingnian Wu. Cooperative training of descriptor and generator networks. *IEEE transactions on pattern analysis and machine intelligence*, 2016.

[44] Lei Xu, Maria Skoularidou, Alfredo Cuesta-Infante, and Kalyan Veeramachaneni. Modeling tabular data using conditional gan. In *Advances in Neural Information Processing Systems*, 2019.

[45] Ruslan Salakhutdinov Yuri Burda, Roger Grosse. Importance weighted autoencoders. In *ICML*, 2015.

[46] Zhanpeng Zhang, Ping Luo, Chen Change Loy, and Xiaoou Tang. From facial expression recognition to interpersonal relation prediction. *International Journal of Computer Vision*, 2018.

[47] Song Zhu, Yingnian Wu, and David Mumford. Filters, random fields and maximum entropy (frame): Towards a unified theory for texture modeling. *International Journal of Computer Vision*, 1998.