

# Generative models with kernel distance in data space

Szymon Knop, Marcin Mazur, Przemysław Spurek, Jacek Tabor, Igor Podolak

Faculty of Mathematics and Computer Science  
Jagiellonian University, Krakw, Poland

## Abstract

Generative models dealing with modeling a joint data distribution are generally either autoencoder or GAN based. Both have their pros and cons, generating blurry images or being unstable in training or prone to mode collapse phenomenon, respectively. The objective of this paper is to construct a model situated between above architectures, one that does not inherit their main weaknesses. The proposed LCW generator (Latent Cramer-Wold generator) resembles a classical GAN in transforming Gaussian noise into data space. What is of utmost importance, instead of a discriminator, LCW generator uses kernel distance. No adversarial training is utilized, hence the name generator. It is trained in two phases. First, an autoencoder based architecture, using kernel measures, is built to model a manifold of data. We propose a Latent Trick mapping a Gaussian to latent in order to get the final model. This results in very competitive FID values.

## Introduction

*Generative modeling* is a fast-growing area of machine learning which deals with modeling a joint distribution of data. Generative modeling's key task is to train a generator network to produce new samples from a given data space by transforming samples from noise distribution. Training usually minimizes the dissimilarity between real and generated data distributions.

Widely used approaches to generative modeling are GANs and models based on autoencoder architecture. Autoencoder consists of an encoder  $\mathcal{E}: \mathcal{X} \rightarrow \mathcal{Z}$  and a decoder  $\mathcal{D}: \mathcal{Z} \rightarrow \mathcal{X}$  (a generator) networks. Training boils down to minimizing a tuned sum of a reconstruction error and some measure of similarity between the distribution of encoded data  $P_{\mathcal{E}(X)}$  and a given prior (noise) distribution  $P_{\mathcal{Z}}$  on the latent  $\mathcal{Z}$ . GAN consists of a generator  $\mathcal{G}: \mathcal{Z} \rightarrow \mathcal{X}$  and a discriminator that distinguishes between samples from real  $P_X$  and "fake"  $P_{\mathcal{G}(Z)}$  data distributions. It learns adversarially by utilizing a *minimax* game rule.

Both approaches have their pros and cons. Autoencoder based generative methods produce theoretically elegant generative models with the drawback that they tend to generate blurry samples when applied to natural images. On the other hand, their main advantage over GAN based models is that they allow fitting manifold of data and approximate probability distribution simultaneously (Goodfellow et al. 2014).

Contrary to GANs, in autoencoder based models, each data point can easily be directly mapped into latent space, needed in conditional data generation.

GAN based methods' main advantage is their ability to produce sharp images, almost indistinguishable from real ones. On the other hand, GANs are harder to train, unstable and may suffer from the mode collapse problem, where the resulting model is unable to capture all the true data distribution variability. GAN based models imitate real data well but frequently do not cover the training data set's entire space (Heusel et al. 2017).

The objective of this paper is to show that it is achievable to effectively train a model that does not inherit the weaknesses of the above models, such as blurry images or complex adversarial training, and provides for better results. The proposed LCW generator<sup>1</sup> is obtained in a two-stage training and uses kernel measures in all cost functions. The discriminator network is no longer necessary.

Our contributions can be summarized as follows:

- we introduce a new LCW generator, which resembles a classical GAN in transforming Gaussian noise into data space and use kernel distance instead of adversarial training,
- we show that kernel based metric can be used as reconstruction error in classical autoencoder based generative models by introducing a CW<sup>2</sup> model that achieves state-of-the-art FID scores in it's category,
- we propose a Latent Trick, which can be applied to any model with latent space to construct density-based interpolations.

## Motivation

In the case of simple data sets, e.g. MNIST or Fashion-MNIST, it is possible to effectively train a generator network by directly minimizing an estimator of some distance between data and model distributions (Deshpande, Zhang, and Schwing 2018; Tabor et al. 2018). Unfortunately, such models may give poor results for the CelebA data set (see Tab. 1).

At first glance, it seems that kernel distance cannot be effectively applied in high dimensional spaces. But it turns out

<sup>1</sup> The code is available <https://github.com/gmum/lcw-generator>

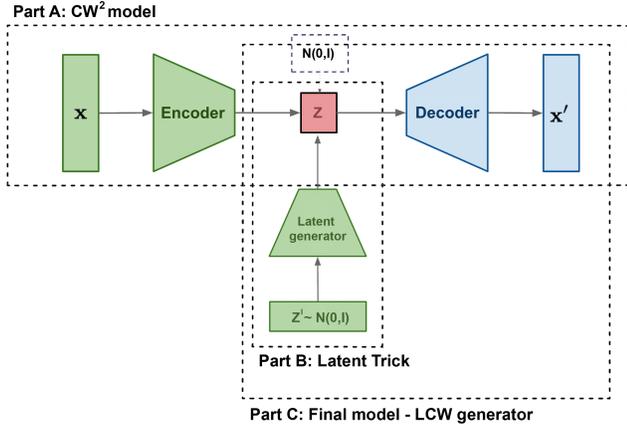


Figure 1: The LCW generator is trained in a two-stage procedure. First, a generative autoencoder (CW<sup>2</sup> model) is trained (see Part A). In the second stage, a new Latent generator neural network is introduced and trained to transport a standard Gaussian noise distribution  $\mathcal{N}(0, I)$  into the autoencoder latent space (see Part B). The final model is a concatenation of the latent generator and the decoder (see Part C).

that the CW distance used in the CW autoencoder (CWAE) model introduced in Tabor et al. (2018), to measure dissimilarity between distribution of encoded data and the Gaussian prior on the latent space, can also be efficiently utilized as a measure of reconstruction error (see the CW<sup>2</sup> model introduction below). This autoencoder model obtains state-of-the-art results, suggesting that the problem does not lie in an objective function but in a training procedure (see Tab. 1).

Inferior results in previous to CW<sup>2</sup> models were due to mini-batch training in high dimensional spaces with high noise. The model needs to minimize both the reconstruction error and the latent distribution distance from the prior, which might be complicated for non-simple data sets. In the proposed model, both are computed using kernel distances, and the training is two-step (see Fig. 2), which partly separates the responsibility for reconstruction/distribution training. Hence much better results are possible.

### General idea

The proposed LCW generator is trained in a two-stage procedure (see Fig. 1). First, the CW<sup>2</sup> model, a modification of CWAE (Tabor et al. 2018), is pretrained (see next section for details). Such architecture gives state-of-the-art FID score in class of generative autoencoder models but does not generate images as sharp as GANs do. To solve this problem the second stage of the construction is applied. There, the Latent Trick (defined thoroughly in later section), in which the current autoencoder architecture is fixed and the new latent generator is trained to transport standard Gaussian noise distribution  $\mathcal{N}(0, I)$  into the autoencoder latent space. This better approaches the encoded data distribution. The Latent Trick is implemented as a fully connected neural network. Conse-

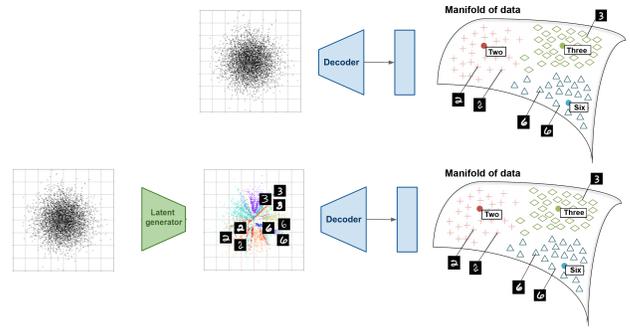


Figure 2: A classical decoder must simultaneously render manifold of data and transport a given prior into data space to model data distribution (see top row). In the proposed modification, we first model manifold of data with autoencoder and then model probability distribution of data using a Latent generator (see bottom row). This increases the generativity of the whole model. Visualization of data manifold inspired by (Socher et al. 2013).

quently, the final model is given as concatenation of the latent generator and previously trained decoder. The CW<sup>2</sup> autoencoder may still be used on its own.

	CWAE	CW-Generator	SW-Generator	CW <sup>2</sup>
MNIST	23.63	<b>16.48</b>	<b>17.85</b>	<b>17.44</b>
F-MNIST	49.95	<b>31.05</b>	48.02	<b>33.34</b>
CelebA	<b>49.69</b>	87.01	141.03	<b>47.09</b>

Table 1: FID scores on MNIST, Fashion-MNIST and CelebA data sets obtained with CWAE (state-of-the-art autoencoder based generative model), classical CW-Generator (see Tabor et al. (2018) for both), SW-Generator (Deshpande, Zhang, and Schwing 2018) and the proposed CW<sup>2</sup>.

We illustrate the concept of the model in Fig. 2. Classical generative autoencoder’s decoder (see the top row in Fig. 2) must simultaneously render manifold of data and transport a given prior distribution on the latent into data space to model data distribution. In the proposed solution, the data manifold is modeled first and then the probability distribution of data model is expanded using the Latent generator (see the bottom row in Fig. 2). Above is the crucial aspect of the proposed solution which inherits the positive properties of both autoencoder and GAN based generative methods. A stable model is obtained, with precise autoencoder latent space producing high-quality images without adversarial training.

### CWAE and CW<sup>2</sup> models

In general, all autoencoder based generative models are trained to minimize an objective function of the form

$$\mathcal{J}(X; \mathcal{E}, \mathcal{D}) = \text{Err}(\mathbf{X}, \mathcal{D}(\mathcal{E}(X))) + \lambda \text{DM}(P_{\mathcal{E}(X)}, P_Z),$$



Figure 3: Results of CWAE+LT and LCW generator models trained on CelebA data set. As we can see, we obtain state-of-the-art samples and interpolations in autoencoder based generative models.

where  $\text{Err}$  is a reconstruction error term,  $\lambda$  is a hyperparameter and  $\text{DM}$  denotes any dissimilarity, not necessarily non-negative, measure between probability distributions on the latent  $\mathcal{Z}$ . CWAE uses the mean squared error  $\text{MSE}$ , logarithm of the square of the CW distance  $d_{\text{CW}}$  and  $\mathcal{N}(0, I)$  as the prior  $P_{\mathcal{Z}}$ , which leads to the following formula

$$\text{MSE}(X, \mathcal{D}(\mathcal{E}(X))) + \lambda \log d_{\text{CW}}^2(P_{\mathcal{E}(X)}, \mathcal{N}(0, I)).$$

Following Tabor et al. (2018), we emphasize that  $d_{\text{CW}}$  can be calculated analytically and approximated as a distance between a latent sample  $Z = (z_i)_{i=1}^n$  and the standard Gaussian prior. Specifically,

$$2\sqrt{\pi} d_{\text{CW}}^2(Z, \mathcal{N}(0, I)) \approx \frac{1}{n^2} \sum_{ij} (\gamma_n + \frac{\|z_i - z_j\|^2}{2D_{\mathcal{Z}} - 3})^{-\frac{1}{2}} + (1 + \gamma_n)^{-\frac{1}{2}} - \frac{2}{n} \sum_i (\gamma_n + \frac{1}{2} + \frac{\|z_i\|^2}{2D_{\mathcal{Z}} - 3})^{-\frac{1}{2}},$$

consequently giving the following CWAE cost function

$$\mathcal{J}_{\text{CW}} = \text{MSE}(X, \mathcal{D}(\mathcal{E}(X))) + \lambda \log d_{\text{CW}}^2(\mathcal{E}(X), \mathcal{N}(0, I)),$$

where  $D_{\mathcal{Z}} = \dim \mathcal{Z}$  and  $\gamma_n$  bandwidth is chosen using Silverman’s rule of thumb, i.e.  $\gamma_n = \hat{\sigma}(\frac{4}{3n})^{2/5}$ , where  $\hat{\sigma}$  denotes a standard deviation that we assume to be equal to 1 as we deal with the standard Gaussian distribution.

We want to point out that the CW distance is given by a characteristic kernel with a closed-form for spherical Gaussians. Moreover, to the best of our knowledge, CWAE model was the first kernel distance based concept that required no sampling from the prior distribution.

Taking into consideration the results of Tabor et al. (2018), it is also possible to obtain the following approximate analytical formula that expresses  $d_{\text{CW}}$  for two given samples  $X = (x_i)_{i=1}^n$  and  $Y = (y_i)_{i=1}^n$  in  $\mathcal{X}$

$$2\sqrt{\pi} d_{\text{CW}}^2(X, Y) \approx \frac{1}{n^2} \sum_{ij} (\gamma_n + \frac{\|x_i - x_j\|^2}{2D_{\mathcal{X}} - 3})^{-\frac{1}{2}} + \frac{1}{n^2} \sum_{ij} (\gamma_n + \frac{\|y_i - y_j\|^2}{2D_{\mathcal{X}} - 3})^{-\frac{1}{2}} - \frac{2}{n} \sum_{ij} (\gamma_n + \frac{\|x_i - y_j\|^2}{2D_{\mathcal{X}} - 3})^{-\frac{1}{2}},$$

where  $D_{\mathcal{X}} = \dim \mathcal{X}$  and  $\gamma_n$  is calculated using standard deviation of joined  $X$  and  $Y$  samples.

This suggests using the CW distance not only in the latent, but also in the data space as  $\text{Err}(\mathbf{X}, \mathcal{D}(\mathcal{E}(X)))$ . Consequently we introduce the  $\text{CW}^2$  model with the following objective function<sup>2</sup>

$$\mathcal{J}_{\text{CW}^2} = d_{\text{CW}}^2(X, \mathcal{D}(\mathcal{E}(X))) + \lambda \log d_{\text{CW}}^2(\mathcal{E}(X), \mathcal{N}(0, I)),$$

minimized in LCW generator construction’s first phase. The square sign in  $\text{CW}^2$  denotes that it uses the CW distance twice.

### Latent Trick

Note that in the first stage (see Part A in Fig. 1) of model construction, the latent distribution is forced to be as close as possible to the Gaussian prior. Hence the obtained  $\text{CW}^2$  model is certainly generative.

<sup>2</sup>It should be mentioned here that this is not a classical autoencoder function, because a reconstruction error does not depend directly on differences between  $x_i$  and  $\mathcal{D}(\mathcal{E}(x_i))$  as in MSE, but is based on a distance between distributions.

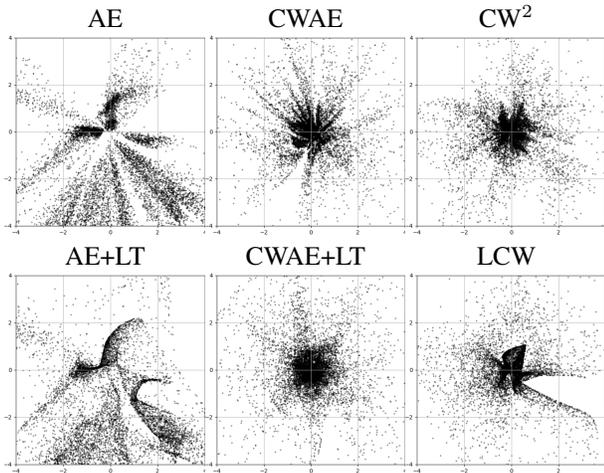


Figure 4: The top row shows encoded data distributions learned by AE, CWAE and  $CW^2$  models on Fashion-MNIST’s validation data set. Note that they are close to but different from the standard Gaussian noise. However, they are very similar to bottom row figures, which present the LCW generator’s output, applied for the above models, i.e the standard Gaussian distribution transported to respective latent spaces.

However, there remain two fundamental problems. One is because there are empty holes/spaces in the latent space , i.e. parts with very low data density being mapped. In the classical autoencoder based approach, there is a need for a compromise between reconstructing and generating abilities. A possible solution is to use an appropriate hyperparameter to balance the two terms (Higgins et al. 2017) but, in practice, it increases the generativity of the model at the expense of reconstruction.

Second problem seems to be more fundamental (Li, Swersky, and Zemel 2015; Dziugaite, Roy, and Ghahramani 2015; Li et al. 2017). It is related to using mini-batch training in high dimensional data space with much lower intrinsic data dimension. Each batch contains only a small, typically unbalanced, subset of a data set. We hypothesize that it affects kernel methods’ effectiveness because they need to learn representation and data distribution simultaneously, see Fig. 2.

To solve the above problems, we introduce the Latent Trick, which is the second stage of our procedure. It involves the creation of a latent generator

$$\mathcal{L}\mathcal{G}: (\mathcal{Z}', \mathcal{N}(0, I)) \longrightarrow (\mathcal{Z}, P_{\mathcal{E}(X)}),$$

which is a neural network trained to transform the standard Gaussian distribution (on the new  $\mathcal{Z}'$  space) so that it resembles the distribution of encoded data on  $CW^2$  model’s latent  $\mathcal{Z}$ . To be precise, the objective is to minimize the following function

$$\mathcal{J}_{\text{LT}}(X, \mathcal{Z}'; \mathcal{L}\mathcal{G}) = d_{\text{CW}}^2(\mathcal{E}(X), \mathcal{L}\mathcal{G}(\mathcal{Z}')),$$

where  $X$  and  $\mathcal{Z}'$  denote a data sample and a sample from  $\mathcal{N}(0, I)$  distribution, respectively. Consistently, to express

Data set	Method	Learn. rate	$\lambda$	Latent dim	Noise dim	Rec. error	FID score
MNIST	AE	1.e-3	-	8	-	11	52
	AE+LT	5.e-4	-	8	8	-	22
	CWAE	1.e-3	1	8	-	11	23
	CWAE+LT	5.e-4	-	8	8	-	20
	$CW^2$	1.e-3	1	8	-	14	<b>17</b>
	LCW	5.e-4	-	8	8	-	<b>14</b>
F-MNIST	AE+LT	5.e-4	-	8	8	-	41
	CWAE	1.e-3	10	8	-	10	49
	CWAE+LT	5.e-4	-	8	8	-	38
	$CW^2$	1.e-3	1	8	-	13	<b>33</b>
	LCW	5.e-4	-	8	8	-	<b>28</b>
CelebA	AE	1.e-3	-	64	-	66	328
	AE+LT	5.e-4	-	128	32	-	45
	CWAE	5.e-4	5	64	-	68	49
	CWAE+LT	5.e-4	-	128	32	-	<b>31</b>
	$CW^2$	5.e-4	0.2	64	-	71	47
	LCW	5.e-4	-	128	32	-	<b>33</b>

Table 2: Comparison of different architectures on MNIST, Fashion-MNIST and CelebA data sets. All models outputs except AE are similarly close to the normal distribution. CWAE achieves the best value of FID score (lower is better). All hyperparameters were found using a grid search (see appendix).

the CW distance an appropriate approximation formula is used (analogous to that in the previous section).

Note that the Latent Trick phase allows to rectify model’s generative power without losing it’s reconstruction quality. Consequently, the final LCW generator is provided as a function  $\mathcal{L}\mathcal{G}_{\text{CW}}: \mathcal{Z}' \longrightarrow \mathcal{X}$ , which transport a Gaussian noise sample  $\mathcal{Z}'$  into the data space, via concatenation of  $\mathcal{L}\mathcal{G}$  and  $\mathcal{D}$  (compare diagram in Fig. 1), i.e.

$$\mathcal{L}\mathcal{G}_{\text{CW}}(\mathcal{Z}') = \mathcal{D}(\mathcal{L}\mathcal{G}(\mathcal{Z}')).$$

The influence of the Latent Trick is visualized in Fig 4. Top row shows the mapping of the whole data set on the latent (in  $\mathbb{R}^2$  here for simplicity) for AE, CWAE, and  $CW^2$  models are shown, while the bottom row shows mapping of the Gaussian noise through Latent Trick onto the same space. The corresponding figures seem to differ slightly, i.e. the  $\mathcal{L}\mathcal{G}$  mappings can be understood as generalizations of the encoded data samples. Hence the small discrepancies.

Latent Trick can be used not only for the above construction but to any classical autoencoder based generative model that uses MSE as a reconstruction error. We examine some of these models in the Experiments section.

## Related work

We divided the related work section into two parts. First, we describe existing approaches to train GAN style models (Generators) with kernel measures in data space. Then we discuss existing solutions improving autoencoder properties by adding a neural network in latent space.



Figure 5: Reconstructions quality of AE, CWAE and  $CW^2$  (first stage of LCW generator) models trained on CelebA data set. Odd columns correspond to the real test points, while even to their reconstructions.

**Generators** Generative Moment Matching Network (GMMN) (Li, Swersky, and Zemel 2015; Dziugaite, Roy, and Ghahramani 2015) is a deep generative model that differs from Generative Adversarial Network (GAN) (Goodfellow et al. 2014) in replacing GAN’s discriminator with a two-sample test based on kernel maximum mean discrepancy (MMD) (Gretton et al. 2012). Unfortunately, these models work only for reasonable simple data sets like MNIST and Fashion-MNIST.

To solve such problems (Li et al. 2017) propose an MMD GAN. Authors improve the model expressiveness of GMMN and its computational efficiency by introducing adversarial kernel learning to replace a fixed Gaussian kernel in the original GMMN. The proposed algorithm is similar to GAN, aiming to optimize two neural networks in a minimax setting, but the objective’s meaning is different. In GAN we train a discriminator (binary) classifier to distinguish two distributions. MMD-GAN discriminates two distributions with a two-sample test via MMD, but with an adversarially learned kernel. It is similar to our approach in using kernel distance but, contrary to LCW generator proposed, still uses adversarial training.

Deshpande, Zhang, and Schwing (2018) show that it is possible to train a GAN like architecture (namely an SW-Generator) by substituting discriminator with kernel measure in data space. In practice, the authors use a Sliced Wasserstein distance. Since the sliced method can reduce data dimensionality by random projections, a model can be trained effectively using kernel measures in high dimensional spaces. Unfortunately, their model works only for reasonable simple data sets like MNIST and Fashion-MNIST.

**Autoencoder based generative models** First and one of the most popular approaches to autoencoder based generative models is VAE (Kingma and Welling 2014). In practice, such models give correct results, but geometry of latent space might have some very low probability areas (empty spaces) (Tolstikhin et al. 2018). An idea to join the training algorithms and, following it, the merits of autoencoders and

GANs, was presented early in a widely cited paper (Larsen et al. 2016). Another method to solve such problems can be obtained by adding additional architecture to the latent, in order to obtain better representation. For example, Ziegler and Rush (2019); Xiao, Yan, and Amit (2019) apply normalizing flows (Kingma and Dhariwal 2018) in the latent space. Thanks to this the latent distribution (which is similar to the standard Gaussian) is transformed into the Gaussian prior. Dai and Wipf (2018), on the other hand, present similar solution based on adding additional autoencoder in the latent space (TwoStageVAE model). They present theoretical results stating that VAE model does not properly approximate properly prior distribution in the latent space. But, as they propose, the second VAE model is able to correct distribution in the latent. Deja et al. (2020) use Sinkhorn autoencoder with Noise Generator (e2e SAE), a simple fully connected architecture which transfers Gaussian noise into the latent space. This model is trained end-to-end and gives similar results to WAE-MMD (Tolstikhin et al. 2018).

## Experiments

In this section, we empirically validate the proposed LCW generator model on standard benchmarks for autoencoder based generative models: CelebA, MNIST and Fashion-MNIST. Since MNIST and Fashion-MNIST are relatively simple, we use them rather as toy examples. CelebA data sets will be used to show real difference between models, see Tab. 3.

It should be mentioned that GAN models obtain essentially better results. But our goal is not to outperform GANs but rather to reduce the gap between the generative quality of the GAN based and non-adversarial AE based models. At the same time, as a “by-product”, we propose the  $CW^2$  model. Still being an autoencoder based model, it obtains much better FID values than other autoencoders (see, e.g., Tab. 4).

**Quantitative tests** To quantitatively compare LCW generator with other models, in the first experiment we compare



Figure 6: Results of AE and AE+LT models trained on CelebA data set. Interpolation in AE (top row) are constructed linearly in the latent between “endpoint” images from AE+LT, while in AE+LT a linear interpolation between samples from  $\mathcal{N}(0, I)$  is performed to be mapped as  $\mathcal{L}\mathcal{G}(z')$  points. This corresponds to two types of interpolations in Fig. 7. As we can see, classical AE is not a generative model (top right). By applying Latent Trick we produce generative model AE+LT (bottom right).

	CW-Generator	SW-Generator	LCW generator
MNIST	16.48	17.85	<b>14.92</b>
F-MNIST	31.05	48.02	<b>28.04</b>
CelebA	87.01	141.03	<b>33.07</b>

Table 3: FID scores on MNIST, Fashion-MNIST and CelebA data sets obtained with CW-Generator, SW-Generator and LCW generator.

different models which use kernel distance in data space, i.e. CW-, SW- and LCW generators. To assess results the Fréchet Inception Distance FID is used (Heusel et al. 2017). As it can be easily seen, LCW generator outperform other approaches. In the case of MNIST and Fashion-MNIST data set the difference is smaller than in the case of CelebA, since the two first data sets are quite simple.

In the other experiment we follow the classical competition setup for generative models on Fashion-MNIST data set. In agreement with the qualitative studies, we observe FID values of LCW generator to be better (lower) than those for VAE, CWAE, SWAE (Kolouri et al. 2019), WAE-MMD and WAE-GAN (Tolstikhin et al. 2018).

We stress here that LCW generator on CelebA achieves FID score 33, compared to 52 and 42 achieved by WAE-

Model	FID	Model	FID	Model	FID
VAE	60	<b>2Stage-VAE*</b>	<b>34</b>	<b>CWAE + LT</b>	<b>31</b>
CWAE	49	CW <sup>2</sup>	47	<b>LCW</b>	<b>33</b>
WAE	52	WAE-GAN*	42	e2e SAE*	54
DCGAN*	12.5	WGAN-GP*	20.3	WGAN-div*	17.5

Table 4: FID values for several compared models, computed on the CelebA data set. The Latent Trick addition expands the model generative capabilities. Models marked with an asterisk have results taken from literature. The DCGAN value was trained with a special two time scale algorithm (Heusel et al. 2017). WGAN have convolutional architectures and the results are taken from (Wu et al. 2017).

MMD and WAE-GAN (see Tab. 4). It should be mentioned that LCW generator gives similar results to 2Stage-VAE and CWAE+LT, which also use two stage training but in first part use MSE as a reconstruction cost.

**Ablation study** In this paragraph, we show how the proposed two-stage training procedure works in various autoencoder-based models. As it was already mentioned, the Latent Trick is a possible solution to the problem with si-

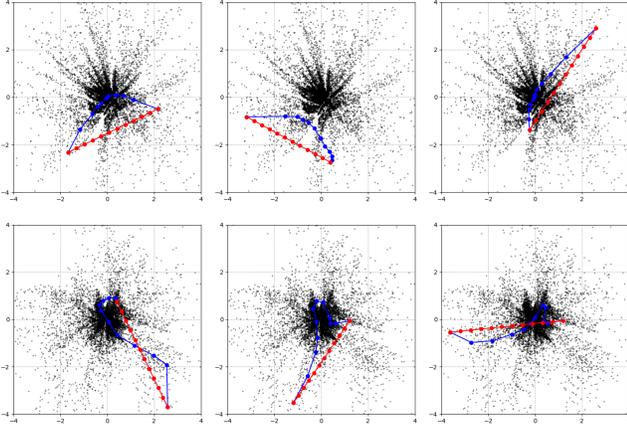


Figure 7: Interpolations in CWAE+LT (top row) and LCW generator (bottom row) models’ latent spaces, trained on CelebA data set. The red dots show a GAN like linear interpolation  $\alpha_i \mathcal{L}\mathcal{G}(z'_1) + (1 - \alpha_i) \mathcal{L}\mathcal{G}(z'_k)$  for  $z'_1, z'_k \sim \mathcal{N}(0, I)$ , while the blue ones are obtained by mapping a linear interpolation  $\alpha_i z'_1 + (1 - \alpha_i) z'_k$  via  $\mathcal{L}\mathcal{G}$ . In the latter case curves follow the latent distribution and avoid empty spaces.

multaneous modeling of a manifold of data and generativity of the model. We claim that we model the manifold of data in the first stage, while the latent distribution in the second stage .

To test whether this hypothesis is true, we apply the Latent Trick to some models. We start with a standard autoencoder (AE), which models only the manifold of data. It turns out that such models work surprisingly well (see the difference between AE and AE+LT methods in Tab. 2). It should be mentioned that AE+LT gets FID score of 45 on CelebA, which is comparable to that of WAE-GAN’s (42) and WAE-MMD’s (52). It is still an open question why the Latent Trick added to CWAE and CW<sup>2</sup> models performs essentially better score than AE+LT (see Tab. 2). It seems that in AE we are able to model data manifold but latent representation is not constrained, so the Latent Trick is not able to describe distribution of data efficiently. In Fig. 6 we present qualitative tests obtained by interpolation and sampling.

Summarizing, Latent Trick can be used to transform any autoencoder based architecture into generative model.

**Qualitative tests** Since the LCW generator, similarly to GAN models, does not produce reconstructions, we compare only those of CW<sup>2</sup> model (trained as the first stage of LCW generator) to ones produced by CWAE and vanilla AE, see Fig. 5. We want to stress that CW<sup>2</sup> produces quite accurate reconstructions even though it is not explicitly required (point to point) by its cost function.

**Interpolations** It is challenging to see the different interpolation abilities of the proposed LCW generator model and its parts. The least interesting is that using the CW<sup>2</sup> model only since it is a simple, autoencoder like, interpolation in

the latent  $\mathcal{Z}$  space. Thus we shall skip it.

On the other hand, the construction of the Latent Trick allows to make two interesting interpolations. First, let us draw  $z'_1, z'_k \sim \mathcal{N}(0, I)$ , i.e. two points from Latent Trick’s input space  $\mathcal{Z}'$ . Mapping  $z_1 = \mathcal{L}\mathcal{G}(z'_1), z_k = \mathcal{L}\mathcal{G}(z'_k)$  it is possible to perform a GAN like type linear interpolation  $z_1, z_2, \dots, z_k$  in the latent  $\mathcal{Z}$ . On the other hand, there is a possibility to produce a linear interpolation  $z'_1, z'_2, \dots, z'_k$  in  $\mathcal{Z}'$  and map all those points to  $\mathcal{L}\mathcal{G}(z'_1), \mathcal{L}\mathcal{G}(z'_2), \dots, \mathcal{L}\mathcal{G}(z'_k)$ . This produces a density-based interpolation since  $\mathcal{L}\mathcal{G}$  network is trained to generate the latent distribution. The results can be seen in Fig. 7, where red dots show the standard like interpolation, while blue ones are obtained using the Latent Trick. Note that in the latter case the model can follow areas of high data density.

Additionally, using all considered models we can construct interpolation and sampling, see Fig. 3. In the case of sampling we can see that LCW generator gives state-of-the-art autoencoder based faces.

The above clearly shows that the proposed LCW generator model can better map the latent space. Several methods for interpolation in latent spaces provide sophisticated approaches to computing the latent space data mapping density, frequently using Riemannian curvatures, see e.g. (Agustsson et al. 2017; Arvanitidis, Hansen, and Hauber 2018; Lesniak, Sieradzki, and Podolak 2019). The proposed model generates density-based interpolations thanks only to applying the Latent Trick.

## Architecture

For CWAE model we used hyperparameters reported in Tabor et al. (2018). For CW<sup>2</sup> model, SW-Generator and CW-Generator we performed a grid search over batch-size in  $\{64, 128, 256, 512\}$  and learning rate values in  $\{0.005, 0.0025, 0.001, 0.0005, 0.00025\}$ . For every model, we reported results for configuration that achieved the lowest value of FID Score.

Autoencoder feed-forward architecture for MNIST or Fashion-MNIST ( $28 \times 28$  images)

**encoder** three feed-forward ReLU layers, 200 neurons each,

**decoder** three feed-forward ReLU layers, 200 neurons each followed by feed-forward sigmoid layer.

Autoencoder architecture for CelebA (with images centered and cropped to  $64 \times 64$  with 3 color layers):

**encoder**

four convolution layers with  $4 \times 4$  filters, each layer was followed by a batch normalization (consecutively 128, 256, 512, and 1024 channels) and ReLU activation,

**decoder**

dense 1024 neuron layer,

three transposed-convolution layers with  $4 \times 4$  filters, and each layer followed by a batch normalization with ReLU activation (consecutively 512, 256, and 128 channels),

transposed-convolution layer with  $3 \times 3$  filter, 3 channels and tanh activation.

For LCW generator architecture we used five feed-forward RELU layers with batch normalization and 512 neurons. The final layer was a feed-forward layer with a linear activation.

## Conclusions

In this paper, we presented the new LCW generator. According to our knowledge, it is the first model that can be effectively trained using a kernel distance in high dimensional data sets. It needs no adversarial training and does not suffer from “mode collapse”. Consequently, the proposed approach might be situated between autoencoder and GAN models, while not inheriting their main weaknesses.

Our experiments show that LCW generator gives superior FID values for generated samples. Another interesting feature is that it is a natural generator of density-based interpolations in the latent space. Its ability to omit all low-density areas might give a robust tool for generating new samples with a smooth changing of features.

Finally, we want to stress that the Latent Trick may be applied to any already trained autoencoder based model, to increase its generative capabilities.

## References

- Agustsson, E.; Sage, A.; Timofte, R.; and Van Gool, L. 2017. Optimal Transport Maps for Distribution Preserving Operations on Latent Spaces of Generative Models. *arXiv:1711.01970*.
- Arvanitidis, G.; Hansen, L. K.; and Hauber, S. 2018. Latent Space Oddity: on the Curvature of Deep Generative Models. In *Proc. of the Int. Conf. on Learning Representations*, ICLR’2018.
- Dai, B.; and Wipf, D. 2018. Diagnosing and Enhancing VAE Models. In *Proc. of the Int. Conf. on Learning Representations*, ICLR’2018.
- Deja, K.; Dubiński, J.; Nowak, P.; Wenzel, S.; and Trzcíński, T. 2020. End-to-end Sinkhorn Autoencoder with Noise Generator. *arXiv:2006.06704*.
- Deshpande, I.; Zhang, Z.; and Schwing, A. 2018. Generative Modeling Using the Sliced Wasserstein Distance. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition*, CVPR’2018, 3483–3491.
- Dziugaite, G. K.; Roy, D. M.; and Ghahramani, Z. 2015. Training Generative Neural Networks via Maximum Mean Discrepancy Optimization. In *Proc. of the Conf. on Uncertainty in Artificial Intelligence*, UAI’15.
- Goodfellow, I.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; and Bengio, Y. 2014. Generative Adversarial Nets. In *Adv. in Neural Information Processing Systems*, NeurIPS’2014, 2672–2680.
- Gretton, A.; Borgwardt, K. M.; Rasch, M. J.; Schölkopf, B.; and Smola, A. 2012. A Kernel Two-sample Test. *The Journal of Machine Learning Research* 13(1): 723–773.
- Heusel, M.; Ramsauer, H.; Unterthiner, T.; Nessler, B.; and Hochreiter, S. 2017. GANs Trained by a Two Time-scale Update Rule Converge to a Local Nash Equilibrium. In *Adv. in Neural Information Processing Systems*, NeurIPS’2017, 6626–6637.
- Higgins, I.; Matthey, L.; Pal, A.; Burgess, C.; Glorot, X.; Botvinick, M.; Mohamed, S.; and Lerchner, A. 2017.  $\beta$ -VAE: Learning Basic Visual Concepts with a Constrained Variational Framework. In *Proc. of the Int. Conf. on Learning Representations*, ICLR’2017.
- Kingma, D.; and Welling, M. 2014. Auto-encoding Variational Bayes. *arXiv:1312.6114*.
- Kingma, D. P.; and Dhariwal, P. 2018. Glow: Generative Flow with Invertible 1x1 Convolutions. In *Adv. in Neural Information Processing Systems*, NeurIPS’2018, 10236–10245.
- Kolouri, S.; Pope, P. E.; Martin, C. E.; and Rohde, G. K. 2019. Sliced Wasserstein Auto-Encoders. In *Proc. of the Int. Conf. on Learning Representations*, ICLR’2019.
- Larsen, A. B. L.; Sønderby, S. K.; Larochelle, H.; and Winther, O. 2016. Autoencoding Beyond Pixels Using a Learned Similarity Metric. In *Proc. of the Int. Conf. on Machine Learning*, ICML’16, 15581566. JMLR.org.
- Lesniak, D.; Sieradzki, I.; and Podolak, I. 2019. Distribution Interpolation Trade-off in Generative Models. In *Proc. of the Int. Conf. on Learning Representations*, ICLR’2019.
- Li, C.-L.; Chang, W.-C.; Cheng, Y.; Yang, Y.; and Póczos, B. 2017. MMD GAN: Towards Deeper Understanding of Moment Matching Network. In *Adv. in Neural Information Processing Systems*, NeurIPS’2017, 2203–2213.
- Li, Y.; Swersky, K.; and Zemel, R. 2015. Generative Moment Matching Networks. In *Proc. of the Int. Conf. on Machine Learning*, ICML’2015, 1718–1727.
- Socher, R.; Ganjoo, M.; Manning, C. D.; and Ng, A. 2013. Zero-shot Learning Through Cross-modal Transfer. In *Adv. in Neural Information Processing Systems*, NeurIPS’2013, 935–943.
- Tabor, J.; Knop, S.; Spurek, P.; Podolak, I.; Mazur, M.; and Jastrzębski, S. 2018. Cramer-Wold Autoencoder. *arXiv:1805.09235*.
- Tolstikhin, I.; Bousquet, O.; Gelly, S.; and Schoelkopf, B. 2018. Wasserstein Auto-encoders. In *Proc. of the Int. Conf. on Learning Representations*, ICLR’2018.
- Wu, J.; Huang, Z.; Thoma, J.; Acharya, D.; and Gool, L. V. 2017. Wasserstein Divergence for GANs.
- Xiao, Z.; Yan, Q.; and Amit, Y. 2019. Generative Latent Flow. *arXiv:1905.10485*.
- Ziegler, Z. M.; and Rush, A. M. 2019. Latent Normalizing Flows for Discrete Sequences. In *Proc. of the Int. Conf. on Machine Learning*, ICML’2019.