

# Mind2Mind : Transfer Learning for GANs

Yael Fregier, Jean-Baptiste Gouray \*

Univ. Artois, UR2462

Laboratoire de Mathématiques de Lens (LML)

F-62300 Lens, France

yael.fregier@univ-artois.fr, jeanbaptiste.gouray@gmail.com

## Abstract

Training generative adversarial networks (GANs) on high quality (HQ) images involves important computing resources. This requirement represents a bottleneck for the development of applications of GANs. We propose a transfer learning technique for GANs that significantly reduces training time. Our approach consists of freezing the low-level layers of both the critic and generator of the original GAN. We assume an auto-encoder constraint in order to ensure the compatibility of the internal representations of the critic and the generator. This assumption explains the gain in training time as it enables us to bypass the low-level layers during the forward and backward passes. We compare our method to baselines and observe a significant acceleration of the training. It can reach two orders of magnitude on HQ datasets when compared with StyleGAN. We prove rigorously, within the framework of optimal transport, a theorem ensuring the convergence of the learning of the transferred GAN. We moreover provide a precise bound for the convergence of the training in terms of the distance between the source and target dataset.

## 1 Introduction

The recent rise of deep learning as a leading paradigm in AI mostly relies on computing power (with generalized use of GPUs) and massive datasets. These requirements represent bottlenecks for most practitioners outside of big labs in industry or academia and are the main obstacles to the generalization of the use of deep learning. Therefore, methods that can bypass such bottlenecks are in strong demand. Transfer learning is one of them and various methods of transfer learning (for classification tasks) specific to deep neural networks have been developed (Tan et al. 2018).

A *generative problem* is a situation in which one wants to be able to produce elements that could belong to a given data set  $\mathcal{D}$ . Generative Adversarial Networks (GANs) were introduced in 2014 (Goodfellow et al. 2014) (Salimans et al. 2016) to tackle generative tasks with deep learning architectures and improved in (Arjovsky, Chintala, and Bottou 2017; Gulrajani et al. 2017) (Wasserstein GANs). They immediately took a leading position in the world of generative models, comforted by progress with HQ images (ProGAN github

repository), (Karras, Laine, and Aila 2019). The goal of our work (see section 4) is to develop in a generative setting, i.e., for GAN architectures, the analog of the *cut-and-paste* approach.

The main idea of our method is to reuse, for the training of a GAN, the weights of an autoencoder already trained on a source dataset  $\mathcal{D}$ . The weights of the low level layers of the generator (resp. critic) will be given by those of the decoder (resp. encoder). We call *MindGAN* the high level layers of the generator. It is a subnetwork that is trained as a GAN on the encoded features of the target dataset  $\mathcal{D}'$ .

We prove in section 5 a theorem that controls the convergence of the transferred GAN in terms of the quality of the autoencoder, the domain shift between  $\mathcal{D}$  and  $\mathcal{D}'$  and the quality of the MindGAN. As a consequence, our experimental results in section 6 rely heavily on the choice of the autoencoder. ALAE autoencoders (Pidhorskyi, Adjeroh, and Doretto 2020) are extremely good autoencoders that turned out to be crucial in our experiments with HQ images. Their use, in conjunction with our transfer technique, enables an **acceleration of the training by a factor of 656**, while keeping a good quality.

## 2 Preliminaries

A GAN consists of two networks trained adversarially. The *generator*  $g : Z \rightarrow \chi$  associates to a vector  $z$  sampled from a latent vector space  $Z$  a vector  $g(z)$  in another vector space  $\chi$  while the *discriminator*  $c : \chi \rightarrow \mathbb{R}$  learns to associate a value close to 1 if the vector  $g(z)$  belongs to  $\mathcal{D}$  and zero otherwise. Their respective loss functions,  $L_g$  and  $L_c$  are recalled in section 4.

One can assume that elements of  $\mathcal{D}$  can be sampled from an underlying probability distribution  $\mathbb{P}_{\mathcal{D}}$  on a space  $\chi$  and try to approximate it by  $\mathbb{P}_{\theta}$ , another distribution on  $\chi$  that depends on some learnable parameters  $\theta$ . *Generating* then means sampling from the distribution  $\mathbb{P}_{\theta}$ . The main idea behind a Wasserstein GAN is to use the *Wasserstein distance* (see appendix A, and (Villani 2008) definition 6.1) to define by  $W(\mathbb{P}_{\mathcal{D}}, \mathbb{P}_{\theta})$  the loss function for this optimisation problem. More precisely, the Wasserstein distance is a distance on Borel probability measures on  $\chi$  (when compact metric space). In particular, the quantity  $W(\mathbb{P}_{\mathcal{D}}, \mathbb{P}_{\theta})$  gives a number which depends on the parameter  $\theta$  since  $\mathbb{P}_{\theta}$  depends itself on  $\theta$ . The main result of (Arjovsky, Chintala, and Bot-

\*both authors contributed equally

ou 2017) asserts that if  $\mathbb{P}_\theta$  is of class  $\mathcal{C}^k$  as a function of  $\theta$  almost everywhere,  $W(\mathbb{P}_D, \mathbb{P}_\theta)$  is also of class  $\mathcal{C}^k$  with respect to  $\theta$ . As a consequence, one can solve this optimization problem by doing gradient descent for the parameters  $\theta$  (using a concrete gradient formula provided by the same theorem) until the two probability distributions coincide.

Among the distributions on  $\chi$ , some can be obtained from a prior distribution  $\mathbb{P}_Z$  on an auxiliary latent space  $Z$  and a map  $g : Z \rightarrow \chi$  as follows. The *push-forward* of  $\mathbb{P}_Z$  under  $g$  (Bogachev 2007) is defined so that a sample is given by  $g(z)$  the image through  $g$  of a sample  $z$  from the distribution  $\mathbb{P}_Z$ . We will denote this pushforward by  $g_\# \mathbb{P}_Z$  and when  $g$  depends on parameters  $\theta$  use instead the notation  $\mathbb{P}_\theta := g_\# \mathbb{P}_Z$ . In practice, one can choose for  $\mathbb{P}_Z$  a uniform or Gaussian distribution, and for  $g$  a (de)convolution deep neural network. In our applications, we will consider for instance  $Z := \mathbb{R}^{128}$  equipped with a multivariate gaussian distribution and  $\chi = [-1, 1]^{28 \times 28}$ , the space of gray level images of resolution  $28 \times 28$ . Hence, sampling from  $\mathbb{P}_\theta$  will produce images.

In order to minimise the function  $W(\mathbb{P}_D, \mathbb{P}_\theta)$ , one needs a good estimate of the Wasserstein distance. The *Rubinstein-Kantorovich duality* (See (Villani 2008) theorem 5.9) states that  $W(\mathbb{P}, \mathbb{P}') = \max_{|c|_L \leq 1} \mathbb{E}_{x \sim \mathbb{P}} c(x) - \mathbb{E}_{x \sim \mathbb{P}'} c(x)$ , where  $\mathbb{E}_{x \sim \mathbb{P}} f(x)$  denotes the expected value of the function  $f$  for the probability measure  $\mathbb{P}$ , while the max is taken on the unit ball for the Lipschitz semi-norm. Concretely, this max is obtained by gradient ascent on a function  $c_\theta$  encoded by a deep convolution neural network.

In our case, when  $\mathbb{P}' := \mathbb{P}_\theta$ , the term  $\mathbb{E}_{x \sim \mathbb{P}'} c(x)$  takes the form  $\mathbb{E}_{z \sim \mathbb{P}_Z} c_\theta(g_\theta(z))$ . One recovers the diagram

$$Z \xrightarrow{g_\theta} \chi \xrightarrow{c_\theta} \mathbb{R} \quad (1)$$

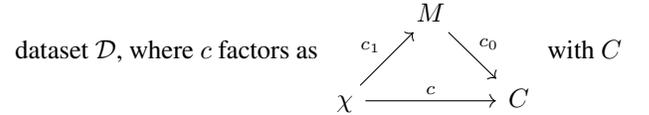
familiar in the adversarial interpretation of GANs. With this observation, one understands that one of the drawbacks of GANs is that there are two networks to train. They involve many parameters during the training, and the error needs to backpropagate through all the layers of the two networks combined, i.e., through  $c_\theta \circ g_\theta$ . This process is computationally expensive and can trigger the vanishing of the gradient. Therefore, specific techniques need to be introduced to deal with very deep GANs, such as in (Karras et al. 2018a). The approach we present in section 4 can circumvent these two problems.

*Transfer learning* is a general approach in machine learning to overcome the constraints of the volume of data and computing power needed for training models. It leverages a priori knowledge from a learned task  $\mathcal{T}$  on a source data set  $\mathcal{D}$  in order to learn more efficiently a task  $\mathcal{T}'$  on a target data set  $\mathcal{D}'$ . It applies in deep learning in at least two ways: *Cut and Paste* and *Fine tuning*.

**Cut and Paste** takes advantage of the difference between high and low-level layers of the network  $c$ . It assumes that the network  $c$  is composed of two networks  $c_0$  and  $c_1$  stacked one on each other, i.e., mathematically that  $c = c_0 \circ c_1$  (one understands the networks as maps). While the low-level layers  $c_1$  process low-level features of the data, usually common for similar datasets, the high-level layers  $c_0$  are in charge of the high-level features which are very specific to

each dataset. Hence, instead of retraining all the weights of an auxiliary network pre-trained on  $\mathcal{D}$ , one can retrain only the parameters of the last layers of the network while keeping the other parameters untouched. This approach boils down to the following steps :

1. identify a dataset  $\mathcal{D}$  similar to the data set  $\mathcal{D}'$  we are interested in, both in the same space  $\chi$ ,
2. import a network  $c = c_0 \circ c_1$ , already trained on the



the space of classes.

3. pass the new dataset  $\mathcal{D}'$  through  $c_1$  and train  $c'_0$  on  $c_1(\mathcal{D}')$ , and
4. use  $c' := c'_0 \circ c_1$  as the new classifier on  $\mathcal{D}'$ .

The main advantages of this approach are the following :

- a. much fewer parameters to train (only the parameters of  $c'_0$ , which in practice correspond to a few dense layers),
- b. need to pass the data  $\mathcal{D}'$  only once through  $c_1$ , and
- c. no need to backpropagate the error through  $c_1$ .

Since the low-level features often represent the most time-consuming part of the training, eliminating the need to train their weights will accelerate the process. If the datasets are similar, only training the last layers generally leads to good results in a much shorter time and with fewer data.

**Fine tuning** is based on the same first two steps than *Cut and Paste*, but instead of steps 3 and 4, uses the weights of  $c$  to initialise the training of the new network  $c'$  on  $\mathcal{D}'$ . It is assumed that  $c$  and  $c'$  share the same architecture.

In both approaches, the network  $c$  must have been previously trained by a third party during a very long time on the source dataset  $\mathcal{D}$ , which is potentially much more massive than  $\mathcal{D}'$ . It turns out in practice that these approaches enable to train a network on a new task with much less computing power and data (see (Donahue et al. 2014)).

### 3 Related works

In the following section, we survey the works to which this paper can be associated. We postpone the analysis of the main differences with these works in section 7.

**Wasserstein autoencoders** A version of autoencoders in conjunction with GANs has been considered in (Makhzani et al. 2015) and later generalized with Wasserstein distance in (Tolstikhin et al. 2017) and subsequent works. In short, in this approach, one trains an adversarially an autoencoder. Moreover, one adds a regularizer term to get a generative model. More details in section A.

**Adversarial learned inference** A second stream of papers, (Dumoulin et al. 2017), (Donahue, Krähenbühl, and Darrell 2017), (Belghazi et al. 2018), (Haidar and Rezagholizadeh 2019) and (Zhao et al. 2018) to cite a few, uses another blend of autoencoders with GANs. Their key idea is to learn adversarially an encoder  $g_x : \chi \rightarrow M$  together with a decoder  $g_m : M \rightarrow \chi$  against a discriminator  $c : \chi \times M \rightarrow \mathbb{R}$  in a way that the distributions given

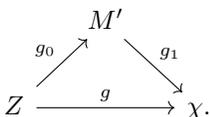
by the couples  $(g_m(m), m)$  and  $(x, g_x(x))$  are indistinguishable from the discriminator point of view. Note that one does not explicitly train  $g_x$  and  $g_m$  to be inverses to each other. However theorem 2 of (Donahue, Krähenbühl, and Darrell 2017) shows that at optimality they are indeed.

**Adversarial Latent Autoencoders** These autoencoders (Pidhorskyi, Adjeroh, and Doretto 2020) have the property of learning the latent distribution to match the encoded distribution. This is very different from other traditional approaches that assume an a priori target latent distribution and learn the encoder to match the encoded distribution with this a priori target. From this perspective, this method is very close to the architecture we use in our work, though the objectives are very different : their point is to disentangle representations in order to be able to control the features, whereas our objective is to do transfer. The similarity in architecture probably explains why ALAE autoencoders are very suited to our method.

**Fine-Tuning** On the side of papers addressing transfer for GANs, we are aware of (Wang et al. 2018), (Shan et al. 2018). Both apply to GANs *fine-tuning*, one of the techniques of transfer learning. It consists in initializing the training of a network on a target dataset  $\mathcal{D}'$  with weights from another network with the same architecture, but already trained on a similar source dataset  $\mathcal{D}$ . The two papers seem to have been written independently. While (Shan et al. 2018) is mainly targeting a specific application of de-noising in medical imagery, (Wang et al. 2018) is rather interested in understanding fine-tuning for GANs per se. Both report a faster convergence and a better quality, though (Wang et al. 2018) also observes that fine-tuning enables training with smaller datasets and that the distance between the source and target datasets influences the quality of the training.

## 4 Mind to mind algorithm

We now adapt to GANs the cut-and-paste procedure described in 2. The difference is that in addition to the classifier  $c$  (or *critic* in the language of WGAN), one also has a generator  $g$ . Let us consider a factorisation of the form  $g = g_1 \circ g_0$




---

### Algorithm 1 Mind2Mind transfer learning.

---

**Require:**  $(c_1, g_1)$ , an autoencoder trained on a source dataset  $\mathcal{D}$ ,  $\alpha$ , the learning rate,  $b$ , the batch size,  $n$ , the number of iterations of the critic per generator iteration,  $\mathcal{D}'$ , a target dataset,  $\varphi'$  and  $\theta'$  the initial parameters of the critic  $c'_0$  and of the generator  $g'_0$ .  
 Compute  $c_1(\mathcal{D}')$ .  
**while**  $\theta'$  has not converged **do**  
   **for**  $t = 0, \dots, n$  **do**  
     Sample  $\{m^{(i)}\}_{i=1}^b \sim c_{1\#}\mathbb{P}_{\mathcal{D}'}$  a batch from  $c_1(\mathcal{D}')$ .  
     Sample  $\{z^{(i)}\}_{i=1}^b \sim \mathbb{P}_Z$  a batch of prior samples.  
     Update  $c'_0$  by descending  $L_c$ .  
   **end for**  
   Sample  $\{z^{(i)}\}_{i=1}^b \sim \mathbb{P}_Z$  a batch of prior samples.  
   Update  $g'_0$  by descending  $-L_g$ .  
**end while**  
**return**  $g_1 \circ g'_0$ .

---

Our algorithm assumes that  $g_1$  comes from an autoencoder  $(c_1, g_1)$  that has been trained on a source dataset  $\mathcal{D}$ . The algorithm passes the second data set  $\mathcal{D}'$  through the encoder  $c_1$  and trains a MindGAN  $(g'_0, c'_0)$  on the encoded data  $c_1(\mathcal{D}')$ . One obtains the final generator as the composition  $g_1 \circ g'_0$  of  $g_1$ , the decoder of the autoencoder, with  $g'_0$ , the generator of the MindGAN. We denote by  $L_c$  and  $L_g$  the losses of the discriminator and the generator.

In the remainder of the paper, we use for  $L_g$  and  $L_c$  the losses of a WGAN with gradient penalty (Gulrajani et al. 2017) :  $L_g := \mathbb{E}_{z \sim \mathbb{P}_Z} c'_0(g'_0(z))$ ,  $L_c := -\mathbb{E}_{m \sim c_{1\#}\mathbb{P}_{\mathcal{D}'}} c'_0(m) + L_g + \lambda \mathbb{E}_{m \sim c_{1\#}\mathbb{P}_{\mathcal{D}'}, z \sim \mathbb{P}_Z, \alpha \sim (0,1)} \{(\|\nabla c'_0(\alpha m + (1-\alpha)g'_0(z))\|_2 - 1)^2\}$ .

**Remark 1.** This algorithm can be applied (with minor modifications) to conditional GANs. We refer to Appendix B for more details.

**Motivation for the approach** The architectures of a generator and a critic of a GAN are symmetric to one another. The high-level features appear in  $g_0$ , the closest to the prior vector (resp.  $c_0$ , the closest to the prediction), while the low-level features are in  $g_1$ , the closest to the generated sample (resp.  $c_1$ , the closest to the input image). Therefore, the analogy with cut and paste is to keep  $g_1$  (the low level features of  $\mathcal{D}$ ) and only learn the high level features  $g'_0$  of the target data set  $\mathcal{D}'$ . However, the only way a generator can access to information from  $\mathcal{D}'$  is through the critic  $c$  via the value of  $c \circ g = c'_0 \circ c_1 \circ g_1 \circ g'_0$  (Gulrajani et al. 2017). Hence, the information needs to back-propagate through  $c_1 \circ g_1$  to reach the weights of  $g'_0$ . Our main idea is to bypass this computational burden and train directly  $g'_0$  and  $c'_0$  on  $c_1(\mathcal{D}')$ . But this requires that the source of  $c'_0$  coincides with the target of  $g'_0$ . Therefore, we assume that  $M = M'$ , a first hint that autoencoders are relevant for us.

A second hint comes from an analogy with humans learning a task, like playing tennis, for instance. One can model a player as a function  $Z \xrightarrow{g} \chi$ , where  $\chi$  is the space of physical actions of the player. His/her coach can be understood as a function  $\chi \xrightarrow{c} \mathbb{R}$ , giving  $c(g(z))$  as a feedback for an action  $g(z)$  of  $g$ . The objective of the player can be understood as to be able to generate instances of the distribution  $\mathcal{D}'$  on  $\chi$

corresponding to the “tennis moves”. However, in practice, a coach rarely gives his/her feedback as a score. He instead describes what the player has done and should do instead. We can model this description as a vector  $c_1(g(z))$  in  $M$ , the mind of  $c = c_0 \circ c_1$ . In this analogy,  $c_1$  corresponds to the coach analyzing the action, while  $c_0$  corresponds to the coach giving a score based on this analysis. One can also decompose the player itself as  $g = g_1 \circ g_0$ . Here  $g_0$  corresponds to the player conceiving the set of movements he/she wants to perform and  $g_1$  to the execution of these actions. Therefore, two conditions are needed for the coach to help efficiently his/her student :

1. they must speak the same language in order to understand one each other,
2. the player must already have a good command of his/her motor system  $g_1$ .

In particular, the first constraint implies that they must share the same feature space, i.e.,  $M = M'$ . A way to ensure that both constraints are satisfied is to check whether the player can reproduce a task described by the coach, i.e., that

$$g_1(c_1(x)) = x \quad (2)$$

holds. One recognizes in (2) the expression of an autoencoder. It is important to remark that usually, based on previous learning, a player already has a good motor control and he/she and his/her coach know how to communicate together. In other words  $g_1$  and  $c_1$  satisfy (2) before the training starts. Then the training consists only in learning  $g'_0$  and  $c'_0$  on the high level feature interpretations of the possible tennis movements, i.e., on  $c_1(\mathcal{D}')$ .

## 5 Theoretical guarantee for convergence

The following theorem (cf. Appendix A) enables a very precise control of the convergence of the generated distribution towards the true target distribution. It gives an upper bound on the convergence error  $er_{conv}$  in terms of the domain shift  $er_{shift}$ , the autoencoder quality  $er_{AE}$  and the quality of the mindGAN  $er_{mind}$ .

**Theorem 1.** *There exist two positive constants  $a$  and  $b$  such that*

$$er_{conv} \leq a \cdot er_{shift} + er_{AE} + b \cdot er_{mind}. \quad (3)$$

To be more precise, with the notations  $er_{conv} = W(\mathbb{P}_{\mathcal{D}'}, \mathbb{P}'_{\theta})$ ,  $er_{shift} = W(\mathbb{P}_{\mathcal{D}}, \mathbb{P}_{\mathcal{D}'})$ ,  $er_{AE} = W(AE(\mathbb{P}_{\mathcal{D}}), \mathbb{P}_{\mathcal{D}})$  and  $er_{mind} = W(c_{1\#}\mathbb{P}_{\mathcal{D}'}, g'_{0\#}\mathbb{P}_Z)$ .

Very concretely, theorem 1 tells us that in order to control the convergence of the transferred GAN towards the distribution of the target dataset  $\mathcal{D}'$ , we need the exact analogues of steps 1-3 of 2 :

1. choose two datasets  $\mathcal{D}'$  and  $\mathcal{D}$  very similar, i.e.,  $W(\mathbb{P}_{\mathcal{D}}, \mathbb{P}_{\mathcal{D}'})$  small,
2. choose a good autoencoder  $(c_1, g_1)$ , i.e.,  $W(\mathbb{P}_{\mathcal{D}}, AE(\mathbb{P}_{\mathcal{D}}))$  small,
3. train well the MindGAN  $(g_0, c_0)$  on  $c_1(\mathcal{D}')$ , i.e.,  $W(c_{1\#}\mathbb{P}_{\mathcal{D}'}, \mathbb{P}'_{\theta})$  small.

**Remark 2.** *The main theorem of (Patrini et al. 2019), theorem 3.1, guarantees the convergence of a Wasserstein autoencoder (WAE). We show in Appendix A that it is a direct consequence of our theorem.*

## 6 Evaluation

**Datasets.** We have tested our algorithm at resolution  $28 \times 28$  in grey levels (scaled in range  $[-1; 1]$ ) on MNIST (LeCun and Cortes 2010), KMNIST (Clanuwat et al. 2018), FashionMNIST (Xiao, Rasul, and Vollgraf 2017) (60 000 images each) and at resolution  $1024 \times 1024$  in color on CelabAHQ (Karras et al. 2018a) (30 000 images) from models trained on the 60 000 first images of FFHQ (Karras, Laine, and Aila 2019) (which consists of 70 000 images), using the **library** Pytorch. The **hardware** for our experiments with  $28 \times 28$  images consisted of a desktop with 16 Go of RAM and a GPU Nvidia GTX 1080 TI. Most of our experiments with HD color images used a node-gpu (Jean Zay website) with two CPU Intel Cascade Lake 6248 and a GPU Nvidia V100 SXM2 32 Go. We have also benchmarked the running time on entry level GPU GTX 1060. Despite its limitations (Borji 2019), we have used *FID* (Frechet Inception Distance) (Heusel et al. 2017) as **metric**. It is the current standard for evaluations of GANs. Our **code** is available at (Mind2mind github repository).

**At resolution  $28 \times 28$ .** The encoder  $c_1$  has three convolutional layers with instance normalisation (in) and relu : 32+in+relu, 64+in+relu, 128+in+relu, followed by a single dense layer 256+tanh. The decoder  $g_1$  has a single dense layer  $4*4*64 +$  relu and three deconvolutional layers with batch normalisation (bn) : 64+bn+relu, 64+bn+relu, 32+bn+relu, 1+tanh. The MindGAN is a MLP WGAN whose generator  $g_0$  consists in three dense layers : 512+bn+relu, 512+bn+relu, 256+tanh and the critic  $c_0$  consists also in three dense layers : 256+relu, 256+relu, 1. Our **hyper-parameters** : learning rate of  $10^{-3}$ , batch size of 128 for all the networks, 80 epochs for  $(g_1, c_1)$  and 100 for the other networks, gradient penalty with  $\lambda = 10$ , beta parameters in Adam optimizer (.9, .9) for  $(g_1, c_1)$ , (.1, .5) for the other networks.

We report the results with  $\mathcal{D}' = \text{MNIST}$  for  $c_1$  trained on each dataset (see Appendix C for other  $\mathcal{D}'$ ). We compare our results to a Vanilla WGAN with architecture  $(g_1 \circ g_0, c_0 \circ c_1)$ , so that the number of parameters agrees, for fair comparison.

**Baseline 1** Vanilla WGAN with gradient penalty trained on MNIST. We have used for the Vanilla WGAN a model similar to the one used in (Gulrajani et al. 2017). However, this model would not converge properly, due to a problem of “magnitude race.” We have therefore added an  $\epsilon$ -term (Karras et al. 2018a), (Aigner and Körner 2018) to ensure its convergence. Our results appear in the first graph on the *l.h.s* of figure 1, with time in seconds in abscissa. One can observe extremely fast convergence of the mindGAN to good scores, in comparison with the Vanilla WGAN. Note that we have not smoothed any of the curves. This observation suggests that our approach, beyond the gain in training time, provides a regularising mechanism. The stability of the training confirms this hypothesis. Indeed, **statistics over 10 runs**

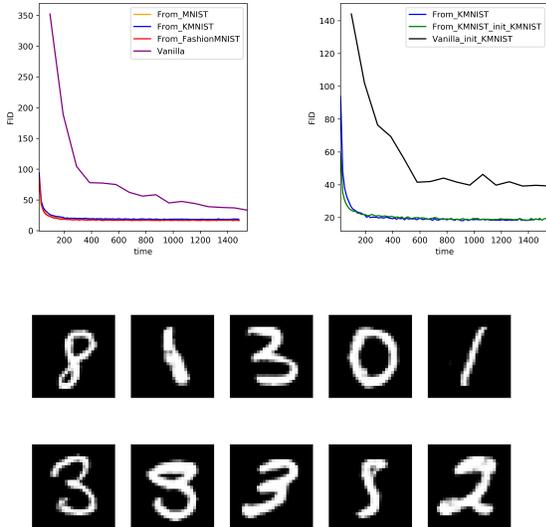


Figure 1: Mind2Mind training and samples in  $28 \times 28$ .

demonstrate a very small standard deviation, as shown in figure 4 in the supplementary material. In particular, this regularization enabled us to use a much bigger learning rate ( $10^{-3}$  instead of  $10^{-4}$ ), adding to the speed of convergence. In terms of epochs, the MindGAN and the Vanilla WGAN learn similarly (cf Appendix C).

**Baseline 2 Fine tuning** studied in (Wang et al. 2018). We have trained a Vanilla WGAN with gradient penalty on  $\mathcal{D} = \text{KMNIST}$ , the dataset the closest to  $\mathcal{D}' = \text{MNIST}$ . We have then *fine-tuned* it on  $\mathcal{D}'$ , i.e., trained a new network on  $\mathcal{D}'$ , initialized with the weights of this previously trained Vanilla WGAN. We display it on the *r.h.s* of figure 1 under the name *Vanilla init Kmnist*, together with our best result, namely a Mind2Mind transfer on  $\mathcal{D}' = \text{MNIST}$  from  $\mathcal{D} = \text{KMNIST}$ . One can observe that the Mind2Mind approach achieves significantly better performances in FID.

The bottom of figure 1 displays samples of images produced by a MindGAN trained on MNIST images encoded using a KMNIST autoencoder.

**At resolution  $1024 \times 1024$ .** We have worked with the encoder and decoder of the ALAE model (Pidhorskyi, Adjeroh, and Doretto 2020) pre-trained on FFHQ available at (ALAE github repository). The generator of our MindGAN has an input dimension of 128, three hidden dense layers with relu activation (128, 256, 512) followed by a dense output layer with 512 units (no activation). The critic has an input dimension of 512, three hidden dense layers with relu activation (512, 256, 128) followed by a dense output layer with one unit (no activation). Its hyperparameters were  $\text{lr} = 1e^{-3}$ ,  $\text{betas} = [0., 0.5]$ ,  $\text{gradient penalty} = 10$ ,  $\text{epsilon penalty} = 1e^{-2}$ ,  $\text{batch size} = 256$ ,  $\text{critic iteration} = 5$ ,  $\text{epochs} = 300$ .

We have encoded the dataset CelebaHQ (Karras et al. 2018b) and then trained a mindGAN on a V100 for 300



Figure 2: Mind2Mind on CelebaHQ transferred from FFHQ.

epochs at 16.95 s/epoch during 1h24m. On a GTX 1060 the same training takes 30.67 s/epoch. We have reached (over 5 runs) an average FID of 15.18 with an average standard deviation of 0.8. This is a better result than the FID score (19.21) of an ALAE directly trained from scratch (see table 5 from (Pidhorskyi, Adjeroh, and Doretto 2020)). Samples are displayed on figure 2. Compared to the results reported on the ProGAN and StyleGAN official repositories (ProGAN github repository), (StyleGAN github repository), our training (1.5 hour) on **1 GPU V100** is roughly **224 times faster** than the training of a **proGAN** (2 weeks) and **656 times faster** than the training of a **StyleGAN** (41 days). The training of a MindGAN on a **GTX 1060** is about **112 times faster** than the training of a proGAN and **328 times faster** than the training of a StyleGAN, both on a V100. Note that a GTX 1060 costs around 200 \$ while a V100 is around 8000 \$. One has to mention however that, on CelebaHQ, the FID of a ProGAN is 8.03 (see table 5 from (Pidhorskyi, Adjeroh, and Doretto 2020)), while the FID of a StyleGAN is 4.40 (see table 4 of (Karras, Laine, and Aila 2019)), so both are significantly better than ours. We see two factors that can explain the acceleration of the training. The first one is that there are much less parameters to train. Indeed, our mindGAN in HD has around 870K parameters, while the ALAE model (based on a StyleGAN architecture) has 51M parameters. So this already represents a difference of almost two orders of magnitude. One can suspect that the rest of the

difference comes from the fact that we bypass 18 layers in the computation of the backpropagation. We believe that the validation of this hypothesis deserves a careful experimental study.

## 7 Comparison to other works

**Wasserstein autoencoders.** WAEs do not provide a solution to the question we address here. Indeed, (Tolstikhin et al. 2017) do not consider at all transfer learning and work only with a single data set  $\mathcal{D}$  at a time. Its goal is rather to give a new approach to Variational Auto Encoders based on the use of the Wasserstein distance.

**Adversarial learned inference.** The works on adversarial learned inference offer an alternative way to train autoencoders. They do not address transfer learning, however, it is possible to choose these types of auto-encoders as building block for our transfer method. We haven't yet conducted experiments with such auto-encoders. In particular, we do not know if theoretical results similar to theorem 1 can be obtained in this setting.

**Adversarial Latent Autoencoders** The architecture of this method is very close to the architecture we use in our work, though the objectives are very different : their point is to disentangle representations in order to be able to control the features, whereas our objective is to do transfer. The similarity in architecture probably explains why ALAE are very suited to our method. We did not have enough time to test whether the disentanglement properties of ALAE are preserved via transfer. We plan to investigate further this question.

**Fine-Tuning.** Our approach is different but can be, in theory, combined with fine-tuning. Indeed, one can initialise the training of our MindGAN  $(g_0, c_0)$  of algorithm 1 on  $c_1(\mathcal{D}')$  with another MindGAN trained on  $c_1(\mathcal{D})$ . We have tried this on the MNISTs datasets; see figure 1, but no significant improvement has been observed. We have compared fine-tuning against Mind2Mind transfer for WGANs and report better results in terms of FID. Note that our theorem 1 gives a theoretical justification, in our setting, of the observation of (Wang et al. 2018) of the influence of the domain shift (distance between  $\mathcal{D}$  and  $\mathcal{D}'$ ) on the convergence of the learning.

## 8 Disadvantages

Our first limitation is that we do not learn the transferred layers during the transfer. On the one hand, it is a feature as it enables faster learning. On the other hand, it is possible that at the asymptote, i.e., after the training has converged, the transferred mindGAN offers a worse quality than a Vanilla WGAN. We have not observed this phenomenon with the ALAE architecture. To the contrary, we obtained a better asymptote since we got a FID of 15.18 for the transferred model compared to the FID of 19.21 of the ALAE model trained from scratch on CelebA HQ. Maybe this is due to the fact that we have focussed the training on a more significant part of the network (the MindGAN).

The second limitation of our approach is that it does not provide improvement of the training on limited data. How-

ever, such a training was not the objective of our work as we believe that acceleration of the training is a useful goal on itself.

Our results, compared to state of the art baselines, show a worse performance in terms of quality (FID of 15.18 compared to 8.03 for ProGAN and 4.40 for StyleGAN). However, this should not be an obstacle to the use of this algorithm. Indeed, the targeted users are practitioners without significant computing capacity. Their need is to reach a reasonable quality in a short time. With this regard, in terms of wall clock time, this algorithm learns roughly 224 times faster than a ProGAN and 656 times faster than a StyleGAN on CelebA HQ. So merits of our algorithm will depend on the tradeoff between the needs of the users in terms of quality and their constraints in terms of computing capacity.

The bottleneck in our approach is the lack of a zoo of models of autoencoders trained on diverse datasets in high resolution. We hope that our method, in conjunction with the ALAE architecture, will lead the main players in the industry to train such models and make them accessible to the community.

## 9 Conclusion

We introduced a method that enables transfer learning for GANs. Given an autoencoder trained on a source dataset, one passes the target dataset through the encoder and uses the encoded features to train a GAN, called a MindGAN, in the latent space of the autoencoder. Composing the MindGAN with the decoder provides the transferred GAN, a generator for the target dataset. We provided theoretical results that guaranty that the transferred GAN converges to the target data. We have demonstrated that our method enables to train GANs much faster (between 6 and 656 times faster, depending on the size of the network) than state of the art methods, in both  $28 \times 28$  gray scale and  $1024 \times 1024$  HQ color datasets.

## Broader impact

The main impact of our work will be a democratization of the use of GANs. Indeed, without the barrier of computing time/costs, GANs for high quality images will become accessible to a pool of practitioners much broader than the researchers/engineers of big industry and academic labs. As a corollary, one can expect a development of research in this field proportional to the increase of number of researchers who will gain access to these tools.

Affordable computing time will offer the possibility to customize models for specific needs and datasets. This can lead to new applications in movie industry, virtual assistants or even telecommunications (it may be more efficient to learn a model of a person, send it and use it to reconstruct a video signal than to transmit the signal itself) to name a few. But this can also lead to malicious uses, in particular for the generation fake profiles and deep fakes that can be used for phishing or disinformation.

## References

- Aigner, S.; and Körner, M. 2018. FutureGAN: Anticipating the Future Frames of Video Sequences using Spatio-Temporal 3d Convolutions in Progressively Growing Autoencoder GANs. *CoRR* abs/1810.01325. URL <http://arxiv.org/abs/1810.01325>.
- ALAE github repository. 2020. URL <https://github.com/podgorskiy/ALAE>.
- Arjovsky, M.; Chintala, S.; and Bottou, L. 2017. Wasserstein Generative Adversarial Networks. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, Australia, 7-9 August, 2017*. URL <http://leon.bottou.org/papers/arjovsky-chintala-bottou-2017>.
- Belghazi, M. I.; Rajeswar, S.; Mastropietro, O.; Rostamzadeh, N.; Mitrovic, J.; and Courville, A. C. 2018. Hierarchical Adversarially Learned Inference. *CoRR* abs/1802.01071.
- Bogachev, V. I. 2007. *Measure Theory*. Springer Verlag.
- Borji, A. 2019. Pros and cons of GAN evaluation measures. *Computer Vision and Image Understanding* 179: 41 – 65. ISSN 1077-3142. doi:<https://doi.org/10.1016/j.cviu.2018.10.009>. URL <http://www.sciencedirect.com/science/article/pii/S1077314218304272>.
- Clanuwat, T.; Bober-Irizar, M.; Kitamoto, A.; Lamb, A.; Yamamoto, K.; and Ha, D. 2018. Deep Learning for Classical Japanese Literature. *CoRR* abs/1812.01718.
- Donahue, J.; Jia, Y.; Vinyals, O.; Hoffman, J.; Zhang, N.; Tzeng, E.; and Darrell, T. 2014. Decaf: A deep convolutional activation feature for generic visual recognition. In *International conference on machine learning*, 647–655.
- Donahue, J.; Krähenbühl, P.; and Darrell, T. 2017. Adversarial Feature Learning. *CoRR* abs/1605.09782.
- Dumoulin, V.; Belghazi, I.; Poole, B.; Lamb, A.; Arjovsky, M.; Mastropietro, O.; and Courville, A. C. 2017. Adversarially Learned Inference. *CoRR* abs/1606.00704.
- Goodfellow, I.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; and Bengio, Y. 2014. Generative Adversarial Nets. In Ghahramani, Z.; Welling, M.; Cortes, C.; Lawrence, N. D.; and Weinberger, K. Q., eds., *Advances in Neural Information Processing Systems 27*, 2672–2680. Curran Associates, Inc. URL <http://papers.nips.cc/paper/5423-generative-adversarial-nets.pdf>.
- Gulrajani, I.; Ahmed, F.; Arjovsky, M.; Dumoulin, V.; and Courville, A. C. 2017. Improved Training of Wasserstein GANs. In *NIPS*.
- Haidar, M. A.; and Rezagholizadeh, M. 2019. TextKD-GAN: Text Generation using KnowledgeDistillation and Generative Adversarial Networks.
- 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 Guyon, I.; Luxburg, U. V.; Bengio, S.; Wallach, H.; Fergus, R.; Vishwanathan, S.; and Garnett, R., eds., *Advances in Neural Information Processing Systems 30*, 6626–6637. Curran Associates, Inc.
- Jean Zay website. 2020. URL <http://www.idris.fr/jean-zay/jean-zay-presentation.html>.
- Karras, T.; Aila, T.; Laine, S.; and Lehtinen, J. 2018a. Progressive Growing of GANs for Improved Quality, Stability, and Variation. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. URL <https://openreview.net/forum?id=Hk99zCeAb>.
- Karras, T.; Aila, T.; Laine, S.; and Lehtinen, J. 2018b. Progressive Growing of GANs for Improved Quality, Stability, and Variation. *ArXiv* abs/1710.10196.
- Karras, T.; Laine, S.; and Aila, T. 2019. A Style-Based Generator Architecture for Generative Adversarial Networks. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 4396–4405.
- LeCun, Y.; and Cortes, C. 2010. MNIST handwritten digit database. <http://yann.lecun.com/exdb/mnist/>. URL <http://yann.lecun.com/exdb/mnist/>.
- Makhzani, A.; Shlens, J.; Jaitly, N.; and Goodfellow, I. J. 2015. Adversarial Autoencoders. *CoRR* abs/1511.05644.
- Mind2mind github repository. 2020. URL <https://github.com/maskedforreview/mindgan>.
- Patrini, G.; Carioni, M.; Forré, P.; Bhargav, S.; Welling, M.; van den Berg, R.; Genewein, T.; and Nielsen, F. 2019. Sinkhorn AutoEncoders. In *UAI*.
- Pidhorskyi, S.; Adjeroh, D. A.; and Doretto, G. 2020. Adversarial Latent Autoencoders. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*. [to appear].
- ProGAN github repository. 2018. URL [https://github.com/tkarras/progressive\\_growing\\_of\\_gans](https://github.com/tkarras/progressive_growing_of_gans).
- Salimans, T.; Goodfellow, I. J.; Zaremba, W.; Cheung, V.; Radford, A.; and Chen, X. 2016. Improved Techniques for Training GANs. In *NIPS*.
- Shan, H.; Zhang, Y.; Yang, Q.; Kruger, U.; Kalra, M. K.; Sun, L.; Cong, W.; and Wang, G. 2018. 3-D Convolutional Encoder-Decoder Network for Low-Dose CT via Transfer Learning From a 2-D Trained Network. *IEEE Transactions on Medical Imaging* 37: 1522–1534.
- StyleGAN github repository. 2019. URL <https://github.com/NVLabs/stylegan>.
- Tan, C.; Sun, F.; Kong, T.; Zhang, W.; Yang, C.; and Liu, C. 2018. A Survey on Deep Transfer Learning. In *ICANN*.
- Tolstikhin, I. O.; Bousquet, O.; Gelly, S.; and Schölkopf, B. 2017. Wasserstein Auto-Encoders. *CoRR* abs/1711.01558.
- Villani, C. 2008. *Optimal transport: old and new*, volume 338. Springer Science & Business Media.
- Wang, Y.; Wu, C.; Herranz, L.; van de Weijer, J.; Gonzalez-Garcia, A.; and Raducanu, B. 2018. Transferring GANs: Generating Images from Limited Data. In *ECCV*.
- Xiao, H.; Rasul, K.; and Vollgraf, R. 2017. Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms. *CoRR* abs/1708.07747.
- Zhao, J. J.; Kim, Y.; Zhang, K.; Rush, A. M.; and LeCun, Y. 2018. Adversarially Regularized Autoencoders. In *ICML*.

## A Proof of things

We first recall notions that will be needed. We then give the proof of our main theoretical result.

### A.1 Wasserstein distance and Lipschitz functions

In the following, all the metric spaces considered will be subsets of normed vector spaces, with the metric on the subset induced by the norm.

First, one recalls some definitions (more details can be found in (Villani 2008)).

**Definition** (transference plan). *Let  $(X, \mathbb{P}_X)$  and  $(Y, \mathbb{P}_Y)$  be two probability spaces. A transference plan  $\gamma$  is a measure on  $X \times Y$  such that :*

$$\int_{A \times Y} d\gamma = \mathbb{P}_X(A),$$

and,

$$\int_{X \times B} d\gamma = \mathbb{P}_Y(B).$$

$\mathbb{P}_X$  and  $\mathbb{P}_Y$  are called the **marginals** of  $\gamma$ . The set of transference plans with marginals  $\mathbb{P}_X$  and  $\mathbb{P}_Y$  is denoted by  $\Pi(\mathbb{P}_X, \mathbb{P}_Y)$ .

**Definition** (p-Wasserstein distance). *Let  $(X, \|\cdot\|)$  be a metric space and  $p \in [1, +\infty)$ . For two probability measure  $\mathbb{P}_1, \mathbb{P}_2$  on  $X$ , the p-Wasserstein distance between  $\mathbb{P}_1$  and  $\mathbb{P}_2$  is defined by the following*

$$W_p(\mathbb{P}_1, \mathbb{P}_2) = \left( \inf_{\gamma \in \Pi(\mathbb{P}_1, \mathbb{P}_2)} \mathbb{E}_{(x,y) \sim \gamma} \|x - y\|^p \right)^{\frac{1}{p}}.$$

In this paper, we used the notation  $W(\mathbb{P}_1, \mathbb{P}_2)$  instead of  $W_1(\mathbb{P}_1, \mathbb{P}_2)$ .

**Definition** (Lipschitz function). *Let  $\phi : X \rightarrow Y$  be a map between metric spaces  $X$  and  $Y$ . It is called a C-Lipschitz function if there exists a constant C such that :*

$$\forall x \text{ and } y \in X, \|\phi(x) - \phi(y)\|_Y \leq C\|x - y\|_X.$$

**Lemma 1.** *Let  $\phi : X \rightarrow Y$  be a locally Lipschitz map, with  $X$  compact, then there exists a constant C such that*

$$W_Y(\phi_{\#}\mu, \phi_{\#}\nu) \leq CW_X(\mu, \nu).$$

*Proof.* Let  $\gamma$  be a transference plan realising  $W_X(\mu, \nu)$ . Define  $\gamma' := (\phi \times \phi)_{\#}\gamma$ . One can check that  $\gamma'$  defines a transference plan between  $\phi_{\#}\mu$  and  $\phi_{\#}\nu$ . Therefore, one has the following relation

$$\begin{aligned} W_Y(\phi_{\#}\mu, \phi_{\#}\nu) &\leq \int \|x - y\| d\gamma'(x, y) \\ &= \int \|\phi(x) - \phi(y)\| d\gamma(x, y) \\ &\leq \int C\|x - y\| d\gamma(x, y) \\ &= CW_X(\mu, \nu), \end{aligned}$$

where the first inequality comes from the fact that  $\gamma'$  is a transference plan, the first equality from the definition of the push forward of a measure by a map (recalled in section 2), the last inequality from lemma 2, and the last equality from the choice of  $\gamma$ .  $\square$

**Lemma 2.** *Let  $\phi : X \rightarrow Y$  be a locally Lipschitz map, and  $X$  a compact metric space. Then there exists C such that  $\phi$  is a C-Lipschitz function.*

*Proof.* By definition of a locally Lipschitz map, for all  $x$  in  $X$ , there exists  $U_x$  a neighbourhood of  $x$  and a constant  $C_x$  such that  $\phi$  is  $C_x$ -Lipschitz on  $U_x$ .

So  $\bigcup_{x \in X} U_x$  is a cover of  $X$ . Since  $X$  is compact, there exists a finite set  $I$  such that  $\bigcup_{i \in I} U_i$  is a cover of  $X$ .

One can check that  $\phi$  is C-Lipschitz on  $X$ , with  $C := \max_{i \in I} (C_{x_i})$ .  $\square$

### A.2 Proof of theorem 1

We can finally turn to the proof of theorem 1 :

**Theorem.** *There exist two positive constants a and b such that*

$$\begin{aligned} W(\mathbb{P}_{\mathcal{D}'}, \mathbb{P}'_{\theta}) &\leq a W(\mathbb{P}_{\mathcal{D}}, \mathbb{P}_{\mathcal{D}'}) + W(\mathbb{P}_{\mathcal{D}}, AE(\mathbb{P}_{\mathcal{D}})) \\ &\quad + b W(c_{1\#}\mathbb{P}_{\mathcal{D}'}, \mathbb{P}'_{\theta}). \end{aligned}$$

*Proof.* From the triangle inequality property of the Wasserstein metric and the definition of  $\mathbb{P}'_{\theta}$ , one has :

$$W(\mathbb{P}_{\mathcal{D}'}, \mathbb{P}'_{\theta}) \leq W(\mathbb{P}_{\mathcal{D}'}, AE(\mathbb{P}_{\mathcal{D}'})) + W(AE(\mathbb{P}_{\mathcal{D}'}), g_{1\#}\mathbb{P}'_{\theta}).$$

One concludes with lemma 3 and lemma 1 with  $\phi = g_1$ .  $\square$

**Lemma 3.** *There exist a positive constant a such that*

$$W(\mathbb{P}_{\mathcal{D}'}, AE(\mathbb{P}_{\mathcal{D}'})) \leq aW(\mathbb{P}_{\mathcal{D}}, \mathbb{P}_{\mathcal{D}'}) + W(\mathbb{P}_{\mathcal{D}}, AE(\mathbb{P}_{\mathcal{D}})).$$

*Proof.* Applying twice the triangle inequality, one has :

$$\begin{aligned} W(\mathbb{P}_{\mathcal{D}'}, AE(\mathbb{P}_{\mathcal{D}'})) &\leq W(\mathbb{P}_{\mathcal{D}'}, \mathbb{P}_{\mathcal{D}}) + W(\mathbb{P}_{\mathcal{D}}, AE(\mathbb{P}_{\mathcal{D}})) \\ &\quad + W(AE(\mathbb{P}_{\mathcal{D}}), AE(\mathbb{P}_{\mathcal{D}'})). \end{aligned}$$

One concludes with lemma 1 with  $\phi = g_1 \circ c_1$ .  $\square$

**Remark 3.** *It is important to remark that in order to be able to apply lemma 1 in the proof of theorem 1, one needs the assumption that  $\mathbb{P}'_{\theta}$  and  $c_{1\#}\mathbb{P}_{\mathcal{D}'}$  have compact support. But as  $\chi$  is itself compact, this is not a problem for  $c_{1\#}\mathbb{P}_{\mathcal{D}'}$  since the image of a compact  $\chi$  by a continuous function  $c_1$  is compact. However, the compactity of the support of  $\mathbb{P}'_{\theta}$  is not a priori granted. An easy fix is to choose a prior  $\mathbb{P}_Z$  with compact support. Therefore, we choose this setting in our applications.*

**Remark 4.** *Our proof of theorem 1 implicitly assumed that neural networks are locally Lipschitz maps (see lemmata 1 and 3). This assumption is justified by the following lemma.*

**Lemma.** *Let  $g : Z \rightarrow X$  be a neural network and  $\mathbb{P}_Z$  a prior over  $Z$  such that  $\mathbb{E}_{z \sim \mathbb{P}_Z} (\|z\|) < \infty$  (such as Gaussian) then  $g$  is locally Lipschitz and  $\mathbb{E}_{z \sim \mathbb{P}_Z} (L_z) < \infty$ , where  $L_z$  are the local Lipschitz constants.*

*Proof.* See Corollary 1. of (Arjovsky, Chintala, and Bottou 2017)  $\square$

### A.3 Application to Wasserstein autoencoders.

The main theorem of (Patrini et al. 2019), theorem 3.1, guarantees the convergence of a Wasserstein autoencoder (WAE). We recall this theorem and show that it is a direct consequence of our theorem 1.

#### Theorem 2.

$$W(\mathbb{P}_{\mathcal{D}}, g_{1\#}\mathbb{P}_{\mathcal{Z}}) \leq W(\mathbb{P}_{\mathcal{D}}, AE(\mathbb{P}_{\mathcal{D}})) \quad (4)$$

$$+ bW(c_{1\#}\mathbb{P}_{\mathcal{D}}, \mathbb{P}_{\mathcal{Z}}). \quad (5)$$

*Proof.* Since WAEs do not involve transfer, one has  $\mathcal{D} = \mathcal{D}'$ , i.e.  $\mathbb{P}_{\mathcal{D}} = \mathbb{P}_{\mathcal{D}'}$ . Then it suffices to replace  $\mathbb{P}_{\theta}^{g_0}$  by  $\mathbb{P}_{\mathcal{Z}}$  and  $\mathbb{P}'_{\theta}$  becomes  $g_{1\#}\mathbb{P}_{\mathcal{Z}}$ .  $\square$

**Remark 5.** *Our proof of theorem 1 is very similar to the proof of theorem 2 given in (Patrini et al. 2019). Therefore, our contribution here consists rather in finding a versatile statement that applies to both problems (transfer and WAE) than in the originality of the tools used in the proofs.*

**Remark 6.** *When one restricts our approach to the case when  $\mathcal{D} = \mathcal{D}'$ , it does not coincide with WAE. Indeed, with the notations of our paper, WAE work with a fixed prior  $\mathbb{P}_M$  on  $M$  that one tries to approximate by  $c_{1\#}\mathbb{P}_{\mathcal{D}}$ , while constraining  $c_1$  to be a right inverse (in measure) of  $g_1$ , and  $g_{1\#}\mathbb{P}_M$  to approximate (in measure)  $\mathbb{P}_{\mathcal{D}}$ . On the other hand, our approach involves an extra auxiliary latent space  $Z$ . Therefore we can consider  $g_{0\#}\mathbb{P}_Z$  as a replacement of  $\mathbb{P}_M$ . Via the flexibility of the learnable weights of  $g_0$ , we use  $g_{0\#}\mathbb{P}_Z$  to approximate  $c_{1\#}\mathbb{P}_{\mathcal{D}}$ , instead of using  $c_{1\#}\mathbb{P}_{\mathcal{D}}$  to approximate  $\mathbb{P}_M$  as in (Makhzani et al. 2015). This is fundamental, because in a setting where  $\mathcal{D} \neq \mathcal{D}'$ , this decoupling permits to train  $c_1$  and  $g_1$  on  $\mathcal{D}$  and  $c_0$  and  $g_0$  on  $c_1(\mathcal{D}')$ , enabling us to do transfer.*

## B Mind2Mind conditional GANs

As suggested to us by L. Cetinsoy, the Mind2Mind approach also applies to conditional GANs. However, one needs to implement the following modifications : replace  $M$  by  $M \times L$  and  $Z$  by  $Z \times L$  in the diagram

$$\begin{array}{ccccc} & M & & M & \\ g_0 \nearrow & & g_1 \searrow & c_1 \nearrow & c_0 \searrow \\ Z & \xrightarrow{g} & \chi & \xrightarrow{c} & \mathbb{R}, \end{array} \quad (6)$$

where  $L$  stands for the space of conditions, in order to get

$$\begin{array}{ccccc} & M \times L & & M \times L & \\ g_0^c \nearrow & & g_1 \times \mathbb{I}_L \searrow & c_1 \times \mathbb{I}_L \nearrow & c_0^c \searrow \\ Z \times L & \xrightarrow{g^c} & \chi \times L & \xrightarrow{c^c} & \mathbb{R}. \end{array} \quad (7)$$

Here,  $(g_0^c, c_0^c)$  and  $(g^c, c^c)$  are conditional GANs, with the generators of the form  $g_0^c(z, l) = (m(z, l), l)$  and  $g^c(z, l) = (x(z, l), l)$ . The autoencoder  $(c_1 \times \mathbb{I}_L, g_1 \times \mathbb{I}_L)$  can be trivially deduced from an autoencoder  $(c_1, g_1)$  via the formulas  $c_1 \times \mathbb{I}_L(x, l) := (c_1(x), l)$  and  $g_1 \times \mathbb{I}_L(m, l) := (c_1(m), l)$ .

In practice, the algorithm 1 becomes a classical conditional GAN algorithm :

---

#### Algorithm 2 Conditional-MindGAN transfer learning.

---

**Require:**  $(c_1, g_1)$ , an autoencoder trained on a source dataset  $\mathcal{D}$ ,  $\alpha$ , the learning rate,  $b$ , the batch size,  $n$ , the number of iterations of the critic per generator iteration,  $\mathcal{D}' \subset \chi \times L$ , a dataset with conditions,  $\varphi'$  and  $\theta'$  the initial parameters of the critic  $c_0^c$  and of the generator  $g_0^c$ .  
 Compute  $(c_1 \times \mathbb{I}_L)(\mathcal{D}')$ .  
**while**  $\theta'$  has not converged **do**  
   **for**  $t = 0, \dots, n_{\text{critic}}$  **do**  
     Sample  $\{(m^{(i)}, l^{(i)})\}_{i=1}^b \sim (c_1 \times \mathbb{I}_L)_{\#}\mathbb{P}_{\mathcal{D}'}$  a batch from  $(c_1 \times \mathbb{I}_L)(\mathcal{D}')$ .  
     Sample  $\{(z^{(i)}, l^{(i)})\}_{i=1}^b \sim \mathbb{P}_{Z \times L}$  a batch of prior samples with conditions.  
     Update  $c_0^c$  by descending  $L_c$ .  
   **end for**  
   Sample  $\{(z^{(i)}, l^{(i)})\}_{i=1}^b \sim \mathbb{P}_{Z \times L}$  a batch of prior samples with conditions.  
   Update  $g_0^c$  by descending  $-L_g$ .  
**end while**  
**return**  $g_1 \circ g_0^c$ .

---

## C Supplementary experiments

In figure 3 we display additional samples from a MindGAN on CelebA HQ transferred from FFHQ. We also display in figure 4 the mean and standard deviation over 10 runs of the training of a MindGAN in  $28 \times 28$ . The source dataset in figure 5 is  $\mathcal{D}' = \text{FashionMNIST}$ , while in figure 6,  $\mathcal{D}' = \text{MNIST}$ . For comparison, we display in figure 7 samples from a vanilla WGAN.



Figure 3: Mind2Mind on CelebaHQ transferred from FFHQ

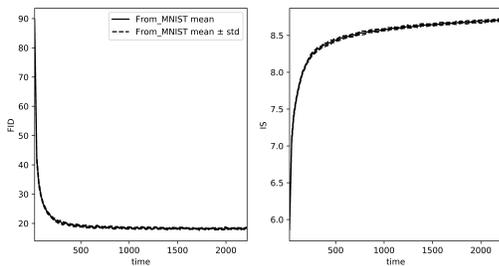


Figure 4: Mean and standard deviation of the training of a MindGAN.



Figure 5: MindGAN from FashionMNIST.



Figure 6: MindGAN from MNIST.



Figure 7: Vanilla WGAN.