
Learning Deep-Latent Hierarchies by Stacking Wasserstein Autoencoders

Anonymous Authors¹

Abstract

Probabilistic models with hierarchical-latent-variable structures provide state-of-the-art results amongst non-autoregressive, unsupervised density-based models. However, the most common approach to training such models based on Variational Autoencoders (VAEs) often fails to leverage deep-latent hierarchies; successful approaches require complex inference and optimisation schemes. Optimal Transport is an alternative, non-likelihood-based framework for training generative models with appealing theoretical properties, in principle allowing easier training convergence between distributions. In this work we propose a novel approach to training models with deep-latent hierarchies based on Optimal Transport, without the need for highly bespoke models and inference networks. We show that our method enables the generative model to fully leverage its deep-latent hierarchy, avoiding the well known "latent variable collapse" issue of VAEs; therefore, providing qualitatively better sample generations as well as more interpretable latent representation than the original Wasserstein Autoencoder with Maximum Mean Discrepancy divergence.

1. Introduction

Probabilistic latent-variable modelling is widely applicable in machine learning as a method for discovering structure directly from large, unlabelled datasets. Variational Autoencoders (VAEs) (Kingma & Welling, 2014; Rezende et al., 2014) have proven to be effective for training generative models parametrised by powerful neural networks, by mapping the data into a low-dimensional embedding space. While allowing the training of expressive models, VAEs often fail when using deeper hierarchies with several stochastic latent layers.

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

In fact, many of the most successful probabilistic latent-variable models use only one stochastic latent layer. Auto-regressive models (Larochelle & Murray, 2011; Van Den Oord et al., 2016a;b; Salimans et al., 2017), for example, produce state-of-the-art samples and likelihood scores. However, auto-regressive models suffer from poor scalability to high-dimensional data. Flow-based models (Rezende & Mohamed, 2015; Kingma et al., 2016; Dinh et al., 2016; Kingma & Dhariwal, 2018) are another class of generative models providing competitive sample quality and are able to scale to higher-dimensional data, but they still lag behind auto-regressive models in terms of likelihood scores.

Deep-latent-variable models are highly expressive models that aim to capture the structure of data in a hierarchical manner, and could thus potentially compete with auto-regressive models for state-of-the-art performance. However, they remain hard to train. Many explanations have been proposed for this, from the use of dissimilarity measure directly in the pixel space (Larsen et al., 2016) resulting in poor sample quality, to the lack of efficient representation in the latent (Zhao et al., 2017), to simply the poor expressiveness of the models used (Zha et al., 2017; Maaløe et al., 2019).

Solutions for training deep-latent-variable models range from introducing auxiliary variables to increase the expressiveness of the posterior distribution (Maaløe et al., 2016) to replacing the simple spherical Gaussian prior with richer prior distributions which can be learned jointly with the generative model (Tomczak & Welling, 2018; Klushyn et al., 2019). Other approaches focus on building and training complex generative models and inference networks introducing latent skip connections in the generative model and inference network (Bachman, 2016), sharing generative model and inference network parameters (Sønderby et al., 2016), as well as bidirectional inference networks (Maaløe et al., 2019). Maaløe et al. (2019) managed to train very deep-hierarchical-latent models achieving near state-of-the-art sample generations, both in term of likelihood score and sample quality. However, in order to leverage their latent hierarchy (working in the VAE framework), they needed both complex, tailored inference networks, and deterministic skip connections in the generative model.

Optimal Transport (OT) (Villani, 2008) is a mathematical framework to compute distances between distributions. Re-

055 recently, it has been used as a training method for generative
 056 models (Genevay et al., 2018; Bousquet et al., 2017). Tol-
 057 stikhin et al. (2018) introduced Wasserstein Autoencoders
 058 (WAEs), where as with VAEs, an encoding distribution
 059 maps the data into a latent space, aiming to learn a low-
 060 dimensional representation of the data. WAEs provide a
 061 non-likelihood based framework with appealing topological
 062 properties (Arjovsky et al., 2017; Bousquet et al., 2017),
 063 in theory allowing for easier training convergence between
 064 distributions. Gaujac et al. (2018) trained a two-latent-layer
 065 generative model using a WAE, showing promising results
 066 in the capacity of the WAE framework to leverage a latent
 067 hierarchy relative to the equivalent VAE.
 068

069 Following these early successes, we propose to train deep-
 070 hierarchical-latent models using the WAE framework, with-
 071 out the need for complex dependency paths in both the
 072 generative model and inference network. As in the works of
 073 Sønderby et al. (2016) and Maaløe et al. (2019), we believe
 074 that a deep-latent hierarchy offers the potential to improve
 075 generative models, if they could be trained properly. In order
 076 to leverage the deep-latent hierarchy, we derive a novel
 077 objective function by stacking WAEs, introducing an infer-
 078 ence network at each level, and encoding the information up
 079 to the deepest layer in a bottom-up way. For convenience,
 080 we refer to our method as STACKEDWAE.
 081

082 Our contributions are two-fold: first, we introduce
 083 STACKEDWAE, a novel objective function based on OT,
 084 designed specifically for generative modelling with latent
 085 hierarchies, and show that it is able to fully leverage its hier-
 086 archy of latents. Second, we show that STACKEDWAE per-
 087 forms significantly better when training hierarchical-latent
 088 models than the original WAE framework regularising the
 089 latent with the Maximum Mean Discrepancy (MMD) diver-
 090 gence (Gretton et al., 2012).
 091

2. Stacked WAE

093 STACKEDWAEs are probabilistic-latent-variable models
 094 with a deep hierarchy of latents. They can be minimalist-
 095 ically simple in their inference and generative models, but
 096 are trained using OT in a novel way. We start by defining the
 097 probabilistic models considered in this work, then introduce
 098 OT, and finally discuss how to train probabilistic models
 099 with deep-latent hierarchies using OT, the method that we
 100 refer to as STACKEDWAE.
 101

102 Throughout this paper, we use calligraphic letters (e.g. \mathcal{X})
 103 for sets, capital letters (e.g. X) for random variables, and
 104 lower case letters (e.g. x) for their realisation. We denote
 105 probability distributions with capital letters (e.g. $P(X)$) and
 106 their densities with lower case letters (e.g. $p(x)$).
 107

2.1. Generative models with deep latent hierarchies

We will consider deep-generative models with Markovian
 hierarchies in their latent variables. Namely, where each
 latent variable depends exclusively on the previous one.
 Denoting by P_Θ the parametric model with N latent layers,
 where $\Theta = \{\theta_1, \dots, \theta_N\}$, we have:

$$p_\Theta(x, z_{1:N}) = p_{\theta_1}(x|z_1) \left[\prod_{n=2}^N p_{\theta_n}(z_{n-1}|z_n) \right] p_{\theta_N}(z_N) \quad (1)$$

where the data is $X \in \mathcal{X}$ and the latent variables are
 $Z_n \in \mathcal{Z}_n$ and we chose $p_{\theta_N}(z_N) = \mathcal{N}(z_N; 0_{\mathcal{Z}_N}, I_{\mathcal{Z}_N})$. The corresponding graphical model for $N = 3$ is given
 Figure 1a.

We will be using variational inference through the WAE
 framework of Tolstikhin et al. (2018), introducing vari-
 ational distributions, $q_\Phi(z_1, \dots, z_N|x)$, to approximate the
 intractable posterior, analogous to VAEs. It will be shown in
 Section 2.3 that without loss of generality, $q_\Phi(z_1, \dots, z_N|x)$
 can have a Markovian latent hierarchy when following the
 STACKEDWAE approach. That is, without loss of general-
 ity,

$$q_\Phi(z_{1:N}|x) = q_{\phi_1}(z_1|x) \prod_{n=2}^N q_{\phi_n}(z_n|z_{n-1}) \quad (2)$$

where each $q_{\phi_n}(z_n|z_{n-1})$ is introduced iteratively by stack-
 ing WAEs at each latent layer. The corresponding graphical
 model for $N = 3$ is given Figure 1b.

We focus on this simple Markovian latent-variable structure
 for the generative model as a proof point for STACKED-
 WAE. This simple modelling setup is famously difficult to
 train, as is discussed extensively in the VAE framework (see
 for example Burda et al. (2015); Sønderby et al. (2016);
 Zhao et al. (2017)). The difficulty in training such models
 comes from the Markovian latent structure of the genera-
 tive model; in particular, the difficulty of learning struc-
 ture in the deeper latents. This is because, to generate
 samples $x \sim p_\Theta(x)$, only the joint $p_\Theta(x, z_1)$ is needed
 as $p_\Theta(x) = \int_{\mathcal{Z}_1} p_{\theta_1}(x|z_1)p_\Theta(z_1)dz_1$. Learning a smooth
 structure in each latent layer is not a strict requirement for
 learning a good generative model, however, it is sought after
 if the latent is to be used downstream or interpreted. We find
 empirically (see Section 3.1) that a better generative model
 is also achieved when the latent hierarchy is well learnt all
 the way down.

Sønderby et al. (2016) sought to overcome the difficulty
 of learning a deep-latent hierarchy by using determinis-
 tic bottom-up inference, followed by top-down inference
 that shared parameters with the generative model. With
 additional optimisation tricks (e.g. KL annealing), their

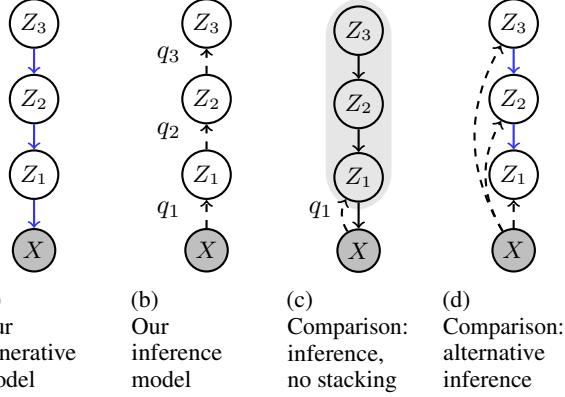


Figure 1: (a) Generative model (blue lines represent generative model parameters). (b) Inference model used in STACKEDWAE. (c) Standard WAE; the generative model has only one latent with prior $p(z_1) = \int p(z_1|z_2)p(z_2|z_3)p(z_3)dz_1dz_2dz_3$. (d) Inference model with skips connections and parameter sharing with the generative model, as in Sønderby et al. (2016).

deeper-latent distributions would still go unused for sufficiently deep-latent hierarchies (discussed in Maaløe et al. (2019)). In order to get deeper hierarchies of latents, Maaløe et al. (2019) introduced additional deterministic connections in the generative model as well as bidirectional inference network to facilitate the deep information flow needed to ensure the usage of the deeper-latent layers.

While the approach in Maaløe et al. (2019) is well motivated and achieves excellent results, we choose the OT framework for training deep-latent models due to its topological properties (see Arjovsky et al. (2017)). Still, the standard WAE encounters the same difficulties as the VAE in learning deep-latent hierarchies. We thus modify the original WAE objective, effectively stacking WAEs, to improve the learning of both the generative model and the inference distribution throughout the latent hierarchy.

2.2. Wasserstein Autoencoders

The Kantorovich formulation of the OT between the true-but-unknown data distribution P_D and the model distribution P_Θ , for a given cost function c , is defined by:

$$\text{OT}_c(P_D, P_\Theta) = \inf_{\Gamma \in \mathcal{P}(P_D, P_\Theta)} \int_{\mathcal{X} \times \mathcal{X}} c(x, \tilde{x}) d\Gamma(x, \tilde{x}) \quad (3)$$

where $\mathcal{P}(P_D, P_\Theta)$ is the space of all couplings of P_D and P_Θ ; namely, the space of joint distributions Γ on $\mathcal{X} \times \mathcal{X}$ whose densities γ have marginals p_D and p_Θ :

$$\begin{aligned} \mathcal{P}(P_D, P_\Theta) &= \left\{ \Gamma \mid \int_{\mathcal{X}} \gamma(x, \tilde{x}) d\tilde{x} = p_D(x), \int_{\mathcal{X}} \gamma(x, \tilde{x}) dx = p_\Theta(\tilde{x}) \right\} \\ &= \left\{ \Gamma \mid \int_{\mathcal{X}} \gamma(x, \tilde{x}) d\tilde{x} = p_D(x), \int_{\mathcal{X}} \gamma(x, \tilde{x}) dx = p_\Theta(\tilde{x}) \right\} \end{aligned} \quad (4)$$

In the WAE framework, the space of couplings is constrained to joint distributions Γ , with density γ , of the form:

$$\gamma(x, \tilde{x}) = \int_{\mathcal{Z}_{1:N}} p_\Theta(\tilde{x}|z_{1:N}) q_{\phi_1}(z_{1:N}|x) p_D(x) dz_{1:N} \quad (5)$$

where $q_{\phi_1}(z_{1:N}|x)$, for $x \in \mathcal{X}$, plays the same role as the variational distribution in variational inference.

Marginalising over \tilde{x} in Eq. (5) automatically gives $p_D(x)$, thus satisfying the first marginal constraint of Eq. (4). However the second marginal constraint in Eq. (4) (that over x giving p_Θ) is not guaranteed. Due to the Markovian structure of the generative model, a sufficient condition for satisfying the second marginal constraint in Eq. (4) for all $z_1 \in \mathcal{Z}_1$, is (more detail provided in Appendix A.1):

$$\int_{\mathcal{X}} q_{\phi_1}(z_1|x) p_D(x) dx = p_{\Theta_{2:N}}(z_1) \quad (6)$$

where we denote $p_{\Theta_{2:N}}(z_1)$ the prior over Z_1 :

$$p_{\Theta_{2:N}}(z_1) \stackrel{\text{def}}{=} \int_{\mathcal{Z}_{2:N}} \left[\prod_{n=2}^N p_{\theta_n}(z_{n-1}|z_n) \right] p_{\theta_N}(z_N) dz_{2:N} \quad (7)$$

Finally, to get the WAE objective of Tolstikhin et al. (2018), the constraint in Eq. (6) is relaxed using a Lagrange multiplier (more detail provided in Appendix A.2):

$$\begin{aligned} \widehat{W}_c(P_D, P_\Theta) &= \inf_{Q_{\phi_1}(Z_1|X=x)} \left[\int_{\mathcal{X} \times \mathcal{X}} c(x, \tilde{x}) \gamma(x, \tilde{x}) dx d\tilde{x} \right. \\ &\quad \left. + \lambda_1 \mathcal{D}_1(Q_1^{\text{agg}}(Z_1), P_{\Theta_{2:N}}(Z_1)) \right] \end{aligned} \quad (8)$$

where \mathcal{D}_1 is any divergence function, λ_1 a relaxation parameter, γ is defined in Eq. (5), and $Q_1^{\text{agg}}(Z_1)$ is the aggregated posterior:

$$Q_1^{\text{agg}}(Z_1) \stackrel{\text{def}}{=} \int_{\mathcal{X}} Q_{\phi_1}(Z_1|x) p_D(x) dx \quad (9)$$

Note that because of the Markovian latent structure of the generative model,

$$\begin{aligned} \gamma(x, \tilde{x}) &= \int_{\mathcal{Z}_{1:N}} p_\Theta(\tilde{x}|z_{1:N}) q_{\phi_1}(z_{1:N}|x) p_D(x) dz_{1:N} \\ &= \int_{\mathcal{Z}_1} p_{\theta_1}(\tilde{x}|z_1) q_{\phi_1}(z_1|x) p_D(x) dz_1 \end{aligned} \quad (10)$$

Eq. (9) and (10) do not depend on $z_{>1}$, so the infimum in Eq. (8) is taken only over $q_{\phi_1}(z_1|x)$ instead of the full $q_{\phi_1}(z_1, \dots, z_N|x)$.

While Eq. (8) is in-principle tractable (e.g. for Gaussian $q_{\phi_1}(z_1|x)$ and sample-based divergence function such as

MMD), it only depends on the first latent Z_1 . Thus it will learn only a good approximation for the joint $p_\Theta(x, z_1) = p_{\theta_1}(x|z_1)p(z_1)$, rather than a full hierarchy with each latent living on a smooth manifold. We show empirically in Section 3.1 that Eq. (8) is indeed insufficient for learning to use the full hierarchy of latents.

2.3. Stacking WAEs for deep latent-variable modelling

In the limit $\lambda_1 \rightarrow \infty$, Eq. (8) does not depend on the choice of divergence \mathcal{D}_1 . However, given the set of approximations used, a divergence that takes into account the smoothness of the full stack of latents will likely help the optimisation. We now explain how, by using the Wasserstein distance itself for \mathcal{D}_1 , we can derive an objective that naturally pairs up inference and generation at every level in the deep-latent hierarchy. After all, the divergence in Eq. (8) is between the intractable aggregate distribution $Q_1^{\text{agg}}(Z_1)$ from which we can only access samples, and an analytically-known distribution $P_{\Theta_{2:N}}(Z_1)$, which is analogous to where we started with the Wasserstein distance between P_D and the full P_Θ .

Specifically, we choose \mathcal{D}_1 in Eq. (8) to be the relaxed Wasserstein distance \tilde{W}_{c_1} , which following the same arguments as before, requires the introduction of a new variational distribution $Q_{\phi_2}(Z_2|Z_1)$:

$$\begin{aligned} \mathcal{D}_1(Q_1^{\text{agg}}(Z_1), P_{\Theta_{2:N}}(Z_1)) &= \\ &\inf_{Q_{\phi_2}(Z_2|Z_1=z_1)} \left[\lambda_2 \mathcal{D}_2(Q_2^{\text{agg}}(Z_2), P_{\Theta_{3:N}}(Z_2)) \right. \\ &\quad \left. + \int_{\mathcal{Z}_1 \times \mathcal{Z}_1} c_1(z_1, \tilde{z}_1) \gamma_1(z_1, \tilde{z}_1) dz_1 d\tilde{z}_1 \right] \end{aligned} \quad (11)$$

where we re-used the notation of Eq. (6) for the prior $P_{\Theta_{3:N}}$, and similarly with Eq. (9) for the aggregated posterior,

$$Q_2^{\text{agg}}(Z_2) \stackrel{\text{def}}{=} \int_{\mathcal{Z}_1} Q_{\phi_2}(Z_2|z_1) q_1^{\text{agg}}(z_1) dz_1 \quad (12)$$

The joint is given by

$$\gamma_1(z_1, \tilde{z}_1) \stackrel{\text{def}}{=} \int_{\mathcal{Z}_2} p_{\theta_2}(\tilde{z}_1|z_2) q_{\phi_2}(z_2|z_1) q_1^{\text{agg}}(Z_1) dz_2 \quad (13)$$

And as before, $Q_{\phi_2}(Z_2|Z_1)$ in the inf does not need to provide a distribution over the $z_{>2}$ without loss of generality.

The divergence \mathcal{D}_1 that arose in Eq. (8) between two distributions over Z_1 is thus mapped onto the latent at the next level in the latent hierarchy, Z_2 , via Eq. (11). This process can be repeated by using \tilde{W}_{c_2} again for \mathcal{D}_2 in Eq. (11) to get an expression that maps to Z_3 , requiring the introduction of another variational distribution $Q_{\phi_3}(Z_3|Z_2)$. Repeating this process until we get to the final layer of the hierarchical-

latent-variable model gives the STACKEDWAE objective:

$$\begin{aligned} W_{\text{STACKEDWAE}}(P_D, P_\Theta) &= \\ &\inf_{Q_\Phi(Z_1, \dots, Z_N|X=x)} \left[\left[\prod_{i=1}^N \lambda_i \right] \mathcal{D}_N(Q_N^{\text{agg}}(Z_N), P(Z_N)) \right. \\ &\quad \left. + \sum_{n=0}^{N-1} \left[\prod_{i=1}^n \lambda_i \right] \int_{\mathcal{Z}_n \times \mathcal{Z}_n} c_n(z_n, \tilde{z}_n) \gamma_n(z_n, \tilde{z}_n) dz_n d\tilde{z}_n \right] \end{aligned} \quad (14)$$

where we denote $(\mathcal{Z}_0, Z_0, z_0) = (\mathcal{X}, X, x)$ and we define the empty product $\prod_{i=1}^0 \lambda_i \stackrel{\text{def}}{=} 1$. Similarly to Eq. (11), we defined the joints γ_n as

$$\begin{aligned} \gamma_n(z_n, \tilde{z}_n) &\stackrel{\text{def}}{=} \\ &\int_{\mathcal{Z}_{n+1}} p_{\theta_{n+1}}(\tilde{z}_n|z_{n+1}) q_{\phi_{n+1}}(z_{n+1}|z_n) q_n^{\text{agg}}(z_n) dz_{n+1} \end{aligned} \quad (15)$$

Each p_{θ_n} is the n^{th} layer of the generative model given in Eq. (1). The q_{ϕ_n} 's are the inference models introduced each time a WAE is "stacked", which combine to make the overall STACKEDWAE Markovian inference model given in Eq. (2), and the aggregated posterior distributions are defined as

$$Q_n^{\text{agg}}(Z_n) \stackrel{\text{def}}{=} \int_{\mathcal{Z}_{n-1}} Q_{\phi_n}(Z_n|z_{n-1}) q_{n-1}^{\text{agg}}(z_{n-1}) dz_{n-1} \quad (16)$$

with $Q_0^{\text{agg}} \stackrel{\text{def}}{=} P_D$.

Note that the STACKEDWAE objective function in Eq. (14) still requires the specification of a divergence function at the highest latent layer \mathcal{D}_N , which we simply take to be the MMD as originally proposed by Tolstikhin et al. (2018). Other choices can be made, as in Patrini et al. (2018), who choose a Wasserstein distance computed using the Sinkhorn algorithm (Cuturi, 2013). While Patrini et al. (2018) provide a theoretical justification for the minimisation of a Wasserstein distance in the prior space, we found that it did not result in significant improvement and comes at an extra efficiency cost. Similarly, one could choose different cost functions c_n at each layer; for simplicity we take all cost functions to be the squared Euclidean distance in their respective spaces.

3. Experiments

We now show how the STACKEDWAE approach of Section 2 can be used to train deep-latent hierarchies without customizing the generative or inference models (e.g. no skip connections, no parameter sharing). We also show through explicit comparison that the STACKEDWAE approach performs significantly better than the standard WAE when training deep-hierarchical-latent models.

220 **3.1. MNIST**

221 EXPERIMENTAL SETUP

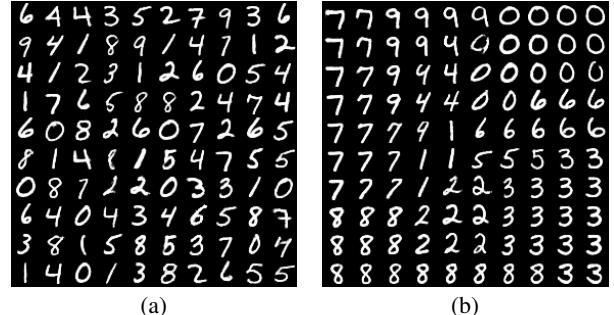
222 We trained a deep-hierarchical-latent variable model with
 223 $N = 5$ latent layers on raw (non-binarised) MNIST (Le-
 224 Cun & Cortes, 2010). The latent layers have dimensions:
 225 $d_{\mathcal{Z}_1} = 32$, $d_{\mathcal{Z}_2} = 16$, $d_{\mathcal{Z}_3} = 8$, $d_{\mathcal{Z}_4} = 4$ and $d_{\mathcal{Z}_5} = 2$.
 226 We chose Gaussian distributions for both the generative and
 227 inference models, except for the bottom layer of the gen-
 228 erative model, which we choose to be deterministic, as in
 229 Tolstikhin et al. (2018). The mean and covariance matrices
 230 used are parametrised by fully connected neural networks
 231 similar to that of Sønderby et al. (2016). Full details of the
 232 experimental setup is given in Appendix B.1.

233 LEARNING A DEEP-LATENT HIERARCHY

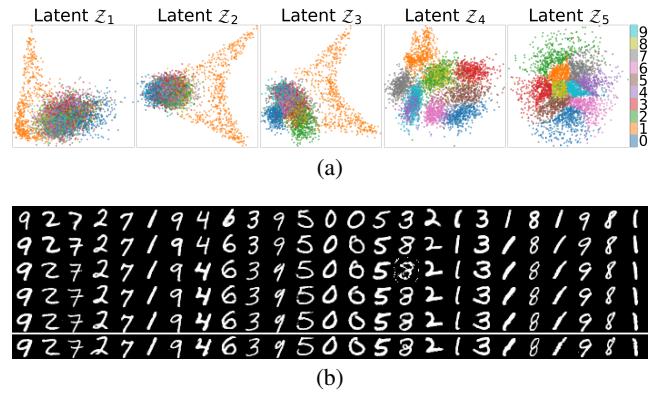
234 Our results are shown in Figure 2, with samples from the
 235 generative model in Figure 2a, and latent space interpo-
 236 lations in Figure 2b where digits are reconstructed from
 237 points in \mathcal{Z}_5 taken evenly in a grid varying between ± 2
 238 standard deviations from the origin. STACKEDWAE gen-
 239 erates compelling samples and learns a smooth manifold (see
 240 Figure 7, top-left, and Figure 8b of Appendix B.1 for addi-
 241 tional samples and latent interpolations respectively). Note
 242 that the choice $d_{\mathcal{Z}_5} = 2$ allows for easy visualisation of the
 243 learned manifold, rather than being the optimal dimension
 244 for capturing all of the variance in MNIST.

245 Figure 2b shows that STACKEDWAE manages to use all of
 246 its latent layers, capturing most of the covariance structure
 247 of the data in the deepest-latent layer, which is something
 248 that VAE methods struggle to accomplish (Sønderby et al.,
 249 2016; Zhao et al., 2017). Figure 3a shows the encoded input
 250 images through the latent layers, with corresponding digit
 251 labels coloured. We see through each layer that STACKED-
 252 WAE leverages the full hierarchy in its latents, with struc-
 253 tured manifolds learnt at each stochastic layer. Note that for
 254 the shallower layers with higher dimensions (*e.g.* $d_{\mathcal{Z}_1} = 32$
 255 or $d_{\mathcal{Z}_2} = 16$), the PCA algorithm results in a poor visualisa-
 256 tion. For such high dimensional spaces, other visualisation
 257 techniques can be used such as UMAP (McInnes & Healy,
 258 2018) as done in Figure 8a of Appendix B.1.

259 An advantage of deep-latent hierarchies is their capacity to
 260 capture information at different levels, augmenting a single-
 261 layer latent space. In Figure 3b, input images, shown in
 262 the bottom-row, are encoded and reconstructed for each
 263 latent layer. More specifically, the inputs are encoded up
 264 to the latent layer i and reconstructed from the encoded z_i
 265 using the generative model $p(x|z_1)p(z_1|z_2)\dots p(z_{i-1}|z_i)$.
 266 Row i in Figure 3b shows the reconstruction obtained from
 267 encoding up to layer i . We can see that each additional en-
 268 coding layer moves slightly farther away from copying the
 269 input image as it moves towards fitting the encoding into the
 270

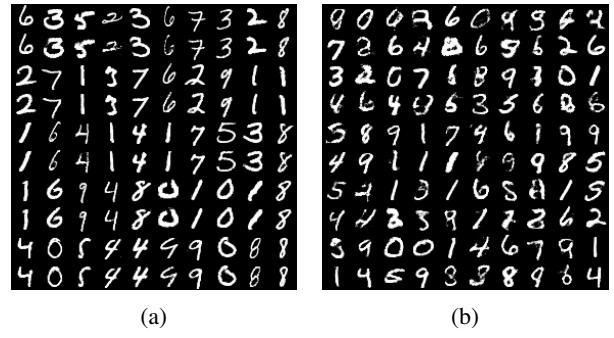


(a) (b)

Figure 2: 5-layer STACKEDWAE. (a) Model samples. (b) \mathcal{Z}_5 latent-space interpolations in a 2-dimensional grid varying between ± 2 standard deviations from the origin.

(a)

(b)

Figure 3: 5-layer STACKEDWAE. (a) Visualisations of latent spaces \mathcal{Z}_i . Each colour corresponds to a digit label. $d_{\mathcal{Z}_5} = 2$ can be directly plotted; for higher dimensions we use PCA. (b) Reconstructions for different encoding layers. The bottom row is data; the i^{th} row from the bottom is generated using the latent codes z_i which are from the i^{th} encoding layer.

(a) (b)

Figure 4: 1-layer implicit-prior WAE. (a) Reconstructions (within pairs of rows, data is above with the corresponding reconstructions below). (b) \mathcal{Z}_5 latent-space interpolations.

2-dimensional, unit-normal prior distribution. The dimensionality of the deeper-latent layers is a modelling choice which determines how much information is preserved; this can be seen through the loss of information from deeper reconstructions in Figure 3b (for a better visualisation of the *real* model reconstructions, corresponding to the top-row in Figure 3b, see Figure 7, bottom-left, of Appendix B.1). Indeed, in each layer, the encoder is asked to map the incoming encodings into a lower-dimensional latent space, filtering the amount of information to keep and pass up to the deeper layer. Thus, if there is a mismatch in the dimensions between the true underlying generative process of the data and the chosen model, the encoder will have to project the encodings into lower-dimensional space, losing information along the way.

ABLATION STUDY: STACKEDWAE VERSUS WAE

In this section, we compare STACKEDWAE with the original WAE framework for training generative models with deep-hierarchical latents. In particular, we train a WAE using the objective defined in Eq. (8) and an inference distribution as in Eq. (2); the corresponding graphical model is shown in Figure 1c. We use the same experimental setup as outlined earlier in this section, and the same parametrised networks. This experiment can be related to the work of Tomczak & Welling (2018); Klushyn et al. (2019) in the VAE framework case, as we can rephrase it as training a plain-vanilla 1-layer WAE whose prior over the latent variable is defined and parametrised as in Eq. (7), and learned alongside the inference network and the generative model. However, we do not use any specific optimisation scheme nor do we constraint the structure of the prior beside the modelling choices made in Section 2.1, as opposed to what is done in Tomczak & Welling (2018); Klushyn et al. (2019).

The results are shown in Figure 4, with reconstructions in Figure 4a, and deepest-latent interpolations in Figure 4b. The latter two should be compared directly with bottom and top-row of Figure 3b and Figure 2b, respectively, for the STACKEDWAE (see also Figure 7 of Appendix B.1 for samples and reconstructions comparisons with STACKEDWAE). The samples generated with the standard WAE when training deep-hierarchical-latent variable models (Figure 7, top-right, of Appendix B.1) are poor in comparison with those of the STACKEDWAE (Figure 7, top-left, of Appendix B.1).

The lack of smooth interpolations in Figure 4b shows that almost no structure has been captured in the deepest-latent layer. This is likely due to the fact that the standard WAE, with objective given in Eq. (8), is independent of the deeper-latent inference distributions, thus weakening the smoothness constraint in the deeper layers. The relatively accurate reconstructions in Figure 4a indicate that the model only needs the first latent layer to capture most of the structure

of the data. This behaviour is similar to that of the Markov HVAE as described in Zhao et al. (2017). They show that, for Markov HVAEs to learn interpretable latents, one needs additional constraints beyond simply maximising the likelihood, or in our case, the WAE objective of Tolstikhin et al. (2018) with MMD divergence.

3.2. Real world Datasets

We now turn to more realistic datasets to show that STACKEDWAE is still able to leverage a deep-latent hierarchy. In particular, we trained a 6-layer and a 10-layer STACKEDWAE on Street View House Numbers (SVHN) (Netzer et al., 2011) and CelebA (Liu et al., 2015), respectively. Our goal in this section, is to train very deep-latent variable models such as the ones in Maaløe et al. (2019), and show how STACKEDWAE still manages to use the deepest-latent layers of the generative model. It is worth stating that we are not aiming to achieve state-of-the-art performance on these datasets, but rather to show that it is easy to learn a deep-hierarchical latent representation of the data.

Our implementation choices resulted in relatively high-dimensional latent spaces (see below and Appendix B.2 for more details). Rubenstein et al. (2018) observed that when training WAEs with high-dimensional latent spaces, the variance of the inference network tends to collapse to 0. The authors argue that this might come either from an optimisation issue or from the failing of the divergence used to regularise the latents. Either way, the collapse to deterministic encoders results in poor sample quality as the deterministic encoder is being asked to map the input manifold into a higher-dimensional space than its intrinsic dimension. They proposed to explicitly include a regulariser that maintains a non-zero variance in Eq. (14). As in Rubenstein et al. (2018), we included a log-variance penalty term given by Eq. (17):

$$\mathcal{L}_{\text{pen}} = \sum_{i=1}^N \lambda_i^\Sigma \sum_{m=1}^{d_{z_i}} |\log \Sigma_i^q[m]| \quad (17)$$

Where Σ_i^q is the covariance matrix of the n^{th} inference network. We find that an exponentially decreasing penalty term λ_i^Σ (see following sections) works well for the dataset at hand with our experimental setup. This choice is motivated by the fact that the bigger the latent dimension (shallower latent layers in the hierarchy), the more likely it is that the latent dimension is larger than the data's intrinsic dimension.

STREET VIEW HOUSE NUMBERS

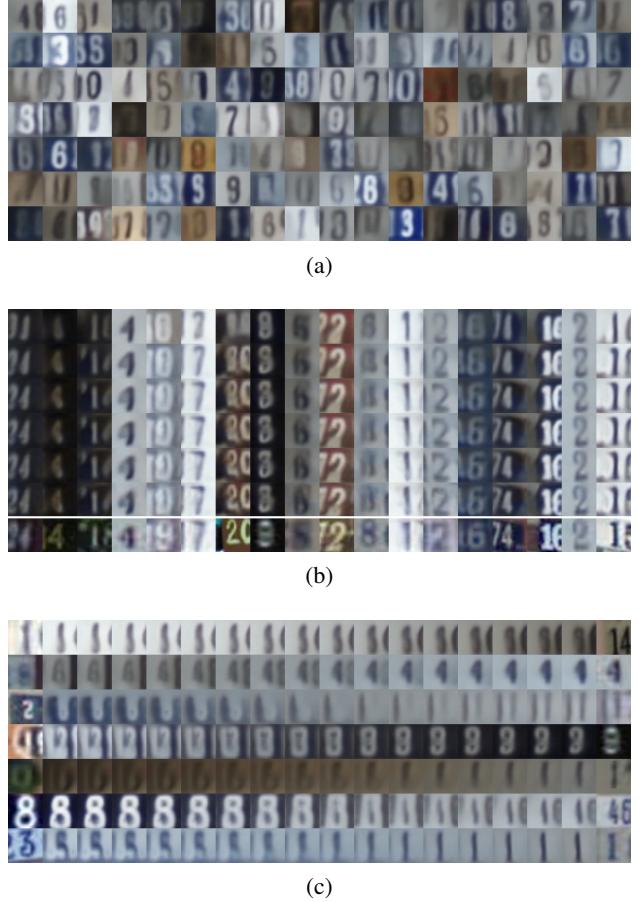
We trained a 6-layer STACKEDWAE with Gaussian inference networks and generative models at each latent layer, with mean and covariance functions parametrised by 3-layer ResNet-style neural networks (Kaiming et al., 2015). The

330 details are given in Appendix B.2. These architectures dictate the dimensionality of the latent spaces, with the latent
 331 dimension given by the size of the output feature map at
 332 that layer times the number of these feature maps. In this
 333 experiment, networks $i = 1, 3, 5$ have stride 2, effectively
 334 performing downsampling and upsampling operations in
 335 the bottom-up path and the top-down path respectively, re-
 336 sulting in latent dimensions (width \times height \times feature map)
 337 of: $16 \times 16 \times 2$, $16 \times 16 \times 1$, $8 \times 8 \times 2$, $8 \times 8 \times 1$, $4 \times 4 \times 2$ and $4 \times 4 \times 1$.
 338 For the regularisation hyperparameters, we use $\prod_{i=1}^n \lambda_i =$
 339 $\lambda_{\text{rec}}^{(n-1)+1}$ for $n = 1, \dots, 5$ (each reconstruction term in the
 340 objective), and $\prod_{i=1}^6 \lambda_i = \lambda_{\text{match}}$ (the final divergence term).
 341 The choice for the reconstruction weights is motivated by
 342 the fact that the effective regularisation hyperparameters
 343 scale exponentially. Thus, to avoid the collapse (or blowing
 344 up for $\lambda_{\text{rec}} > 1$) of the corresponding reconstruction terms,
 345 we choose the weights to scale like $\mathcal{O}(\lambda_{\text{rec}}^n)$. We found
 346 that $(\lambda_{\text{rec}}, \lambda_{\text{match}}) = (0.5, 10)$ worked well with our ex-
 347 perimental setup. As mentioned above, in order to avoid the
 348 collapse to deterministic encoder networks, we penalised the
 349 log-variance of the inference networks (as done in Ruben-
 350 stein et al. (2018)). We found that $\lambda_i^\Sigma = \lambda^\Sigma \cdot e^{-(i-1)}$, for
 351 $i = 1, \dots, 6$, with $\lambda^\Sigma = 2.5$ worked well in our setting.
 352 More details of the experimental setup can be found in Ap-
 353 pendix B.2.
 354

355 Our results on SVHN are given in Figure 5. Samples from
 356 the generative model are shown in Figure 5a. Figure 5b
 357 shows the reconstructions of the data points (along the bot-
 358 tom row) at each latent layer in the hierarchy. Similarly to
 359 the results for MNIST, we can see that the deepest latent
 360 layer may not be large enough to enable high-fidelity recon-
 361 structions. Our intention is to show that the hierarchy of
 362 latents can be learnt, which is clearly the case, rather than to
 363 model SVHN perfectly, so we do not attempt to tune to the
 364 optimal latent dimensionality. Figure 5c represents point
 365 interpolation, with the actual anchor data points shown in
 366 the first and last column, their reconstructions in the sec-
 367 ond, respectively second-to-last, columns, and the latent
 368 interpolation in-between.
 369

CELEBA

370 We trained a 10-layers STACKEDWAE. Similarly to
 371 SVHN, the inference and generative networks are fully-
 372 convolutional ResNet. The inference networks and gener-
 373 ative models, $i = 1, 4, 7, 10$ have stride 2. These down-
 374 sampling, respectively upsampling, operations resulted in
 375 4 groups of latent spaces sharing features maps of the
 376 same spatial dimensions: latent layers $i = 1, 2, 3$ with
 377 output features maps of size 32×32 , layers $i = 4, 5, 6$
 378 with features maps of size 16×16 , layers $i = 7, 8, 9$ with
 379 size 8×8 and layer $i = 10$ with size 4×4 . Within each
 380 group, the number of feature maps is decreased by two
 381 as we go up the hierarchy, with the biggest layer having
 382



383 Figure 5: 6-layer STACKEDWAE. (a) Model samples. (b)
 384 Reconstructions for different encoding layers, as in Fig-
 385 ure 3b. (c) Points interpolations; the first and last columns
 386 are actual data points, with corresponding reconstruc-
 387 tions shown in the second, respectively second-to-last, columns.
 388

385 8 features maps and the smallest 4. As with SVHN, we
 386 used an exponentially decreasing regularisation parame-
 387 ters. More precisely, we chose $\prod_{i=1}^p \lambda_i = \lambda_{\text{rec}}^{(p-1)/3+1}$ for
 388 $p = 1, \dots, 9$ and $\prod_{i=1}^{10} \lambda_i = \lambda_{\text{match}}$. For this dataset, choos-
 389 ing $(\lambda_{\text{rec}}, \lambda_{\text{match}}) = (10^{-1}, 10^0)$ worked well. We also
 390 added a L_1 penalty on the log-covariance matrices of the
 391 encoders to avoid any variance collapse. We used the same
 392 penalisation scheme than in the previous experiment, that is
 393 $\lambda_i^\Sigma = \lambda^\Sigma \cdot e^{-(i-1)}$, for $i = 1, \dots, 10$, with $\lambda^\Sigma = 2.5$. See
 394 Appendix B.2 for more details about the experimental setup.
 395

396 Results are shown Figure 6, with model samples given in
 397 Figure 6a and layer-wise reconstructions in Figure 6b. As
 398 with the 6-layer STACKEDWAE trained on SVHN, Figure
 399 6b shows that STACKEDWAE manages to use all its
 400 latent layers, up to the deepest ones. While the full recon-
 401 structions, shown in the top-row of Figure 6b, are relatively
 402 close to the original data points, shown in the bottom-row,
 403 we can notice, as in the MNIST and SVHN experiments, a
 404 loss of information as we go deeper in the hierarchy. In other
 405 words, the deeper the encoding, the smoother or blurrier the
 406 reconstructions. Here again, one possible explanation is
 407 the excessive filtering of information by the encoder when
 408 going up in the hierarchy due to the miss-match between
 409 the intrinsic dimension and the latent dimension.

4. Conclusion

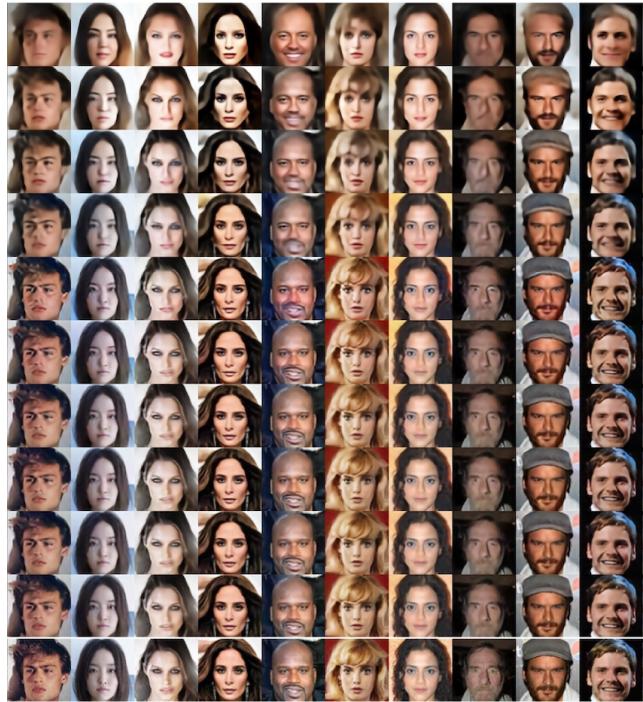
410 In this work we introduced a novel objective function for
 411 training generative models with deep hierarchies of latent
 412 variables using Optimal Transport. Our approach recur-
 413 sively applies the Wasserstein distance as the regularisation
 414 divergence, allowing the stacking of WAEs for arbitrarily
 415 deep-latent hierarchies. We showed that this approach en-
 416 ables the learning of smooth latent distributions even in
 417 deep latent hierarchies, which otherwise requires extensive
 418 model design and tweaking of the optimisation procedure
 419 to train. We also showed that our approach is significantly
 420 more effective at learning smooth hierarchical latents than
 421 the standard WAE.
 422

References

- 423 Arjovsky, M., Chintala, S., and Bottou, L. Wasserstein gen-
 424 erative adversarial networks. In *International Conference
 425 on Machine Learning*, 2017.
- 426 Bachman, P. An architecture for deep, hierarchical genera-
 427 tive models. In *Advances in Neural Information Process-
 428 ing Systems*, 2016.
- 429 Bousquet, O., Gelly, S., Tolstikhin, I., Simon-Gabriel, C. J.,
 430 and Schölkopf, B. From optimal transport to generative
 431 modeling: the VEGAN cookbook. 2017.
- 432 Burda, Y., R., G., and Salakhutdinov, R. Importance
 433



(a)



(b)

Figure 6: 10-layer STACKEDWAE. (a) Model samples. (b) Reconstructions for different encoding layers, as in Figure 3b.

- 440 weighted autoencoders. In *International Conference on*
 441 *Learning Representations*, 2015.
- 442 Cuturi, M. Sinkhorn distances: lightspeed computation of
 443 optimal transport. In *Advances in Neural Information*
 444 *Processing Systems*, 2013.
- 445 Dinh, L., Sohl-Dickstein, J., and Bengio, S. Density esti-
 446 mation using real NVP. In *International Conference on*
 447 *Learning Representations*, 2016.
- 448 Gaujac, B., Feige, I., and Barber, D. Gaussian mixture
 449 models with Wasserstein distance. *arXiv:1806.04465*,
 450 2018.
- 451 Genevay, A., Peyre, G., and Cuturi, M. Learning genera-
 452 tive models with Sinkhorn divergences. In *International*
 453 *Conference on Artificial Intelligence and Statistics*, 2018.
- 454 Glorot, X., Bordes, A., and Bengio, Y. Deep sparse rectifier
 455 neural networks. In *International Conference on Artificial*
 456 *Intelligence and Statistics*, 2011.
- 457 Gretton, A., Borgwardt, K. M., Rasch, M. J., Schölkopf, B.,
 458 and Smola, A. A kernel two-sample test. In *Journal of*
 459 *Machine Learning Research*, 2012.
- 460 Ioffe, S. and Szegedy, C. Batch normalization: Accelerating
 461 deep network training by reducing internal covariate shift.
 462 In *International Conference on Machine Learning*, 2015.
- 463 Kaiming, H., Xiangyu, Z., Shaoqing, R., and Jian, S. Deep
 464 residual learning for image recognition. In *IEEE Confer-
 465 ence on Computer Vision and Pattern Recognition*, 2015.
- 466 Kingma, D. P. and Ba, J. Adam: a method for stochastic
 467 optimization. In *International Conference on Learning*
 468 *Representations*, 2015.
- 469 Kingma, D. P. and Dhariwal, P. Glow: Generative flow
 470 with invertible 1x1 convolutions. In *Advances in Neural*
 471 *Information Processing Systems*, 2018.
- 472 Kingma, D. P. and Welling, M. Auto-encoding variational
 473 Bayes. In *International Conference on Learning Repre-
 474 sentations*, 2014.
- 475 Kingma, D. P., Salimans, T., Jozefowicz, R., Chen, X.,
 476 Sutskever, I., and Welling, M. Improved variational in-
 477 ference with inverse autoregressive flow. In *Advances in*
 478 *Neural Information Processing Systems*, 2016.
- 479 Klushyn, A., Chen, N., Kurle, R., Cseke, B., and van ver
 480 Smagt, P. Learning hierarchical priors in VAEs. In *Ad-
 481 vances in Neural Information Processing Systems*, 2019.
- 482 Larochelle, H. and Murray, I. The neural autoregressive
 483 distribution estimator. In *International Conference on*
 484 *Artificial Intelligence and Statistics*, 2011.
- 485 Larsen, A. B. L., Sønderby S., K., Larochelle, H., and
 486 Winther, O. Autoencoding beyond pixels using a learned
 487 similarity metric. In *International Conference on Ma-
 488 chine Learning*, 2016.
- 489 LeCun, Y. and Cortes, C. MNIST handwritten digit database.
 490 2010.
- 491 Liu, Z., Luo, P., Wang, X., and Tang, X. Deep learning face
 492 attributes in the wild. In *Proceedings of International*
 493 *Conference on Computer Vision*, 2015.
- 494 Maaløe, L., Fraccaro, M., Liévin, V., and Winther, O. BIVA:
 495 a very deep hierarchy of latent variables for generative
 496 modeling. In *Advances in neural information processing*
 497 *systems*, 2019.
- 498 Maaløe, L., Sønderby, C. K., Sønderby, S. K., and Winther,
 499 O. Auxiliary deep generative models. In *International*
 500 *Conference on Machine Learning*, 2016.
- 501 McInnes, L. and Healy, J. UMAP: uniform manifold
 502 approximation and projection for dimension reduction.
 503 *arXiv:1802.03426*, 2018.
- 504 Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., and
 505 Ng, A. Y. Reading digits in natural images with unsuper-
 506 vised feature learning. In *Advances in Neural Information*
 507 *Processing Systems Workshop on Deep Learning and Un-
 508 supervised Feature Learning*, 2011.
- 509 Patrini, G., Carioni, M., Forré, P., Bhargav, S., Welling, M.,
 510 Van Den Berg, R., Genewein, T., and Nielsen, F. Sinkhorn
 511 autoencoders. *arXiv:1810.01118*, 2018.
- 512 Rezende, D. J. and Mohamed, S. Variational inference
 513 with normalizing flows. In *International Conference on*
 514 *Machine Learning*, 2015.
- 515 Rezende, D. J., Mohamed, S., and Wierstra, D. Stochastic
 516 backpropagation and approximate inference in deep gen-
 517 erative models. In *International Conference on Machine*
 518 *Learning*, 2014.
- 519 Rubenstein, P. K., Schölkopf, B., and Tolstikhin, I. On the
 520 latent space of Wasserstein auto-encoders. In *Workshop*
 521 *track - International Conference on Learning Represen-
 522 tations*, 2018.
- 523 Salimans, T., Karpathy, A., Chen, X., and Kingma, D. P.
 524 Pixelcnn++: Improving the PixelCNN with discretized
 525 logistic mixture likelihood and other modifications. In
 526 *International Conference on Learning Representations*,
 527 2017.
- 528 Sønderby, C. K., Raiko, T., Maaløe, L., Sønderby, S. K., and
 529 Winther, O. Ladder variational autoencoders. In *Advances*
 530 *in neural information processing systems*, 2016.

- 495 Tolstikhin, I., Bousquet, O., Gelly, S., and Schoelkopf, B.
496 Wasserstein auto-encoders. In *International Conference*
497 *on Learning Representations*, 2018.
- 498 Tomczak, J. M. and Welling, W. VAE with VampPrior. In
499 *International Conference on Artificial Intelligence and*
500 *Statistics*, 2018.
- 501
- 502 Van Den Oord, A., Kalchbrenner, N., and Kavukcuoglu,
503 K. Pixel recurrent neural networks. In *International*
504 *Conference on Machine Learning*, 2016a.
- 505
- 506 Van Den Oord, A., Kalchbrenner, N., Vinyals, O., Espeholt,
507 L., Graves, A., and Kavukcuoglu, K. Conditional image
508 generation with PixelCNN decoders. In *Advances in*
509 *Neural Information Processing Systems*, 2016b.
- 510
- 511 Villani, C. *Optimal Transport: Old and New*. Springer
512 Berlin Heidelberg, 2008.
- 513 Zha, S., Song, J., and Ermon, S. Towards deeper
514 understanding of variational autoencoding models.
515 *arxiv:1702.08658*, 2017.
- 516
- 517 Zhao, S., Song, J., and Ermon, S. Learning hierarchical
518 features from deep generative models. In *International*
519 *Conference on Machine Learning*, 2017.
- 520
- 521
- 522
- 523
- 524
- 525
- 526
- 527
- 528
- 529
- 530
- 531
- 532
- 533
- 534
- 535
- 536
- 537
- 538
- 539
- 540
- 541
- 542
- 543
- 544
- 545
- 546
- 547
- 548
- 549

A. StackedWAE derivation details

A.1. Marginal constraint

The space of couplings $\Gamma \in \mathcal{P}(P_D, P_\Theta)$ that defines the OT distance in Eq. (3) is constrained according to Eq. (4). WAE assumes a joint density of the form given in Eq. (5), which automatically satisfies the p_D marginal constraint, but requires the further sufficient condition given in Eq. (6) in order to satisfy the p_Θ marginal constraint. To see that Eq. (6) is indeed a sufficient condition for the p_Θ marginal constraint, note that from the Markovian assumption of the generative model (see Eq. (1)), we can write γ as

$$\forall(x, \tilde{x}) \in \mathcal{X} \times \mathcal{X},$$

$$\begin{aligned}\gamma(x, \tilde{x}) &= \int_{\mathcal{Z}_{1:N}} p_{\theta_1}(\tilde{x}|z_1) q_{\phi_1}(z_{1:N}|x) p_D(x) dz_{1:N} \\ &= \int_{\mathcal{Z}_1} p_{\theta_1}(\tilde{x}|z_1) q_{\phi_1}(z_1|x) p_D(x) dz_1\end{aligned}\quad (18)$$

The constraint in Eq. (4) on the p_D marginal is trivially true as the integral over the second variable can be brought inside the integral over \mathcal{Z}_1 , after which all the integrals simply integrate to unity leaving p_D .

The constraint on the second marginal is obtained by integrating Eq. (18) over the first variable,

$$\forall x \in \mathcal{X},$$

$$\begin{aligned}\int_{\mathcal{X}} \gamma(x, \tilde{x}) dx &= \int_{\mathcal{X}} \int_{\mathcal{Z}_1} p_{\theta_1}(\tilde{x}|z_1) q_{\phi_1}(z_1|x) p_D(x) dz_1 dx \\ &= \int_{\mathcal{Z}_1} p_{\theta_1}(\tilde{x}|z_1) \int_{\mathcal{X}} q_{\phi_1}(z_1|x) p_D(x) dx dz_1\end{aligned}\quad (19)$$

Thus, to satisfy Eq. (4), we need:

$$\forall \tilde{x} \in \mathcal{X},$$

$$\begin{aligned}\int_{\mathcal{Z}_1} p_{\Theta}(\tilde{x}|z_1) \overbrace{\int_{\mathcal{X}} q_{\phi_1}(z_1|x) p_D(x) dx}^{\stackrel{\text{def}}{=} q_1^{\text{agg}}(z_1)} dz_1 &\stackrel{\text{need}}{=} \\ &\int_{\mathcal{Z}_1} p_{\theta_1}(\tilde{x}|z_1) p_{\Theta_{2:N}}(z_1) dz_1\end{aligned}\quad (20)$$

where we use the generative model defined in Eq. (1) and we introduced:

$$\forall z_1 \in \mathcal{Z}_1,$$

$$p_{\Theta_{2:N}}(z_1) \stackrel{\text{def}}{=} \int_{\mathcal{Z}_{2:N}} \prod_{n=2}^N p_{\theta_n}(z_{n-1}|z_n) p(z_N) dz_{2:N}\quad (21)$$

To satisfy Eq. (20), one obvious sufficient condition on the aggregated posterior distribution $Q_1^{\text{agg}}(Z_1)$ defined in Eq. (9) is Eq. (6), namely that

$$\forall z_1 \in \mathcal{Z}_1, \quad q_1^{\text{agg}}(z_1) = p_{\Theta_{2:N}}(z_1) \quad (22)$$

which is what we sought out to show. However, Eq. (22) is a sufficient condition, not a necessary one: indeed Eq. (22) must only hold under $\int_{\mathcal{Z}_1} P_\Theta(\tilde{X}|z_1) dz_1$. So for example, if $P_{\theta_1}(\tilde{X}|z_1) = P_{\theta_1}(\tilde{X})$, then Eq. (20) would boil down to a constraint only on the expectations of $Q_1^{\text{agg}}(Z_1)$ and $P_{\Theta_{2:N}}(Z_1)$.

A.2. WAE objective

Starting from the definition of the OT distance given in Eq. (3), and using the WAE approach with density γ written as Eq. (18), we find:

$$\begin{aligned}\text{OT}_c(P_D, P_\Theta) &= \inf_{\Gamma \in \mathcal{P}(P_D, P_\Theta)} \int_{\mathcal{X} \times \mathcal{X}} c(x, \tilde{x}) d\Gamma(x, \tilde{x}) \\ &\leq \inf_{\substack{Q_{\phi_1}(Z_1, Z_2, \dots, Z_n|X) \\ \text{Eq.(20) satisfied}}} \int_{\mathcal{X} \times \mathcal{X}} c(x, \tilde{x}) \gamma(x, \tilde{x}) dx d\tilde{x}\end{aligned}$$

where γ is defined in Eq. (18). Given that the above does not depend on $z_{>1}$, the inf can be written over $Q_{\phi_1}(Z_1|X)$ rather than the full $Q_{\phi_1}(Z_1, Z_2, \dots, Z_n|X)$. Replacing the constraint in the inf with the sufficient condition according to Eq. (6), which amounts to replacing Eq. (20) with Eq. (22), we obtain:

$$\text{OT}_c(P_D, P_\Theta) \leq \inf_{\substack{Q_{\phi_1}(Z_1|X) \\ Q_1^{\text{agg}} = P_{\Theta_{2:N}}}} \int_{\mathcal{X} \times \mathcal{X}} c(x, \tilde{x}) \gamma(x, \tilde{x}) dx d\tilde{x}\quad (23)$$

with γ still as in Eq. (18). Eq. (8) is then obtained by relaxing constraint in Eq. (23); replacing the hard constraint by a soft constraint via a penalty term added to the objective, weighted by a λ_1 :

$$\begin{aligned}\widehat{W}_c(P_D, P_\Theta) &= \inf_{Q_{\phi_1}(z_1|x)} \left[\int_{\mathcal{X} \times \mathcal{X}} c(x, \tilde{x}) \gamma(x, \tilde{x}) dx d\tilde{x} \right. \\ &\quad \left. + \lambda_1 \mathcal{D}_1(Q_1^{\text{agg}}(Z_1), P_{\Theta_{2:N}}(Z_1)) \right]\end{aligned}\quad (24)$$

where \mathcal{D}_1 is any divergence function between distributions on \mathcal{Z}_1 .

B. Experiments

B.1. MNIST experiments

EXPERIMENTAL SETUP

We train a deep-hierarchical latent-variable model with $N = 5$ latent layers whose dimensions are $d_{\mathcal{Z}_1} = 32$, $d_{\mathcal{Z}_2} =$

605 $d_{\mathcal{Z}_3} = 8$, $d_{\mathcal{Z}_4} = 4$ and $d_{\mathcal{Z}_5} = 2$, respectively. We
 606 parametrise the generative and inference models as:

$$\begin{aligned} q_{\phi_i}(z_i|z_{i-1}) &= \mathcal{N}(z_i; \mu_i^q(z_{i-1}), \Sigma_i^q(z_{i-1})), \quad i = 1, \dots, 5 \\ p_{\theta_i}(z_{i-1}|z_i) &= \mathcal{N}(z_{i-1}; \mu_i^p(z_i), \Sigma_i^p(z_i)), \quad i = 2, \dots, 5 \\ p_{\theta_1}(x|z_1) &= \delta(x - f_{\theta_1}(z_1)) \end{aligned} \quad (25)$$

612 For both the encoder and decoder, the mean and diagonal
 613 covariance functions μ_i, Σ_i are fully-connected networks
 614 with 2 same-size hidden layers (consider f_{θ_1} as μ_1^p). For $i =$
 615 $1, 2, 3, 4, 5$, the number of units is 2048, 1024, 512, 256, 128
 616 respectively.

617 For the regularisation hyperparameters, we use $\prod_{i=1}^n \lambda_i =$
 618 $\lambda_{\text{rec}}^n / d_{\mathcal{Z}_n}$ for $n = 1, \dots, 4$ (for each reconstruction term
 619 in the objective), and $\prod_{i=1}^5 \lambda_i = \lambda_{\text{match}}$ (for the final
 620 divergence term). We then perform a grid search over
 621 the 25 pairs $(\lambda_{\text{rec}}, \lambda_{\text{match}}) \in \{0.005, 0.01, 0.05, 0.1, 0.5\} \otimes$
 622 $\{10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}\}$ and find the best result
 623 (smallest Eq. (14)) is obtained with $(\lambda_{\text{rec}}, \lambda_{\text{match}}) =$
 624 $(0.01, 10^{-4})$.

625 Finally, we choose the squared Euclidean distance as the
 626 cost function: $c_n(z_n, \tilde{z}_n) = \|z_n - \tilde{z}_n\|_{L^2}^2$. The expectations
 627 in Eq. (14) are computed analytically whenever possible,
 628 and with Monte Carlo sampling otherwise. We use
 629 batch normalisation (Ioffe & Szegedy, 2015) after each hidden
 630 fully-connected layer, followed by a ReLU activation
 631 (Glorot et al., 2011). We train the models over 5,000 epochs
 632 using Adam optimiser (Kingma & Ba, 2015) with default
 633 parameters and batch size of 128.

636 ADDITIONAL RESULTS

637 B.2. Real word dataset

638 In the following experiments, the inference networks and
 639 generative models at each latent layer are taken to be
 640 Gaussian distributions with mean and covariance functions
 641 parametrised by 3-layer ResNet (similar to those of Eq. 25).
 642 A M -layer residual network is composed of $M - 1$ convolutional
 643 blocks followed by a resampling convolution, and a
 644 residual connection. The outputs of the two are added and a
 645 last operation (either fully connected or convolution layer) is
 646 applied on the result. A convolutional block is composed of
 647 a convolution layer followed by batch normalisation (Ioffe
 648 & Szegedy, 2015) and a ReLU non-linearity (Glorot et al.,
 649 2011). We also use batch normalisation and ReLU after
 650 the sum of the convolutional blocks output and the residual
 651 connection. See Figure 9 for an example of a 3-layers
 652 residual network with a last convolution operation. When
 653 resampling, we use a convolution layer with stride 2 in both
 654 the skip connection and the resampling convolution for the
 655 inference networks and a deconvolution layer with stride 2
 656 in both the skip connection and the resampling convolution
 657 in the generative models. If no resampling is performed,

9 6 6 2 3 8 1 7 5 3 6 7	2 1 0 2 1 5 0 9 6 8 2 2
1 5 1 3 9 0 4 7 9 0 4 4	4 9 5 7 7 8 1 9 7 8 9 9
6 1 9 3 8 7 3 3 7 2 5 3	8 2 6 9 8 0 3 6 3 4 7 5
9 5 9 0 0 5 8 0 7 6 0 8	4 3 9 8 8 3 1 1 8 4 1 1
7 7 9 6 9 0 4 6 6 4 9 4	5 8 1 8 3 1 9 9 1 1 2 6
0 2 0 1 8 5 3 5 6 7 7 0	2 1 2 0 3 3 0 8 8 1 9 0
9 1 9 3 8 6 1 3 7 5 9 8	6 8 3 4 1 9 7 9 1 8 9 8
9 0 1 1 2 1 0 9 2 1 3 6	7 9 6 8 7 8 3 5 1 1 6 1
6 2 1 1 6 2 1 6 0 3 0 0	2 0 8 1 1 0 8 8 4 1 6 2
5 4 2 1 3 0 9 4 7 2 4 9	2 4 1 6 0 0 9 1 9 7 8
1 5 8 5 4 9 9 2 1 1 2 0	4 9 6 4 8 3 9 4 0 8 1 3
0 1 8 1 7 7 1 1 8 0 0	3 6 9 7 0 8 0 9 4 8 2 5
6 3 5 2 3 6 7 3 2 8 6 9	6 3 5 2 3 6 7 3 2 8 6 9
6 3 5 7 3 4 7 3 2 8 6 9	6 3 5 2 3 6 7 3 2 8 6 9
2 7 1 3 7 6 2 9 1 1 9 0	2 7 1 3 7 6 2 9 1 1 9 0
2 7 1 3 7 6 2 9 1 1 9 0	2 7 1 3 7 6 2 9 1 1 9 0
1 6 4 1 4 1 7 5 3 8 4 3	1 6 4 1 4 1 7 5 3 8 4 3
1 5 4 1 4 1 7 5 3 8 4 3	1 6 4 1 4 1 7 5 3 8 4 3
1 6 9 4 8 0 1 0 1 8 4 6	1 6 9 4 8 0 1 0 1 8 4 6
1 6 7 4 3 6 1 0 1 8 4 6	1 6 9 4 8 0 1 0 1 8 4 6
4 0 5 4 9 9 0 8 8 5 3	4 0 5 4 9 9 0 8 8 5 3
4 0 8 9 7 9 9 0 8 8 5 3	4 0 5 4 9 9 0 8 8 5 3
0 0 1 9 1 9 4 5 1 7 1 5	0 0 1 9 1 9 4 5 1 7 1 5
0 0 1 9 1 9 4 5 1 7 1 5	0 0 1 9 1 9 4 5 1 7 1 5

Figure 7: 5-layer STACKEDWAE (left column) *versus* 1-layer implicit-prior WAE (right column). First row: Model samples. Second row: Reconstructions (within pairs of rows, data is above with the corresponding reconstructions below).

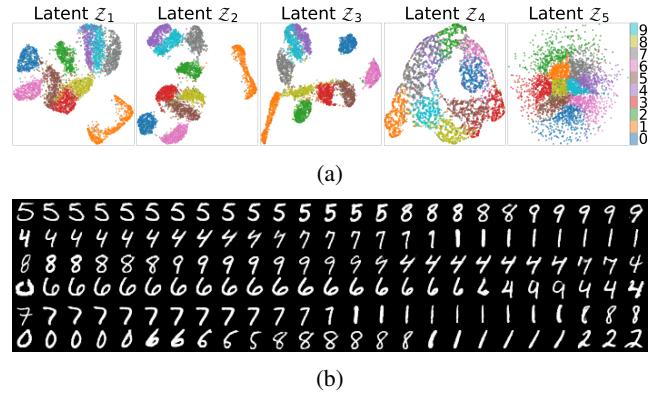


Figure 8: 5-layer STACKEDWAE. (a) UMAP (McInnes & Healy, 2018) visualisations of latent spaces \mathcal{Z}_i . Each colour corresponds to a digit label. $d_{\mathcal{Z}_5} = 2$ can be directly plotted; for higher dimensions we use UMAP. (b) Points interpolations; the first and last columns are actual data points with corresponding reconstructions shown in the second, respectively second-to-last, columns.

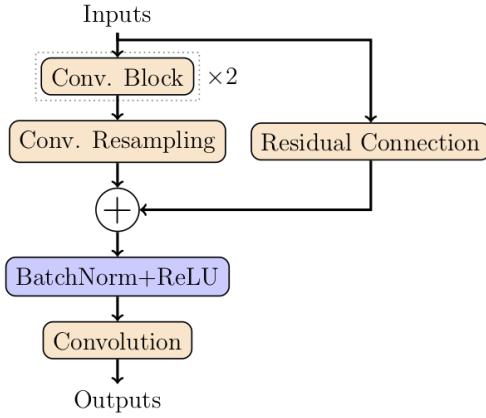


Figure 9: Residual network with 3 hidden convolutions.

then the resampling convolution is a simple convolution layer with stride 1 and the skip connection performs the identity operation. The latent dimensions are then given by the dimensions and the number of features in the last convolutional operation.

The ground cost is chosen to be the L_2 -squared norm given by $c_n(z_n, \tilde{z}_n) = \|z_n - \tilde{z}_n\|_{L^2}^2$, and the expectations in Eq. (14) are computed analytically whenever possible, otherwise using Monte Carlo sampling.

SVHN

We train a 6-layers STACKEDWAE on the SVHN dataset using both the training dataset (73,257 digits) and the additional training dataset (531,131 digits). We train the models over 1000 epochs using Adam optimiser (Kingma & Ba, 2015) with default parameters and batch size of 100.

As previously mentioned, we use 3-layer ResNet for the mean and covariance functions of the inference networks and generative models. More specifically, for each ResNet, we use $M = 2$ convolutional blocks with the dimensions of the filters specified in the top-table of Figure 10. Networks in layers $i = 1, 2$ have 64 convolution filters, layers $i = 3, 4$ 64 and layers $5, 6$ 128, each filters having the same size within each residual network, and doubling (in the inference networks) or divide by 2 (in the generative models) their number in the resampling convolution if any resampling is performed. The latent layers $i = 1, 3, 5$ have a stride of 2 and we choose the number of features to be 2, 1, 2, 1, 2 and 1 for the latent layers $i = 1 \dots 6$ (top-table of Figure 10 for the full details).

SVHN			
Layer i	Filters dim.	Resampling	Output dim.
Layer 1	$5 \times 5 \times 64$	down / up	$16 \times 16 \times 2$
Layer 2	$3 \times 3 \times 64$	None	$16 \times 16 \times 1$
Layer 3	$3 \times 3 \times 96$	down / up	$8 \times 8 \times 2$
Layer 4	$3 \times 3 \times 96$	None	$8 \times 8 \times 1$
Layer 5	$3 \times 3 \times 128$	down / up	$4 \times 4 \times 2$
Layer 6	$3 \times 3 \times 128$	None	$4 \times 4 \times 1$

CelebA			
Layer i	Filters dim.	Resampling	Output dim.
Layer 1	$7 \times 7 \times 64$	down / up	$32 \times 32 \times 8$
Layer 2	$5 \times 5 \times 64$	None	$32 \times 32 \times 6$
Layer 3	$3 \times 3 \times 64$	None	$32 \times 32 \times 4$
Layer 4	$3 \times 3 \times 64$	down / up	$16 \times 16 \times 8$
Layer 5	$3 \times 3 \times 96$	None	$16 \times 16 \times 6$
Layer 6	$3 \times 3 \times 96$	None	$16 \times 16 \times 4$
Layer 7	$3 \times 3 \times 96$	down / up	$8 \times 8 \times 8$
Layer 8	$3 \times 3 \times 128$	None	$8 \times 8 \times 6$
Layer 9	$3 \times 3 \times 128$	None	$8 \times 8 \times 4$
Layer 10	$3 \times 3 \times 128$	down / up	$4 \times 4 \times 8$

Figure 10: Details of the architectures used in Section 3.2. Top-table: architecture of the 6-layer STACKEDWAE trained on SVHN. Bottom-table: architecture of the 10-layer STACKEDWAE trained on CelebA.

715 CELEBA

716 We train a 10-layers STACKEDWAE on the CelebA dataset
717 over 100 epochs using Adam optimiser ([Kingma & Ba,](#)
718 [2015](#)) with default parameters and batch size of 64.

720 We use the same ResNet building blocks than previously,
721 mainly, $M = 2$ convolutional block, each filter within a
722 ResNet block having the same size, and doubling their num-
723 ber in the last convolutional operation. Networks in layers
724 $i = 1, 2, 3, 4$ have 64 filters, 96 in layers $i = 5, 6, 7$ and
725 128 in the remaining layers. The latent layers $i = 1, 4, 7, 10$
726 have stride 2 and we choose the number of features to be
727 8, 4, 2, 8, 4, 2, 8 (See bottom-table of Figure 10 for the full
728 details).

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769