702 A APPENDIX

A.1 CONVERSION BETWEEN CARTESIAN AND HYPERSPHERICAL COORDINATES

For reference, we note here the standard formulas for converting between cartesian and spherical coordinates (as they appear in https://en.wikipedia.org/wiki/N-sphere).

⁷⁰⁸ ⁷⁰⁹ ⁷⁰⁹ ⁷⁰⁰ ⁷¹⁰ ⁷¹⁰ ⁷¹¹ In *n* dimensions, given a set of cartesian coordinates x_k with $k \in \{1, ..., n\}$, the hyperspherical ⁷⁰⁹ ⁷⁰⁹ ⁷¹⁰ ⁷¹⁰ ⁷¹¹ ⁷¹¹ ⁷¹¹ ⁷¹¹ ⁷¹¹ ⁷¹¹ ⁷¹¹ ⁷¹¹ ⁷¹² ⁷¹² ⁷¹² ⁷¹³ ⁷¹⁴ ⁷¹⁵ ⁷¹⁵ ⁷¹⁶ ⁷¹⁷ ⁷¹⁷ ⁷¹⁷ ⁷¹⁸ ⁷¹⁸ ⁷¹⁹ ⁷¹⁹ ⁷¹⁰ ⁷¹¹

(9)

From hyperspherical to cartesian conversion:

716

 $x_2 = r \sin(\varphi_1) \cos(\varphi_2)$

- $x_2 = r\sin(\varphi_1)\sin(\varphi_2)\cos(\varphi_3)$
- 717
- 718

 $x_1 = r\cos(\varphi_1)$

719 720 $x_{n-1} = r \sin(\varphi_1) \sin(\varphi_2) \dots \sin(\varphi_{n-2}) \cos(\varphi_{n-1})$ $x_n = r \sin(\varphi_1) \sin(\varphi_2) \dots \sin(\varphi_{n-2}) \sin(\varphi_{n-1})$

721 722

723

741 742

743 744 From cartesian to hyperspherical conversion:

 $r = \sqrt{x_n^2 + x_{n-1}^2 + \ldots + x_2^2 + x_1^2}$ $\cos(\varphi_1) = \frac{x_1}{\sqrt{x_n^2 + x_{n-1}^2 + \ldots + x_2^2 + x_1^2}}$ $\cos(\varphi_2) = \frac{x_2}{\sqrt{x_n^2 + x_{n-1}^2 + \ldots + x_2^2}}$ 724 725 726 727 728 729 730 731 (10)732 : $\cos(\varphi_{n-2}) = \frac{x_{n-2}}{\sqrt{x_n^2 + x_{n-1}^2 + x_{n-2}^2}}$ 733 734 735 $\cos(\varphi_{n-1}) = \frac{x_{n-1}}{\sqrt{x_n^2 + x_{n-1}^2}}$ 736 737 738 739 A.2 VECTORIZED CODE FOR CONVERTING BETWEEN CARTESIAN AND HYPERSPHERICAL 740

A.2 VECTORIZED CODE FOR CONVERTING BETWEEN CARTESIAN AND HYPERSPHERICAL COORDINATES

This code is accessible here and provided below for reference.

745 import torch 746 747 def r (x): 748 749 r = torch.linalq.norm (x, dim=1) 750 return r 751 752 753 def cart_to_cos_sph (x, device): 754 m = x.size(0)

```
756
          n = x.size(1)
757
758
         mask = torch.triu(torch.ones(n, n)).to(device)
759
         mask = torch.unsqueeze(mask, dim=0)
760
761
         mask = mask.expand(m, n, n)
762
763
          X = torch.unsqueeze(x, dim=1).expand(m, n, n)
764
         X_squared = torch.square(X)
765
766
          X_squared_masked = X_squared * mask
767
768
          denom = torch.sqrt(torch.sum(X_squared_masked, dim=2)+0.001)
769
         cos_phi = x / denom
770
771
          return cos_phi[:, 0:n-1]
772
773
774
      def cart_to_sin_sph (x, device):
775
          return torch.sqrt (1 - cart_to_cos_sph (x, device).pow(2))
776
777
778
      def cart_to_sph (x, device):
779
         m = x.size(0)
780
781
          n = x.size(1)
782
783
         mask = torch.triu(torch.ones(n, n)).to(device)
784
         mask = torch.unsqueeze(mask, dim=0)
785
786
         mask = mask.expand(m, n, n)
787
788
         X = torch.unsqueeze(x, dim=1).expand(m, n, n)
789
          X_squared = torch.square(X)
790
791
         X_squared_masked = X_squared * mask
792
793
          denom = torch.sqrt(torch.sum(X_squared_masked, dim=2)+0.001)
794
         phi_plus = torch.arccos (x / denom)
795
796
          phi_minus = 2*3.141592654 - phi_plus
797
798
         phi = phi_plus
799
          phi[:, n-2] = torch.where (x[:, n-1] >= 0, phi_plus[:, n-2],
800
             phi_minus[:, n-2])
801
802
          return phi[:, 0:n-1]
803
804
      def sph_to_cart (R, phi, device):
805
806
         m = phi.size(0)
807
808
          n = phi.size(1)+1
809
         mask = torch.tril(torch.ones(n-1, n-1)).to(device)
```

mask = torch.unsqueeze(mask, dim=0)

sin_PHI_masked = sin_PHI * mask + mask_

sin_prod = torch.prod (sin_PHI_masked, dim=2)

sin_PROD = torch.column_stack((ones, sin_prod))

PHI = torch.unsqueeze(phi, dim=1).expand(m, n-1, n-1)

mask_ = torch.unsqueeze(torch.triu(torch.ones(n-1, n-1),

diagonal=1).to(device), dim=0).expand(m, n-1, n-1)

mask = mask.expand(m, n-1, n-1)

ones = torch.ones(m).to(device)

torch.unsqueeze(R, dim=1))

x = torch.mul (sin_PROD, cos_R)

return x

sin_PHI = torch.sin(PHI)

810

832 833

834 835

836

837

A.3 HYPERVOLUME ELEMENT IN HYPERSPHERICAL COORDINATES

The hypervolume element of the hypersphere \mathbb{S}_R^{n-1} is given by the following expression when using hyperspherical coordinates (see https://en.wikipedia.org/wiki/N-sphere):

cos_R = torch.mul (torch.column_stack((torch.cos (phi), ones)),

854

$$dV_{\mathbb{S}_R^{n-1}} = R^{n-1} \sin^{n-2}(\varphi_1) \sin^{n-3}(\varphi_2) \cdots \sin(\varphi_{n-2}) d\varphi_1 d\varphi_2 \cdots d\varphi_{n-1}$$
(11)

841 In the small angle regime, where $\sin \varphi \approx \varphi$, we can approximately integrate this expression for an angular coordinate hypercube $[0, \varphi_0]^{n-1}$, and the result is proportional to $v_0 = R^{n-1} \varphi_0^{n(n-1)/2}$. If 842 843 now we reduce the size of the angular coordinate hypercube by a schedule of the form $\varphi_t = \varphi_0 (1 - \varphi_t)$ 844 t), $t \in [0,1]$, then we can compare the percentage of hypervolume being reduced from the initial 845 value, while keeping R fixed, to the percentage obtained by reducing the size of the hypersphere by an schedule of the form $R_t = R(1-t), t \in [0,1]$, while keeping φ_0 fixed (this second case is 846 equivalent to reducing all the Cartesian coordinates at once, because $r^2 = x_n^2 + x_{n-1}^2 + \ldots + x_2^2 + x_1^2$. 847 Indeed, we get, respectively, $v_t = v_0(1-t)^{n(n-1)/2}$ and $v_t = v_0(1-t)^{n-1}$. In Fig. 3 we plot the 848 behavior of v_t/v_0 in terms of the reduction of the coordinate, given by (1-t), for three, increasing 849 values of dimension n. As we can see, already in dimension 20 (bottom figure in the panel), there 850 is a sharp decrease in volume in the angular case as soon as one decreases the angular coordinates 851 by a minimal amount; in comparison to the radial/Cartesian coordinate case, the abrupt decrease in 852 volume looks almost discontinuous. 853

A.4 CONCENTRATION OF MEASURE EFFECTS 855

856 In this appendix, we collect the results of simple experiments that clearly show the concentration of measure effects that occur in high dimensions. In Fig. 4a), we show the distribution of a simple 858 Normal distribution in 2 dimensions (left), and the histogram for the norm of the samples (right). In 859 b), the same but for a Normal distribution in 100 dimensions. In Fig. 5a), we show the histogram for the angle between two random samples from a Normal distribution in 2 dimensions (left), and the same but for a Normal distribution in 100 dimensions (right). In b), we display a schematic diagram 861 of the mass concentration of the uniform measure of the hypersphere in very high dimensions. The 862 intuition in this diagram comes from the more precise result (Wainwright, 2019) which states that, 863 for any given $y \in \mathbb{R}^n$, if we define on the hypersphere an 'equatorial' slice of width $\epsilon > 0$ as



Figure 3: Hypervolume element reduction comparison: $(1 - t)^{n-1}$ vs. $(1 - t)^{n(n-1)/2}$.

 $T_y(\epsilon) \doteq \{z \in \mathbb{S}^{n-1} | (z,y) | \le \epsilon/2\}$, then its volume according to the uniform measure satisfies the following concentration inequality:

$$\mathbb{P}\left[T_y(\epsilon)\right] \ge 1 - \sqrt{2\pi} \exp(-\frac{n\epsilon^2}{2}). \tag{12}$$

The previous inequality shows that, in very high dimensions, the equatorial slice $T_y(\epsilon)$ occupies a huge portion of the total volume, even for a very small width.

Finally, with this in place, we can understand the peculiar shape that a high dimensional Normal distribution takes when expressed in hyperspherical coordinates (Fig.6).





Figure 6: High dimensional Normal distribution in hyperspherical coordinates. For the first three images from the left, each horizontal slice at some vertical index value shows the color coded histogram (red, high density; blue, low density) for the range of the coordinate of that index; the vertical axis stacks all the histograms for all the dimensions (in this example, 40). The white dots represent the mean and the black dots represent the standard deviation of the corresponding histogram. The numbers on top are the total mean and standard deviation of all these previous values taken together.

1026 A.5 ADDITIONAL ANALYSIS OF CIFAR10 RESULTS

Here, we continue the analysis of the experimental results that we obtained for CIFAR10. Fig.7 shows the same results as Fig.2, but we now make a more detailed breakdown of the dependencies of both the MSE and self-FID wrt both the number of latent space dimensions and the total gain β (as in Fig.2, solid lines correspond to the compression model, while dashed lines to the standard one). In the standard VAE, for a fixed β , as we increase the latent dimension, the self-FID increases (worse generation), but the MSE decreases (more sharp, less blurry images); for a fixed latent dimension, as we increase β , the self-FID decreases (better generation), but the MSE increases (less sharp, more blurry images).

Fig.8 shows the typical training of a standard VAE in one of our experimental rounds. In the upper panel we show, from left to right, the histograms of μ , σ , and z, respectively, using the same con-ventions as in Fig.6. The fourth histogram in this panel shows the norm histograms of μ and z, as well as the 'replica angle' (dashed red lines) between the testing samples and the mean for all the test set (this value should give an idea about the angular size of the island as well as to signal if there is an overall replica symmetry breaking in our model; in this particular example, there is no such phase transition, since the mean value of the replica angle is close to $\pi/2$). The second, middle panel shows the behavior of the MSE and KLD loses during training for the test set. The bottom panel corresponds to the histogram of the cosine of the hyperspherical coordinates of μ (cf. Fig.6).

Fig.9 shows the typical training of our compression VAE in one of our experimental rounds. The conventions are the same as in Fig.8. Of note is that the replica angle value in this case shows the desired phase transition. The middle panel shows the annealing schedule used for training our model. Finally, we can see how in the histogram of the cosine of the hyperspherical coordinates all of them are shifted towards a cosine value of 1, which corresponds to an angle equal to 0, as expected.







1242 A.6 The different regimes of the standard β VAE in HD

1244 In this appendix, we illustrate with several examples from our experiments the different regimes in which a β VAE can operate according to the value of the parameter β , while maintaining the 1245 dimension of the latent space fixed but high enough. This is important, since it is known (see, 1246 e.g., Cinelli et al. (2021), section 5.5.2) that HD VAEs are prone to exhibit a phenomenon known 1247 as posterior collapse when β is too high: "[...] [a] state where the variational posterior and true 1248 model posterior collapse to the prior, the posterior encodes no information about the input x, and no 1249 useful latent representation was learned" (quoted from the mentioned reference). This of course, is 1250 a problem, since the collapsed latent dimensions become inoperative for the model and in-utilizable 1251 for other tasks. Furthermore, if used, they can introduce errors in those analysis. 1252

In practice, a simple solution to avoid this issue that often works is to simply reduce the value of β , which acts as a gain for the KLD term in the VAE loss function. One can check for any collapse by inspecting the histograms of the means μ of the latent encoding and making sure that the standard deviation (std) there is appreciably away from zero for each latent dimension. A threshold value can be implemented, but we will keep the discussion qualitative in that aspect.

In Fig.10 we show a standard VAE trained with a high β (= 1.00) in HD (n = 200), it has more 1258 than half of its dimensions collapsed yet the self-FID remains the lowest for the examples (for the 1259 standard VAE, that is) in this dimension as we decrease the β (cf. Fig.7, second row, right; this 1260 is the case for all the dimensions we checked except the lowest, n = 50; see A.12 for this latter 1261 case). Thus, posterior collapse here acts as an effective dimensional reduction mechanism for the 1262 generation, since the collapse actually improves the self-FID profile (we believe that what happens 1263 here is that the weights of the network corresponding to these dimensions are inactive or close to 1264 0 and, therefore, the decoder simply ignores the dimensions in question; see also Dai et al. (2018); 1265 Rolínek et al. (2019)). Nevertheless, since many dimensions are ignored, the model's latent space 1266 lacks representation capacity, which translates into poor reconstructions (MSE = 9.92): the model 1267 works similarly to a non-collapsed one with a much more lower latent dimension.

1268 In Fig.11 we show a standard VAE trained with a medium/balanced β (= 0.20). In this case, 1269 there are more functional dimensions than collapsed or almost collapsed ones. Thus, the model has 1270 more representation capacity and this is reflected in a lower reconstruction error (MSE = 7.12). 1271 Nevertheless, since the decoder now actually operates with a much higher number of dimensions, 1272 then the sparsity and high hypervolume of HD spaces becomes an issue, and this is reflected in 1273 a worse generative performance (higher self-FID than the previous case). In Fig.12, we show a standard VAE trained with a low β (= 0.09) VAE. In this example, the mentioned trends continue 1274 and intensify, now with a much better reconstruction (MSE = 6.32), but very poor generation. 1275

1276

1277

1278

1279

1280

1281 1282

1283

1284

1285

1286

1288

1289

1290

1291

1201

1292

1293

1322



Figure 10: Results of a high β (= 1.00) VAE training (MSE = 9.92, poor). Notice the collapsed dimensions in the histograms for μ (the variance, black dots, for each of those dimensions is very close to 0). Good generation.



Figure 11: Results of a medium/balanced β (= 0.20) VAE training (MSE = 7.12, regular). There are more functional dimensions than collapsed or almost collapsed ones. Regular generation.

Under review as a conference paper at ICLR 2025



Figure 12: Results of a low β (= 0.09) VAE training (MSE = 6.32, good). There are no collapsed dimensions, but the model becomes almost an autoencoder (i.e., the VAE's σ is close to 0). Bad generation.

A.7 AVOIDING POSTERIOR COLLAPSE IS NOT ENOUGH TO IMPROVE GENERATION IN A HD VAE VAE

In this appendix, we show an example in which we encourage the mean of the radial coordinate r^{μ} of the encoded means μ to lie on the hypersphere of radius \sqrt{n} , i.e., $a_{\mu,r} = \sqrt{n}$, and the means of the (cosine) hyperspherical angles ϕ_k^{μ} to lie in the equators, i.e., $a_{\mu,k} = 0$, $\forall k$; furthermore, we also balance the variance of the (cosine) angles ϕ_k^{μ} by encouraging it to be in the same direction as the vector whose Cartesian coordinates are $(1, \ldots, 1)$, i.e., $b_{\mu,k} = 1/\sqrt{k+1}$, $\forall k$. With this setup, our experiments show that posterior collapse is avoided (in both the Cartesian and hyperspherical coordinates representations), while the distribution of μ is still similar to a uniform distribution on the hypersphere, like in the standard VAE (cf. Bardes et al. (2021)). Nevertheless, as expected from the discussion in the previous section, this is not enough to guarantee good generation (Fig.13).



Figure 13: Results of a non-collapsed, non-compressed VAE training. We repeated the experiments for several target values for σ and β , but the results were qualitatively the same as in the present figure.

A.8 HIGH HYPERVOLUME COMPRESSION REDUCES SPARSITY AND IMPROVES GENERATION IN HD VAES

Continuing the analysis of the previous section and figure, then, now that we are sure that we don't 1461 have any collapsed latent dimensions and thus are using the full representation capacity of the HD 1462 space, we can try to improve the poor generation. If our hypothesis about the sparsity introduced by 1463 the exponentially (wrt the dimension) diverging hypervolume in the equators being the root cause of 1464 this issue is true, then by implementing our compression via hyperspherical coordinates we should 1465 be able to improve this generation while remaining on the hyphersphere, un-collapsed and thus re-1466 taining the full expressive capacity of the HD space (unless we compress too much and the excessive 1467 overlap hinders the reconstruction). 1468

In Fig.14 we start with a moderate amount of compression by encouraging the mean of the (cosine)

1470 angles $\dot{\varphi}_k$ to be in the same direction as the vector whose Cartesian coordinates are $(1, \ldots, 1)$, i.e., 1471 $a_{\mu,k} = 1/\sqrt{k+1}, \forall k$. Indeed, recall from appendix A.3 that the closer we get to the north pole, the lower the volume. Nevertheless, this moderate compression is not enough to significantly improve 1472 the generation. Thus, in Fig.15 we go to full compression mode by setting $a_{\mu,k} = 1, \forall k$, which 1473 encourages all the points to converge and condense at the north pole. it is only in this regime of very 1474 high compression that we get a significantly appreciably improvement in the generation. Further-1475 more, we consider this a direct proof of our hypothesis regarding the sparsity of HD spaces and their 1476 impact on generation. In Fig.2 of the main text we showed our experimental results for the more 1477 challenging dataset CIFAR10 regarding how we can use this to systematically take advantage of the 1478 better representation capacity of un-collapsed HD latent spaces to maintain a good and stable recon-1479 struction, while we use our method of volume compression to improve at the same time the quality 1480 of the generation. This allowed us to reach more valuable zones of the MSE-self-FID plane which 1481 are not accessible via the standard VAE in any combination of the parameters n (latent dimension) 1482 and β .

1483 As an additional comment, by looking at the histogram for μ in Cartesian coordinates in Fig.15, 1484 one may think that the lower (in coordinate index) latent dimensions seem heavily collapsed. But 1485 this is not the case: the latent data distribution lies exactly on the hypersphere, and this forces cor-1486 relations in the Cartesian coordinates, reason by which the fact that one or many more Cartesian 1487 coordinates (and their variance) are close to 0 is not conclusive of the irrelevance of many of the la-1488 tent dimensions; indeed, if we now check the histogram for the (cosine) angles φ_k^{μ} in hyperspherical 1489 coordinates (which are a set of uncorrelated coordinates on the hypersphere, by construction), then 1490 we see that there is no collapse in any dimension there. Adding to this point, we can see in Fig.8 that, in the standard VAE, the collapse in Cartesian coordinates (e.g., around index 20 in the first 1491 histogram to the left in the first row) translates into a collapse in the (cos) hyperspherical coordinates 1492 (third row histogram, same index), while this is not the case in our compression VAE in Fig.9, where 1493 the apparent collapse in Cartesian coordinates around, e.g., index 20, doesn't translate into an anal-1494 ogous collapse in the (cos) hyperspherical coordinates: we believe that the reason for this is that, in 1495 the standard VAE, we are not exactly on the hypersphere (in the fourth histogram to the right in the 1496 first row in Fig.8, we can see that the norm of μ , in orange, has a non-zero variance, since the prior 1497 is still a multivariate Gaussian, not exactly a uniform distribution on the hypersphere), while our 1498 compression VAE is indeed exactly on the hypersphere (analogous norm histogram in Fig.9), since 1499 we explicitly encourage the variance of the radial coordinate of μ to be 0. Thus, we emphasize that 1500 the improvements in generation by our compression method cannot be explained by selective pos-1501 terior collapse (as in Fig.10), where the HD collapsed latent representation is effectively equivalent 1502 to a non-collapsed one in lower dimensions, since this comes at the cost of loosing reconstruction quality; but our method is able to improve generation while retaining some amount of better recon-1503 struction, and this is why some of the best performative examples in Fig.2 cannot be re-obtained by a 1504 standard VAE with a different combination of parameters n and β (possibly in a selective collapsed 1505 mode). The improvement in our method is coming from the reduction of the sparsity by compression 1506 of the latent hypervolume and by performing this in a key angular way due to the peculiar equatorial 1507 nature of the volume in HD spaces. 1508

1509

1510



Figure 14: Results of a moderately compressed VAE training.



Figure 15: Results of a fully compressed VAE training.

1620 A.9 THE SPIN GLASS ANALOGY DURING TRAINING

1622 We described in Section 4.1 of the main text the following training schedule and the reasons behind this choice: "[...] we use an annealing schedule for the gain β of the KLD-like loss, consisting of 1623 an initial stage which increases proportionally with $\sqrt{\text{epoch}}$ for the first 100 epochs, and is constant 1624 afterwards. This was necessary because we observed that too much compression of the volume 1625 was detrimental to the performance, while a strong compression was still necessary at the initial 1626 stage[...]". The gain β here has the role² of the inverse temperature, $\beta = 1/T$. In spin glasses and 1627 complex systems, the energy function has exponentially many local minima in the equatorial region 1628 of the hypersphere. To overcome them, a very strong signal or bias towards the desired region is 1629 necessary at the beginning, together with a rapid cooling or quenching. Thus, our initial high β (i.e., 1630 very low temperature T) setting, and in the presence of the high intensity (regulated by the β^{-1} factor in front of the MSE) hyperspherical external magnetic fields as bias in directions away from the equator, should make the gradient descent dynamics to quickly tend towards a low temperature 1633 distribution with replica symmetry breaking. Indeed, this is what we observed in our experiments, since we check for the replica angle, as mentioned before. This initial strong compression helps 1634 escaping those undesirable equatorial minima (Fig.16). Nevertheless, the obtained state shows too 1635 much overlapping between samples, so we then perform the annealing (i.e., lower the β , or increase the temperature T, and also lower the intensity of the magnetic fields) in order to allow the system to 1637 relax the strong order introduced by the initial bias and, in this way, transition to a replica symmetry 1638 breaking state with a bigger angle between replicas (that is, to go back up a bit in the ultrametricity 1639 tree/hierarchy of the replica angle values; cf. Mourrat (2024)). This decreases the MSE and makes 1640 the decoded images more sharp, at the cost of some generation quality (Fig.17). Note how the replica 1641 angle (red dashed lines in fourth histogram to the left in second row) doesn't fully go back to $\pi/2$, 1642 even when the KLD term (where the external magnetic fields are) stops optimizing at this stage of 1643 the training process (red line in third row), but instead jumps to a different value, higher than the 1644 initial one but still below $\pi/2$. This is fully consistent with the spin glass analogy in a quenched and then annealed system, where the glass, always in the replica symmetry breaking phase, jumps 1645 from one so-called 'pure state' to a different pure state, i.e., goes back up a bit in the ultrametricity 1646 tree/hierarchy of the replica angle values, as mentioned before. But the system has escaped the zone 1647 with exponentially many local minima in the equator. 1648 1649

- 1649 1650
- 16

1652

1654

165

1656 1657

1658

1659 1660

10

1663 1664

100

1667

1669

1670

1671

1672 1673 $\frac{{}^{2}\mathcal{L} = \beta \left(\beta^{-1} \text{MSE}(x, x_{z}) + \text{KLD}_{\text{HSphCoords}}^{w/Prior}(\varphi_{k}, r) \right) = \beta \mathcal{H}. \text{ cf. footnote 1, where } \nabla \mathcal{L} = \beta \nabla \mathcal{H} \text{ for the gradient descent dynamics on } \mathcal{L}.$







A.10 RESULTS ON CELEBA64

In this appendix we include additional experimental results conducted on the dataset CelebA (Liu et al., 2015), resized to a 64×64 image size.

The analysis is of the same type as the one we performed on CIFAR10 (cf. Figs.2, 7), and the results show qualitatively the same trends (Figs.19, 20).



Figure 18: Effect of latent dimension and β on the trade off between reconstruction and generation on CelebA64 (as in CIFAR10, solid lines closer to the bottom left corner than the dashed lines).





Figure 20: Results of standard VAE training with a balanced β (in the 3-D embedding diagram, the samples are normalized by the overall mean of the radial coordinate, rather than set exactly to the sphere; thus, rather than looking like a uniform-like distribution on the 2-D sphere, it looks like a normal distribution in 3-D, but this difference is only merely in the convention being used regarding the radial normalization for the 3-D embedding). MSE = 5.43 and self - FID = 34.94, n = 600.



n = 600.

1998 A.11 INTERPOLATIONS ON MNIST

Here we complement the results and associated claims of Fig.1 with interpolations experiments on the same models. They highlight the lack of continuity in the standard VAE case (Fig.22), while they show the gained continuity and how densely packed the clusters are in our compressed version (Fig.23).



Figure 22: Interpolations on the standard VAE. a) from 0 to 1; b) from 7 to 2; c) from 0 to 4.



Figure 23: Interpolations on the compression VAE. a)-b) Idem as previous figure.

A.12 The different regimes of the standard β VAE in low dimensions

In this appendix, we perform a similar analysis as the one in A.6, but now for the model with the lowest latent dimension (n = 50).

In this situation, the trends actually reverse: de-collapsing the model (that is, going from Figs.24 to 25 and so on) improves the generation as measured by the self-FID (cf. Figs.10 to 11 and so on, in the HD case, where it becomes worse). See also Fig.7, second row, right.

Nevertheless, these cases are pathological and not very useful, since all of them have very high MSE, that is, the images are too 'blurry'. Thus, both the collapsed and the non-collapsed cases fall into the bottom far right of Fig.2, way outside the more useful area of the MSE-self-FID plane.





Figure 24: Results of a high β (= 1.00) VAE training (MSE = 12.27, very poor). Notice the collapsed dimensions in the histograms for μ (the variance, black dots, for each of those dimensions is very close to 0). Good generation.



Figure 25: Results of a medium/balanced β (= 0.20) VAE training (MSE = 9.97, poor). There are more functional dimensions than collapsed or almost collapsed ones. Better generation than previous figure.



Figure 26: Results of a low β (= 0.09) VAE training (MSE = 9.59, poor). There are no collapsed dimensions, but the model becomes almost an autoencoder (i.e., the VAE's σ is close to 0). Even better generation than previous figure.

A.13 RE-WRITING OF THE KLD TERM

In this appendix, we make explicit the steps to go from the standard form of the KLD term in the VAE to the one we used as a starting point for our own KLD in hyperspherical coordinates.

In Cartesian coordinates, the KLD divergence between the estimated posterior defined by μ_k and σ_k and the prior defined by μ_k^p and σ_k^p is (Odaibo, 2019):

$$\operatorname{KLD}_{\operatorname{CartCoords}}^{w/Prior} = \frac{1}{2} \sum_{k=1}^{n} \left[\left(\frac{\sigma_k}{\sigma_k^p} \right)^2 - \log \left(\frac{\sigma_k}{\sigma_k^p} \right)^2 - 1 + \frac{\left(\mu_k - \mu_k^p \right)^2}{\left(\sigma_k^p \right)^2} \right]$$
(13)

A Taylor approximation (up to second order) of the part for sigma around its prior yields for some constants γ_k and $\widetilde{\gamma}_k$:

$$\operatorname{KLD}_{\operatorname{CartCoords}}^{w/Prior} \approx \sum_{k=1}^{n} \left[\gamma_k \left(\sigma_k - \sigma_k^p \right)^2 + \widetilde{\gamma}_k \left(\mu_k - \mu_k^p \right)^2 \right]$$
(14)

In practice, the optimization is performed over mini batches of data (of size N_b), using the objective below:

$$\operatorname{KLD}_{\operatorname{CartCoords}}^{w/Prior} \approx \frac{1}{N_b} \sum_{l=1}^{N_b} \sum_{k=1}^n \left(\gamma_k \left(\sigma_{k,l} - \sigma_k^p \right)^2 + \widetilde{\gamma}_k \left(\mu_{k,l} - \mu_k^p \right)^2 \right)$$
(15)

If we denote the corresponding batch statistics as \mathbb{E}_b and σ_b , then, by using the basic formula,

$$\mathbb{E}_{b}[X^{2}] = \mathbb{E}_{b}[X]^{2} + \sigma_{b}[X]^{2}, \tag{16}$$

we can write this objective as (we omit the constants for ease of reading)

$$\operatorname{KLD}_{\operatorname{CartCoords}}^{w/Prior} \approx \sum_{k=1}^{n} \left(\left(\mathbb{E}_{b}[\sigma_{k}] - \sigma_{k}^{p} \right)^{2} + \sigma_{b}[\sigma_{k}]^{2} + \left(\mathbb{E}_{b}[\mu_{k}] - \mu_{k}^{p} \right)^{2} + \sigma_{b}[\mu_{k}]^{2} \right)$$
(17)



This is a complete analogue of A.6 but for the CelebA64 dataset with n = 1000.



Figure 27: Results of a high β (= 1.00) VAE training (MSE = 5.22, good). Regular to bad generation.



Figure 28: Results of a medium/balanced β (= 0.20) VAE training (MSE = 3.74, good). Bad generation.



Figure 29: Results of a low β (= 0.09) VAE training (MSE = 3.31, good to very good). Very poor generation.

A.15 IT IS THE SPARSITY!

2323

In this apppendix, we collect further experimental evidence, from different directions, that complements our findings in Fig.1, as well as the discussion in A.7, A.8, and A.11, that supports our hypothesis about the sparsity of HD latent spaces as the root cause of the poor generation for the standard VAE.

Continuing the analysis of A.7 and Fig.13, to further confirm that this poor generation is caused by the sparsity, we started with a very low value for the radial mean of the VAE's σ , from $a_{\sigma,r} = 0.1 \times \sqrt{n}$ (this is similar to an autoencoder) in Fig.30, to $a_{\sigma,r} = 1.0 \times \sqrt{n}$ (this is similar to a VAE regime) in Fig.13, and finally to $a_{\sigma,r} = 2.0 \times \sqrt{n}$ (this is similar to a VAE regime with very wide encoded latent distributions for each encoded data point x) in Fig.31. Actually, in this latter case, after the reparameterization trick and normalization of z_x to the hypersphere, z_x could be, for each x, in *any* location on the hypersphere, not just a portion of it.

The logic behind this is that the bigger $a_{\sigma,r}$ is, then more of the latent space is covered during 2336 training because of the reparameterization trick; thus, if sparsity is the cause of the poor generation, 2337 we should be able to see some improvement as we increase $a_{\sigma,r}$. And this is indeed what our 2338 experiments show. Nevertheless, the amount of improvement in generation by this method quickly 2339 reaches a limit for the extreme case $a_{\sigma,r} = 2.0 \times \sqrt{n}$, since there is probably too much overlapping 2340 between the encoded distributions for all the different datapoints. On the other hand, our method 2341 which consists in the angular coordinates compression of the distributions (rather than the previous 2342 radial expansion of σ) gives better, more consistent and efficient results (we believe this is because 2343 of the natural angular behavior of the hypervolume in HD; cf. A.3, A.4), cf. Fig.15.

2344 In Fig.32, we show a standard VAE with decoded random samples from the aggregate approximate 2345 posterior distribution, rather than the prior. We can see that there is no much improvement in the 2346 generation, and thus this cannot be attributed to a mismatch between the prior and aggregate ap-2347 proximate posterior distribution in this case (as in some of the references mentioned in 2). In Fig.14 2348 of A.8, we take a moderately compressed model (by our method) and show the 3D-embedding of 2349 the approximate aggregate posterior, which is a von Mises-Fisher-like distribution, and it shows no holes or cracks (since our method, by compressing the data, makes the clusters more packed and 2350 closer to each other; thus, it could also be applied to solve or mitigate the prior hole problem, but we 2351 will leave it at that for now, since the main point of this paper is the sparsity issue in HD). If the phi-2352 losophy of the mismatch between prior and posterior as root cause of bad generation were true, then 2353 sampling from such an approximate aggregate posterior von Mises-Fisher-like distribution should 2354 produce better generation. But we do exactly that in that figure and we don't see any improvement. 2355 It's only when we do a full compression mode that we get the improvement in generation. 2356

In Fig.33, we computed the *persistence diagrams* corresponding to the experiments in Figs.13, 14, 15 (*a*), *b*), *c*) in Fig.33, respectively). We used the GUDHI library (Maria et al.; https://gudhi.inria.fr/) for producing the persistence diagrams presented here.

2360 This is a key tool from Topological Data Analysis (TDA) (Carlsson, 2009), which allows us to 2361 explicitly and directly assess the precise amount of sparsity in our HD latent spaces. Consider the encoded means μ of our dataset in latent space as a point cloud, then we take a cover of n-balls of a 2362 fixed radius ρ around each point in the cloud. At 'time' of radius $\rho = 0$, all the balls are obviously 2363 non-intersecting with each other, this is called the 'birth time' B_{μ} of the points μ (thus, all points 2364 have birth time $B_{\mu} = 0$). We then start increasing the value of ρ with 'time'. If at some moment of 2365 radius ρ_{μ} the ball around a point μ intersects other balls, this is called the 'death time' of that point, 2366 D_{μ} , and has the value $D_{\mu} = \varrho_{\mu}$. In this way, each point μ is characterized here by its corresponding 2367 pair $(B_{\mu}, D_{\mu}) = (0, \varrho_{\mu})$, and it is these pairs that we plot in the persistence diagrams, for all the 2368 points in the cloud. 2369

Thus, it is not difficult to see why this directly assess the sparsity of the cloud: the longer the death times for a given cloud, this means that the points are more separated from each other and that the point cloud is more sparse.

In this way, we can see a very interesting correlation between the sparsity measured in this manner and the quality of the generation. When we go from Fig.13 to 14, we see that there is indeed some compression (this can be seen in the replica angle and the 3D-embedding in those figures), but there is no appreciable improvement in the generation. In the corresponding persistence diagrams of Fig.33 (*a*) and *b*), respectively), we can see that there is a small diminution of the death times (of around $\Delta \rho_{\mu} \approx 1000$ for each point) between diagrams. Nevertheless, the difference in the death times between the shortest living point and the longest one in the same diagram remains almost unchanged, around $\Delta \rho \approx 9000$ for each diagram. Thus, this moderate compression only achieves a reduction in death times by a constant translation factor of $\Delta \rho_{\mu} \approx 1000$ on all death times.

It is only when we pass to the full compression mode, Fig.15, that we see an appreciable improve-ment in the generation. In the corresponding persistence diagram of Fig.33 (c), we can see that there is a considerable greater reduction in the life time of the longest living point, from around $D_{\mu_{max}} \approx 10000$ in the previous two diagrams, to around $D_{\mu_{max}} \approx 2000$. Furthermore, we also get a considerable decrease of the difference in the death times between the shortest living point and the longest one in the same diagram, from around $\Delta \rho \approx 9000$ in the previous diagrams, to around $\Delta \rho \approx 2000$ in the diagram in c). Thus, our analysis here shows that this dramatic reduction in sparsity is strongly correlated to the appreciable improvement in the quality of the generation. It was the sparsity!



Figure 30: Hyperspherical VAE with $a_{\sigma,r} = 0.1 \times \sqrt{n}$.



Figure 31: Hyperspherical VAE with $a_{\sigma,r} = 2.0 \times \sqrt{n}$ (histograms for σ out of scale here).



Figure 32: Standard VAE with decoded random samples from the aggregate approximate posterior distribution.



47

2538 A.16 DIFFERENCES IN TRAINING SPEED

2549

We provide here data regarding the differences in the training speeds between the standard VAE and our compression VAE via hyperspherical coordinates. The origin of this difference mainly lies in the extra calculations needed for the coordinate transformations in A.1, which are implemented via the script in A.2.

The measurements were done during typical trainings in a NVIDIA H100 GPU. In Fig.34 we show the results for the case of trainings with CIFAR10, with a batch size of 200 samples, and the changes in training speed (measured as how many batches per second are being processed) in terms of the dimension n of the latent space. After n = 200, until n = 800, the decay is almost linear in n, with a decay rate in the speed of 20 batch/s every 200 latent dimensions.

