
Appendix: A Wrapped Normal Distribution on Hyperbolic Space for Gradient-Based Learning

A. Derivations

A.1. Inverse Exponential Map

As we mentioned in the main text, the exponential map from $T_{\boldsymbol{\mu}}\mathbb{H}^n$ to \mathbb{H}^n is given by

$$\boldsymbol{z} = \exp_{\boldsymbol{\mu}}(\boldsymbol{u}) = \cosh(\|\boldsymbol{u}\|_{\mathcal{L}})\boldsymbol{\mu} + \sinh(\|\boldsymbol{u}\|_{\mathcal{L}})\frac{\boldsymbol{u}}{\|\boldsymbol{u}\|_{\mathcal{L}}}. \quad (1)$$

Solving (1) for \boldsymbol{u} , we obtain

$$\boldsymbol{u} = \frac{\|\boldsymbol{u}\|_{\mathcal{L}}}{\sinh(\|\boldsymbol{u}\|_{\mathcal{L}})}(\boldsymbol{z} - \cosh(\|\boldsymbol{u}\|_{\mathcal{L}})\boldsymbol{\mu}).$$

We still need to obtain the evaluable expression for $\|\boldsymbol{u}\|_{\mathcal{L}}$. Using the characterization of the tangent space (main text, (2)), we see that

$$\begin{aligned} \langle \boldsymbol{\mu}, \boldsymbol{u} \rangle_{\mathcal{L}} &= \frac{\|\boldsymbol{u}\|_{\mathcal{L}}}{\sinh(\|\boldsymbol{u}\|_{\mathcal{L}})} \left(\langle \boldsymbol{\mu}, \boldsymbol{z} \rangle_{\mathcal{L}} - \cosh(\|\boldsymbol{u}\|_{\mathcal{L}}) \langle \boldsymbol{\mu}, \boldsymbol{\mu} \rangle_{\mathcal{L}} \right) = 0, \\ \cosh(\|\boldsymbol{u}\|_{\mathcal{L}}) &= -\langle \boldsymbol{\mu}, \boldsymbol{z} \rangle_{\mathcal{L}}, \\ \|\boldsymbol{u}\|_{\mathcal{L}} &= \operatorname{arccosh}(-\langle \boldsymbol{\mu}, \boldsymbol{z} \rangle_{\mathcal{L}}). \end{aligned}$$

Now, defining $\alpha = -\langle \boldsymbol{\mu}, \boldsymbol{z} \rangle_{\mathcal{L}}$, we can obtain the inverse exponential function as

$$\boldsymbol{u} = \exp_{\boldsymbol{\mu}}^{-1}(\boldsymbol{z}) = \frac{\operatorname{arccosh}(\alpha)}{\sqrt{\alpha^2 - 1}}(\boldsymbol{z} - \alpha\boldsymbol{\mu}). \quad (2)$$

A.2. Inverse Parallel Transport

The parallel transportation on the Lorentz model along the geodesic from $\boldsymbol{\nu}$ to $\boldsymbol{\mu}$ is given by

$$\begin{aligned} \operatorname{PT}_{\boldsymbol{\nu} \rightarrow \boldsymbol{\mu}}(\boldsymbol{v}) &= \boldsymbol{v} - \frac{\langle \exp_{\boldsymbol{\nu}}^{-1}(\boldsymbol{\mu}), \boldsymbol{v} \rangle_{\mathcal{L}}}{d_{\ell}(\boldsymbol{\nu}, \boldsymbol{\mu})^2} (\exp_{\boldsymbol{\nu}}^{-1}(\boldsymbol{\mu}) + \exp_{\boldsymbol{\mu}}^{-1}(\boldsymbol{\nu})) \\ &= \boldsymbol{v} + \frac{\langle \boldsymbol{\mu} - \alpha\boldsymbol{\nu}, \boldsymbol{v} \rangle_{\mathcal{L}}}{\alpha + 1} (\boldsymbol{\nu} + \boldsymbol{\mu}), \end{aligned} \quad (3)$$

where $\alpha = -\langle \boldsymbol{\nu}, \boldsymbol{\mu} \rangle_{\mathcal{L}}$. Next, likewise, for the exponential map, we need to be able to compute the inverse of the parallel transform. Solving (3) for \boldsymbol{v} , we get

$$\boldsymbol{v} = \boldsymbol{u} - \frac{\langle \boldsymbol{\mu} - \alpha\boldsymbol{\nu}, \boldsymbol{v} \rangle_{\mathcal{L}}}{\alpha + 1} (\boldsymbol{\nu} + \boldsymbol{\mu}).$$

Now, observing that

$$\begin{aligned} \langle \boldsymbol{\nu} - \alpha\boldsymbol{\mu}, \boldsymbol{u} \rangle_{\mathcal{L}} &= \langle \boldsymbol{\nu}, \boldsymbol{v} \rangle_{\mathcal{L}} + \frac{\langle \boldsymbol{\mu} - \alpha\boldsymbol{\nu}, \boldsymbol{v} \rangle_{\mathcal{L}}}{\alpha + 1} (\langle \boldsymbol{\nu}, \boldsymbol{\nu} \rangle_{\mathcal{L}} + \langle \boldsymbol{\mu}, \boldsymbol{\nu} \rangle_{\mathcal{L}}) \\ &= -\langle \boldsymbol{\mu}, \boldsymbol{v} \rangle_{\mathcal{L}} = -\langle \boldsymbol{\mu} - \alpha\boldsymbol{\nu}, \boldsymbol{v} \rangle_{\mathcal{L}}, \end{aligned}$$

we can write the inverse parallel transport as

$$\mathbf{v} = \text{PT}_{\nu \rightarrow \mu}^{-1}(\mathbf{u}) = \mathbf{u} + \frac{\langle \nu - \alpha \mu, \mathbf{u} \rangle_{\mathcal{L}}}{\alpha + 1} (\nu + \mu). \quad (4)$$

The inverse of parallel transport from ν to μ coincides with the parallel transport from μ to ν .

A.3. Determinant of exponential map

We provide the details of the computation of the log determinant of \exp_{μ} . Let $\mu \in \mathbb{H}^n$, let $\mathbf{u} \in T_{\mu}(\mathbb{H}^n)$, and let $\mathbf{v} = \exp_{\mu}(\mathbf{u})$. Then the derivative is a map from the tangent space of $T_{\mu}(\mathbb{H}^n)$ at \mathbf{u} to the tangent space of \mathbb{H}^n at \mathbf{v} . The determinant of this derivative will not change by any orthogonal change of basis. Let us choose an orthonormal basis of $T_{\mathbf{u}}(T_{\mu}(\mathbb{H}^n)) \cong T_{\mu}(\mathbb{H}^n)$ containing $\bar{\mathbf{u}} = \mathbf{u}/\|\mathbf{u}\|_{\mathcal{L}}$:

$$\{\bar{\mathbf{u}}, \mathbf{u}'_1, \mathbf{u}'_2, \dots, \mathbf{u}'_{n-1}\}$$

The desired determinant can be computed by tracking how much each element of this basis grows in magnitude under the transformation.

The derivative in the direction of each basis element can be computed as follows:

$$\begin{aligned} d \exp_{\mu}(\bar{\mathbf{u}}) &= \left. \frac{\partial}{\partial \epsilon} \right|_{\epsilon=0} \exp_{\mu}(\mathbf{u} + \epsilon \bar{\mathbf{u}}) \\ &= \left. \frac{\partial}{\partial \epsilon} \right|_{\epsilon=0} \left[\cosh(r + \epsilon) \mu + \sinh(r + \epsilon) \frac{\mathbf{u} + \epsilon \bar{\mathbf{u}}}{\|\mathbf{u} + \epsilon \bar{\mathbf{u}}\|_{\mathcal{L}}} \right] \\ &= \sinh(r) \mu + \cosh(r) \bar{\mathbf{u}}. \end{aligned} \quad (5)$$

$$\begin{aligned} d \exp_{\mu}(\mathbf{u}'_k) &= \left. \frac{\partial}{\partial \epsilon} \right|_{\epsilon=0} \exp_{\mu}(\mathbf{u} + \epsilon \mathbf{u}'_k) \\ &= \left. \frac{\partial}{\partial \epsilon} \right|_{\epsilon=0} \left[\cosh(r) \mu + \sinh(r) \frac{\mathbf{u} + \epsilon \mathbf{u}'_k}{r} \right] \\ &= \frac{\sinh r}{r} \mathbf{u}'_k. \end{aligned} \quad (6)$$

In the second line of the computation of the directional derivative with respect to \mathbf{u}'_k , we used the fact that $\|\mathbf{u} + \epsilon \mathbf{u}'_k\|_{\mathcal{L}} = \sqrt{\langle \mathbf{u}, \mathbf{u} \rangle_{\mathcal{L}} + \epsilon \langle \mathbf{u}, \mathbf{u}'_k \rangle_{\mathcal{L}} + \epsilon^2 \langle \mathbf{u}'_k, \mathbf{u}'_k \rangle_{\mathcal{L}}} = \|\mathbf{u}\|_{\mathcal{L}} + \mathcal{O}(\epsilon^2)$ and that $\mathcal{O}(\epsilon^2)$ in the above expression will disappear in the $\epsilon \rightarrow 0$ limit of the finite difference. All together, the derivatives computed with respect to our choice of the basis elements are given by

$$\left(\sinh(r) \mu + \cosh(r) \bar{\mathbf{u}}, \frac{\sinh r}{r} \mathbf{u}'_1, \frac{\sinh r}{r} \mathbf{u}'_2, \dots, \frac{\sinh r}{r} \mathbf{u}'_{n-1} \right)$$

The desired determinant is the product of the Lorentzian norms of the vectors of the set above. Because all elements of $T_{\mu}(\mathbb{H}^n)$ are orthogonal with respect to the Lorentzian inner product and because $\|\sinh(r) \mu + \cosh(r) \bar{\mathbf{u}}\|_{\mathcal{L}} = 1$ and $\|\sinh(r)/r \cdot \mathbf{u}'_k\|_{\mathcal{L}} = \sinh(r)/r$, we get

$$\det \left(\frac{\partial \exp_{\mu}(\mathbf{u})}{\partial \mathbf{u}} \right) = \left(\frac{\sinh r}{r} \right)^{n-1}. \quad (7)$$

A.4. Determinant of parallel transport

Next, let us compute the determinant of the parallel transport. Let $\mathbf{v} \in T_{\mu_0} \mathbb{H}^n$, and let $\mathbf{u} = \text{PT}_{\mu_0 \rightarrow \mu}(\mathbf{v}) \in T_{\mu} \mathbb{H}^n$. The derivative of this map is a map from $T_{\mathbf{v}}(T_{\mu_0}(\mathbb{H}^n))$ to $T_{\mathbf{u}}(\mathbb{H}^n)$. Let us choose an orthonormal basis ξ_k (In Lorentzian sense). Likewise above, we can compute the desired determinant by tracking how much each element of this basis grows in magnitude under the transformation.

Denoting $\alpha = -\langle \mu_0, \mu \rangle_{\mathcal{L}}$, we get

$$\begin{aligned}
 d \text{PT}_{\mu_0 \rightarrow \mu}(\xi) &= \left. \frac{\partial}{\partial \epsilon} \right|_{\epsilon=0} \text{PT}_{\mu_0 \rightarrow \mu}(\mathbf{v} + \epsilon \xi) \\
 &= \left. \frac{\partial}{\partial \epsilon} \right|_{\epsilon=0} \left[(\mathbf{v} + \epsilon \xi) + \frac{\langle \mu - \alpha \mu_0, \mathbf{v} + \epsilon \xi \rangle_{\mathcal{L}}}{\alpha + 1} (\mu_0 + \mu) \right] \\
 &= \xi + \frac{\langle \mu - \alpha \mu_0, \xi \rangle_{\mathcal{L}}}{\alpha + 1} (\mu_0 + \mu) = \text{PT}_{\mu_0 \rightarrow \mu}(\xi).
 \end{aligned} \tag{8}$$

and see that each basis element ξ_k is mapped by $d \text{PT}_{\mu_0 \rightarrow \mu}$ to

$$(\text{PT}_{\mu_0 \rightarrow \mu}(\xi_1), \text{PT}_{\mu_0 \rightarrow \mu}(\xi_2) \cdots, \text{PT}_{\mu_0 \rightarrow \mu}(\xi_n))$$

Because parallel transport is a norm preserving map, $\|\text{PT}_{\mu_0 \rightarrow \mu}(\xi)\|_{\mathcal{L}} = 1$. That is,

$$\det \left(\frac{\partial \text{PT}_{\mu_0 \rightarrow \mu}(\mathbf{v})}{\partial \mathbf{v}} \right) = 1. \tag{9}$$

B. Visual Examples of Hyperbolic Wrapped Distribution $\mathcal{G}(\mu, \Sigma)$

Figure 1 shows examples of hyperbolic wrapped distribution $\mathcal{G}(\mu, \Sigma)$ with various μ and Σ . We plotted the log-density of these distributions by heatmaps. We designate the μ by the \times mark. The right side of these figures expresses their log-density on the Poincaré ball model, and the left side expresses the same one on the corresponding tangent space.

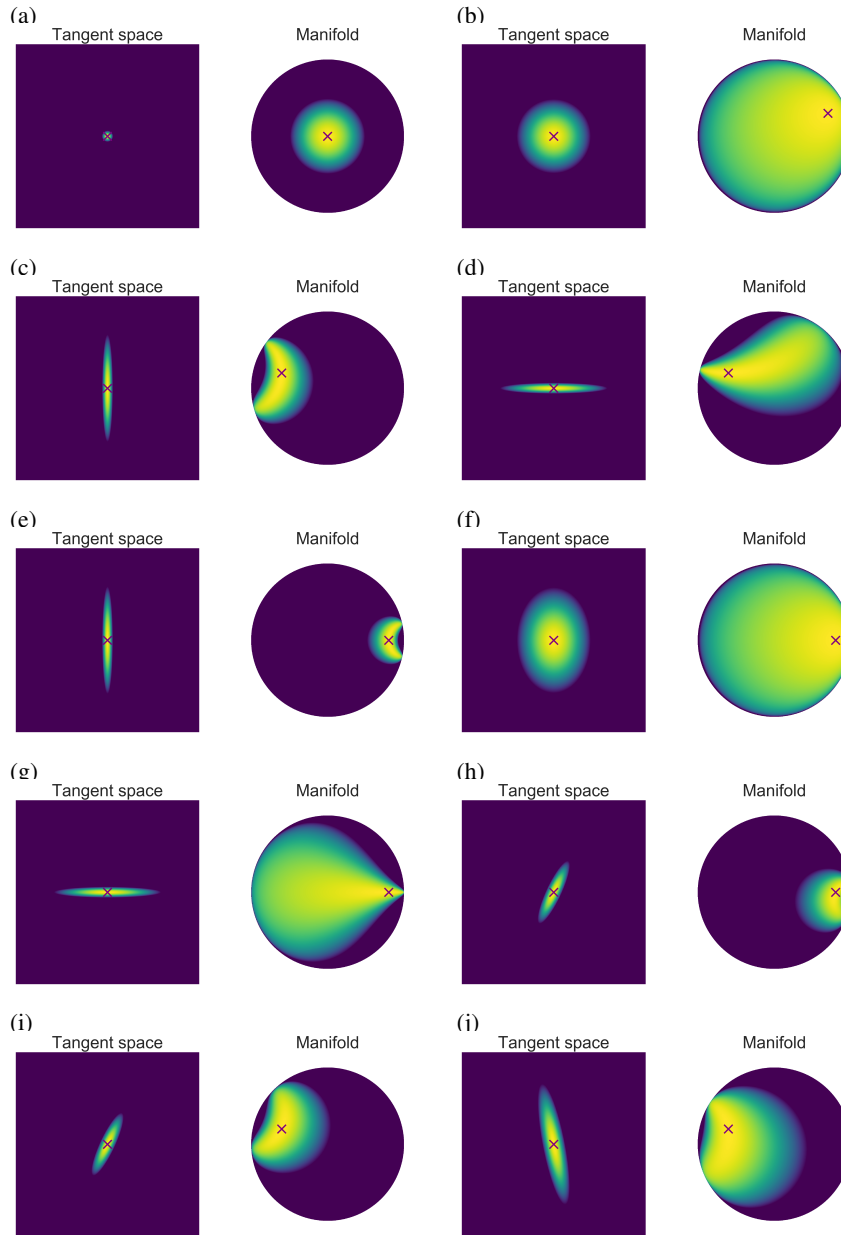


Figure 1: Visual examples of hyperbolic wrapped distribution on \mathbb{H}^2 . Log-density is illustrated on \mathcal{B}^2 by translating each point from \mathbb{H}^2 for clarity. We designate the origin of hyperbolic space by the \times mark.

C. Additional Numerical Evaluations

C.1. Synthetic Binary Tree

We qualitatively compared the learned latent space of Vanilla and Hyperbolic VAEs. Figure 2 shows the embedding vectors of the synthetic binary tree dataset on the two-dimensional latent space. We evaluated the latent space of Vanilla VAE with $\beta = 0.1, 1.0, 2.0,$ and $3.0,$ and Hyperbolic VAE. Note that the hierarchical relations in the original tree were **not** used during the training phase. Red points are the embeddings of the noiseless observations. As we mentioned in the main text, we evaluated the correlation coefficient between the Hamming distance on the data space and the hyperbolic (Euclidean for Vanilla VAEs) distance on the latent space. Consistently with this metric, the latent space of the Hyperbolic VAE captured the hierarchical structure inherent in the dataset well. In the comparison between Vanilla VAEs, the latent space captured the hierarchical structure according to increase the β . However, the posterior distribution of the Vanilla VAE with $\beta = 3.0$ collapsed and lost the structure. Also, the blue points are the embeddings of noisy observation, and pink \times represents the origin of the latent space. In latent space of Vanilla VAEs, there was bias in which embeddings of noisy observations were biased to the center side.

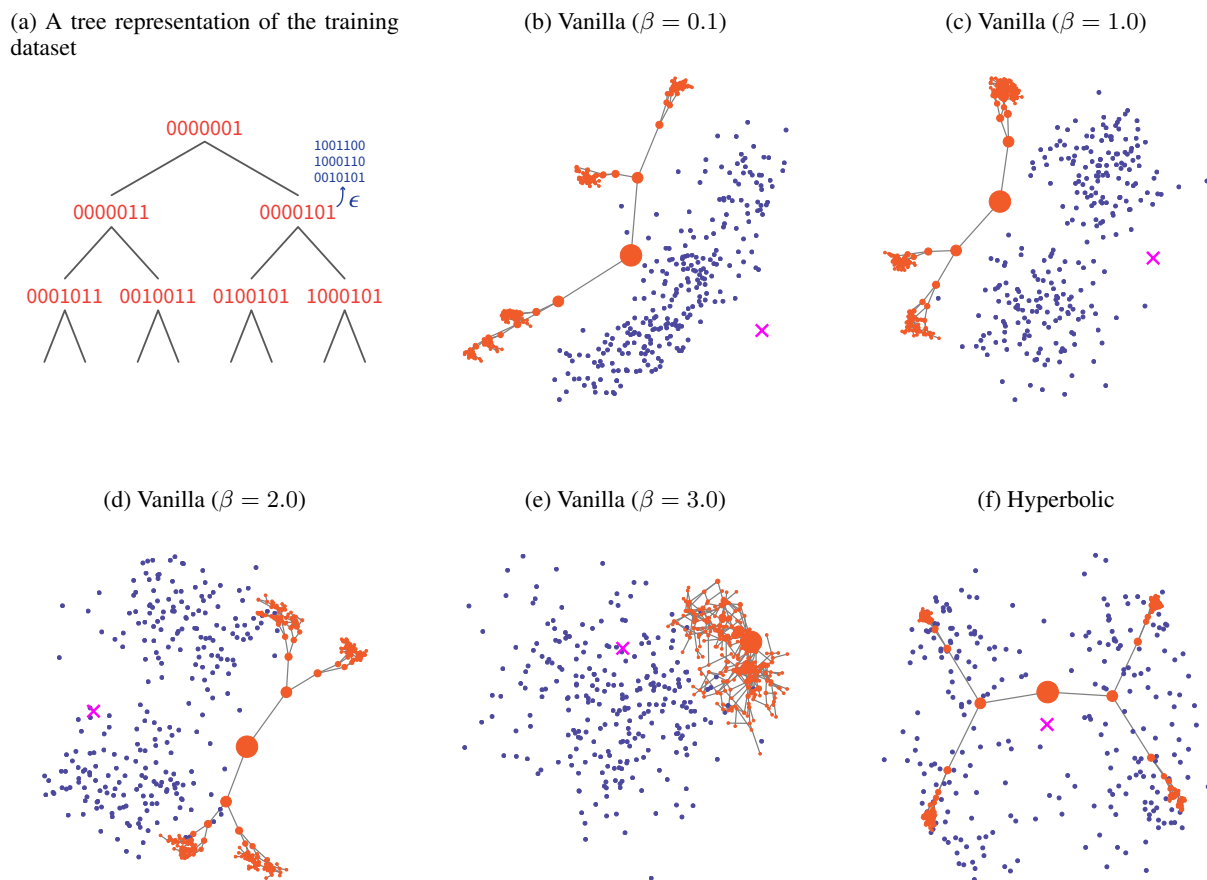


Figure 2: The visual results of Vanilla and Hyperbolic VAEs applied to an artificial dataset generated by applying a random perturbation to a binary tree. The visualization is being done in the Poincaré ball. Red points are the embeddings of the original tree, and the blue points are the embeddings of all other points in the dataset. Pink \times represents the origin of hyperbolic space. Note that the hierarchical relations in the original tree was **not** used during the training phase.

C.2. MNIST

n	Vanilla VAE		Hyperbolic VAE	
	ELBO	LL	ELBO	LL
2	-145.53±.65	-140.45±.47	-143.23±0.63	-138.61±0.45
5	-111.32±.38	-105.78±.51	-111.09±0.39	-105.38±0.61
10	-92.49±.52	-86.25±.52	-93.10±0.26	-86.40±0.28
20	-85.17±.40	-77.89±.36	-88.28±0.34	-79.23±0.20

Table 1: Quantitative comparison of Hyperbolic VAE against Vanilla VAE on the MNIST dataset in terms of ELBO and log-likelihood (LL) for several values of latent space dimension n . LL was computed using 500 samples of latent variables. We calculated the mean and the ± 1 SD with five different experiments.

We showed the numerical performance of Vanilla and Hyperbolic VAEs for MNIST data in the main text in terms of the log-likelihood. In this section, we also show the evidence lower bound for the same dataset in Table 1.

C.3. Atari 2600 Breakout

To evaluate the performance of Hyperbolic VAE for hierarchically organized dataset according to time development, we applied our Hyperbolic VAE to a set of trajectories that were explored by an agent with a trained policy during multiple episodes of Breakout in Atari 2600. We used a pretrained Deep Q-Network to collect trajectories, and Figure 3 shows examples of observed screens.

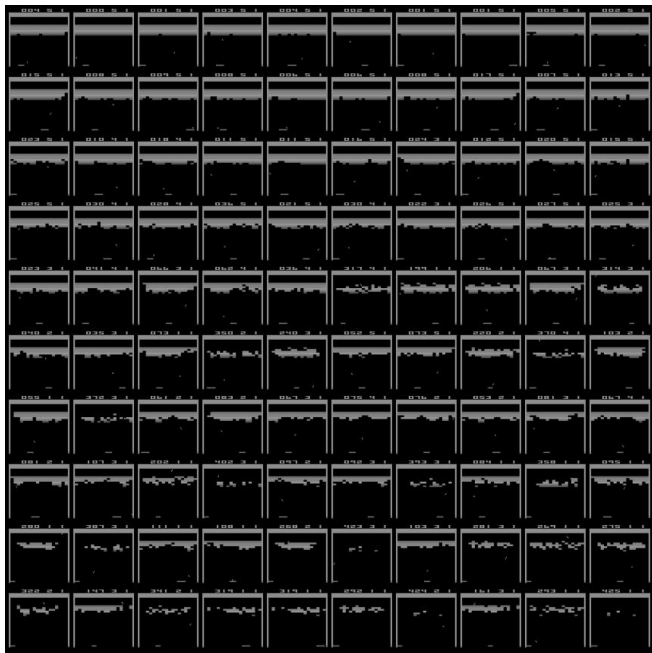


Figure 3: Examples of observed screens in Atari 2600 Breakout.

We showed three trajectories of samples from the prior distribution with the scaled norm for both models in the main text. We also visualize more samples in Figure 4 and 5. For both models, we generated samples with $\|\tilde{v}\|_2 = 0, 1, 2, 3, 5,$ and 10.

Vanilla VAE tended to generate oversaturated images when the norm $\|\tilde{v}\|$ was small. Although the model generated several images which include a small number of blocks as the norm increases, it also generated images with a constant amount of blocks even $\|\tilde{v}\| = 10$. On the other hand, the number of blocks contained in the generated image of Hyperbolic VAE gradually decreased according to the norm.

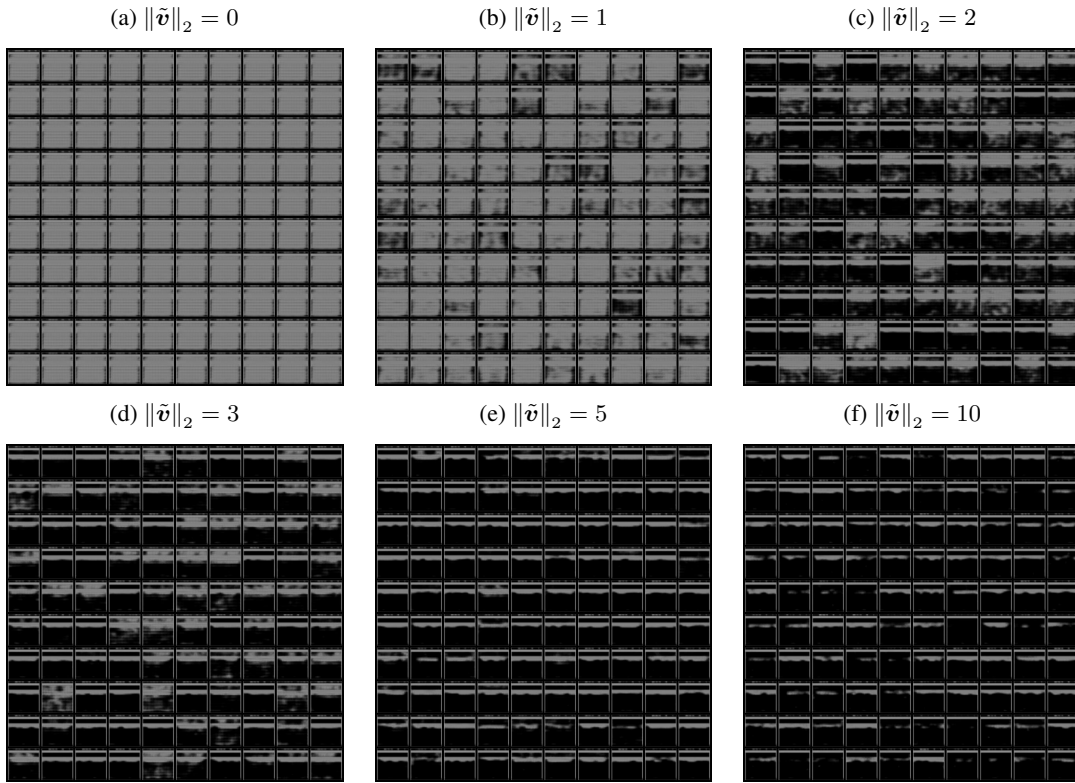


Figure 4: Images generated by Vanilla VAE with constant norm $\|\tilde{\mathbf{v}}\|_2 = a$.

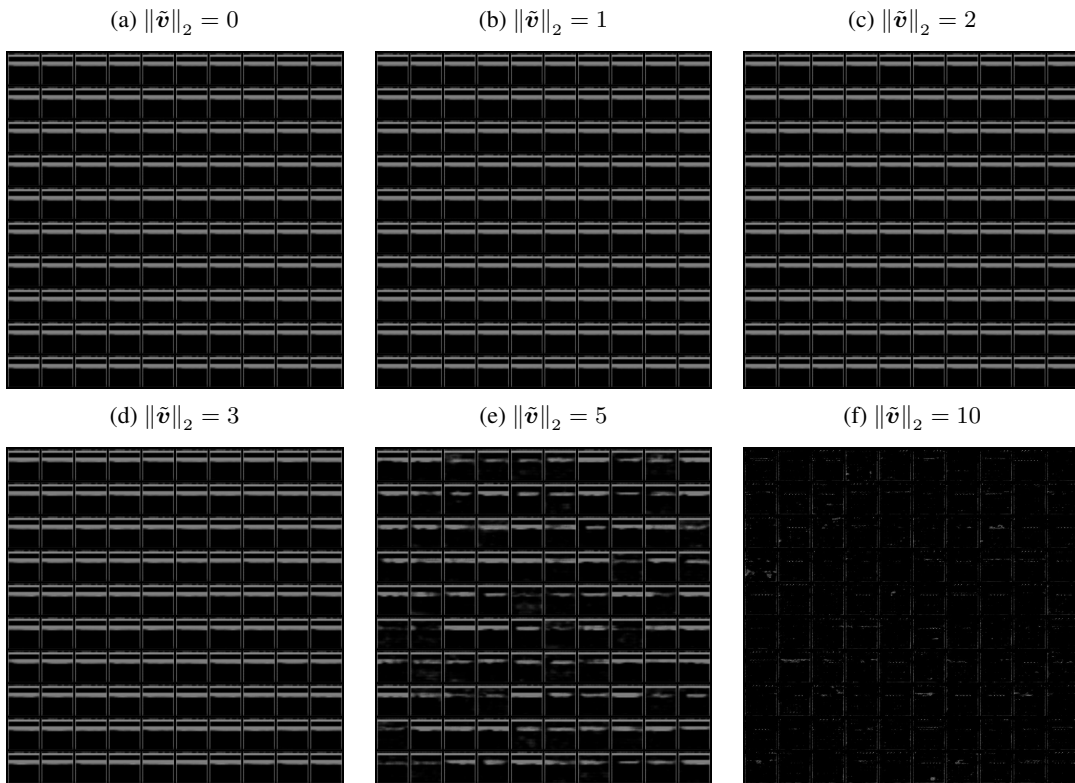


Figure 5: Images generated by Hyperbolic VAE with constant norm $\|\tilde{\mathbf{v}}\|_2 = a$.

C.4. Word Embeddings

We showed the experimental results of probabilistic word embedding models with diagonal variance in the main text. Table 2 shows the same comparison with the reference model by Nickel & Kiela (2017). We also show the results with unit variance (Table 3). When the dimensions of the latent variable are small, the performance of the model on hyperbolic space did not deteriorate much by changing the variance from diagonal to unit. However, the same change dramatically worsened the performance of the model on Euclidean space.

n	Euclid		Hyperbolic		Nickel & Kiela (2017)	
	MAP	Rank	MAP	Rank	MAP	Rank
5	0.296±.006	25.09±.80	0.506±.017	20.55±1.34	0.823	4.9
10	0.778±.007	4.70±.05	0.795±.007	5.07±.12	0.851	4.02
20	0.894±.002	2.23±.03	0.897±.005	2.54±.20	0.855	3.84
50	0.942±.003	1.51±.04	0.975±.001	1.19±.01	0.86	3.98
100	0.953±.002	1.34±.02	0.978±.002	1.15±.01	0.857	3.9

Table 2: Experimental results of the reconstruction performance on the transitive closure of the WordNet noun hierarchy for several latent space dimension n . We calculated the mean and the ± 1 SD with three different experiments.

n	Euclid		Hyperbolic	
	MAP	Rank	MAP	Rank
5	0.217±.008	55.28±3.54	0.529±.010	22.38±.70
10	0.698±.030	6.54±.65	0.771±.006	5.89±.29
20	0.832±.016	3.08±.16	0.862±.002	2.80±.13
50	0.910±.006	1.78±.071	0.903±.003	1.94±.03
100	0.882±.008	4.75±2.01	0.884±.003	2.57±.09

Table 3: Experimental results of the word embedding models with unit variance on the WordNet noun dataset. We calculated the mean and the ± 1 SD with three different experiments.

D. Network Architecture

Table 4 shows the network architecture that we used in Breakout experiments. We evaluated Vanilla and Hyperbolic VAEs with a DCGAN-based architecture (Radford et al., 2015) with the kernel size of the convolution and deconvolution layers as 3. We used leaky ReLU nonlinearities for the encoder and ReLU nonlinearities for the decoder. We set the latent space dimension as 20. We gradually increased β from 0.1 to 4.0 linearly during the first 30 epochs. To ensure the initial embedding vector close to the origin, we initialized γ for the batch normalization layer (Ioffe & Szegedy, 2015) of the encoder as 0.1. We modeled the probability distribution of the data space $p(\mathbf{x}|\mathbf{z})$ as Gaussian, so the decoder output a vector twice as large as the original image.

Encoder		Decoder	
Layer	Size	Layer	Size
Input	$80 \times 80 \times 1$	Linear	$10 \times 10 \times 64$
Convolution	$80 \times 80 \times 16$	BatchNormalization	
BatchNormalization		Deconvolution	$20 \times 20 \times 32$
Convolution	$40 \times 40 \times 32$	BatchNormalization	
BatchNormalization		Convolution	$20 \times 20 \times 32$
Convolution	$40 \times 40 \times 32$	BatchNormalization	
BatchNormalization		Deconvolution	$40 \times 40 \times 16$
Convolution	$20 \times 20 \times 64$	BatchNormalization	
BatchNormalization		Convolution	$40 \times 40 \times 16$
Convolution	$20 \times 20 \times 64$	Deconvolution	$80 \times 80 \times 2$
BatchNormalization		Convolution	$80 \times 80 \times 2$
Convolution	$10 \times 10 \times 64$		
Linear	$2n$		

Table 4: Network architecture for Atari 2600 Breakout dataset.

References

- Ioffe, S. and Szegedy, C. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *Proceedings of the 32nd International Conference on International Conference on Machine Learning*, volume 37, pp. 448–456, 2015.
- Nickel, M. and Kiela, D. Poincaré embeddings for learning hierarchical representations. In *Advances in Neural Information Processing Systems 30*, pp. 6338–6347. 2017.
- Radford, A., Metz, L., and Chintala, S. Unsupervised representation learning with deep convolutional generative adversarial networks. *CoRR*, abs/1511.06434, 2015.