

Continuous Hierarchical Representations with Poincaré Variational Auto-Encoders

Emile Mathieu[†], Charline Le Lan[†], Chris J. Maddison^{†*}, Ryota Tomioka[‡] and Yee Whye Teh^{†*}

[†] Department of Statistics, University of Oxford, United Kingdom, ^{*} DeepMind, London, United Kingdom, [‡] Microsoft Research, Cambridge, United Kingdom



UNIVERSITY OF
OXFORD

Overview & Contributions

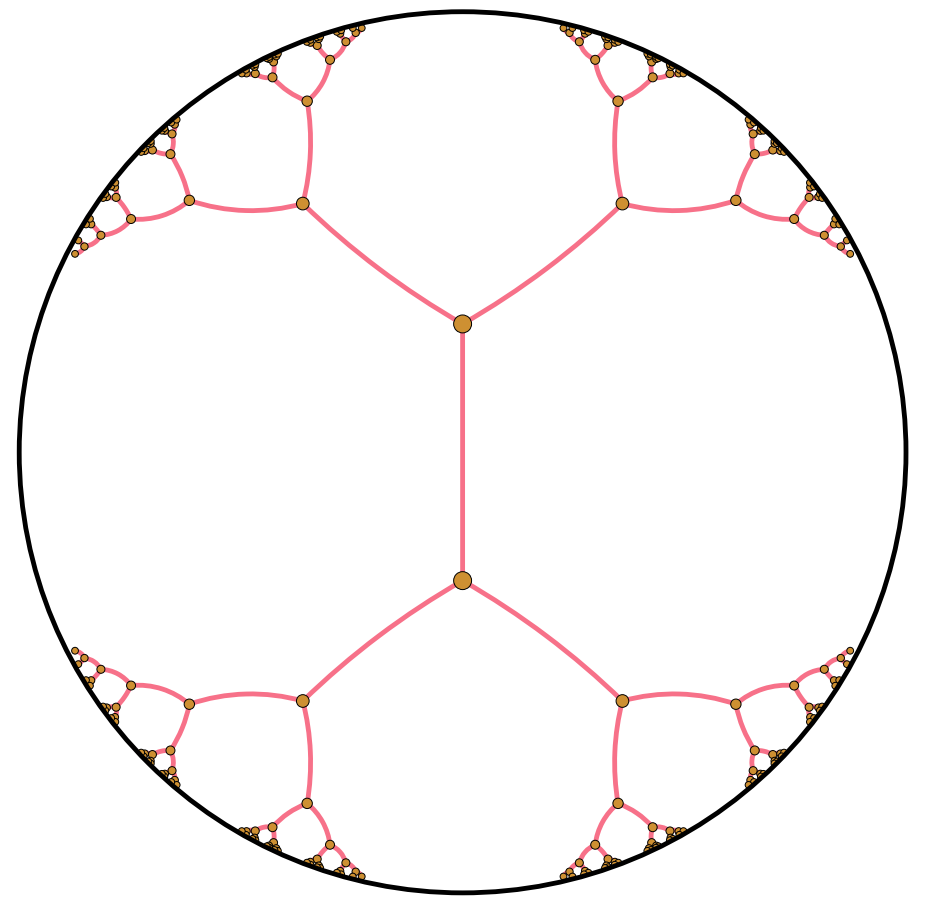
Many real datasets are **hierarchically structured**. However, traditional variational auto-encoders (VAEs) [1, 2] map data in a **Euclidean latent space** which cannot efficiently embed tree-like structures. **Hyperbolic spaces** with negative curvature can [3] and their smoothness is well-suited for gradient based approaches [4].

1. We empirically demonstrate that endowing a VAE with a **Poincaré ball latent space** can be beneficial in terms of model generalisation and can yield more interpretable representations.
2. We propose **efficient and reparametrisable sampling schemes**, and calculate the **probability density functions**, for two canonical Gaussian generalisations defined on the Poincaré ball, namely the **maximum-entropy** and **wrapped normal** distributions.
3. We introduce a **decoder architecture** taking into account the hyperbolic geometry, which we empirically show to be crucial.

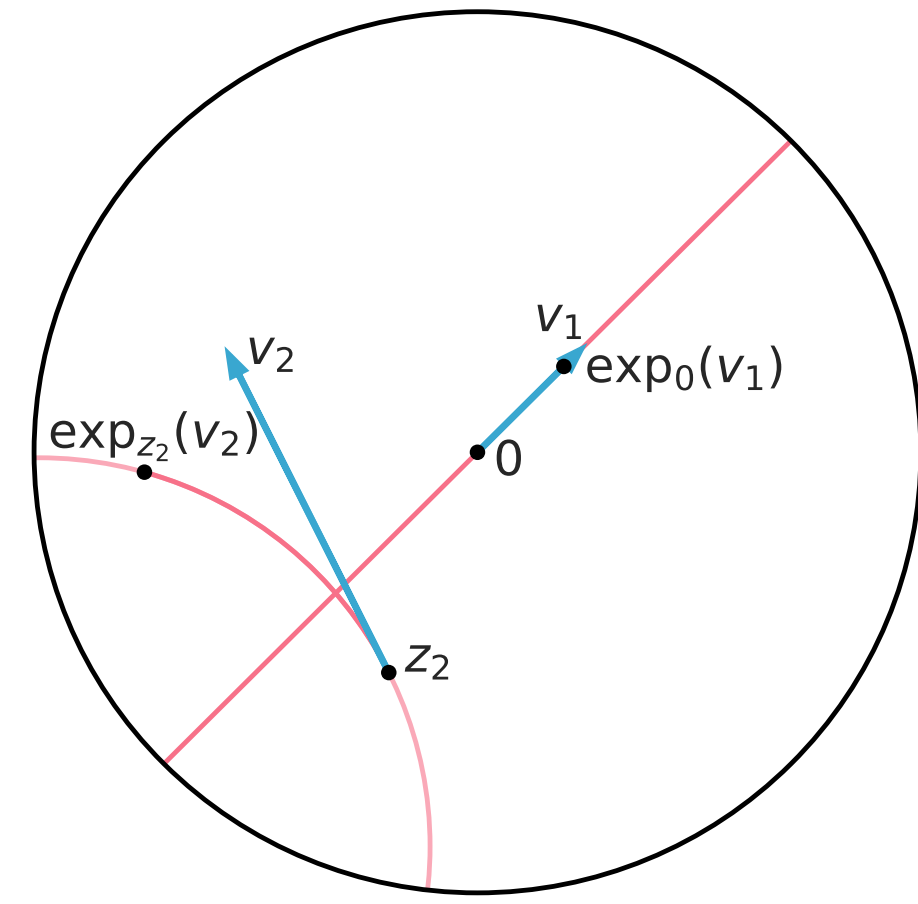
The Poincaré ball model of hyperbolic geometry

The d-dimensional Poincaré ball with curvature $-c$ is the Riemannian manifold $\mathbb{B}_c^d = (\mathcal{B}_c^d, \mathfrak{g}_p^c)$ [5], where $\mathcal{B}_c^d = \{z \in \mathbb{R}^d \mid \|z\|_2 \leq 1/\sqrt{c}\}$, and \mathfrak{g}_p^c its *metric tensor*,

$$\mathfrak{g}_p^c(z) = (\lambda_z^c)^2 \mathfrak{g}_e(z) = \left(\frac{2}{1 - c\|z\|^2} \right)^2 \mathfrak{g}_e(z) \quad (1)$$



(a) Isometric embedding of tree.



(b) Geodesics and exponential maps.

We are extremely grateful to Adam Foster, Phillippe Gagnon and Emmanuel Chevallier for their help. EM, YWT's research leading to these results received funding from the European Research Council under the European Union's Seventh Framework Programme (FP7/2007- 2013) ERC grant agreement no. 617071 and they acknowledge Microsoft Research and EPSRC for funding EM's studentship, and EPSRC grant agreement no. EP/N509711/1 for funding CL's studentship.

- [1] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2014.
- [2] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic Backpropagation and Approximate Inference in Deep Generative Models. In *ICML*, 2014.
- [3] Rik Sarkar. Low distortion delaunay embedding of trees in hyperbolic plane. In Marc van Kreveld and Bettina Speckmann, editors, *Graph Drawing*, pages 355–366. Springer Berlin Heidelberg, 2012.
- [4] Maximilian Nickel and Douwe Kiela. Poincaré embeddings for learning hierarchical representations. In *Advances in Neural Information Processing Systems*, pages 6341–6350, 2017.
- [5] E. Beltrami. *Teoria fondamentale degli spazii di curvatura costante: memoria*. F. Zanetti, 1868.
- [6] Xavier Pennec. Intrinsic statistics on riemannian manifolds: Basic tools for geometric measurements. *Journal of Mathematical Imaging and Vision*, 25(1):127, Jul 2006.
- [7] Salem Said, Lionel Bombrun, and Yannick Berthoumieu. New riemannian priors on the univariate normal model. *Entropy*, 16(7):4015–4031, 2014.
- [8] Octavian-Eugen Ganeva, Gary Bécigneul, and Thomas Hofmann. Hyperbolic neural networks. In *International Conference on Neural Information Processing Systems*, pages 5350–5360, 2018.
- [9] Michael Figurnov, Shakir Mohamed, and Andriy Mnih. Implicit reparameterization gradients. In *International Conference on Neural Information Processing Systems*, pages 439–450, 2018.
- [10] Thomas N Kipf and Max Welling. Variational graph auto-encoders. *Workshop on Bayesian Deep Learning, NIPS*, 2016.

Hyperbolic normal distributions reparametrisation

Invariant measure: $d\mathcal{M}(z) = \sqrt{|G(z)|}dz = (\lambda_z^c)^d dz$ with dz the Lebesgue measure

Riemannian normal: $\mathcal{N}_{\mathbb{B}_c^d}^R(z|\mu, \sigma^2) \propto \exp(-d_p^c(\mu, z)^2/2\sigma^2)$ (maximum-entropy [6, 7])

Wrapped normal: $z = \exp_{\mu}^c(v/\lambda_{\mu}^c)$ with $v \sim \mathcal{N}(\cdot|0, \Sigma)$ (push-forward)

$$\mathcal{N}_{\mathbb{B}_c^d}^W(z|\mu, \Sigma) = \mathcal{N}(\lambda_{\mu}^c \log_{\mu}^c(z)|0, \Sigma) (\sqrt{c} d_p^c(\mu, z)/\sinh(\sqrt{c} d_p^c(\mu, z)))^{d-1}$$

Reparametrisation through the exponential map: $z \sim \mathcal{N}_{\mathbb{B}_c^d}(z|\mu, \sigma^2) d\mathcal{M}(z)$

$$z = \exp_{\mu}^c(G(\mu)^{-\frac{1}{2}} v) = \exp_{\mu}^c(v/\lambda_{\mu}^c) \quad (2)$$

Isotropic: $v = r \alpha$, with direction $\alpha \sim \mathcal{U}(\mathbb{S}^{d-1})$, and hyperbolic radius $r = d_p^c(\mu, x)$ density

$$\rho^R(r) \propto 1_{\mathbb{R}_+}(r) e^{-\frac{r^2}{2\sigma^2}} \left(\frac{\sinh(\sqrt{c}r)}{\sqrt{c}} \right)^{d-1} \xrightarrow{c \rightarrow 0} \rho^W(r) \propto 1_{\mathbb{R}_+}(r) e^{-\frac{r^2}{2\sigma^2}} r^{d-1} \quad (3)$$

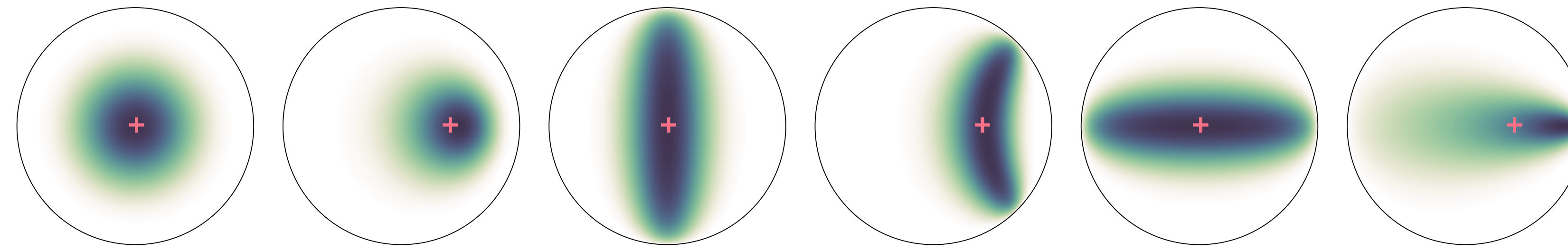


Figure 2: Wrapped normal probability measures for Fréchet means μ , concentrations $\Sigma = \text{diag}(\sigma)$ and $c = 1$.

Decoder & encoder architectures

Decoder Compute geodesic distance to hyperplanes (i.e. *gyroplanes*)

$$f_{a,p}(z) = \langle a, z - p \rangle = \text{sign}(\langle a, z - p \rangle) \|a\| d_E(z, H_{a,p}), \text{ with, } H_{a,p} = p + \{a\}^{\perp}.$$

$$f_{a,p}^c(z) = \text{sign}(\langle a, \log_p^c(z) \rangle_p) \|a\|_p d_p^c(z, H_{a,p}^c), \text{ with, } H_{a,p}^c = \exp_p^c(\{a\}^{\perp}) \text{ [8].}$$

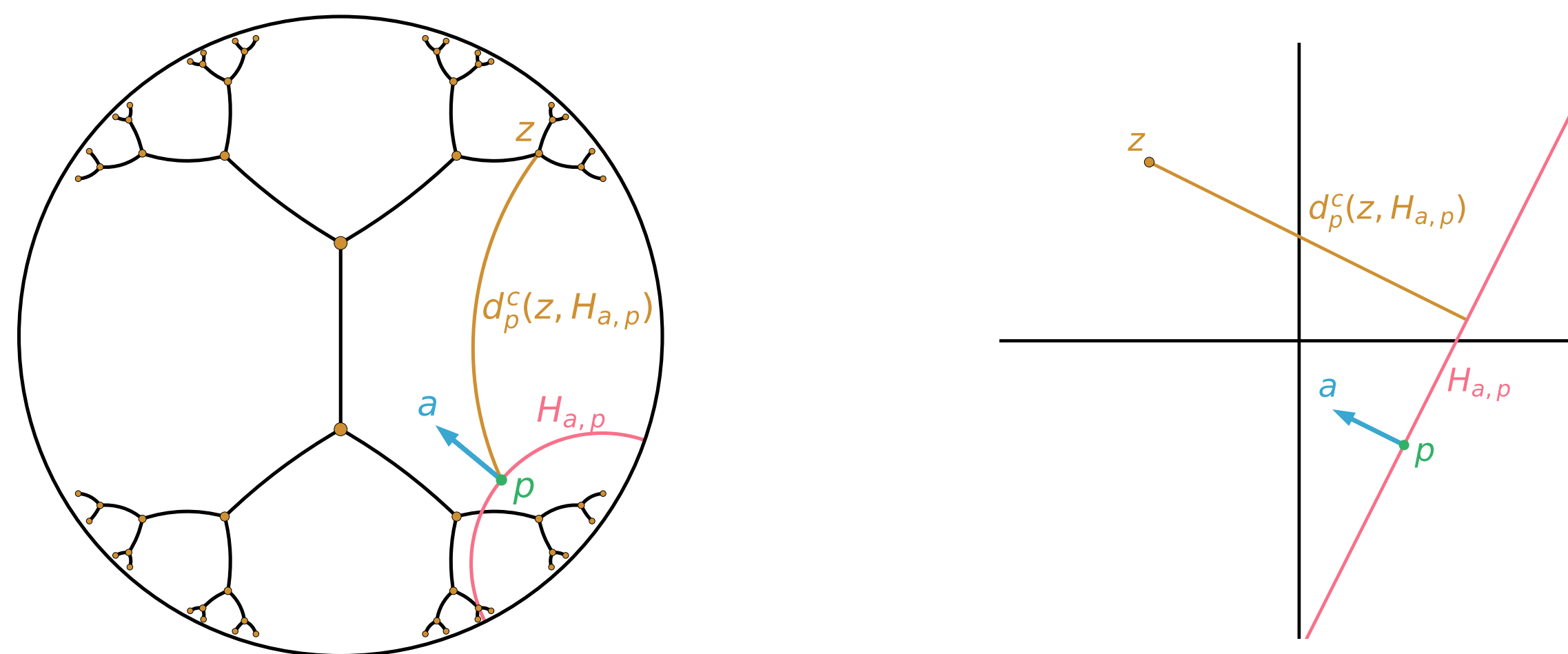


Figure 3: Orthogonal projection on a hyperplane in \mathbb{B}_c^2 (a) and \mathbb{R}^2 (b).

Encoder $\mu = \exp_{\mu}(\text{enc}_{\phi}^{\mu}(x)) \in \mathbb{B}_c^d$ and $\sigma = \text{softplus}(\text{enc}_{\phi}^{\sigma}(x)) \in \mathbb{R}_*^+$.

Training

Model $p_{\theta}(x|z) = p(x|\text{dec}_{\theta}(z))$, $p(z) = \mathcal{N}_{\mathbb{B}_c^d}(z|0, \sigma_0^2)$ and $q_{\phi}(z|x) = \mathcal{N}_{\mathbb{B}_c^d}(z|\text{enc}_{\phi}^{\mu}(x), \text{enc}_{\phi}^{\sigma}(x)^2)$.

$$\text{ELBO } \log p(x) \geq \mathcal{L}_{\mathcal{M}}(x; \theta, \phi) \triangleq \int_{\mathcal{M}} \ln \left(\frac{p_{\theta}(x|z)p(z)}{q_{\phi}(z|x)} \right) q_{\phi}(z|x) d\mathcal{M}(z) \text{ via Monte Carlo} \quad (4)$$

Gradients $\nabla_{\mu} z$ via reparametrisation. $\nabla_{\sigma} z$ via reparametrisation for the *wrapped* normal and via implicit reparametrisation [9] of ρ^R via its cdf $F^R(r; \sigma)$ for the *Riemannian* normal.

Branching diffusion process

Nodes $(x_1, \dots, x_N) \in \mathbb{R}^n$ are hierarchically sampled as follow

$$x_i \sim \mathcal{N}(\cdot | x_{\pi(i)}, \gamma^2) \quad \forall i \in 1, \dots, N, \pi(i) : \text{parent of } i\text{th node}$$

Table 1: Negative test marginal likelihood estimates $\mathcal{L}_{\text{IWAE}}$ (with 5000 samples).

		Models					
	σ_0	\mathcal{N} -VAE	$\mathcal{P}^{0.1}$ -VAE	$\mathcal{P}^{0.3}$ -VAE	$\mathcal{P}^{0.8}$ -VAE	$\mathcal{P}^{1.0}$ -VAE	$\mathcal{P}^{1.2}$ -VAE
$\mathcal{L}_{\text{IWAE}}$	1	57.1 \pm 0.2	57.1 \pm 0.2	57.2 \pm 0.2	56.9 \pm 0.2	56.7 \pm 0.2	56.6 \pm 0.2
$\mathcal{L}_{\text{IWAE}}$	1.7	57.0 \pm 0.2	56.8 \pm 0.2	56.6 \pm 0.2	55.9 \pm 0.2	55.7 \pm 0.2	55.6\pm0.2

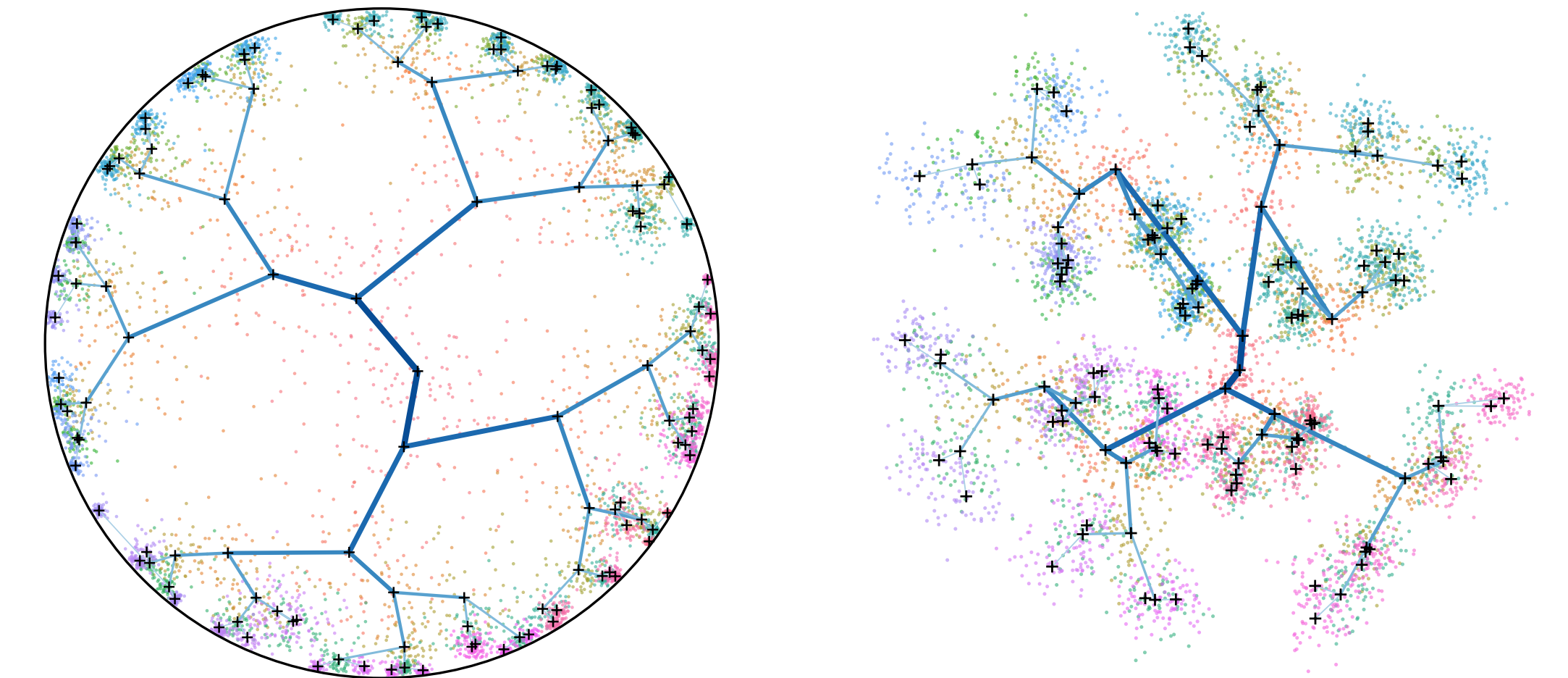


Figure 4: Latent representations learned by \mathcal{P}^1 -VAE (a), \mathcal{N} -VAE (b). Embeddings are represented by black crosses, and colour dots are posterior samples. Blue lines represent true hierarchy.

MNIST digits

One can view the natural clustering in MNIST images as a hierarchy with each of the 10 classes being internal nodes of the hierarchy.

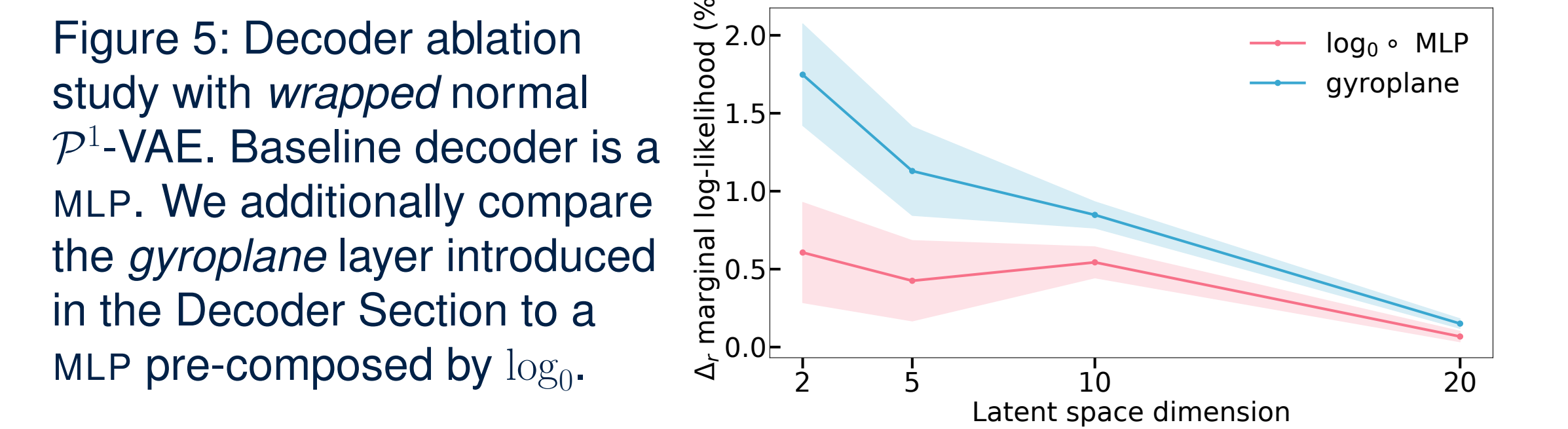


Table 2: Per digit accuracy of a classifier trained on the 2-d latent embeddings. Results are averaged over 10 sets of embeddings and 5 classifier trainings.

Digits	0	1	2	3	4	5	6	7	8	9	Avg
\mathcal{N} -VAE	89	97	81	75	59	43	89	78	68	57	73.6
$\mathcal{P}^{1.4}$ -VAE	94	97	82	79	69	47	90	77	68	53	75.6

Graph embeddings

We evaluate the performance of a variational graph auto-encoder (VGAE) [10] for link prediction in networks.

Table 3: Results on network link prediction. 95% confidence intervals are computed over 40 runs.

	Phylogenetic		CS PhDs		Diseases	
	AUC	AP	AUC	AP	AUC	AP
\mathcal{N} -VAE	54.2 \pm 2.2	54.0 \pm 2.1	56.5 \pm 1.1	56.4 \pm 1.1	89.8 \pm 0.7	91.8 \pm 0.7
\mathcal{P} -VAE	59.0\pm1.9	55.5\pm1.6	59.8\pm1.2	56.7\pm1.2	92.3\pm0.7	93.6\pm0.5