

# Geometry and Deep generative modelling

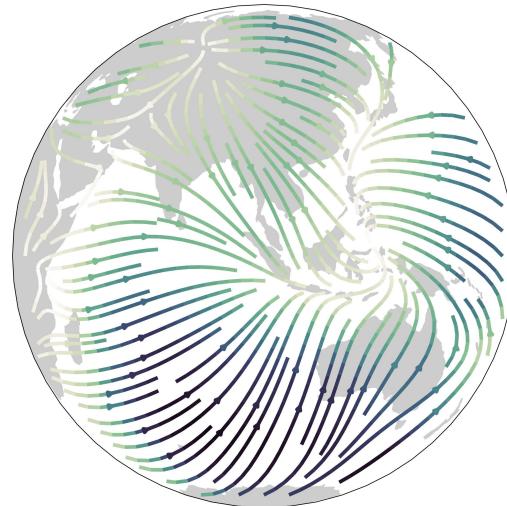
*Geometry as an inductive bias and a manifold constraint*

*Emile Mathieu*

*Data-Centric Engineering Reading Group,  
@ The Alan Turing Institute.*



UNIVERSITY OF  
**OXFORD**

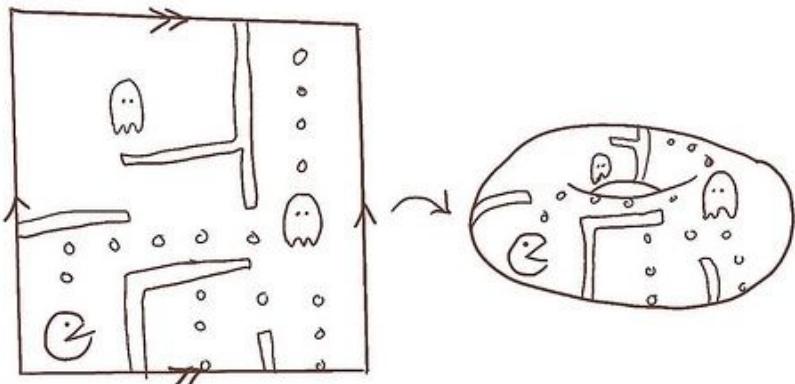


# Overview: Geometry and Deep generative models

- *Geometry as an inductive bias*
  - Variational Auto-Encoder for hierarchical data
- *Geometry as a manifold constraint*
  - Normalizing Flow for manifold data

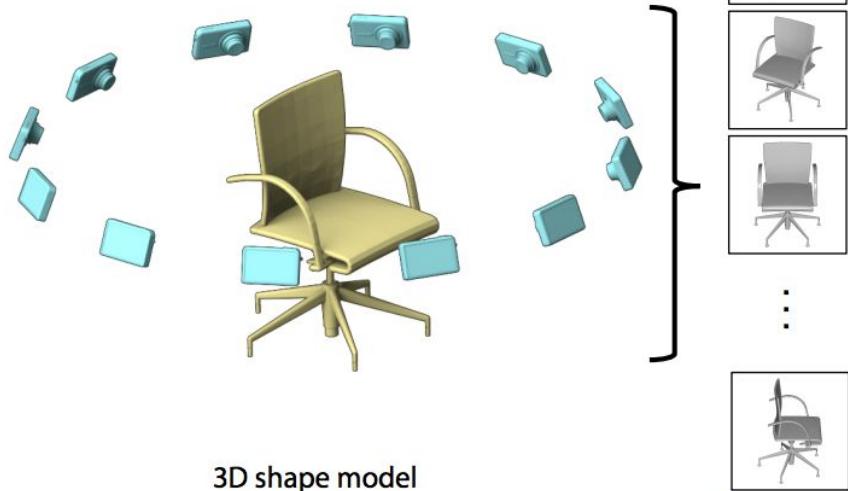
# Geometry as an inductive bias

# Data with loops



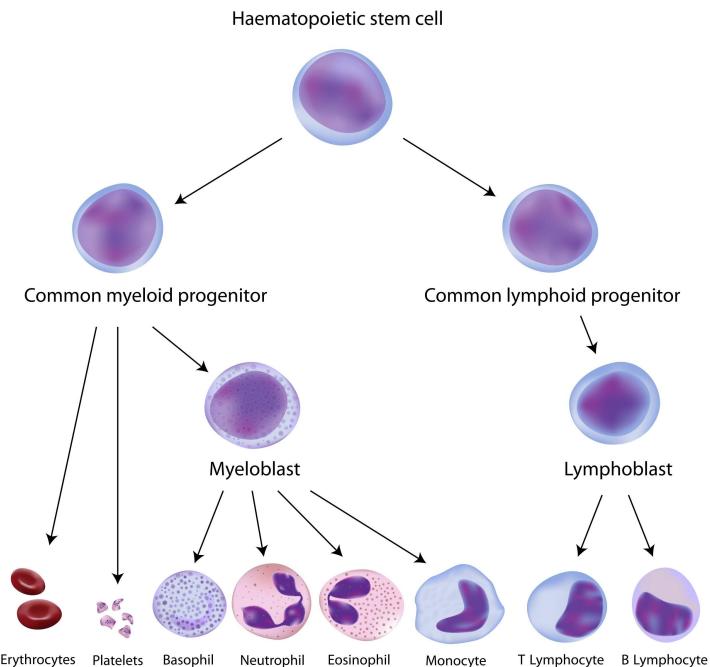
The torus explained.

From Human Mathematics Blog, 2009

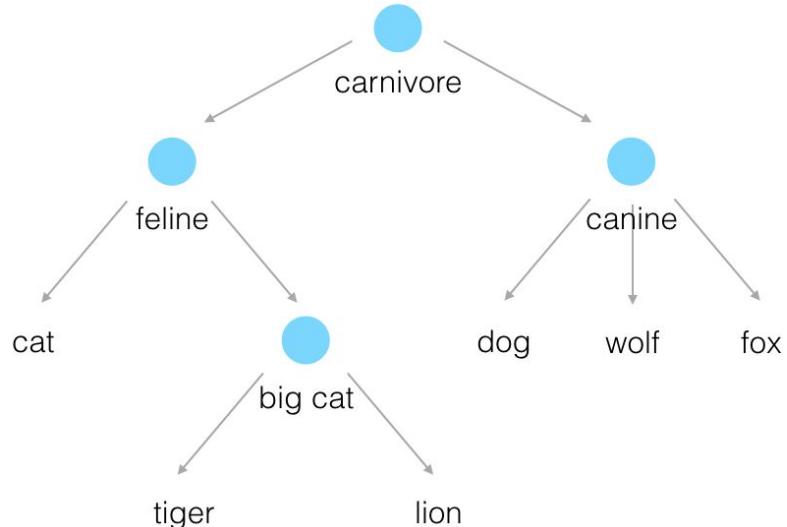


From Connelly Barnes, 3D Reconstruction and Understanding, 2017

# Data with hierarchy



From Tom Ulrich et al. 2014



Fragment of WordNet taxonomy graph

# Poincaré Variational Auto-Encoders

A generative model for data with underlying hierarchical structure

- Latent space endowed with hyperbolic geometry
- Derive necessary methods for two main Gaussian generalisations on hyperbolic space
- Decoder architecture that is geometry-aware

---

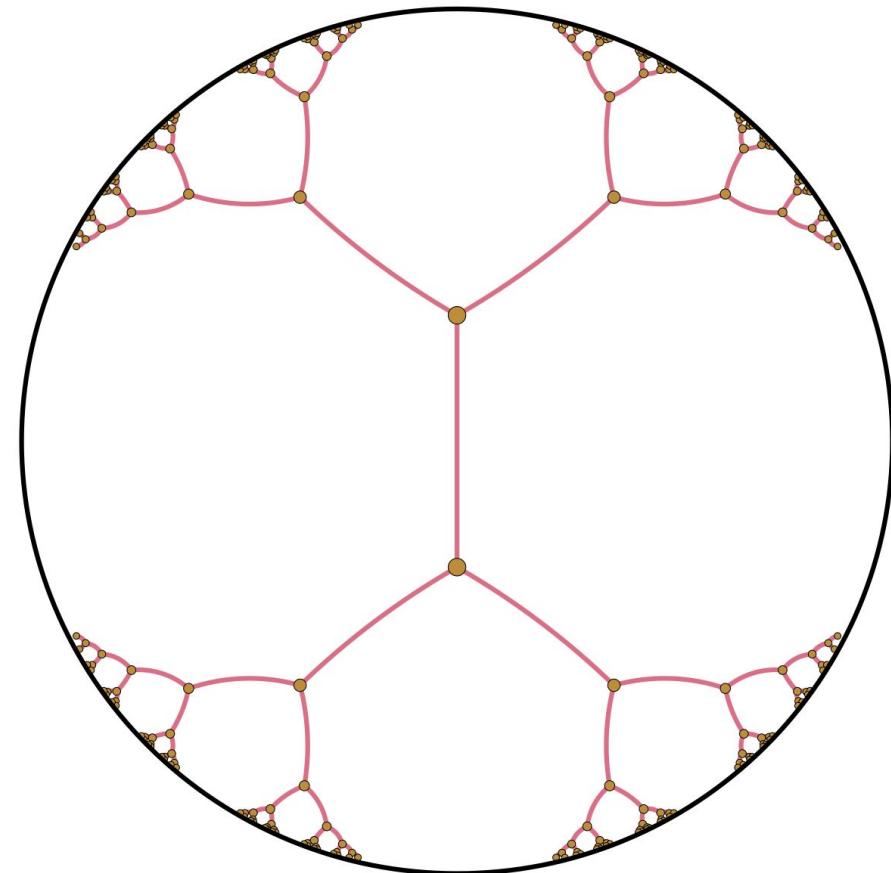
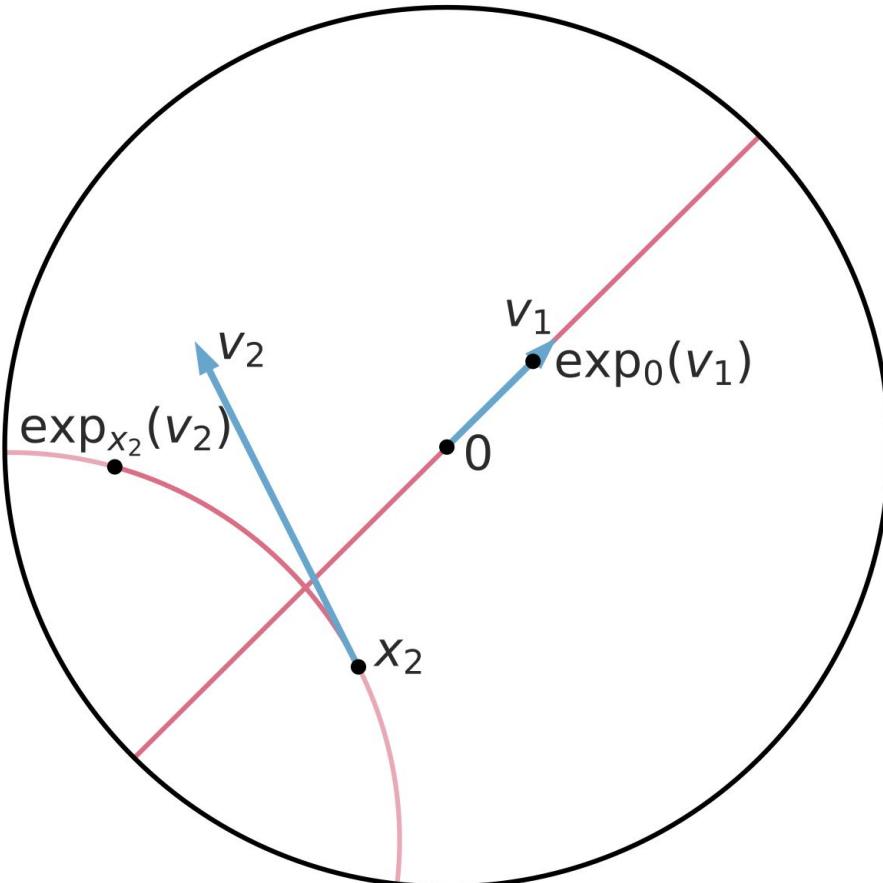
With C. Le Lan, C. Maddison,  
R. Tomioka, Y.W Teh

# Motivation

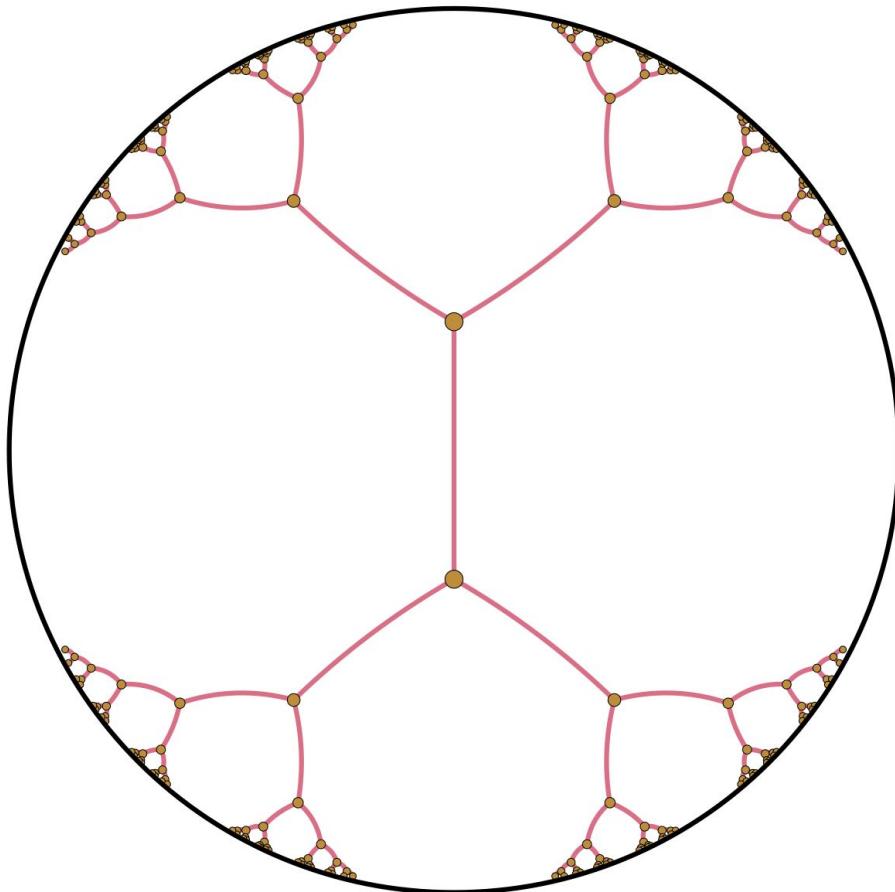
**Theorem [Bourgain]** A Euclidean space is unable to obtain low distortion for tree embeddings—even using an unbounded number of dimensions.

**Theorem [R. Sarkar, 2011]** Trees can be embedded with arbitrarily low distortion into the Poincaré disk.

# Poincaré Ball: A model of Hyperbolic geometry

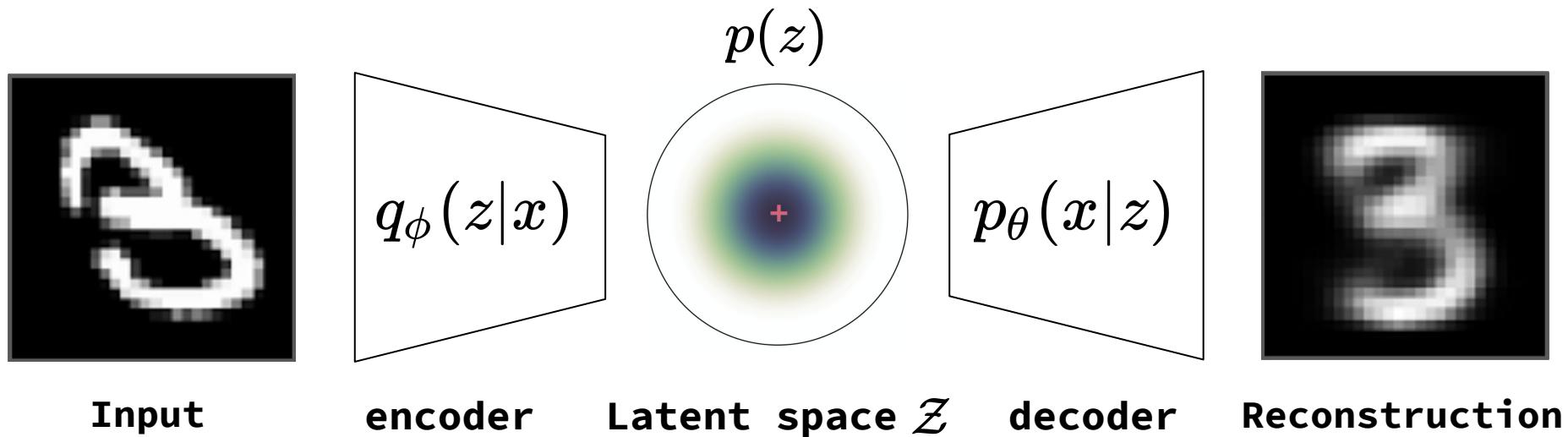


# Poincaré Ball: A model of Hyperbolic geometry

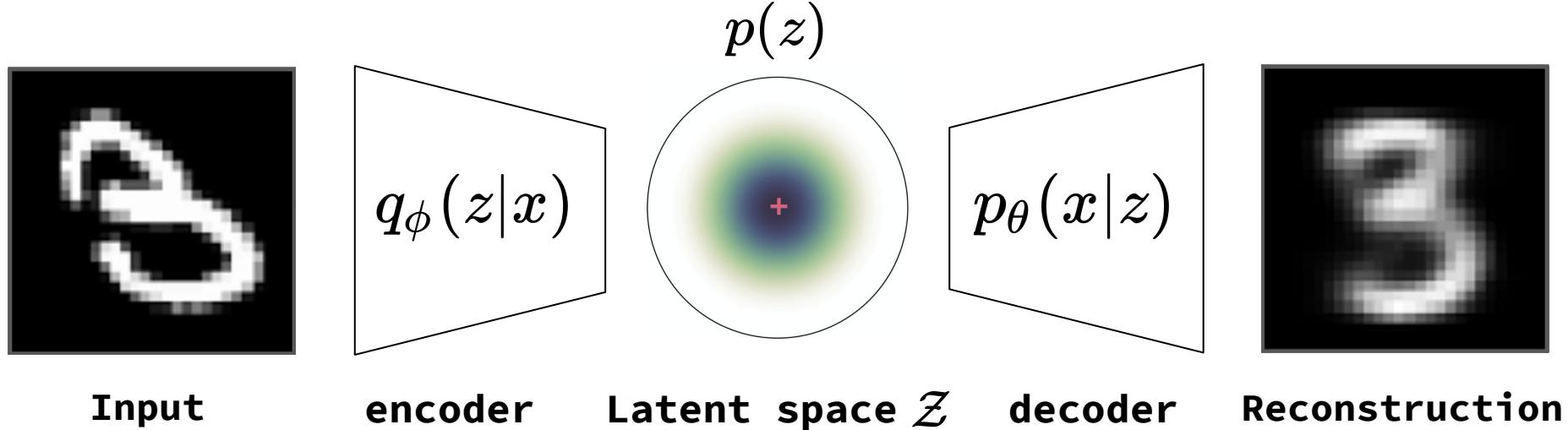


# Model

# Poincaré Variational Auto-Encoder

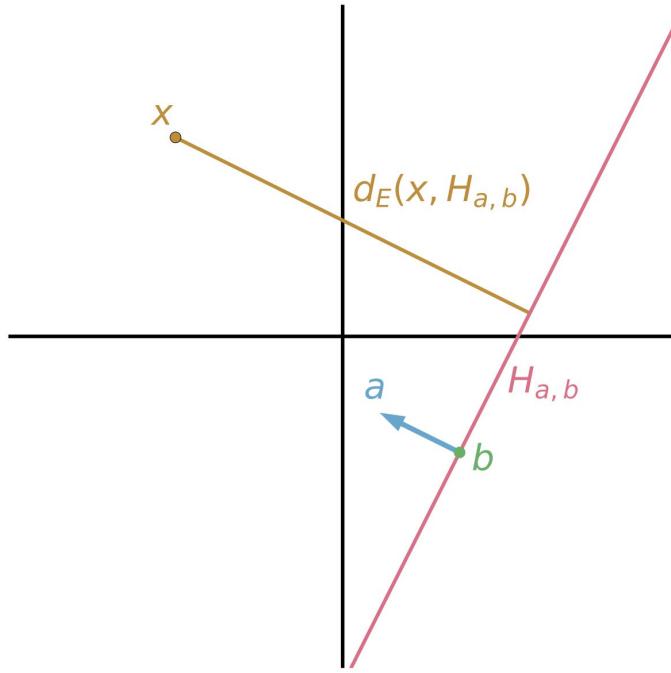


# Poincaré Variational Auto-Encoder



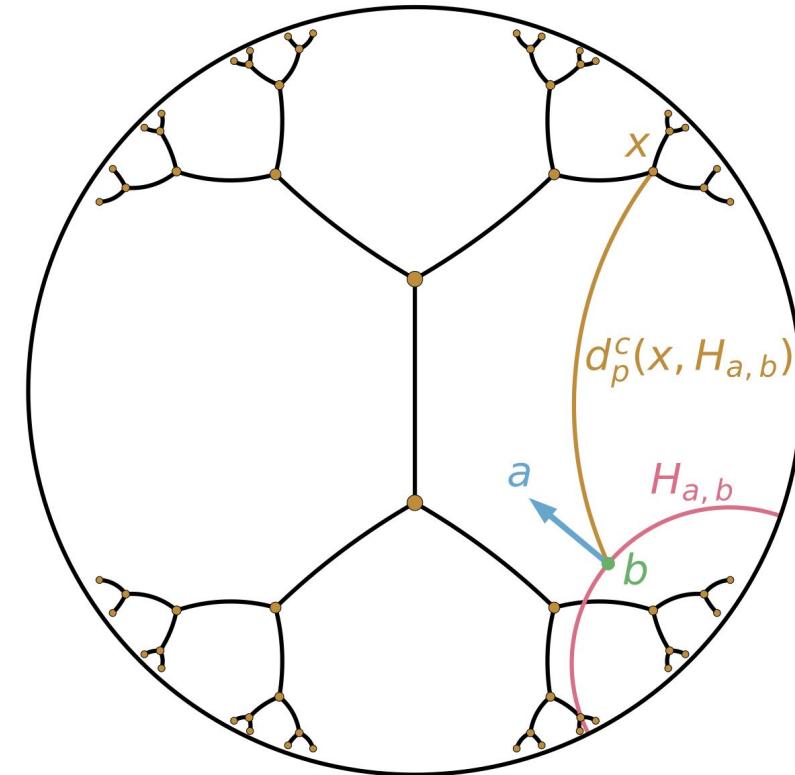
$$\log p(\mathbf{x}) \geq \int_{\mathcal{M}} \ln \left( \frac{p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right) q_\phi(\mathbf{z}|\mathbf{x}) d\mathcal{M}(\mathbf{z}).$$

# Decoder: From latents to observations



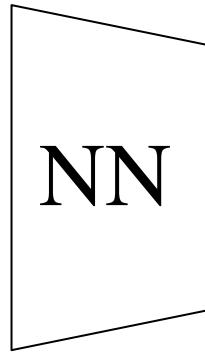
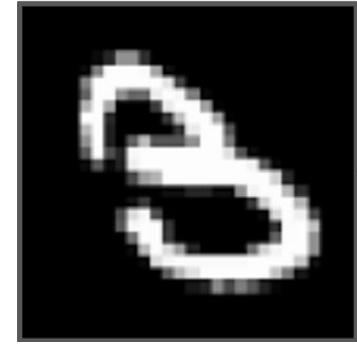
$$f_{\mathbf{a}, \mathbf{b}}(\mathbf{x}) = \langle \mathbf{a}, \mathbf{x} - \mathbf{b} \rangle$$

$$f_{\mathbf{a}, \mathbf{b}}(\mathbf{x}) = \text{sign}(\langle \mathbf{a}, \mathbf{x} - \mathbf{b} \rangle) \|\mathbf{a}\| d_E(\mathbf{x}, H_{\mathbf{a}, \mathbf{b}}^c)$$



$$f_{\mathbf{a}, \mathbf{b}}^c(\mathbf{x}) = \text{sign}(\langle \mathbf{a}, \log_x^c(\mathbf{b}) \rangle_{\mathbf{b}}) \|\mathbf{a}\|_{\mathbf{b}} d_p^c(\mathbf{x}, H_{\mathbf{a}, \mathbf{b}}^c)$$

# Encoder: From observations to latents



MLP<sub>σ</sub>

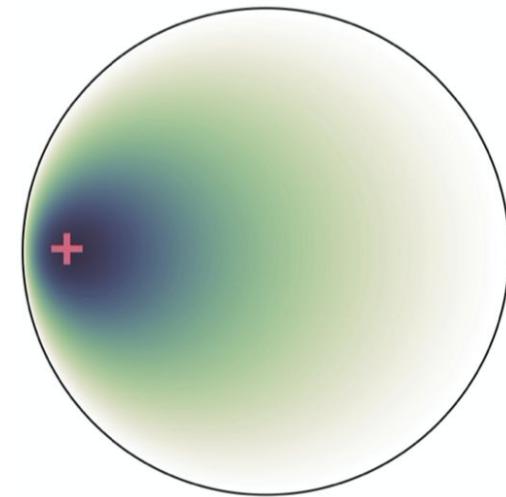
softplus

σ

MLP<sub>μ</sub>

exp<sub>0</sub><sup>c</sup>

μ



$$\text{exp}_0^c : \mathbb{R}^d \rightarrow \mathbb{B}_c^d$$

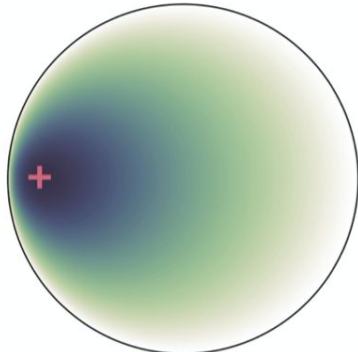
$$u \mapsto \tanh\left(\frac{\sqrt{c}}{2} \|u\|_2\right) \frac{u}{\sqrt{c}\|u\|_2}$$

$$q_\phi(z|x) = q(z; \mu, \sigma)$$

# Generalizing the Normal distribution

## Riemannian Normal

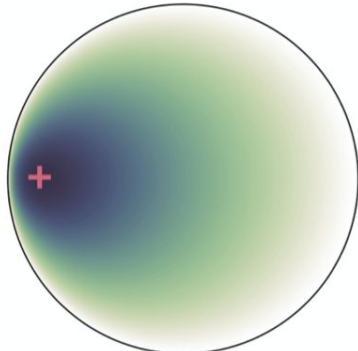
$$\frac{d\nu^R(x|\mu, \sigma^2)}{d\mathcal{M}(x)} \propto \exp\left(-\frac{d_p^c(\mu, x)^2}{2\sigma^2}\right)$$



# Generalizing the Normal distribution

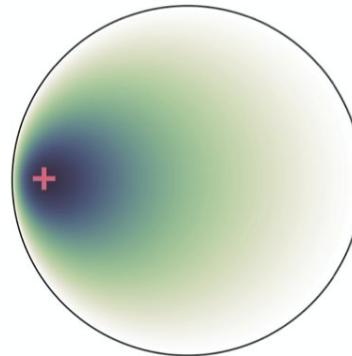
## Riemannian Normal

$$\frac{d\nu^R(x|\mu, \sigma^2)}{d\mathcal{M}(x)} \propto \exp\left(-\frac{d_p^c(\mu, x)^2}{2\sigma^2}\right)$$



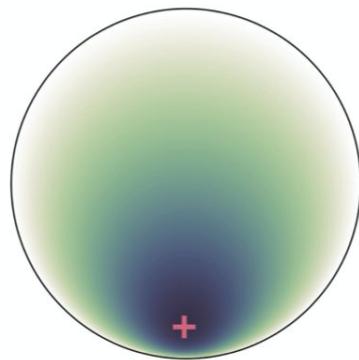
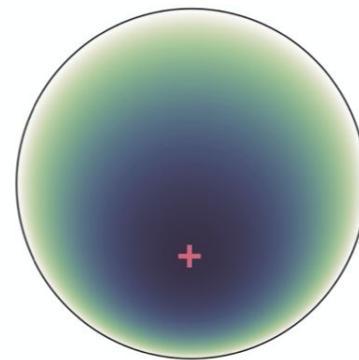
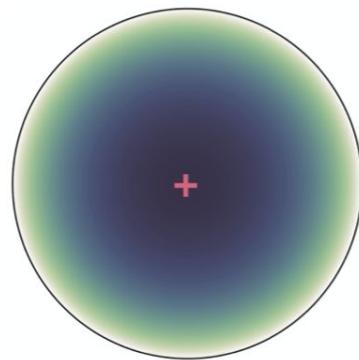
## Wrapped Normal

$$x = \exp_{\mu}^c \left( \frac{v}{\lambda_{\mu}^c} \right) \quad v \sim \mathcal{N}(\cdot | \mathbf{0}, \sigma^2)$$

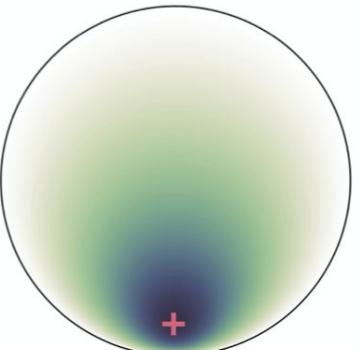
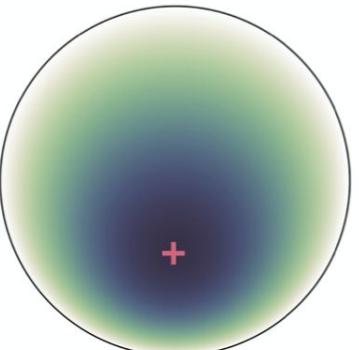
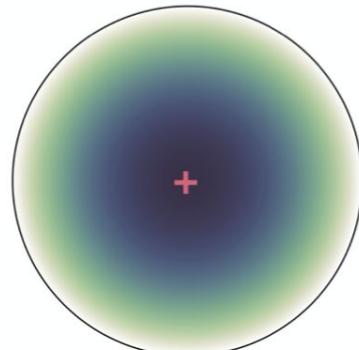


# Generalizing the Normal distribution

Riemannian



Wrapped



$$\sqrt{c}\|\boldsymbol{\mu}\|_2 = 0$$

$$\sqrt{c}\|\boldsymbol{\mu}\|_2 = 0.4$$

$$\sqrt{c}\|\boldsymbol{\mu}\|_2 = 0.8$$

# Reparametrization “trick”

Riemannian Normal

$$\frac{d\nu^R(x|\mu, \sigma^2)}{d\mathcal{M}(x)} \propto \exp\left(-\frac{d_p^c(\mu, x)^2}{2\sigma^2}\right)$$

Wrapped Normal

$$x = \exp_{\mu}^c \left( \frac{v}{\lambda_{\mu}^c} \right) \quad v \sim \mathcal{N}(\cdot | \mathbf{0}, \sigma^2)$$

Reparametrization 

$$x = \exp_{\mu}^c \left( \frac{r}{\lambda_{\mu}^c} \alpha \right) \quad \text{with} \quad r = d_p^c(\mu, x)$$

# Reparametrization “trick”

Riemannian Normal

$$\frac{d\nu^R(x|\mu, \sigma^2)}{d\mathcal{M}(x)} \propto \exp\left(-\frac{d_p^c(\mu, x)^2}{2\sigma^2}\right)$$

Wrapped Normal

$$x = \exp_{\mu}^c \left( \frac{v}{\lambda_{\mu}^c} \right) \quad v \sim \mathcal{N}(\cdot | \mathbf{0}, \sigma^2)$$

Reparametrization 

$$x = \exp_{\mu}^c \left( \frac{r}{\lambda_{\mu}^c} \alpha \right) \quad \text{with} \quad r = d_p^c(\mu, x)$$

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

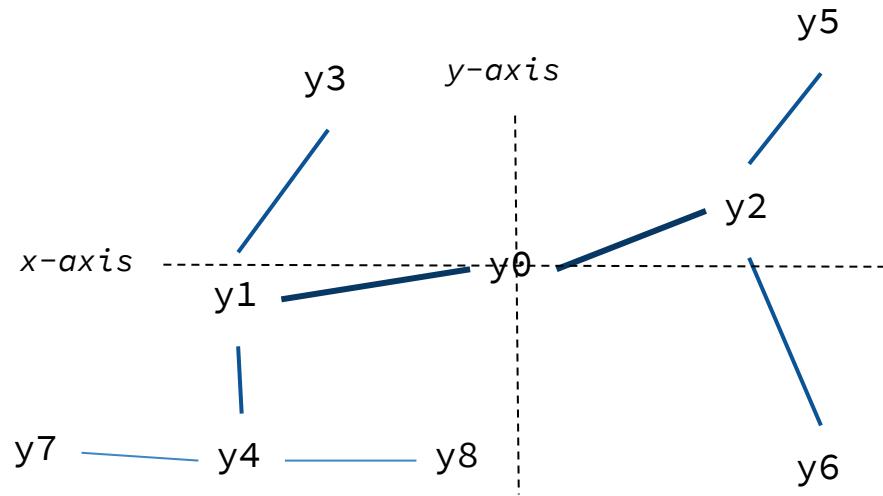
$$\rho^W(r) \propto \mathbb{1}_{\mathbb{R}_+}(r) e^{-\frac{r^2}{2\sigma^2}} r^{d-1}$$

# Experiments

# Synthetic dataset: Branching diffusion process

Nodes  $(\mathbf{y}_1, \dots, \mathbf{y}_N) \in \mathbb{R}^n$

$\mathbf{y}_i \sim \mathcal{N}(\cdot | \mathbf{y}_{\pi(i)}, \sigma_0^2)$   $\forall i \in 1, \dots, N$

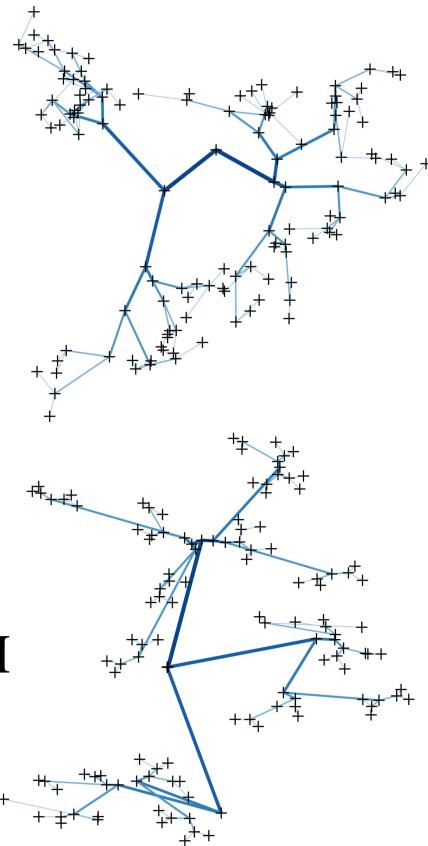


## Models

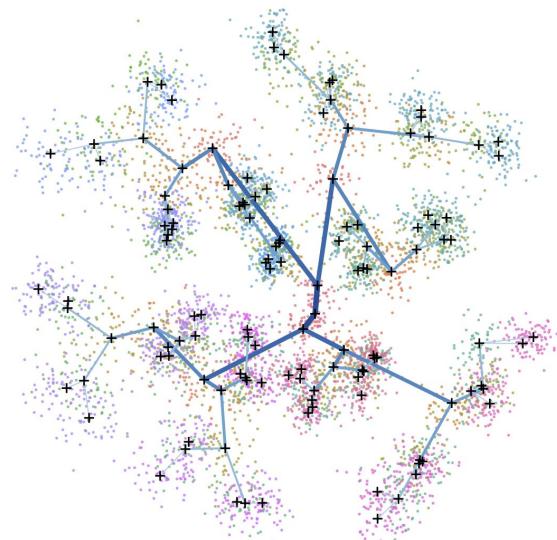
$\sigma_0$	$\mathcal{N}$ -VAE	$\mathcal{P}^{0.1}$ -VAE	$\mathcal{P}^{0.3}$ -VAE	$\mathcal{P}^{0.8}$ -VAE	$\mathcal{P}^{1.0}$ -VAE	$\mathcal{P}^{1.2}$ -VAE	
$\mathcal{L}_{\text{IWAE}}$	1	$57.1 \pm 0.2$	$57.1 \pm 0.2$	$57.2 \pm 0.2$	$56.9 \pm 0.2$	$56.7 \pm 0.2$	$56.6 \pm 0.2$
$\mathcal{L}_{\text{IWAE}}$	1.7	$57.0 \pm 0.2$	$56.8 \pm 0.2$	$56.6 \pm 0.2$	$55.9 \pm 0.2$	$55.7 \pm 0.2$	<b><math>55.6 \pm 0.2</math></b>

# Synthetic dataset: Branching diffusion process

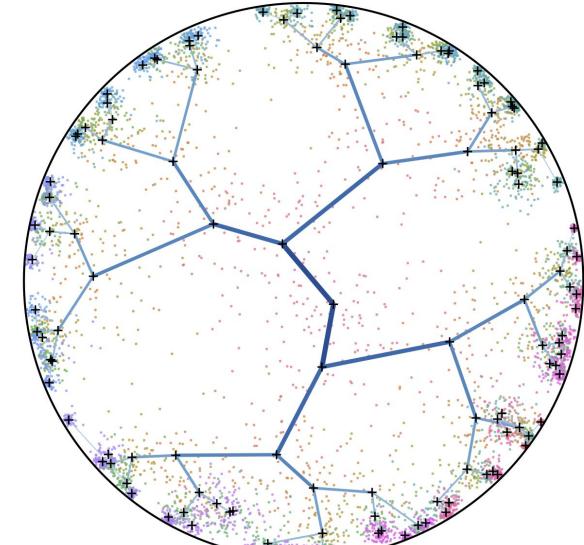
PCA



GPLVM

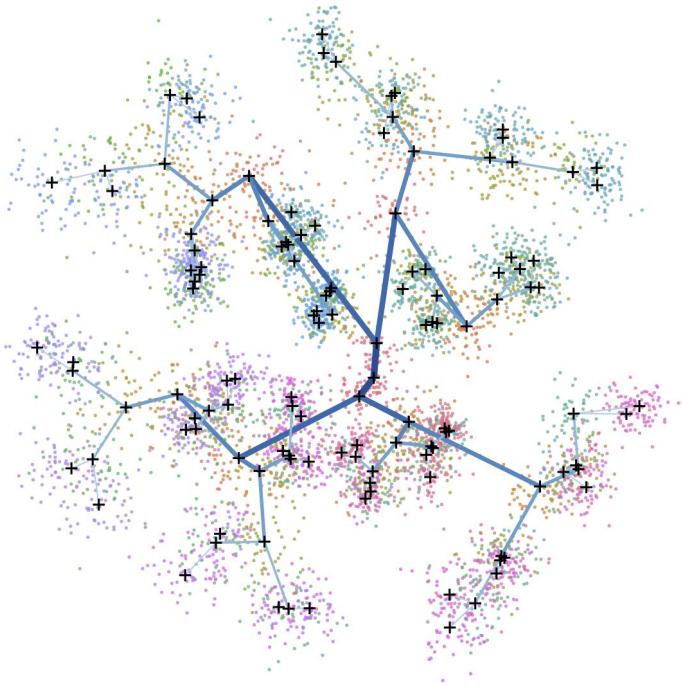


$\mathcal{N}$ -VAE

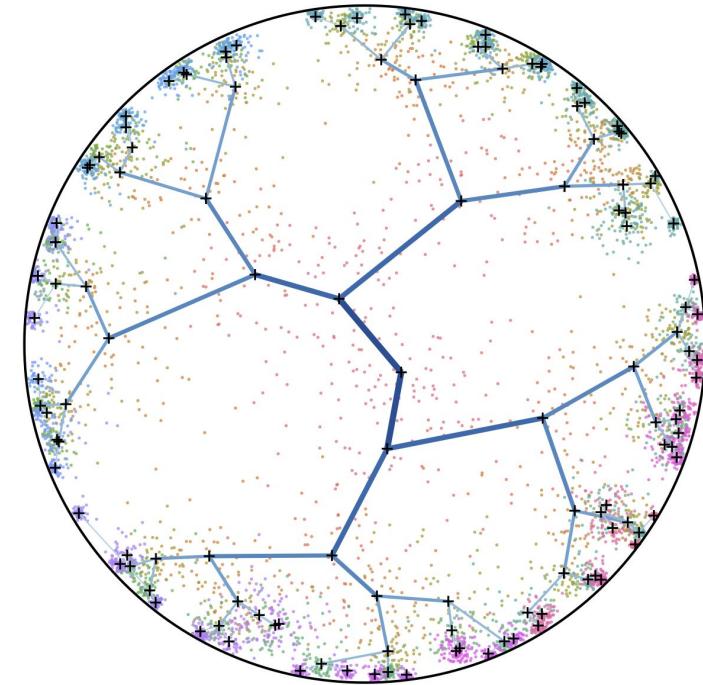


$\mathcal{P}^1$ -VAE

# Synthetic dataset: Branching diffusion process

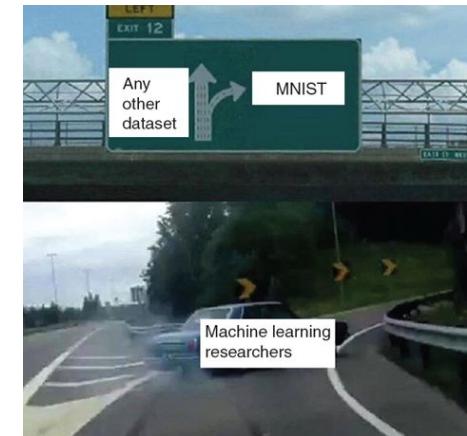
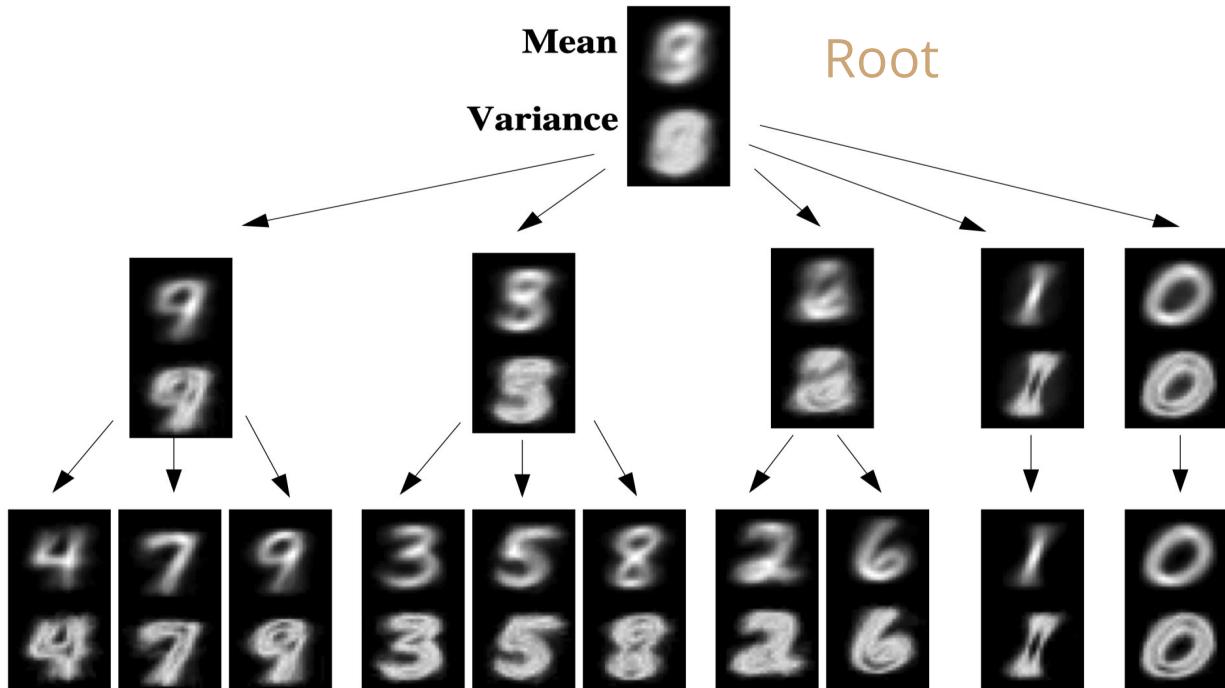


$\mathcal{N}$ -VAE



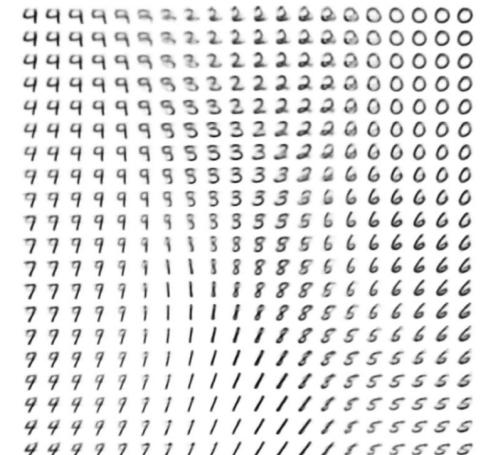
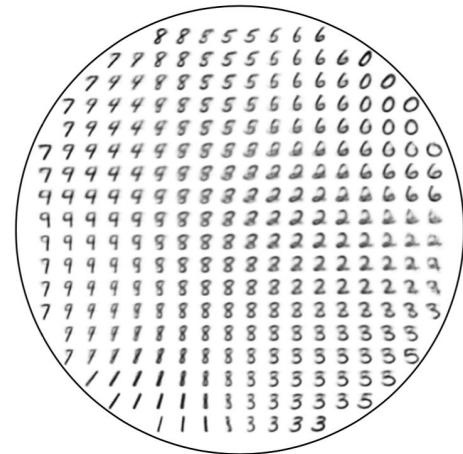
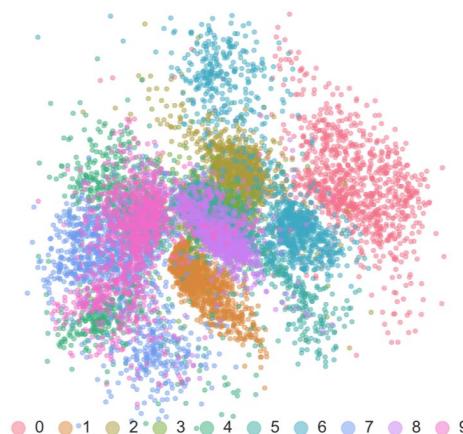
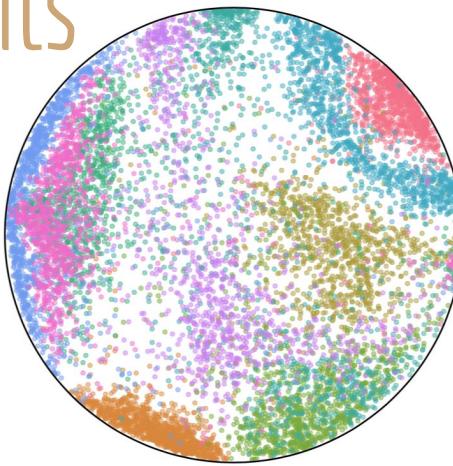
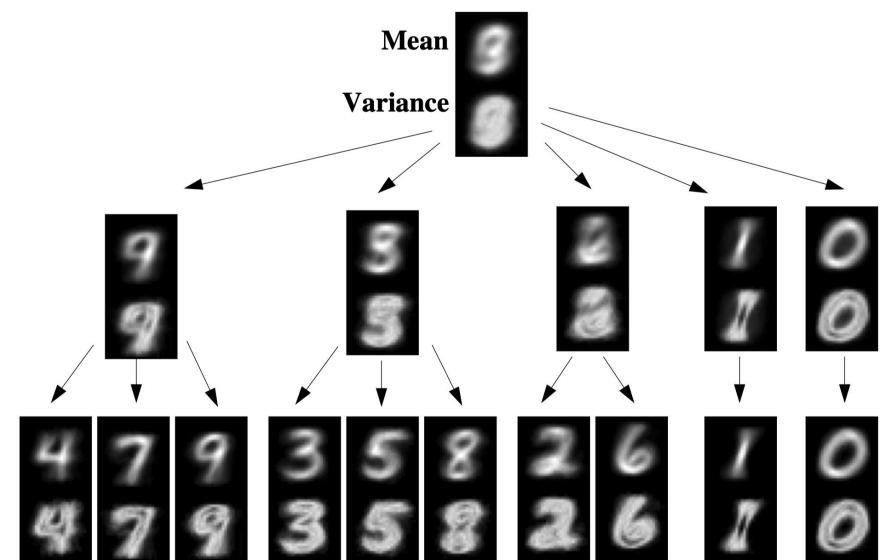
$\mathcal{P}^1$ -VAE

# MNIST: Handwritten digits



From Salakhutdinov et al., One-Shot Learning with a Hierarchical Nonparametric Bayesian Model

# MNIST: Handwritten digits



# MNIST: Handwritten digits

		Dimensionality										
		c	2	5	10	20						
<b><math>\mathcal{N}</math>-VAE</b>		(0)	$144.5 \pm 0.4$	$114.7 \pm 0.1$	$100.2 \pm 0.1$	$97.6 \pm 0.1$						
<b><math>\mathcal{P}</math>-VAE (Wrapped)</b>	0.1		$143.9 \pm 0.5$	$115.5 \pm 0.3$	$100.2 \pm 0.1$	$97.2 \pm 0.1$						
	0.2		$144.2 \pm 0.5$	$115.3 \pm 0.3$	$100.0 \pm 0.1$	$97.1 \pm 0.1$						
	0.7		$143.8 \pm 0.6$	$115.1 \pm 0.3$	$100.2 \pm 0.1$	$97.5 \pm 0.1$						
	1.4		$144.0 \pm 0.6$	$114.7 \pm 0.1$	$100.7 \pm 0.1$	$98.0 \pm 0.1$						
<b><math>\mathcal{P}</math>-VAE (Riemannian)</b>	0.1		$143.7 \pm 0.6$	$115.2 \pm 0.2$	$99.9 \pm 0.1$	$97.0 \pm 0.1$						
	0.2		$143.8 \pm 0.4$	$114.7 \pm 0.3$	$99.7 \pm 0.1$	$97.4 \pm 0.1$						
	0.7		$143.1 \pm 0.4$	$114.1 \pm 0.2$	$101.2 \pm 0.2$	*						
	1.4		<b><math>142.5 \pm 0.4</math></b>	$115.5 \pm 0.3$	*	*						
Digits		0	1	2	3	4	5	6	7	8	9	Avg
$\mathcal{N}$ -VAE		89	97	81	75	59	43	89	<b>78</b>	68	<b>57</b>	73.6
$\mathcal{P}^{1.4}$ -VAE		<b>94</b>	97	<b>82</b>	<b>79</b>	<b>69</b>	<b>47</b>	<b>90</b>	77	68	53	<b>75.6</b>

# MNIST: Handwritten digits

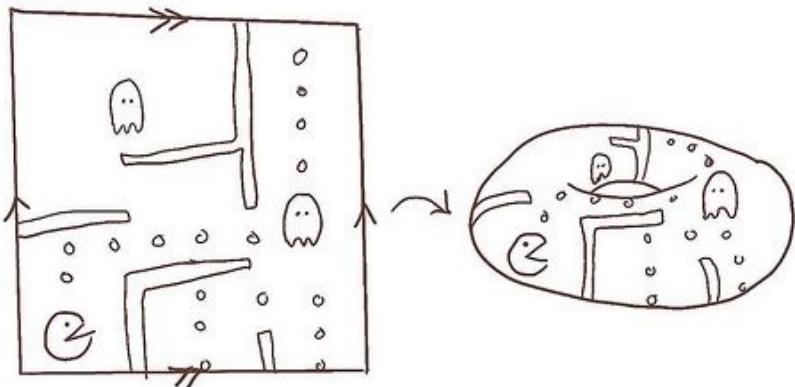
		Dimensionality				
		c	2	5	10	20
<b><math>\mathcal{N}</math>-VAE</b>		(0)	$144.5 \pm 0.4$	$114.7 \pm 0.1$	$100.2 \pm 0.1$	$97.6 \pm 0.1$
<b><math>\mathcal{P}</math>-VAE (Wrapped)</b>	0.1	$143.9 \pm 0.5$	$115.5 \pm 0.3$	$100.2 \pm 0.1$	$97.2 \pm 0.1$	
	0.2	$144.2 \pm 0.5$	$115.3 \pm 0.3$	$100.0 \pm 0.1$	$97.1 \pm 0.1$	
	0.7	$143.8 \pm 0.6$	$115.1 \pm 0.3$	$100.2 \pm 0.1$	$97.5 \pm 0.1$	
	1.4	$144.0 \pm 0.6$	$114.7 \pm 0.1$	$100.7 \pm 0.1$	$98.0 \pm 0.1$	
<b><math>\mathcal{P}</math>-VAE (Riemannian)</b>	0.1	$143.7 \pm 0.6$	$115.2 \pm 0.2$	$99.9 \pm 0.1$	<b><math>97.0 \pm 0.1</math></b>	
	0.2	$143.8 \pm 0.4$	$114.7 \pm 0.3$	<b><math>99.7 \pm 0.1</math></b>	$97.4 \pm 0.1$	
	0.7	$143.1 \pm 0.4$	<b><math>114.1 \pm 0.2</math></b>	$101.2 \pm 0.2$	*	
	1.4	<b><math>142.5 \pm 0.4</math></b>	$115.5 \pm 0.3$	*	*	

# Recap

- Deep generative model for hierarchical data
  - Leverage a hyperbolic latent space to induce hierarchical representations
  - Mostly useful for low latent dimensionality
-

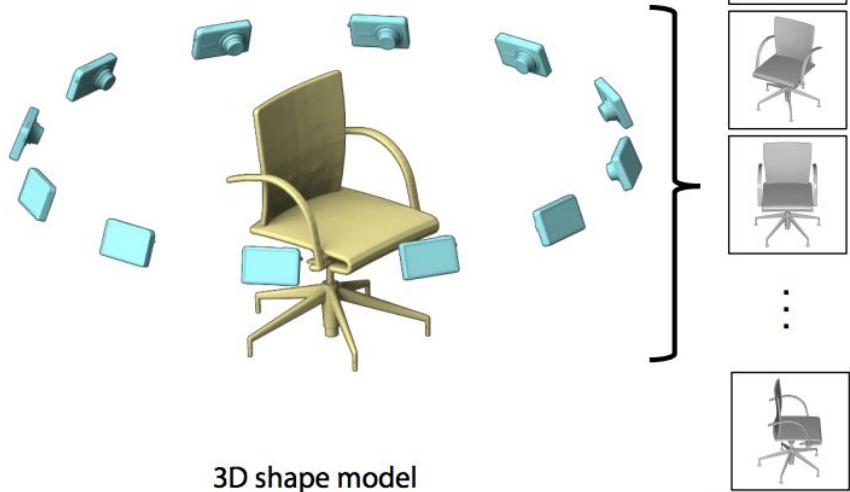
# Geometry as a manifold constraint

# Data with loops



The torus explained.

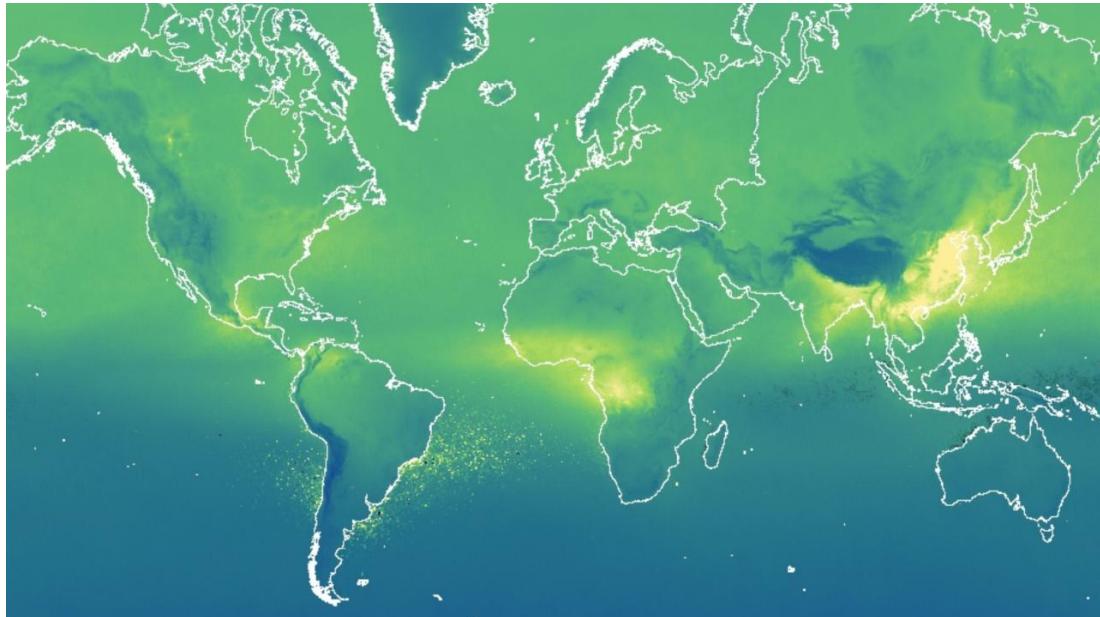
From Human Mathematics Blog, 2009



From Connelly Barnes, 3D Reconstruction and Understanding, 2017

# Earth science

$$\mathcal{M} = \mathbb{S}^2$$



*Concentrations of Carbon monoxide (CO) and water vapor.  
From The Sentinel-5 Precursor, European Space Agency.*

# Riemannian Continuous Normalizing Flows

Learning expressive probability  
distributions on manifolds

- Flows on manifolds as solution of ODEs
- Continuous change of variables for manifold-valued variables
- Vector field parametrization



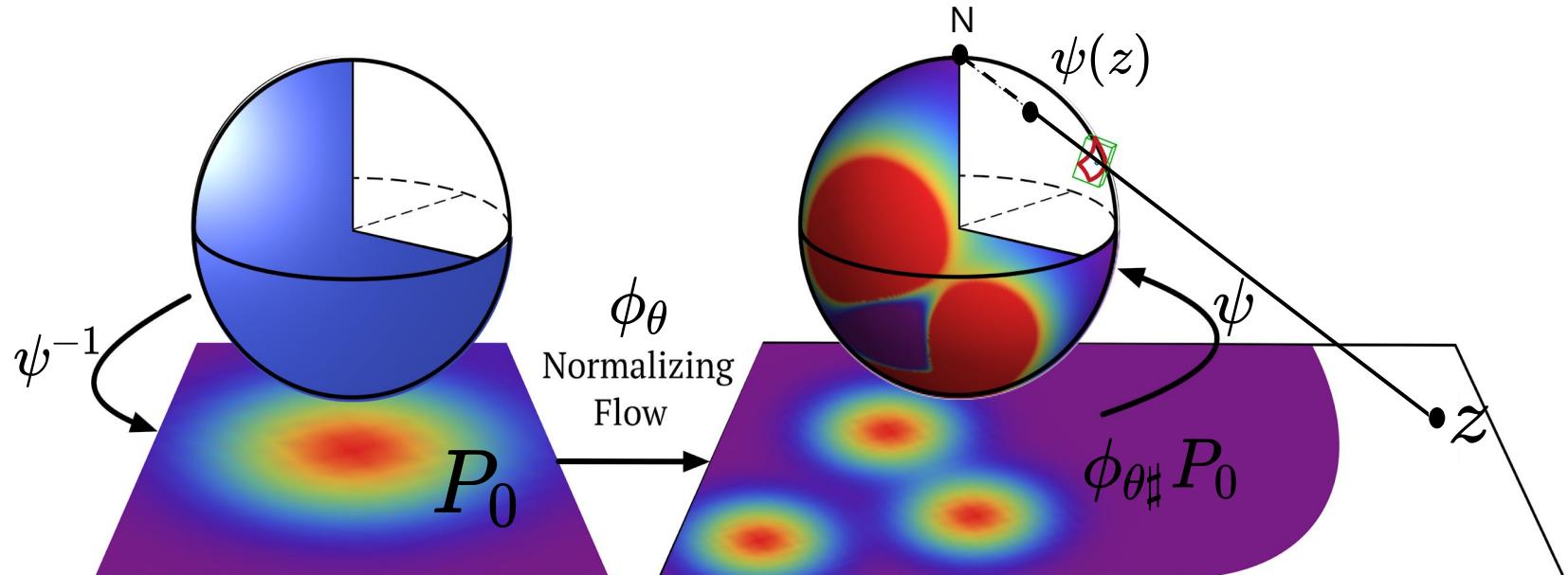
With Max Nickel

# Motivation

# Stereographic projection: An injection

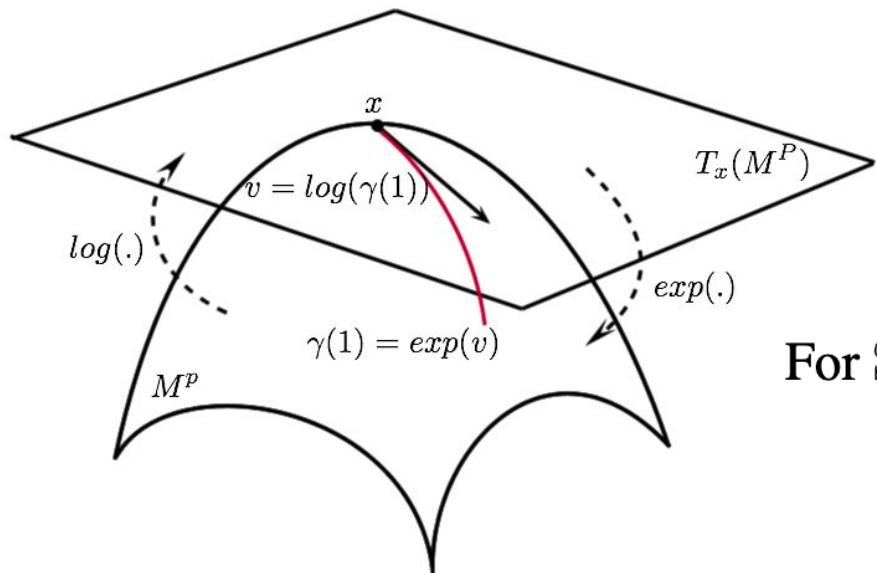
$$\psi : \mathbb{R}^d \rightarrow \mathcal{M}$$

$$P_\theta = (\psi \circ \phi_\theta)_\sharp P_0$$



From Gemici et al, 2016

# Exponential map: A surjection



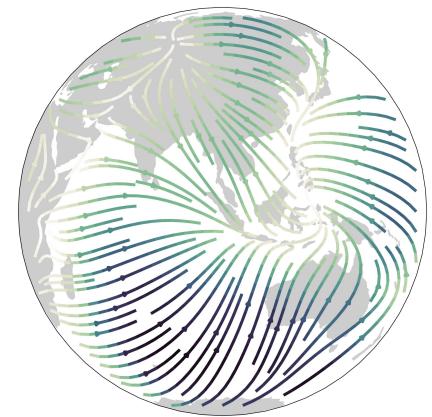
$$\exp_x : \mathcal{T}_x \mathcal{M} \cong \mathbb{R}^d \rightarrow \mathcal{M}$$

For  $\mathbb{S}^d$ ,  $\exp_x(v + 2k\pi) = \exp_x(v) \forall k \in \mathbb{Z}$

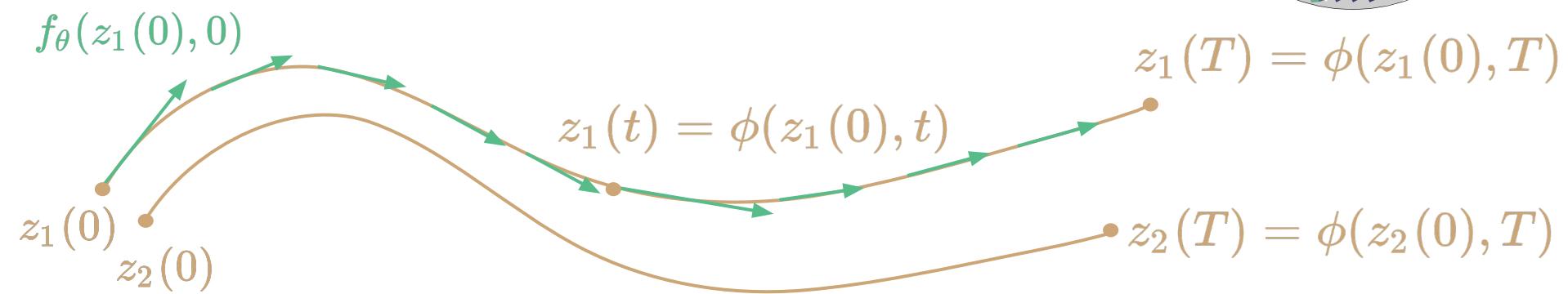
From Matthieu Simeoni, 2013

# Model

# Ordinary Differential Equation & Flow



**ODE**  $\frac{dz(t)}{dt} = f_\theta(z(t), t)$  vector field



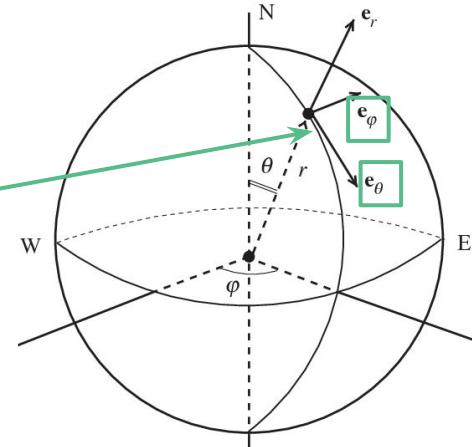
**Flow**  $\phi(z) := \phi(z, T)$   $\phi : \mathcal{M} \rightarrow \mathcal{M}$   $\mathcal{C}^1$ -diffeomorphism if  $f_\theta$  is  $\mathcal{C}^1$  and bounded

**Model**  $P_\theta = \phi_{\sharp} P_0$

# Vector field parametrization: local coordinates

Basis

$$f_\theta(z, t) = \sum_i^d f_\theta^i(z, t) \frac{\partial}{\partial z^i}(z)$$



Divergence

$$\operatorname{div}(f_\theta(z, t)) = \frac{1}{\sqrt{|G(z)|}} \sum_i^d \frac{\partial}{\partial z^i} \left( \sqrt{|G(z)|} f_\theta^i(z, t) \right)$$

Example

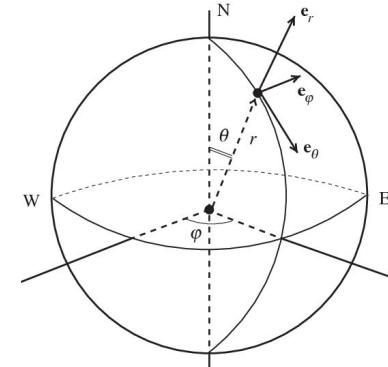
e.g.  $z = (\theta, \phi)$  polar coordinates for  $\mathbb{S}^2$

$$\sqrt{|G(\theta, \varphi)|} = \sin(\theta)$$

# Vector field parametrization: local coordinates

Basis

$$f_\theta(z, t) = \sum_i^d f_\theta^i(z, t) \frac{\partial}{\partial z^i}(z)$$



Divergence

$$\operatorname{div}(f_\theta(z, t)) = \frac{1}{\sqrt{|G(z)|}} \sum_i^d \frac{\partial}{\partial z^i} \left( \sqrt{|G(z)|} f_\theta^i(z, t) \right)$$

Example

$$z(\theta, \varphi) = (\sin(\theta) \cos(\varphi), \sin(\theta) \sin(\varphi), \cos(\theta))$$

$$\sqrt{|G(\theta, \varphi)|} = \sin(\theta)$$

$$\operatorname{div}(f(\theta, \varphi)) = \frac{1}{\sin(\theta)} \frac{\partial}{\partial \theta} (\sin(\theta) f^\theta(\theta, \varphi)) + \frac{1}{\sin(\theta)} \frac{\partial}{\partial \varphi} (f^\varphi(\theta, \varphi))$$

# Likelihood

## Continuous change of variable

$$\frac{\partial \log p_\theta(\mathbf{z}(t))}{\partial t} = -\operatorname{div}(\mathbf{f}_\theta(\mathbf{z}(t), t)) = -|G(\mathbf{z}(t))|^{-\frac{1}{2}} \operatorname{tr}\left(\frac{\partial \sqrt{|G(\mathbf{z}(t))|} \mathbf{f}_\theta(\mathbf{z}(t), t)}{\partial \mathbf{z}}\right) \quad (1)$$

## Stochastic estimator

$$\operatorname{div}(\mathbf{f}_\theta(\mathbf{z}(t), t)) = |G(\mathbf{z}(t))|^{-\frac{1}{2}} \mathbb{E}_{p(\boldsymbol{\epsilon})} \left[ \boldsymbol{\epsilon}^\top \frac{\partial \sqrt{|G(\mathbf{z}(t))|} \mathbf{f}_\theta(\mathbf{z}(t), t)}{\partial \mathbf{z}} \boldsymbol{\epsilon} \right]$$

# Vector field parametrization: divergence-free basis

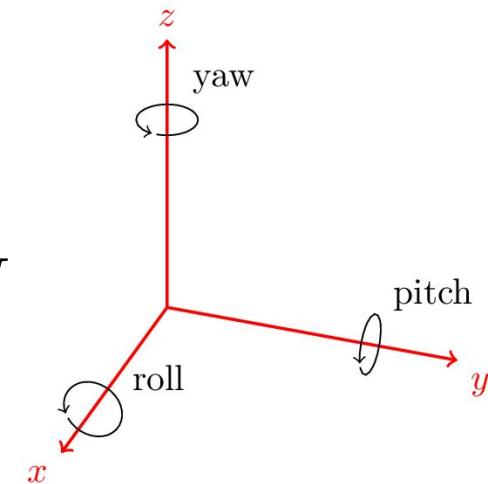
**Basis**  $f_i(z) = \frac{d}{dt} \Big|_{t=0} \exp(t \xi_i) \cdot z$

→  $\operatorname{div}(f_i(z)) = 0$

→  $(f_i(z); i = 1, \dots, n)$  span the tangent space  $T_z M$

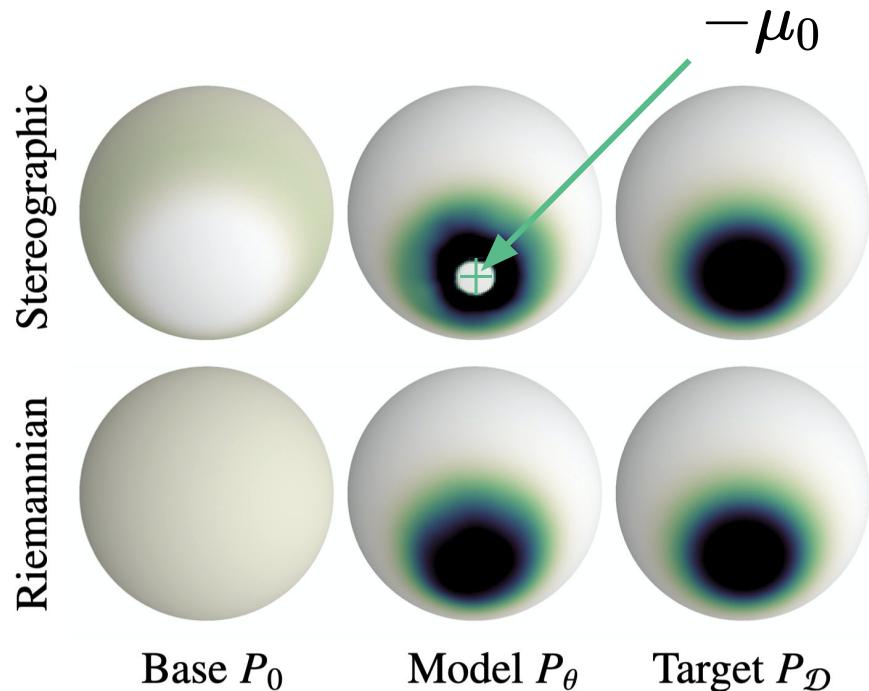
→  $f_\theta(z, t) = \sum_i^n f_\theta^i(z, t) f_i(z)$

**Divergence**  $\operatorname{div}(f_\theta(z, t)) = \sum_i^n \langle \frac{\partial}{\partial z} f^i(z, t), f_i(z) \rangle$



# Experiments

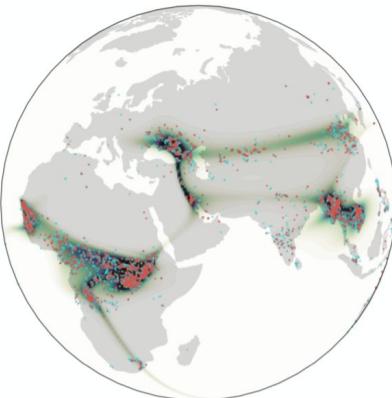
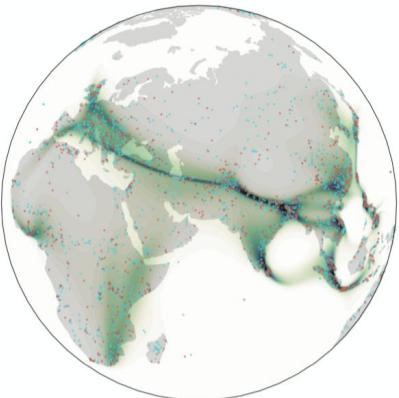
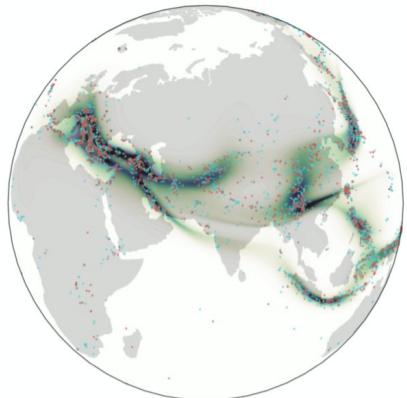
# Experiments: Von Mises-Fisher



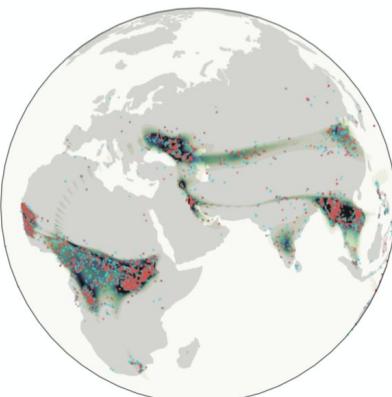
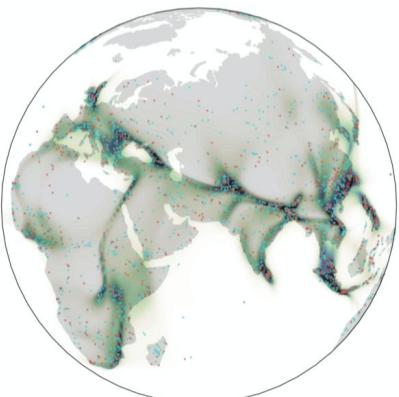
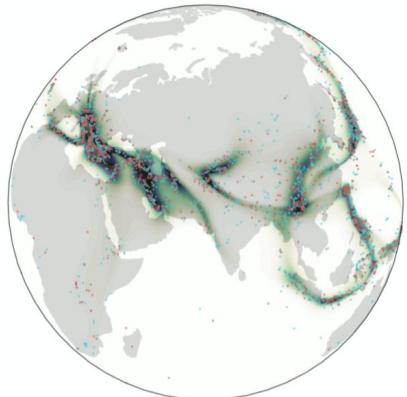
Loss	model $\kappa$	Stereographic	Riemannian
$\mathcal{L}^{\text{Like}}$	100	$63.60 \pm 3.56$	$-1.78 \pm 0.01$
	50	$32.68 \pm 3.15$	$-1.09 \pm 0.01$
	10	$6.45 \pm 2.42$	$0.52 \pm 0.01$
$\mathcal{L}^{\text{KL}}$	100	$1.56 \pm 0.34$	$0.04 \pm 0.02$
	50	$0.68 \pm 0.16$	$0.03 \pm 0.02$
	10	$0.12 \pm 0.01$	$0.01 \pm 0.00$

# Experimental Results: Earth Science

Stereographic



Riemannian

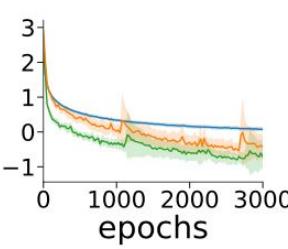
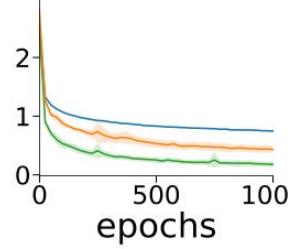
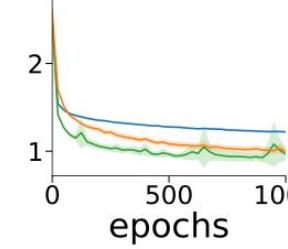
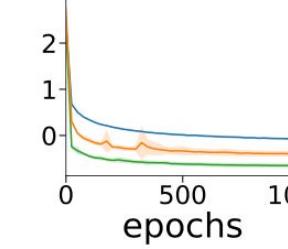


Earthquake

Flood

Fire

# Experimental Results: Earth Science

	Volcano	Earthquake	Flood	Fire
<b>Mixture vMF</b> ■	$-0.31 \pm 0.07$	$0.59 \pm 0.01$	$1.09 \pm 0.01$	$-0.23 \pm 0.02$
<b>Stereographic</b> □	$-0.64 \pm 0.20$	$0.43 \pm 0.06$	$0.99 \pm 0.04$	$-0.40 \pm 0.06$
<b>Riemannian</b> ▨	$-0.97 \pm 0.15$	$0.18 \pm 0.05$	$0.90 \pm 0.03$	$-0.66 \pm 0.05$
Learning curves				
Data size	829	6124	4877	12810

# Recap

- Extended neural ODEs to smooth manifolds
  - The generated flow is a diffeomorphism
  - Generalize better than the non-surjective stereographic approach
-

# Recap: Geometry and Deep generative models

- *Geometry as an inductive bias*
  - Hyperbolic geometry for hierarchical data
- *Geometry as a manifold constraint*
  - Normalizing Flows for manifold-valued data

Thank you for your attention!

Questions?