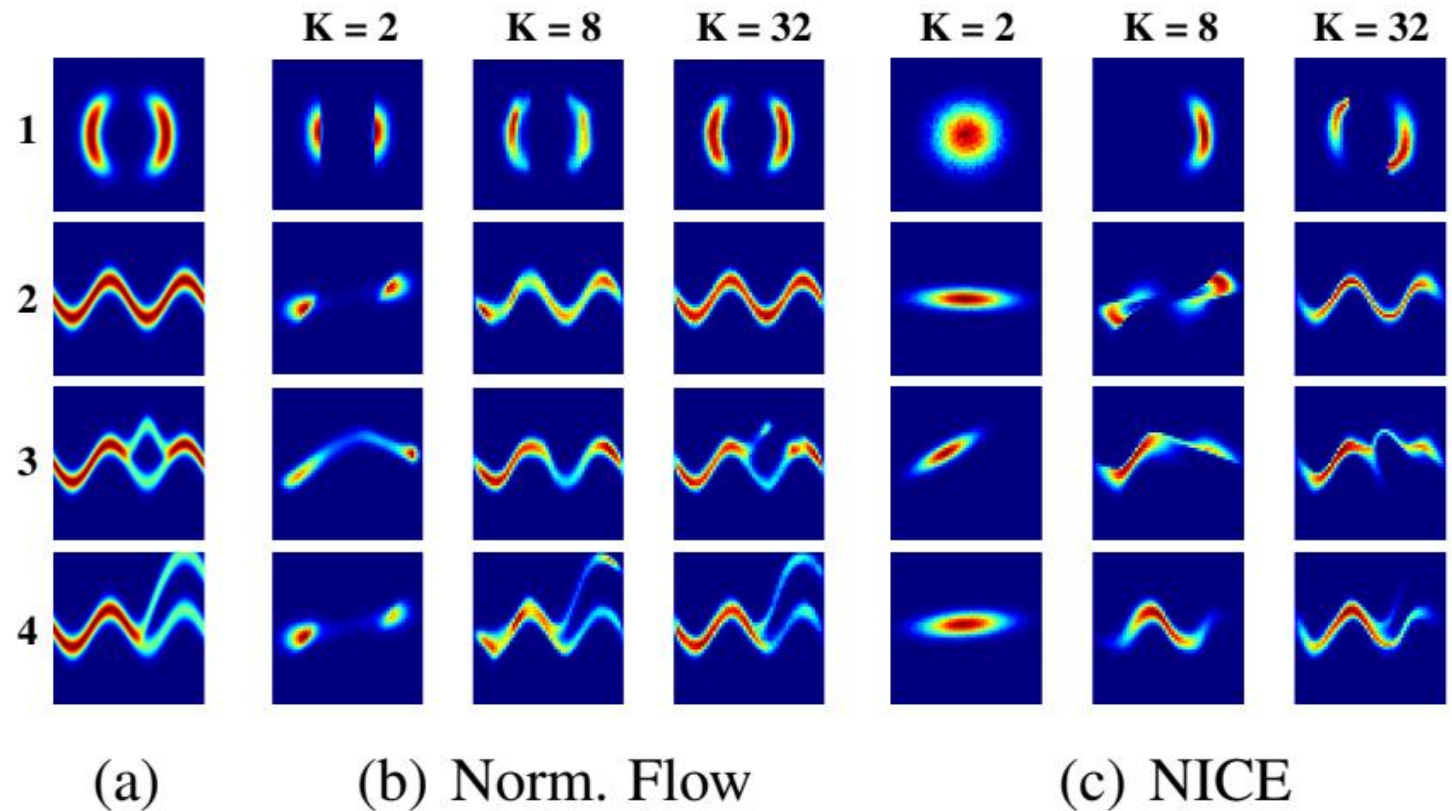
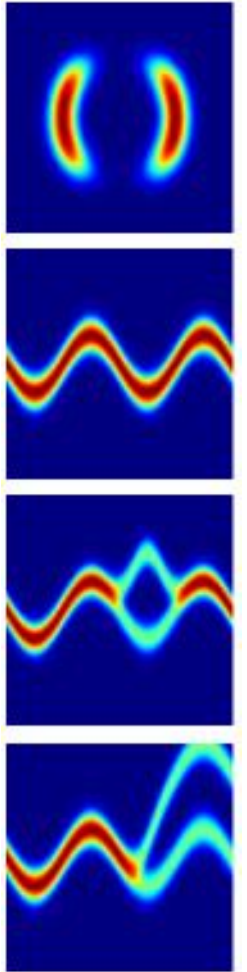
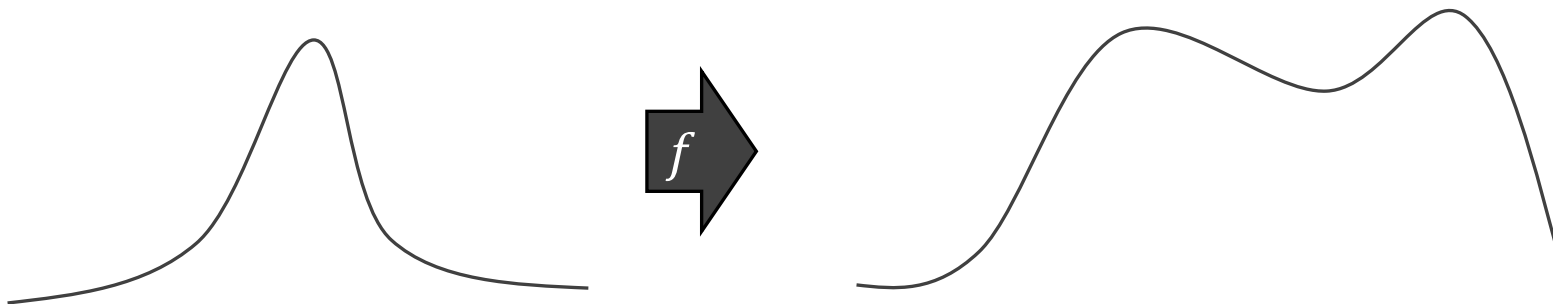


Normalizing flows



Intuition

- Often our posterior approximation is not enough
 - Imagine data generated by two modes
→ approximating with a standard Gaussian would be problematic
- If we applying a transformation to a simple density input
 - *e.g.*, a Gaussian
 - We can morph it into a more complicated density function
- Doing it many times → model any complex density



Change of variables

- For 1-d variables we know that

$$\int f(g(x))g'(x)dx = \int f(u)du, \text{ where } u = g(x)$$

- This is called change of variables (or [integration by substitution](#))
- The density is $du = g'(x)dx$

- For multivariate cases

$$\int f(g(\mathbf{x})) \left| \det \frac{dg}{d\mathbf{x}} \right| d\mathbf{x} = \int f(\mathbf{u}) d\mathbf{u}, \text{ where } \mathbf{u} = g(\mathbf{x})$$

And \mathbf{u} and \mathbf{x} have the same dimensionality

Normalizing flows

- Our model is an encoder $f: \mathbf{x} \rightarrow \mathbf{z}$
 - It maps the input \mathbf{x} with density $p(\mathbf{x})$ to the latent \mathbf{z} with density $p(\mathbf{z})$
- The inverse model is the decoder $f^{-1}: \mathbf{z} \rightarrow \mathbf{x}$
 - It maps the latent \mathbf{z} back to the input \mathbf{x}
- For our forward model we have

$$p(\mathbf{z}) \left| \det \frac{df}{d\mathbf{x}} \right| = p(\mathbf{x}) \Leftrightarrow \log p(\mathbf{z}) + \log \left| \det \left(\frac{df}{d\mathbf{x}} \right) \right| = \log p(\mathbf{x})$$

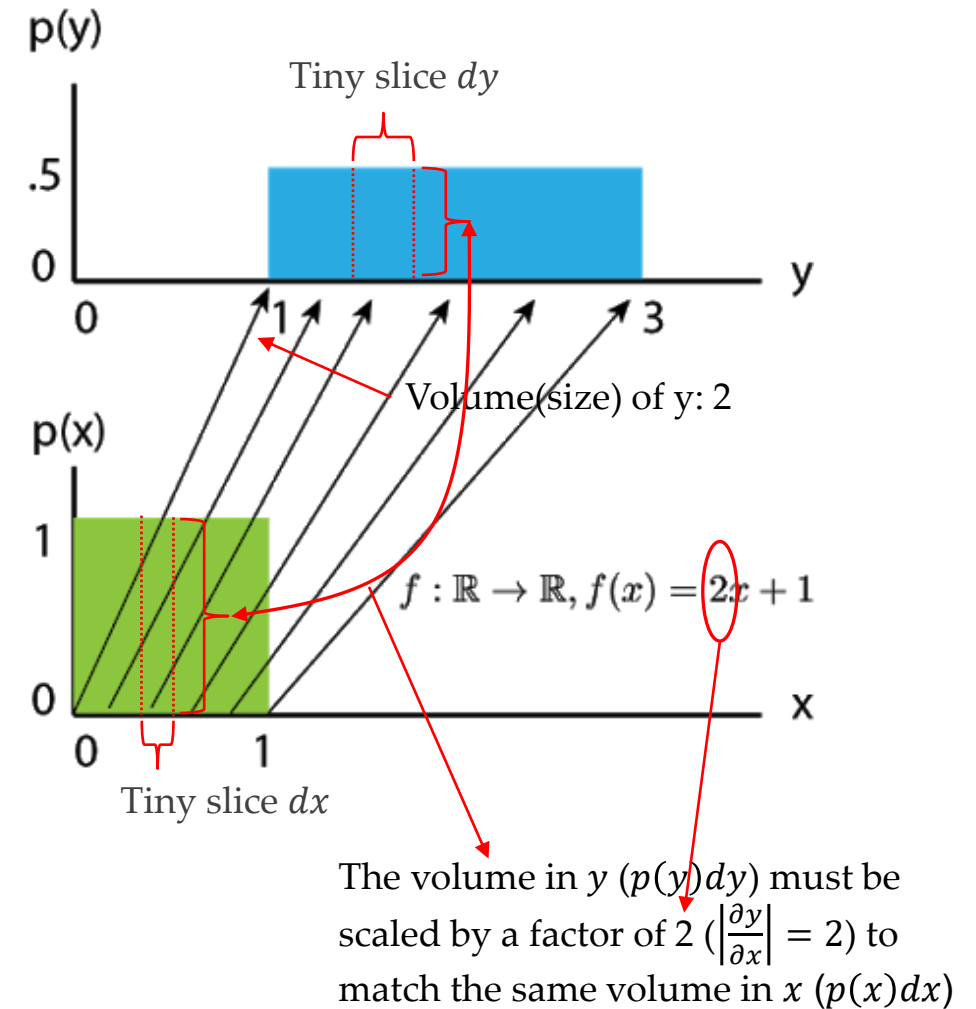
Normalizing flows geometrically

- The determinant shows the change in the volume of the input and output probability spaces

- 1-d case

$$p(y)dy = p(x)dx \Rightarrow p(y) \left| \frac{\partial y}{\partial x} \right| = p(x)$$

- $\left| \frac{dy}{dx} \right|$ indicates how much I must rescale a similarly tiny slice in y so the two densities are the same
- The volumes (sizes) rescale so that $\int p(\mathbf{z}) = 1$ and $\int p(\mathbf{x}) = 1$
- Normalizing flows expand or contract the density



Stacking normalizing flows

- The change of variables can be applied recursively

$$\mathbf{x} = \mathbf{z}_K = f_K^{-1} \circ \dots \circ f_2^{-1} \circ f_1^{-1}(\mathbf{z}_0)$$

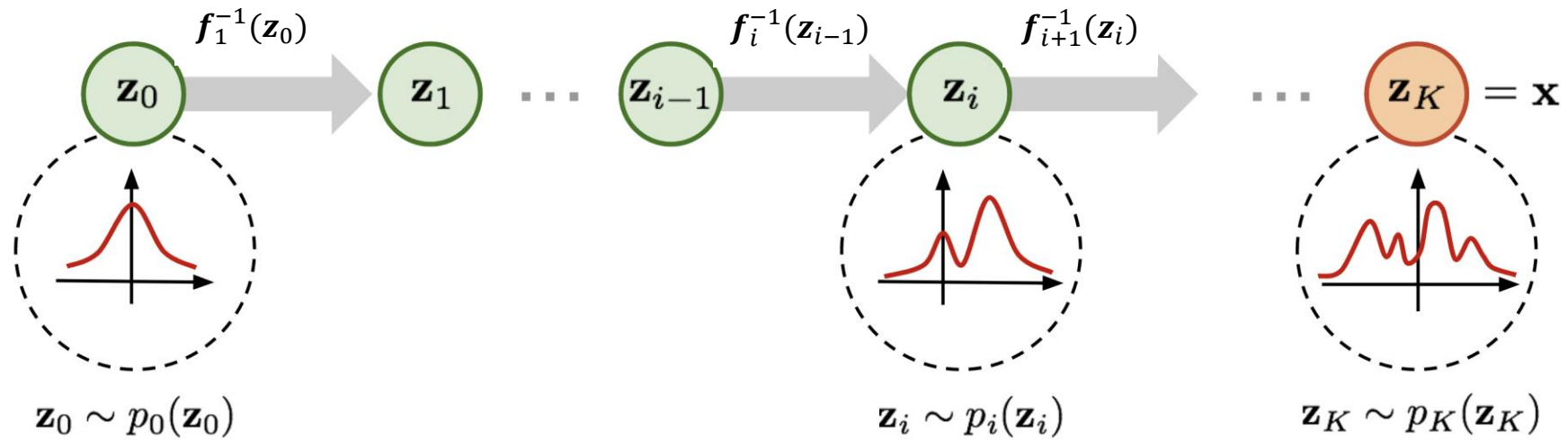
(e.g., $\mathbf{z}_1 = f_1^{-1}(\mathbf{z}_0) \Leftrightarrow \mathbf{z}_0 = f_1(\mathbf{z}_1)$)

- The log density of our data is

$$\log p(\mathbf{x}) = \log p_K(\mathbf{z}_K) = \log p_0(\mathbf{z}_0) - \sum_{k=1}^K \log \det \left(\frac{df_k}{d\mathbf{x}_k} \right)$$

- Optimize with maximum likelihood

Stacking normalizing flows



What transformations?

$$\log p_K(\mathbf{z}_K) = \log p_0(\mathbf{z}_0) - \sum_{k=1}^K \log \det \left(\frac{df}{d\mathbf{x}} \right)$$

- We want smooth, differentiable transformations f_k
 - For which it is easy to compute inverse f_k^{-1}
 - and determinant of the Jacobian $\det \frac{df_k}{dz_k}$
- Example transformations
 - Planar flows
 - Radial flows
 - Coupling layers

Planar flow

- The transformation is

$$f(\mathbf{z}) = \mathbf{z} + \mathbf{u}h(\mathbf{w}^T \mathbf{z} + b)$$

- $\mathbf{u}, \mathbf{w}, b$ are free parameters
- h is an element-wise non-linearity (element-wise so that it is easy to invert)
- The log-determinant of the Jacobian is

$$\psi(\mathbf{z}) = h'(\mathbf{w}^T \mathbf{z} + b)\mathbf{w}$$
$$\det\left(\frac{df}{d\mathbf{z}}\right) = |1 + u^T \psi(\mathbf{z})|$$

Radial flow

- The transformation is

$$f(\mathbf{z}) = \mathbf{z} + \beta h(\alpha, r)(\mathbf{z} - \mathbf{z}_0)$$

Where $h(\alpha, r) = 1/(\alpha + r)$

- The log-determinant of the Jacobian is

$$\det\left(\frac{df}{d\mathbf{z}}\right) = [1 + \beta h(\alpha, r)]^{d-1} [1 + \beta h(\alpha, r) + h'(\alpha, r)r]$$

VAE with Normalizing Flows

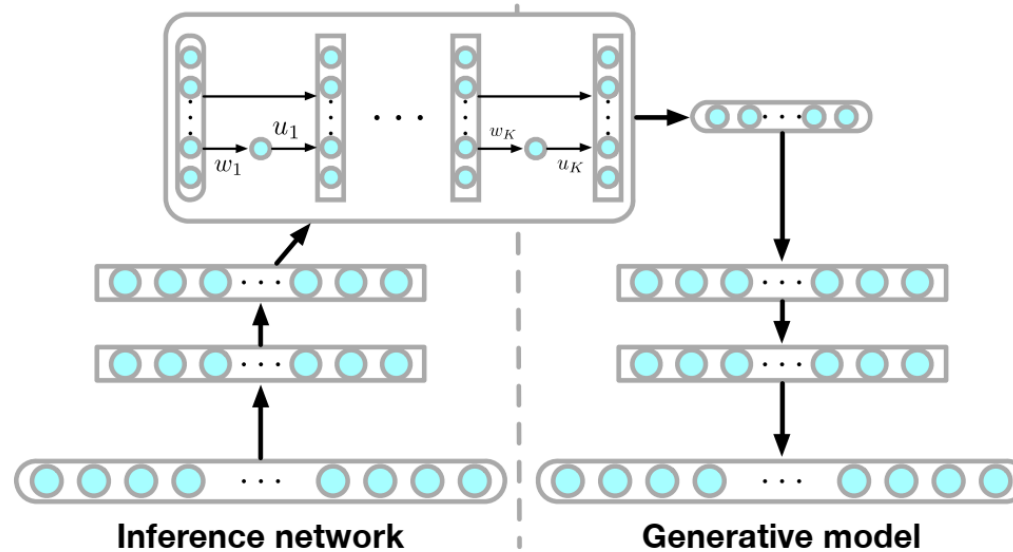


Figure 2. Inference and generative models. Left: Inference network maps the observations to the parameters of the flow; Right: generative model which receives the posterior samples from the inference network during training time. Round containers represent layers of stochastic variables whereas square containers represent deterministic layers.

Learning better posteriors with variational inference

- Again: the evidence lower bound is

$$\text{ELBO}_{\theta, \phi}(\mathbf{x}) = \log p(\mathbf{x}) - \text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}))$$

- Replace the simple approximate posterior by normalizing flows

$$\mathbb{E}_{q_0(\mathbf{z}_0|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z}_K)] - \text{KL}(q_0(\mathbf{z}_0|\mathbf{x}) \parallel p(\mathbf{z})) + \mathbb{E}_{q_0(\mathbf{z}_0|\mathbf{x})} \left[\sum_{k=1}^K \log \left| \det \frac{df_k}{dz_k} \right|^{-1} \right]$$

Algorithm 1 Variational Inf. with Normalizing Flows

Parameters: ϕ variational, θ generative

while not converged **do**

$\mathbf{x} \leftarrow \{\text{Get mini-batch}\}$

$\mathbf{z}_0 \sim q_0(\bullet|\mathbf{x})$

$\mathbf{z}_K \leftarrow f_K \circ f_{K-1} \circ \dots \circ f_1(\mathbf{z}_0)$

$\mathcal{F}(\mathbf{x}) \approx \mathcal{F}(\mathbf{x}, \mathbf{z}_K)$

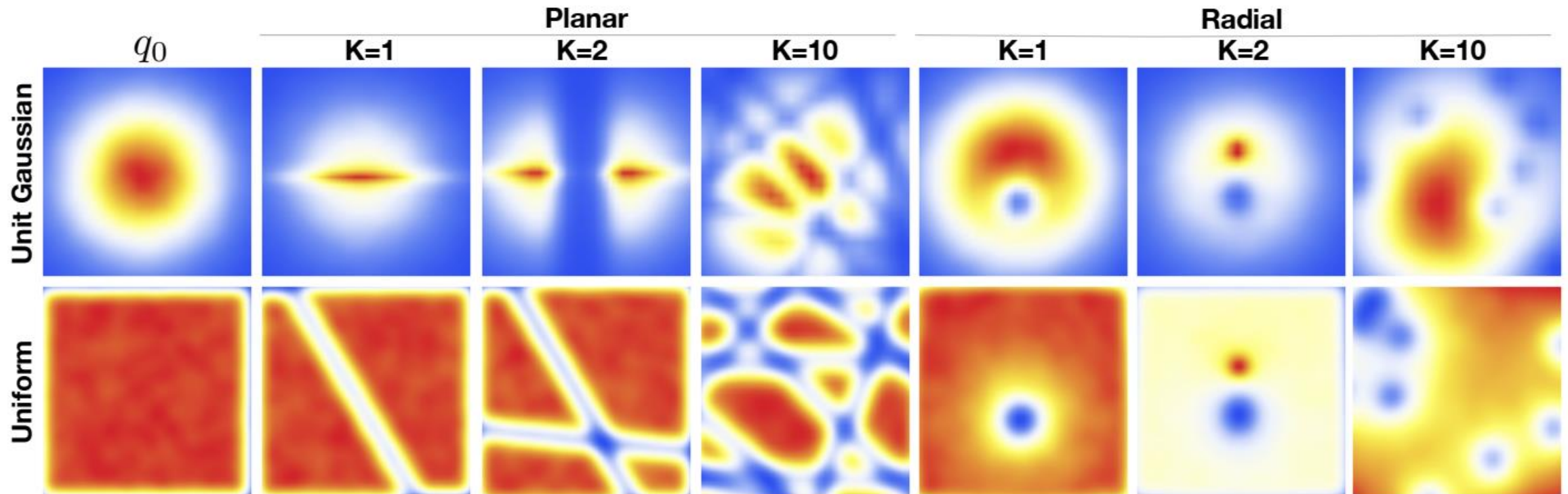
$\Delta\theta \propto -\nabla_{\theta}\mathcal{F}(\mathbf{x})$

$\Delta\phi \propto -\nabla_{\phi}\mathcal{F}(\mathbf{x})$

end while

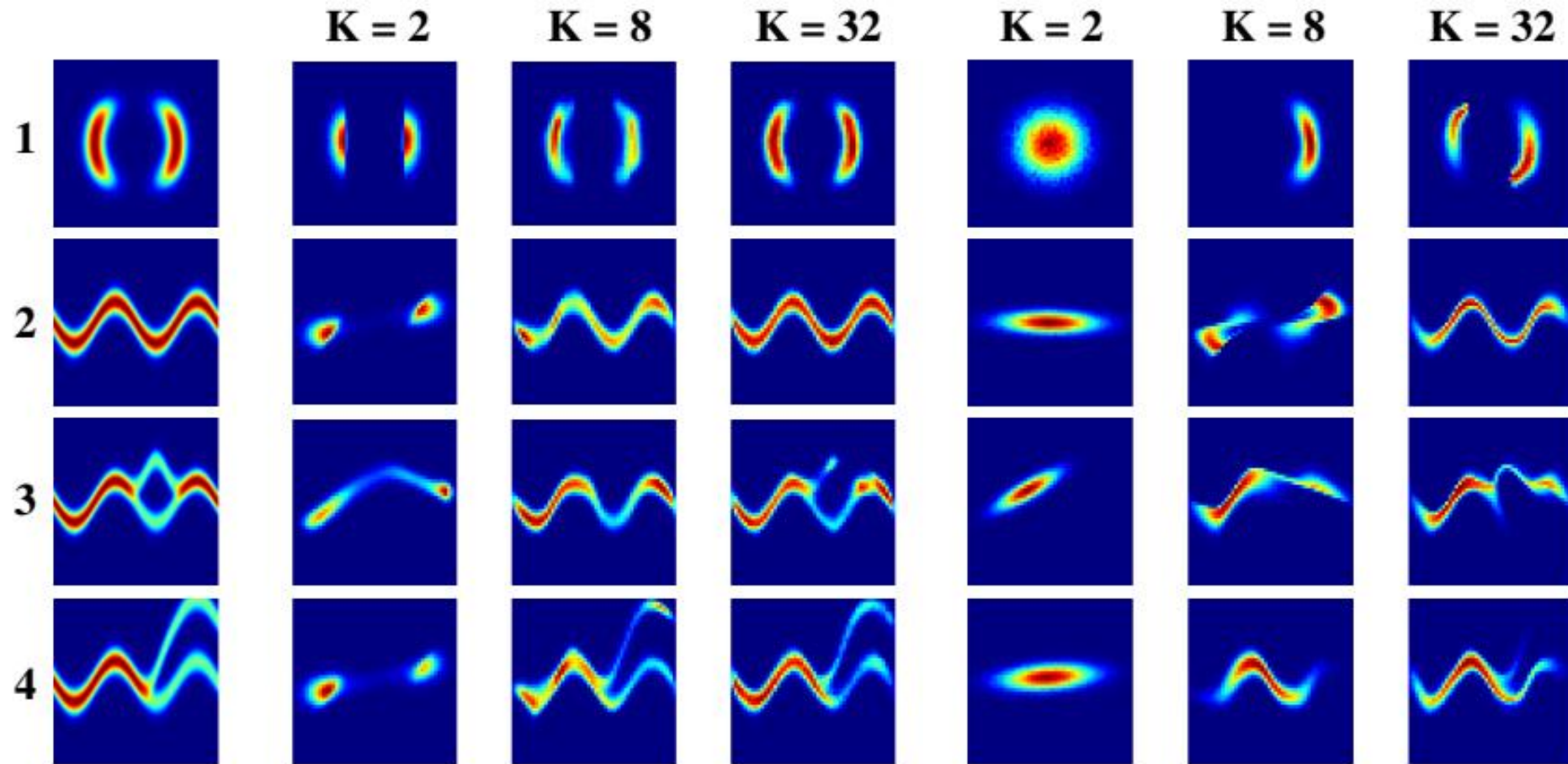
Rezende and Mohamed, Variational Inference with Normalizing Flows

The effect of number of transformations/flows



Rezende and Mohamed, Variational Inference with Normalizing Flows

Some results



(a)

(b) Norm. Flow

(c) NICE

Rezende and Mohamed, Variational Inference with Normalizing Flows