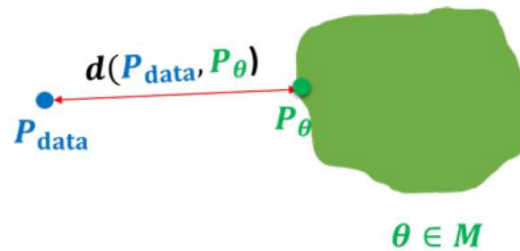


Normalizing flows

Recap of likelihood-based learning so far:

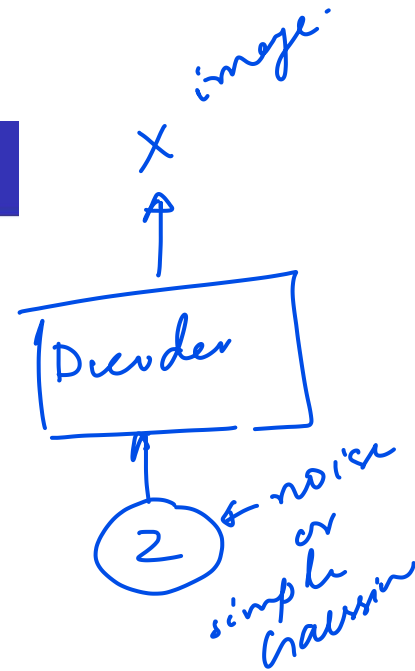


$$\mathbf{x}_i \sim P_{\text{data}} \\ i = 1, 2, \dots, n$$



Model family

- Model families: *Graphical model to enforce joint distribution.*
 - Autoregressive Models: $p_{\theta}(\mathbf{x}) = \prod_{i=1}^n p_{\theta}(x_i | \mathbf{x}_{<i})$ ✗
 - Variational Autoencoders: $p_{\theta}(\mathbf{x}) = \int_{\mathbf{z}} p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z}$
- Autoregressive models provide tractable likelihoods but no direct mechanism for learning features : *actual dependency is best captured latent in a space.*
- Variational autoencoders can learn feature representations (via latent variables \mathbf{z}) but have intractable marginal likelihoods
- Key question:** Can we design a latent variable model with tractable likelihoods? Yes!



Simple Prior to Complex Data Distributions

- Desirable properties of any model distribution:
 - Analytic density. —
 - Easy-to-sample — + Latent space of features.
- Many simple distributions satisfy the above properties e.g., Gaussian, uniform distributions, or a Bayesian network. — //
- Unfortunately, data distributions could be much more complex (multi-modal) —
- **Key idea**: Map simple distributions (easy to sample and evaluate densities) to complex distributions (learned via data) using **change of variables**. —

Normalizing Flows for Probabilistic Modeling and Inference

George Papamakarios*

GPAPAMAK@GOOGLE.COM

Eric Nalisnick*

ENALISNICK@GOOGLE.COM

Danilo Jimenez Rezende

DANILOR@GOOGLE.COM

Shakir Mohamed

SHAKIR@GOOGLE.COM

Balaji Lakshminarayanan

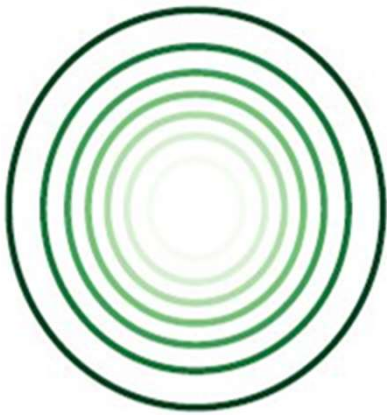
BALAJILN@GOOGLE.COM

DeepMind

Abstract

Normalizing flows provide a general mechanism for defining expressive probability distributions, only requiring the specification of a (usually simple) base distribution and a series of bijective transformations. There has been much recent work on normalizing flows, ranging from improving their expressive power to expanding their application. We believe the field has now matured and is in need of a unified perspective. In this review, we attempt to provide such a perspective by describing flows through the lens of probabilistic modeling and inference. We place special emphasis on the fundamental principles of flow design, and discuss foundational topics such as expressive power and computational trade-offs. We also broaden the conceptual framing of flows by relating them to more general probability transformations. Lastly, we summarize the use of flows for tasks such as generative modeling, approximate inference, and supervised learning.

Basics of Normalizing Flows



$$\underline{\mathbf{u}} \sim p(\mathbf{u}) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

latent space \mathbf{z}

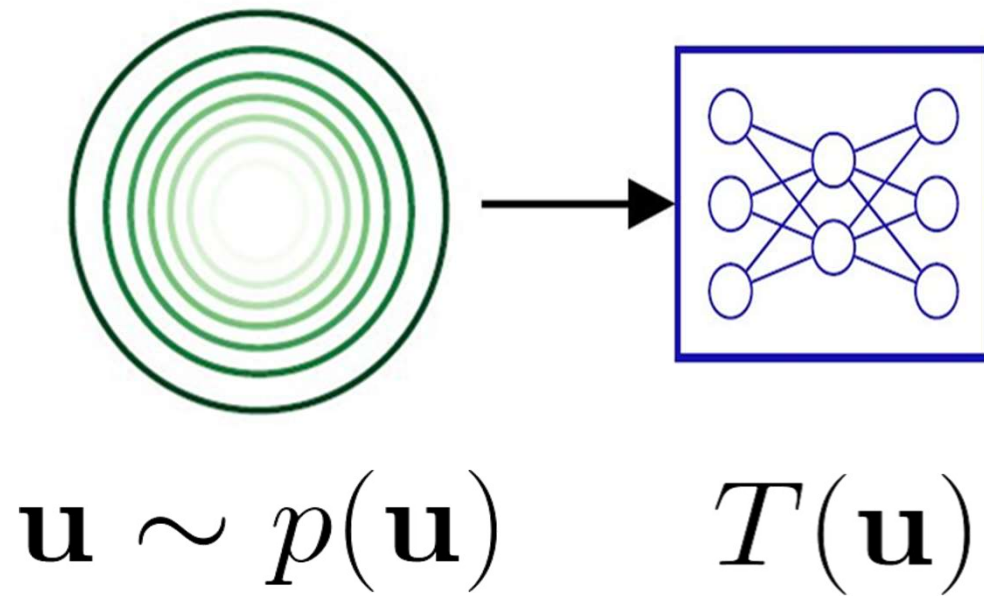
/

.

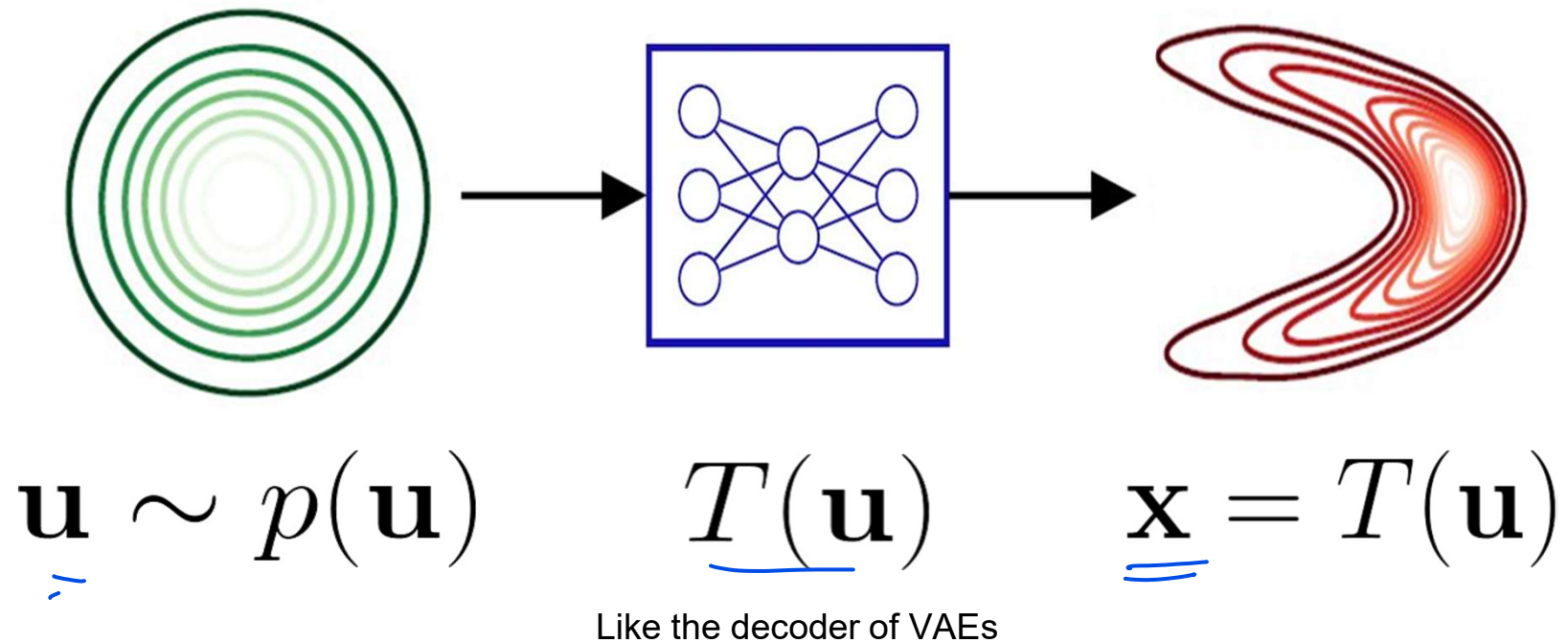
.

.

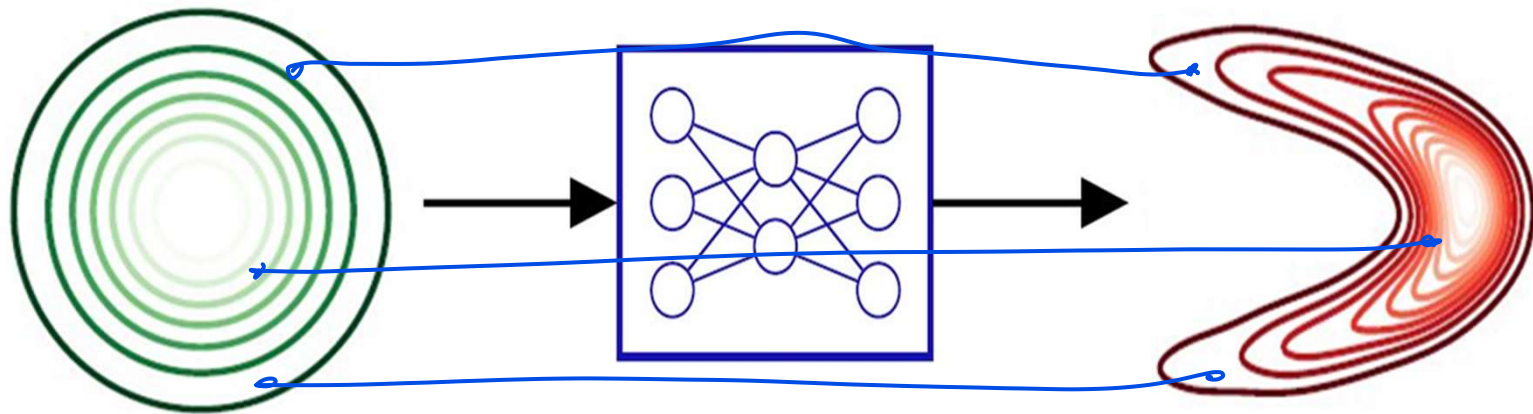
Basics of Normalizing Flows



Basics of Normalizing Flows



Basics of Normalizing Flows



$$\mathbf{u} \sim p(\mathbf{u}) \xrightarrow{\text{bijection } T} \mathbf{x} = T(\mathbf{u})$$

$\tilde{\mathbf{u}} = T^{-1}(\mathbf{x})$ T must be invertible

$$\frac{p(\tilde{\mathbf{x}})}{p(\tilde{\mathbf{x}})} = p(\tilde{\mathbf{u}})$$

Calculating Density $p(x)$ using change of variables

$$x = T(u) \quad x \in \mathbb{R}, u \in \mathbb{R} \quad x \quad u$$

What is $p(x)$?

Consider 1-d data $x \in \mathbb{R}$

$$\underline{F_x(x)} = P(\underline{x} \leq \underline{x})$$

$$= P(u \leq T^{-1}(x)) = F_u(T^{-1}(x))$$

← CDF of U

chain rule of differentiation

$$\begin{aligned} \underline{p(x)} &= \frac{\partial}{\partial x} \underline{F_x(x)} = \frac{\partial}{\partial x} \underline{F_u(T^{-1}(x))} = \frac{\frac{\partial}{\partial u} F_u(u)}{\frac{\partial u}{\partial x}} \frac{\partial T^{-1}(x)}{\partial x} \\ &= \frac{p(u)}{u} \frac{\partial T^{-1}(x)}{\partial x} \end{aligned}$$

$$\boxed{\frac{\partial T^{-1}(x)}{\partial x} = \left(\frac{\partial x}{\partial T^{-1}(x)} \right)^{-1} = \left(\frac{\partial T(u)}{\partial u} \right)^{-1}}$$

Change of variable formula for 1-D variables.

$$\begin{aligned} \log \underline{p_x(x)} &= \log p_u(u) \left[\frac{\partial T(u)}{\partial u} \right]^{-1} \\ &= \log \underline{p_u(u)} = \log \frac{\partial T(u)}{\partial u} \end{aligned}$$

Jacobians

$x \in \mathbb{R}^D$, $u \in \mathbb{R}^D$ T is invertible.

$$x_1 = T_1(u) = T_1(u_1, \dots, u_D)$$

$$x_2 = T_2(u) = T_2(u_1, \dots, u_D)$$

$$\vdots$$

$$x_D = T_D(u) = T_D(u_1, \dots, u_D)$$

Jacobian derivative

$$J_T = \begin{bmatrix} \frac{\partial T_1}{\partial u_1} & \dots & \frac{\partial T_1}{\partial u_D} \\ \vdots & & \vdots \\ \frac{\partial T_D}{\partial u_1} & \dots & \frac{\partial T_D}{\partial u_D} \end{bmatrix}$$

$$\underline{J_{T^{-1}}}(x) = \begin{bmatrix} \frac{\partial (T^{-1})_1}{\partial x_1} & \dots & \frac{\partial (T^{-1})_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial (T^{-1})_n}{\partial x_1} & \dots & \frac{\partial (T^{-1})_n}{\partial x_n} \end{bmatrix}$$

$$\underline{J_T}(u) = \left[\underline{J_{T^{-1}}}(x) \right]^{-1} \text{ where } x = T(u)$$


Change-of-Variables Formula

$$\log \underline{p_X(\mathbf{x})} = \log \underline{p_u(\mathbf{u})} - \log \underline{|\det J_T(\mathbf{u})|}$$

Jacobian matrix of T

$$J_T(\mathbf{u}) = \begin{bmatrix} \frac{\partial T_1}{\partial u_1} & \cdots & \frac{\partial T_1}{\partial u_D} \\ \vdots & \ddots & \vdots \\ \frac{\partial T_D}{\partial u_1} & \cdots & \frac{\partial T_D}{\partial u_D} \end{bmatrix}$$

Change-of-Variables Formula

$$\log p_{\mathbf{x}}(\mathbf{x}) = \log p_{\mathbf{u}}(\mathbf{u}) - \log |\det J_T(\mathbf{u})|$$


Jacobian matrix of T



In practice, two requirements on T :

1. Invertible ←
2. Easy-to-compute determinant of Jacobian

Flows Support Two Core Operations

Sampling:

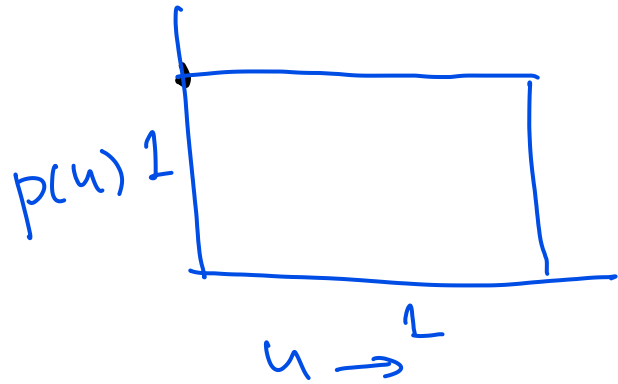
$$\hat{\mathbf{u}} \sim p(\mathbf{u}), \quad \hat{\mathbf{x}} = T(\hat{\mathbf{u}})$$

Density Evaluation:

$$\log p_{\mathbf{x}}(\mathbf{x}) = \log p_{\mathbf{u}}(T^{-1}(\mathbf{x})) + \log |\det J_{T^{-1}}(\mathbf{x})|$$

Example Uniform distribution

$$u \sim U(0, 1)$$



$$x = 4u = T(u) \quad T^{-1}(x) = \frac{x}{4}$$
$$p(x) = p(u) \frac{\partial T^{-1}(x)}{\partial x}$$

$$1 \cdot \frac{\partial (x/4)}{\partial x} = \frac{1}{4}$$

$$p(x) = \frac{1}{4}$$

$$x \sim U(0, \underline{\underline{\frac{1}{4}}}) \quad p(x) = \underline{\underline{\frac{1}{4}}}$$

Normalizing flows - simplest case - linear transform in 1D

Normalizing Flows allow for defining complex densities by transforming simple one by invertible mappings i.e. bijections.

Let's build a simple model for the simple example considered in the previous slide.

- we start from the simplest univariate gaussian distribution

$$u \sim \mathcal{N}(0, 1) \quad p(u) = \frac{e^{-\frac{u^2}{2}}}{\sqrt{2\pi}}$$

- we define a affine transformation (a forward pass)

$$x = \mu + u\sigma \quad x = \mu + u\sigma = T(u)$$

- and inverse (with constraint sigma > 0):

$$\varepsilon = (x - \mu) / \sigma \quad u = \frac{T^{-1}(x)}{\sigma} = \frac{(x - \mu)}{\sigma}$$

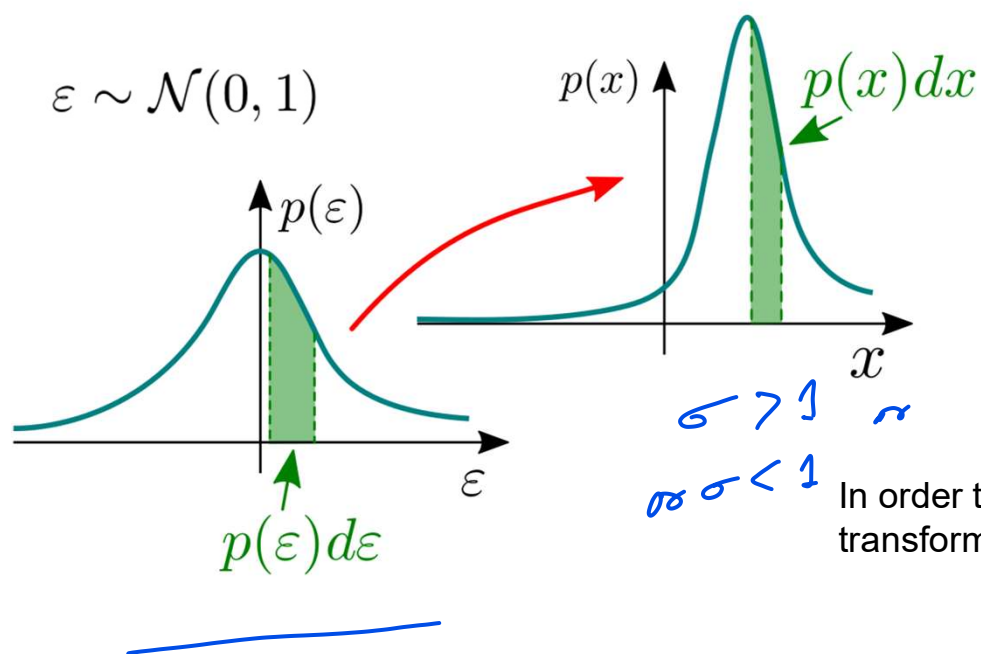
- we know $p(u)$ and we want to find $p(x)$

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$\begin{aligned} p(x) &= p(u) \left| \frac{\partial T^{-1}(x)}{\partial x} \right| \\ &= \frac{e^{-\frac{(x-\mu)^2}{2\sigma^2}}}{\sqrt{2\pi}} \left[\frac{\partial}{\partial x} \left(\frac{x-\mu}{\sigma} \right) \right] = \end{aligned}$$

Normalizing flows - transforming scalars

The volume change, but density must be preserved:



The necessary condition for this is:

$$p(\varepsilon)d\varepsilon = p(x)dx$$

The transformed density is then:

$$x = \mu + \varepsilon\sigma \quad \text{change of volume}$$

$$p(x) = p(\varepsilon) \left| \frac{dx}{d\varepsilon} \right|^{-1} = \frac{p(\varepsilon)}{\sigma}$$

In order to estimate the density at x we have to apply inverse transform:

$$\varepsilon = (x - \mu) / \sigma$$

Which will result in:

$$p_{\text{model}}(x; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{(x - \mu)^2}{2\sigma^2} \right]$$

$$D=2$$

$$T_1(\vec{u}) = T_1(u_1, u_2)$$

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} \frac{2u_1 + 3u_2}{u_1 - u_2} \end{bmatrix}$$

$T(u)$

$$\underline{T^{-1}(x)} = \begin{bmatrix} (x_1 + 3x_2)/5 \\ (x_1 - 2x_2)/5 \end{bmatrix}$$

$$p(\underline{u}_1, \underline{u}_2) \sim \underline{N}(0, I) = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$$

$$p(\underline{x}_1, \underline{x}_2) = \frac{e^{-\frac{u_1^2 + u_2^2}{2}}}{\sqrt{2\pi}}$$

Gaussian Copulas

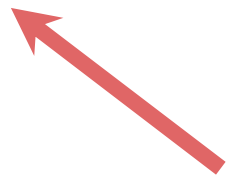
Straightforward to Fit via Divergence of Choice

Example: *Maximum Likelihood Estimation*

$$D_{\text{KL}} [p^*(\mathbf{x}) \parallel p(\mathbf{x})] :$$

Data generating
process

Flow model



Straightforward to Fit via Divergence of Choice

Example: Maximum Likelihood Estimation

$$\frac{1}{N} \sum_{i=1}^N \log p(x^i)$$

$$D_{\text{KL}} [p^*(\mathbf{x}) || p(\mathbf{x})] = -\mathbb{E}_{p^*} [\log p(\mathbf{x})] + \text{const.}$$

Data generating
process

Flow model

Straightforward to Fit via Divergence of Choice

Example: *Maximum Likelihood Estimation*

$$D_{\text{KL}} [p^*(\mathbf{x}) \parallel p(\mathbf{x})] = -\mathbb{E}_{p^*} [\log p(\mathbf{x})] + \text{const.}$$

$$= -\mathbb{E}_{p^*} \left[\log p_{\text{u}}(T^{-1}(\mathbf{x})) + \log |\det J_{T^{-1}}(\mathbf{x})| \right] + \text{const.}$$

Universal Representation

Can any distribution be represented as a flow?

*modest conditions apply

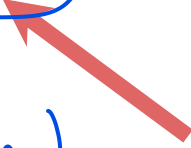
Universal Representation

Can any distribution be represented as a flow?

Yes! The intuition is given by inverse transform sampling...

$$\begin{aligned} \underline{u} &\sim \text{Uniform}(0, 1) \\ \underline{x} &= \text{CDF}^{-1}(u) \end{aligned}$$

$T(u)$



*modest conditions apply

Don't know this in practice.

Key Properties for Tractability

Key Properties for Tractability

#1 Compositionality

#2 Link Between Chain Rule and Triangular Jacobians

#3 Free to Choose Directionality

Key Properties for Tractability

#1 Compositionality

#2 Link Between Chain Rule and Triangular Jacobians

#3 Free to Choose Directionality

#1 Compositionality

Transformations can be composed without violating invertibility:

$$\underline{T} = T_K \circ \dots \circ T_1$$

#1 Compositionality

Transformations can be composed without violating invertibility:

$$T = T_K \circ \dots \circ T_1$$

The det-Jacobian decomposes locally over sub-flows:

$$\log |J_T(\mathbf{u})| = \log \prod_k |J_{T_k}(\mathbf{u})| = \sum_k \log |J_{T_k}(\mathbf{u})|$$

Handwritten notes:
- A blue arrow points from the word "determinant" to $|J_T(\mathbf{u})|$.
- A blue line underlines $J_{T_k}(\mathbf{u})$.
- A blue line underlines the entire right-hand side of the equation.

#1 Compositionality

Transformations can be composed without violating invertibility:

$$T = \underline{T_K} \circ \dots \circ \underline{T_1}$$

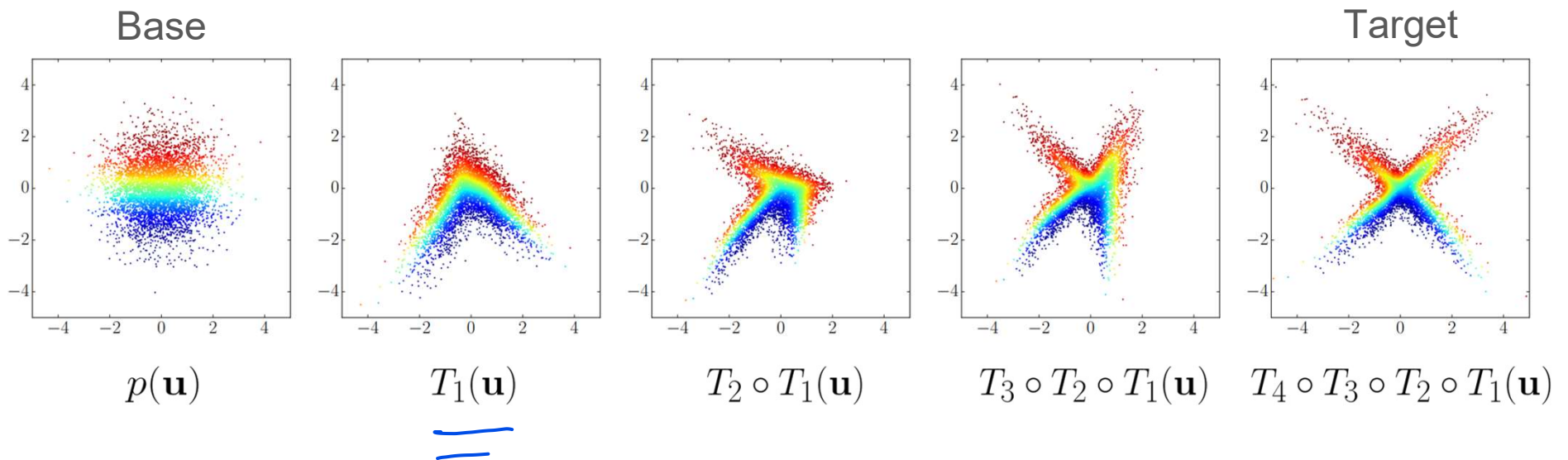
The det-Jacobian decomposes locally over sub-flows:

$$\log |J_T(\mathbf{u})| = \log \prod_k |J_{T_k}(\mathbf{u})| = \sum_k \log |J_{T_k}(\mathbf{u})|$$

Pay only $O(K)$ cost for composition!

#1 Compositionality

Linear cost allows expressive transforms to be defined by composing many simple sub-flows...



Key Properties for Tractability

#1 Compositionality

#2 Link Between Chain Rule and Triangular Jacobians

#3 Free to Choose Directionality