

# Normalizing Flows

---

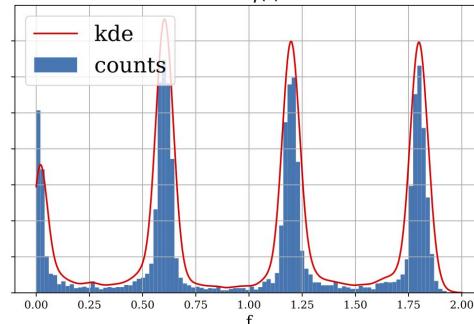
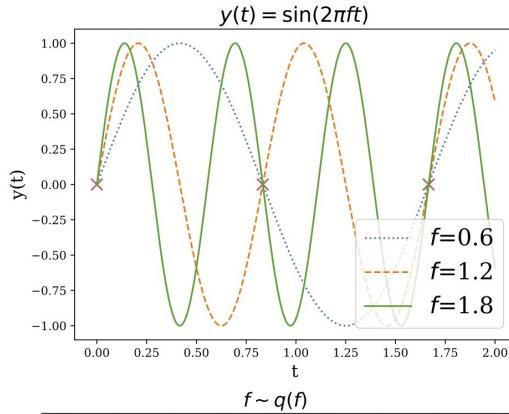
Chin-Wei Huang  
(2020/03/16)

# Invertible Neural network and why is it appealing?

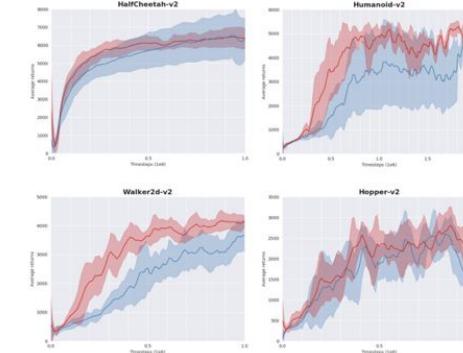
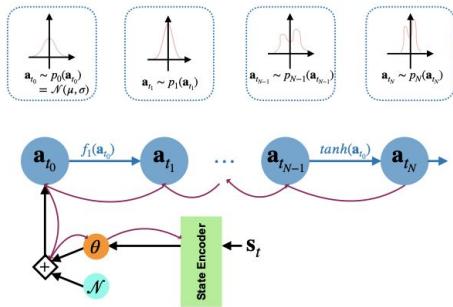
Generative modeling



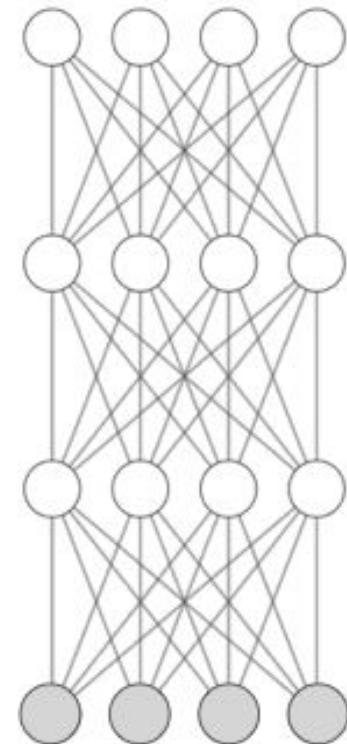
Bayesian inference



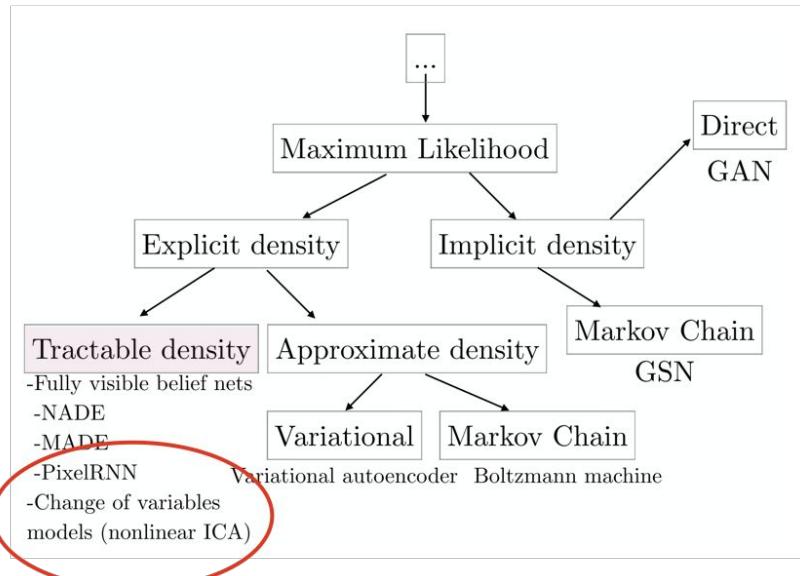
Exploration in RL



O(1) mem backprop



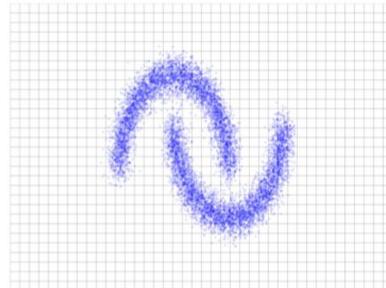
# Explicit Density Model Without Approximation



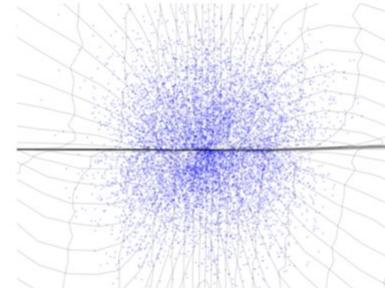
### Inference

$$x \sim \hat{p}_X$$
$$z = f(x)$$

Data space  $\mathcal{X}$

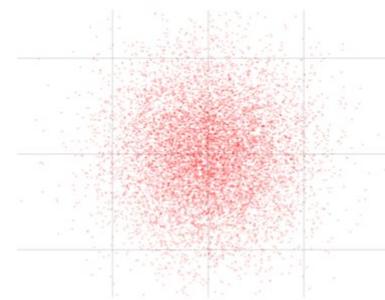
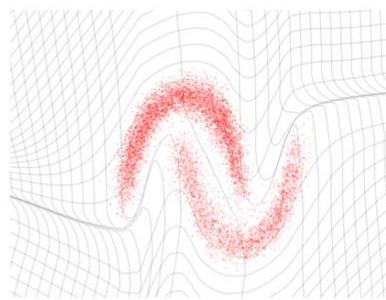


Latent space  $\mathcal{Z}$



### Generation

$$z \sim p_Z$$
$$x = f^{-1}(z)$$



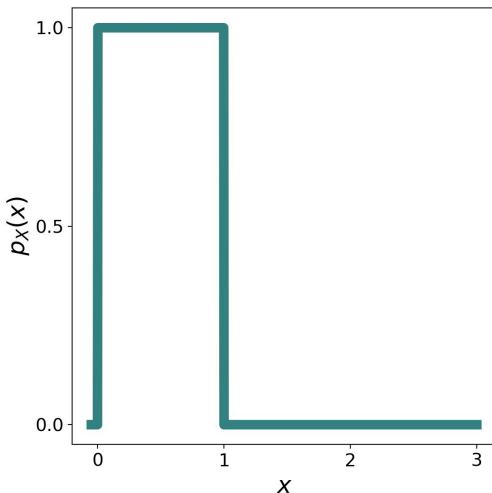
# Normalizing Flows

Translating probability distributions

# Change of Variable Density Needs to Be Normalized

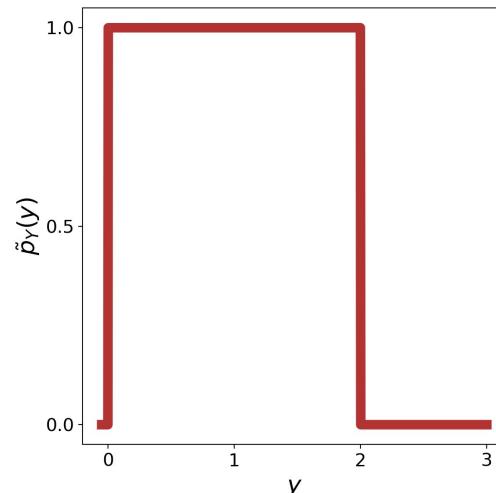
$$X \sim p_X$$

$$p_X(x) = \begin{cases} 1 & \text{for } 0 \leq x \leq 1 \\ 0 & \text{else} \end{cases}$$

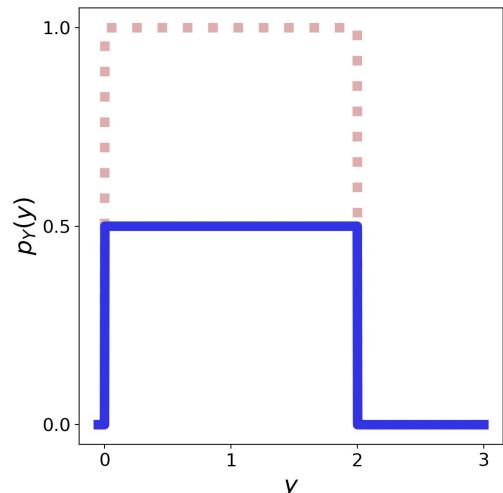


$$Y := 2X$$

$$\tilde{p}_Y(y) = p_X(y/2)$$



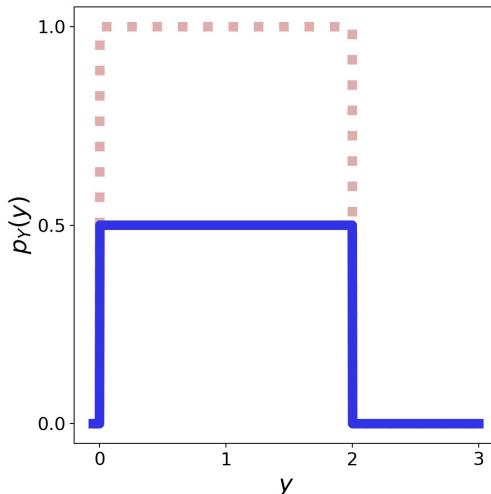
$$p_Y(y) = p_X(y/2)/2$$



# Change of Variable Density (m-Dimensional)

For a multivariable invertible mapping  $f : \mathbb{R}^m \rightarrow \mathbb{R}^m$   $X \sim p_X$   $Y := f(X)$

$$p_Y(y) = p_X(f^{-1}(y)) \left| \det \frac{\partial f^{-1}(y)}{\partial y} \right|$$



$$Y := 2X$$

$$p_Y(y) = p_X(y/2)/2$$

Local change of volume

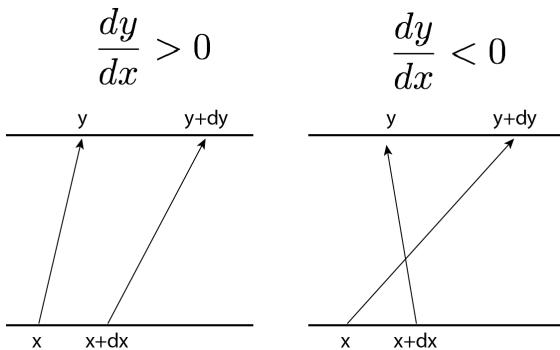
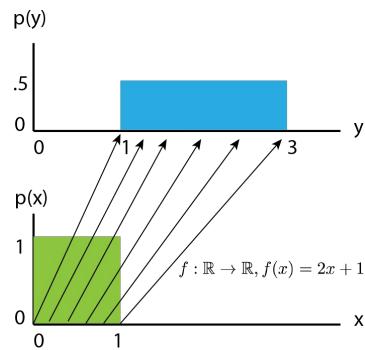
mass = density \* **volume**

# Change of Variable Density (m-Dimensional)

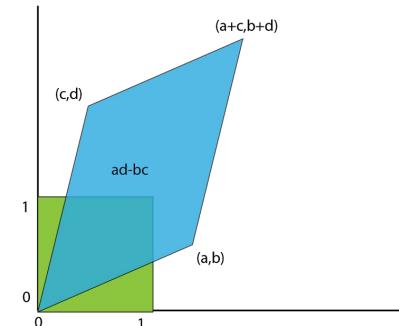
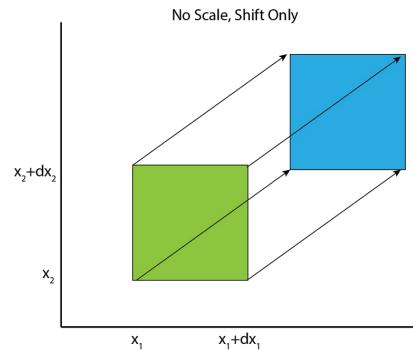
For a multivariable invertible mapping  $f : \mathbb{R}^m \rightarrow \mathbb{R}^m$   $X \sim p_X$   $Y := f(X)$

$$p_Y(y) = p_X(f^{-1}(y)) \left| \det \frac{\partial f^{-1}(y)}{\partial y} \right|$$

## 1-D

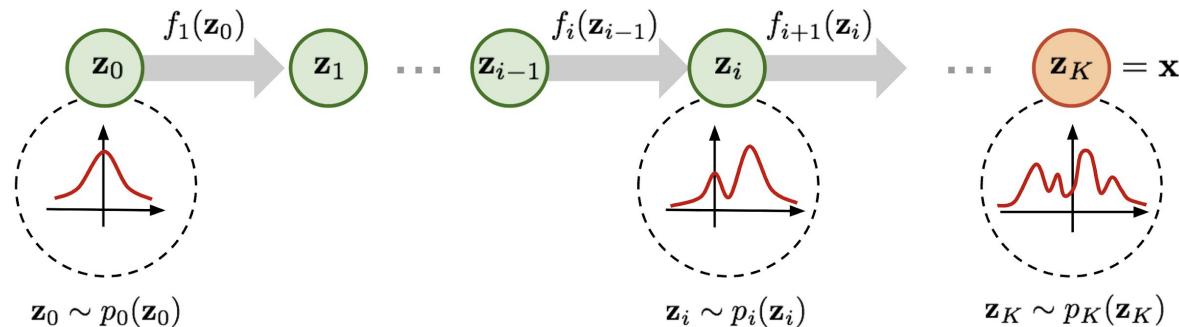


## 2-D



# Chaining Invertible Mappings (Composition)

$$f = f_S \circ \cdots \circ f_2 \circ f_1 \quad f(x) = f_S(\cdots f_2(f_1(x)))$$



$$\frac{\partial f(x)}{\partial x} = \frac{f_S(x_{S-1})}{\partial x_{S-1}} \cdots \frac{f_2(x_1)}{\partial x_1} \frac{f_1(x_0)}{\partial x_0} \quad x_s = f_s(x_{s-1}) \quad x_0 = x$$

Chain rule

$$\det \left( \frac{\partial f(x)}{\partial x} \right) = \det \left( \frac{f_S(x_{S-1})}{\partial x_{S-1}} \right) \cdots \det \left( \frac{f_2(x_1)}{\partial x_1} \right) \det \left( \frac{f_1(x_0)}{\partial x_0} \right)$$

Determinant of matrix product

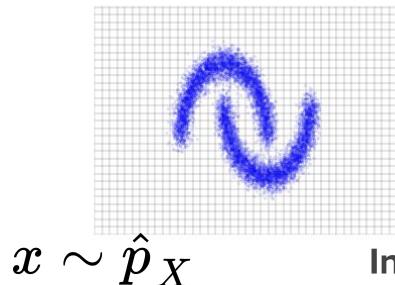
# Training with Maximum Likelihood Principle

$$Z \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad X = g(Z)$$

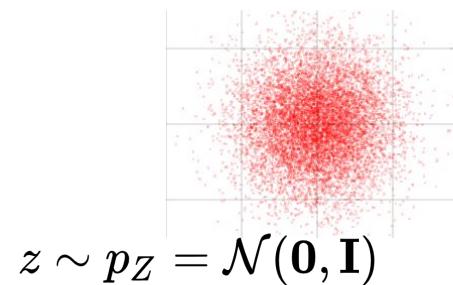
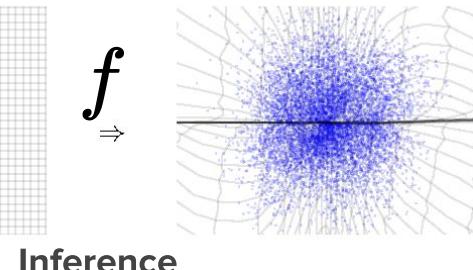
$g = f^{-1}$  bijective

$$\mathbb{E}_x [\log p(x)] = \mathbb{E}_x \left[ \log \mathcal{N}(f(x); \mathbf{0}, I) \left| \det \frac{\partial f(x)}{\partial x} \right| \right]$$

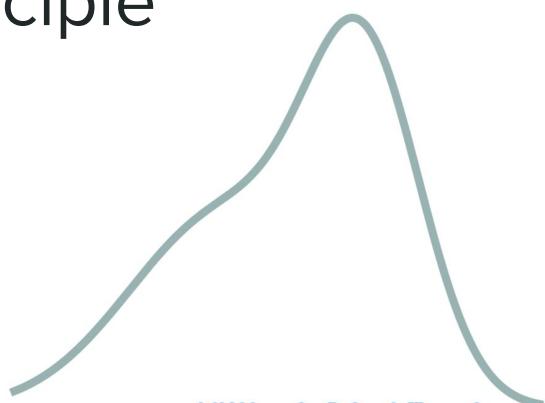
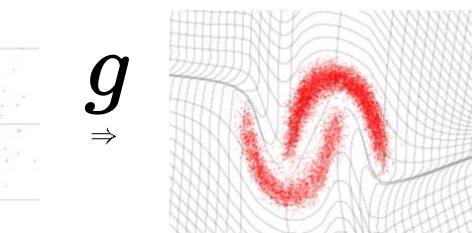
Regularizes the entropy



$$f \Rightarrow$$



$$g \Rightarrow$$



Higher likelihood

# Pathways to Designing a Normalizing Flow

1. Require an invertible architecture.
  - Coupling layers, autoregressive, etc.
2. Require efficient computation of a change of variables equation.

$$\log p(x) = \log p(f(x)) + \log \left| \det \frac{df(x)}{dx} \right|$$

<Model distribution>                          <Base distribution>

(or a continuous version)  $\log p(x(t_N)) = \log p(x(t_0)) + \int_{t_0}^{t_N} \text{tr} \left( \frac{\partial f(x(t), t)}{\partial x(t)} \right) dt$



$\mathcal{O}(m^3)$

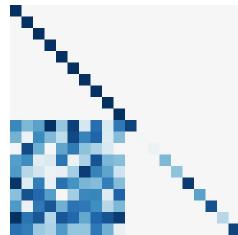
# Architectural Taxonomy

**Sparse connection**

$$f(\mathbf{x})_t = g(\mathbf{x}_{1:t})$$

**1. Block coupling**

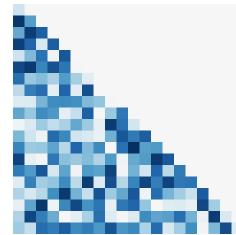
NICE/RealNVP/Glow  
Cubic Spline Flow  
Neural Spline Flow



(Lower triangular +  
structured)

**2. Autoregressive**

IAF/MAF/NAF  
SOS polynomial  
UMNN



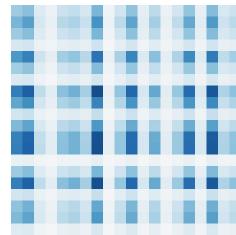
(Lower triangular)

**Residual Connection**

$$f(\mathbf{x}) = \mathbf{x} + g(\mathbf{x})$$

**3. Det identity**

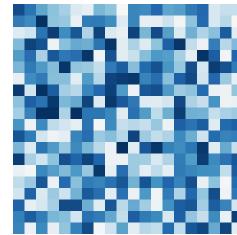
Planar/Sylvester flows  
Radial flow



(Low rank)

**4. Stochastic estimation**

Residual Flow  
FFJORD



(Arbitrary)

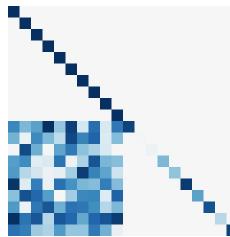
# Architectural Taxonomy

## Sparse connection

$$f(\mathbf{x})_t = g(\mathbf{x}_{1:t})$$

### 1. Block coupling

NICE/RealNVP/Glow  
Cubic Spline Flow  
Neural Spline Flow



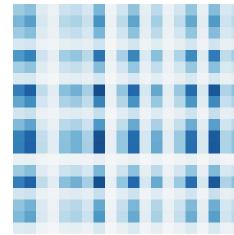
(Lower triangular +  
structured)

## Residual Connection

$$f(\mathbf{x}) = \mathbf{x} + g(\mathbf{x})$$

### 3. Det identity

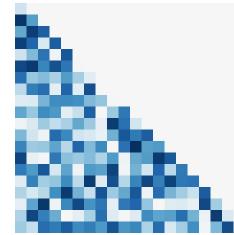
Planar/Sylvester flows  
Radial flow



(Low rank)

### 2. Autoregressive

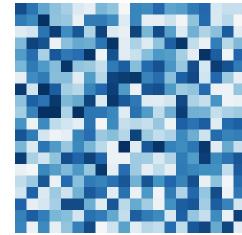
IAF/MAF/NAF  
SOS polynomial  
UMNN



(Lower triangular)

### 4. Stochastic estimation

Residual Flow  
FFJORD



(Arbitrary)

Jacobian

# Coupling Law - NICE

- General form

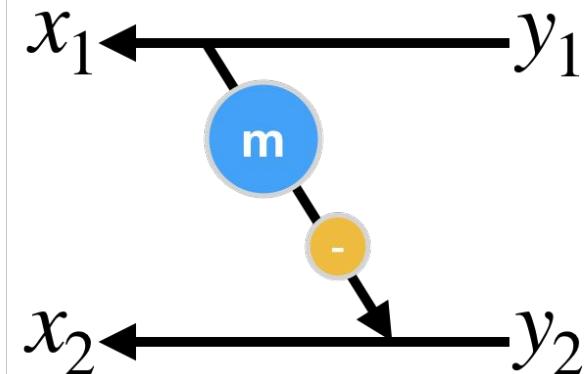
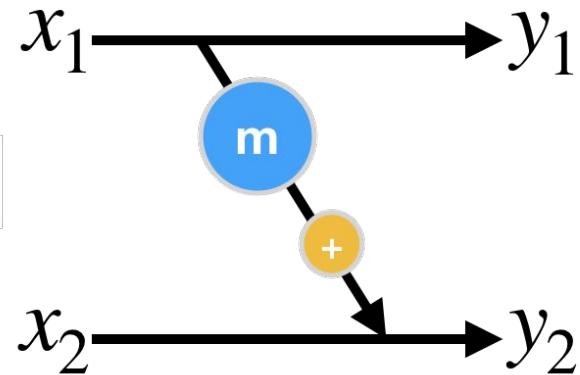
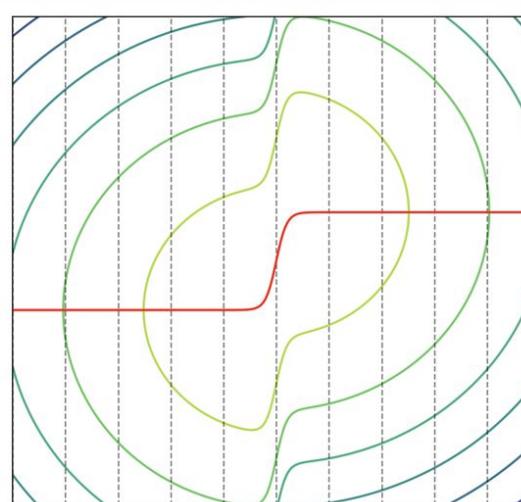
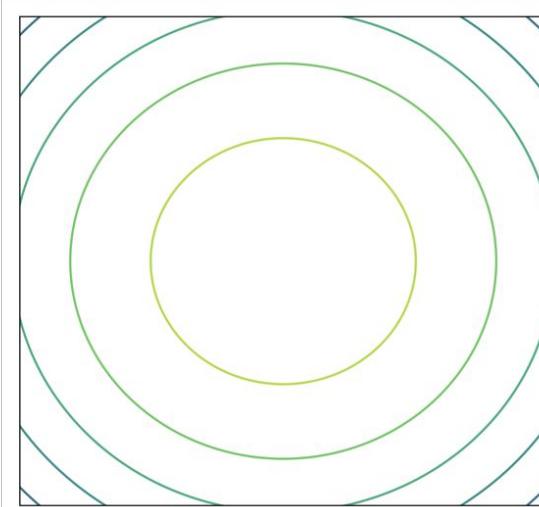
$$f(\mathbf{x})_1 = x_1, \quad f(\mathbf{x})_2 = x_2 + \mathcal{F}(\mathbf{x}_1)$$

- Invertibility

no constraint

- Jacobian determinant

=1 (volume preserving)



# Coupling Law - RealNVP

- General form

$$f(x)_1 = x_1,$$

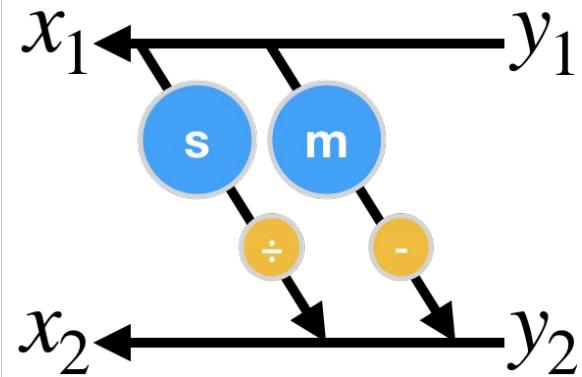
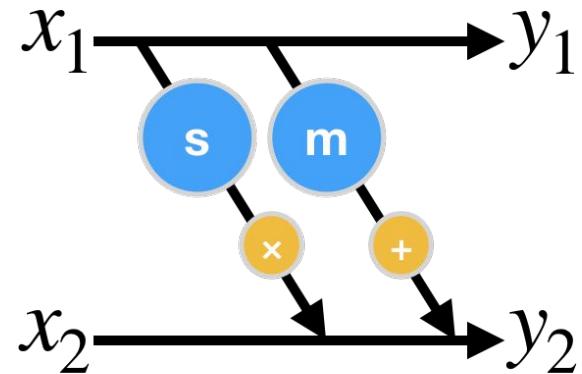
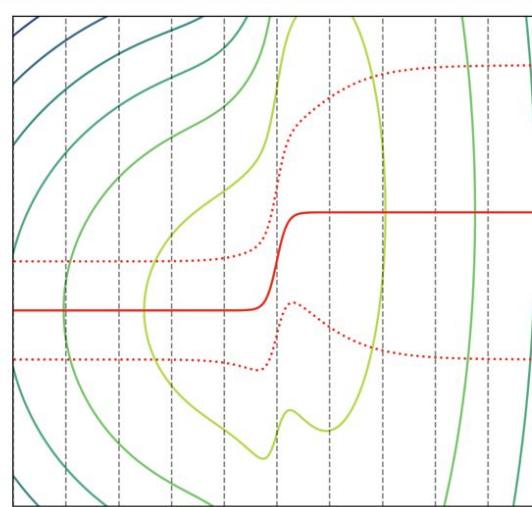
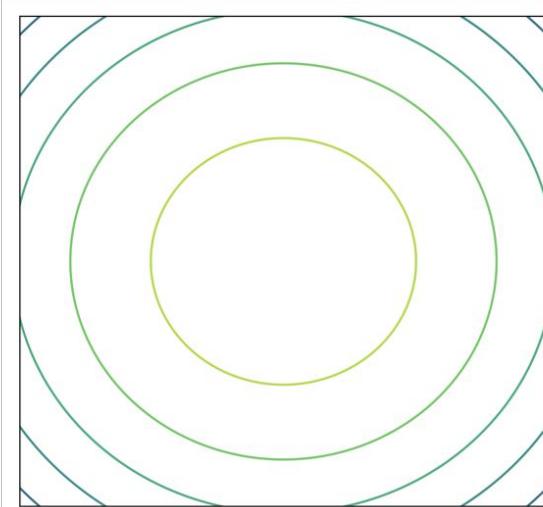
$$f(x)_2 = s(x_1) \odot x_2 + m(x_1)$$

- Invertibility

$s > 0$  (or simply non-zero)

- Jacobian determinant

product of  $s$



# Real NVP via Masked Convolution

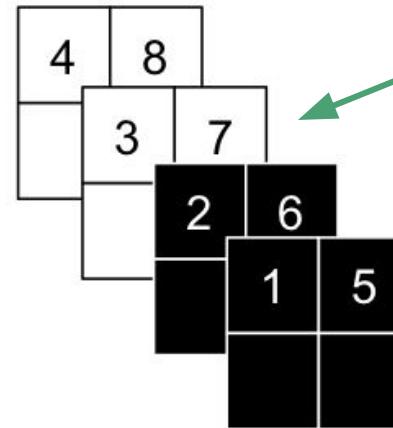
Partitioning can be implemented using a binary mask  $b$ , and using the functional form for  $y$

$$f(x) = b \odot x + (1 - b) \odot (x \odot \exp(s_-(b \odot x))) + m(b \odot x))$$

# Real NVP via Masked Convolution

Partitioning can be implemented using a binary mask  $b$ , and using the functional form for  $y$

$$f(x) = b \odot x + (1 - b) \odot (x \odot \exp(s_-(b \odot x)) + m(b \odot x))$$



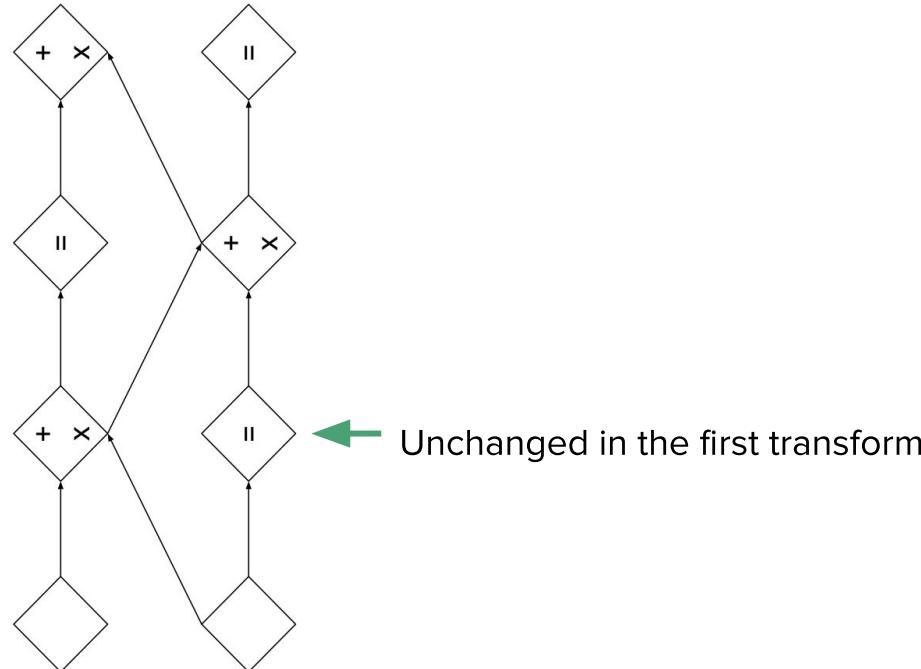


Celeba-64 (left) and LSUN bedroom (right)

Figures from *Density Estimation Using Real NVP* by Dinh et al., 2017

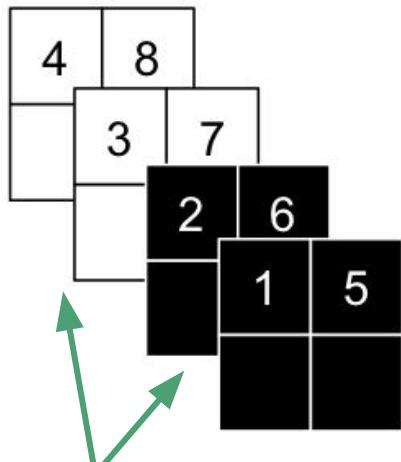
# Glow: Generative Flow with 1x1 Convolutions

Replacing permutation with 1x1 convolution (soft permutation)



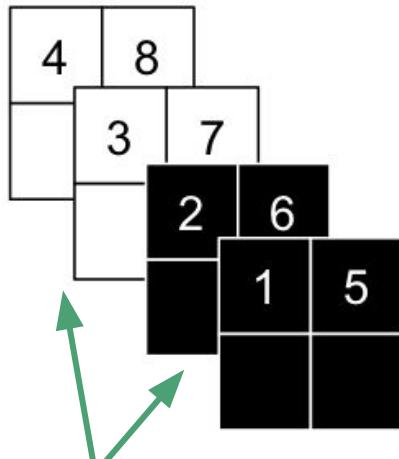
# Glow: Generative Flow with 1x1 Convolutions

Replacing permutation with 1x1 convolution (soft permutation)



# Glow: Generative Flow with 1x1 Convolutions

Replacing permutation with 1x1 convolution (soft permutation)



$$\begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} x_3 \\ x_4 \\ x_1 \\ x_2 \end{bmatrix}$$

Alternating masks

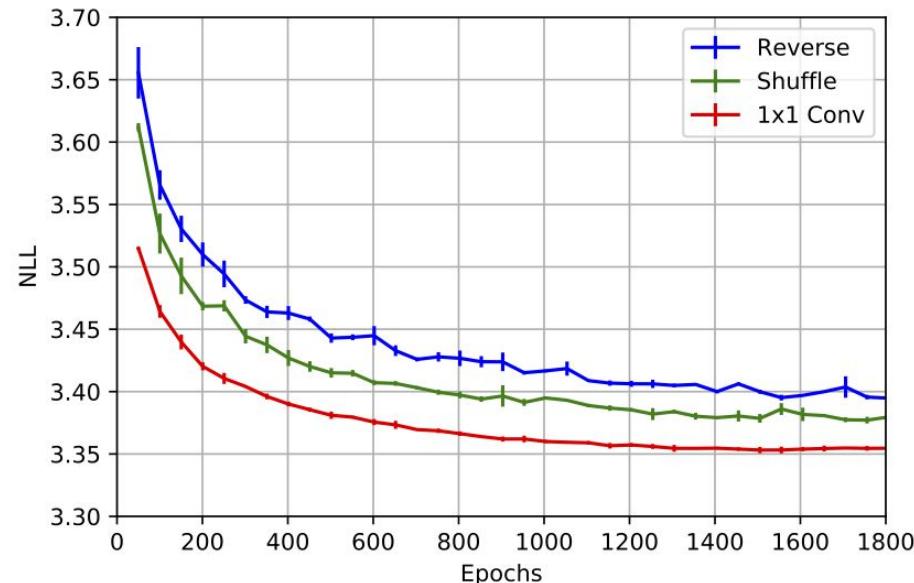
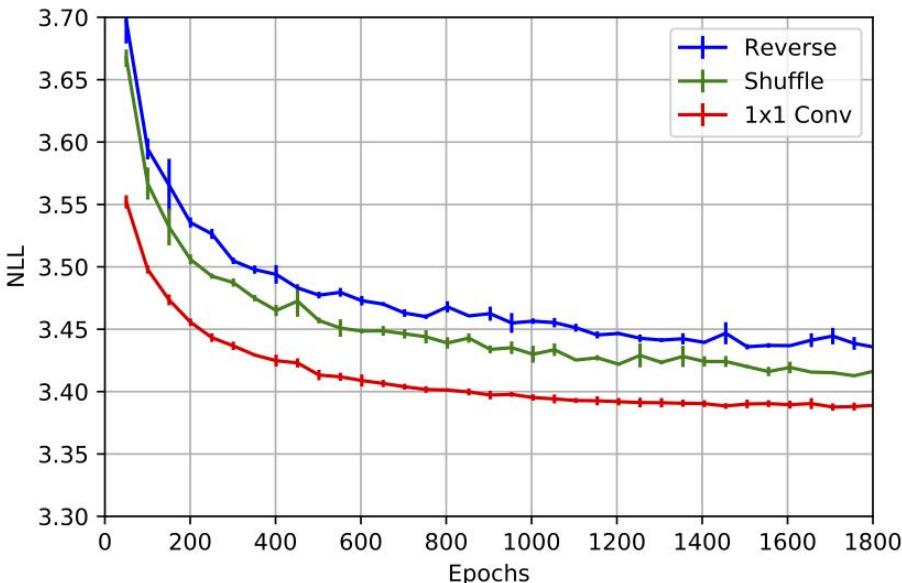
Replace with a general invertible matrix  $W$

Represent  $W$  as a 1x1 convolutional kernel of shape  $[c, c, 1, 1]$ ;  $c$  being # channels

$$\log \left| \det \left( \frac{\partial \text{conv2D}(h; W)}{\partial h} \right) \right| = h \cdot w \cdot \log | \det(W) |$$

# Ablation: Permutation vs 1x1 Convolution

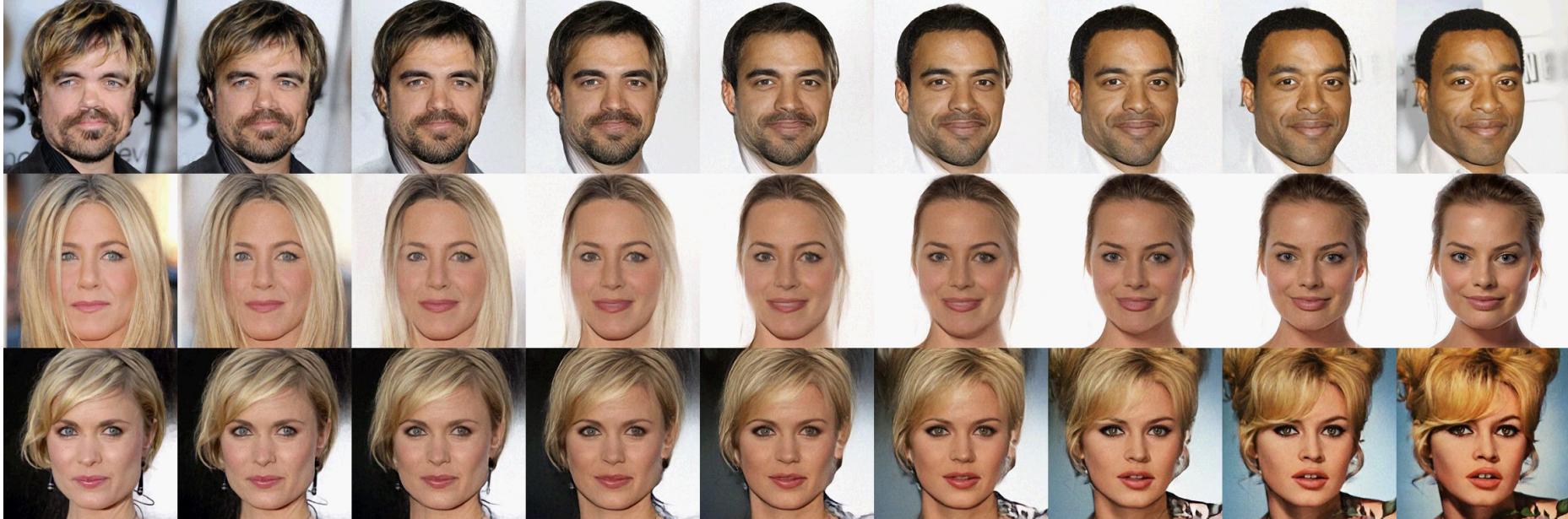
Model	CIFAR-10	ImageNet 32x32	ImageNet 64x64	LSUN (bedroom)	LSUN (tower)	LSUN (church outdoor)
RealNVP	3.49	4.28	3.98	2.72	2.81	3.08
Glow	<b>3.35</b>	<b>4.09</b>	<b>3.81</b>	<b>2.38</b>	<b>2.46</b>	<b>2.67</b>



Bits-per-dim on CIFAR: left: additive, right: affine



Figure from *Glow: Generative Flow with Invertible 1×1 Convolutions* by Kingma and Dhariwal, 2018



# Interpolation with Generative Flows

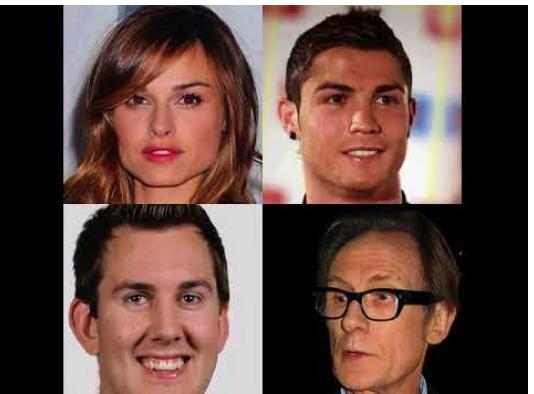


Figure from *Glow: Generative Flow with Invertible 1×1 Convolutions* by Kingma and Dhariwal, 2018  
Video from Durk Kingma's youtube channel

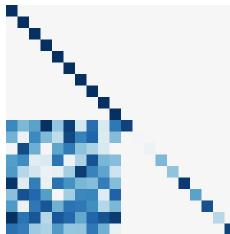
# Architectural Taxonomy

**Sparse connection**

$$f(\mathbf{x})_t = g(\mathbf{x}_{1:t})$$

**1. Block coupling**

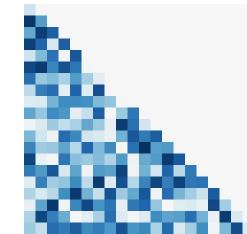
NICE/RealNVP/Glow  
Cubic Spline Flow  
Neural Spline Flow



(Lower triangular +  
structured)

**2. Autoregressive**

IAF/MAF/NAF  
SOS polynomial  
UMNN



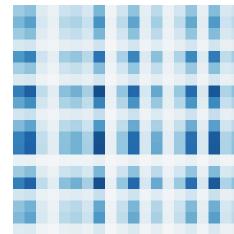
(Lower triangular)

**Residual Connection**

$$f(\mathbf{x}) = \mathbf{x} + g(\mathbf{x})$$

**3. Det identity**

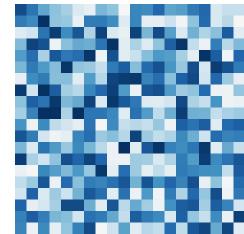
Planar/Sylvester flows  
Radial flow



(Low rank)

**4. Stochastic estimation**

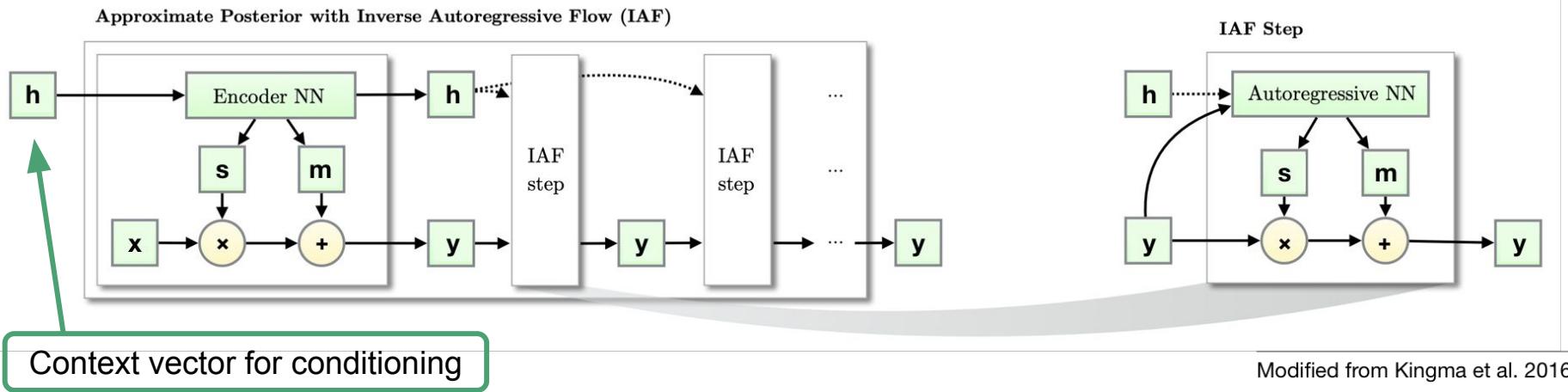
Residual Flow  
FFJORD



(Arbitrary)

# Inverse (Affine) Autoregressive Flows

- General form  $f(\mathbf{x})_t = s(\mathbf{x}_{<t}) \cdot \mathbf{x}_t + m(\mathbf{x}_{<t})$
- Invertibility  $s > 0$  (or simply non-zero)
- Jacobian determinant product of  $s$



Modified from Kingma et al. 2016

# Inverse (Affine) Autoregressive Flows

- General form

$$f(\mathbf{x})_t = s(\mathbf{x}_{<t}) \cdot \mathbf{x}_t + m(\mathbf{x}_{<t})$$

- Invertibility

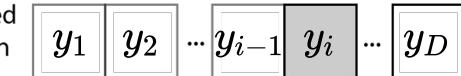
$s > 0$  (or simply non-zero)

- Jacobian determinant

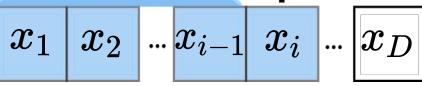
product of  $s$

## Autoregressive NN

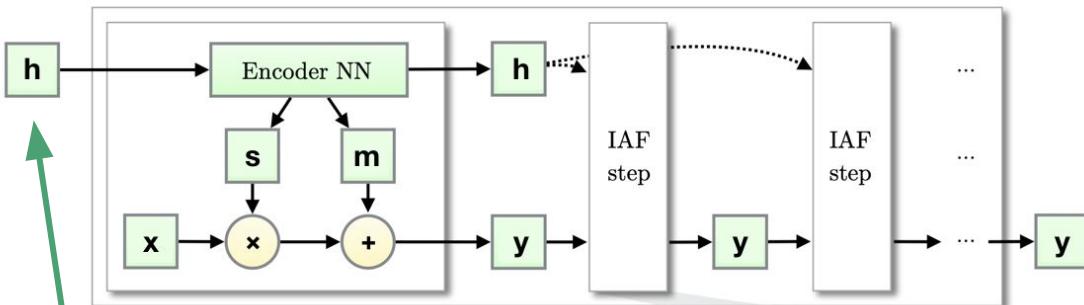
transformed distribution



base distribution

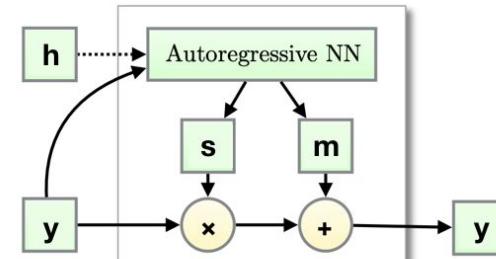


Approximate Posterior with Inverse Autoregressive Flow (IAF)



Context vector for conditioning

IAF Step

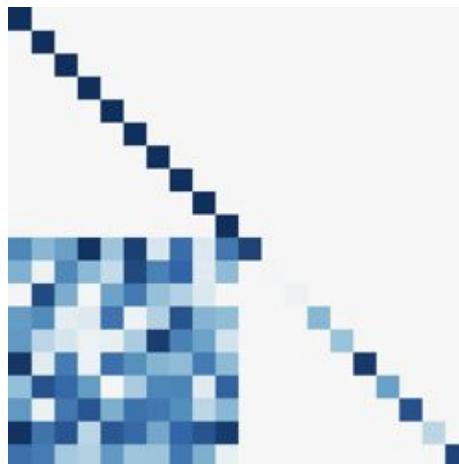


# Trade-off between Expressivity and Inversion Cost

## Block autoregressive

- Limited capacity
- Inverse takes constant time

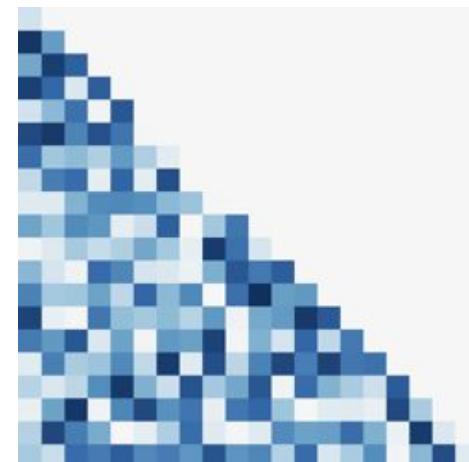
Jacobian



(Block triangular)

## Autoregressive

- Higher capacity
- Inverse takes linear time (dimensionality)



(Triangular)

# Neural Autoregressive Flows

- General form

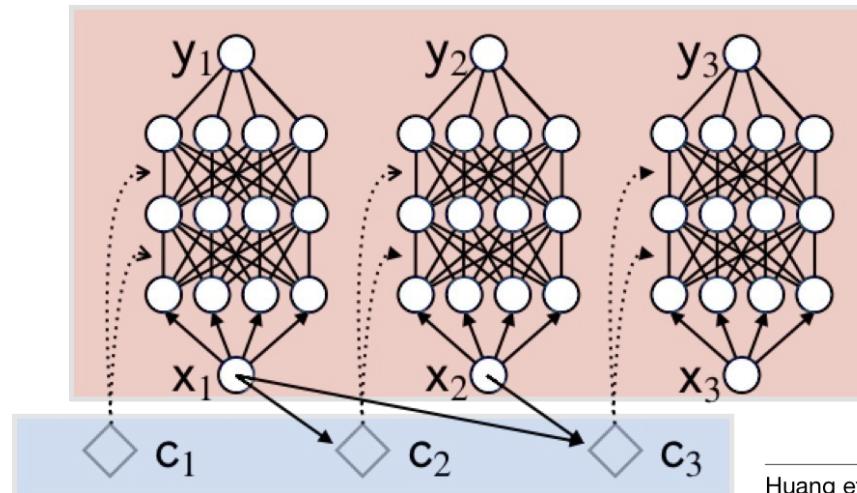
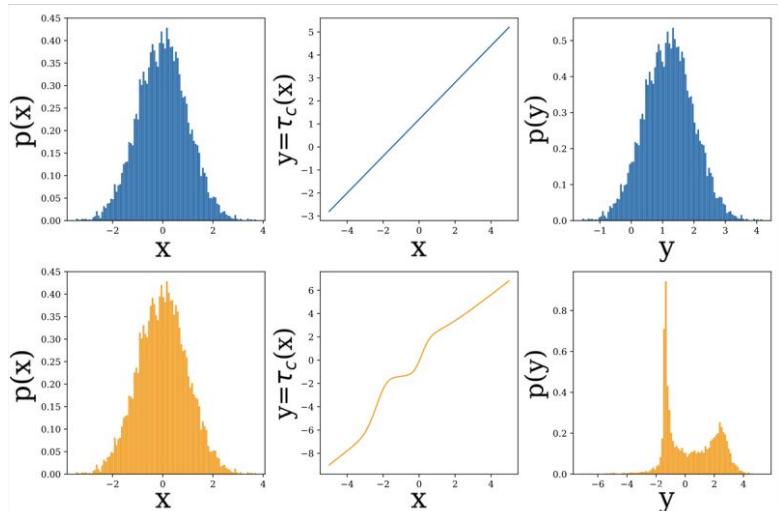
$$f(\mathbf{x})_t = \mathcal{P}(\mathbf{x}_t; \mathcal{H}(\mathbf{x}_{$$

- Invertibility

monotonic activation and positive weight in  $\mathcal{P}$

- Jacobian determinant

product of derivatives (elementwise)



# Inverse Transform Sampling

Let  $F$  be a cumulative distribution function (CDF) on the real space with inverse

$$G(u) := F^{-1}(u) := \inf\{x : F(x) = u\} \quad \forall 0 < u < 1$$

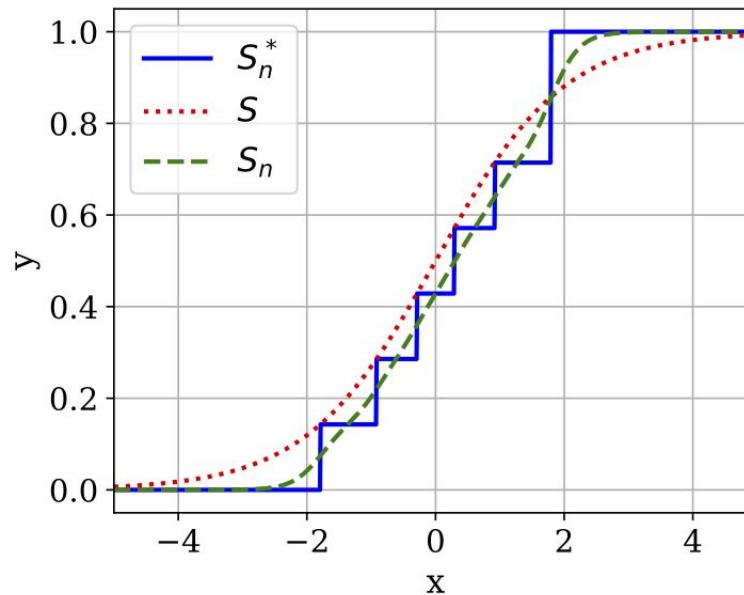
**Theorem (Probability Integral transform & Inverse Probability integral transform)**

If  $F$  is the CDF of  $X$ , then  $Y:=F(X)$  has a uniform distribution on  $[0,1]$ .

Conversely, if  $Y$  follows a uniform distribution on  $[0,1]$ ,  $X:=G(Y)$  has  $F$  as its CDF.

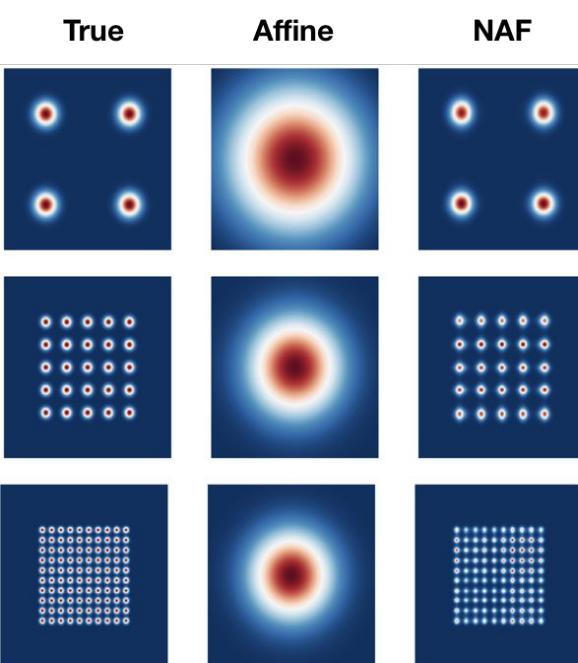
# Universal Flows in 1D

Assume the inverse CDF  $G(u)$  of a random variable  $X$  is continuous, then one can use a monotonic neural network to approximate  $G$  (to an arbitrary precision).

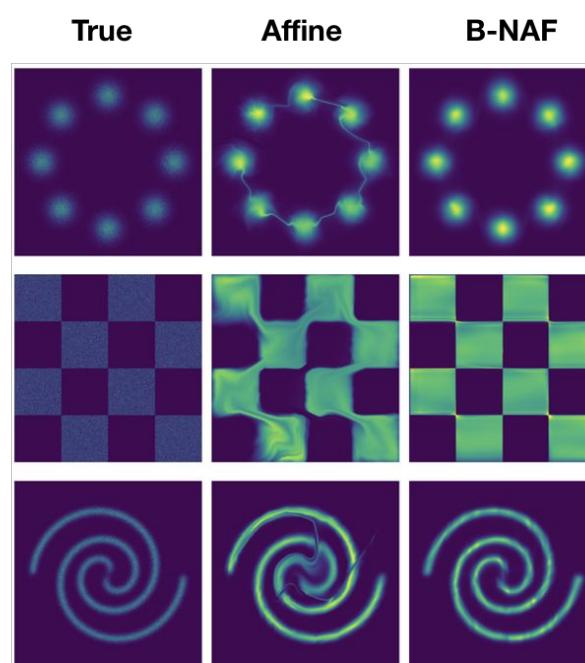


Further reading: the Knothe-Rosenblatt rearrangement for higher dimensionality.

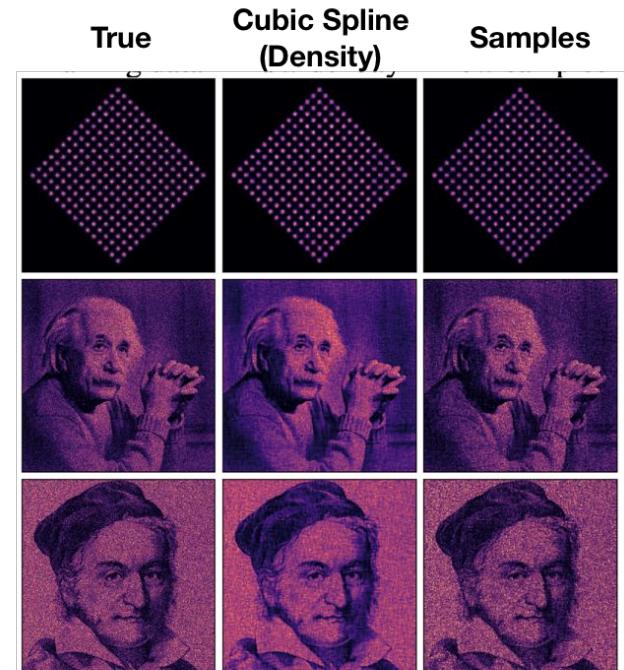
# Examples of Universal Flows



Neural Autoregressive Flows



Block Neural Autoregressive Flows

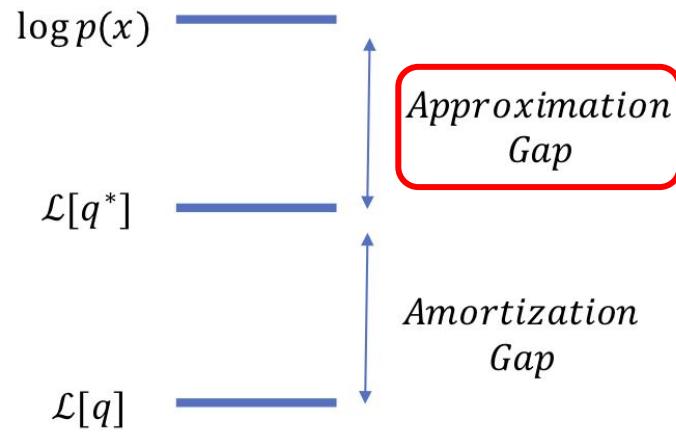
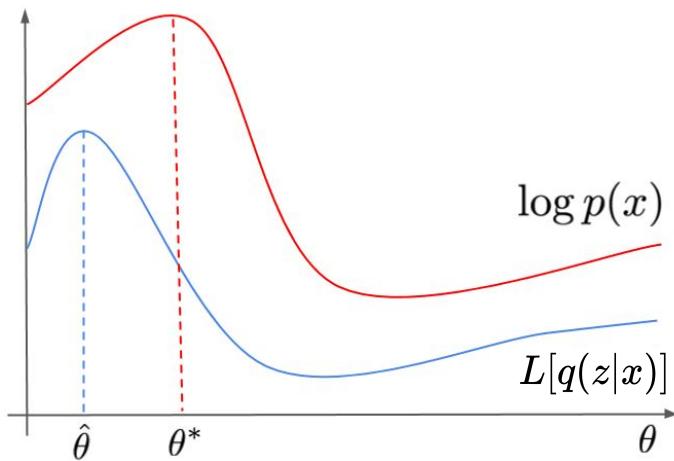


Cubic-Spline Flows

# A Digression on Variational Inference

The variational gap can be reduced by considering a larger family of distributions.

$$\log p(x) - L[q(z|x)] = \underbrace{\log p(x) - L[q^*]}_{\text{Approximation gap}} + \underbrace{L[q^*] - L[q(z|x)]}_{\text{Amortization gap}}$$



# A Digression on Variational Inference

The variational gap can be reduced by considering a larger family of distributions.

Recall the *reparameterization trick*:

Given some function  $g_\phi$  such that for some  $\epsilon \sim q_\epsilon$ ,  $g_\phi(\epsilon) \stackrel{d}{=} z \sim q_\phi(z)$   
 $g_\phi$  depends on  $\phi$  but  $q_\epsilon$  does not.

$$\nabla_\phi \mathbb{E}_{z \sim q_\phi(z)} [\log p_\theta(x, z) - \log q_\phi(z)] = \nabla_\phi \mathbb{E}_{\epsilon \sim q_\epsilon(\epsilon)} [\log p_\theta(x, g_\phi(\epsilon)) - \log q_\phi(g_\phi(\epsilon))]$$

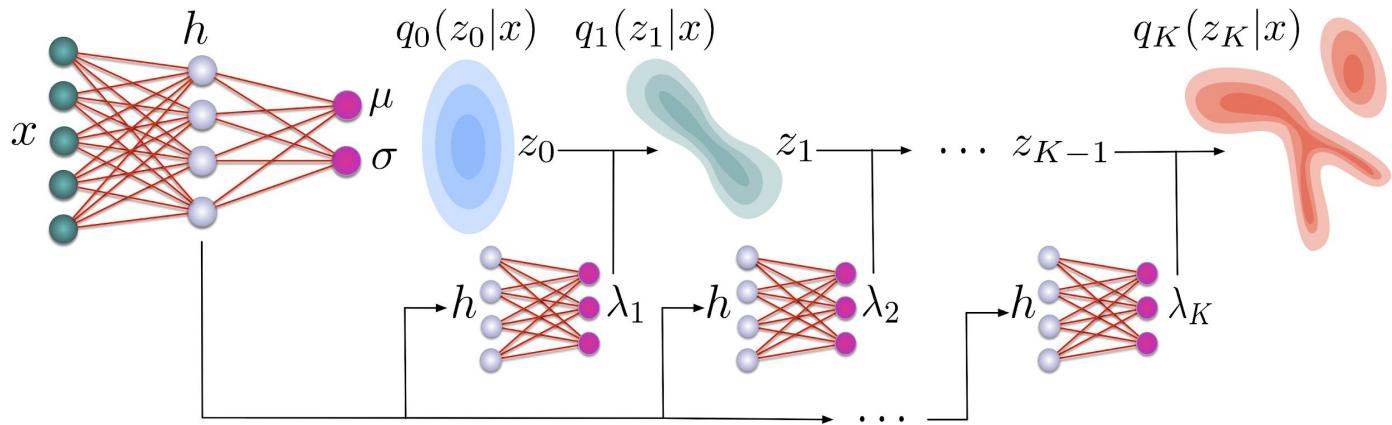
If furthermore,  $g_\phi$  is invertible, the entropy of  $z$  can be estimated by

$$-\mathbb{E}_\epsilon [\log q_\phi(g_\phi(\epsilon))] = -\mathbb{E}_\epsilon \left[ \log q_\epsilon(\epsilon) \left| \det \frac{\partial g_\phi(\epsilon)}{\partial \epsilon} \right|^{-1} \right]$$

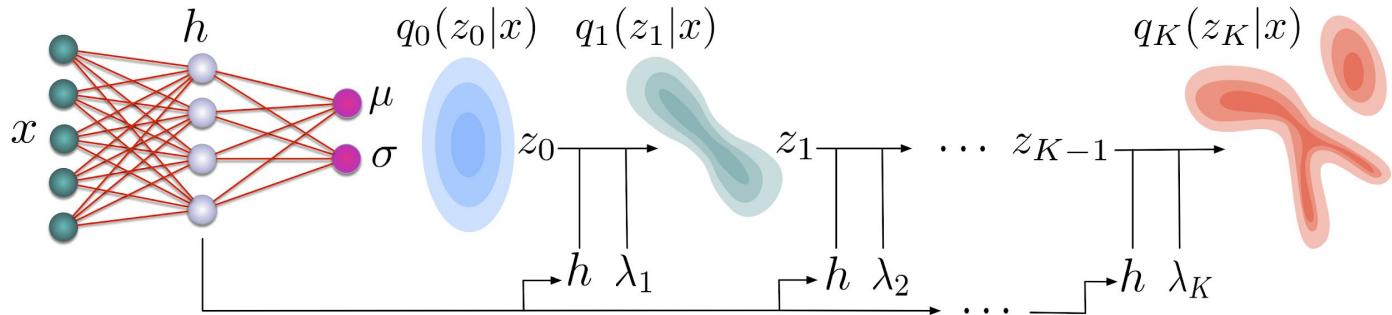
Prevents the distribution of  $z$  from collapsing to a point mass

# Amortization Strategies

Hypernetwork



Concatenation



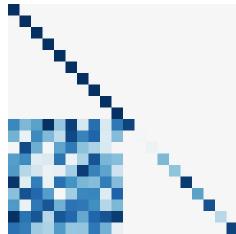
# Architectural Taxonomy

**Sparse connection**

$$f(\mathbf{x})_t = g(\mathbf{x}_{1:t})$$

**1. Block coupling**

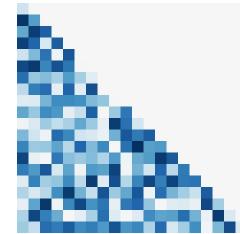
NICE/RealNVP/Glow  
Cubic Spline Flow  
Neural Spline Flow



(Lower triangular +  
structured)

**2. Autoregressive**

IAF/MAF/NAF  
SOS polynomial  
UMNN



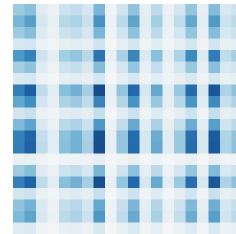
(Lower triangular)

**Residual Connection**

$$f(\mathbf{x}) = \mathbf{x} + g(\mathbf{x})$$

**3. Det identity**

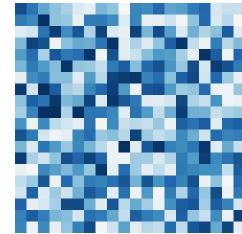
Planar/Sylvester flows  
Radial flow



(Low rank)

**4. Stochastic estimation**

Residual Flow  
FFJORD



(Arbitrary)

# Determinant Identity - Planar Flows

- General form

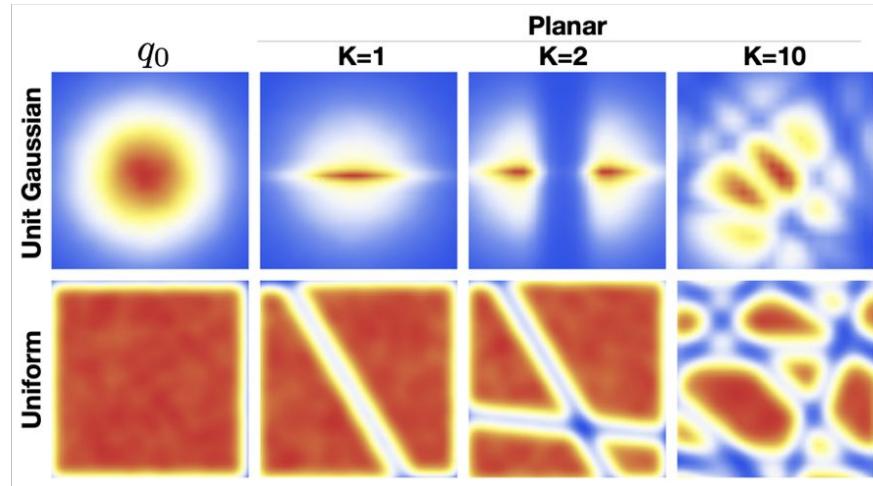
$$f(\mathbf{x}) = \mathbf{x} + \mathbf{u}h(\mathbf{w}^\top \mathbf{x} + b)$$

- Invertibility

$$\mathbf{u}^\top \mathbf{w} > -1 \text{ if } h = \tanh$$

- Jacobian determinant

$$\left| \det \frac{\partial f}{\partial \mathbf{x}} \right| = \left| \det \left( \mathbf{I} + h'(\mathbf{w}^\top \mathbf{x} + b) \mathbf{u} \mathbf{w}^\top \right) \right| = \left| 1 + h'(\mathbf{w}^\top \mathbf{x} + b) \mathbf{u}^\top \mathbf{w} \right|$$



## VAE on binary MNIST

Model	$-\ln p(\mathbf{x})$
DLGM diagonal covariance	$\leq 89.9$
DLGM+NF (k = 10)	$\leq 87.5$
DLGM+NF (k = 20)	$\leq 86.5$
DLGM+NF (k = 40)	$\leq 85.7$
DLGM+NF (k = 80)	$\leq 85.1$

Rezende et al. 2015

# Determinant Identity - Sylvester Flows

- General form

$$f(\mathbf{x}) = \mathbf{x} + \mathbf{A}h(\mathbf{B}\mathbf{x} + \mathbf{b})$$

$\mathbf{A} \in \mathbb{R}^{m \times d}$ ,  $\mathbf{B} \in \mathbb{R}^{d \times m}$ ,  $\mathbf{b} \in \mathbb{R}^d$ , and  $d \leq m$

- Invertibility Similar to planer flows

- Jacobian determinant Using Sylvester's Thm:  $\det(\mathbf{I}_m + \mathbf{AB}) = \det(\mathbf{I}_d + \mathbf{BA})$

Model	Freyfaces		Omniglot		Caltech 101	
	-ELBO	NLL	-ELBO	NLL	-ELBO	NLL
VAE	$4.53 \pm 0.02$	$4.40 \pm 0.03$	$104.28 \pm 0.39$	$97.25 \pm 0.23$	$110.80 \pm 0.46$	$99.62 \pm 0.74$
Planar	<b><math>4.40 \pm 0.06</math></b>	<b><math>4.31 \pm 0.06</math></b>	$102.65 \pm 0.42$	$96.04 \pm 0.28$	$109.66 \pm 0.42$	$98.53 \pm 0.68$
IAF	$4.47 \pm 0.05$	$4.38 \pm 0.04$	$102.41 \pm 0.04$	$96.08 \pm 0.16$	$111.58 \pm 0.38$	$99.92 \pm 0.30$
O-SNF	$4.51 \pm 0.04$	$4.39 \pm 0.05$	$99.00 \pm 0.29$	$93.82 \pm 0.21$	$106.08 \pm 0.39$	$94.61 \pm 0.83$
H-SNF	$4.46 \pm 0.05$	$4.35 \pm 0.05$	<b><math>99.00 \pm 0.04</math></b>	<b><math>93.77 \pm 0.03</math></b>	<b><math>104.62 \pm 0.29</math></b>	<b><math>93.82 \pm 0.62</math></b>
T-SNF	$4.45 \pm 0.04$	$4.35 \pm 0.04$	$99.33 \pm 0.23$	$93.97 \pm 0.13$	$105.29 \pm 0.64$	$94.92 \pm 0.73$

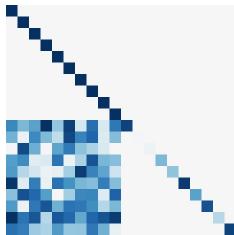
# Architectural Taxonomy

## Sparse connection

$$f(\mathbf{x})_t = g(\mathbf{x}_{1:t})$$

### 1. Block coupling

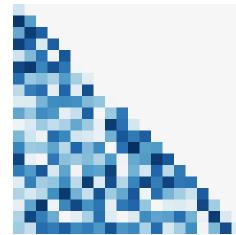
NICE/RealNVP/Glow  
Cubic Spline Flow  
Neural Spline Flow



(Lower triangular +  
structured)

### 2. Autoregressive

IAF/MAF/NAF  
SOS polynomial  
UMNN



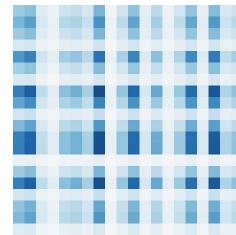
(Lower triangular)

## Residual Connection

$$f(\mathbf{x}) = \mathbf{x} + g(\mathbf{x})$$

### 3. Det identity

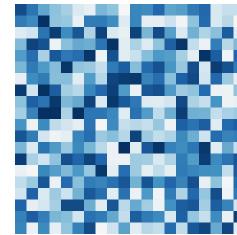
Planar/Sylvester flows  
Radial flow



(Low rank)

### 4. Stochastic estimation

Residual Flow  
FFJORD



(Arbitrary)

# Stochastic Estimation for General Residual Form

- General form
- Invertibility
- Jacobian determinant

$$f(\mathbf{x}) = \mathbf{x} + g(\mathbf{x})$$

$$\left\| \frac{\partial g(\mathbf{x})}{\partial \mathbf{x}} \right\|_2 < 1$$

---

*Algorithm 1.* Inverse of i-ResNet layer via fixed-point iteration.

**Input:** output from residual layer  $y$ , contractive residual block  $g$ , number of fixed-point iterations  $n$   
Init:  $x^0 := y$   
**for**  $i = 0, \dots, n$  **do**  
     $x^{i+1} := y - g(x^i)$   
**end for**

---

**Theorem 1.** If  $g(x)$  is  $L$ -Lipschitz for  $0 < L < 1$ , then  $f(x) := x + g(x)$  is invertible.

**Proof** Fix any  $y$  in the codomain of  $f$ . Define  $x_0 = y$  and  $x_t = y - g(x_{t-1})$  for  $t \geq 1$ , which is a contractive mapping. Thus, by the Banach fixed-point theorem,  $x_t$  converges to a unique fixed point  $x^*$  which satisfies  $f(x^*) = y$ .

# Stochastic Estimation for General Residual Form

- General form

$$f(\mathbf{x}) = \mathbf{x} + g(\mathbf{x})$$

- Invertibility

$$\left\| \frac{\partial g(\mathbf{x})}{\partial \mathbf{x}} \right\|_2 < 1$$

- Jacobian determinant

$$\log \left| \det \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \right| = \text{tr} \left( \log \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \right)$$

Jacobi's formula

# Stochastic Estimation for General Residual Form

- General form

$$f(\mathbf{x}) = \mathbf{x} + g(\mathbf{x})$$

- Invertibility

$$\left\| \frac{\partial g(\mathbf{x})}{\partial \mathbf{x}} \right\|_2 < 1$$

- Jacobian determinant

$$\log \left| \det \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \right| = \text{tr} \left( \log \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \right)$$

Jacobi's formula

$$\text{tr} \left( \log \left( \mathbf{I} + \frac{\partial g(\mathbf{x})}{\partial \mathbf{x}} \right) \right) = \sum_{k=1}^{\infty} \frac{(-1)^{k+1}}{k} \text{tr} \left( \frac{\partial g(\mathbf{x})}{\partial \mathbf{x}}^k \right)$$

Power series expansion

# Stochastic Estimation for General Residual Form

- General form

$$f(\mathbf{x}) = \mathbf{x} + g(\mathbf{x})$$

- Invertibility

$$\left\| \frac{\partial g(\mathbf{x})}{\partial \mathbf{x}} \right\|_2 < 1$$

- Jacobian determinant

$$\log \left| \det \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \right| = \text{tr} \left( \log \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \right)$$

Jacobi's formula

$$\text{tr} \left( \log \left( \mathbf{I} + \frac{\partial g(\mathbf{x})}{\partial \mathbf{x}} \right) \right) = \sum_{k=1}^{\infty} \frac{(-1)^{k+1}}{k} \text{tr} \left( \frac{\partial g(\mathbf{x})}{\partial \mathbf{x}}^k \right)$$

Power series expansion

$$\approx \mathbb{E}_v \left[ \sum_{k=1}^{\textcolor{red}{n}} \frac{(-1)^{k+1}}{k} v^\top \left( \frac{\partial g(\mathbf{x})}{\partial \mathbf{x}}^k \right) v \right]$$

Truncation &  
Hutchinson trace estimator

# Stochastic Estimation for General Residual Form

- General form

$$f(\mathbf{x}) = \mathbf{x} + g(\mathbf{x})$$

- Invertibility

$$\left\| \frac{\partial g(\mathbf{x})}{\partial \mathbf{x}} \right\|_2 < 1$$

- Jacobian determinant

$$\log \left| \det \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \right| = \text{tr} \left( \log \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \right)$$

Jacobi's formula

$$\text{tr} \left( \log \left( \mathbf{I} + \frac{\partial g(\mathbf{x})}{\partial \mathbf{x}} \right) \right) = \sum_{k=1}^{\infty} \frac{(-1)^{k+1}}{k} \text{tr} \left( \frac{\partial g(\mathbf{x})}{\partial \mathbf{x}}^k \right)$$

Power series expansion

$$\approx \mathbb{E}_v \left[ \sum_{k=1}^n \frac{(-1)^{k+1}}{k} v^\top \left( \frac{\partial g(\mathbf{x})}{\partial \mathbf{x}}^k \right) v \right]$$

Truncation &  
Hutchinson trace estimator

Bias

# Stochastic Estimation for General Residual Form

- General form

$$f(\mathbf{x}) = \mathbf{x} + g(\mathbf{x})$$

- Invertibility

$$\left\| \frac{\partial g(\mathbf{x})}{\partial \mathbf{x}} \right\|_2 < 1$$

- Jacobian determinant

$$\log \left| \det \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \right| = \text{tr} \left( \log \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \right)$$

Jacobi's formula

$$\text{tr} \left( \log \left( \mathbf{I} + \frac{\partial g(\mathbf{x})}{\partial \mathbf{x}} \right) \right) = \sum_{k=1}^{\infty} \frac{(-1)^{k+1}}{k} \text{tr} \left( \frac{\partial g(\mathbf{x})}{\partial \mathbf{x}}^k \right)$$

Power series expansion

$$= \mathbb{E}_{v,n} \left[ \sum_{k=1}^n \frac{(-1)^{k+1}}{k \cdot \mathbb{P}(N \geq k)} v^\top \left( \frac{\partial g(\mathbf{x})}{\partial \mathbf{x}}^k \right) v \right]$$

Russian roulette estimator &  
Hutchinson trace estimator

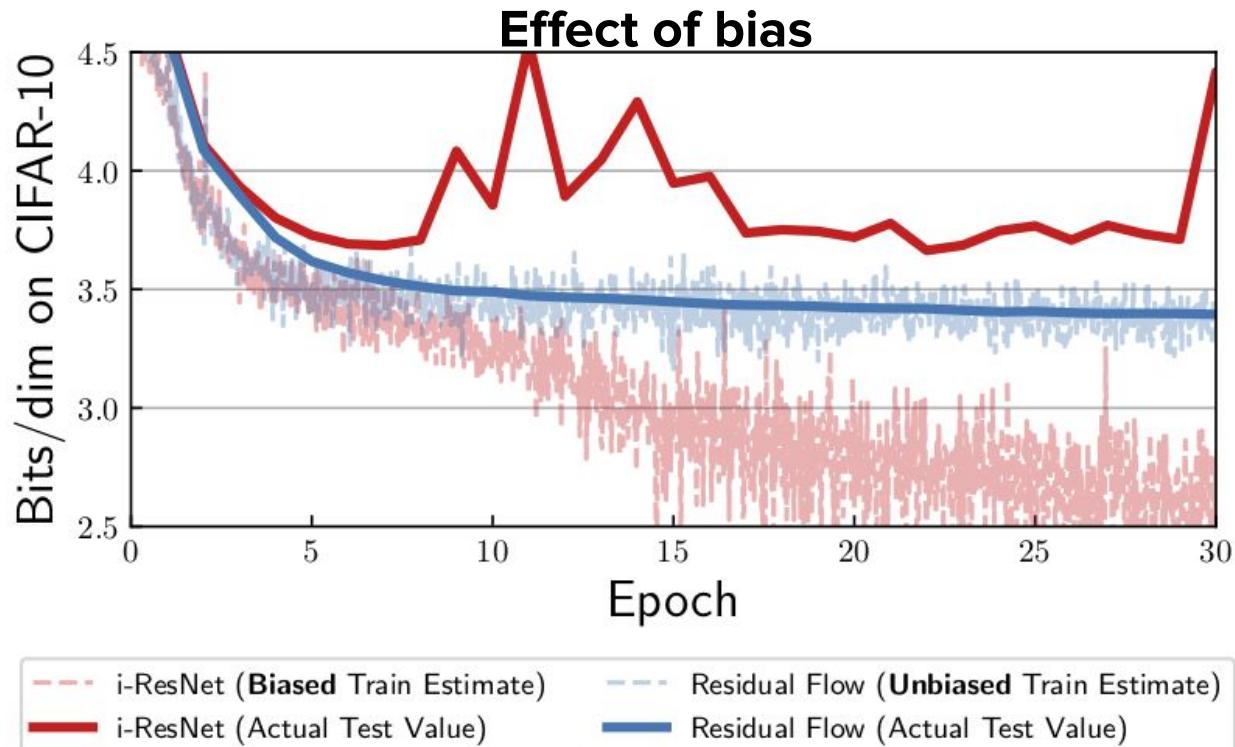
CelebA samples



Cifar10 samples



Imagenet-32 samples



# Continuous Time Flows

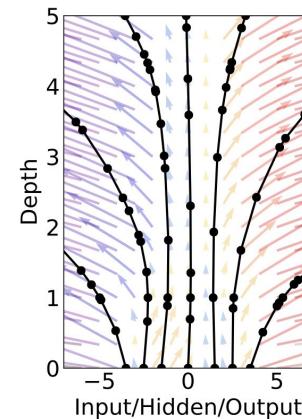
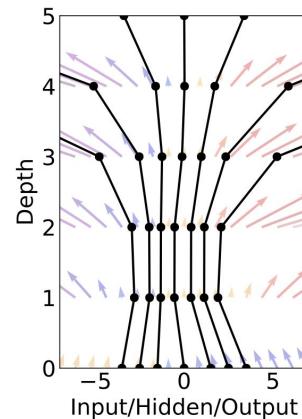
---

# Neural ODE

A neural ordinary differential equation (neural ODE) parameterizes the **time derivative** of a function using a neural network

$$\frac{dz(t)}{dt} = f(z(t), t, \theta) \quad z(t_1) = z(t_0) + \int_{t_0}^{t_1} f(z(t), t, \theta) dt$$

ResNet



ODE Network

# Backpropagation of Neural ODE

**Adjoint sensitivity method** enables reverse-mode differentiation (backprop)

$$a(t) := \frac{\partial L(z(t_1))}{\partial z(t)}$$

Adjoint variable

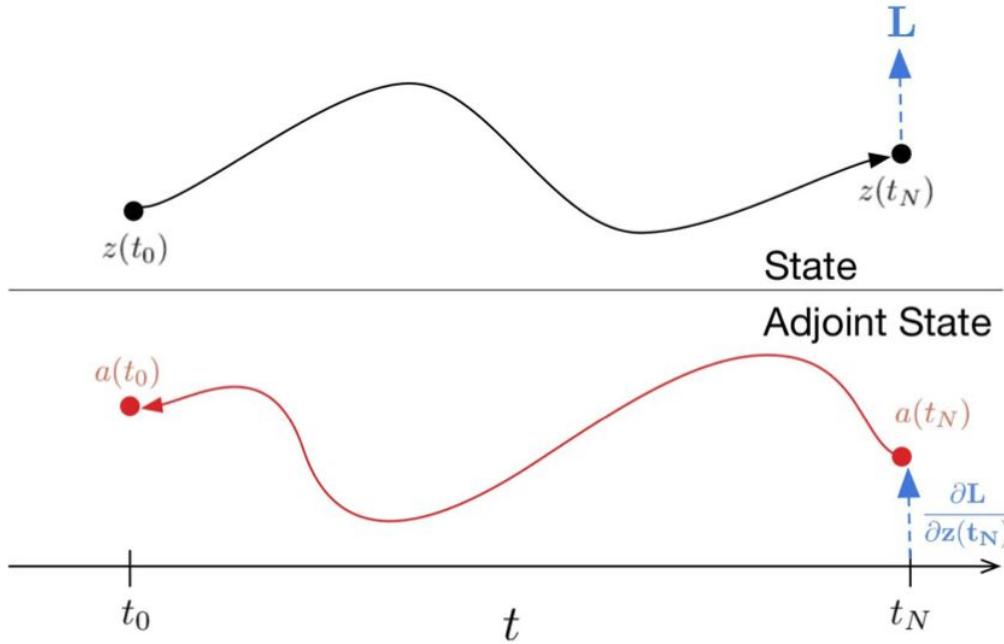
$$\frac{da(t)}{dt} = -a(t)^\top \frac{\partial f(z(t), t, \theta)}{\partial z}$$

Solve for  $a$  in reverse time

Recompute  $z$  backwards  
Constant memory!

# A Differentiable Primitive for AutoDiff

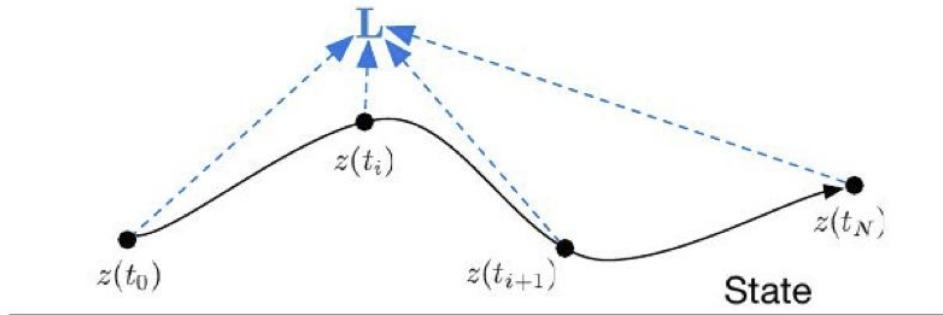
Forward:



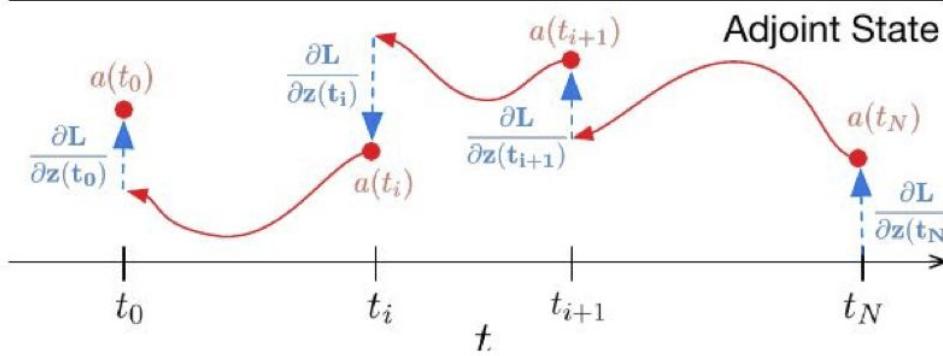
Backward:

# A Differentiable Primitive for AutoDiff

Forward:



Backward:



# Backpropagation of Neural ODE (cont'd)

**Adjoint sensitivity method** enables reverse-mode differentiation (backprop)

$$a(t) := \frac{\partial L(z(t_1))}{\partial z(t)}$$

Adjoint variable

Solve for  $a$  in reverse time

$$\frac{da(t)}{dt} = -a(t)^\top \frac{\partial f(z(t), t, \theta)}{\partial z}$$

1

Recompute  $z$  backwards  
Constant memory!

Gradient wrt parameters

$$\frac{dL}{d\theta} = - \int_{t_1}^{t_0} a(t)^\top \frac{\partial f(z(t), t, \theta)}{\partial \theta} dt$$

3

Solve 1 ODE forward  
Solve 3 ODEs backward

# Uniqueness of Solution $\Rightarrow$ Invertibility

## Theorem (Picard-Lindelof)

The initial value problem  $\frac{dz(t)}{dt} = f(z(t), t, \theta) \quad z(0) = z_0$

has a unique solution if  $f$  is ***uniformly Lipschitz in time***.

Reverse time always has a unique solution  $\Rightarrow$  invertible!

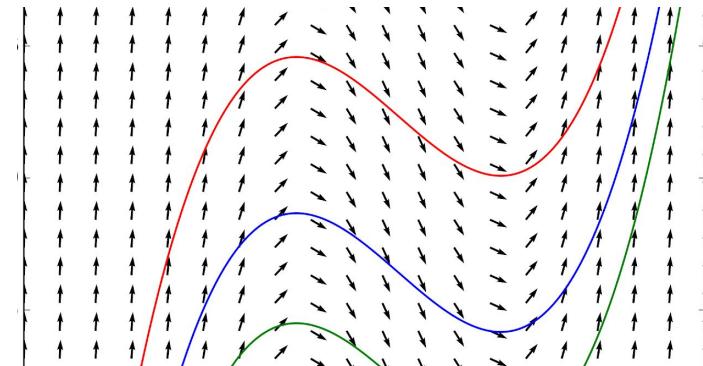


Figure from blog post: *Neural Ordinary Differential Equations and Dynamics Models* by Aidan Abdulali, 2019

# Instantaneous Change of Variables

**Theorem (Instantaneous Change of Variables, Chen et al., 2018)**

Assume  $f$  is ***uniformly Lipschitz in time***, and  $z(t)$  follows the probability density function  $p(z(t), t)$ . Then the change in log probability function follows the ODE

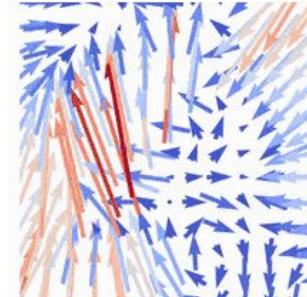
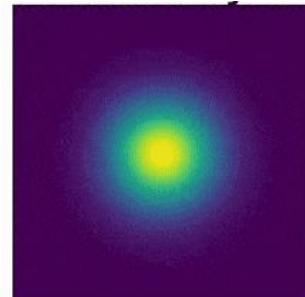
$$\frac{\partial \log p(z(t), t)}{\partial t} = -\text{tr} \left( \frac{df}{dz(t)} \right)$$

# Instantaneous Change of Variables

**Theorem (Instantaneous Change of Variables, Chen et al., 2018)**

Assume  $f$  is ***uniformly Lipschitz in time***, and  $z(t)$  follows the probability density function  $p(z(t), t)$ . Then the change in log probability function follows the ODE

$$\frac{\partial \log p(z(t), t)}{\partial t} = -\text{tr} \left( \frac{df}{dz(t)} \right)$$



1 4 0 1 6 8 1 0 0 5  
0 3 2 7 1 8 8 3 5 9  
9 5 2 9 1 5 0 4 5 6  
1 3 6 2 0 4 1 7 3 2  
9 9 3 5 4 2 4 7 7 2  
4 3 8 4 0 7 2 8 7 9  
4 6 2 8 3 8 9 5 1 8  
8 3 5 6 1 7 1 9 5 4  
1 6 9 8 1 6 6 3 7 5  
0 9 6 9 8 8 9 8 9 9



# Residual Flows as Finite Neural ODE

If  $f$  is Lipschitz, then the Euler method for numerical integration is invertible given that the step size is sufficiently small

$$z_{t+\epsilon} = z_t + \epsilon \cdot f(z_t, t, \theta)$$

