

# Density Estimation Using Real NVP

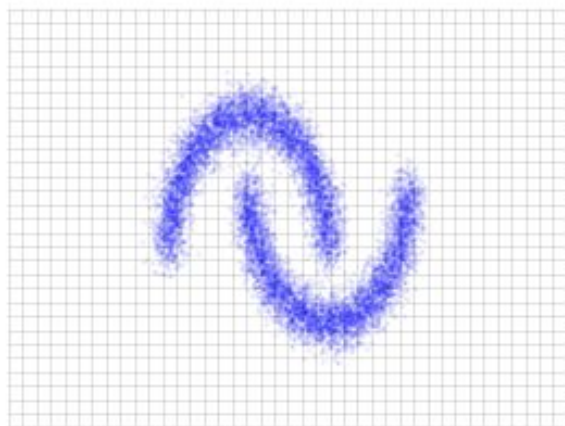
Presenter: Zhengyuan Cui

**Inference**

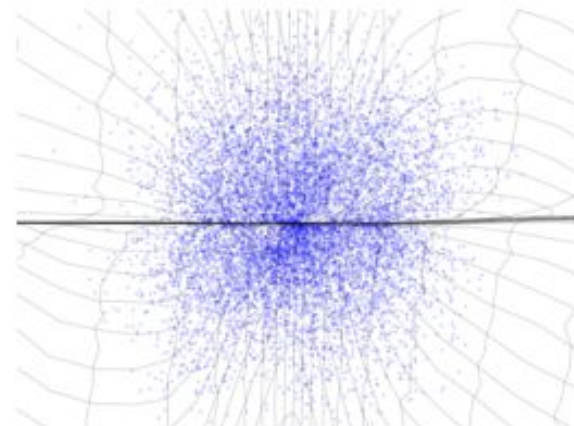
$$x \sim \hat{p}_X$$

$$z = f(x)$$

Data space  $\mathcal{X}$



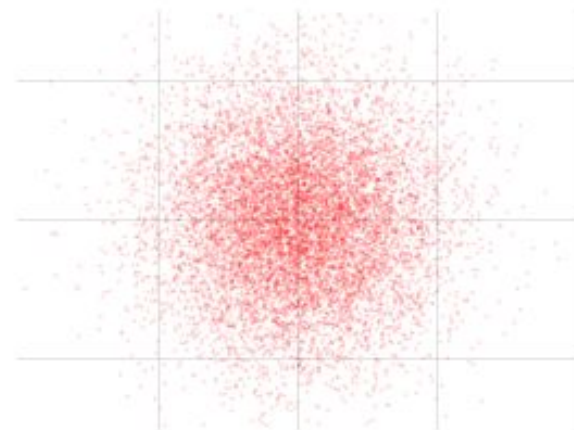
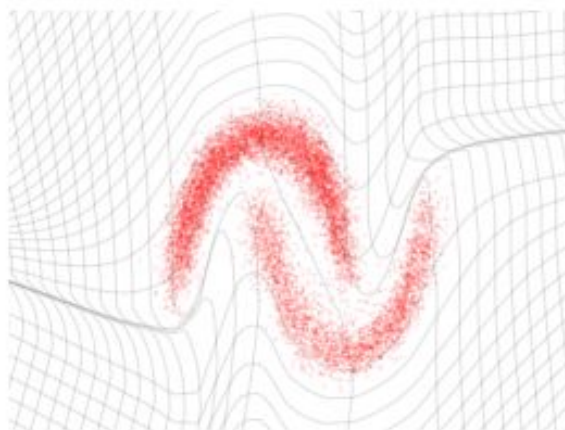
Latent space  $\mathcal{Z}$



**Generation**

$$z \sim p_Z$$

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



$$p_X(x) = p_Z(f(x)) \left| \det \left( \frac{\partial f(x)}{\partial x^T} \right) \right|$$

$$\log(p_X(x)) = \log(p_Z(f(x))) + \log \left( \left| \det \left( \frac{\partial f(x)}{\partial x^T} \right) \right| \right),$$

# Coupling layer

$$y_{1:d} = x_{1:d}$$

$$y_{d+1:D} = x_{d+1:D} \odot \exp(s(x_{1:d})) + t(x_{1:d}),$$

$$\frac{\partial y}{\partial x^T} = \begin{bmatrix} \mathbb{I}_d & 0 \\ \frac{\partial y_{d+1:D}}{\partial x_{1:d}^T} & \text{diag}(\exp[s(x_{1:d})]) \end{bmatrix}$$

# Coupling Layer

$$\begin{cases} y_{1:d} &= x_{1:d} \\ y_{d+1:D} &= x_{d+1:D} \odot \exp(s(x_{1:d})) + t(x_{1:d}) \end{cases}$$

$$\Leftrightarrow \begin{cases} x_{1:d} &= y_{1:d} \\ x_{d+1:D} &= (y_{d+1:D} - t(y_{1:d})) \odot \exp(-s(y_{1:d})), \end{cases}$$

---

## Another expression

$$y = b \odot x + (1 - b) \odot \left( x \odot \exp \left( s(b \odot x) \right) + t(b \odot x) \right).$$

1	2	5	6
3	4	7	8

4	8			
	3	7		
		2	6	
			1	5

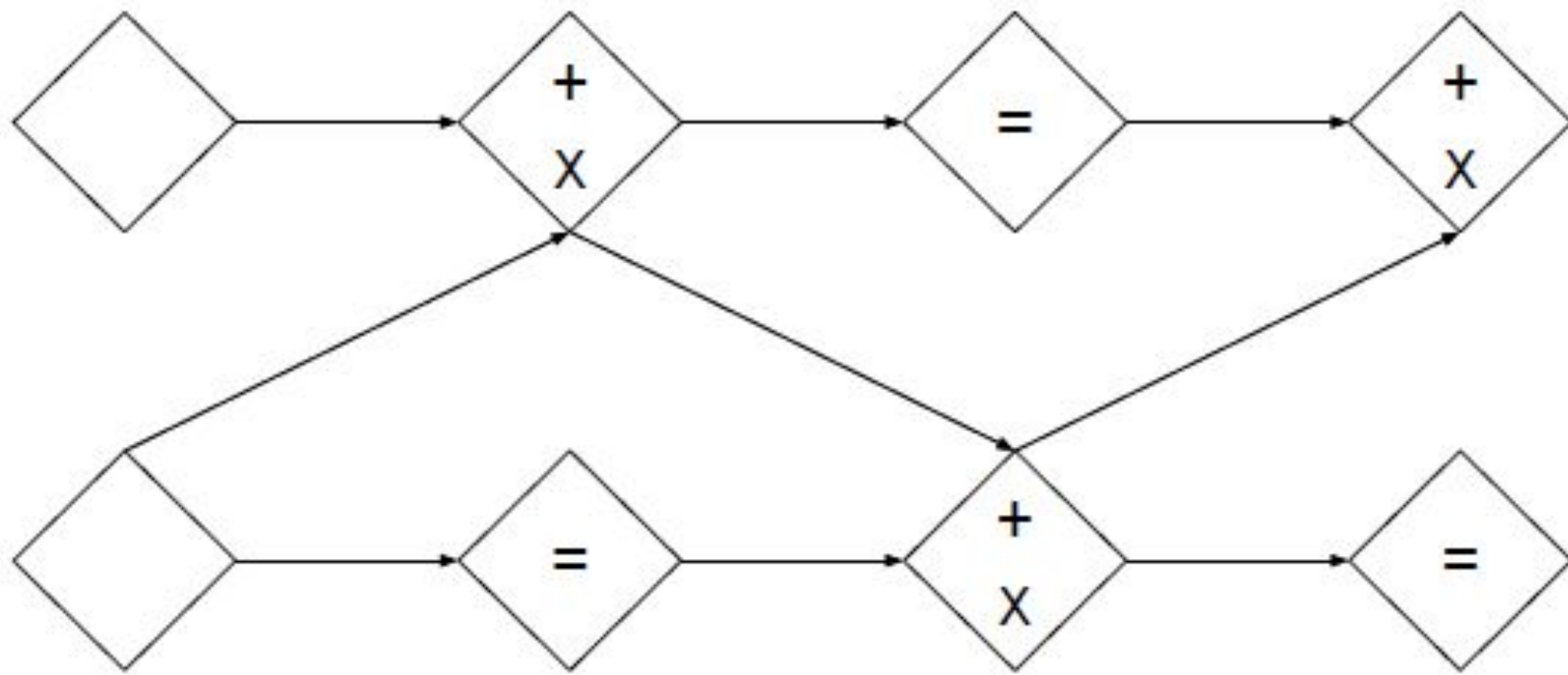
# Combining coupling layers

$$\frac{\partial(f_b \circ f_a)}{\partial x_a^T}(x_a) = \frac{\partial f_a}{\partial x_a^T}(x_a) \cdot \frac{\partial f_b}{\partial x_b^T}(x_b = f_a(x_a))$$

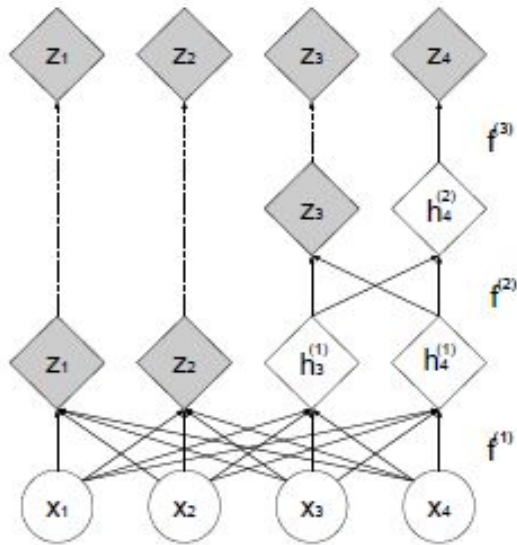
$$\det(A \cdot B) = \det(A) \det(B).$$

$$(f_b \circ f_a)^{-1} = f_a^{-1} \circ f_b^{-1}.$$





# Multi-scale architecture



$$h^{(0)} = x$$

$$(z^{(i+1)}, h^{(i+1)}) = f^{(i+1)}(h^{(i)})$$

$$z^{(L)} = f^{(L)}(h^{(L-1)})$$

$$z = (z^{(1)}, \dots, z^{(L)}).$$

# Batch Normalization

$$x \mapsto \frac{x - \tilde{\mu}}{\sqrt{\tilde{\sigma}^2 + \epsilon}}$$

$$\left( \prod_i (\tilde{\sigma}_i^2 + \epsilon) \right)^{-\frac{1}{2}}$$

# Experiment

- Multi-scale architecture (repeated recursively until the input of the last recursion is a  $4 * 4 * c$  tensor)
- Deep convolutional residue network in coupling layer
- Optimizer: ADAM
- L2 regularization on weight scale parameters

# Result

Dataset	PixelRNN [46]	Real NVP	Conv DRAW [22]	IAF-VAE [34]
<b>CIFAR-10</b>	3.00	3.49	< 3.59	< 3.28
<b>Imagenet</b> ( $32 \times 32$ )	3.86 (3.83)	4.28 (4.26)	< 4.40 (4.35)	
<b>Imagenet</b> ( $64 \times 64$ )	3.63 (3.57)	3.98 (3.75)	< 4.10 (4.04)	
<b>LSUN</b> (bedroom)		2.72 (2.70)		
<b>LSUN</b> (tower)		2.81 (2.78)		
<b>LSUN</b> (church outdoor)		3.08 (2.94)		
<b>CelebA</b>		3.02 (2.97)		

# Results





# Results



# Smooth semantically consistent meaning

$$z = \cos(\phi) (\cos(\phi')z_{(1)} + \sin(\phi')z_{(2)}) + \sin(\phi) (\cos(\phi')z_{(3)} + \sin(\phi')z_{(4)})$$







# Pros and cons

- Performance: Competitive but has lots of room for improvement
- Pros:
  - Allow tractable and exact log likelihood
  - Flexible functional form allows fast and exact sampling from model distribution
  - Does not require the use of fixed form reconstruction cost
  - Able to learn a semantically meaningful latent space
- Cons:
  - Latent space does not reduce number of dimensions
  - Need much more space, too much networks





# Discussion points

- In the multi-scale architecture, the factoring operation factors out lots of positions without going through enough transformation. Would this be a factor that hurts the final performance of the model?
- What other masks we can use besides the checkerboard mask and the channel-wise mask?

# Useful resources

- <https://www.youtube.com/watch?v=6GUSrmo9Qpw&t=458s> The author of this paper talking about the model
- [https://www.youtube.com/watch?v=u3vVyFVU\\_I](https://www.youtube.com/watch?v=u3vVyFVU_I) Introduction to normalizing flow
- <https://arxiv.org/abs/1908.09257> Introduction to lots of flow methods: “Normalizing Flows: An Introduction and Review of Current Methods”