# Self Normalizing Flows

T. Anderson Keller, Jorn Peters, Priyank Jaini,
Emiel Hoogeboom, Patrick Forré, Max Welling

https://arxiv.org/abs/2011.07248

# Self Normalizing Flows

T. Anderson Keller

Jorn Peters

Priyank Jaini

Emiel Hoogeboom

Patrick Forré

Max Welling

$$f = g^{-1}$$

$$p(\mathbf{x}) = p(\mathbf{z}) \, |\mathbf{J}_f|$$

$$p(\mathbf{x}) = p(\mathbf{z}) \left| \frac{\partial \mathbf{z}}{\partial \mathbf{x}} \right|$$

$$p(\mathbf{x}) = p(\mathbf{z}) \, |\mathbf{J}_{g^{-1}}|$$

$$g = f^{-1}$$

$g(\partial \mathbf{z}) = \partial \mathbf{x}$

$f(\partial \mathbf{x}) = \partial \mathbf{z}$

$$f = g^{-1}$$

$$p(\mathbf{x}) = p(\mathbf{z}) \,|\mathbf{J}_f|$$

$$p(\mathbf{x}) = p(\mathbf{z}) \,\left|\frac{\partial \mathbf{z}}{\partial \mathbf{x}}\right|$$

$$g(\partial \mathbf{z}) = \partial \mathbf{x}$$

$$p(\mathbf{x}) = p(\mathbf{z}) \,|\mathbf{J}_{g^{-1}}|$$
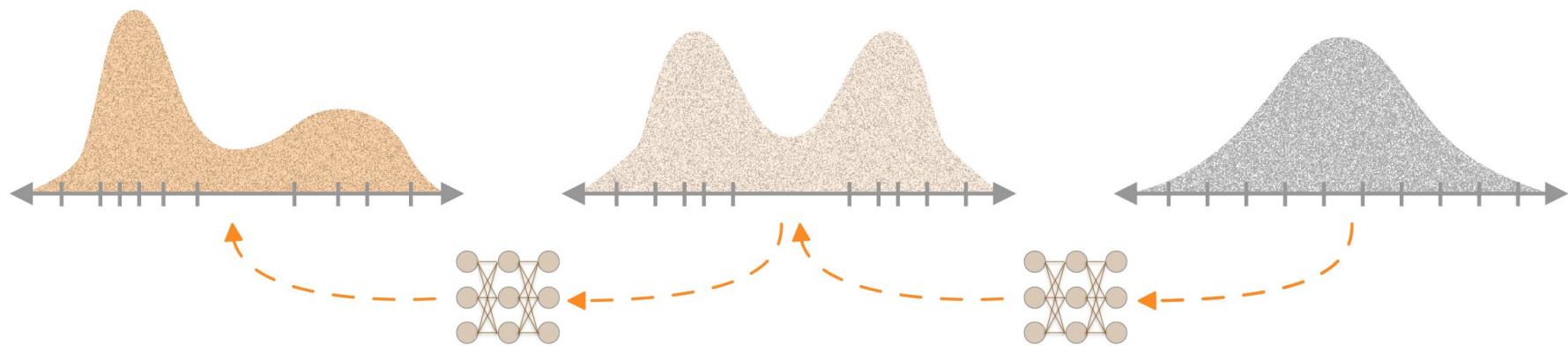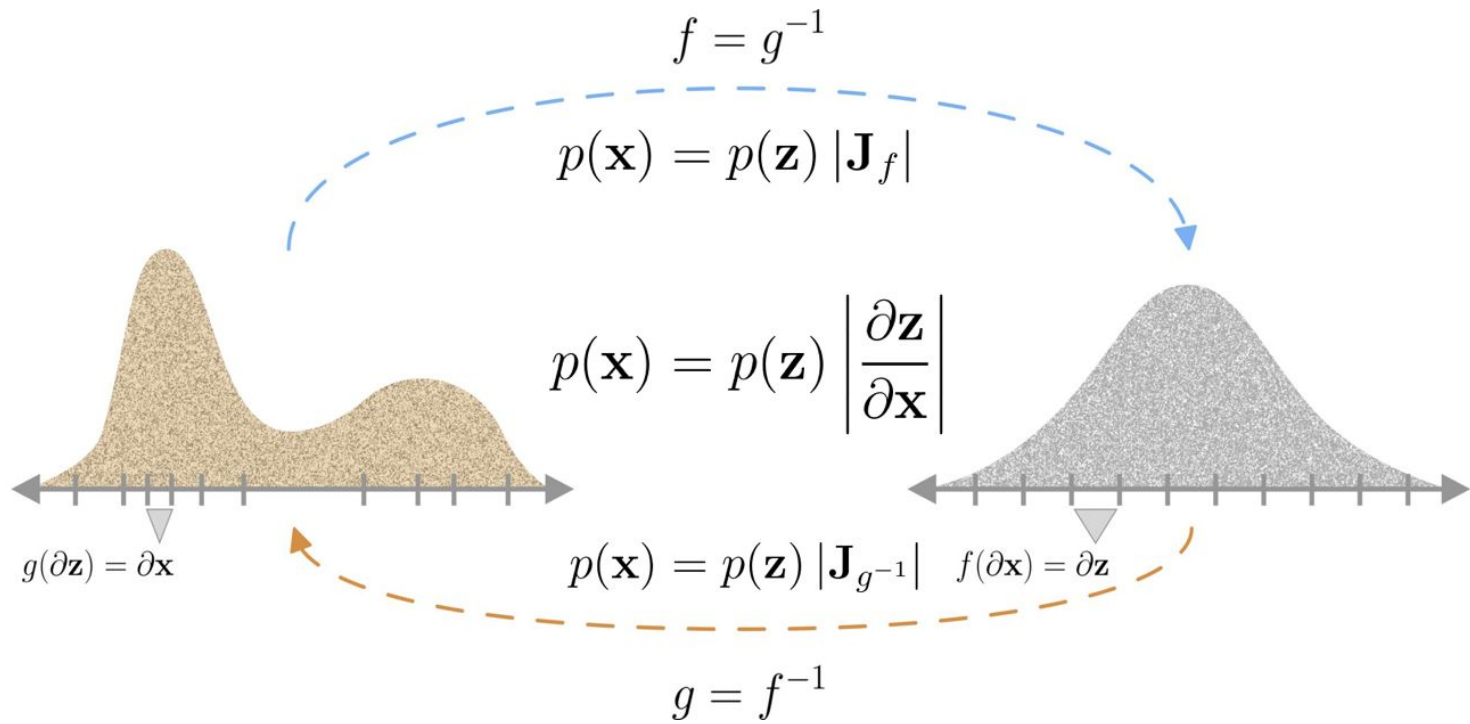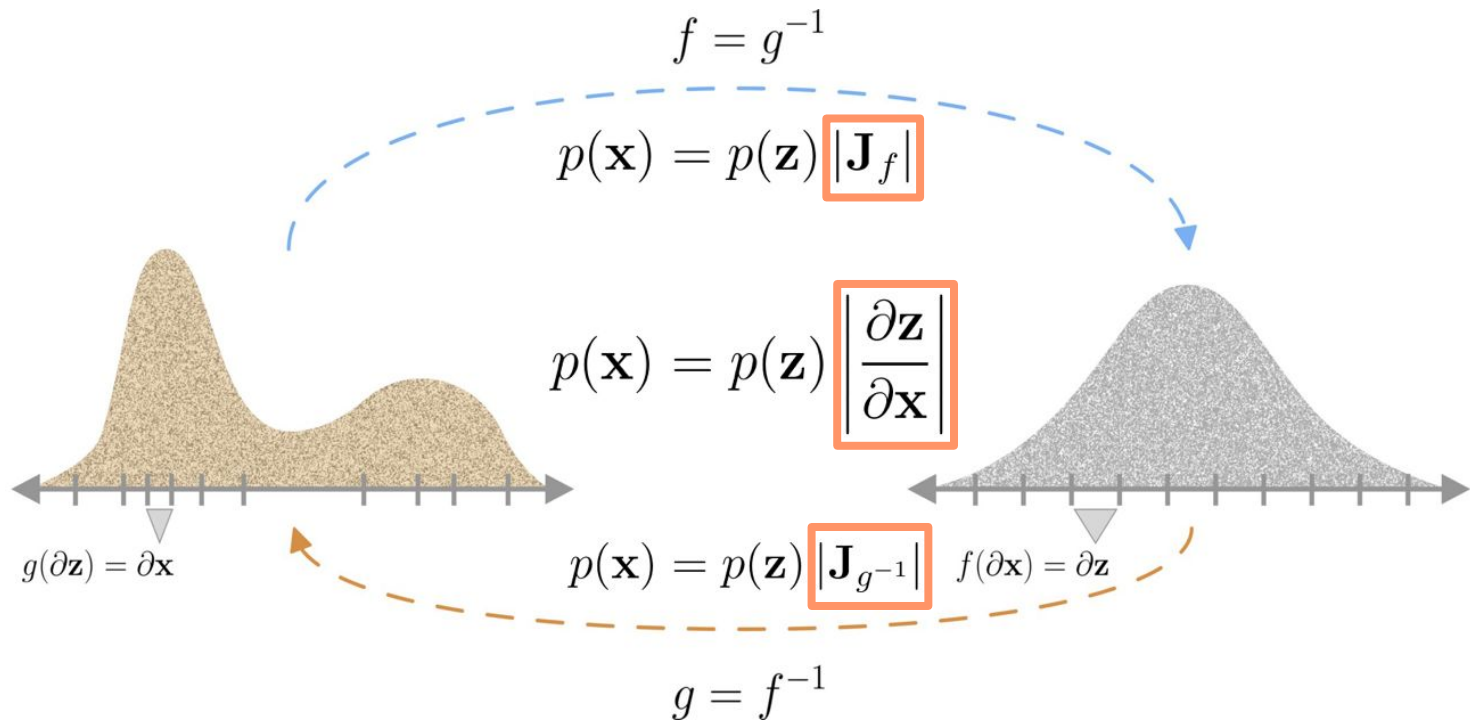
$$f(\partial \mathbf{x}) = \partial \mathbf{z}$$

$$g = f^{-1}$$

# Prior Work

- NICE (Non-linear independent components estimation) (Dinh et al., 2015)
- Real non-volume preserving flow (Real NVP) (Dinh et al., 2017)
- Inverse autoregressive flow (IAF) (Kingma et al., 2016)
- Masked autoregressive flow (MAF) (Papamakarios et al., 2017)
- Glow (Kingma and Dhariwal, 2018)
- Neural Autoregressive Flow (NAF) (Huang et al., 2018)
- block-NAF (B-NAF) (De Cao et al., 2019)
- Flow++ (Ho et al., 2019)
- Sums-of-squares Polynomial transformer (Jaini et al., 2019)

# Prior Work

- Neural Spline Flows (Durkan et al., 2019)
- Residual Flows (Chen et al., 2019)
- Invertible Residual Networks (Jens Behrmann et al., 2018)
- Sylvester Flows (van den Berd et al., 2018)
- Radial Flows (Tabak and Turner, 2013)
- Planar Flows (Rezende and Mohamed, 2015)
- Emerging Convolutions (Hoogeboom et al., 2019)
- Integer Discrete Flows (Hoogeboom et al., 2019)
- The Convolution Exponential (Hoogeboom et al., 2020)

$$\frac{\partial \log |\mathbf{J}_f|}{\partial \mathbf{J}_f}$$

$$\frac{\partial \log |\mathbf{J}_f|}{\partial \mathbf{J}_f} = \mathbf{J}_f^{-T}$$

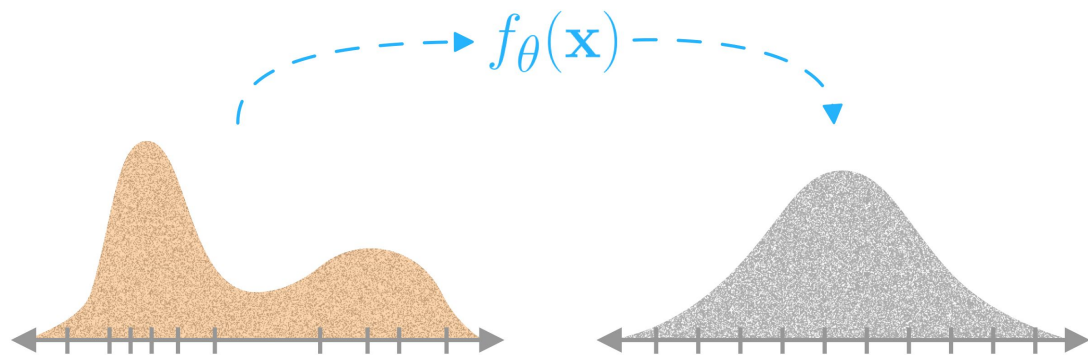$$\frac{\partial \log |\mathbf{J}_f|}{\partial \mathbf{J}_f} = \mathbf{J}_f^{-T} = \mathbf{J}_{f^{-1}}^T$$

$$\frac{\partial \log |\mathbf{J}_f|}{\partial \mathbf{J}_f} = \mathbf{J}_f^{-T} = \mathbf{J}_{f^{-1}}^{T}$$

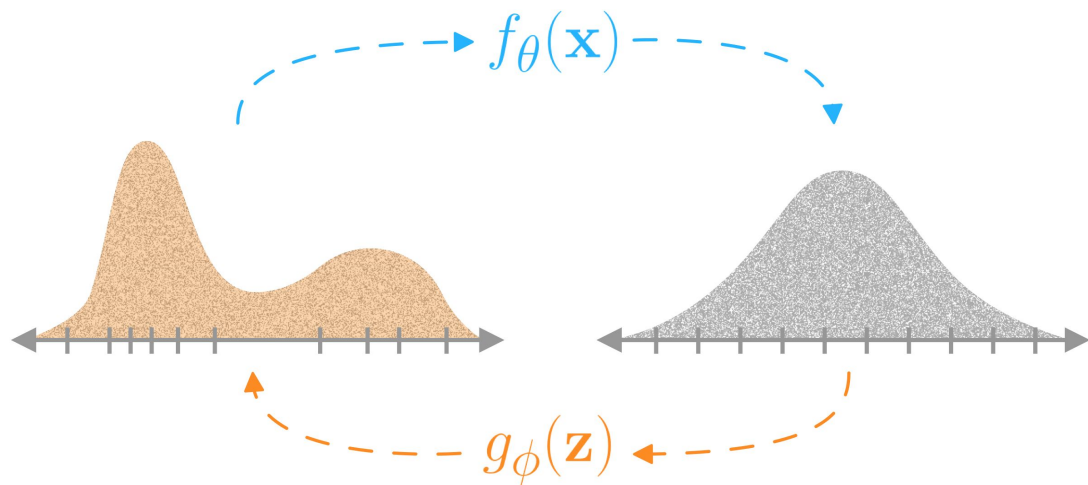$$\text{if } g \approx f^{-1}$$

$$\frac{\partial \log |\mathbf{J}_f|}{\partial \mathbf{J}_f} = \mathbf{J}_f^{-T} = \mathbf{J}_{f^{-1}}^T \approx \mathbf{J}_g^T$$
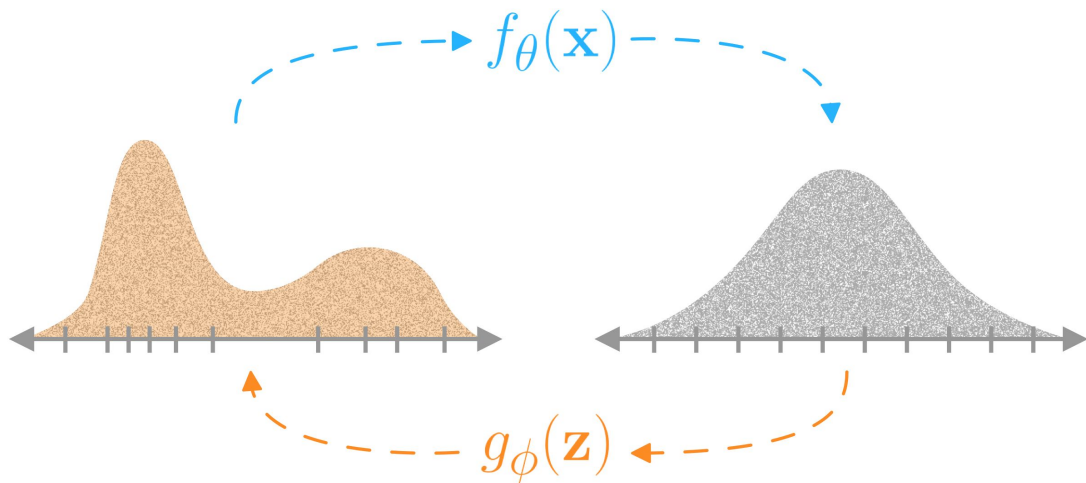
$$\text{if } g \approx f^{-1}$$
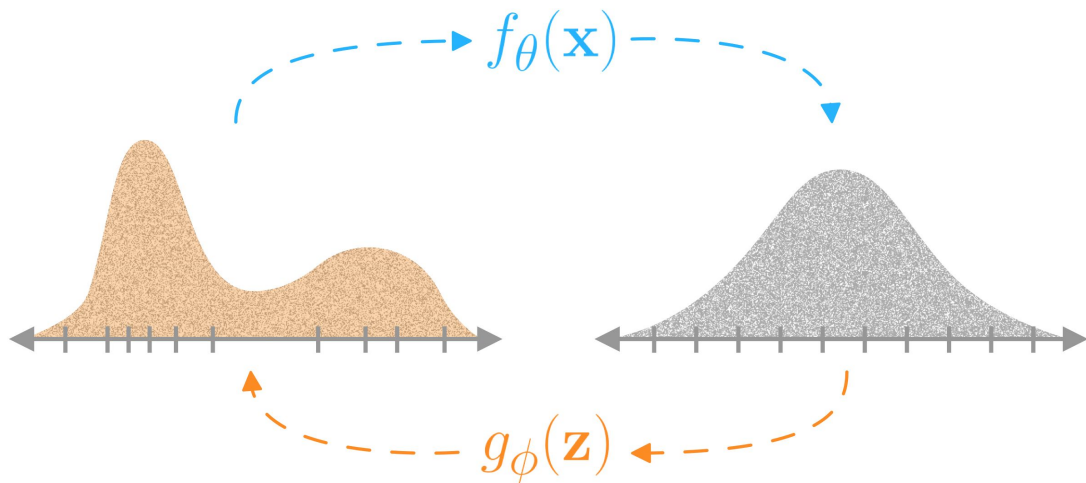
# Self-Normalizing Flows

# Self-Normalizing Flows

# Self-Normalizing Flows



$$\mathcal{L}(\mathbf{x}) = ||g_\phi(f_\theta(\mathbf{x})) - \mathbf{x}||_2^2$$

# Self-Normalizing Flows



$$\mathcal{L}(\mathbf{x}) = ||g_\phi(f_\theta(\mathbf{x})) - \mathbf{x}||_2^2$$

$$\log p_{\mathbf{X}}(\mathbf{x}) = \frac{1}{2}\log p_{\mathbf{X}}^f(\mathbf{x}) + \frac{1}{2}\log p_{\mathbf{X}}^g(\mathbf{x})$$

# Self-Normalizing Flows

$$\log p_{\mathbf{X}}^{f}(\mathbf{x}) = \log p_{\mathbf{Z}}\left(f_{\theta}(\mathbf{x})\right) + \log |\mathbf{J}_{f}|$$

$$\log p_{\mathbf{X}}^{g}(\mathbf{x}) = \log p_{\mathbf{Z}}\left(g_{\phi}^{-1}(\mathbf{x})\right) + \log |\mathbf{J}_{g^{-1}}|$$

# Self-Normalizing Flows

$$\log p_{\mathbf{X}}^{f}(\mathbf{x}) = \log p_{\mathbf{Z}}\left(f_{\theta}(\mathbf{x})\right) + \log |\mathbf{J}_{f}|$$

$$\frac{\partial}{\partial \theta} \log p_{\mathbf{X}}^{f}(\mathbf{x}) = \frac{\partial}{\partial \theta} \log p_{\mathbf{Z}}\left(f_{\theta}(\mathbf{x})\right) + \frac{\partial(\text{vec } \mathbf{J}_{f})^{T}}{\partial \theta}(\text{vec } \mathbf{J}_{f}^{-T})$$

$$\log p_{\mathbf{X}}^{g}(\mathbf{x}) = \log p_{\mathbf{Z}}\left(g_{\phi}^{-1}(\mathbf{x})\right) + \log |\mathbf{J}_{g^{-1}}|$$

$$\frac{\partial}{\partial \phi} \log p_{\mathbf{X}}^{g}(\mathbf{x}) = \frac{\partial}{\partial \phi} \log p_{\mathbf{Z}}\left(g_{\phi}^{-1}(\mathbf{z})\right) + \frac{\partial(\text{vec } \mathbf{J}_{g^{-1}})^{T}}{\partial \phi}(\text{vec } \mathbf{J}_{g^{-1}}^{-T})$$

# Self-Normalizing Flows

$$\log p_{\mathbf{X}}^{f}(\mathbf{x}) = \log p_{\mathbf{Z}}\left(f_{\theta}(\mathbf{x})\right) + \log |\mathbf{J}_{f}|$$

$$\frac{\partial}{\partial \theta} \log p_{\mathbf{X}}^{f}(\mathbf{x}) = \frac{\partial}{\partial \theta} \log p_{\mathbf{Z}}\left(f_{\theta}(\mathbf{x})\right) + \frac{\partial(\text{vec } \mathbf{J}_{f})^{T}}{\partial \theta}(\text{vec } \mathbf{J}_{f}^{-T})$$

$$\approx \frac{\partial}{\partial \theta} \log p_{\mathbf{Z}}\left(f_{\theta}(\mathbf{x})\right) + \frac{\partial(\text{vec } \mathbf{J}_{f})^{T}}{\partial \theta}(\text{vec } \mathbf{J}_{g}^{T})$$

$$\log p_{\mathbf{X}}^{g}(\mathbf{x}) = \log p_{\mathbf{Z}}\left(g_{\phi}^{-1}(\mathbf{x})\right) + \log |\mathbf{J}_{g^{-1}}|$$

$$\frac{\partial}{\partial \phi} \log p_{\mathbf{X}}^{g}(\mathbf{x}) = \frac{\partial}{\partial \phi} \log p_{\mathbf{Z}}\left(g_{\phi}^{-1}(\mathbf{z})\right) + \frac{\partial(\text{vec } \mathbf{J}_{g^{-1}})^{T}}{\partial \phi}(\text{vec } \mathbf{J}_{g^{-1}}^{-T})$$

$$\approx \frac{\partial}{\partial \phi} \log p_{\mathbf{Z}}\left(g_{\phi}^{-1}(\mathbf{z})\right) - \frac{\partial(\text{vec } \mathbf{J}_{g})^{T}}{\partial \phi}(\text{vec } \mathbf{J}_{f}^{T})$$

# Fully-Connected

$$f(\boldsymbol{x}) = \boldsymbol{W}\boldsymbol{x} = \boldsymbol{z}$$

$$g(\boldsymbol{z}) = \boldsymbol{R}\boldsymbol{z}$$

$$\boldsymbol{W}^{-1} \approx \boldsymbol{R}$$

$$\frac{\partial}{\partial \mathbf{W}} \log p_{\mathbf{X}}^f(\mathbf{x}) = \frac{\partial}{\partial \mathbf{W}} \log p_{\mathbf{Z}}(\mathbf{W}\mathbf{x}) + \mathbf{W}^{-T}$$

$$\approx \delta_{\mathbf{z}} \mathbf{x}^T + \mathbf{R}^T$$

$$\frac{\partial}{\partial \mathbf{R}} \log p_{\mathbf{X}}^g(\mathbf{x}) = \frac{\partial}{\partial \mathbf{R}} \log p_{\mathbf{Z}}(\mathbf{R}^{-1}\mathbf{x}) - \mathbf{R}^{-T}$$

$$\approx -\delta_{\mathbf{x}} \mathbf{z}^T - \mathbf{W}^T$$

$$\boldsymbol{\delta}_{\boldsymbol{x}} = \frac{\partial \log p_{\boldsymbol{Z}}(\boldsymbol{z})}{\partial \boldsymbol{x}}$$

$$\boldsymbol{\delta}_{\boldsymbol{z}} = \frac{\partial \log p_{\boldsymbol{Z}}(\boldsymbol{z})}{\partial \boldsymbol{z}}$$

# Fully-Connected

$f(\boldsymbol{x}) = \boldsymbol{W}\boldsymbol{x} = \boldsymbol{z}$

$g(\boldsymbol{z}) = \boldsymbol{R}\boldsymbol{z}$

$\boldsymbol{W}^{-1} \approx \boldsymbol{R}$

$$\frac{\partial}{\partial \mathbf{W}} \log p_{\mathbf{X}}^f(\mathbf{x}) = \frac{\partial}{\partial \mathbf{W}} \log p_{\mathbf{Z}}(\mathbf{W}\mathbf{x}) + \mathbf{W}^{-T}$$

$$\approx \delta_{\mathbf{z}} \mathbf{x}^T + \mathbf{R}^T$$

$$\frac{\partial}{\partial \mathbf{R}} \log p_{\mathbf{X}}^g(\mathbf{x}) = \frac{\partial}{\partial \mathbf{R}} \log p_{\mathbf{Z}}(\mathbf{R}^{-1}\mathbf{x}) - \mathbf{R}^{-T}$$

$$\approx -\delta_{\mathbf{x}} \mathbf{z}^T - \mathbf{W}^T$$

$\boldsymbol{\delta_x} = \frac{\partial \log p_{\boldsymbol{Z}}(\boldsymbol{z})}{\partial \boldsymbol{x}}$

$\boldsymbol{\delta_z} = \frac{\partial \log p_{\boldsymbol{Z}}(\boldsymbol{z})}{\partial \boldsymbol{z}}$

# Convolutional

$f(\boldsymbol{x}) = \boldsymbol{w} \star \boldsymbol{x} = \boldsymbol{z}$

$g(\boldsymbol{z}) = \boldsymbol{r} \star \boldsymbol{z}$

$f^{-1} \approx g$

$$\frac{\partial}{\partial \boldsymbol{w}} \log p_{\boldsymbol{X}}^f(\boldsymbol{x}) = \boldsymbol{\delta}_{\boldsymbol{z}}^f \star \boldsymbol{x} + \frac{\partial (\mathrm{vec}\,\mathcal{T}(\boldsymbol{w}))^T}{\partial \boldsymbol{w}} \left(\mathrm{vec}\,\mathcal{T}(\boldsymbol{w})^{-T}\right)$$

$$\approx \boldsymbol{\delta}_{\boldsymbol{z}} \star \boldsymbol{x} + \mathrm{flip}(\boldsymbol{r}) \odot \boldsymbol{m}$$

$$\frac{\partial}{\partial \boldsymbol{r}} \log p_{\boldsymbol{X}}^g(\boldsymbol{x}) = \frac{\partial (\mathrm{vec}\,\mathcal{T}(\boldsymbol{r}))^T}{\partial \boldsymbol{r}} \left(\mathrm{vec}\,\left[-\mathcal{T}(\boldsymbol{r})^{-T}\boldsymbol{\delta}_{\boldsymbol{z}}^g \boldsymbol{x}^T \mathcal{T}(\boldsymbol{r})^{-T}\right] - \mathrm{vec}\,\mathcal{T}(\boldsymbol{r})^{-T}\right)$$

$$\approx -\boldsymbol{\delta}_{\boldsymbol{x}} \star \boldsymbol{z} - \mathrm{flip}(\boldsymbol{w}) \odot \boldsymbol{m}$$
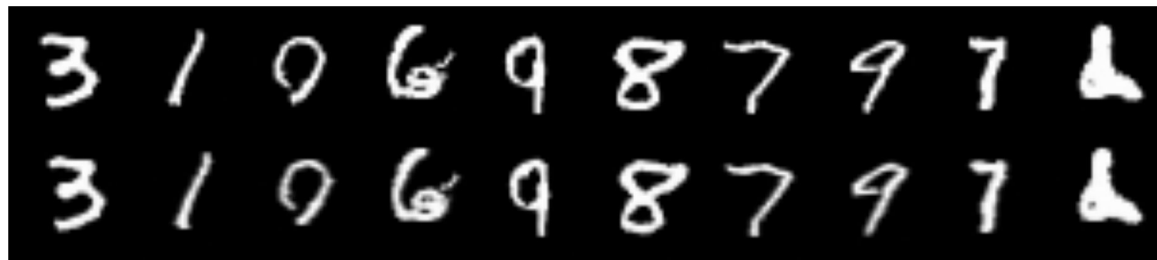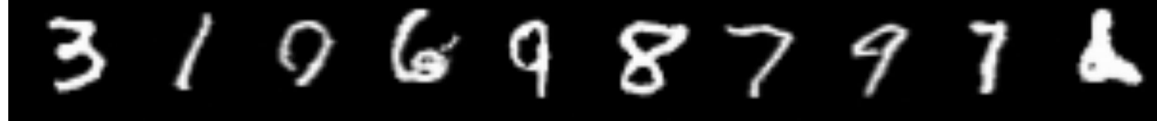
# Experiments

# Experiments

# Experiments



$$\mathbf{z} \sim p_{\mathbf{Z}}(\mathbf{z})$$

$f^{-1}(\mathbf{z})$

$g(\mathbf{z})$

$f^{-1}(\mathbf{z})$

$g(\mathbf{z})$

# Experiments

| Model | $-\log p_{\boldsymbol{X}}(\boldsymbol{x})$ |
|---|---|
| Relative Grad. FC 2-Layer [13] | $1096.5 \pm 0.5$ |
| Exact Gradient FC 2-Layer | $947.6 \pm 0.2$ |
| SNF FC 2-Layer (ours) | $947.1 \pm 0.2$ |
| Emerging Conv. 9-Layer [16] | $645.7 \pm 3.6$ |
| SNF Conv. 9-Layer (ours) | $638.6 \pm 0.9$ |
| Conv. Exponential 9-Layer [15] | $638.1 \pm 1.0$ |
| Exact Gradient Conv. 9-Layer | $637.4 \pm 0.2$ |
| Glow-like 32-Layer [20] | $575.7 \pm 0.8$ |
| SNF Glow 32-Layer (ours) | $575.4 \pm 1.4$ |

# Experiments

| Model | CIFAR-10 | ImageNet32 |
|---|---|---|
| Glow | $3.36 \pm 0.002$ | $4.12 \pm 0.002$ |
| SNF Glow | $3.37 \pm 0.004$ | $4.14 \pm 0.007$ |

# Thank you!

Paper: https://arxiv.org/abs/2011.07248

Code: https://github.com/akandykeller/SelfNormalizingFlows

Blog: http://keller.org/research/2020-10-21-self-normalizing-flows/

Contact: T.Anderson.Keller@Gmail.com