#### Deep Unsupervised Learning using Nonequilibrium Thermodynamics

#### **Jascha Sohl-Dickstein**<sup>1</sup>, Eric Weiss<sup>2</sup>, Niru Maheswaranathan<sup>3</sup>, Surya Ganguli<sup>3</sup>

<sup>1</sup> Google Brain, <sup>2</sup> UC Berkeley, <sup>3</sup> Stanford University





Niru



Surya

Eric

## Outline

- Motivation: The promise of deep unsupervised learning
- Physical intuition: Diffusion processes and time reversal
- Diffusion probabilistic model: Derivation and experimental results
- Other projects: Inverse Ising, non-equilibrium Monte Carlo, stat. mech. of neural networks

Jascha Sohl-Dickstein

### Outline

- Motivation: The promise of deep unsupervised learning
- Physical intuition: Diffusion processes and time reversal
- Diffusion probabilistic model: Derivation and experimental results
- Other projects: Inverse Ising, non-equilibrium Monte Carlo, stat. mech. of neural networks

Jascha Sohl-Dickstein

• Unknown features/labels

- Unknown features/labels
  - Novel modalities

- Unknown features/labels
  - Novel modalities



[Trans Biomed Eng, 2015]

Jascha Sohl-Dickstein

- Unknown features/labels
  - Novel modalities

- Unknown features/labels
  - Novel modalities
  - Exploratory data analysis

- Unknown features/labels
  - Novel modalities
  - Exploratory data analysis

7 exemplar multiunits responding to 40 repeated trials of natural video in cat V1



[PLoS Comp Bio 2014] [Neuron 2013]

Jascha Sohl-Dickstein

- Unknown features/labels
  - Novel modalities
  - Exploratory data analysis

- Unknown features/labels
  - Novel modalities
  - Exploratory data analysis
- Expensive labels

Jascha Sohl-Dickstein

- Unknown features/labels
  - Novel modalities
  - Exploratory data analysis
- Expensive labels





[SPIE 2009] [Med Phys 2014]

- Unknown features/labels
  - Novel modalities
  - Exploratory data analysis
- Expensive labels

Jascha Sohl-Dickstein

- Unknown features/labels
  - Novel modalities
  - Exploratory data analysis
- Expensive labels
- Unpredictable tasks / one shot learning

## Outline

- Motivation: The promise of deep unsupervised learning
- Physical intuition: Diffusion processes and time reversal
- Diffusion probabilistic model: Derivation and experimental results
- Other projects: Inverse Ising, non-equilibrium Monte Carlo, stat. mech. of neural networks

Jascha Sohl-Dickstein

## Outline

- Motivation: The promise of deep unsupervised learning
- Physical intuition: Diffusion processes and time reversal
  - Destroy structure in data
  - Carefully characterize the destruction
  - Learn how to reverse time
- Diffusion probabilistic model: Derivation and experimental results
- Other projects: Inverse Ising, non-equilibrium Monte Carlo, stat. mech. of neural networks

Jascha Sohl-Dickstein



 Dye density represents probability density



- Dye density represents probability density
- Goal: Learn structure of probability density



- Dye density represents probability density
- Goal: Learn structure of probability density
- Observation: Diffusion destroys structure



- Dye density represents probability density
- Goal: Learn structure of probability density
- Observation: Diffusion destroys structure



- Dye density represents probability density
- Goal: Learn structure of probability density
- Observation: Diffusion destroys structure

#### Data distribution

Jascha Sohl-Dickstein

Uniform distribution



• What if we could reverse time?

Jascha Sohl-Dickstein



• What if we could reverse time?

Jascha Sohl-Dickstein



• What if we could reverse time?

Data distribution



#### Uniform distribution

Jascha Sohl-Dickstein



- What if we could reverse time?
- Recover data distribution by starting from uniform distribution and running dynamics backwards

Data distribution

Jascha Sohl-Dickstein

Diffusion Probabilistic Models

Uniform distribution



- What if we could reverse time?
- Recover data distribution by starting from uniform distribution and running dynamics backwards

Data distribution

Jascha Sohl-Dickstein



#### Uniform distribution



Jascha Sohl-Dickstein

Jascha Sohl-Dickstein



© Rutger Saly

- Microscopic view
- Brownian motion

Jascha Sohl-Dickstein



© Rutger Saly

- Microscopic view
- Brownian motion

Jascha Sohl-Dickstein



© Rutger Saly

- Microscopic view
- Brownian motion

Jascha Sohl-Dickstein



© Rutger Saly

- Microscopic view
- Brownian motion

Jascha Sohl-Dickstein



© Rutger Saly

- Microscopic view
- Brownian motion
- Position updates are small Gaussians

Jascha Sohl-Dickstein



© Rutger Saly

- Microscopic view
- Brownian motion
- Position updates are small Gaussians

Jascha Sohl-Dickstein



© Rutger Saly

- Microscopic view
- Brownian motion
- Position updates are small Gaussians

Jascha Sohl-Dickstein
# Observation 2: Microscopic Diffusion is Time Reversible



© Rutger Saly

- Microscopic view
- Brownian motion
- Position updates are small Gaussians
  - Both forwards and backwards in time

Jascha Sohl-Dickstein

Jascha Sohl-Dickstein

Destroy all structure in data distribution using diffusion process

- Destroy all structure in data distribution using diffusion process
- Learn reversal of diffusion process
  - Estimate function for mean and covariance of each step in the reverse diffusion process (binomial rate for binary data)

- Destroy all structure in data distribution using diffusion process
- Learn reversal of diffusion process
  - Estimate function for mean and covariance of each step in the reverse diffusion process (binomial rate for binary data)
- Reverse diffusion process is the model of the data

# Outline

- Motivation: The promise of deep unsupervised learning
- Physical intuition: Diffusion processes and time reversal
- Diffusion probabilistic model: Derivation and experimental results
  - Algorithm
  - Deep convolutional network: Universal function approximator
  - Multiplying distributions: Inputation, denoising, computing posteriors
- Other projects: Inverse Ising, non-equilibrium Monte Carlo, stat. mech. of neural networks

# Outline

- Motivation: The promise of deep unsupervised learning
- Physical intuition: Diffusion processes and time reversal
- Diffusion probabilistic model: Derivation and experimental results
  - Algorithm
  - Deep convolutional network: Universal function approximator
  - Multiplying distributions: Inputation, denoising, computing posteriors
- Other projects: Inverse Ising, non-equilibrium Monte Carlo, stat. mech. of neural networks

Data distribution

$$q\left(\mathbf{x}^{(0)}\right)$$

Jascha Sohl-Dickstein

Data distribution

Forward diffusion

$$q\left(\mathbf{x}^{(0)}\right)$$

$$q\left(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)}\right) = \mathcal{N}\left(\mathbf{x}^{(t)}; \mathbf{x}^{(t-1)}\sqrt{1-\beta_t}, \mathbf{I}\beta_t\right)$$

Jascha Sohl-Dickstein

Data distribution

Forward diffusion

$$q\left(\mathbf{x}^{(0)}\right)$$

$$q\left(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)}\right) = \mathcal{N}\left(\mathbf{x}^{(t)};\mathbf{x}^{(t-1)}\sqrt{1-\beta_t},\mathbf{I}\beta_t\right)$$

Decay towards origin

Jascha Sohl-Dickstein

Data distribution

Forward diffusion

$$q\left(\mathbf{x}^{(0)}\right)$$

$$q\left(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)}\right) = \mathcal{N}\left(\mathbf{x}^{(t)};\mathbf{x}^{(t-1)}\sqrt{1-\beta_t},\mathbf{I}\beta_t\right)$$

Decay towards origin

Add small noise

Jascha Sohl-Dickstein



Jascha Sohl-Dickstein

# Forward Diffusion Process on Swiss Roll









Jascha Sohl-Dickstein

# Forward Diffusion Process on Swiss Roll









Jascha Sohl-Dickstein

Noise distribution

$$p\left(\mathbf{x}^{(T)}\right) = \mathcal{N}\left(\mathbf{x}^{(T)}; 0, \mathbf{I}\right)$$

Jascha Sohl-Dickstein



Jascha Sohl-Dickstein



Jascha Sohl-Dickstein



Jascha Sohl-Dickstein

# Learned Reverse Diffusion Process on Swiss Roll

Gaussian diffusion

w/learned kernel



Jascha Sohl-Dickstein

Data

dist.

**Diffusion Probabilistic Models** 

Isotropic

Gaussian

# Learned Reverse Diffusion Process on Swiss Roll

Gaussian diffusion

w/learned kernel



Jascha Sohl-Dickstein

Data

dist.

**Diffusion Probabilistic Models** 

Isotropic

Gaussian

# Summary of Forward and Reverse Diffusion on Swiss Roll



**Diffusion Probabilistic Models** 

# Summary of Forward and Reverse Diffusion on Swiss Roll



Jascha Sohl-Dickstein

# Summary of Forward and Reverse Diffusion on Swiss Roll



Jascha Sohl-Dickstein

# Training the Reverse Diffusion Process

Model probability

$$p\left(\mathbf{x}^{(0)}\right) = \int d\mathbf{x}^{(1\cdots T)} p\left(\mathbf{x}^{(0\cdots T)}\right)$$

# Training the Reverse Diffusion Process

Model probability

$$p\left(\mathbf{x}^{(0)}\right) = \int d\mathbf{x}^{(1\cdots T)} p\left(\mathbf{x}^{(0\cdots T)}\right)$$

Annealed importance sampling / Jarzynski equality

$$p\left(\mathbf{x}^{(0)}\right) = \int d\mathbf{x}^{(1\cdots T)} q\left(\mathbf{x}^{(1\cdots T)} | \mathbf{x}^{(0)}\right) \frac{p\left(\mathbf{x}^{(0\cdots T)}\right)}{q\left(\mathbf{x}^{(1\cdots T)} | \mathbf{x}^{(0)}\right)}$$

# Training the Reverse Diffusion Process

Model probability

$$p\left(\mathbf{x}^{(0)}\right) = \int d\mathbf{x}^{(1\cdots T)} p\left(\mathbf{x}^{(0\cdots T)}\right)$$

Annealed importance sampling / Jarzynski equality

$$p\left(\mathbf{x}^{(0)}\right) = \int d\mathbf{x}^{(1\cdots T)} q\left(\mathbf{x}^{(1\cdots T)} | \mathbf{x}^{(0)}\right) \frac{p\left(\mathbf{x}^{(0\cdots T)}\right)}{q\left(\mathbf{x}^{(1\cdots T)} | \mathbf{x}^{(0)}\right)}$$

Log Likelihood

$$L = \int d\mathbf{x}^{(0)} q\left(\mathbf{x}^{(0)}\right) \log\left[\int d\mathbf{x}^{(1\cdots T)} q\left(\mathbf{x}^{(1\cdots T)}\right) \frac{p\left(\mathbf{x}^{(0\cdots T)}\right)}{q\left(\mathbf{x}^{(1\cdots T)}\right)}\right]$$

Model probability

$$p\left(\mathbf{x}^{(0)}\right) = \int d\mathbf{x}^{(1\cdots T)} p\left(\mathbf{x}^{(0\cdots T)}\right)$$

Annealed importance sampling / Jarzynski equality

$$p\left(\mathbf{x}^{(0)}\right) = \int d\mathbf{x}^{(1\cdots T)} q\left(\mathbf{x}^{(1\cdots T)} | \mathbf{x}^{(0)}\right) \frac{p\left(\mathbf{x}^{(0\cdots T)}\right)}{q\left(\mathbf{x}^{(1\cdots T)} | \mathbf{x}^{(0)}\right)}$$

Log Likelihood

$$L = \int d\mathbf{x}^{(0)} q\left(\mathbf{x}^{(0)}\right) \log\left[\int d\mathbf{x}^{(1\cdots T)} q\left(\mathbf{x}^{(1\cdots T)}\right) \frac{p\left(\mathbf{x}^{(0\cdots T)}\right)}{q\left(\mathbf{x}^{(1\cdots T)}\right)}\right]$$

Jensen's inequality

$$L \ge \int d\mathbf{x}^{(0\cdots T)} q\left(\mathbf{x}^{(0\cdots T)}\right) \log \left[\frac{p\left(\mathbf{x}^{(0\cdots T)}\right)}{q\left(\mathbf{x}^{(1\cdots T)} | \mathbf{x}^{(0)}\right)}\right]$$

$$p\left(\mathbf{x}^{(0)}\right) = \int d\mathbf{x}^{(1\cdots T)} q\left(\mathbf{x}^{(1\cdots T)} | \mathbf{x}^{(0)}\right) \frac{p\left(\mathbf{x}^{(0\cdots T)}\right)}{q\left(\mathbf{x}^{(1\cdots T)} | \mathbf{x}^{(0)}\right)}$$

Log Likelihood

$$L = \int d\mathbf{x}^{(0)} q\left(\mathbf{x}^{(0)}\right) \log\left[\int d\mathbf{x}^{(1\cdots T)} q\left(\mathbf{x}^{(1\cdots T)}\right) \frac{p\left(\mathbf{x}^{(0\cdots T)}\right)}{q\left(\mathbf{x}^{(1\cdots T)}\right)}\right]$$

Jensen's inequality

$$L \ge \int d\mathbf{x}^{(0\cdots T)} q\left(\mathbf{x}^{(0\cdots T)}\right) \log \left[\frac{p\left(\mathbf{x}^{(0\cdots T)}\right)}{q\left(\mathbf{x}^{(1\cdots T)} | \mathbf{x}^{(0)}\right)}\right]$$

... algebra ...

$$L \ge -\sum_{t=2}^{T} \int d\mathbf{x}^{(0)} d\mathbf{x}^{(t)} q\left(\mathbf{x}^{(0)}, \mathbf{x}^{(t)}\right) D_{KL} \left(q\left(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}\right) \middle| \left| p\left(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}\right) \right. \right. \\ \left. + \operatorname{const} \right.$$

$$L \ge -\sum_{t=2}^{T} \int d\mathbf{x}^{(0)} d\mathbf{x}^{(t)} q\left(\mathbf{x}^{(0)}, \mathbf{x}^{(t)}\right) D_{KL}\left(q\left(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}\right) \middle| \left| p\left(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}\right)\right)$$

 $+ \operatorname{const}$ 

$$L \ge -\sum_{t=2}^{T} \int d\mathbf{x}^{(0)} d\mathbf{x}^{(t)} q\left(\mathbf{x}^{(0)}, \mathbf{x}^{(t)}\right) D_{KL} \left(q\left(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}\right) \middle| \left| p\left(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}\right) \right) + \text{const}$$

$$\mathsf{Gaussian}$$

$$L \ge -\sum_{t=2}^{T} \int d\mathbf{x}^{(0)} d\mathbf{x}^{(t)} q\left(\mathbf{x}^{(0)}, \mathbf{x}^{(t)}\right) D_{KL} \left(q\left(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}\right) \middle| \left| p\left(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}\right) \right\rangle + \text{const}$$

$$+ \text{const}$$
Gaussian

$$L \ge -\sum_{t=2}^{T} \int d\mathbf{x}^{(0)} d\mathbf{x}^{(t)} q\left(\mathbf{x}^{(0)}, \mathbf{x}^{(t)}\right) D_{KL} \left(q\left(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}\right) \middle| \left| p\left(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}\right) \right\rangle + \text{const}$$

$$+ \text{const}$$

$$Gaussian$$

$$p\left(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}\right) = \mathcal{N}\left(\mathbf{x}^{(t-1)}; f_{\mu}\left(\mathbf{x}^{(t)}, t\right), f_{\Sigma}\left(\mathbf{x}^{(t)}, t\right)\right)$$

$$L \ge -\sum_{t=2}^{T} \int d\mathbf{x}^{(0)} d\mathbf{x}^{(t)} q\left(\mathbf{x}^{(0)}, \mathbf{x}^{(t)}\right) D_{KL} \left(q\left(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}\right) \middle| \left| p\left(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}\right) \right) + \text{const}$$

$$p\left(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}\right) = \mathcal{N}\left(\mathbf{x}^{(t-1)}; f_{\mu}\left(\mathbf{x}^{(t)}, t\right), f_{\Sigma}\left(\mathbf{x}^{(t)}, t\right)\right)$$

#### Training

$$\underset{f_{\mu}(\mathbf{x}^{(t)},t),f_{\Sigma}(\mathbf{x}^{(t)},t)}{\operatorname{argmin}} \mathbb{E}\left[D_{KL}\left(q\left(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)},\mathbf{x}^{(0)}\right)\Big|\Big|p\left(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}\right)\right)\right]$$

$$L \ge -\sum_{t=2}^{T} \int d\mathbf{x}^{(0)} d\mathbf{x}^{(t)} q\left(\mathbf{x}^{(0)}, \mathbf{x}^{(t)}\right) D_{KL} \left(q\left(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}\right) \middle| \left| p\left(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}\right) \right) + \text{const}$$

$$p\left(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}\right) = \mathcal{N}\left(\mathbf{x}^{(t-1)}; f_{\mu}\left(\mathbf{x}^{(t)}, t\right), f_{\Sigma}\left(\mathbf{x}^{(t)}, t\right)\right)$$

#### Training

$$\underset{f_{\mu}(\mathbf{x}^{(t)},t),f_{\Sigma}(\mathbf{x}^{(t)},t)}{\operatorname{argmin}} \mathbb{E}\left[D_{KL}\left(q\left(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)},\mathbf{x}^{(0)}\right)\Big|\Big|p\left(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}\right)\right)\right]$$



# Outline

- Motivation: The promise of deep unsupervised learning
- Physical intuition: Diffusion processes and time reversal
- Diffusion probabilistic model: Derivation and experimental results
  - Algorithm
  - Deep convolutional network: Universal function approximator
  - Multiplying distributions: Inputation, denoising, computing posteriors
- Other projects: Inverse Ising, non-equilibrium Monte Carlo, stat. mech. of neural networks

### Use Deep Network as Function Approximator for Images



Jascha Sohl-Dickstein
#### Use Deep Network as Function Approximator for Images



Jascha Sohl-Dickstein

#### Use Deep Network as Function Approximator for Images



Jascha Sohl-Dickstein

### Multiscale Convolution

• Single multi-scale convolutional layer:



Jascha Sohl-Dickstein



Diffusion Probabilistic Models

Jascha Sohl-Dickstein



Jascha Sohl-Dickstein



Jascha Sohl-Dickstein



Jascha Sohl-Dickstein



Jascha Sohl-Dickstein



Jascha Sohl-Dickstein

#### Use Deep Network as Function Approximator for Images



Jascha Sohl-Dickstein



Jascha Sohl-Dickstein



#### Training Data

Jascha Sohl-Dickstein



Training Data



Sample from [Theis *et al*, 2012]

Jascha Sohl-Dickstein



Training Data



Sample from [Theis *et al*, 2012]



Jascha Sohl-Dickstein



Training Data



Sample from [Theis *et al*, 2012]



Sample from diffusion model

Jascha Sohl-Dickstein



Training Data



Sample from [Theis *et al*, 2012]



Sample from diffusion model

Jascha Sohl-Dickstein



Training Data

Jascha Sohl-Dickstein



Training Data

Jascha Sohl-Dickstein



Samples from Generative Adverserial [Goodfellow *et al*, 2014]



Training Data

Samples from DRAW [Gregor *et al*, 2015]

Jascha Sohl-Dickstein



Training Data

Jascha Sohl-Dickstein



Samples from DRAW [Gregor *et al*, 2015]





















Samples from diffusion model

### Outline

- Motivation: The promise of deep unsupervised learning
- Physical intuition: Diffusion processes and time reversal
- Diffusion probabilistic model: Derivation and experimental results
  - Algorithm
  - Deep convolutional network: Universal function approximator
  - Multiplying distributions: Inputation, denoising, computing posteriors
- Other projects: Inverse Ising, non-equilibrium Monte Carlo, stat. mech. of neural networks

Jascha Sohl-Dickstein

Interested in  $\tilde{p}(\mathbf{x}^{(0)}) \propto p(\mathbf{x}^{(0)}) r(\mathbf{x}^{(0)})$ 

- Required to compute posterior distributions
  - Missing data (inpainting)
  - Corrupted data (denoising)

Interested in  $\tilde{p}(\mathbf{x}^{(0)}) \propto p(\mathbf{x}^{(0)}) r(\mathbf{x}^{(0)})$ 

- Required to compute posterior distributions
  - Missing data (inpainting)
  - Corrupted data (denoising)
- Difficult and expensive using competing techniques
  - e.g. VAEs, GSNs, NADEs, GANs, RNVP, most graphical models

Jascha Sohl-Dickstein

Interested in  $\tilde{p}(\mathbf{x}^{(0)}) \propto p(\mathbf{x}^{(0)}) r(\mathbf{x}^{(0)})$ 

Interested in  $\tilde{p}(\mathbf{x}^{(0)}) \propto p(\mathbf{x}^{(0)}) r(\mathbf{x}^{(0)})$ 

Interested in  $\tilde{p}(\mathbf{x}^{(0)}) \propto p(\mathbf{x}^{(0)}) r(\mathbf{x}^{(0)})$ 

Acts as small perturbation to diffusion process

Jascha Sohl-Dickstein

Multiplying Distributions is Straightforward Interested in  $\tilde{p}(\mathbf{x}^{(0)}) \propto p(\mathbf{x}^{(0)}) r(\mathbf{x}^{(0)})$ Acts as small perturbation to diffusion process  $p\left(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}\right) = \mathcal{N}\left(\mathbf{x}^{(t-1)}; f_{\mu}\left(\mathbf{x}^{(t)}, t\right), f_{\Sigma}\left(\mathbf{x}^{(t)}, t\right)\right)$  $\tilde{p}\left(\mathbf{x}^{(t-1)} \mid \mathbf{x}^{(t)}\right) \approx \mathcal{N}\left(x^{(t-1)}; \mathbf{f}_{\mu}\left(\mathbf{x}^{(t)}, t\right) + \mathbf{f}_{\Sigma}\left(\mathbf{x}^{(t)}, t\right) \frac{\partial \log r\left(\mathbf{x}^{(t-1)'}\right)}{\partial \mathbf{x}^{(t-1)'}} \bigg|_{\mathbf{x}^{(t-1)'} = f_{\mu}\left(\mathbf{x}^{(t)}, t\right)}, \mathbf{f}_{\Sigma}\left(\mathbf{x}^{(t)}, t\right)\right)$ 

Jascha Sohl-Dickstein

# Image Denoising by Sampling from Posterior



Holdout Data

Jascha Sohl-Dickstein

# Image Denoising by Sampling from Posterior



Holdout Data

Corrupted (SNR = 1)

Jascha Sohl-Dickstein

# Image Denoising by Sampling from Posterior



Corrupted

(SNR = 1)

Denoised

Jascha Sohl-Dickstein

**Diffusion Probabilistic Models** 

# Image Inpainting by Sampling from Posterior



Inpainted image



True image

Jascha Sohl-Dickstein

# Image Inpainting by Sampling from Posterior



Inpainted image



True image

Jascha Sohl-Dickstein

Jascha Sohl-Dickstein

• Flexible: Diffusion process for any (smooth) distribution

- Flexible: Diffusion process for any (smooth) distribution
  - Binary or continuous state space

- Flexible: Diffusion process for any (smooth) distribution
  - Binary or continuous state space
- Tractable: Training, exact sampling, inference, evaluation
# Flexible and tractable method for deep unsupervised learning

- Flexible: Diffusion process for any (smooth) distribution
  - Binary or continuous state space
- Tractable: Training, exact sampling, inference, evaluation
- Deep networks with thousands of layers (/ time steps)

# Flexible and tractable method for deep unsupervised learning

- Flexible: Diffusion process for any (smooth) distribution
  - Binary or continuous state space
- Tractable: Training, exact sampling, inference, evaluation
- Deep networks with thousands of layers (/ time steps)
- Easy to multiply distributions (e.g. for posterior)

# Flexible and tractable method for deep unsupervised learning

- Flexible: Diffusion process for any (smooth) distribution
  - Binary or continuous state space
- Tractable: Training, exact sampling, inference, evaluation
- Deep networks with thousands of layers (/ time steps)
- Easy to multiply distributions (e.g. for posterior)
- Bounds on entropy production

## Outline

- Motivation: The promise of deep unsupervised learning
- Physical intuition: Diffusion processes and time reversal
- Diffusion probabilistic model: Derivation and experimental results
- Other projects: Inverse Ising, non-equilibrium Monte Carlo, stat. mech. of neural networks

• Train energy based models





[PRL, 2011; ICML, 2011]

**Diffusion Probabilistic Models** 

• Train energy based models

 $\hat{\theta}_{ML} = \operatorname*{argmin}_{\theta} D_{KL} \left( \mathbf{p}^{(\mathbf{0})} || \mathbf{p}^{(\infty)} \left( \theta \right) \right)$ 





[PRL, 2011; ICML, 2011]

**Diffusion Probabilistic Models** 

• Train energy based models

 $\hat{\theta}_{ML} = \operatorname*{argmin}_{\theta} D_{KL} \left( \mathbf{p}^{(\mathbf{0})} || \mathbf{p}^{(\infty)} \left( \theta \right) \right)$ 







[PRL, 2011; ICML, 2011]





 $\hat{\theta}_{ML} = \operatorname{argmin} D_{KL} \left( \mathbf{p}^{(\mathbf{0})} || \mathbf{p}^{(\infty)} \left( \theta \right) \right)$ 



[PRL, 2011; ICML, 2011]

Jascha Sohl-Dickstein



Jascha Sohl-Dickstein



Jascha Sohl-Dickstein

• Train energy based model

 $\mathbf{p^{(0)}}||\mathbf{p^{(\infty)}}\left( heta
ight)
ight)$ 



Progression of learning



 $\theta_{ML} = \operatorname{argmin} D_{KL}$ 

[PRL, 2011; ICML, 2011]

**Diffusion Probabilistic Models** 

• Train energy based model  $\hat{\theta}_{ML} = \underset{\theta}{\operatorname{argmin}} D_{KL} \left( \mathbf{p}^{(\mathbf{0})} || \mathbf{p}^{(\infty)} \left( \theta \right) \right)$ 

$$\hat{\theta}_{MPF} = \operatorname*{argmin}_{\theta} D_{KL} \left( \mathbf{p}^{(\mathbf{0})} || \mathbf{p}^{(\epsilon)} \left( \theta \right) \right)$$

Progression of learning



[PRL, 2011; ICML, 2011]

**Diffusion Probabilistic Models** 

 Train energy based model  $\mathbf{p^{(0)}}||\mathbf{p^{(\infty)}}\left( heta
ight)|$  $\hat{\theta}_{MPF} = \operatorname*{argmin}_{\theta} D_{KL} \left( \mathbf{p}^{(\mathbf{0})} || \mathbf{p}^{(\epsilon)} \left( \theta \right) \right)$  $\theta_{ML} = \operatorname{argmin} D_{KL}$ Progression of learning Intermediate Probability distribution space First learning step Learning complete learning step  $n^{(\infty)}(\theta)$  $p^{(\epsilon)}$  $p^{(\epsilon)}$  $(\infty)_{\ell}$  $p^{(0)}$ 

[PRL, 2011; ICML, 2011]

**Diffusion Probabilistic Models** 

• More rapidly solves inverse Ising problem (estimate Ising couplings from samples)



First 60 seconds

First 800 seconds

First 25,000 seconds

[PRL, 2011; ICML, 2011]

**Diffusion Probabilistic Models** 

• Detailed balance

#### [ICML, 2014]

Jascha Sohl-Dickstein

- Detailed balance
  - + Easy to sample from correct distribution (use any proposal, accept/reject with Metropolis-Hastings)

#### [ICML, 2014]

Jascha Sohl-Dickstein

- Detailed balance
  - + Easy to sample from correct distribution (use any proposal, accept/reject with Metropolis-Hastings)
  - Random walk behavior (by definition, go backwards exactly as often as forwards, explore distribution like sqrt(steps)).

#### [ICML, 2014]

Jascha Sohl-Dickstein

- Detailed balance
  - + Easy to sample from correct distribution (use any proposal, accept/reject with Metropolis-Hastings)
  - Random walk behavior (by definition, go backwards exactly as often as forwards, explore distribution like sqrt(steps)).



Random walk (slow mixing)

#### [ICML, 2014]

Jascha Sohl-Dickstein

- Detailed balance
  - + Easy to sample from correct distribution (use any proposal, accept/reject with Metropolis-Hastings)
  - Random walk behavior (by definition, go backwards exactly as often as forwards, explore distribution like sqrt(steps)).



Random walk (slow mixing)



Asymmetric transitions (faster)

[ICML, 2014]

Jascha Sohl-Dickstein

- Detailed balance
  - + Easy to sample from correct distribution (use any proposal, accept/reject with Metropolis-Hastings)
  - Random walk behavior (by definition, go backwards exactly as often as forwards, explore distribution like sqrt(steps)).

Goal:

 Make Hamiltonian Monte Carlo mix faster by violating detailed balance

#### [ICML, 2014]

Jascha Sohl-Dickstein

• HMC as operators on discrete state space

#### [ICML, 2014]

Jascha Sohl-Dickstein

- HMC as operators on discrete state space
- **Operators:**



#### [ICML, 2014]

Jascha Sohl-Dickstein

- HMC as operators on discrete state space
- **Operators:**





Jascha Sohl-Dickstein

- Greedy leapfrog
  - Try  $\mathbf{L}\zeta$
  - If fail, try  $\mathbf{L}^2 \zeta$  (no rejection)
  - If fail, try  $\mathbf{L}^3 \zeta$  ...



#### [ICML, 2014]

Jascha Sohl-Dickstein

- Greedy leapfrog
  - Try  $\mathbf{L}\zeta$
  - If fail, try  $\mathbf{L}^2 \zeta$  (no rejection)
  - If fail, try  $\mathbf{L}^3 \zeta$  ...
- Tractable, because equilibrium condition now a sum over discrete states, rather than an integral



#### [ICML, 2014]

Jascha Sohl-Dickstein

Improved mixing by violating detailed balance



Jascha Sohl-Dickstein

#### Neural network after random initialization:

$$z_i^l = \sum_j W_{ij}^l y_j^l + b^j \qquad \qquad y_i^{l+1} = \phi(z_i^l)$$
$$W_{ij}^l \sim \mathcal{N}(0, \sigma_w^2/N_{l-1}) \qquad \qquad b_i^l \sim \mathcal{N}(0, \sigma_b^2)$$

[NIPS, 2016] [ICLR, 2017]

Jascha Sohl-Dickstein

#### Neural network after random initialization:

$$\begin{aligned} z_i^l &= \sum_j W_{ij}^l y_j^l + b^j \qquad \qquad y_i^{l+1} = \phi(z_i^l) \\ W_{ij}^l &\sim \mathcal{N}(0, \sigma_w^2/N_{l-1}) \qquad \qquad b_i^l \sim \mathcal{N}(0, \sigma_b^2) \end{aligned}$$

#### Central limit theorem recurrence relation:

$$z_i^l \sim \mathcal{N}(0, q^l)$$
$$q^l = \sigma_w^2 \frac{1}{\sqrt{2\pi}} \int dz e^{-\frac{1}{2}z^2} \phi^2(\sqrt{q^{l-1}}z) + \sigma_b^2$$

Jascha Sohl-Dickstein

**Diffusion Probabilistic Models** 

[NIPS, 2016]

[ICLR, 2017]

[NIPS, 2016] [ICLR, 2017]

Jascha Sohl-Dickstein

#### Phase diagram:



Di

Phase diagram:

#### **Predict trainable depth:**

 $6\xi_{c}$ 



## Thanks!

- Unsupervised Learning using Nonequilibrium Thermodynamics
  - Eric Weiss
  - Niru Maheswaranathan
  - Surya Ganguli
- Minimum Probability Flow
  - Peter Battaglino
  - Michael R. DeWeese
- Hamiltonian Monte Carlo
   without Detailed Balance
  - Mayur Mudigonda
  - Michael R. DeWeese

- Statistical Physics of Deep Networks
  - Sam Schoenholz
  - Ben Poole
  - Jeffrey Pennington
  - Justin Gilmer
  - Surya Ganguli
  - Maithra Raghu
  - Subhaneil Lahiri

### SCRAP SLIDES FROM HERE ON

Jascha Sohl-Dickstein

## Setting Diffusion Rate

Erase constant fraction of stimulus variance each step

$$\beta_t = \frac{1}{T - t + 1}$$

• Can also train  $\beta_t$ 

Jascha Sohl-Dickstein

#### Deep Network Architecture for Diffusion



Jascha Sohl-Dickstein

#### Deep Network Architecture for Diffusion



Jascha Sohl-Dickstein
# Image Inpainting by Sampling from Posterior

• Training data [Lazebnik et al, 2005]



Jascha Sohl-Dickstein

## Diffusion Probabilistic Model Applied to MNLST

Model	Log likelihood estimate*
Stacked CAE	121 ± 1.6 bits
DBN	138 ± 2 bits
Deep GSN	214 ± 1.1 bits
Diffusion	220 ± 1.9 bits
Adversarial net	225 ± 2 bits

\* via Parzen window code from [Goodfellow *et al*, 2014] Jascha Sohl-Dickstein



Jascha Sohl-Dickstein

Continuous time formulation

- Continuous time formulation
- Perturbation around energy based model

- Continuous time formulation
- Perturbation around energy based model
- Binary data (e.g. spike trains)

Jascha Sohl-Dickstein

# Outline

- Other projects: Training energy based models, Monte Carlo, deep learning theory
- Motivation: The promise of deep unsupervised learning
- Physical intuition: Diffusion processes and time reversal
- Diffusion probabilistic model: Derivation and experimental results

### Toy Binary Sequence Learning



Jascha Sohl-Dickstein

# Outline

- Motivation: The promise of deep unsupervised learning
- Physical intuition: Diffusion processes and time reversal
- Diffusion probabilistic model: Derivation and experimental results
  - Algorithm
  - Deep convolutional network: Universal function approximator
  - Multiplying distributions: Inputation, denoising, computing posteriors

Jascha Sohl-Dickstein

• Optimization: Combining SGD and quasi-Newton optimization (SFO optimizer) [IСML 2014]



**Diffusion Probabilistic Models** 

• Optimization: Combining SGD and quasi-Newton optimization (SFO optimizer) [IСML 2014]



Jascha Sohl-Dickstein

• Sampling and evaluation: Hamiltonian Monte Carlo without detailed balance [ICML 2014] and for log likelihood evaluation [Tech Report 2012], fast sampling for natural image models [NIPS 2012]





Jascha Sohl-Dickstein

• Training energy-based models: Minimum Probability Flow learning [ICML 2011] [PRL 2011]



• Model design: capturing dynamics with Lie groups [Under Revision at NECO], bilinear generative models [ICCV 2011]

Horizontal Translation





 Properties of deep networks: Characterization in function space

2d slice through function space for 2-layer network

Jascha Sohl-Dickstein



#### • Online education data



Jascha Sohl-Dickstein

• Medical imaging data [SPIE 2009] [Med Phys 2014]



A. Projection Mammograms





Jascha Sohl-Dickstein

• Neuroscience ele [Neuron 2013]

ata: [PLoS Comp Bio 2014]

#### a) Stimulus frames



b) Example data, 2s of data in 20ms bins





 Human ultrasonic echolocation: Blind assistive device [ТВМЕ 2015]





Jascha Sohl-Dickstein

• Planetary science: multispectral observations [Science 2004a] [Science 2004b]



Jascha Sohl-Dickstein

# Thanks!

#### Collaborators

- Craig Abbey
- Peter Battaglino
- Shaowen Bao
- Matthias Bethge
- Jack Culpepper
- Liberty Hamilton
- Chris Hillar
- Alex Huth
- Kilian Koepsell
- Urs Köster
- Niru Maheswaranathan
- Mayur Mudigonda
- Ben Poole
- Lucas Theis
- Jimmy Wang
- Eric Weiss

#### Jascha Sohl-Dickstein

#### Mentors

- Surya Ganguli
- Bruno Olshausen
- Michael R.
  DeWeese
- James F. Bell III

#### Endless discussion

- The Redwood Center for Theoretical Neuroscience
- The Ganguli Gang



Eric

Weiss



Niru Maheswaranathan



Surya Ganguli

### Differences from Variational Autoencoders

- Can analytically evaluate KL divergence between steps in forward and reverse trajectories.
- Can multiply with other distributions, and compute posteriors
- Erases structure, rather than transforming it
- Thousands of layers or time steps, rather than only a small handful
- Connections to nonequilibrium statistical mechanics

Jascha Sohl-Dickstein

## Continuous Time

$$q\left(\mathbf{x}^{t}|\mathbf{x}^{0},\mathbf{x}^{t+dt}\right) = \mathcal{N}\left(\mathbf{x}^{t};\mathbf{x}^{t+dt} - \mathbf{x}^{t+dt}\frac{\exp\left(-\beta t\right)}{1 - \exp\left(-\beta t\right)}\beta dt - \frac{1}{2}\mathbf{x}^{t+dt}\beta dt + \frac{1}{2}\mathbf{x}^{0}\operatorname{csch}\left(\frac{\beta t}{2}\right)\beta dt,\beta dt\right)$$
$$p\left(\mathbf{x}^{t}|\mathbf{x}^{t+dt}\right) = \mathcal{N}\left(\mathbf{x}^{t};\mathbf{x}^{t+dt} - \mathbf{x}^{t+dt}\frac{\exp\left(-\beta t\right)}{1 - \exp\left(-\beta t\right)}\beta dt - \frac{1}{2}\mathbf{x}^{t+dt}\beta dt + \frac{1}{2}f_{0}\left(\mathbf{x}^{t+dt},t\right)\operatorname{csch}\left(\frac{\beta t}{2}\right)\beta dt,\beta dt\right)$$

$$D_{KL}\left(q\left(\mathbf{x}^{t}|\mathbf{x}^{0},\mathbf{x}^{t+dt}\right)||p\left(\mathbf{x}^{t}|\mathbf{x}^{t+dt}\right)\right) = \frac{1}{2}\frac{\Sigma_{q}}{\Sigma_{p}} + \frac{1}{2}\log\frac{\Sigma_{p}}{\Sigma_{q}} + \frac{1}{2}\frac{\left(\mu_{p}-\mu_{q}\right)^{2}}{\Sigma_{p}} - \frac{1}{2}$$
$$= \frac{1}{8}\left(f_{0}\left(\mathbf{x}^{t+dt},t\right)-\mathbf{x}^{0}\right)^{2}\operatorname{csch}^{2}\left(\frac{\beta t}{2}\right)\beta dt$$

#### Denoising autoencoder penalty

**Diffusion Probabilistic Models** 

# Related Methods

- Generative stochastic networks
- Variational autoencoders
- (Deep) (Recurrent) Neural Autoregressive Distribution Estimators

- Variational Bayesian(e.g. variational autoencoder)
  - Posterior over intermediate layers has analytic form > KL divergence has analytic form
  - Can multiply distributions
  - Generative model is small perturbation around inference model makes learning easier
  - Models have *thousands* of layers (or time steps)

Jascha Sohl-Dickstein