

Deep Unsupervised Learning using Nonequilibrium Thermodynamics

Jascha Sohl-Dickstein¹, Eric Weiss²,
Niru Maheswaranathan³, Surya Ganguli³

¹ Google Brain, ² UC Berkeley, ³ Stanford University



Eric



Niru



Surya

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

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

The Promise of Deep Unsupervised Learning

The Promise of Deep Unsupervised Learning

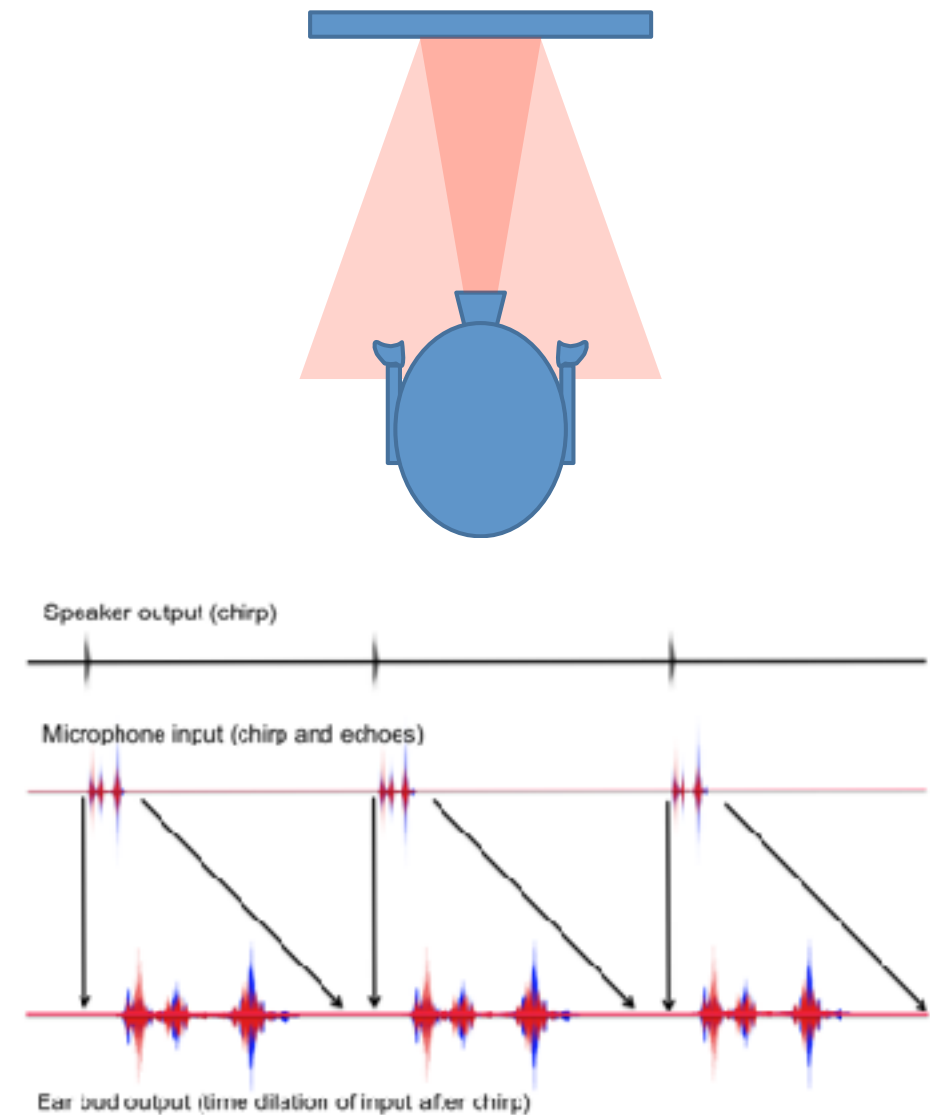
- Unknown features/labels

The Promise of Deep Unsupervised Learning

- Unknown features/labels
 - Novel modalities

The Promise of Deep Unsupervised Learning

- Unknown features/labels
- Novel modalities



[Trans Biomed Eng, 2015]

The Promise of Deep Unsupervised Learning

- Unknown features/labels
 - Novel modalities

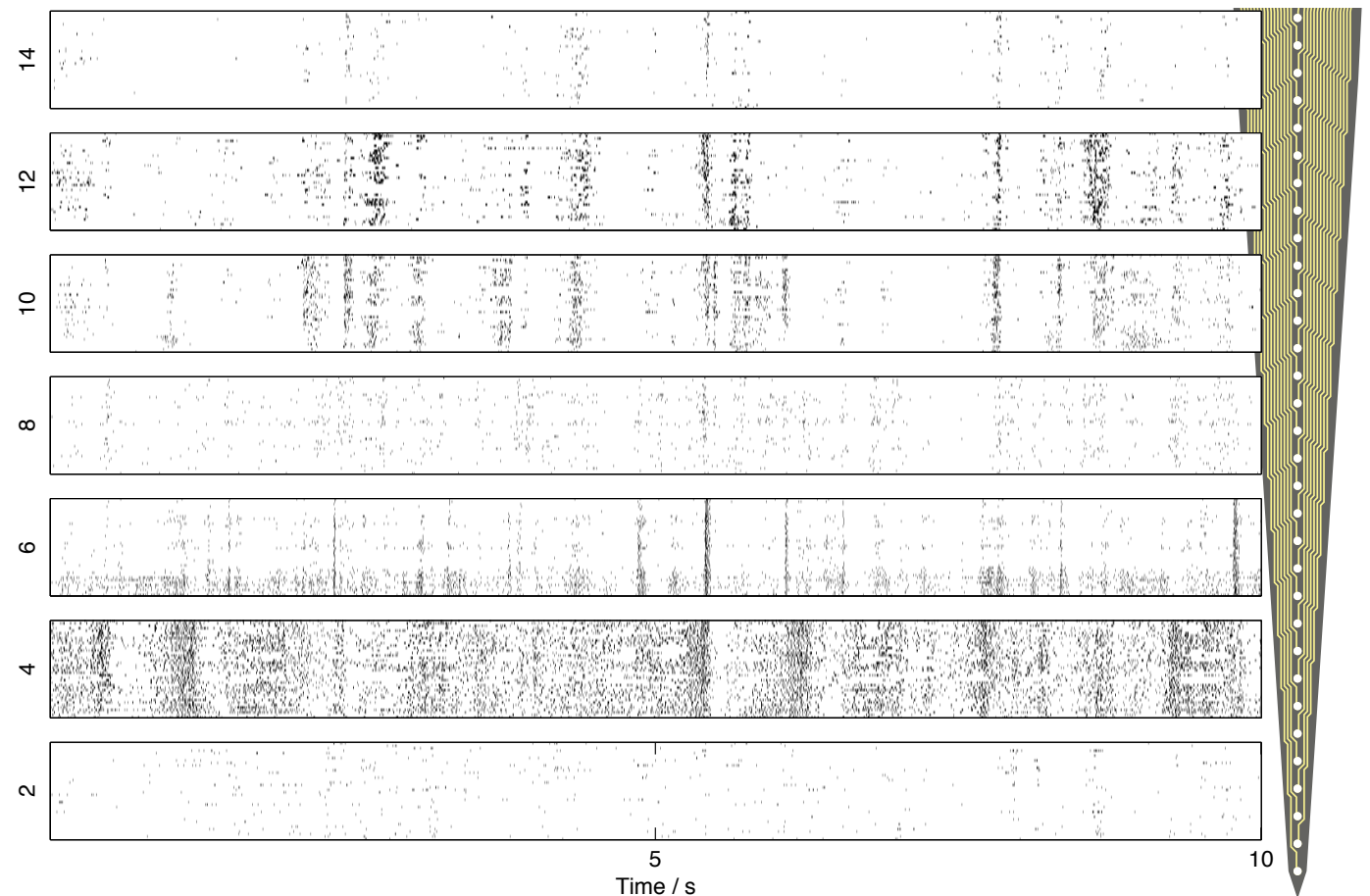
The Promise of Deep Unsupervised Learning

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

The Promise of Deep Unsupervised Learning

- 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]

The Promise of Deep Unsupervised Learning

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

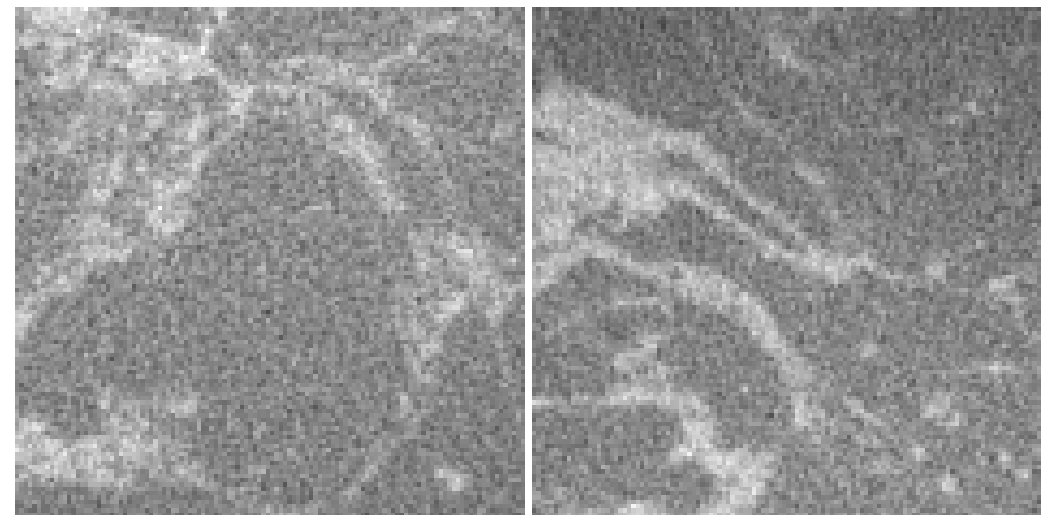
The Promise of Deep Unsupervised Learning

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

The Promise of Deep Unsupervised Learning

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

Coronal breast CT



[SPIE 2009] [Med Phys 2014]

The Promise of Deep Unsupervised Learning

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

The Promise of Deep Unsupervised Learning

- 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

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

Observation 1: Diffusion Destroys Structure



- Dye density represents probability density

Observation 1: Diffusion Destroys Structure



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

Observation 1: Diffusion Destroys Structure



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

Observation 1: Diffusion Destroys Structure



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

Observation 1: Diffusion Destroys Structure



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

Data distribution



Uniform distribution

Core Idea: Recover Structure by Reversing Time



- What if we could reverse time?

Core Idea: Recover Structure by Reversing Time



- What if we could reverse time?

Core Idea: Recover Structure by Reversing Time



- What if we could reverse time?

Data distribution



Uniform distribution

Core Idea: Recover Structure by Reversing Time



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

Data distribution



Uniform distribution

Core Idea: Recover Structure by Reversing Time



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

Data distribution



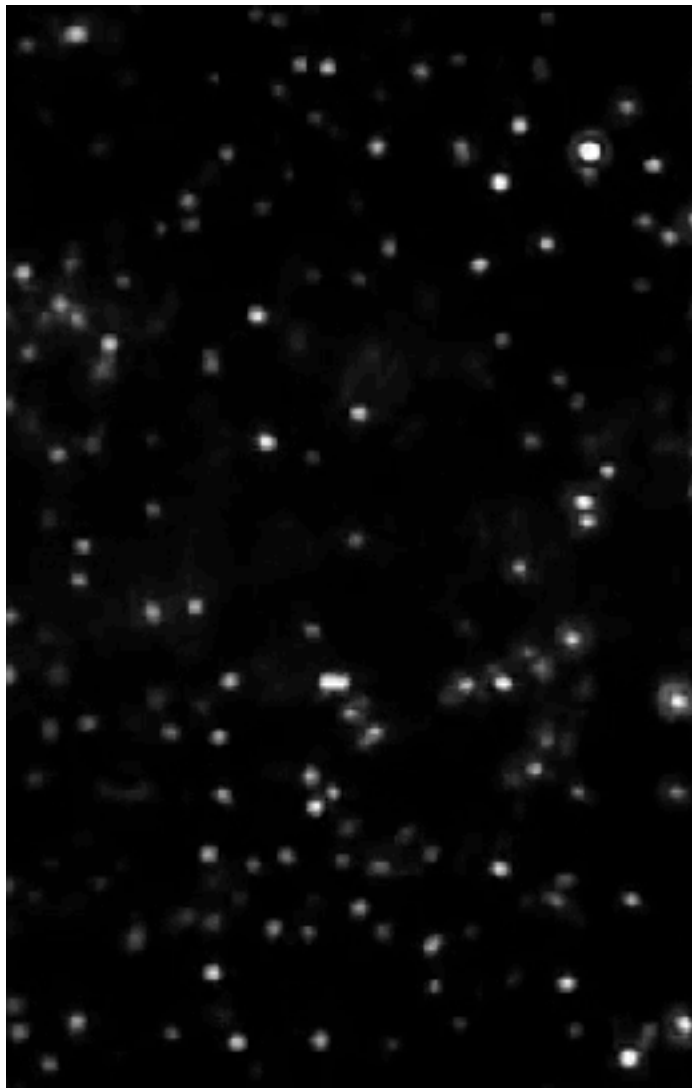
Uniform distribution

Core Idea: Recover Structure by Reversing Time



Core Idea: Recover Structure by Reversing Time

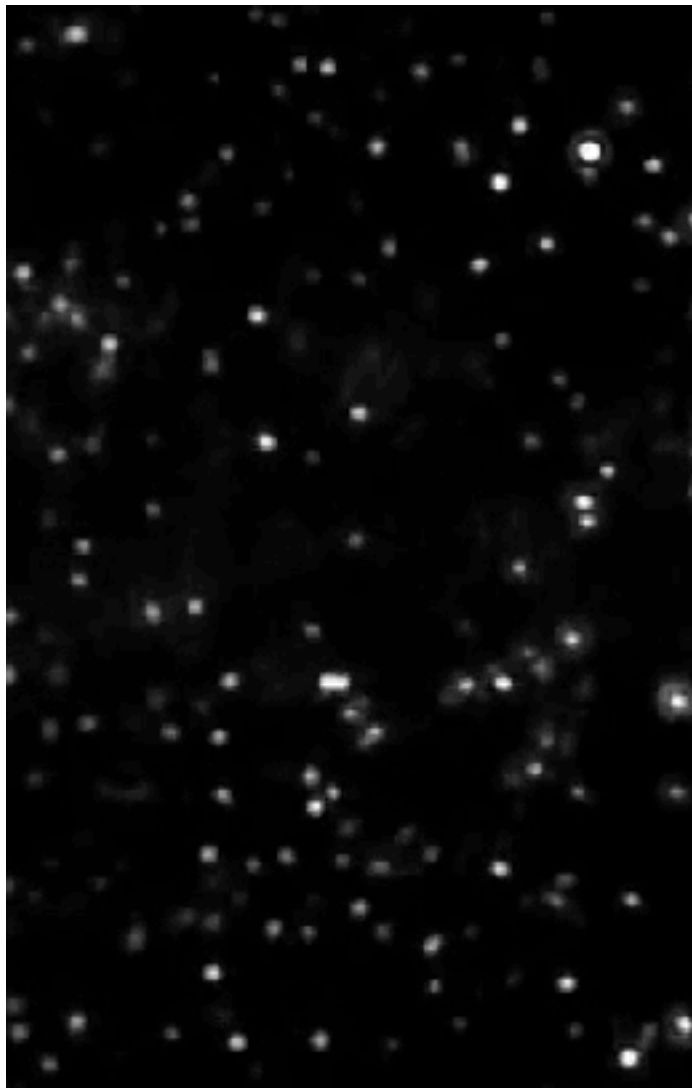
Observation 2: Microscopic Diffusion is Time Reversible



© Rutger Saly

- Microscopic view
- Brownian motion

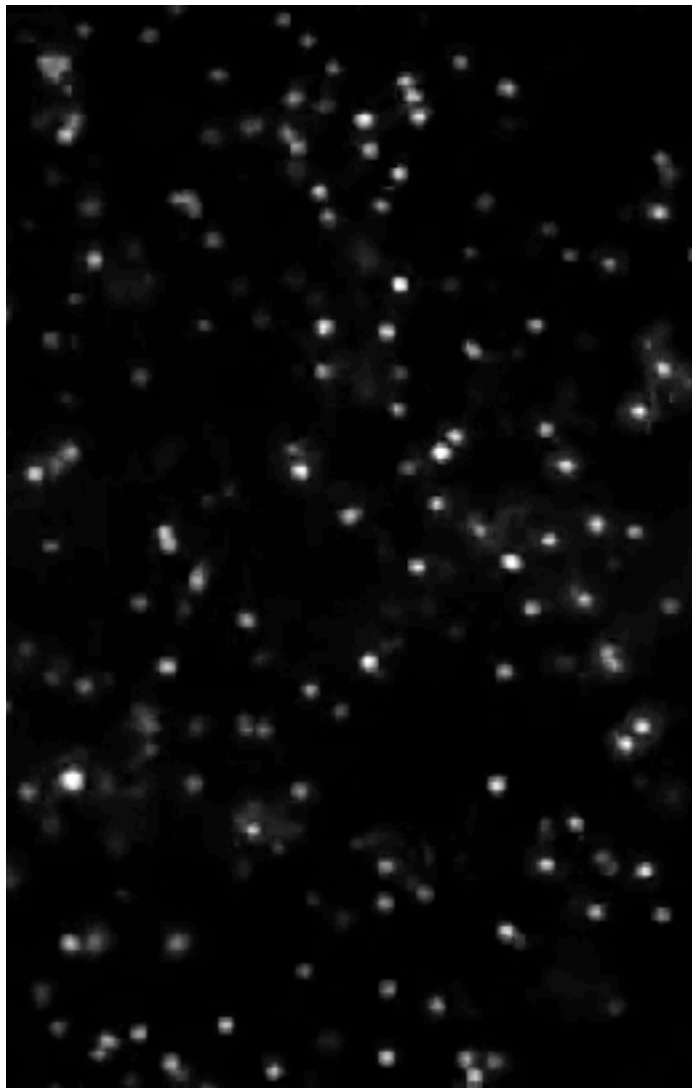
Observation 2: Microscopic Diffusion is Time Reversible



© Rutger Saly

- Microscopic view
- Brownian motion

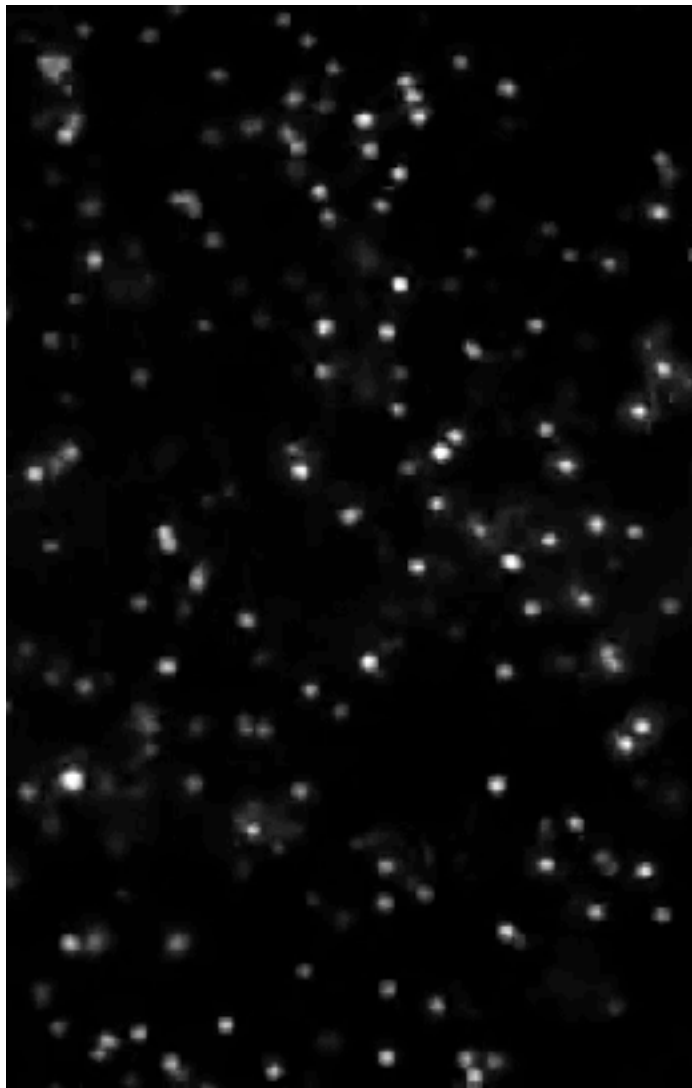
Observation 2: Microscopic Diffusion is Time Reversible



© Rutger Saly

- Microscopic view
- Brownian motion

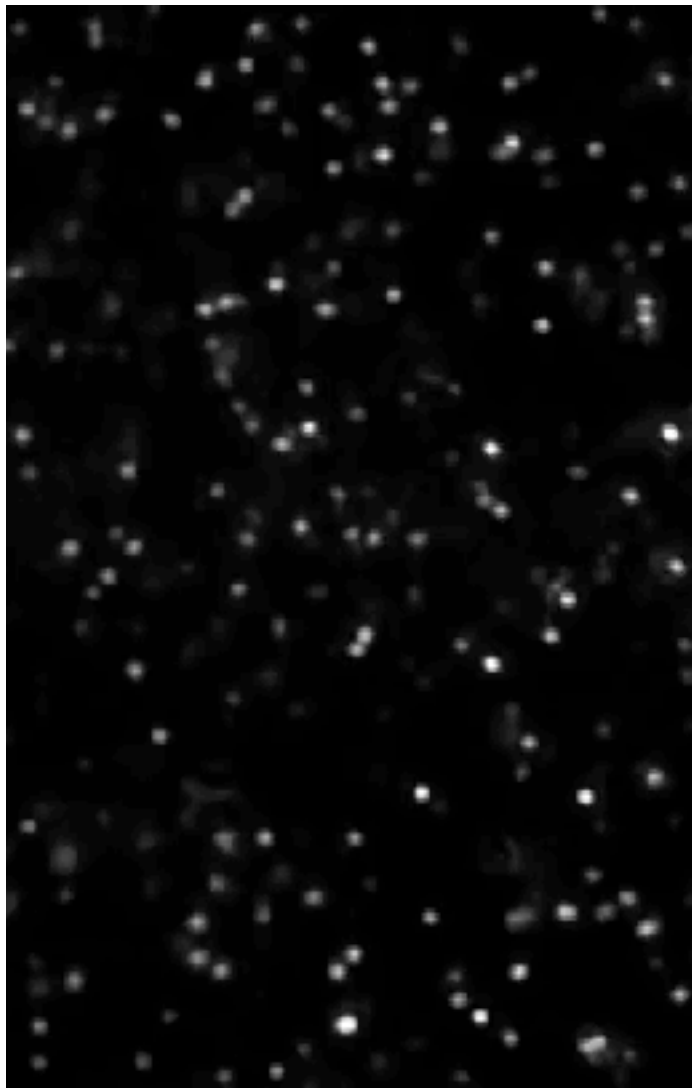
Observation 2: Microscopic Diffusion is Time Reversible



© Rutger Saly

- Microscopic view
- Brownian motion

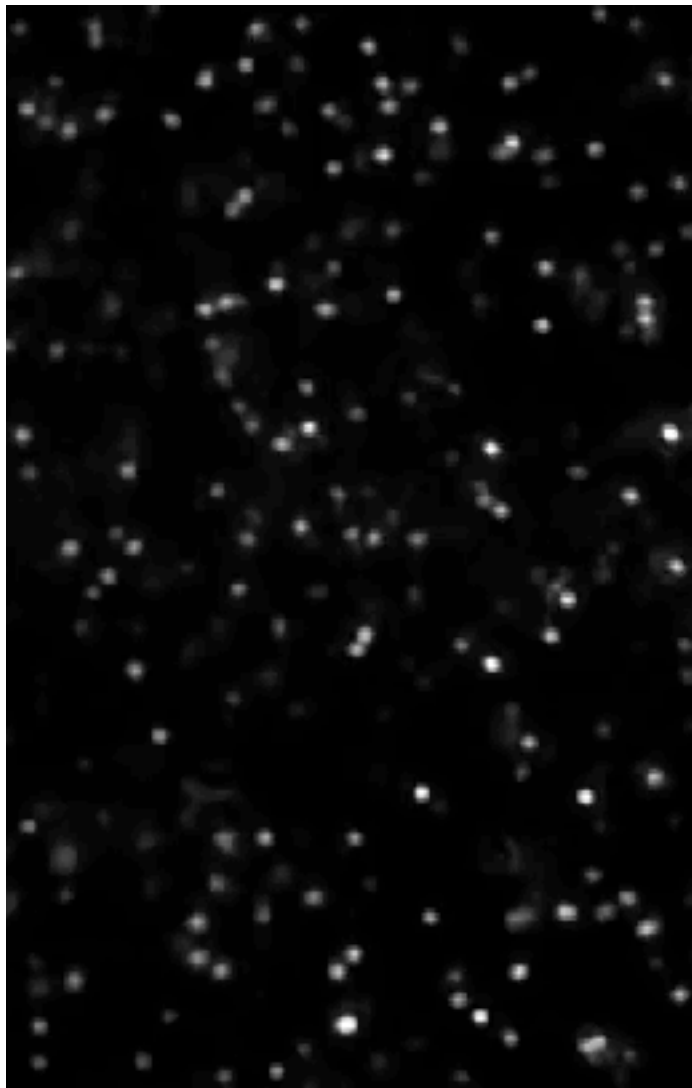
Observation 2: Microscopic Diffusion is Time Reversible



© Rutger Saly

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

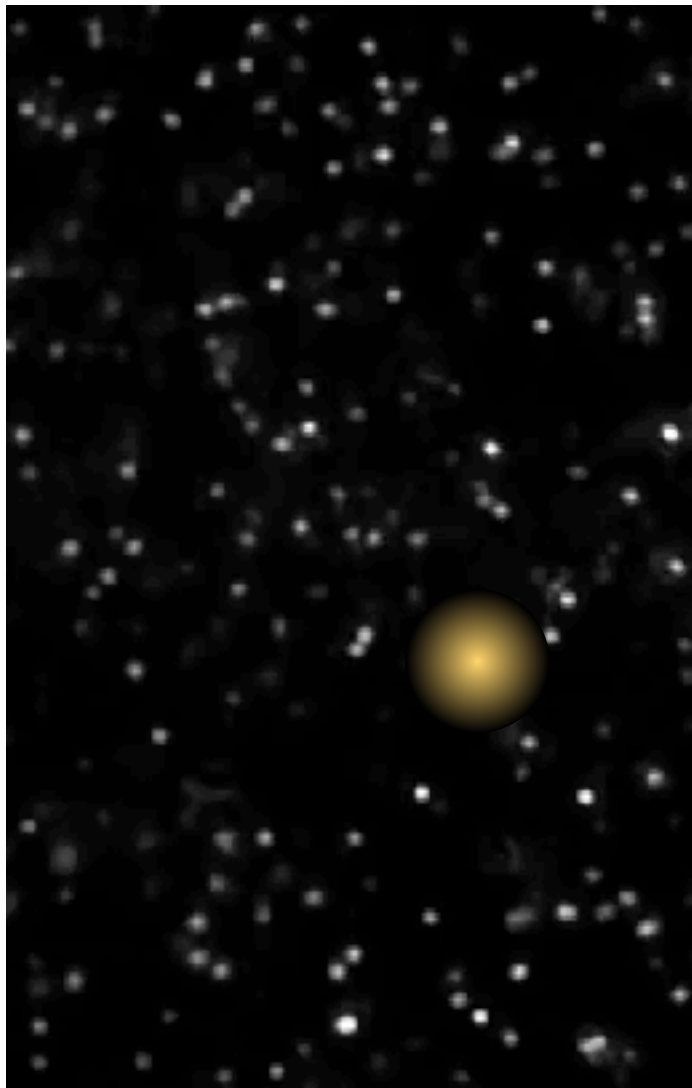
Observation 2: Microscopic Diffusion is Time Reversible



© Rutger Saly

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

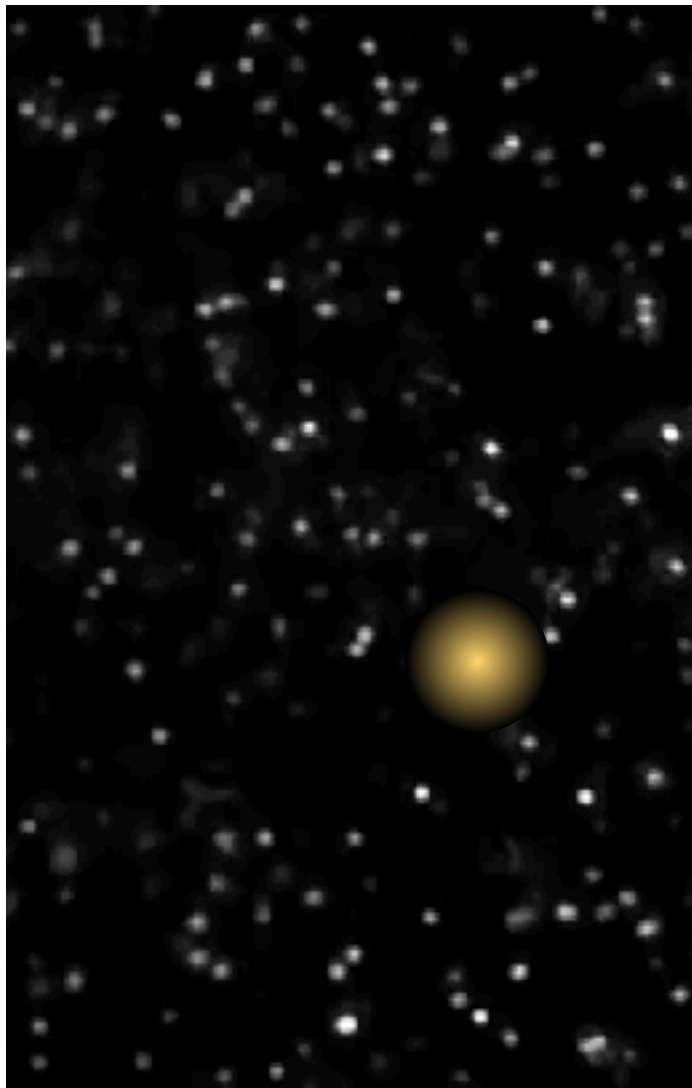
Observation 2: Microscopic Diffusion is Time Reversible



© Rutger Saly

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

Observation 2: Microscopic Diffusion is Time Reversible



© Rutger Saly

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

Overview of Diffusion-based Probabilistic Models

Overview of Diffusion-based Probabilistic Models

- Destroy all structure in data distribution using diffusion process

Overview of Diffusion-based Probabilistic Models

- 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)

Overview of Diffusion-based Probabilistic Models

- 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

Destroy All Structure in Data using Diffusion Process

Data
distribution

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

Destroy All Structure in Data using Diffusion Process

Data
distribution

Forward
diffusion

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



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

Destroy All Structure in Data using Diffusion Process

Data
distribution

Forward
diffusion

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



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

Decay towards origin

Destroy All Structure in Data using Diffusion Process

Data
distribution

Forward
diffusion

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



$$q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)}) = \mathcal{N}(\mathbf{x}^{(t)}; \mathbf{x}^{(t-1)} \underbrace{\sqrt{1 - \beta_t}}_{\text{Decay towards origin}}, \underbrace{\mathbf{I}\beta_t}_{\text{Add small noise}})$$

Decay towards origin

Add small noise

Destroy All Structure in Data using Diffusion Process

Data
distribution

Forward
diffusion

Noise
distribution

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



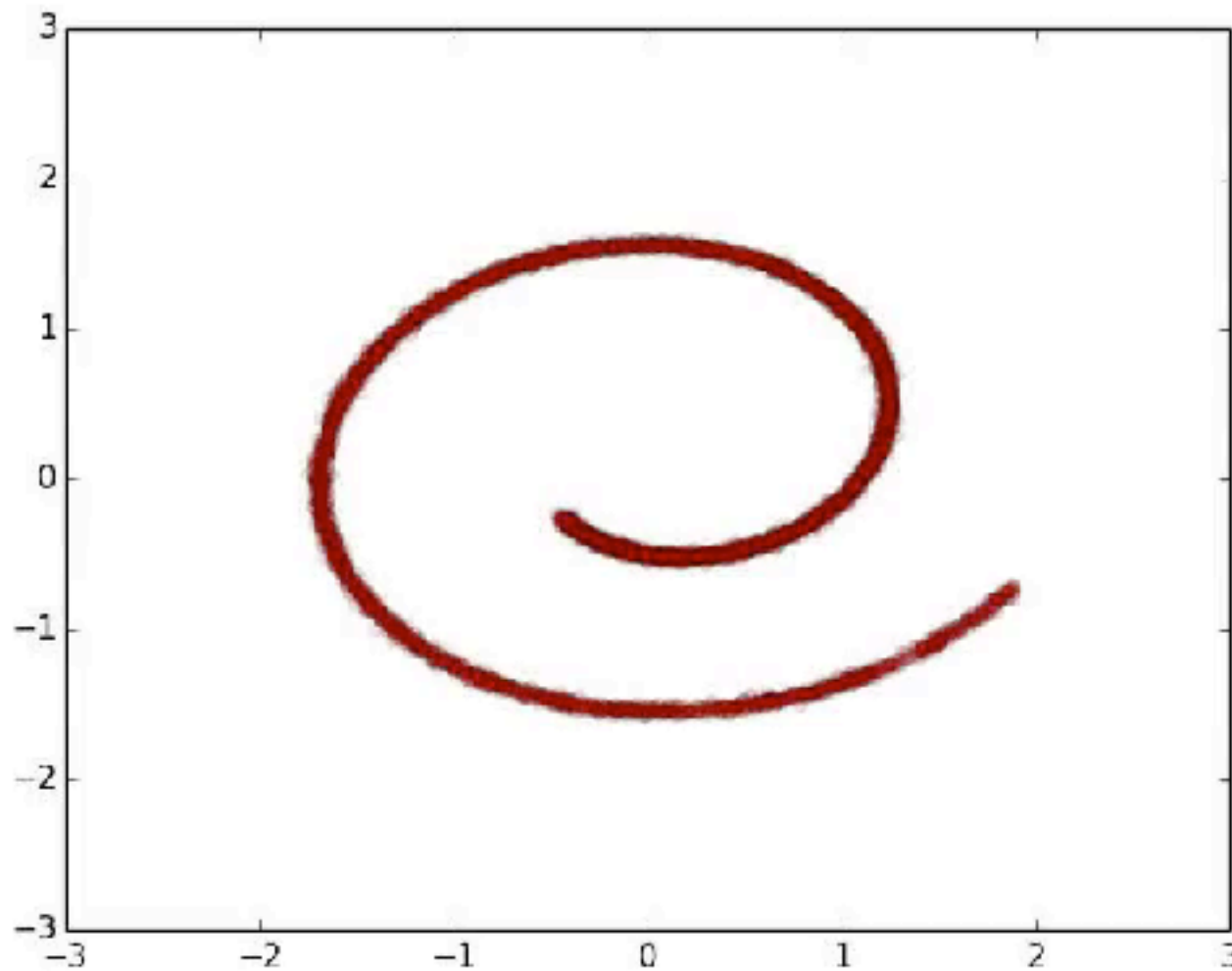
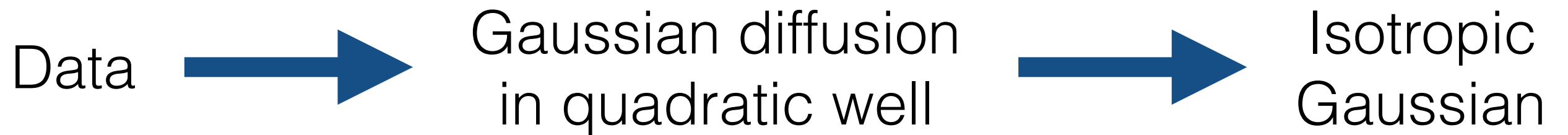
$$q(\mathbf{x}^{(T)}) \approx \mathcal{N}(\mathbf{x}^{(T)}; 0, \mathbf{I})$$

$$q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)}) = \mathcal{N}(\mathbf{x}^{(t)}; \underbrace{\mathbf{x}^{(t-1)} \sqrt{1 - \beta_t}}_{\text{Decay towards origin}}, \underbrace{\mathbf{I} \beta_t}_{\text{Add small noise}})$$

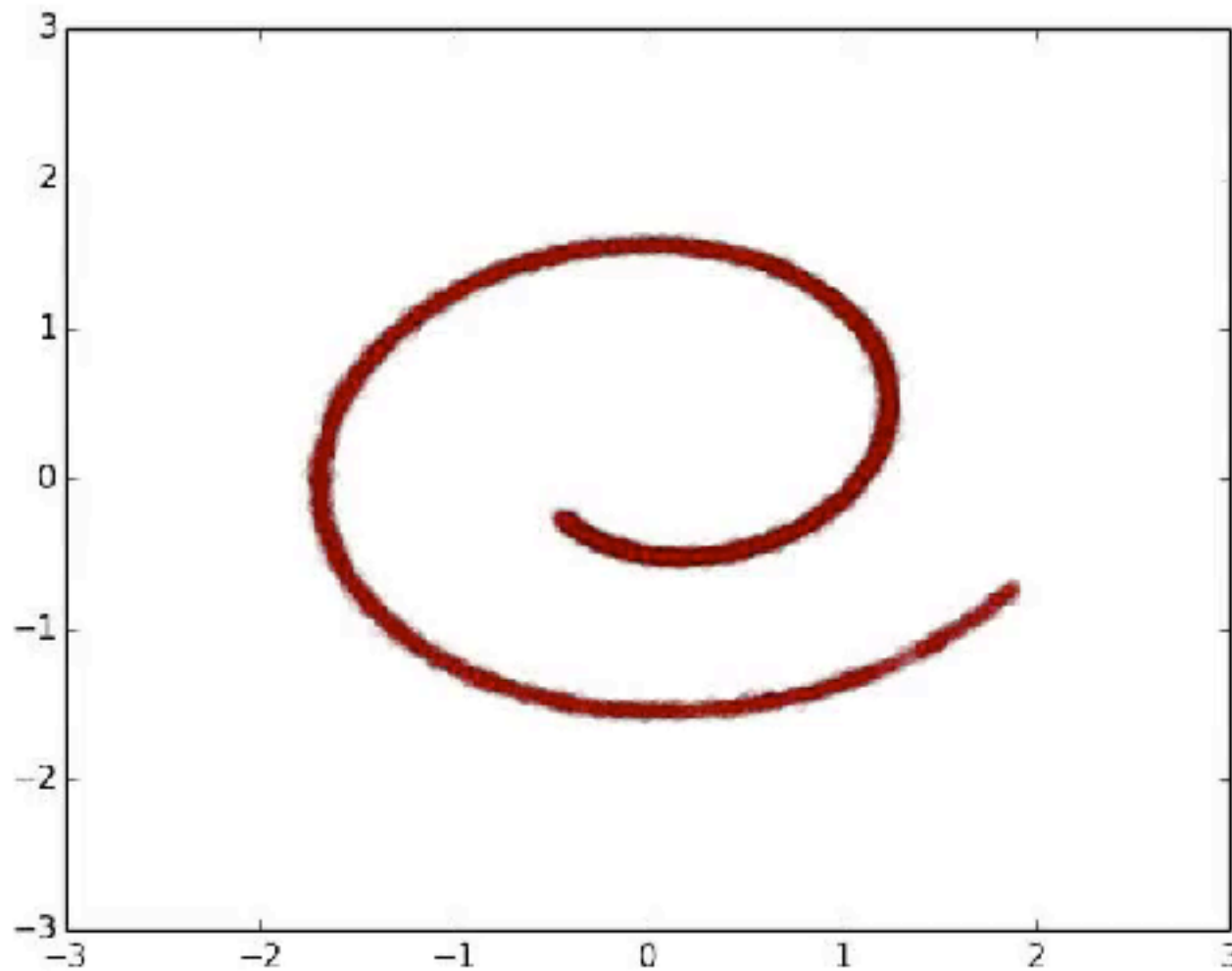
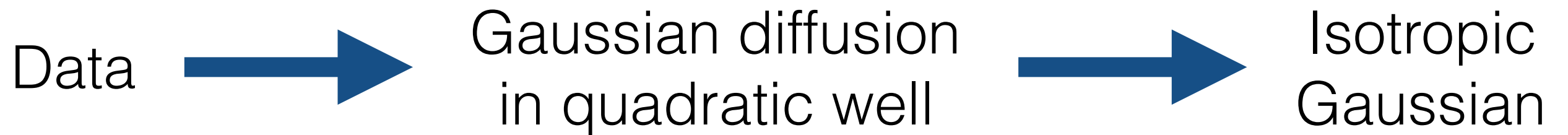
Decay towards origin

Add small noise

Forward Diffusion Process on Swiss Roll



Forward Diffusion Process on Swiss Roll



Recover Structure in Data using Reversal of Diffusion Process

Noise
distribution

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

Recover Structure in Data using Reversal of Diffusion Process

Reverse
diffusion

Noise
distribution



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

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

Recover Structure in Data using Reversal of Diffusion Process

Reverse
diffusion

Noise
distribution



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

$$p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}) = \mathcal{N}(\mathbf{x}^{(t-1)}; \underbrace{f_{\mu}(\mathbf{x}^{(t)}, t), f_{\Sigma}(\mathbf{x}^{(t)}, t)}_{\text{Learned drift and covariance functions}})$$

Learned drift and covariance functions

Recover Structure in Data using Reversal of Diffusion Process

Data
distribution

Reverse
diffusion

Noise
distribution

$$p(\mathbf{x}^{(0)}) \approx q(\mathbf{x}^{(0)})$$



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

$$p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}) = \mathcal{N}(\mathbf{x}^{(t-1)}; \underbrace{f_{\mu}(\mathbf{x}^{(t)}, t), f_{\Sigma}(\mathbf{x}^{(t)}, t)}_{\text{Learned drift and covariance functions}})$$

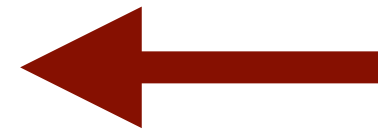
Learned drift and covariance functions

Learned Reverse Diffusion Process on Swiss Roll

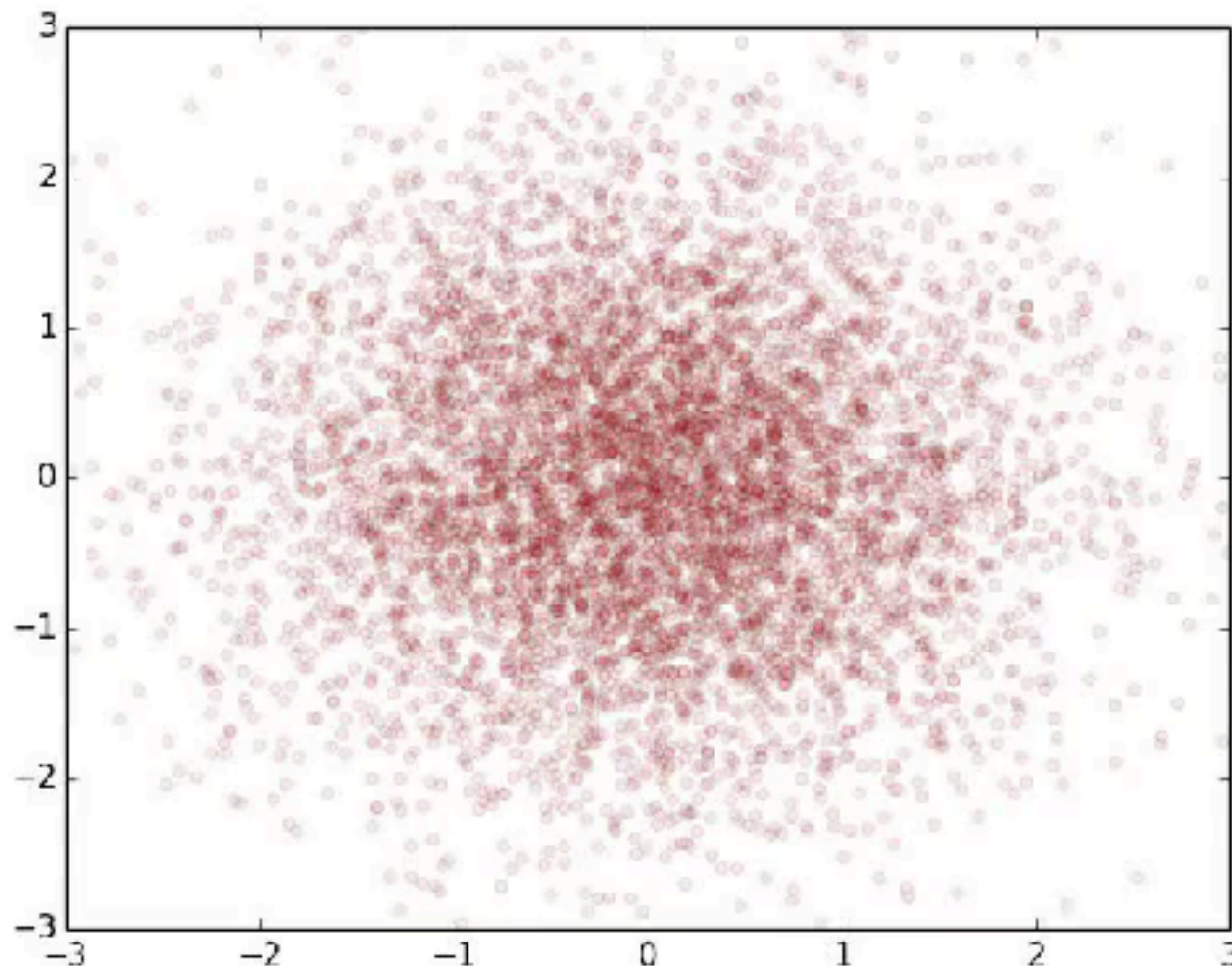
Data
dist.



Gaussian diffusion
w/ learned kernel



Isotropic
Gaussian



Learned Reverse Diffusion Process on Swiss Roll

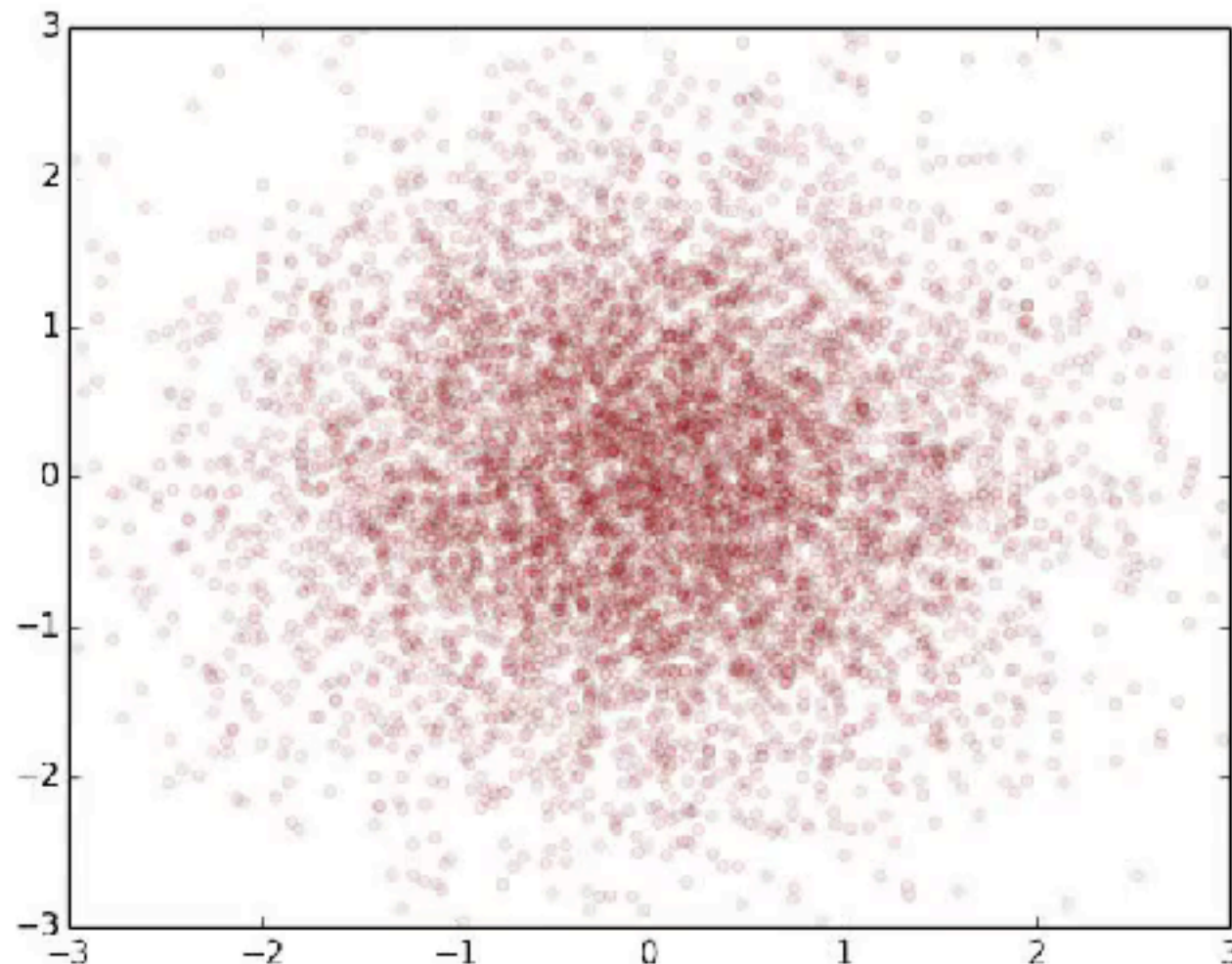
Data
dist.



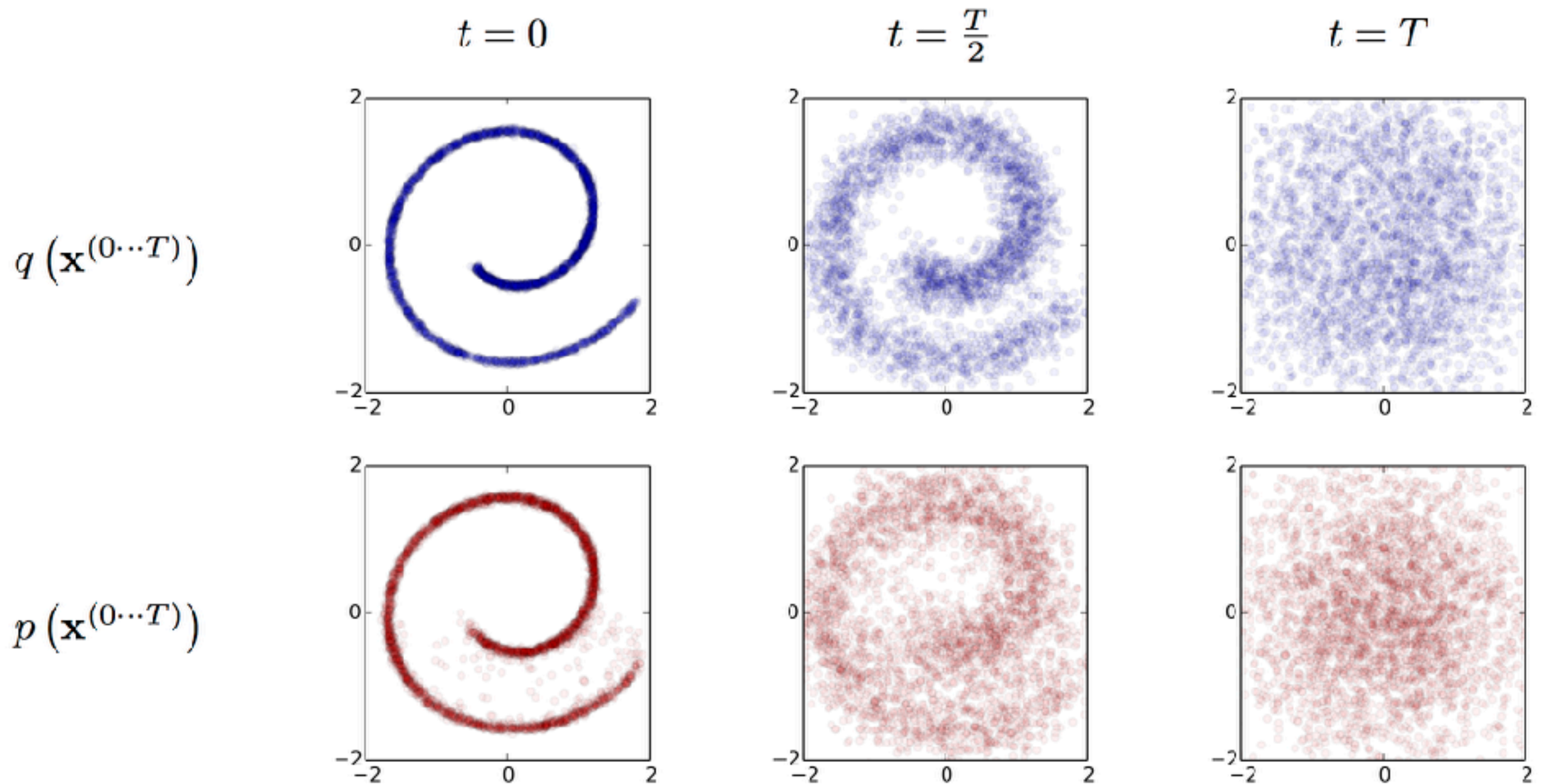
Gaussian diffusion
w/ learned kernel



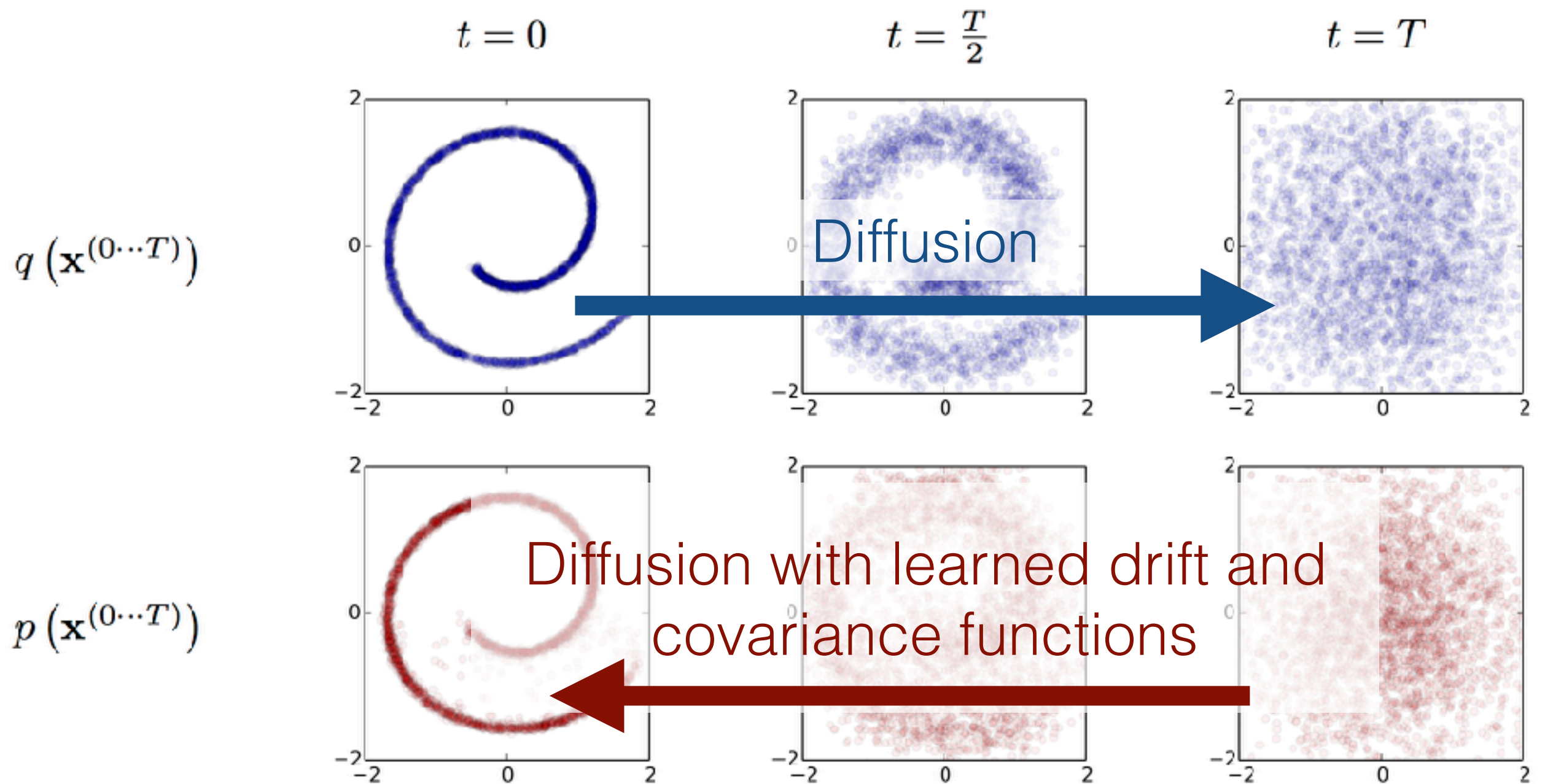
Isotropic
Gaussian



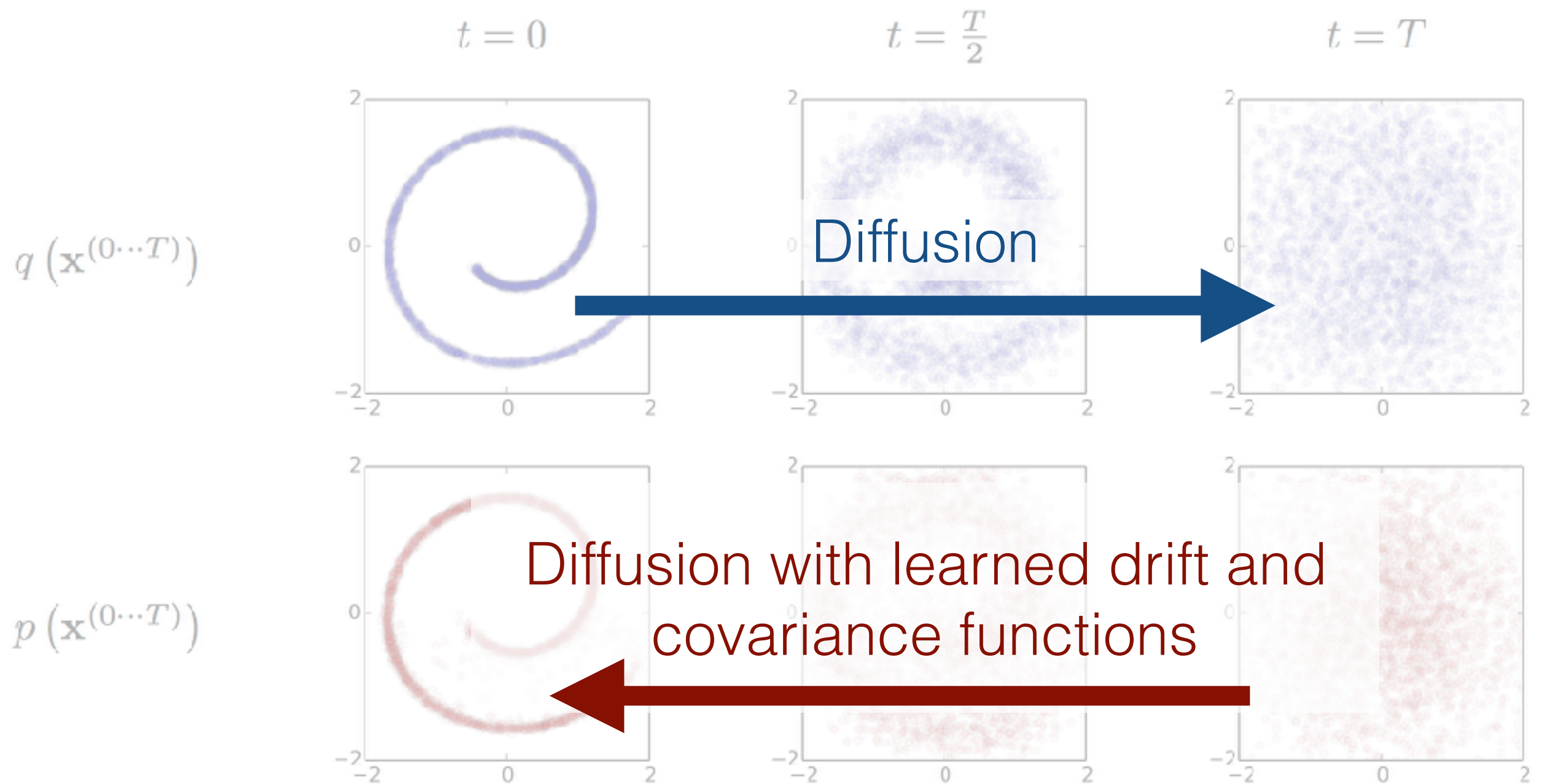
Summary of Forward and Reverse Diffusion on Swiss Roll



Summary of Forward and Reverse Diffusion on Swiss Roll



Summary of Forward and Reverse Diffusion on Swiss Roll



Training the Reverse Diffusion Process

Model probability

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

Training the Reverse Diffusion Process

Model probability

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

Annealed importance sampling / Jarzynski equality

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

Training the Reverse Diffusion Process

Model probability

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

Annealed importance sampling / Jarzynski equality

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

Log Likelihood

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

Model probability

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

Annealed importance sampling / Jarzynski equality

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

Log Likelihood

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

Jensen's inequality

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

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

Log Likelihood

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

Jensen's inequality

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

... algebra ...

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

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

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



Gaussian


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

+ const




Gaussian

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



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

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


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

Training

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

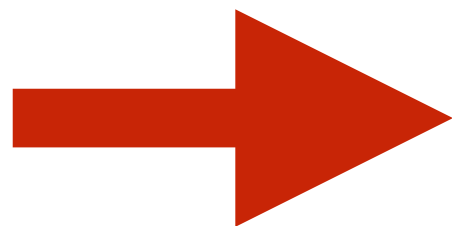
$$L \geq - \sum_{t=2}^T \int d\mathbf{x}^{(0)} d\mathbf{x}^{(t)} q(\mathbf{x}^{(0)}, \mathbf{x}^{(t)}) D_{KL} \left(q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}) \parallel p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}) \right) + \text{const}$$

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

Training

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

Unsupervised
learning

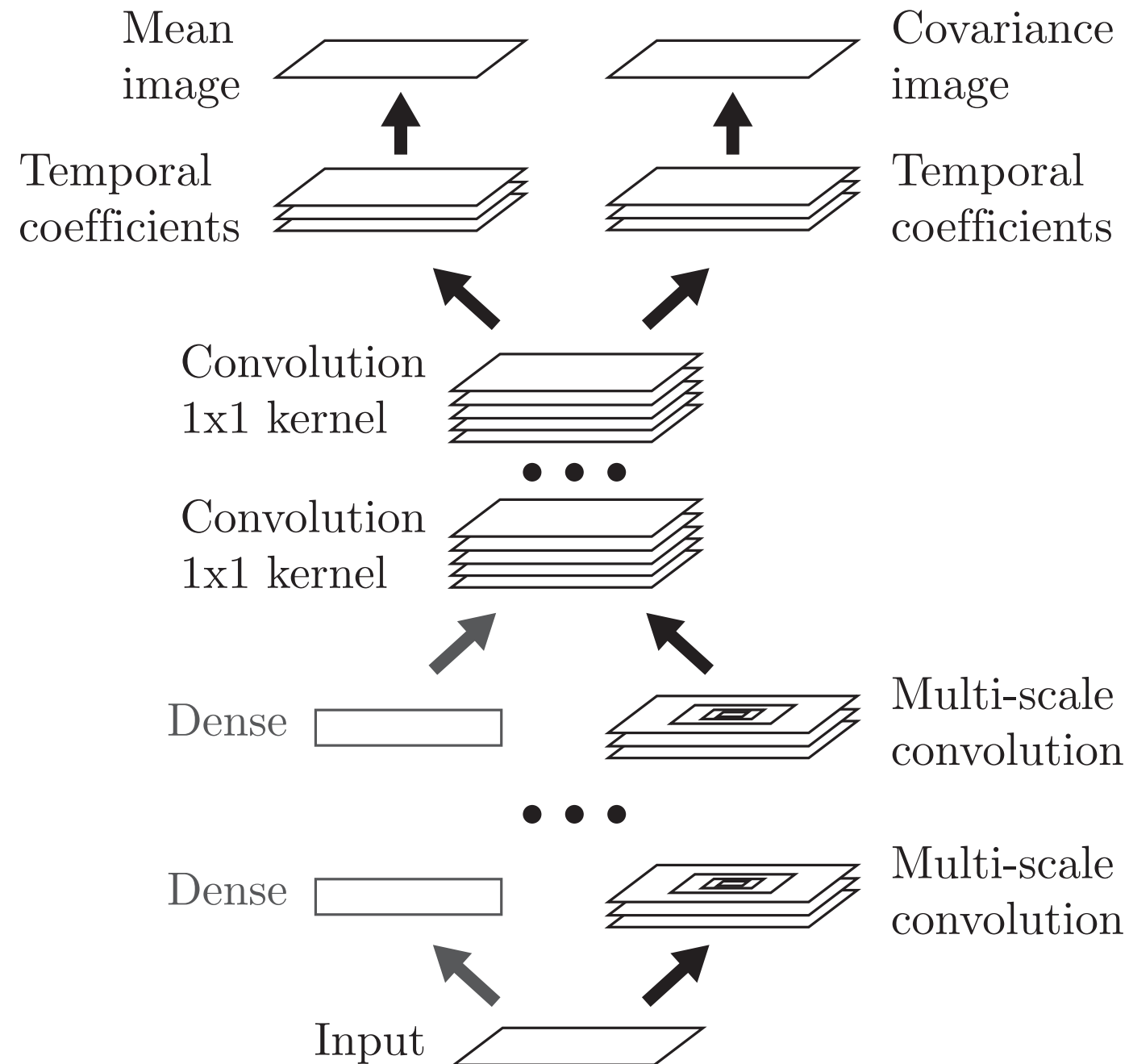


Regression

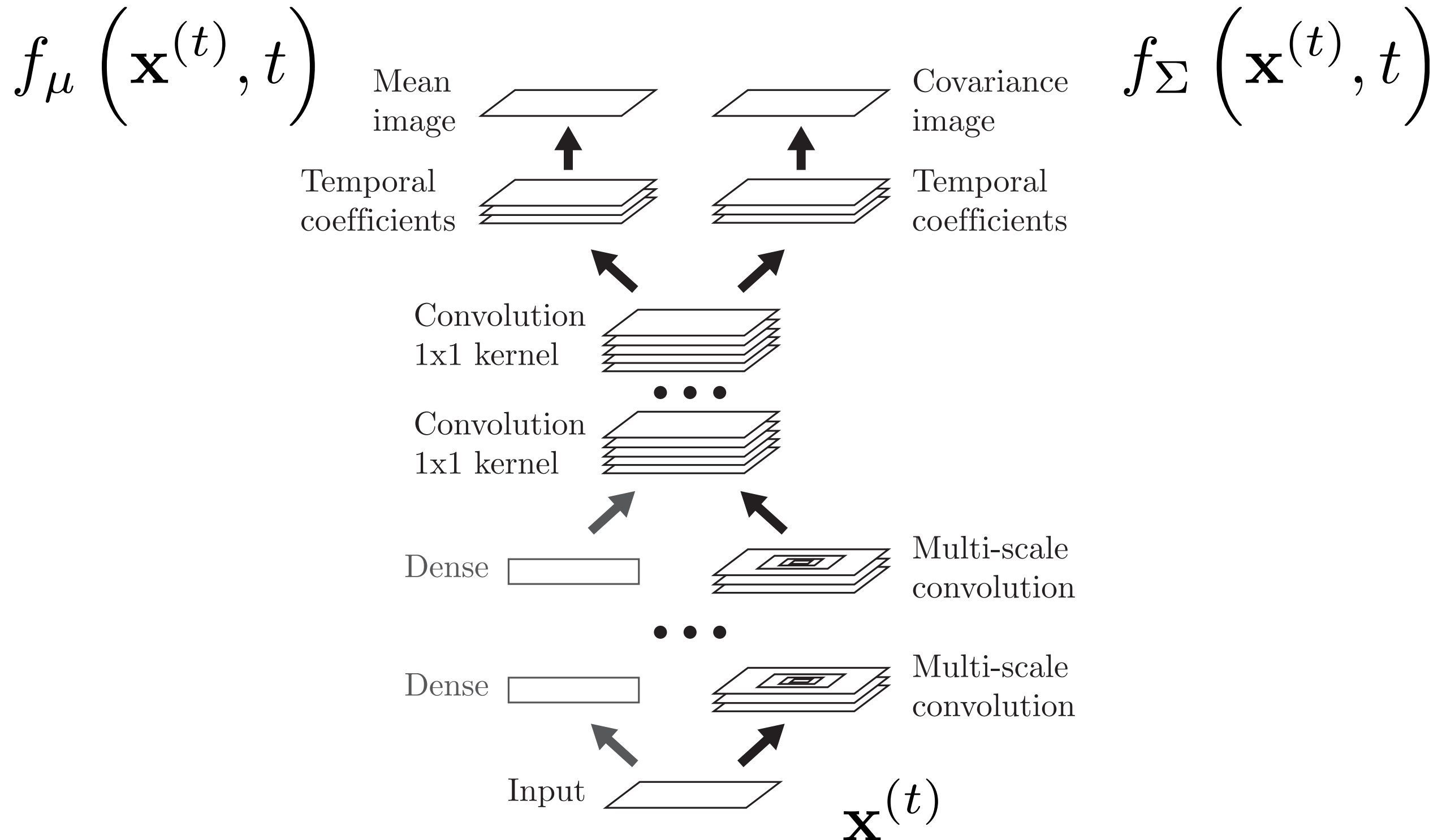
Outline

- **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



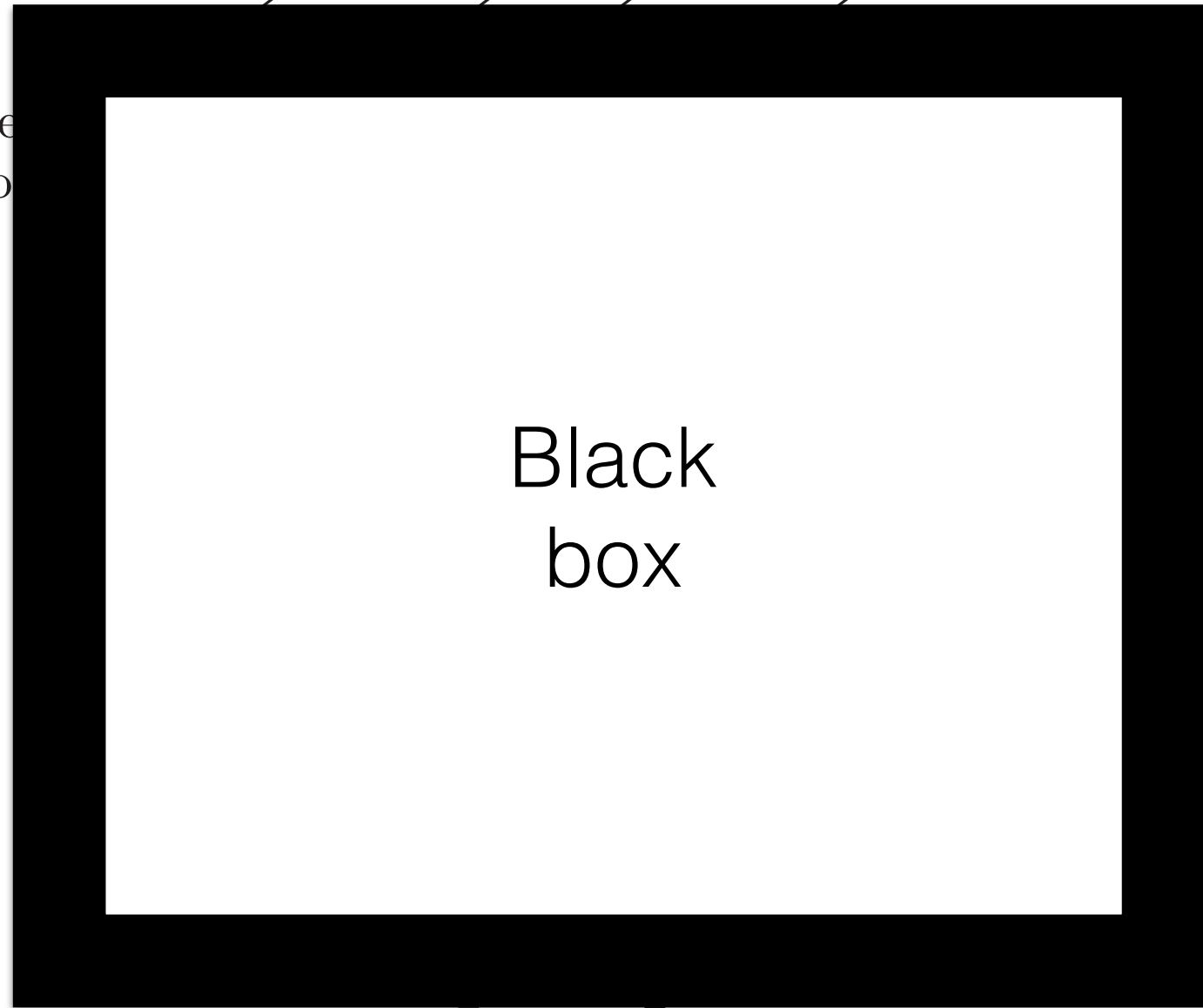
Use Deep Network as Function Approximator for Images



Use Deep Network as Function Approximator for Images

$$f_{\mu} \left(\mathbf{x}^{(t)}, t \right) \quad \text{Mean} \quad \text{Covariance} \quad f_{\Sigma} \left(\mathbf{x}^{(t)}, t \right)$$

Te
co

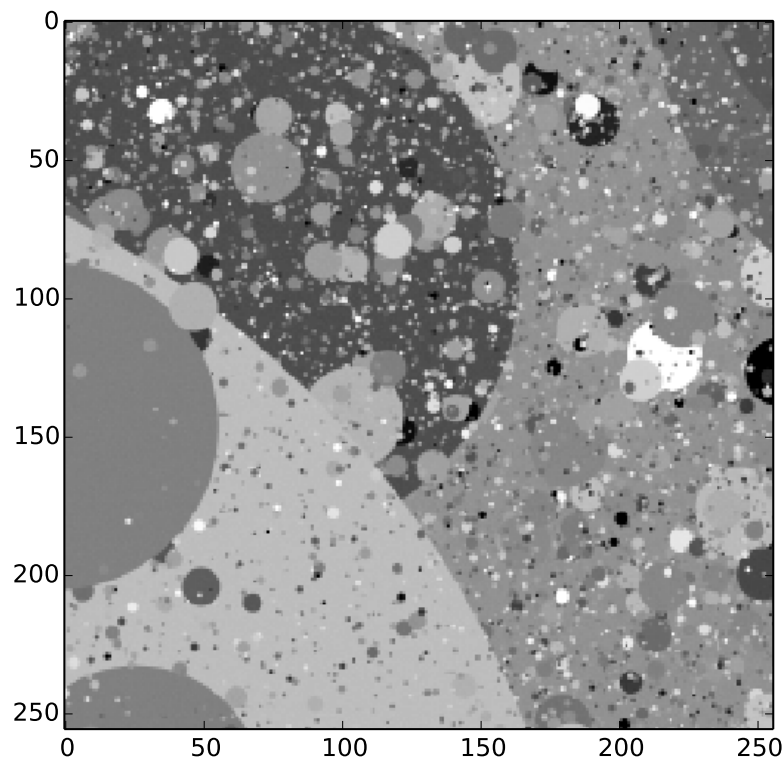


Black
box

Input

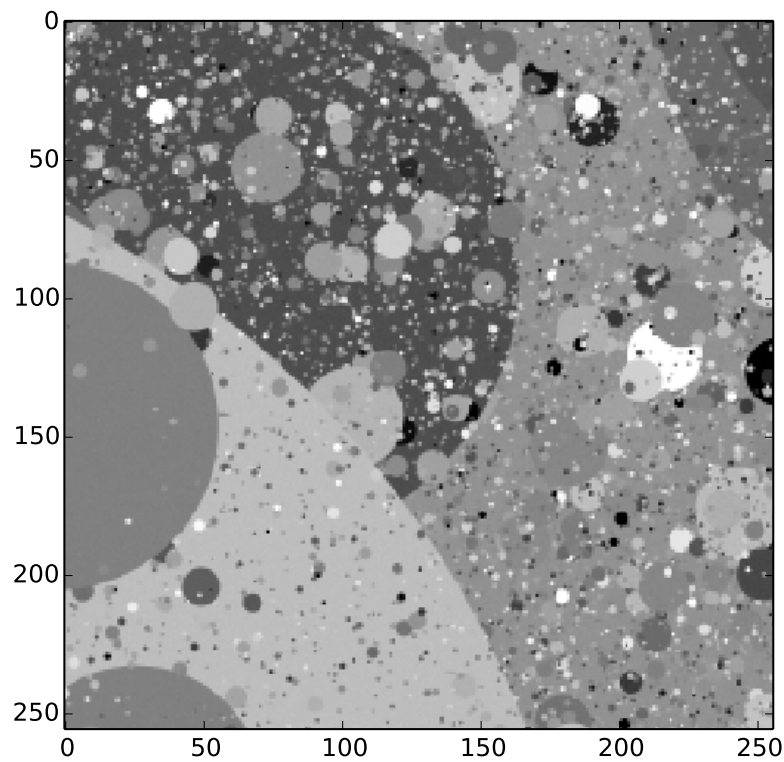
$\mathbf{x}^{(t)}$

Diffusion Probabilistic Model Applied to Dead Leaves

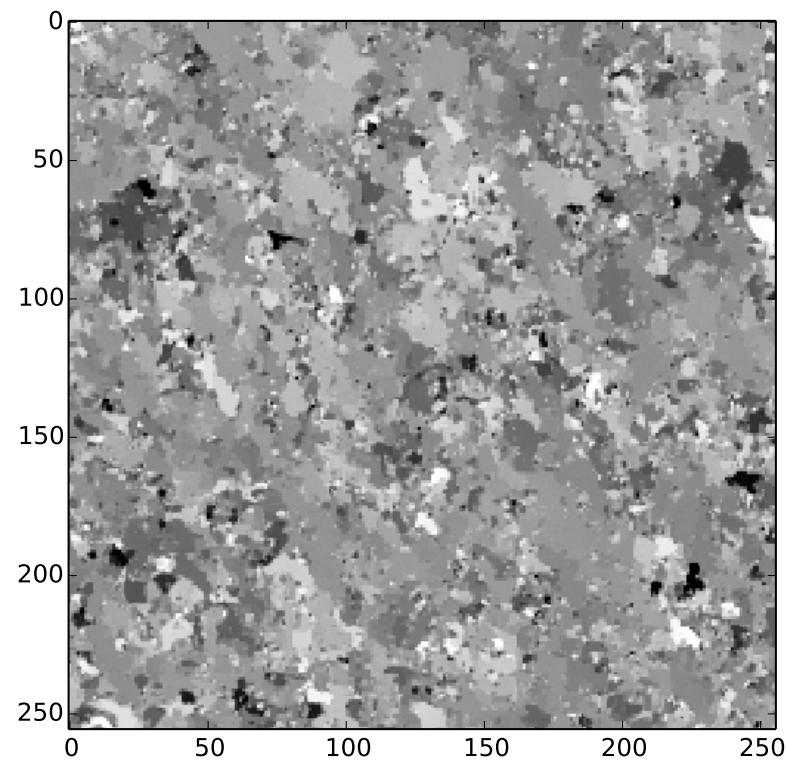


Training Data

Diffusion Probabilistic Model Applied to Dead Leaves

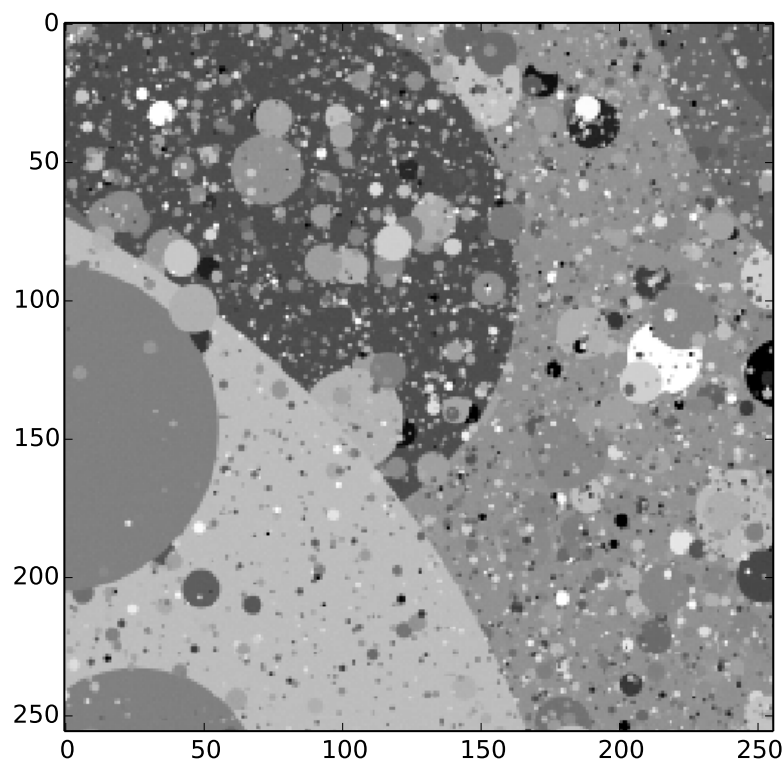


Training Data

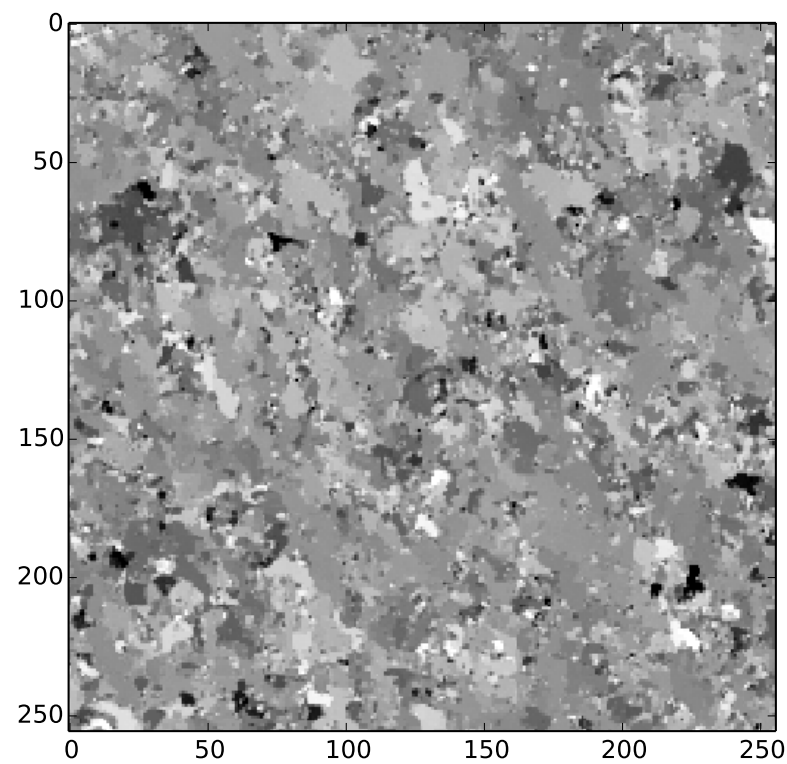


Sample from
[Theis *et al*, 2012]

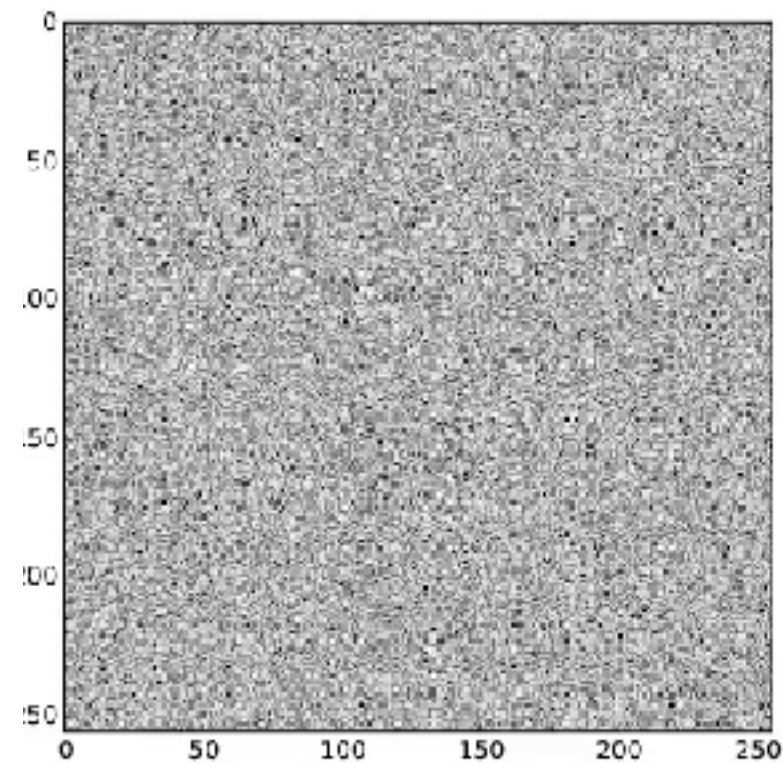
Diffusion Probabilistic Model Applied to Dead Leaves



Training Data

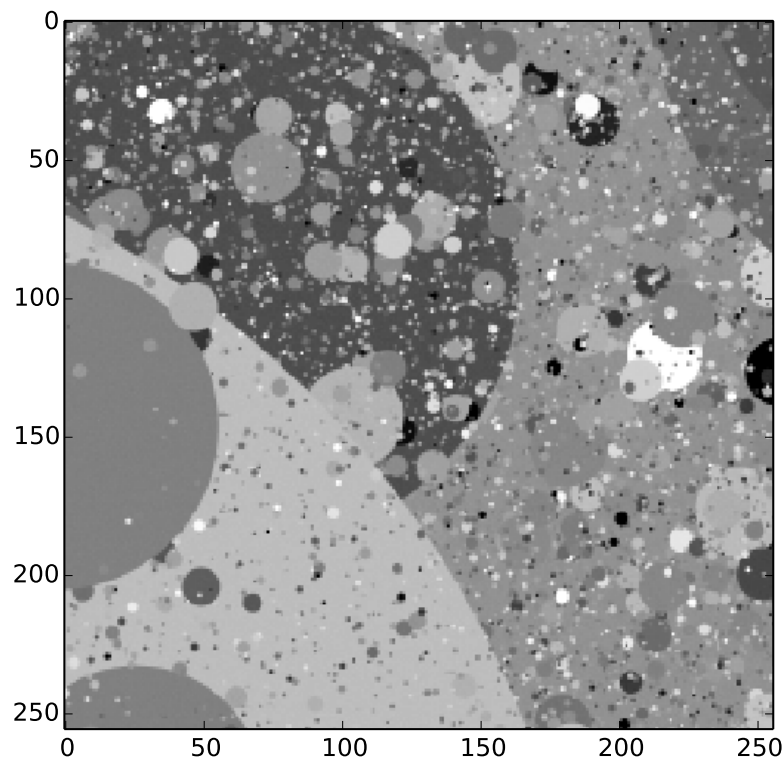


Sample from
[Theis *et al*, 2012]

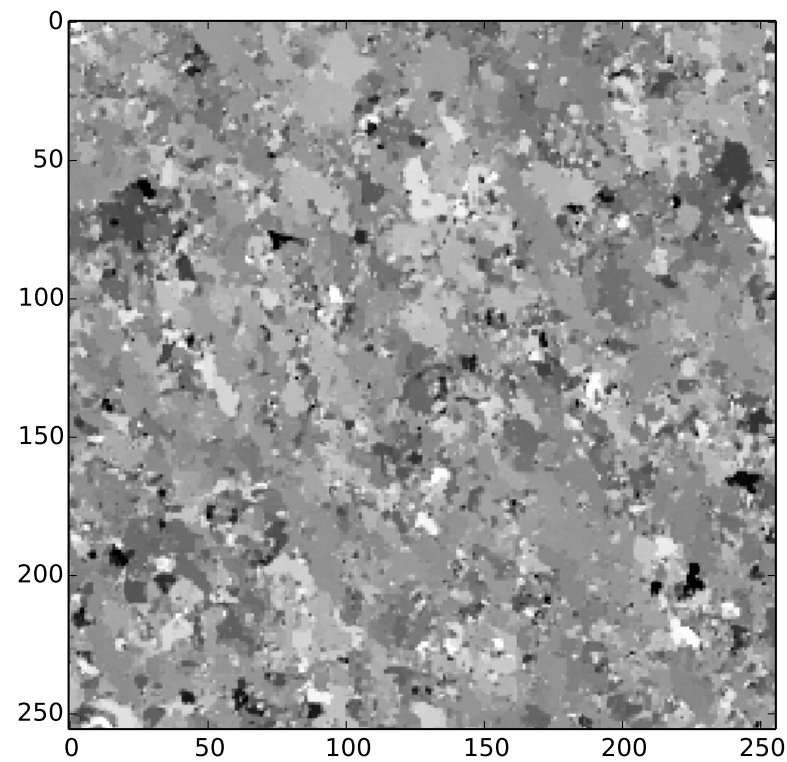


Diffusion Probabilistic Models

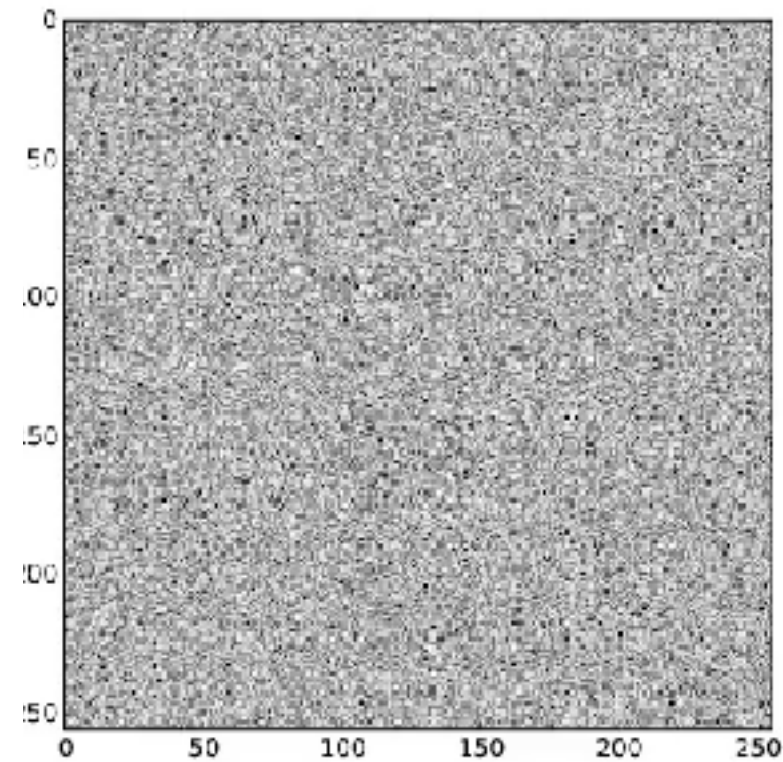
Diffusion Probabilistic Model Applied to Dead Leaves



Training Data

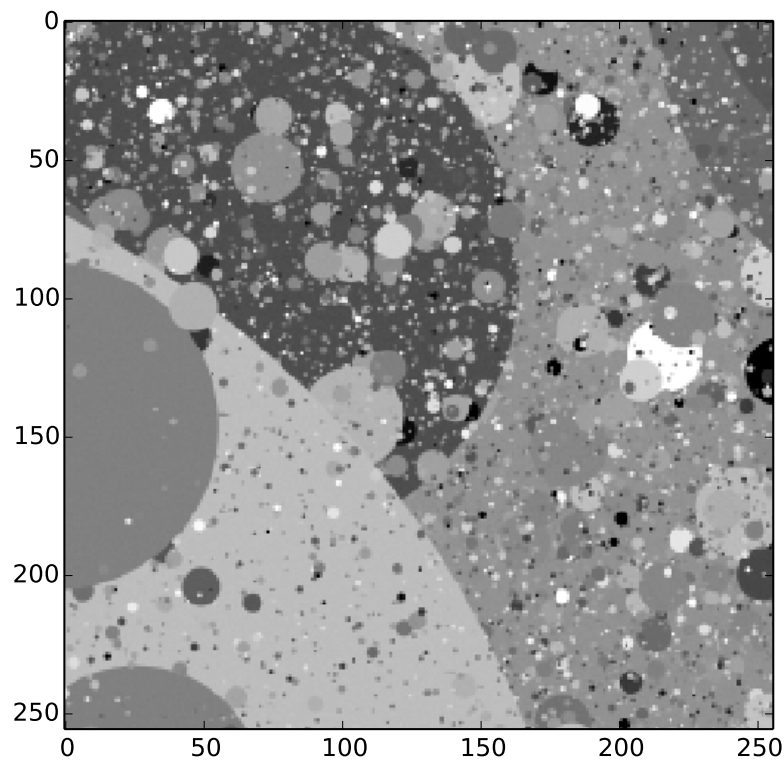


Sample from
[Theis *et al*, 2012]



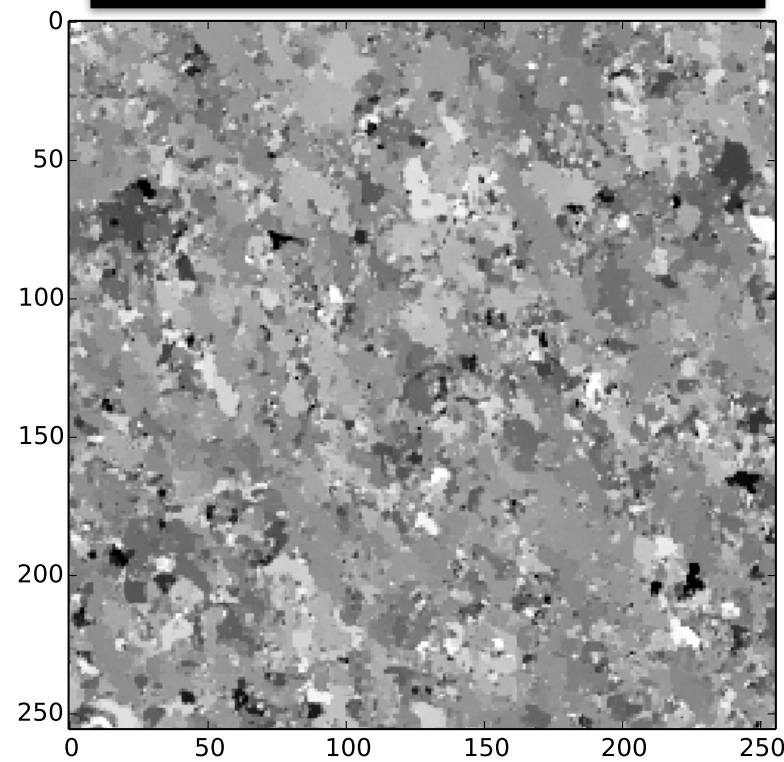
Sample from
diffusion model

Diffusion Probabilistic Model Applied to Dead Leaves



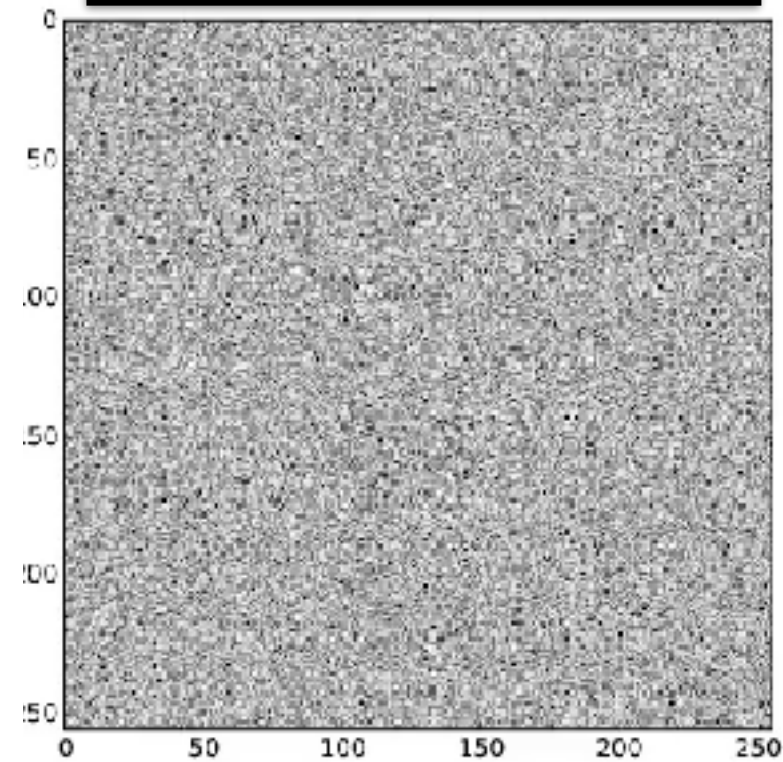
Training Data

Log likelihood
1.24 bits/pixel



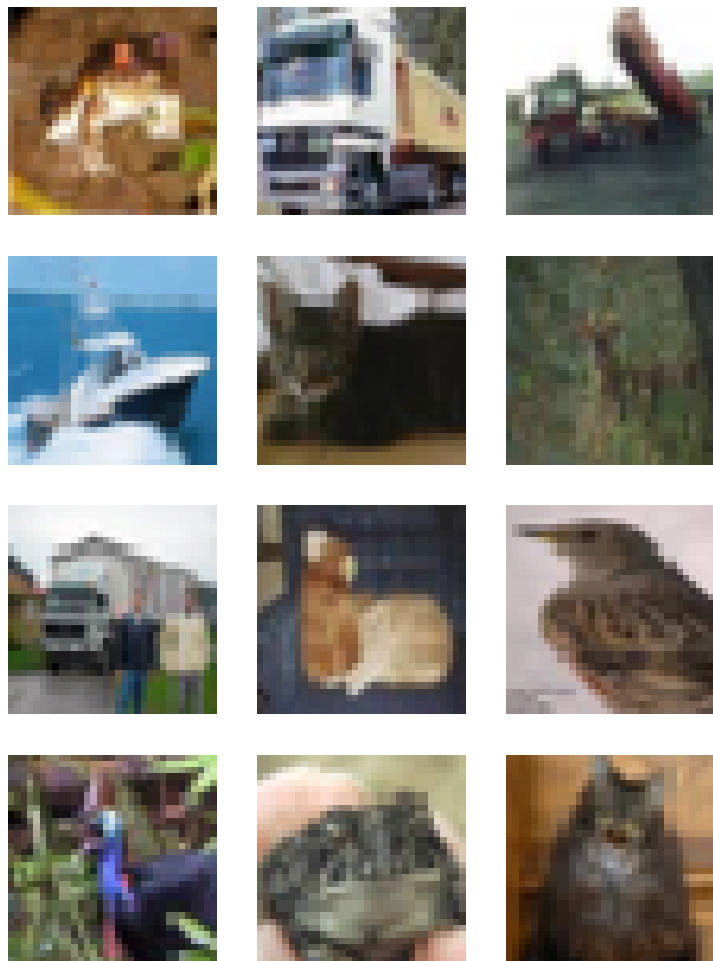
Sample from
[Theis *et al*, 2012]

Log likelihood
1.49 bits/pixel



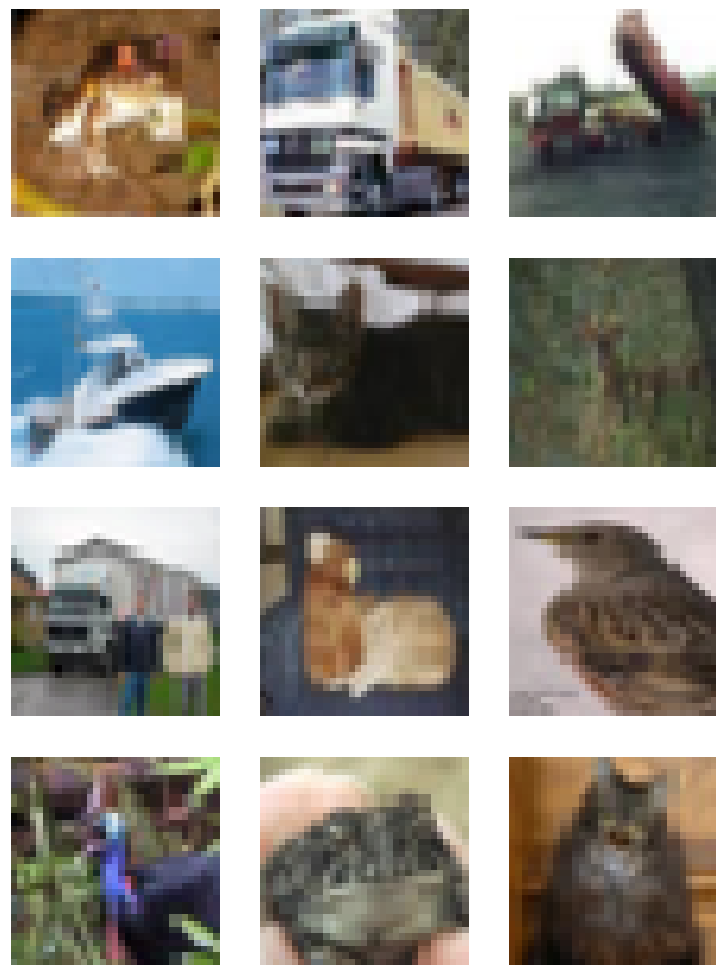
Sample from
diffusion model

Diffusion Probabilistic Model Applied to CIFAR-10

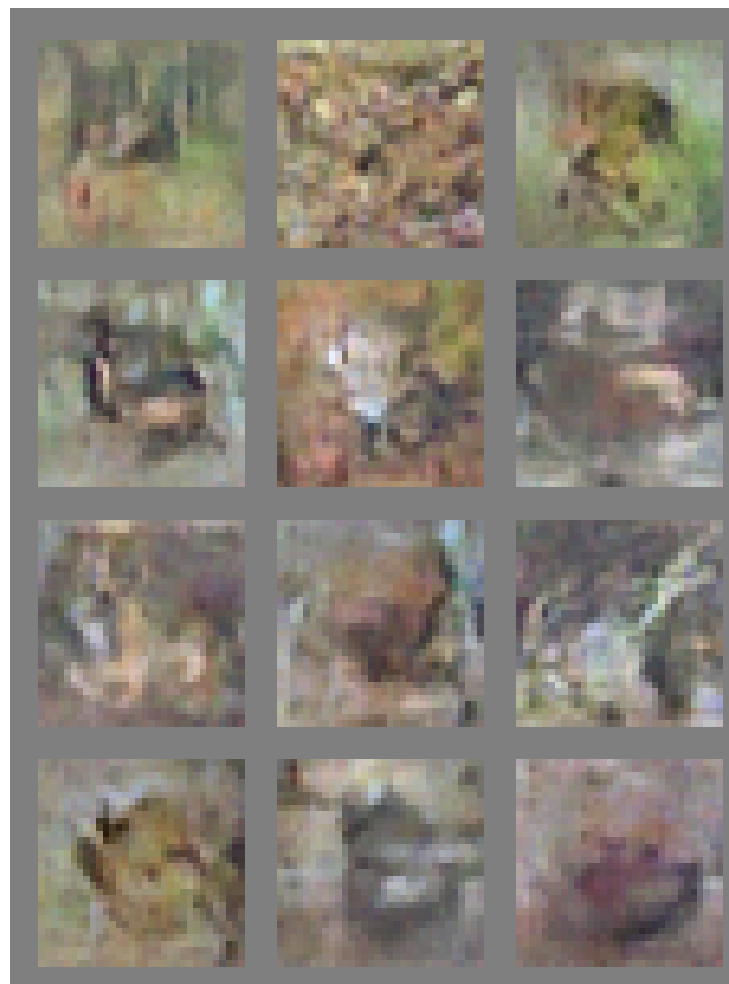


Training Data

Diffusion Probabilistic Model Applied to CIFAR-10

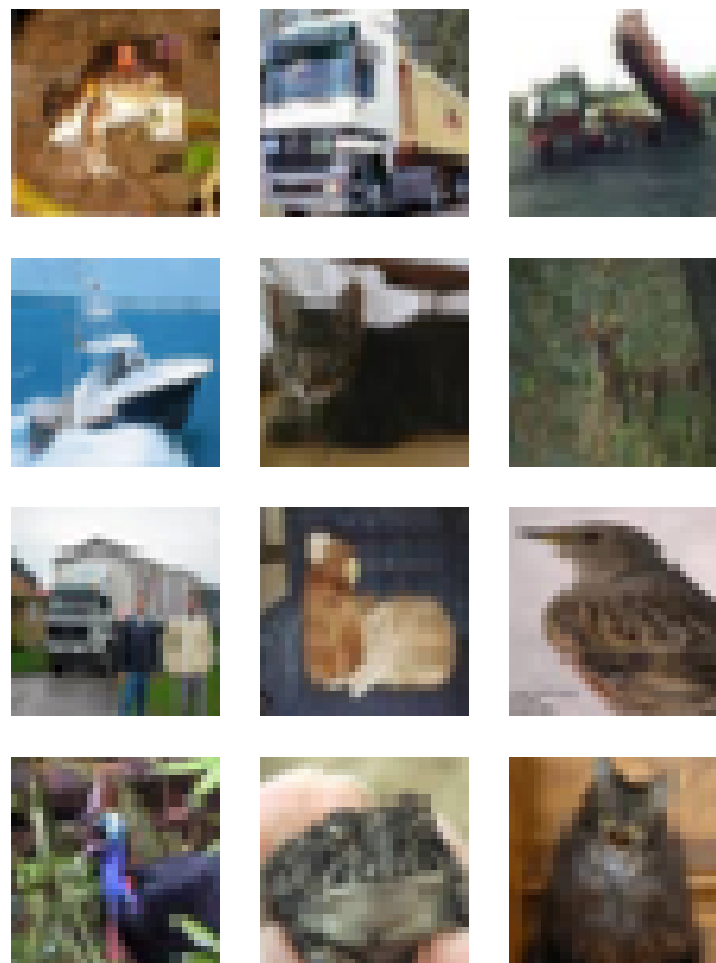


Training Data

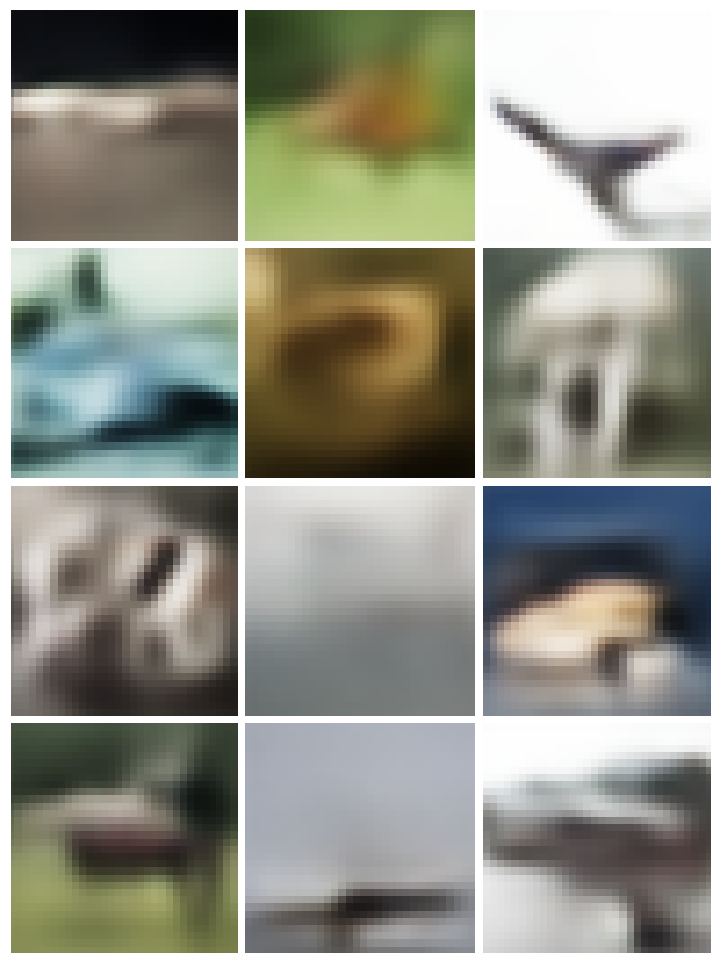


Samples from
Generative Adversarial
[Goodfellow *et al*, 2014]

Diffusion Probabilistic Model Applied to CIFAR-10



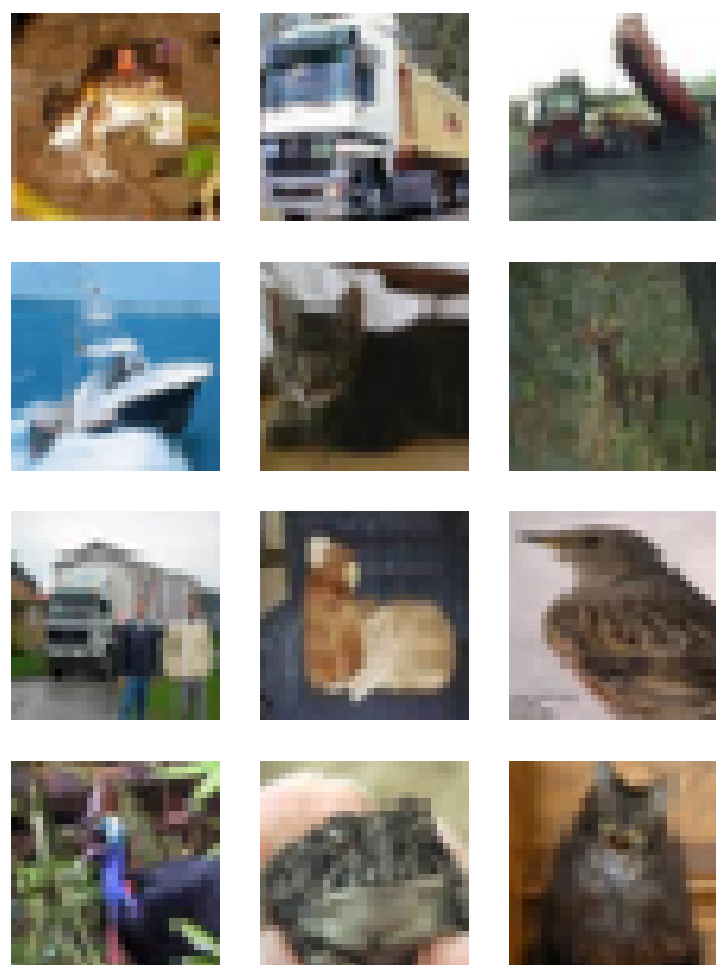
Training Data



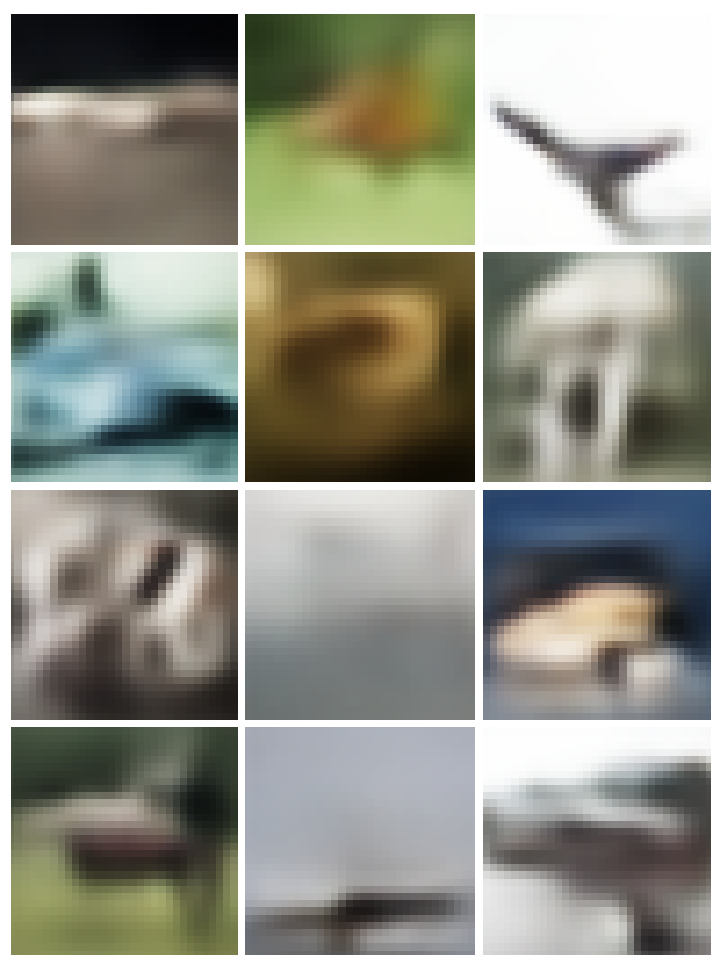
Samples from
DRAW

[Gregor *et al*, 2015]

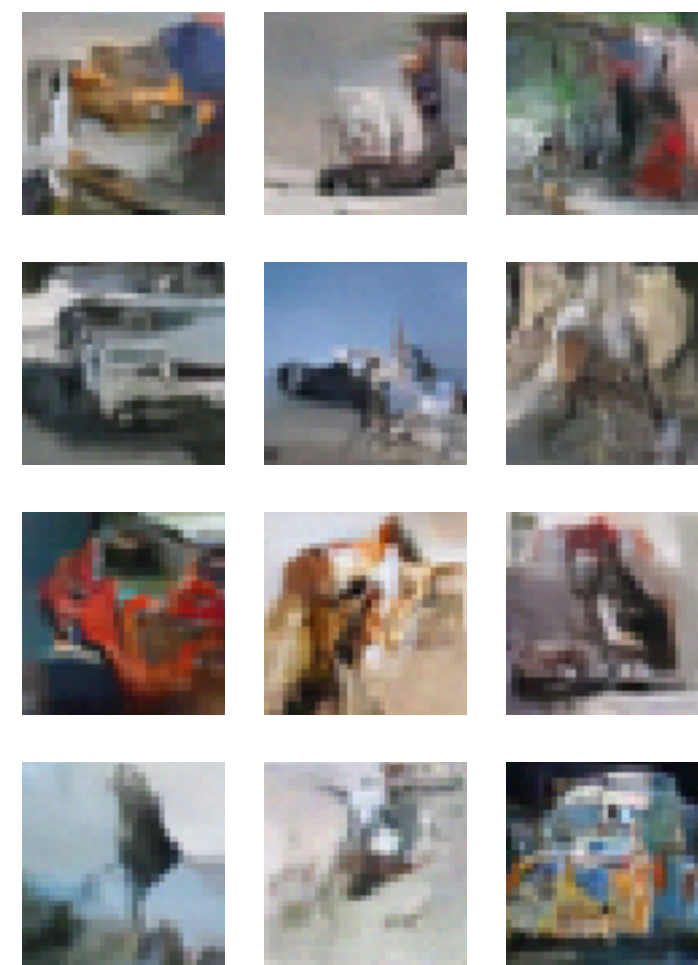
Diffusion Probabilistic Model Applied to CIFAR-10



Training Data



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

Multiplying Distributions is Straightforward

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)

Multiplying Distributions is Straightforward

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

Multiplying Distributions is Straightforward

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

Multiplying Distributions is Straightforward

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

Multiplying Distributions is Straightforward

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

Acts as small perturbation to diffusion process

Multiplying Distributions is Straightforward

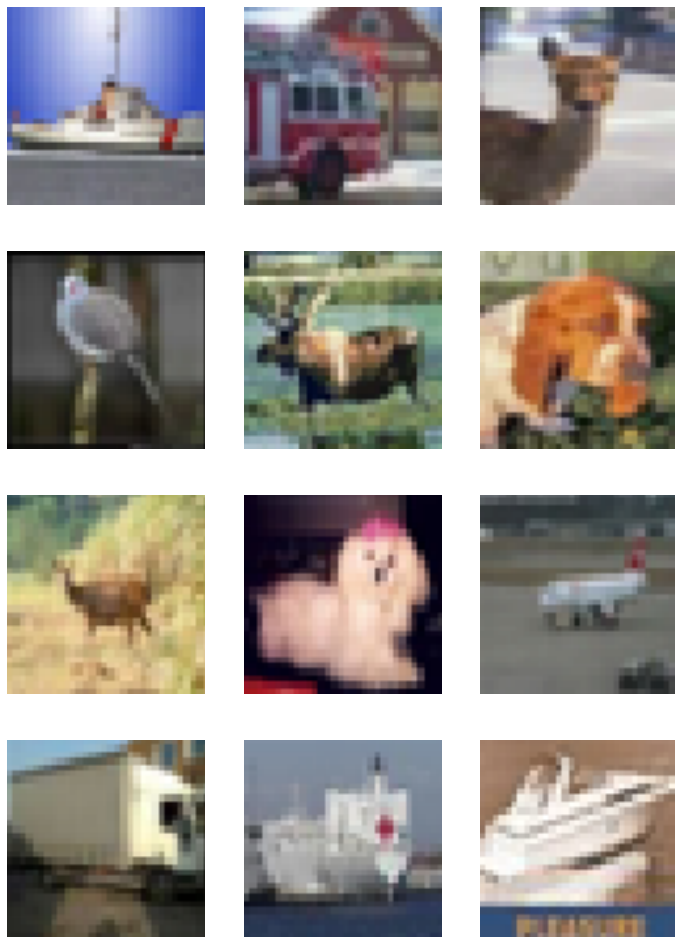
Interested in $\tilde{p}(\mathbf{x}^{(0)}) \propto p(\mathbf{x}^{(0)}) \underline{r(\mathbf{x}^{(0)})}$

Acts as small perturbation to diffusion process

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

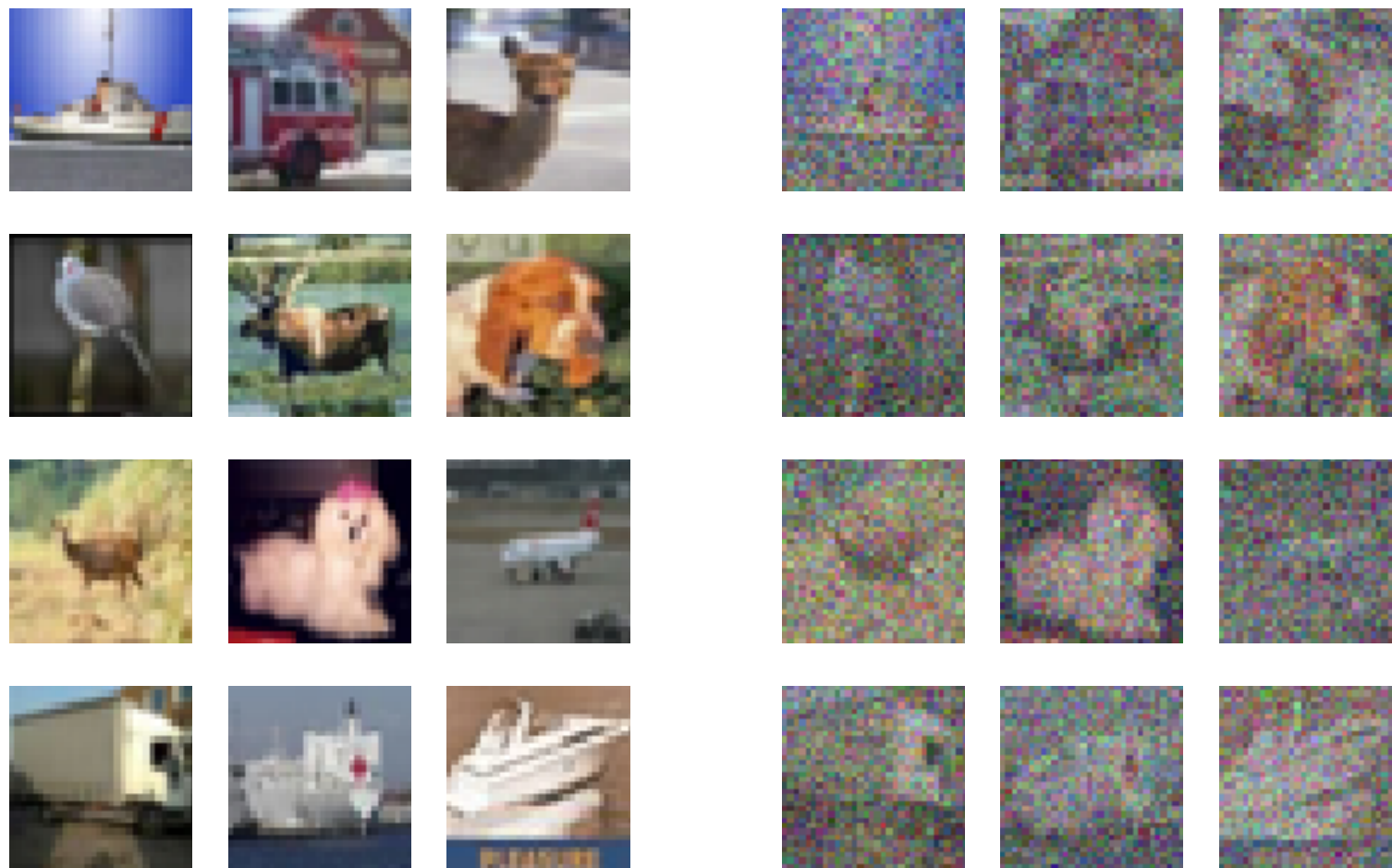
$$\tilde{p}(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}) \approx \mathcal{N}\left(\mathbf{x}^{(t-1)}; \mathbf{f}_{\mu}(\mathbf{x}^{(t)}, t) + \mathbf{f}_{\Sigma}(\mathbf{x}^{(t)}, t) \frac{\partial \log r(\mathbf{x}^{(t-1)'})}{\partial \mathbf{x}^{(t-1)'}} \Big|_{\mathbf{x}^{(t-1)' = \mathbf{f}_{\mu}(\mathbf{x}^{(t)}, t)}, \mathbf{f}_{\Sigma}(\mathbf{x}^{(t)}, t)\right)$$

Image Denoising by Sampling from Posterior



Holdout Data

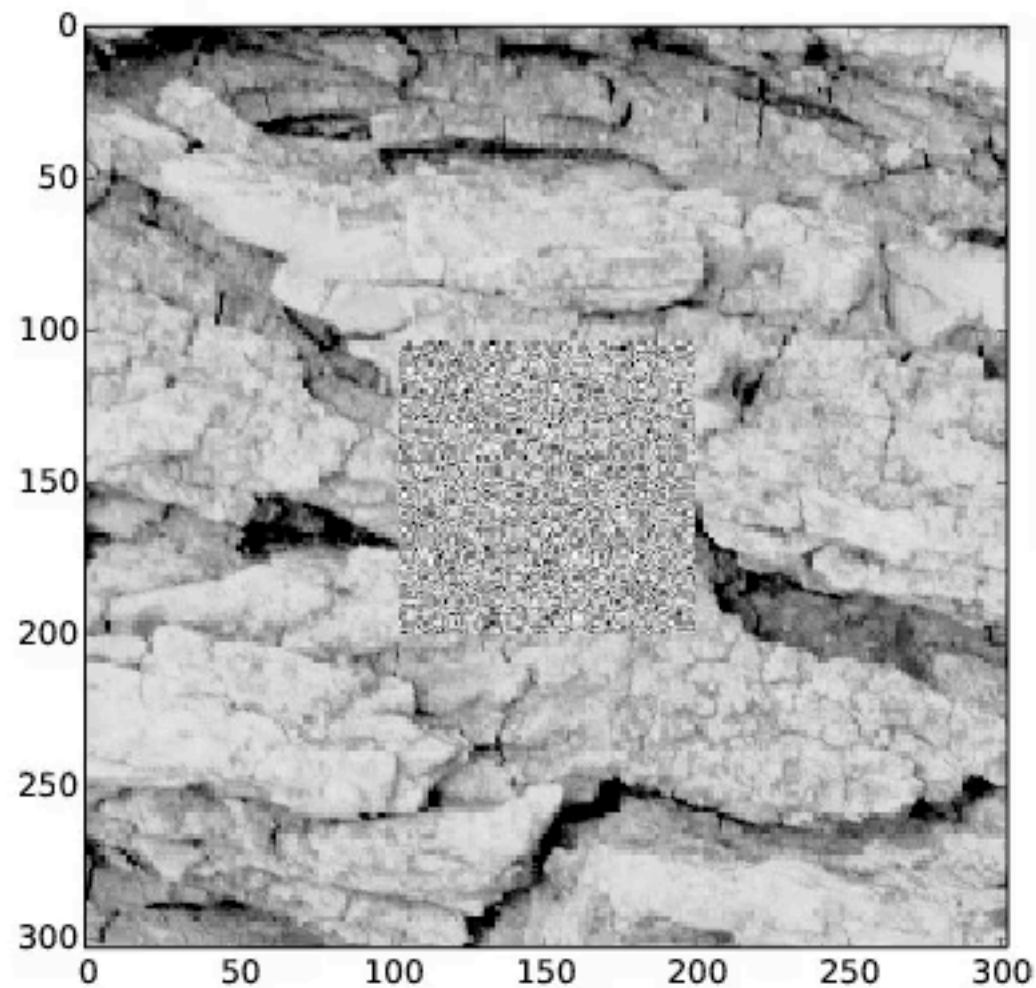
Image Denoising by Sampling from Posterior



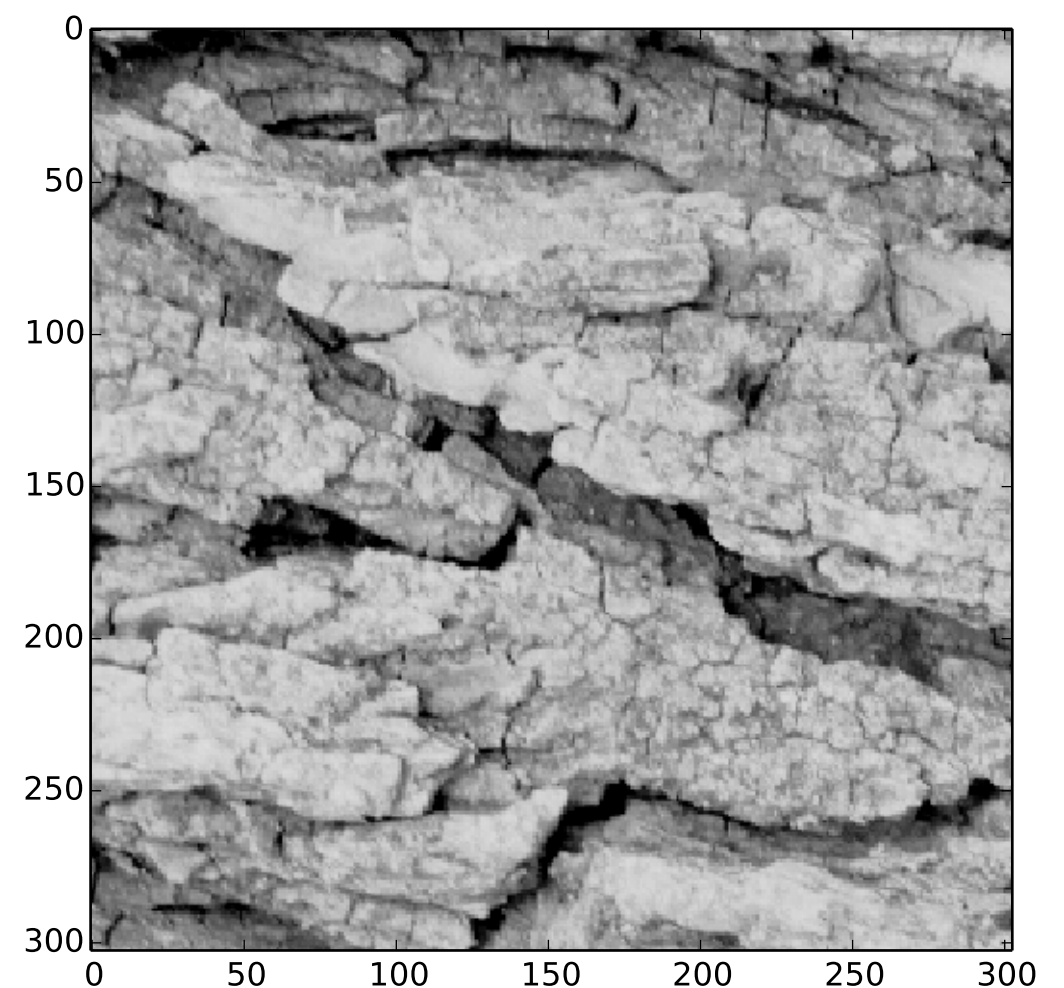
Holdout Data

Corrupted
(SNR = 1)

Image Inpainting by Sampling from Posterior

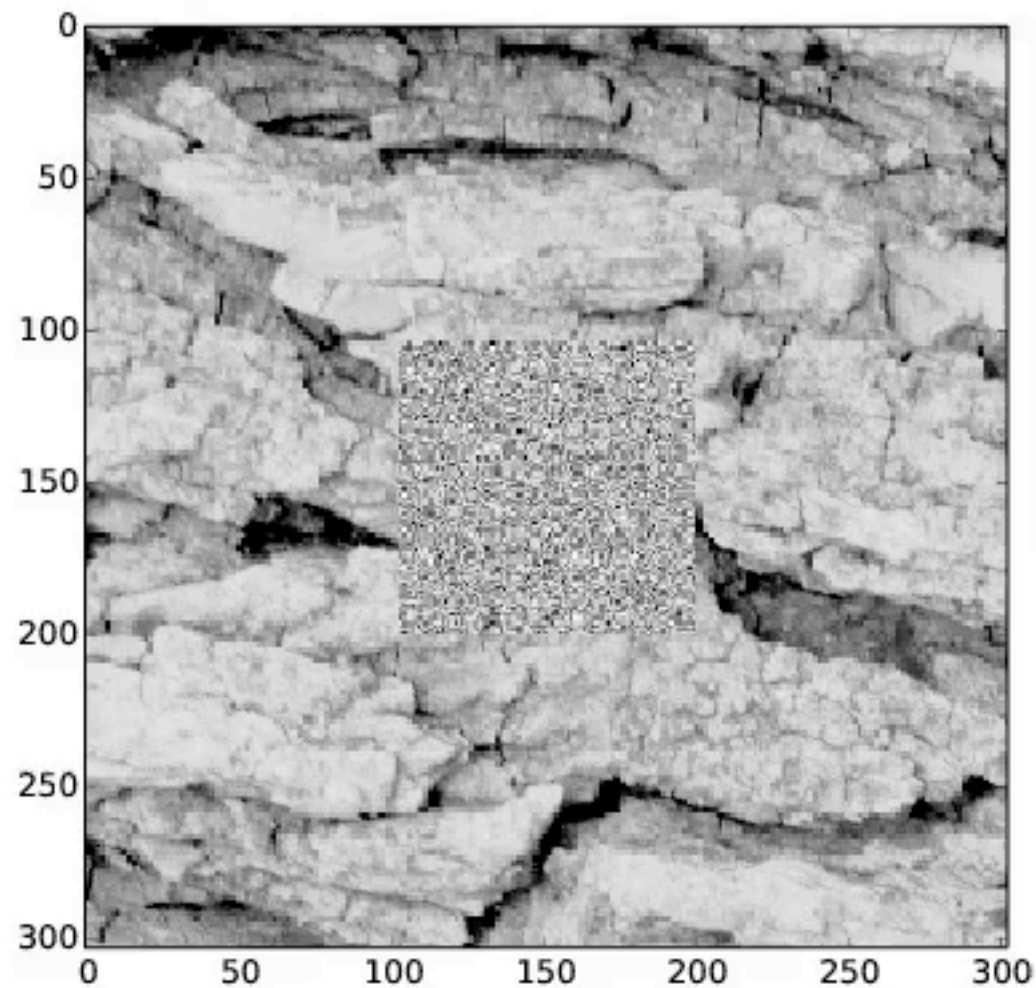


Inpainted image

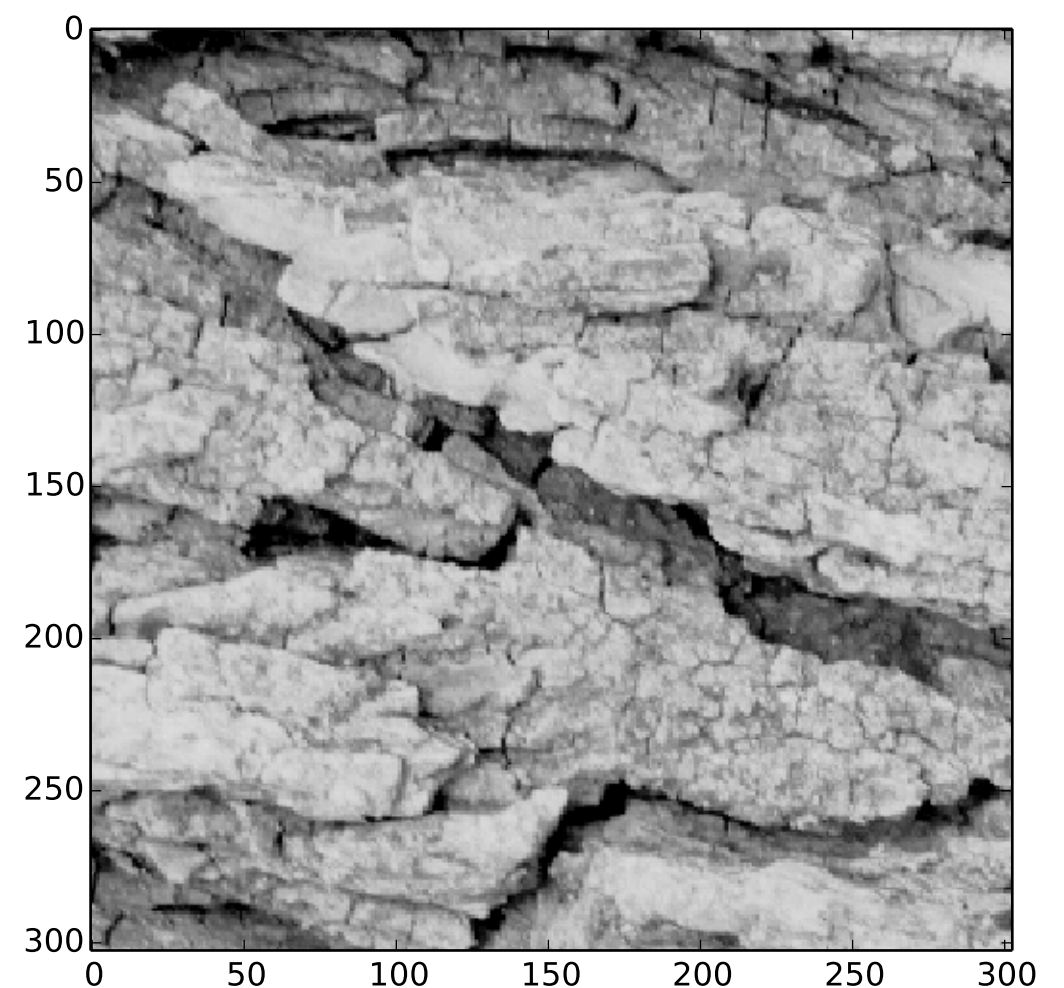


True image

Image Inpainting by Sampling from Posterior



Inpainted image



True image

Flexible and tractable method for deep unsupervised learning

Flexible and tractable method for deep unsupervised learning

- **Flexible:** Diffusion process for any (smooth) distribution

Flexible and tractable method for deep unsupervised learning

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

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

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

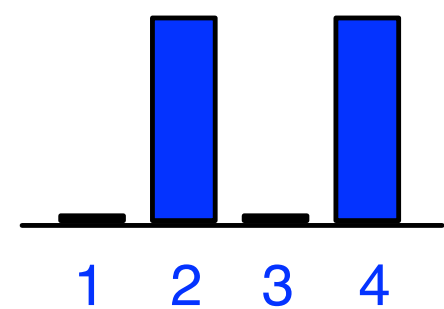
- **Flexible:** Diffusion process for any (smooth) distribution
 - Binary or continuous state space
- **Tractable:** Training, exact sampling, inference, evaluation
- Deep networks with thousands of layers (/ time steps)
- Easy to multiply distributions (e.g. for posterior)
- Bounds on entropy production

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

Minimum Probability Flow Learning

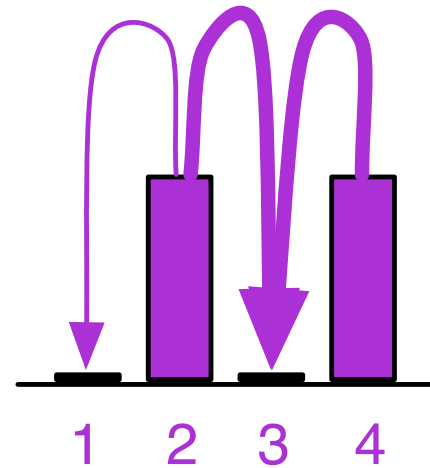
- Estimate parameters in energy based models, by minimizing probability flow under master equation from stat. mech.



data distribution

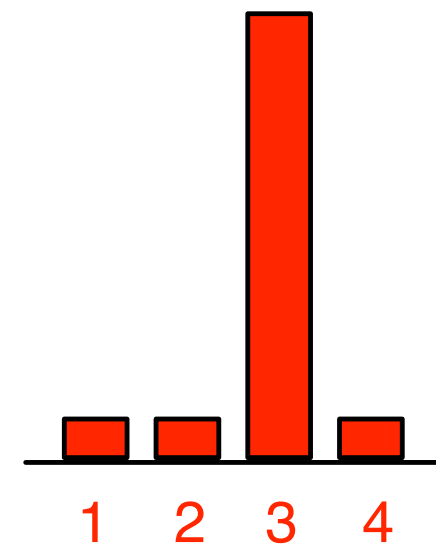
$$p_i^{(0)} = \text{fraction data in state } i$$

$$\dot{p}_i^{(t)} = \sum_{j \neq i} \Gamma_{ij}(\theta) p_j^{(t)} - \sum_{j \neq i} \Gamma_{ji}(\theta) p_i^{(t)}$$



dynamics

$$\dot{p}_i^{(t)} = \sum_{j \neq i} \Gamma_{ij}(\theta) p_j^{(t)} - \sum_{j \neq i} \Gamma_{ji}(\theta) p_i^{(t)}$$



model distribution

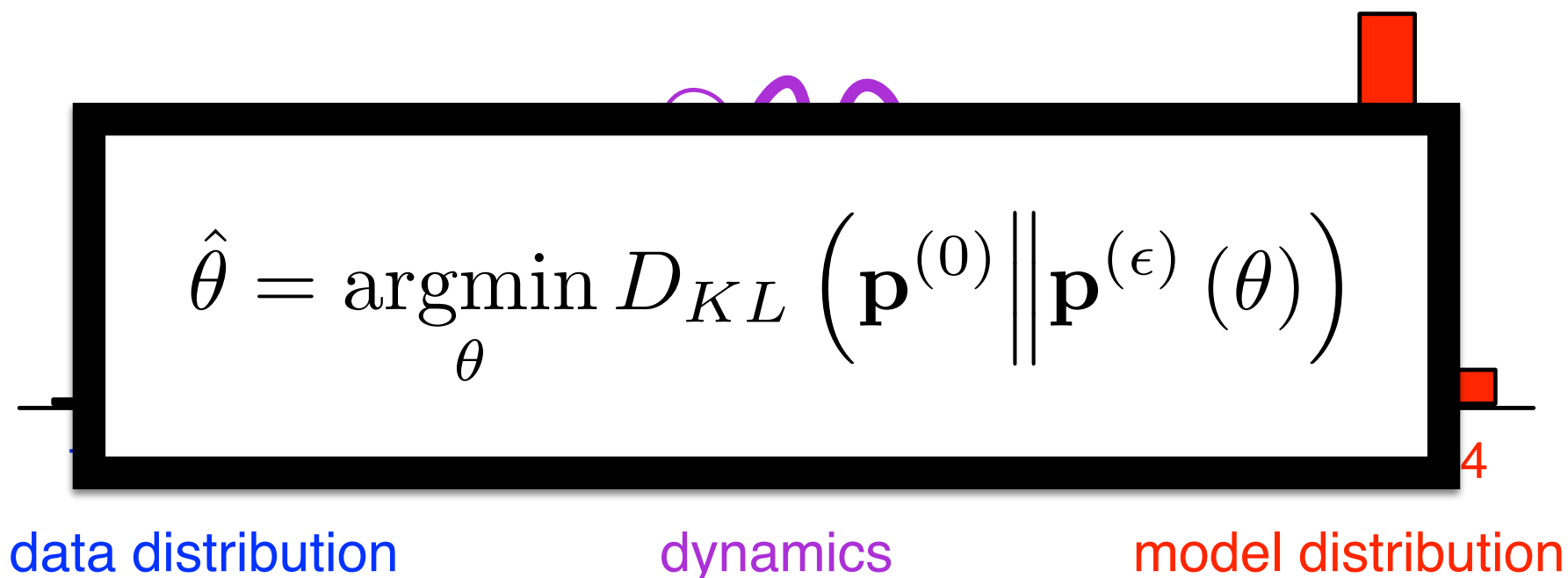
$$p_i^{(\infty)}(\theta) = \frac{e^{-E_i(\theta)}}{Z(\theta)}$$

$$\dot{p}_i^{(\infty)} = 0$$

[PRL, 2011]
[ICML, 2011]

Minimum Probability Flow Learning

- Estimate parameters in energy based models, by minimizing probability flow under master equation from stat. mech.



$p_i^{(0)}$ = fraction data
in state i

$$\dot{p}_i^{(t)} = \sum_{j \neq i} \Gamma_{ij}(\theta) p_j^{(0)} - \sum_{j \neq i} \Gamma_{ji}(\theta) p_i^{(0)}$$



$$\dot{p}_i^{(t)} = \sum_{j \neq i} \Gamma_{ij}(\theta) p_j^{(t)} - \sum_{j \neq i} \Gamma_{ji}(\theta) p_i^{(t)}$$

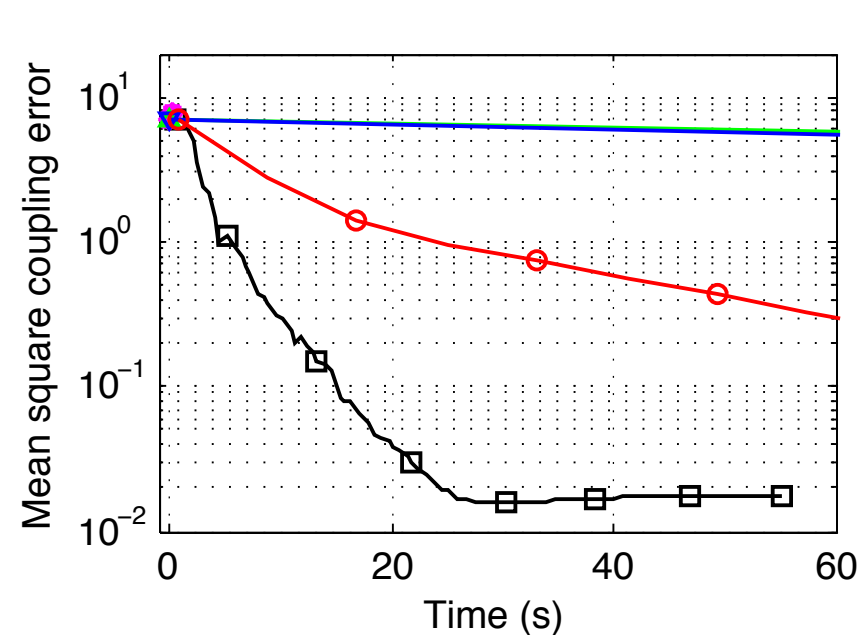
$$p_i^{(\infty)}(\theta) = \frac{e^{-E_i(\theta)}}{Z(\theta)}$$

$$\dot{p}_i^{(\infty)} = 0$$

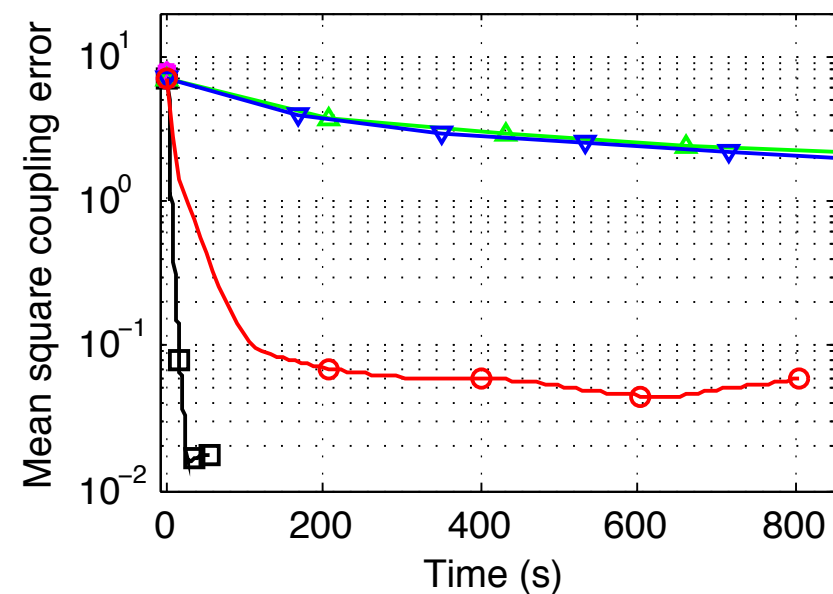
[PRL, 2011]
[ICML, 2011]

Minimum Probability Flow Learning

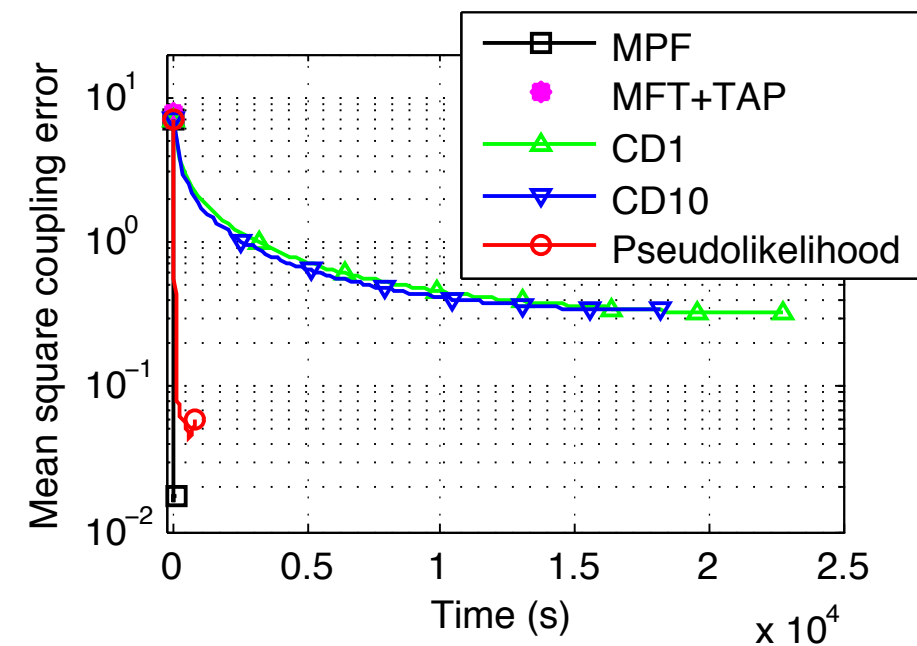
- More rapidly solves inverse Ising problem (estimate Ising couplings from samples)



First 60 seconds



First 800 seconds



First 25,000 seconds

[PRL, 2011]
[ICML, 2011]

Hamiltonian (hybrid) Monte Carlo Without Detailed Balance

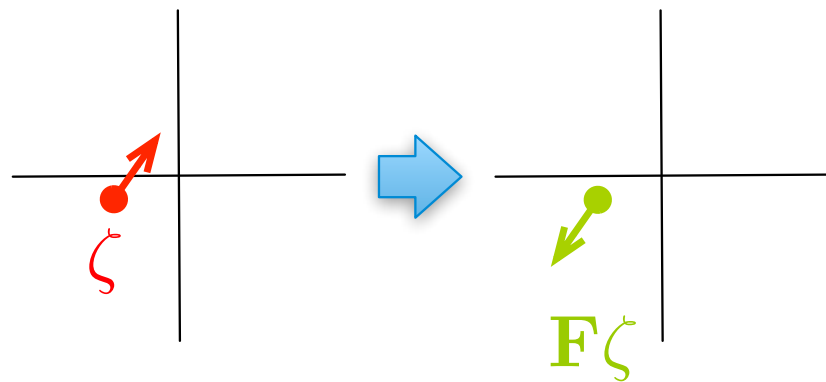
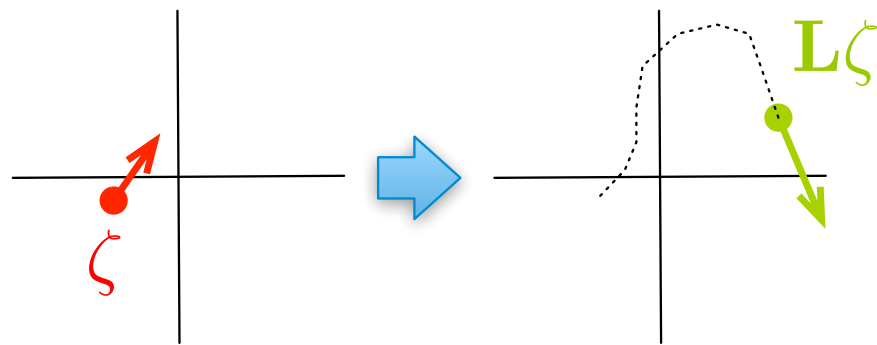
- HMC as operators on discrete state space

[ICML, 2014]

Hamiltonian (hybrid) Monte Carlo Without Detailed Balance

- HMC as operators on discrete state space

Operators:

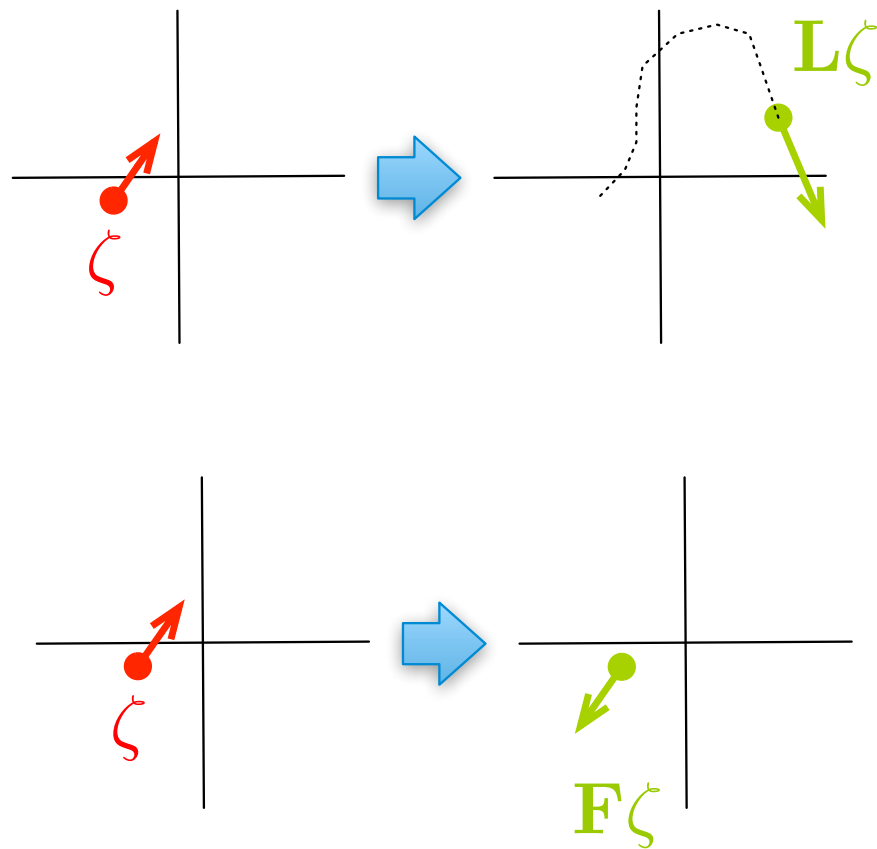


[ICML, 2014]

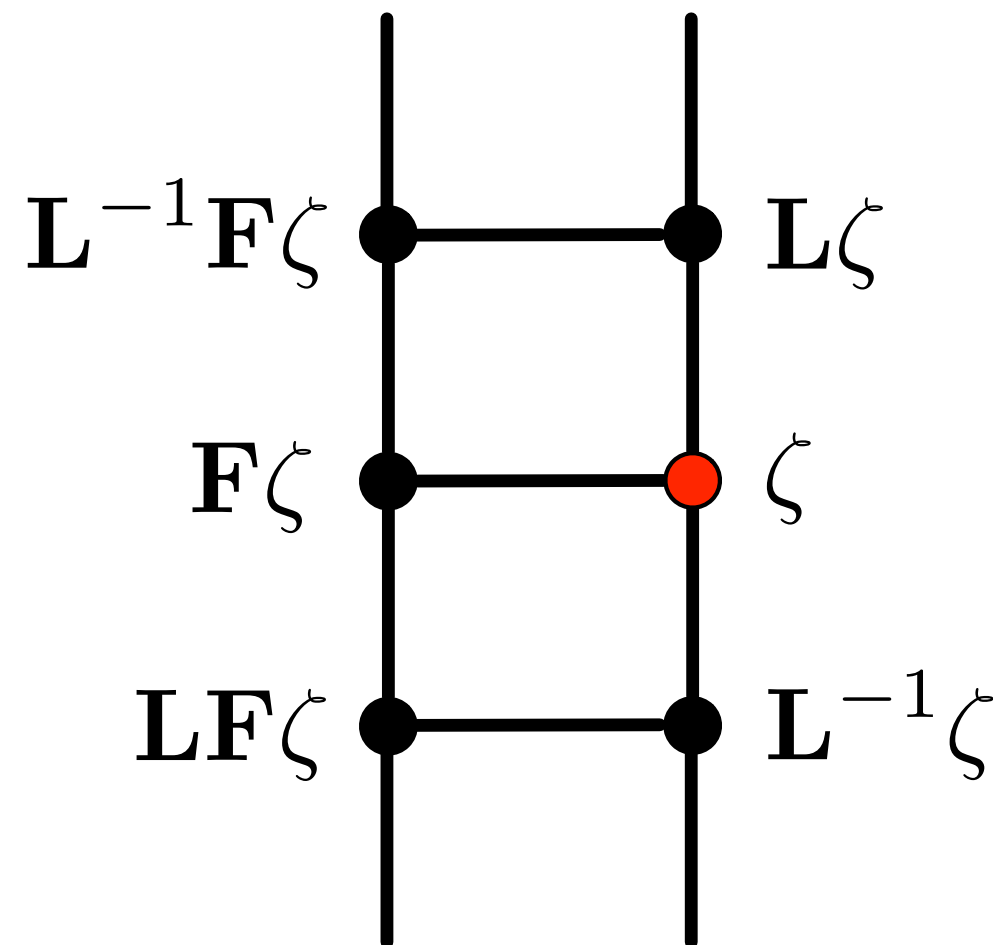
Hamiltonian (hybrid) Monte Carlo Without Detailed Balance

- HMC as operators on discrete state space

Operators:



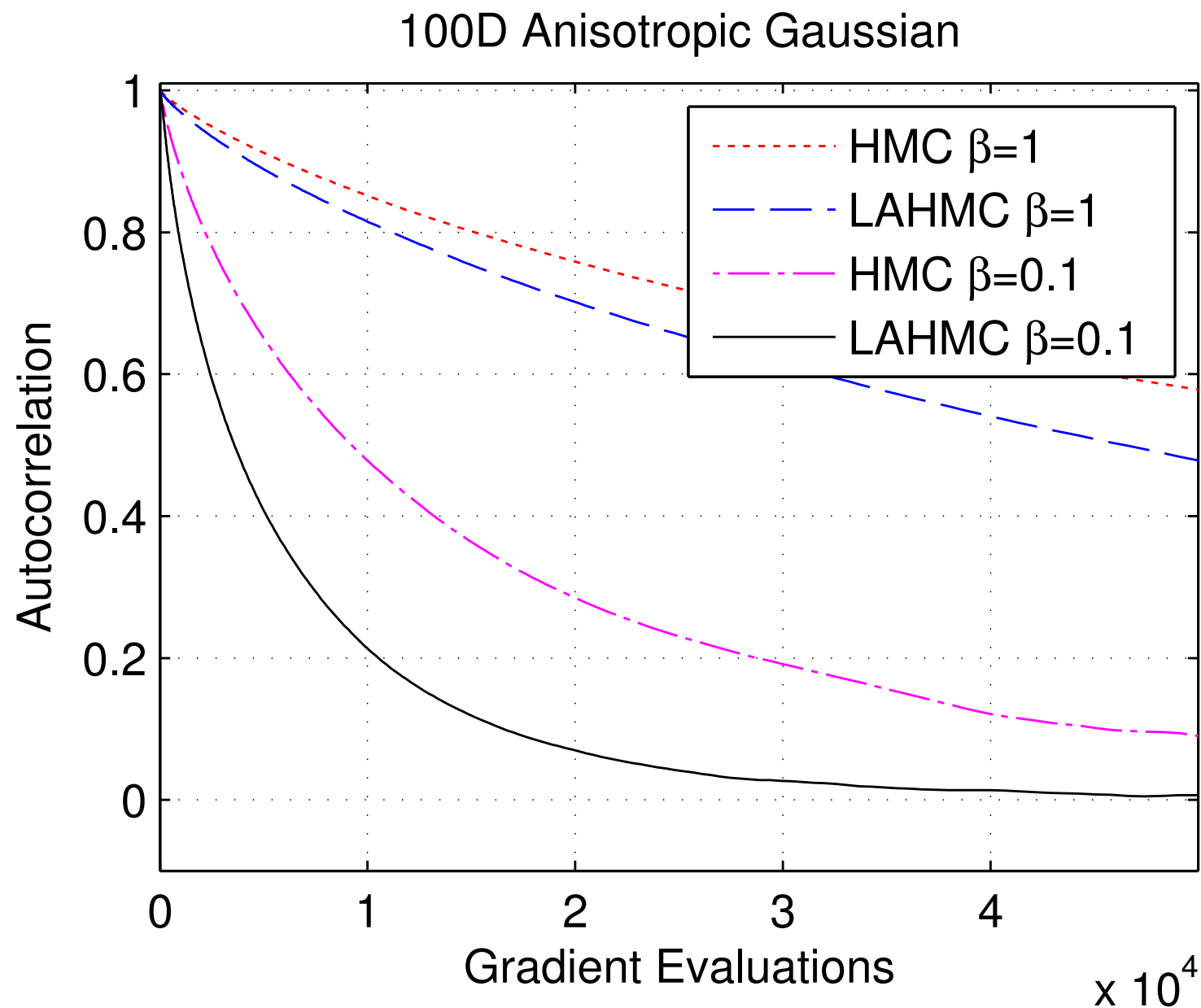
State space:



[ICML, 2014]

Hamiltonian (hybrid) Monte Carlo Without Detailed Balance

- Improved mixing by violating detailed balance



[ICML, 2014]

Statistical Physics of Deep Neural Networks

Neural network after random initialization:

$$z_i^l = \sum_j W_{ij}^l y_j^l + b^j \qquad y_i^{l+1} = \phi(z_i^l)$$

$$W_{ij}^l \sim \mathcal{N}(0, \sigma_w^2 / N_{l-1}) \qquad b_i^l \sim \mathcal{N}(0, \sigma_b^2)$$

**[NIPS, 2016]
[ICLR, 2017]**

Statistical Physics of Deep Neural Networks

Neural network after random initialization:

$$z_i^l = \sum_j W_{ij}^l y_j^l + b^j \quad y_i^{l+1} = \phi(z_i^l)$$

$$W_{ij}^l \sim \mathcal{N}(0, \sigma_w^2 / N_{l-1}) \quad b_i^l \sim \mathcal{N}(0, \sigma_b^2)$$

Central limit theorem \Rightarrow recurrence relation:

$$z_i^l \sim \mathcal{N}(0, q^l)$$

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

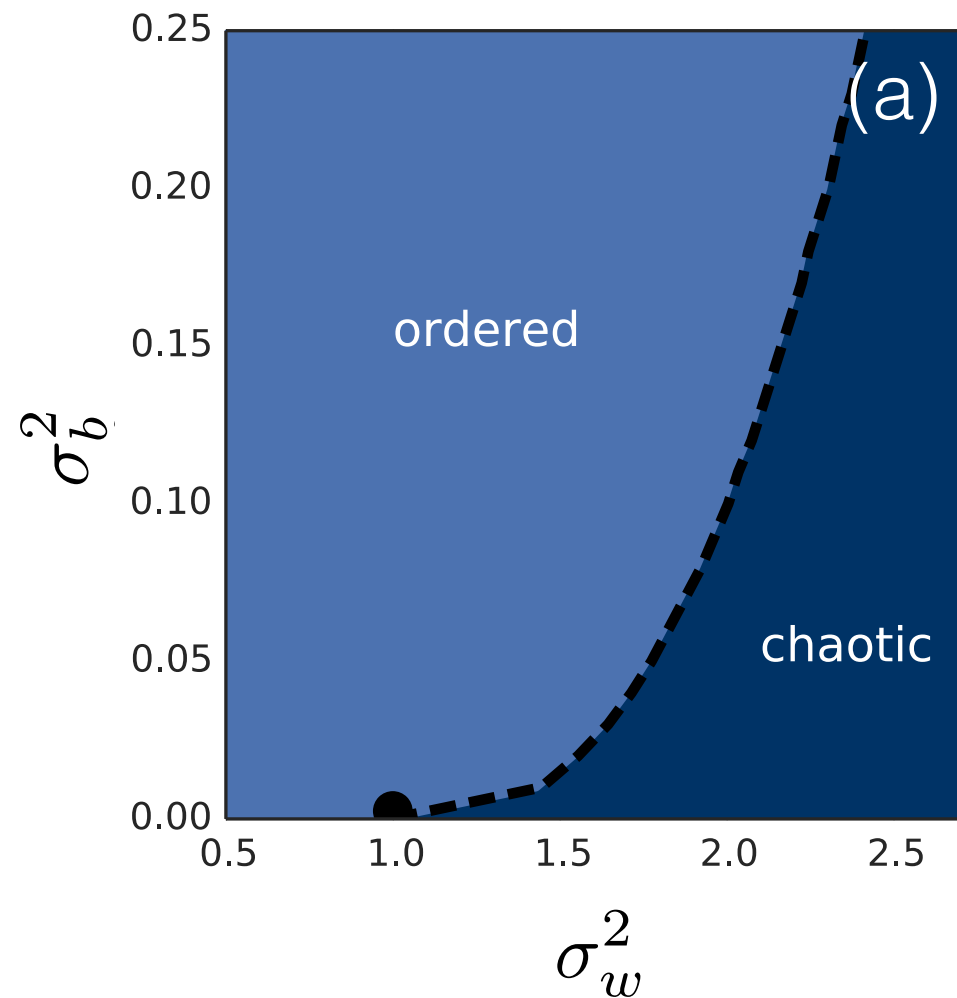
**[NIPS, 2016]
[ICLR, 2017]**

Statistical Physics of Deep Neural Networks

[NIPS, 2016]
[ICLR, 2017]

Statistical Physics of Deep Neural Networks

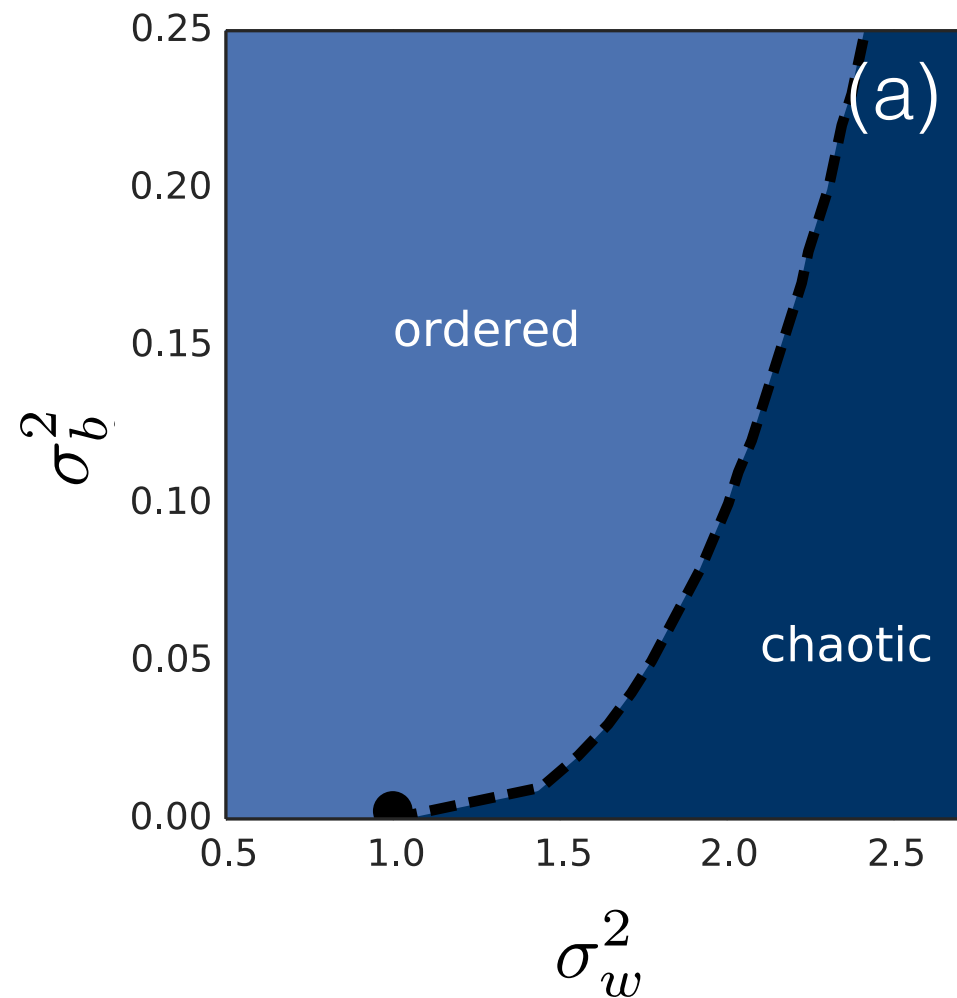
Phase diagram:



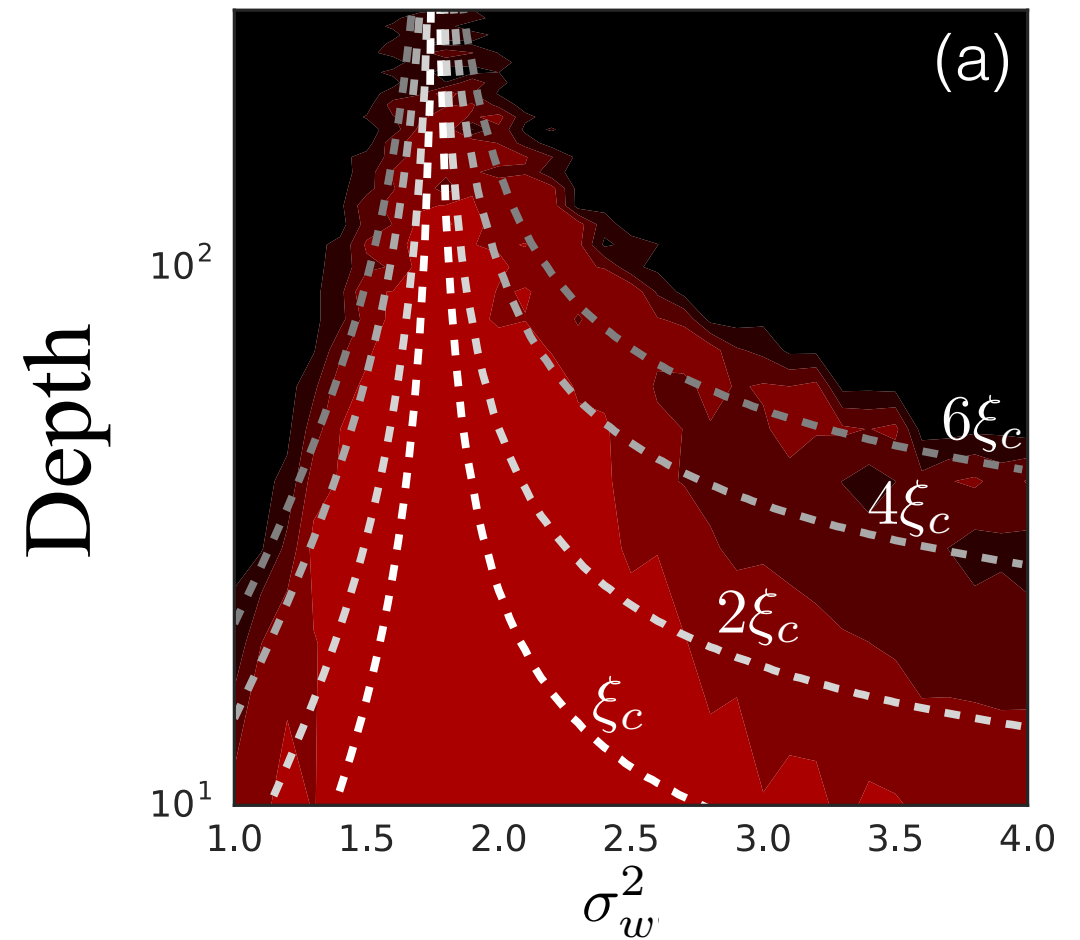
[NIPS, 2016]
[ICLR, 2017]

Statistical Physics of Deep Neural Networks

Phase diagram:



Predict trainable depth:



[NIPS, 2016]
[ICLR, 2017]

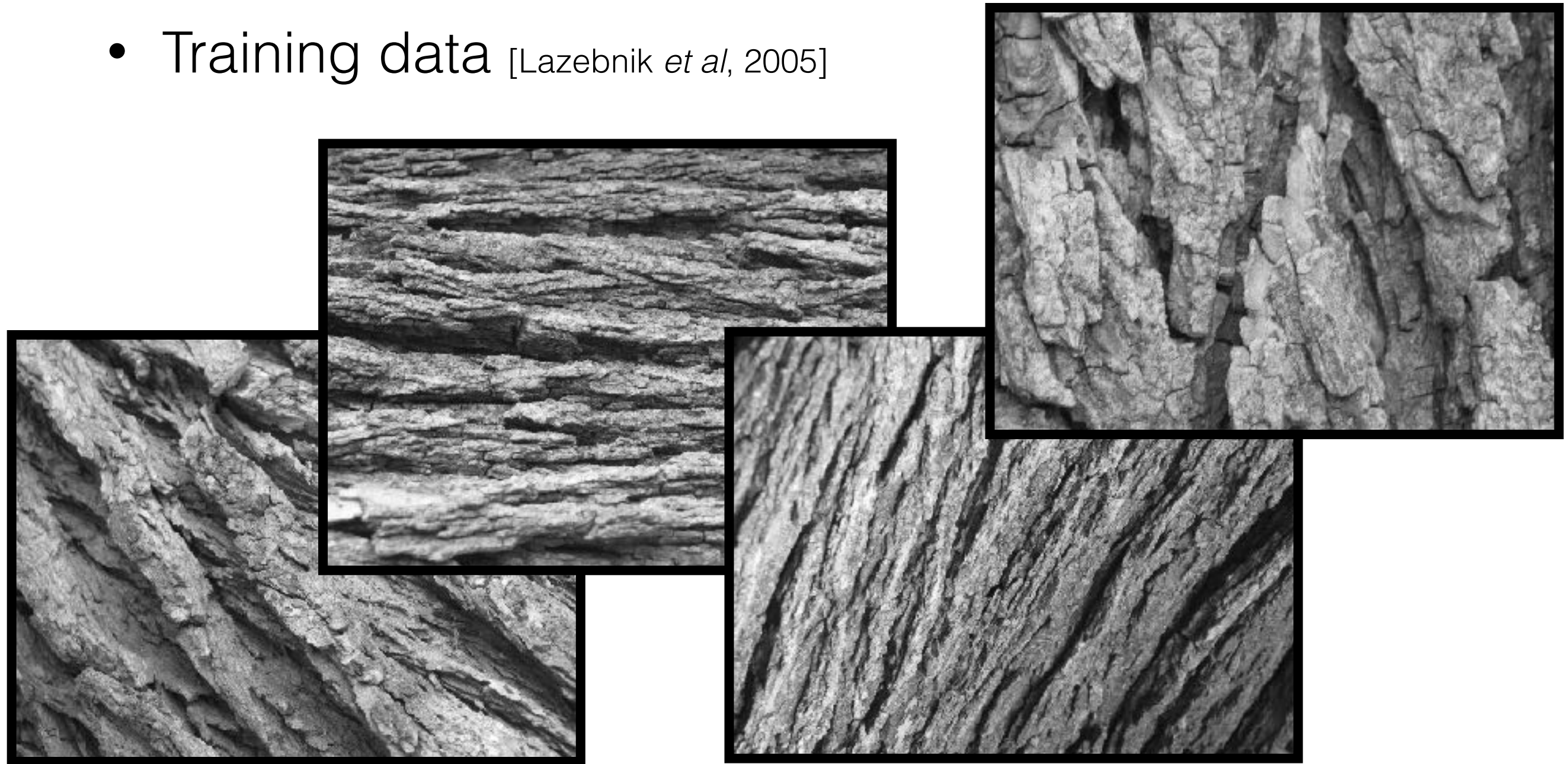
Thanks!

- 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

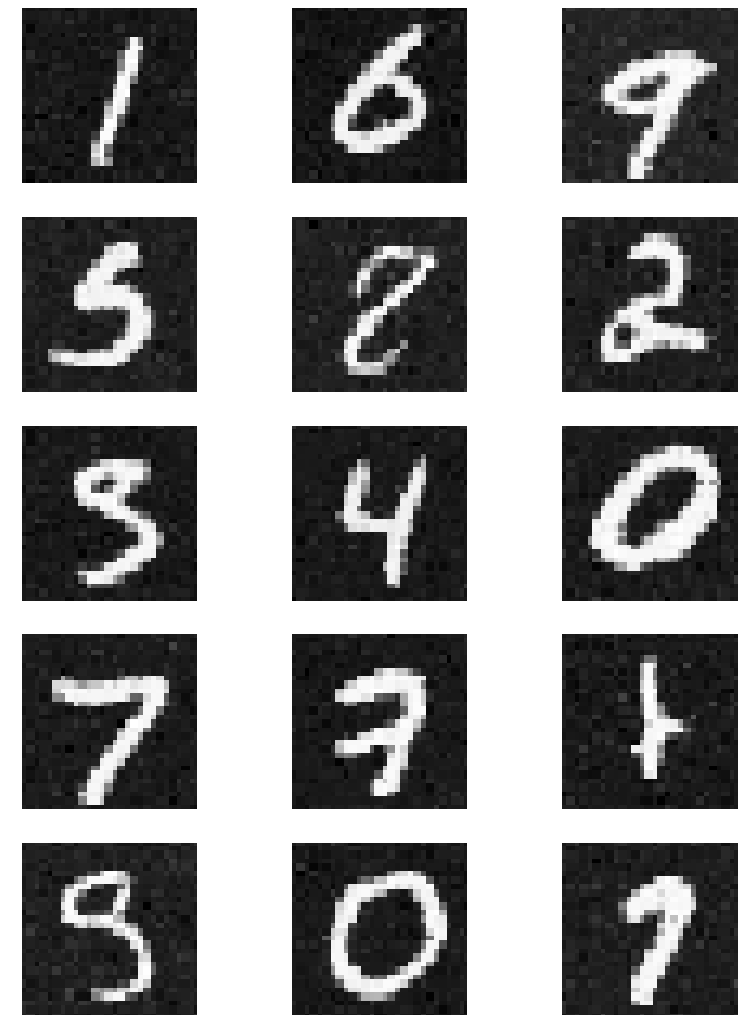
Image Inpainting by Sampling from Posterior

- Training data [Lazebnik *et al*, 2005]



Diffusion Probabilistic Model Applied to MNIST

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



Samples from
diffusion model

* via Parzen window code from [Goodfellow *et al*, 2014]

Future Work

Future Work

- Continuous time formulation

Future Work

- Continuous time formulation
- Perturbation around energy based model

Future Work

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

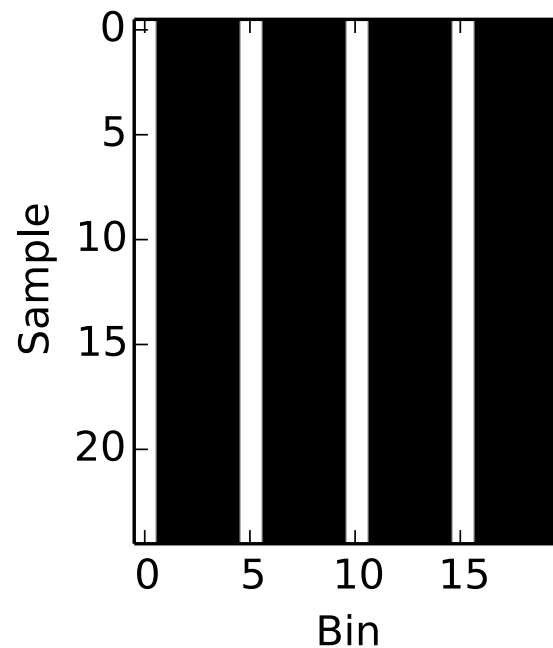
Outline

- **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

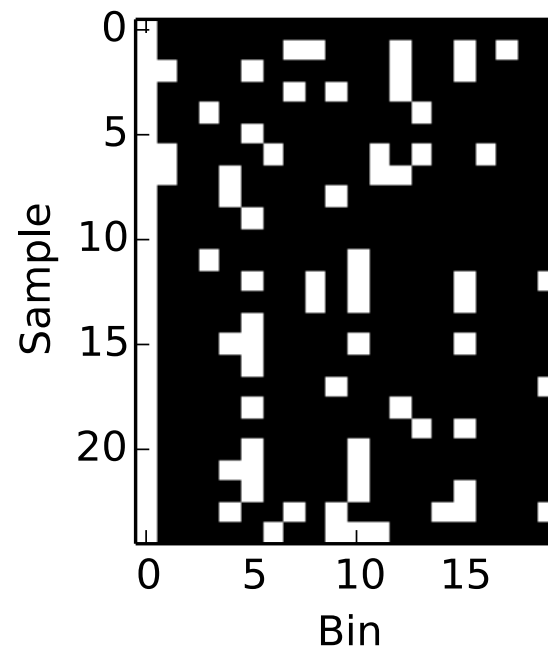
Toy Binary Sequence Learning

$$p(\mathbf{x}^{(0 \dots T)})$$

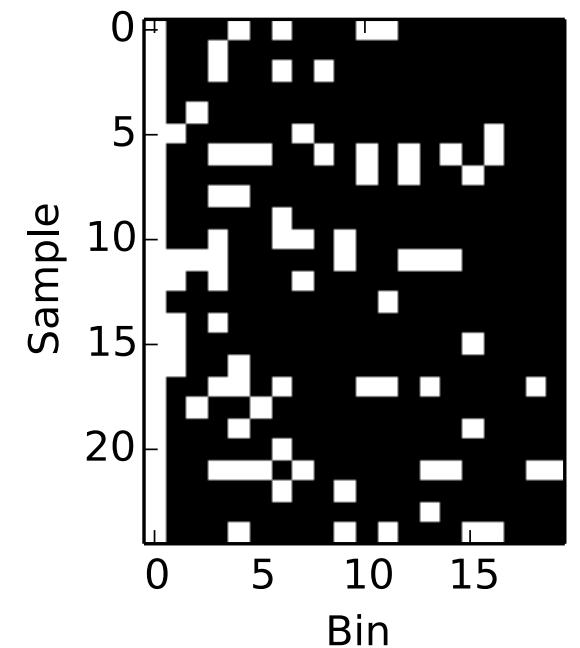
$t = 0$



$t = \frac{T}{2}$



$t = T$



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

Deep Networks

- Extremely flexible, parametric, function approximation

Deep Networks

- Extremely flexible, parametric, function approximation
- **Single layer:** linear transformation, pointwise nonlinearity

Deep Networks

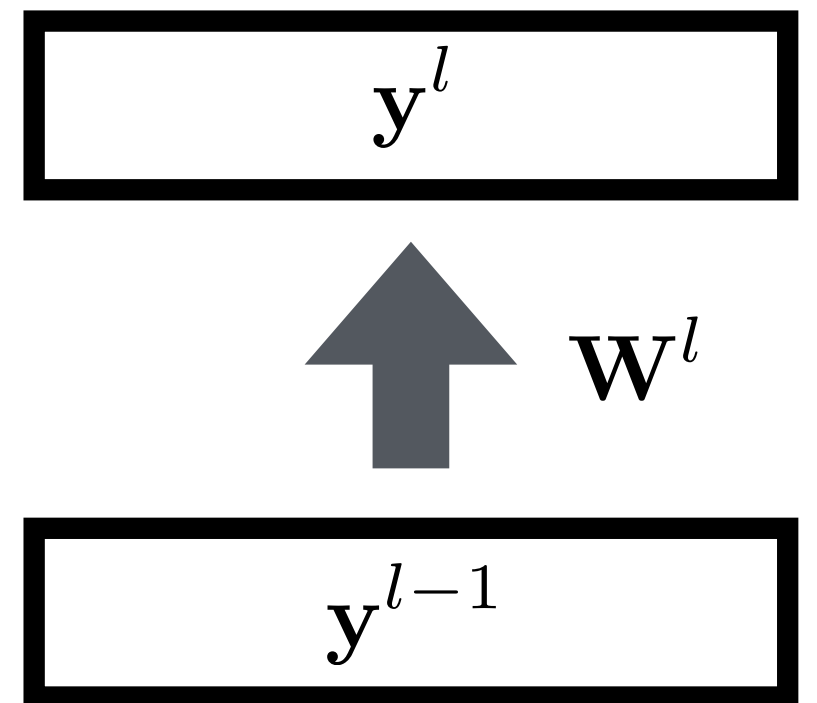
- Extremely flexible, parametric, function approximation
- **Single layer:** linear transformation, pointwise nonlinearity

$$\mathbf{y}^l = \sigma(\mathbf{W}^l \mathbf{y}^{l-1})$$

Deep Networks

- Extremely flexible, parametric, function approximation
- **Single layer:** linear transformation, pointwise nonlinearity

$$\mathbf{y}^l = \sigma(\mathbf{W}^l \mathbf{y}^{l-1})$$



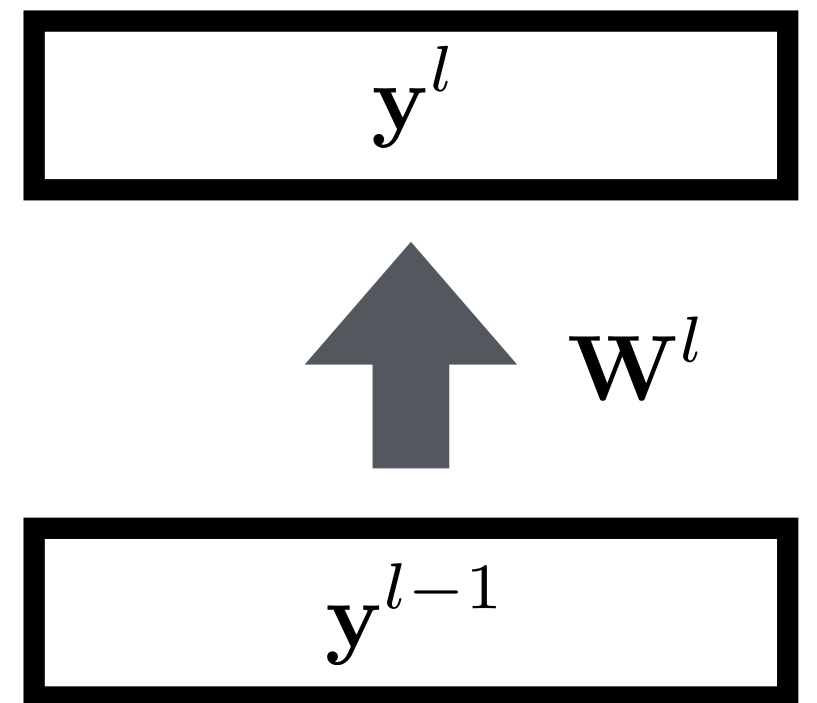
Deep Networks

- Extremely flexible, parametric, function approximation
- **Single layer:** linear transformation, pointwise nonlinearity

$$\mathbf{y}^l = \sigma(\mathbf{W}^l \mathbf{y}^{l-1})$$

$$\sigma(u) \equiv \text{leaky ReLU}$$

$$= \begin{cases} u & u \geq 0 \\ 0.05u & u < 0 \end{cases}$$



Deep Networks

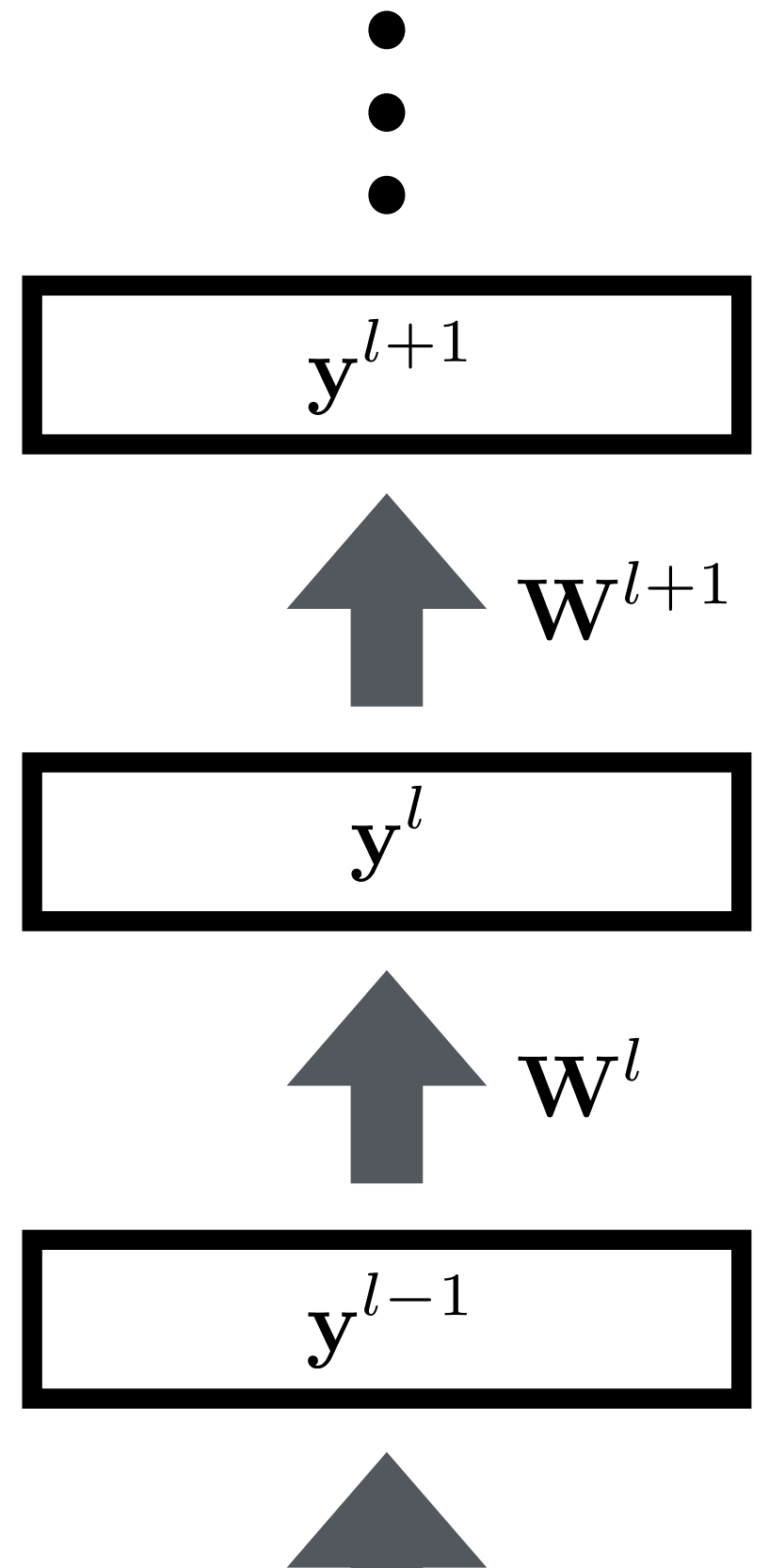
- Extremely flexible, parametric, function approximation
- **Single layer:** linear transformation, pointwise nonlinearity

Deep Networks

- Extremely flexible, parametric, function approximation
- **Single layer:** linear transformation, pointwise nonlinearity
- **Deep network:** stack single layers

Deep Networks

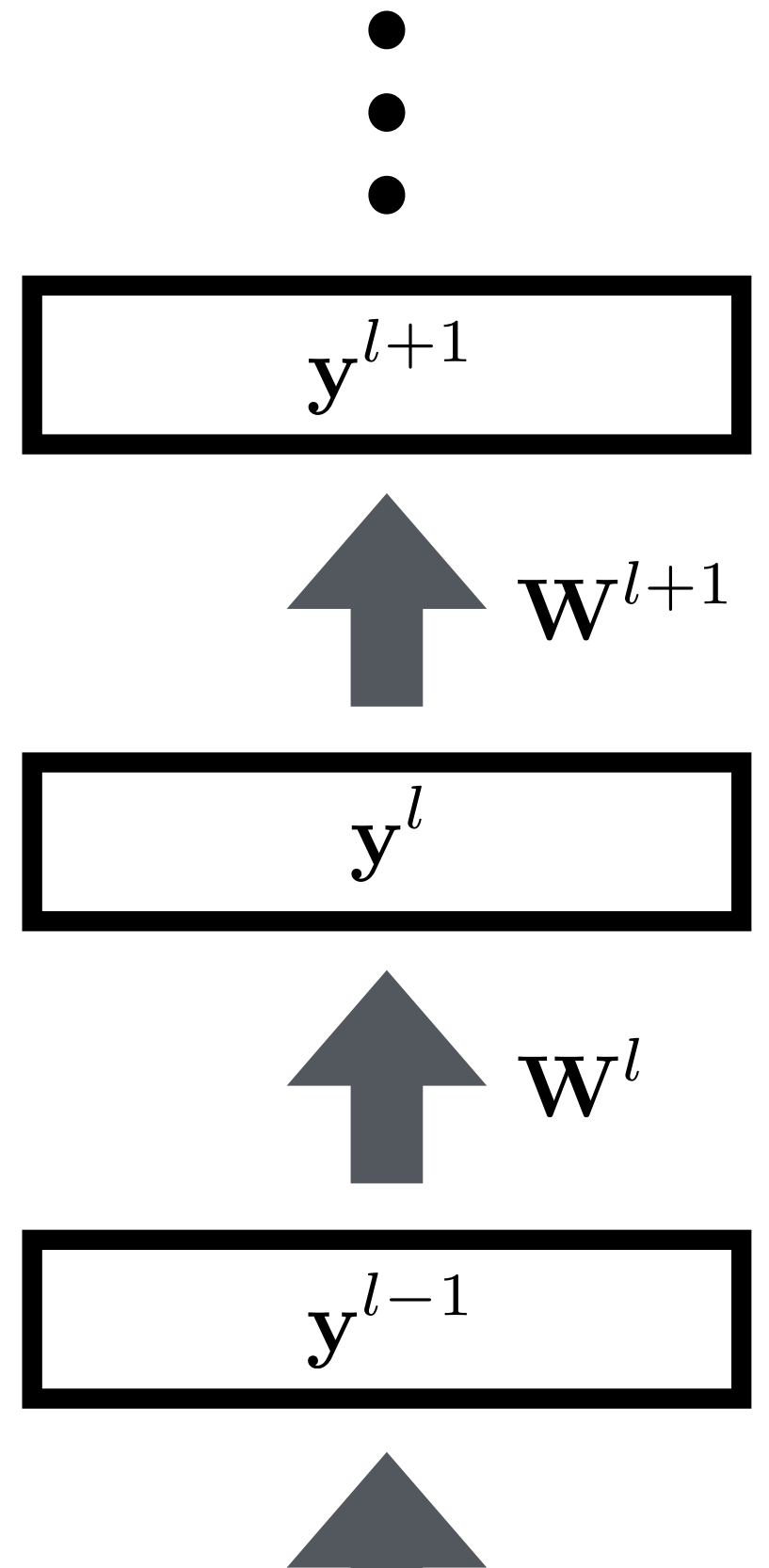
- Extremely flexible, parametric, function approximation
- **Single layer:** linear transformation, pointwise nonlinearity
- **Deep network:** stack single layers



Deep Networks

- Extremely flexible, parametric, function approximation
- **Single layer:** linear transformation, pointwise nonlinearity
- **Deep network:** stack single layers

$$\mathbf{y}^L = \sigma(\mathbf{W}^L \sigma(\mathbf{W}^{L-1} \dots \sigma(\mathbf{W}^1 \mathbf{y}^0)))$$

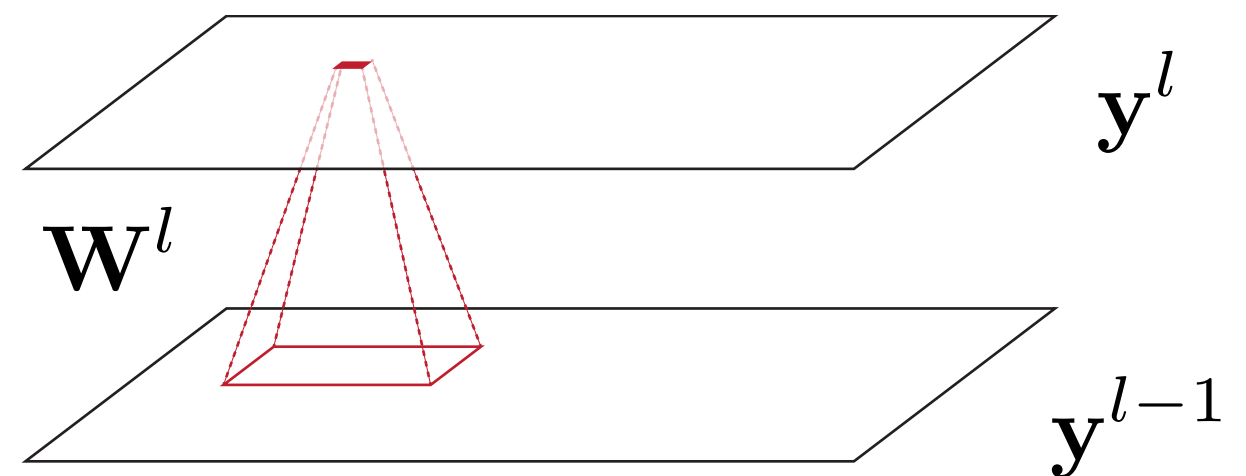


Convolutional Neural Network

- Single convolutional layer:
 - Same linear transform for every pixel
 - Pointwise nonlinearity

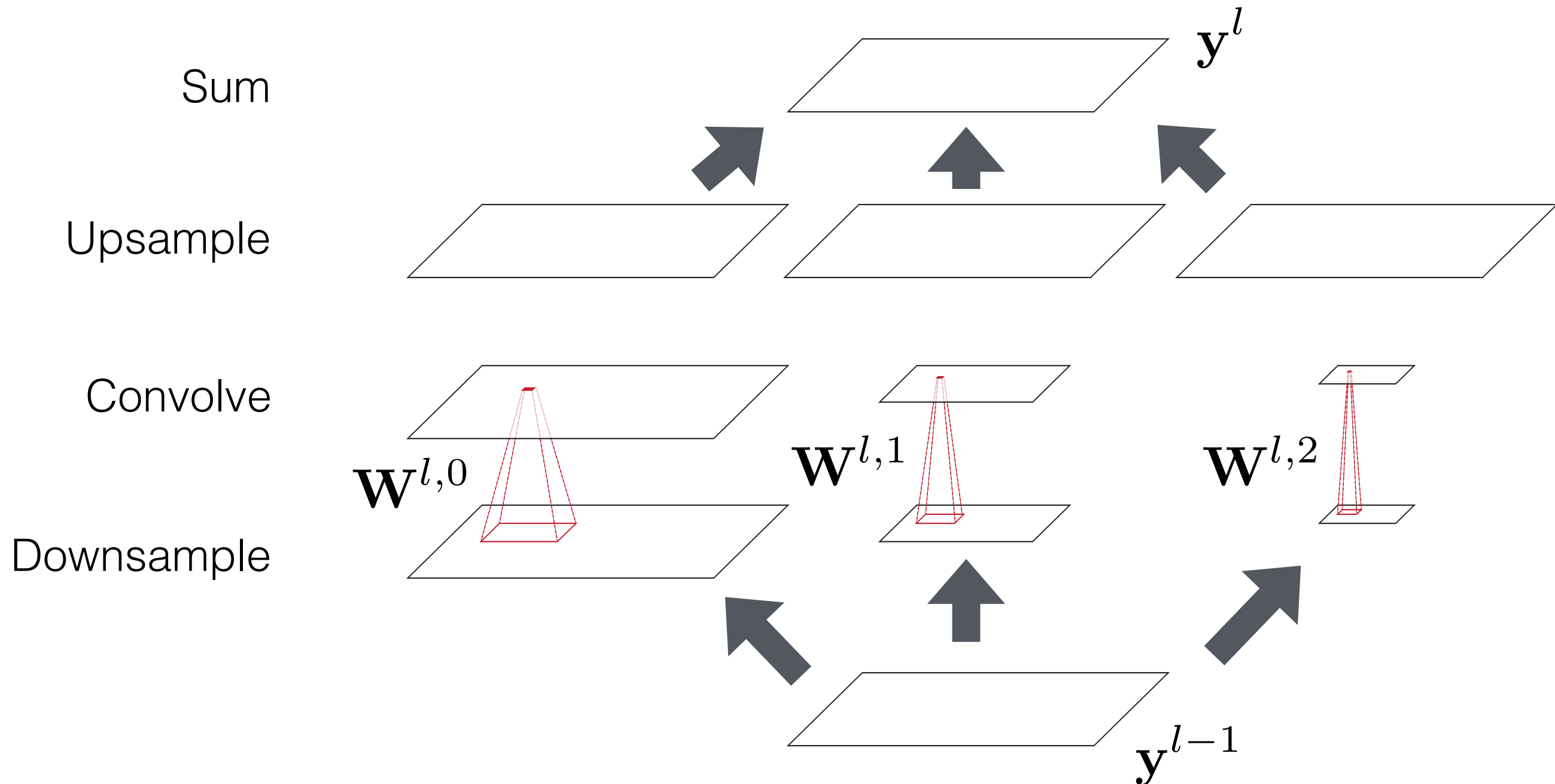
Convolutional Neural Network

- Single convolutional layer:
 - Same linear transform for every pixel
 - Pointwise nonlinearity



Multiscale Convolution

- Single multi-scale convolutional layer:



Deep Network Architecture for Diffusion

$$f_{\mu}(\mathbf{x}^{(t)}, t)$$

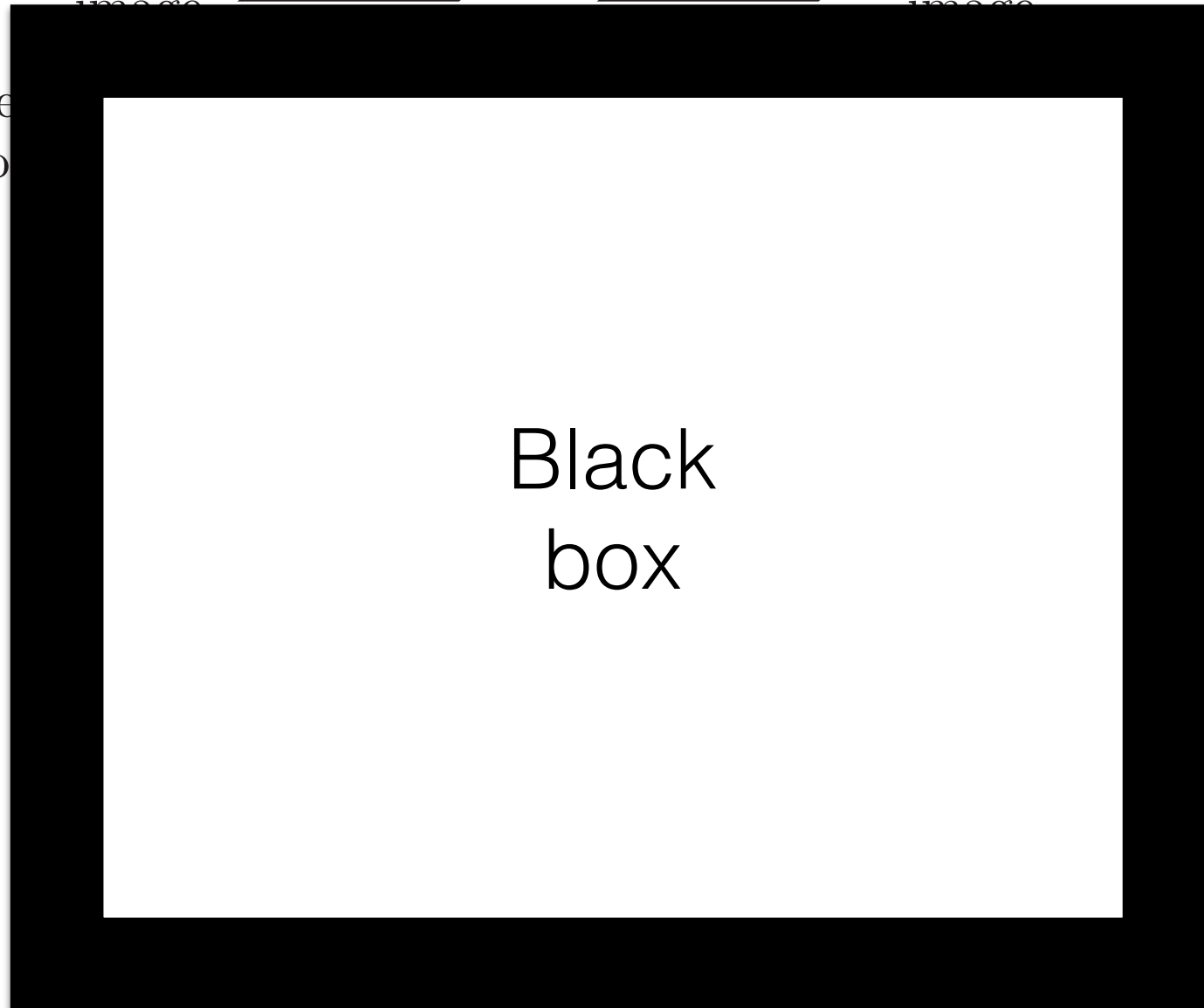
Mean



Covariance

$$f_{\Sigma}(\mathbf{x}^{(t)}, t)$$

Te
co



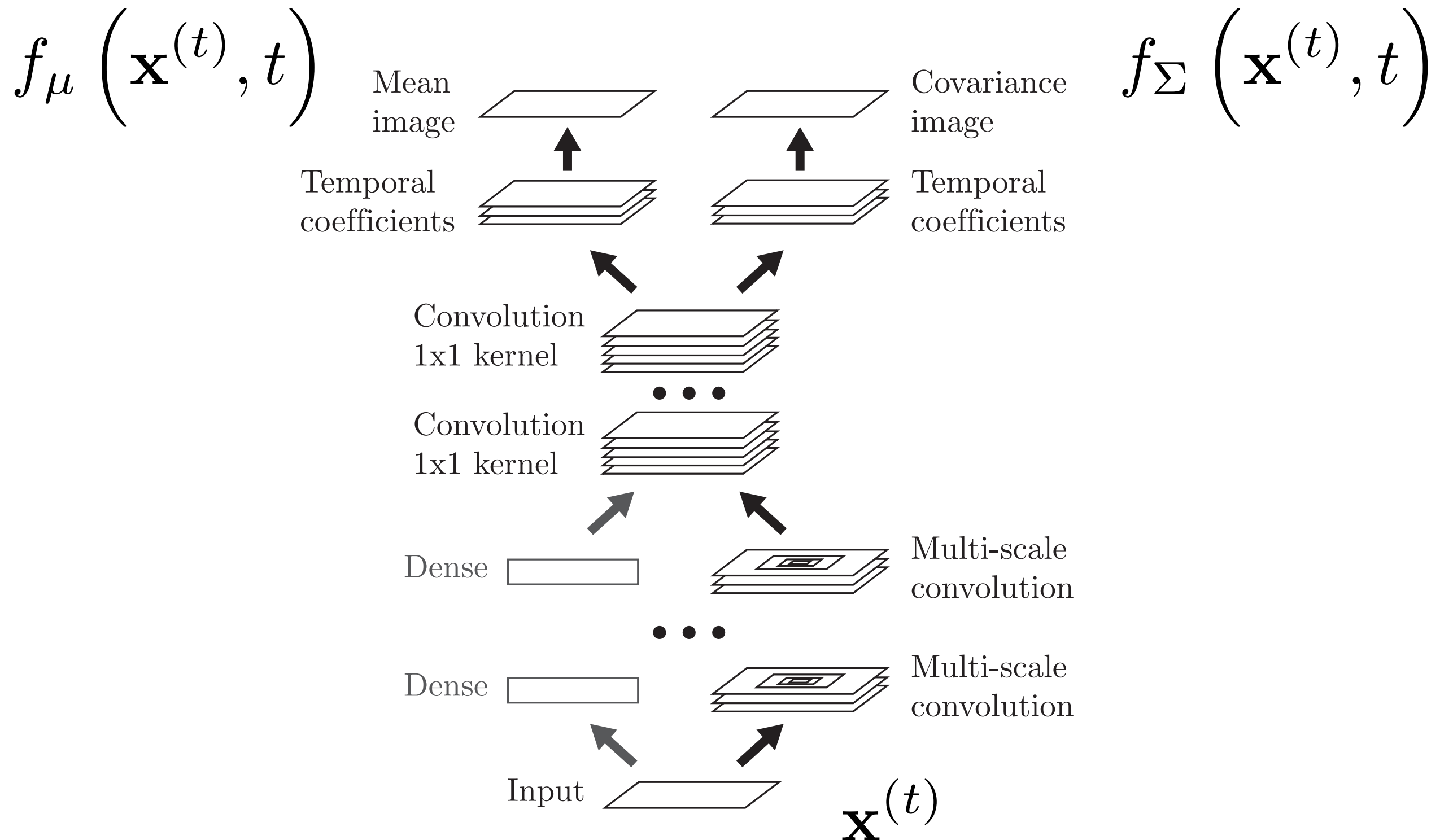
Black
box

Input

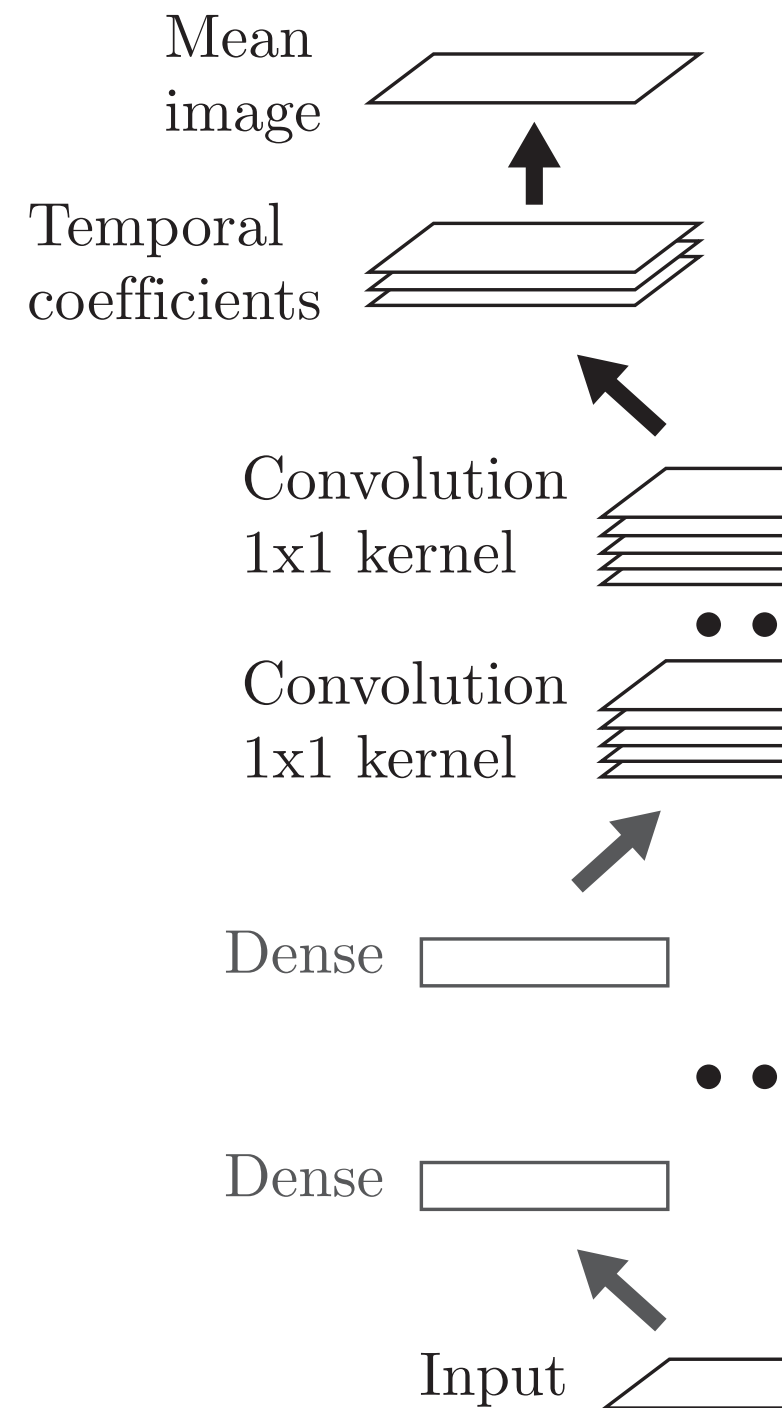


$\mathbf{x}^{(t)}$

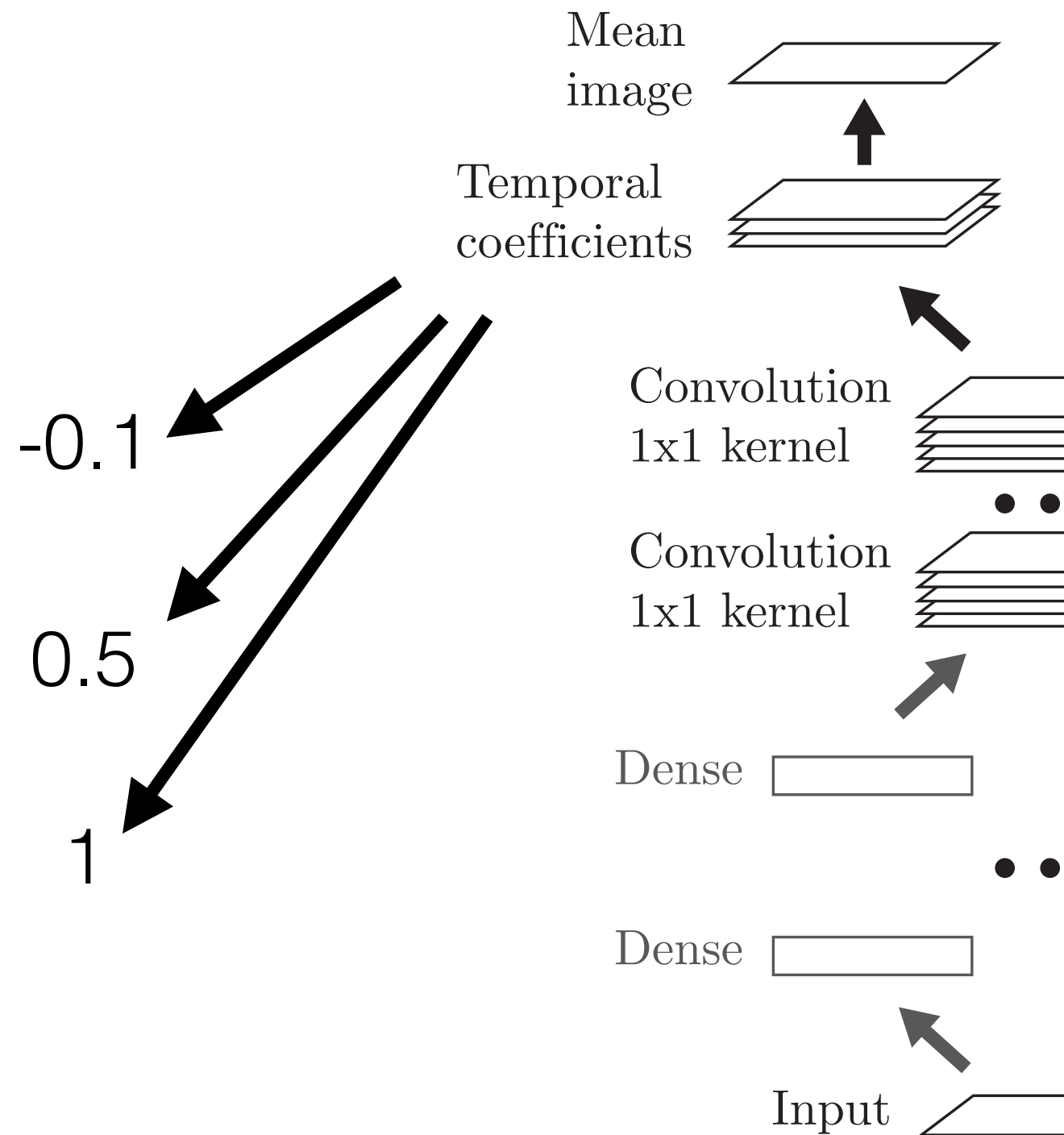
Deep Network Architecture for Diffusion



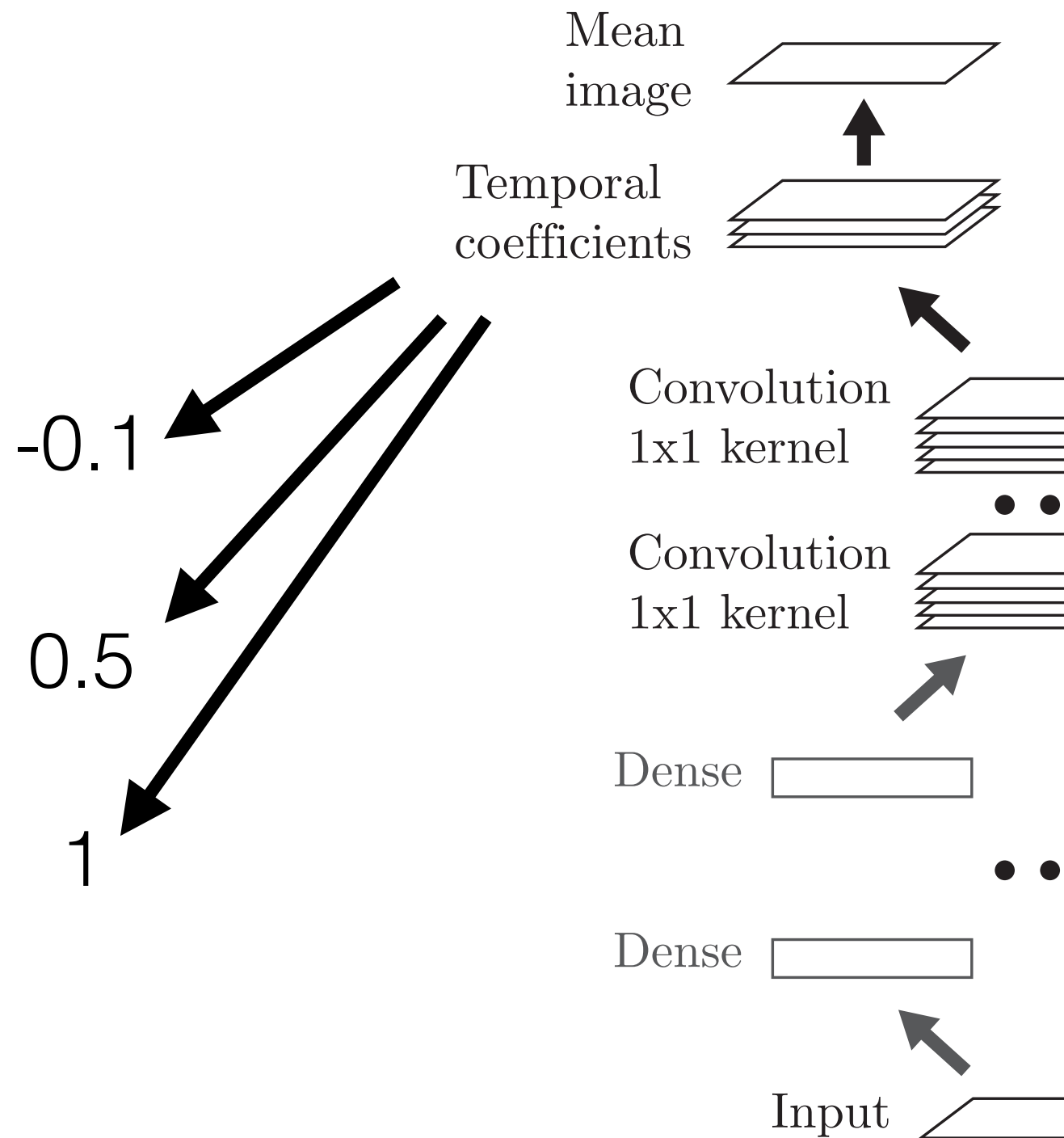
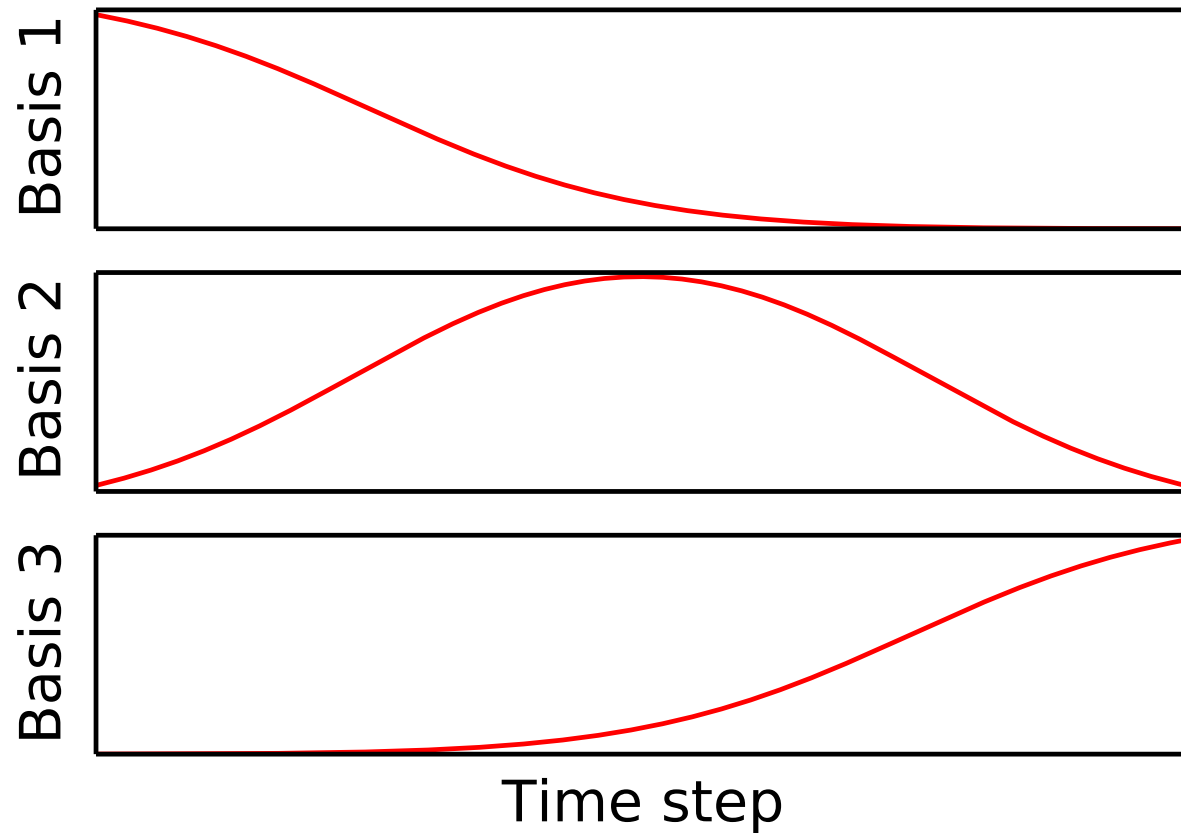
Time Dependence using Temporal Basis



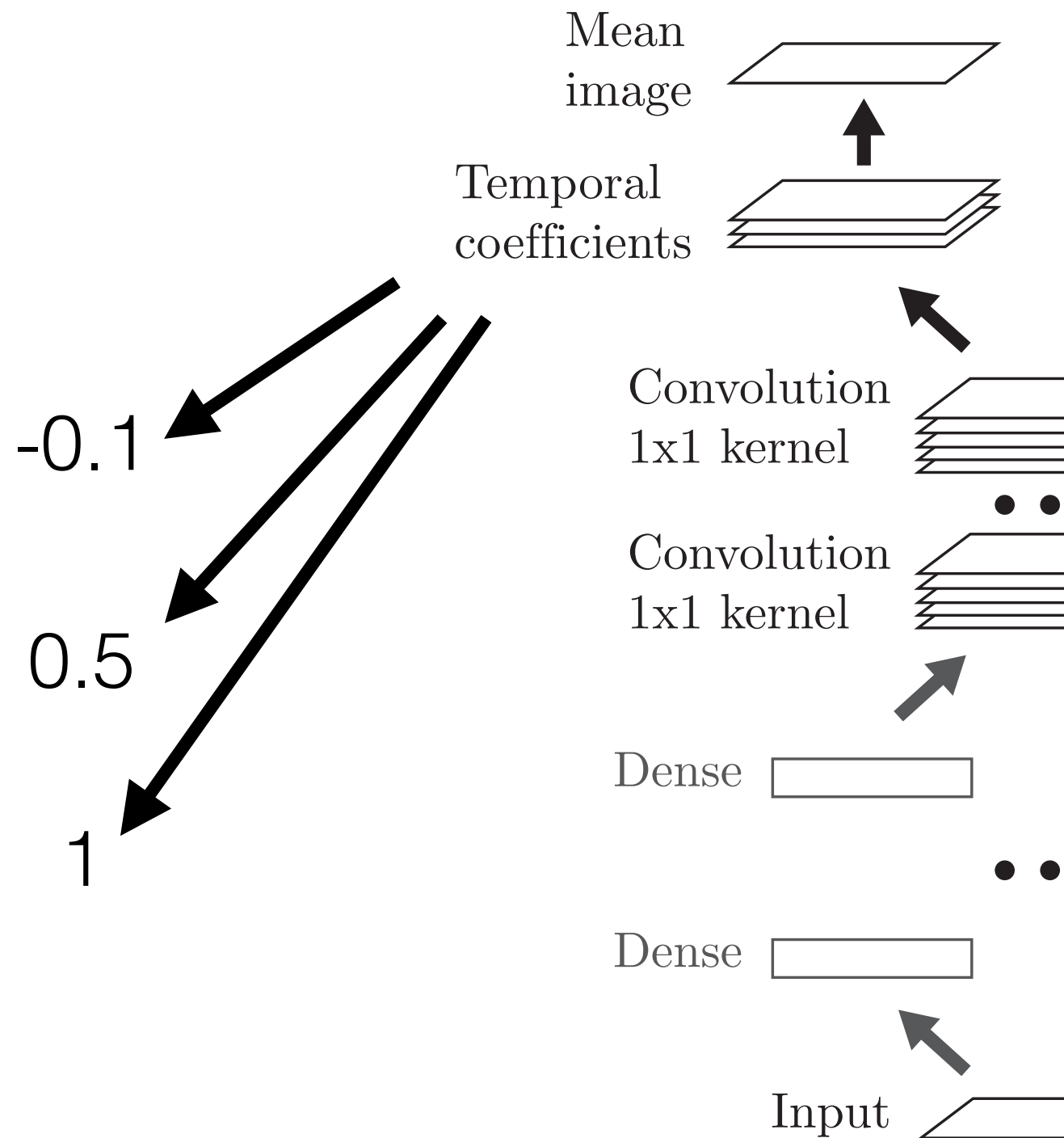
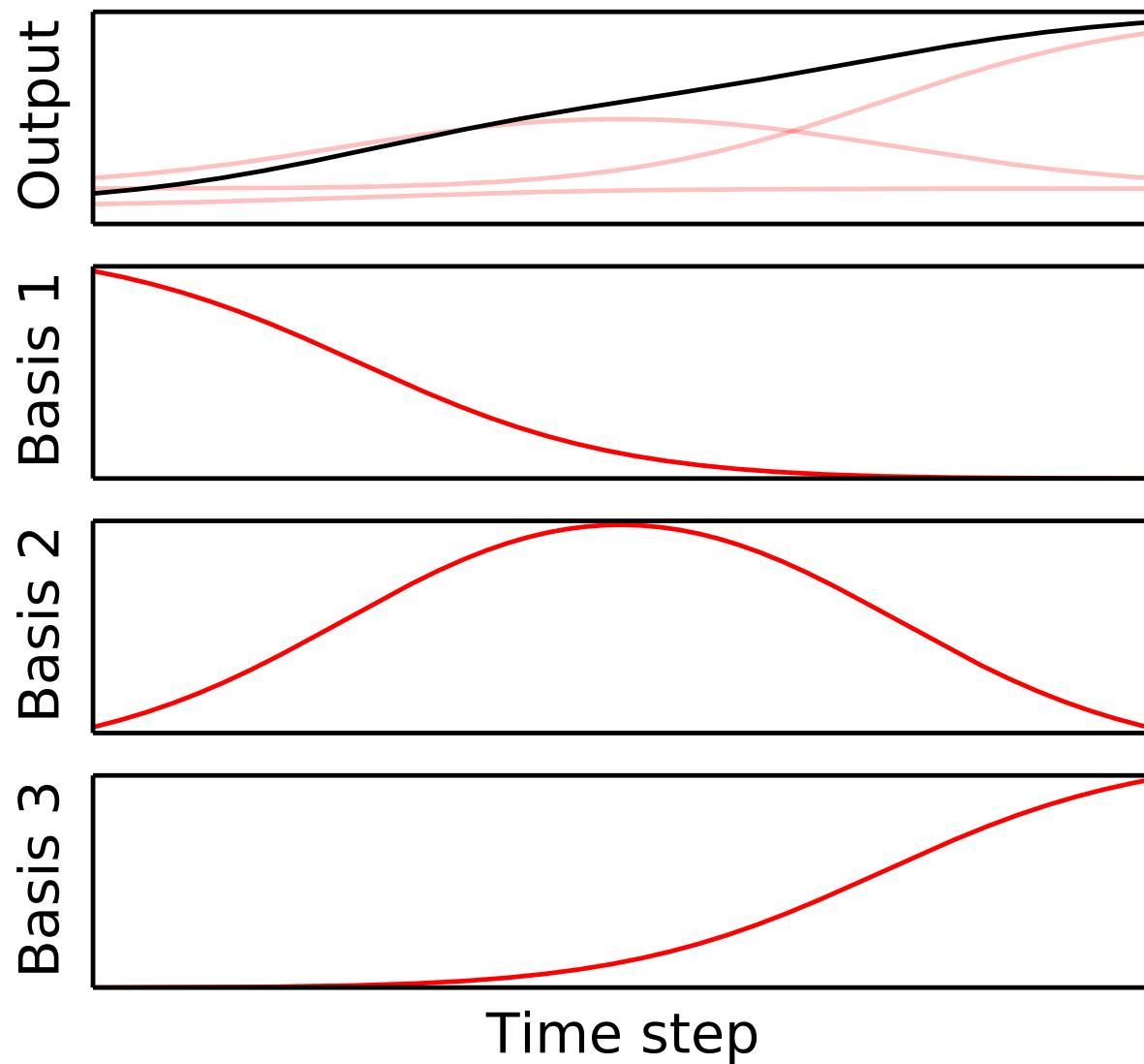
Time Dependence using Temporal Basis



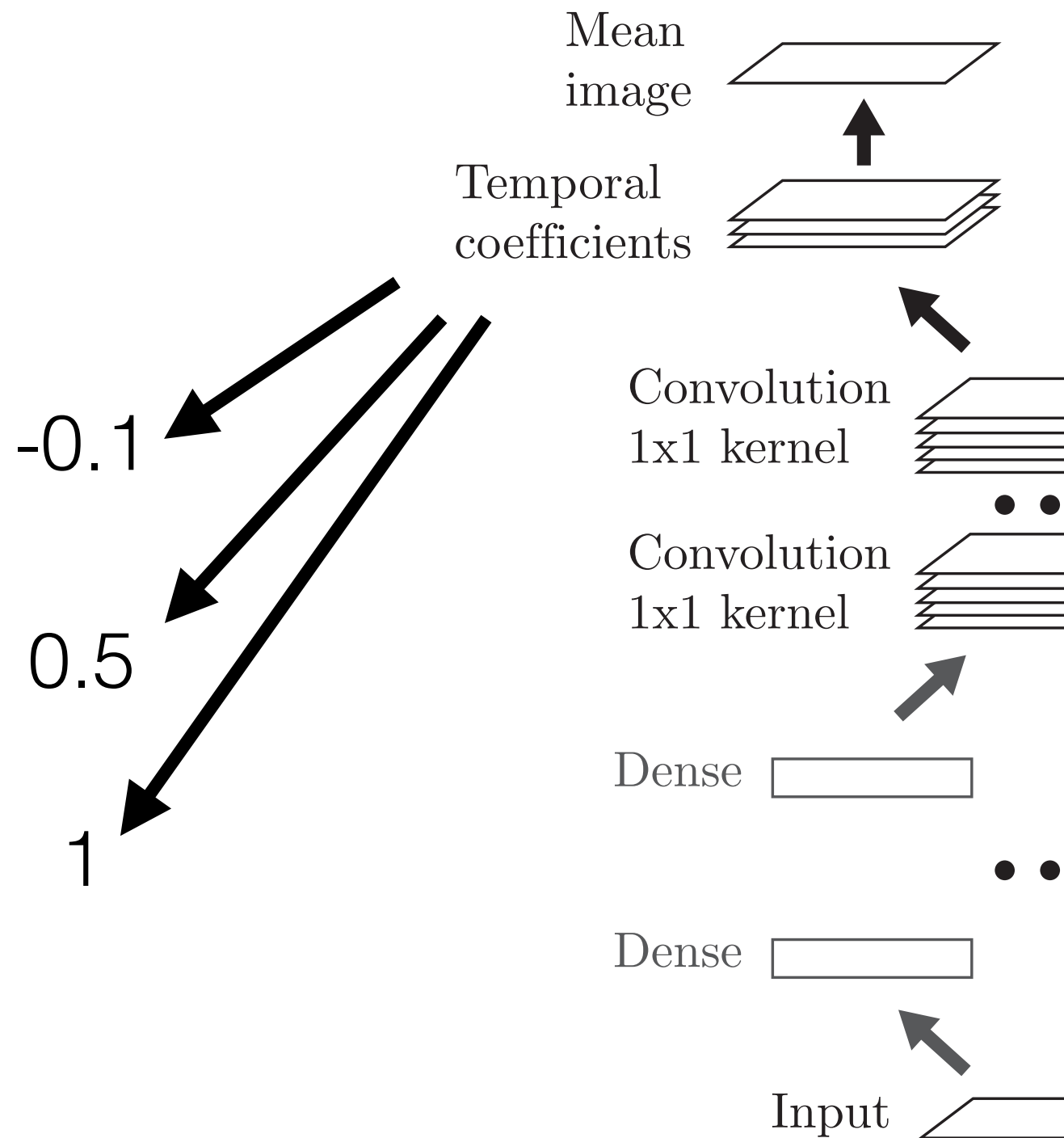
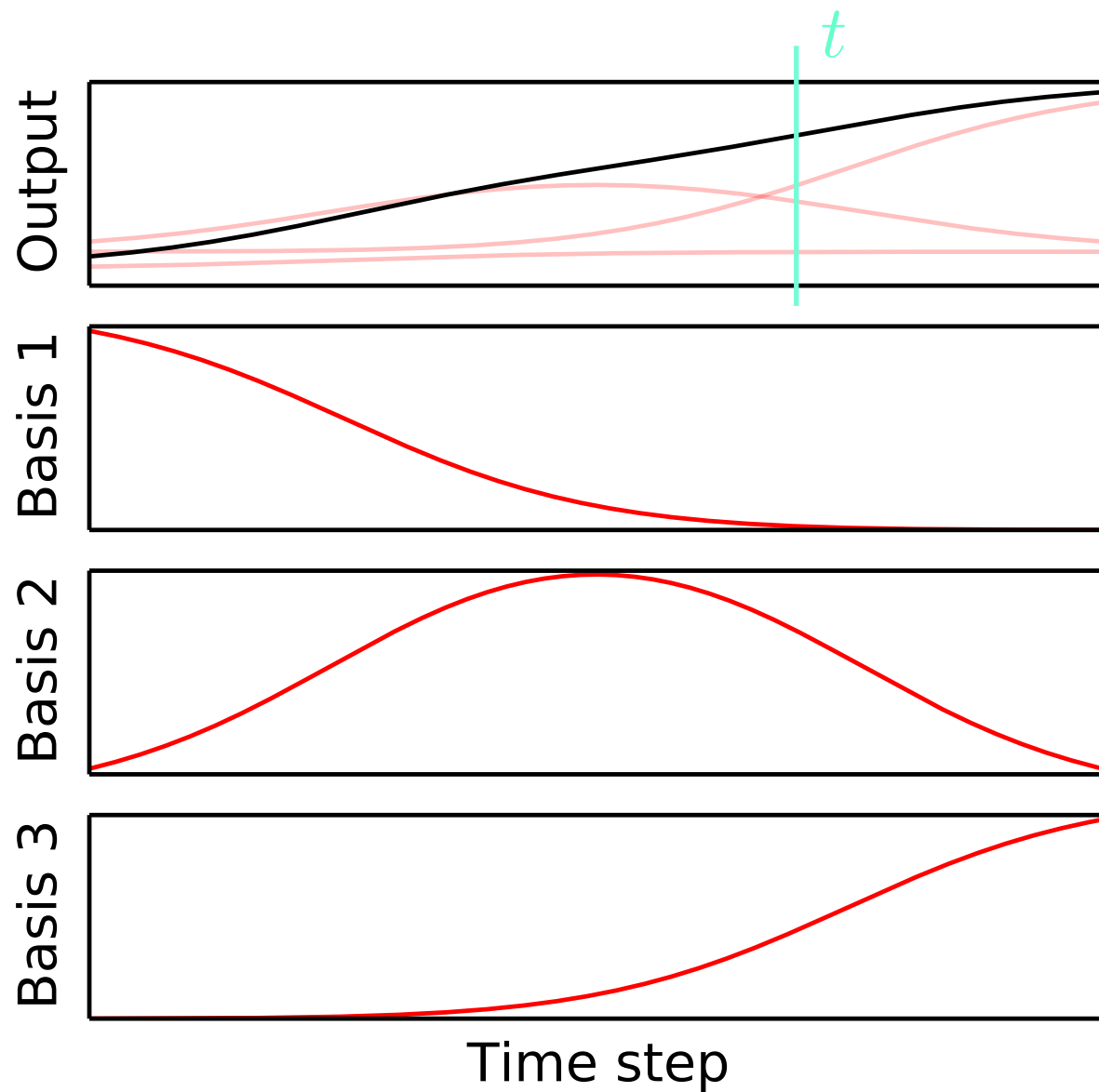
Time Dependence using Temporal Basis



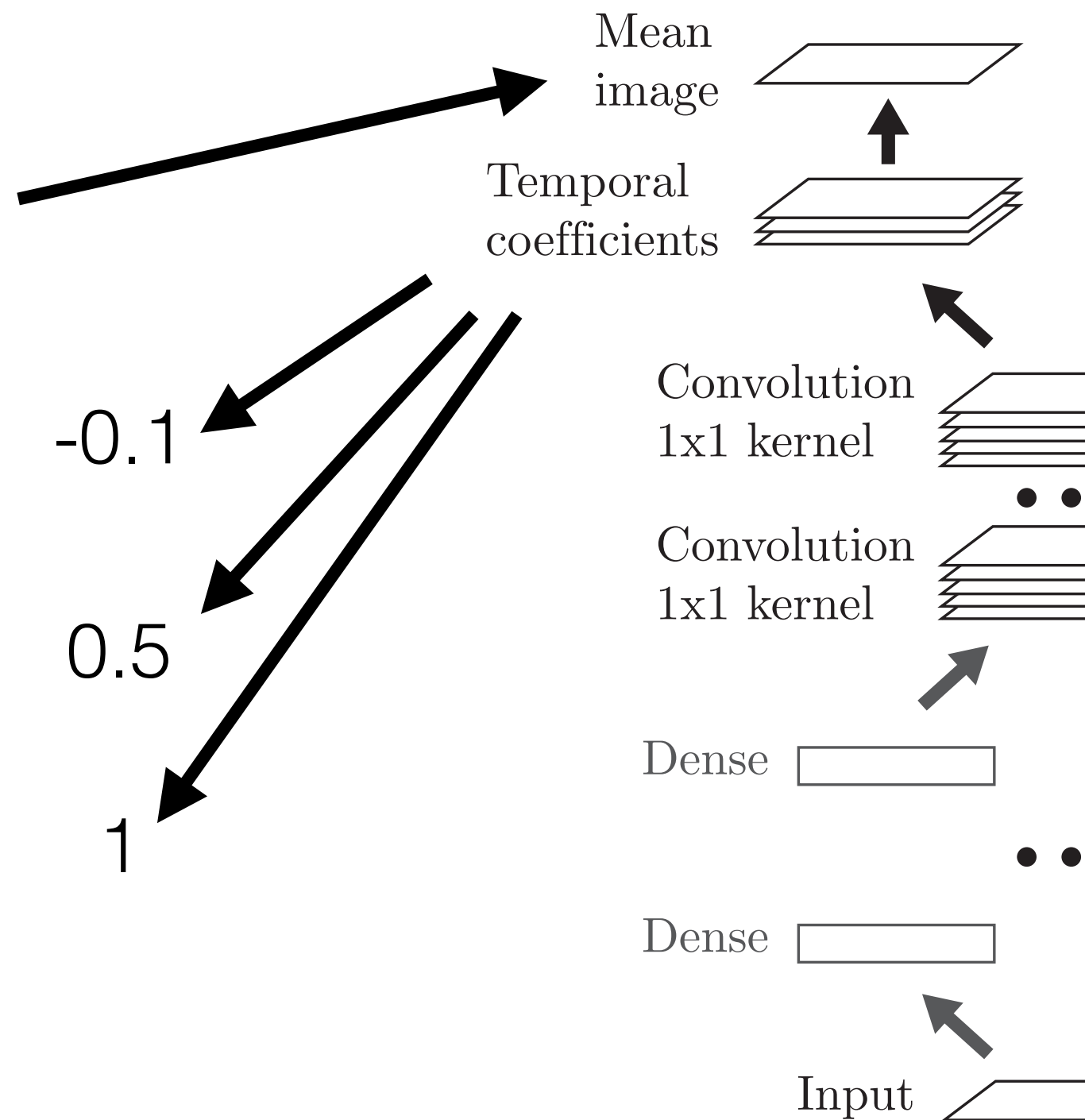
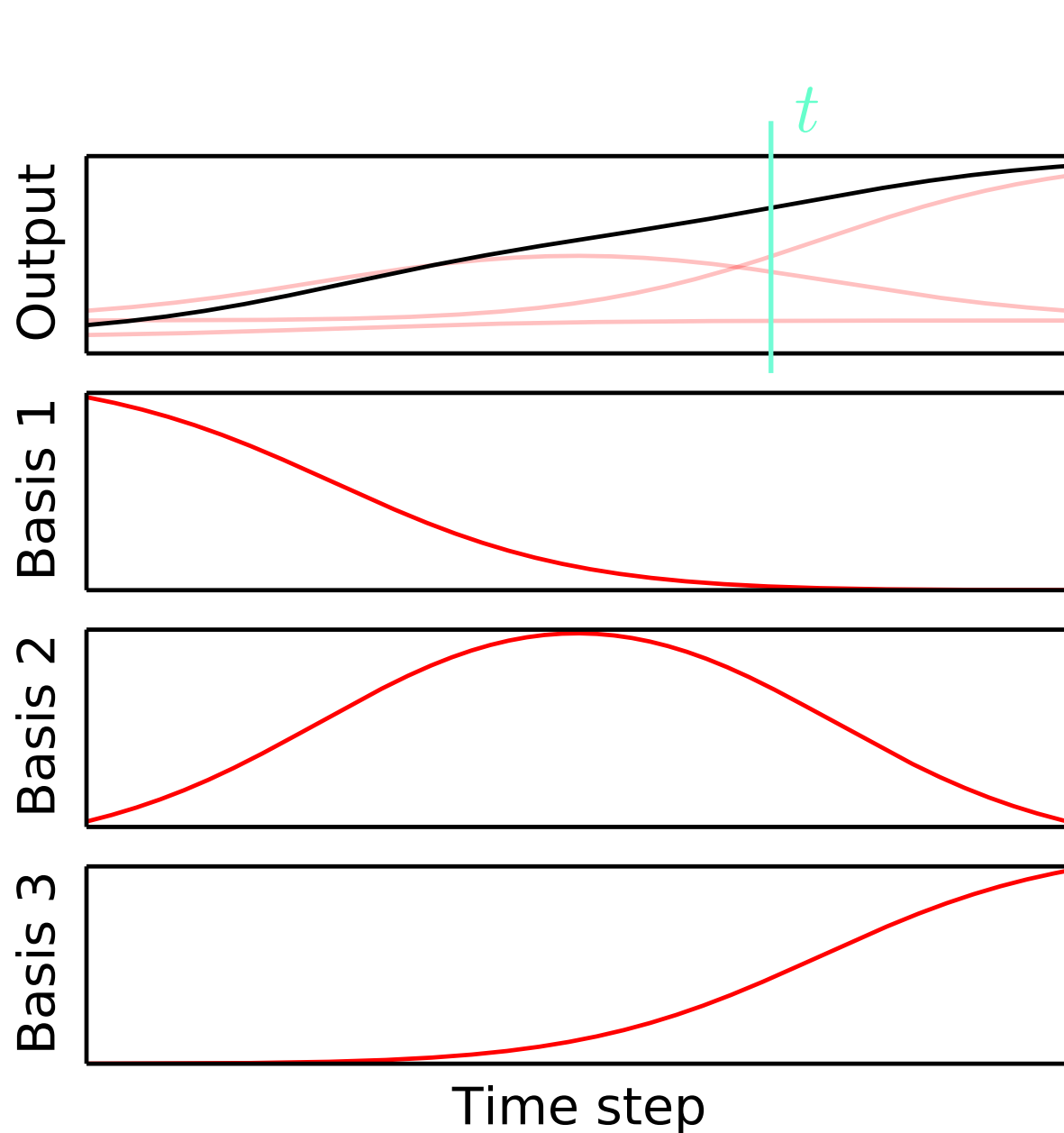
Time Dependence using Temporal Basis



Time Dependence using Temporal Basis



Time Dependence using Temporal Basis



Setting Diffusion Rate

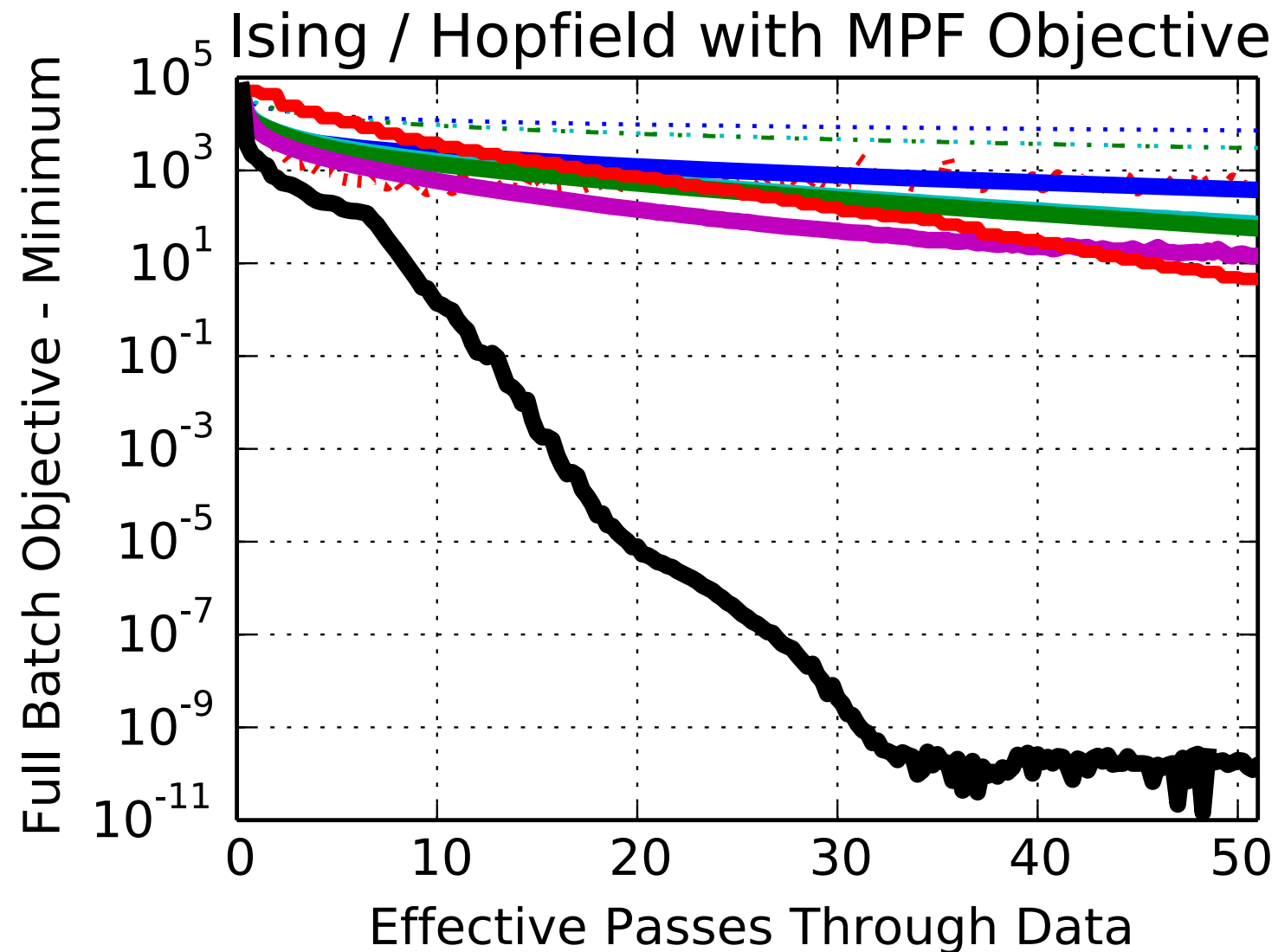
- Erase constant fraction of stimulus variance each step

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

- Can also train β_t

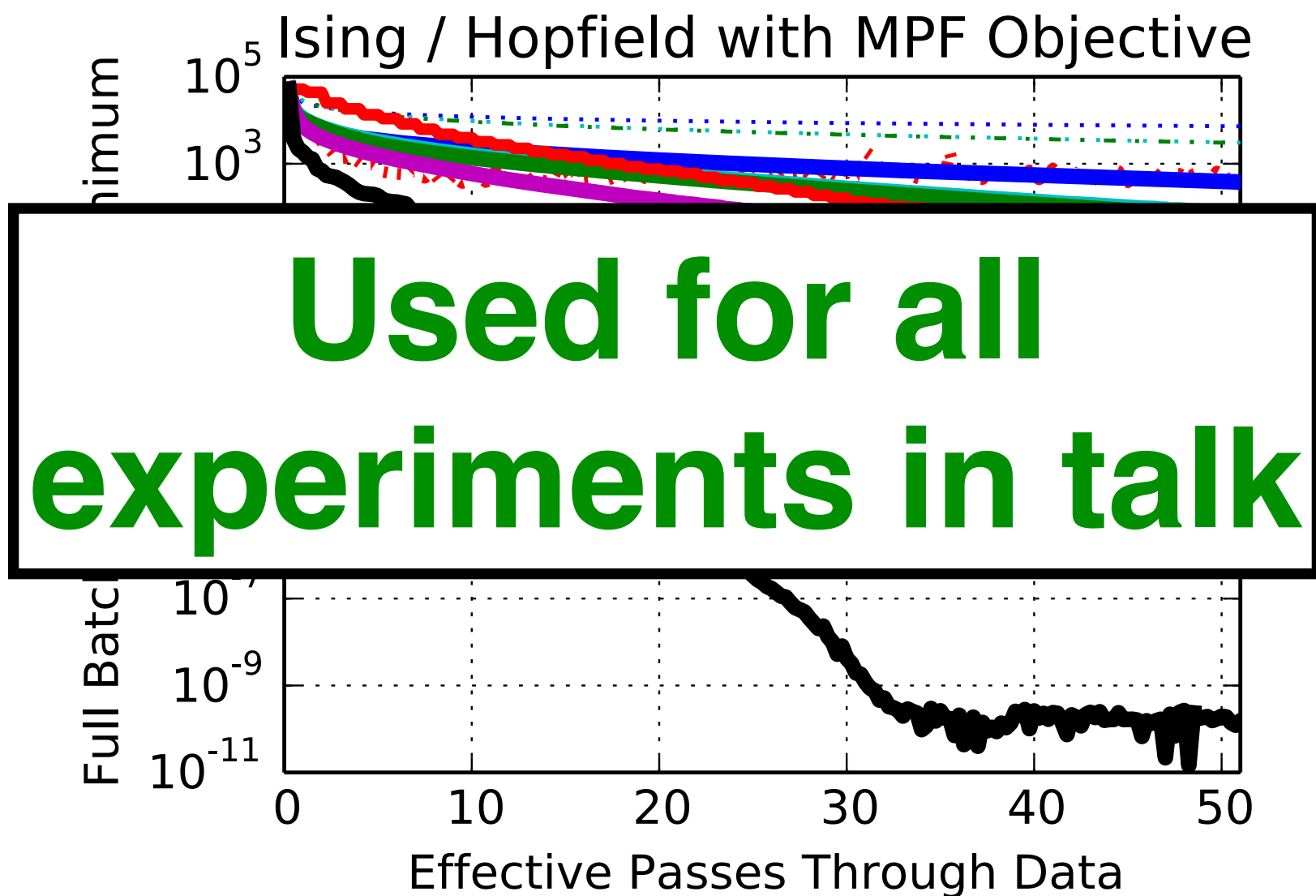
Theoretical Breakthroughs in Machine Learning

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



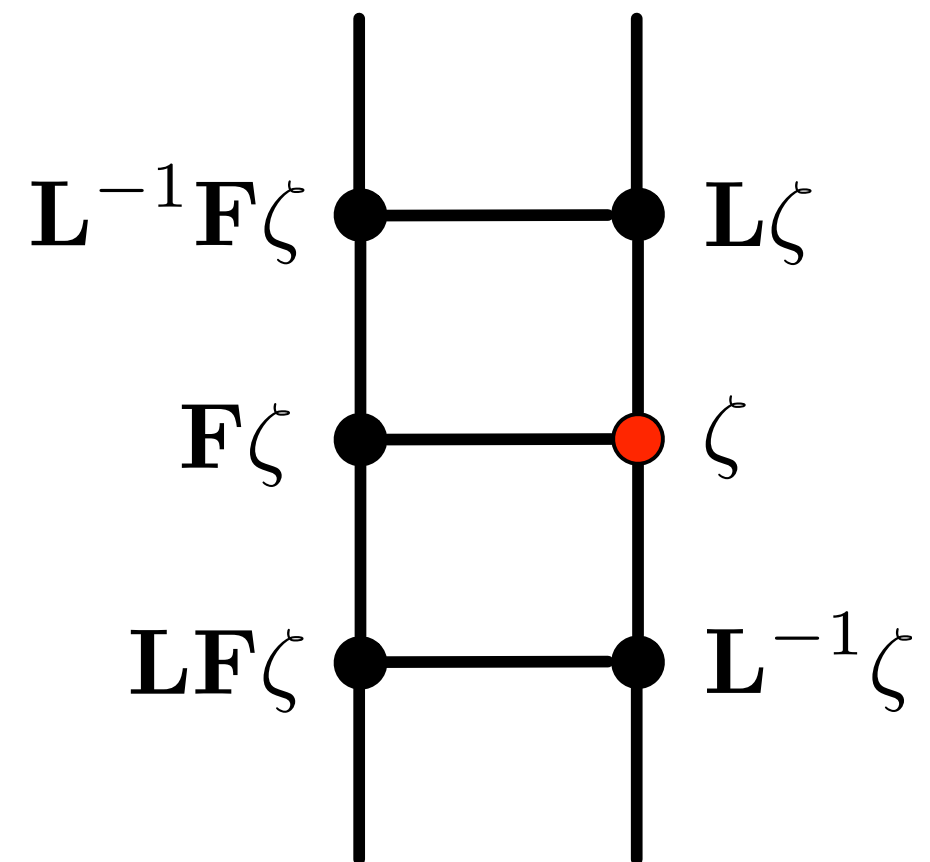
Theoretical Breakthroughs in Machine Learning

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



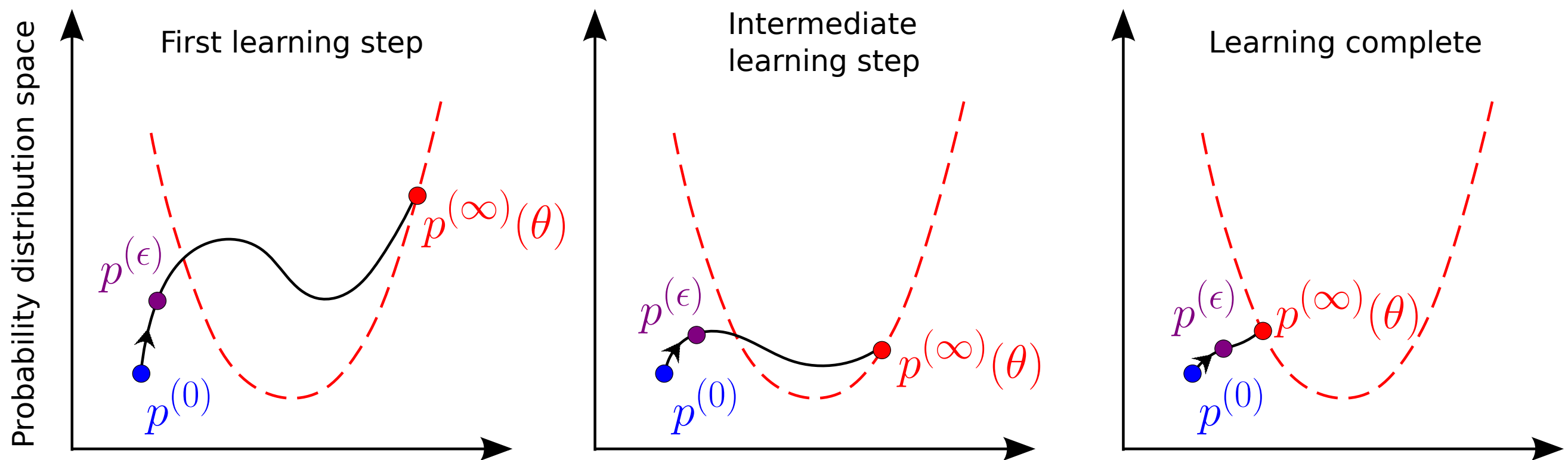
Theoretical Breakthroughs in Machine Learning

- **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]



Theoretical Breakthroughs in Machine Learning

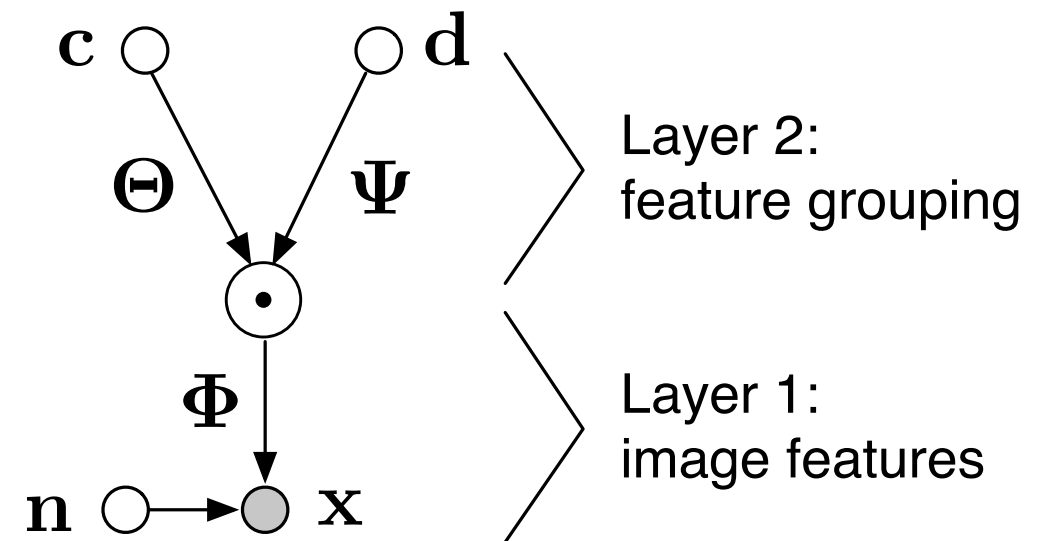
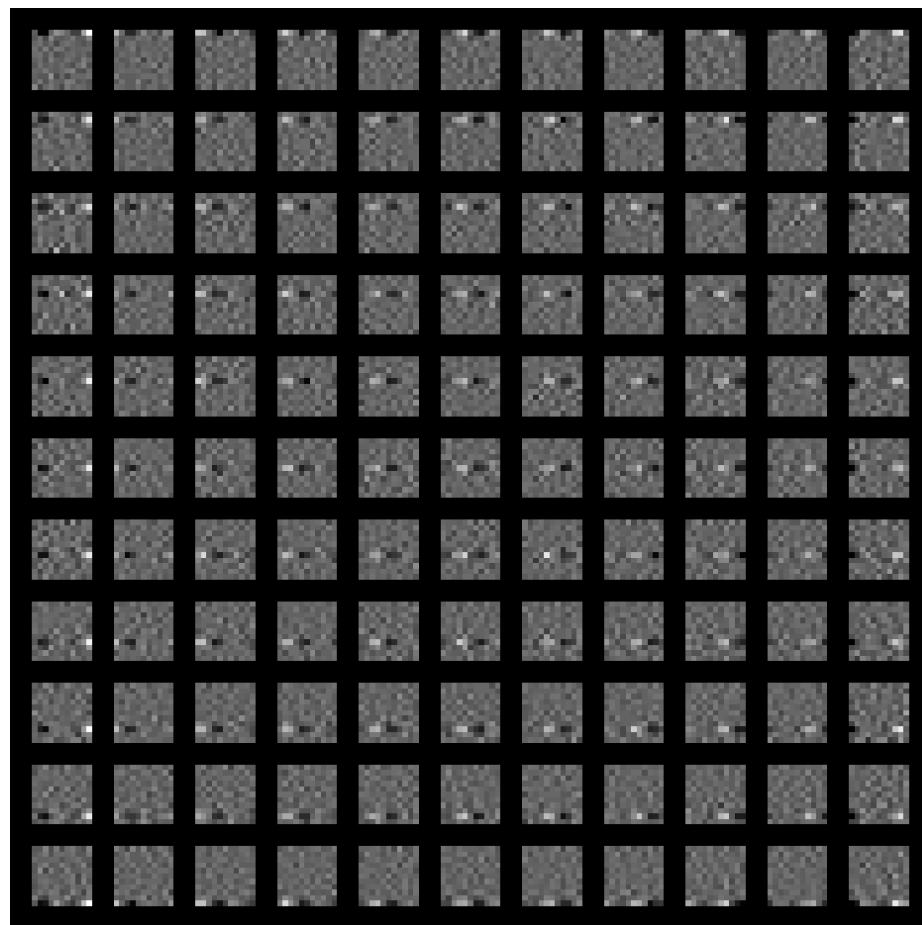
- **Training energy-based models:** Minimum Probability Flow learning [ICML 2011] [PRL 2011]



Theoretical Breakthroughs in Machine Learning

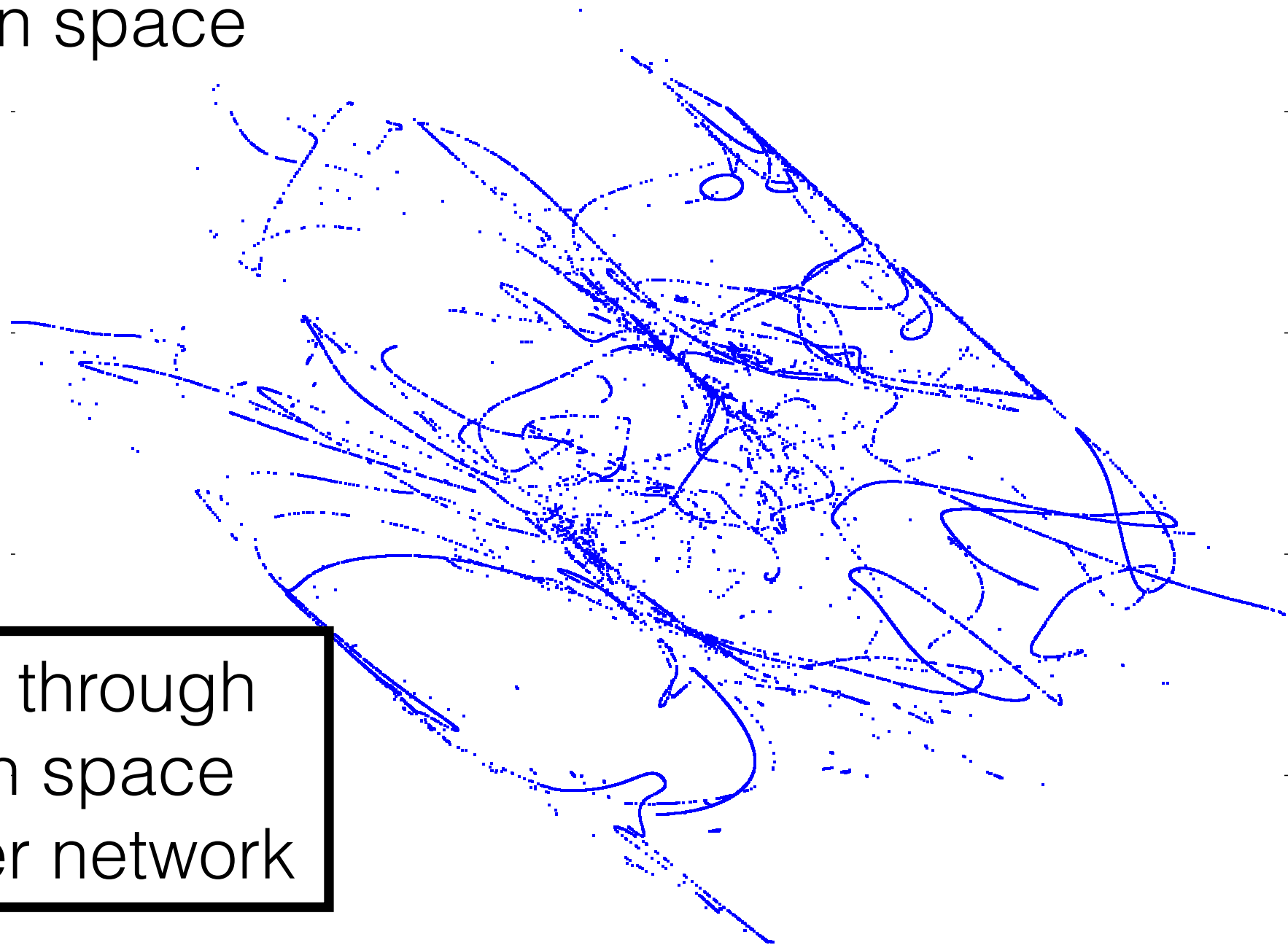
- **Model design:** capturing dynamics with Lie groups [Under Revision at **NECO**], bilinear generative models [ICCV 2011]

Horizontal Translation



Theoretical Breakthroughs in Machine Learning

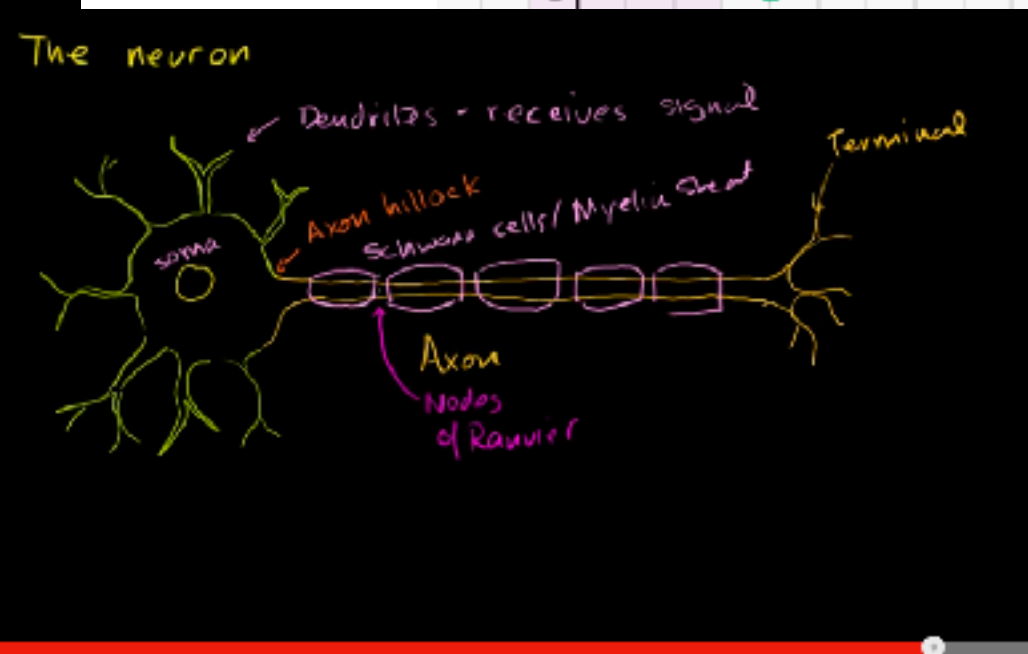
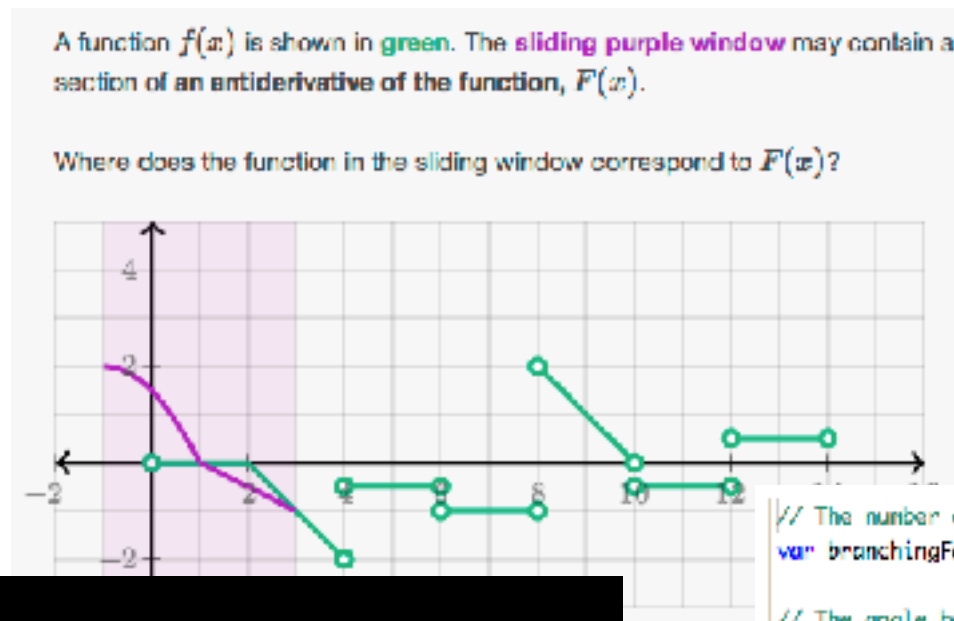
- **Properties of deep networks:** Characterization in function space



2d slice through
function space
for 2-layer network

Understanding Real Data

- Online education data



```
// The number of branches each branch splits into
var branchingFactor = 3;

// The angle between the branches in degrees
var angleBetweenBranches = 30;

// Controls how much smaller each level of the tree gets
var scaleFactor = 0.7;

// The number of levels of the tree drawn
var numLevels = 6;

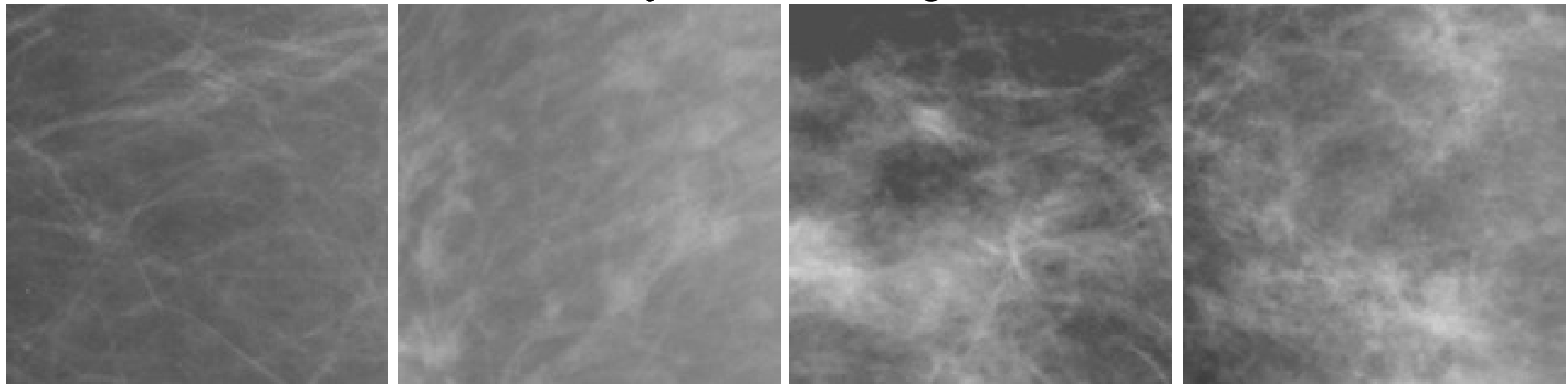
// The length of the branches
var baseBranchLength = 80;

var forward = function(distance) {
  line(0, 0, 0, -distance);
  translate(0, -distance);
}
```

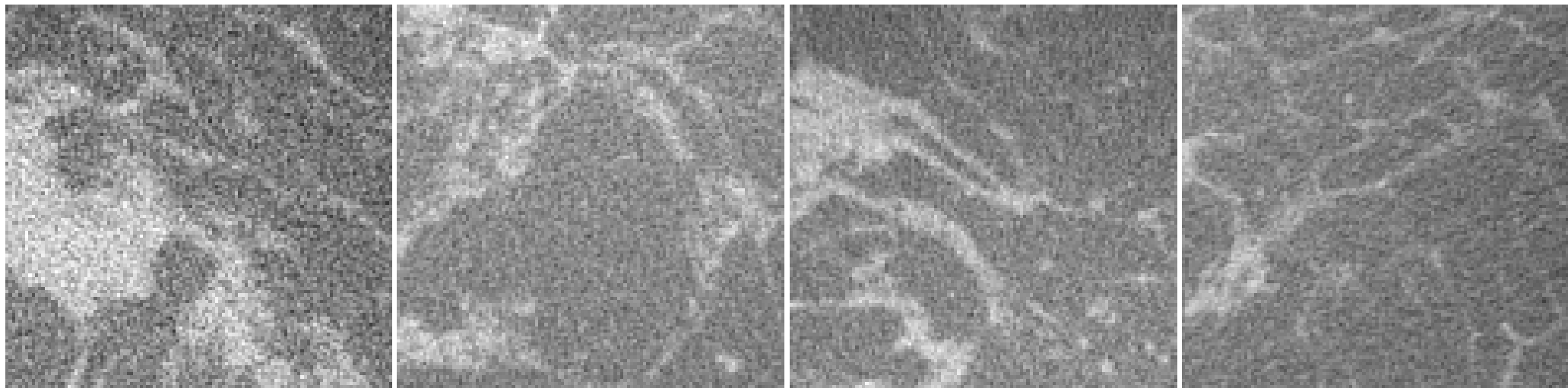
Understanding Real Data

- [Medical imaging data](#) [SPIE 2009] [Med Phys 2014]

A. Projection Mammograms



B. Coronal Breast CT



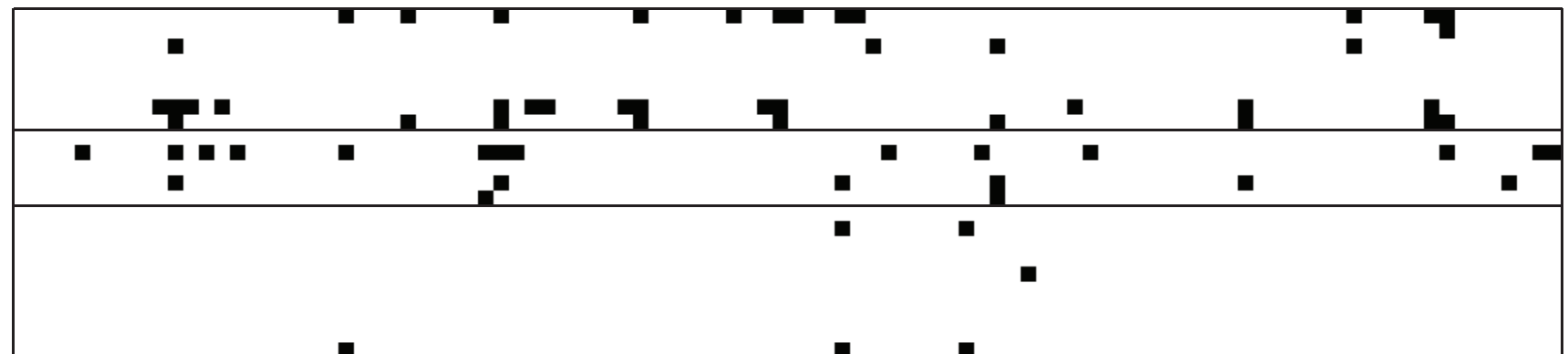
Understanding Real Data

- [Neuroscience electrophysiology data](#): [PLoS Comp Bio 2014]
[Neuron 2013]

a) Stimulus frames

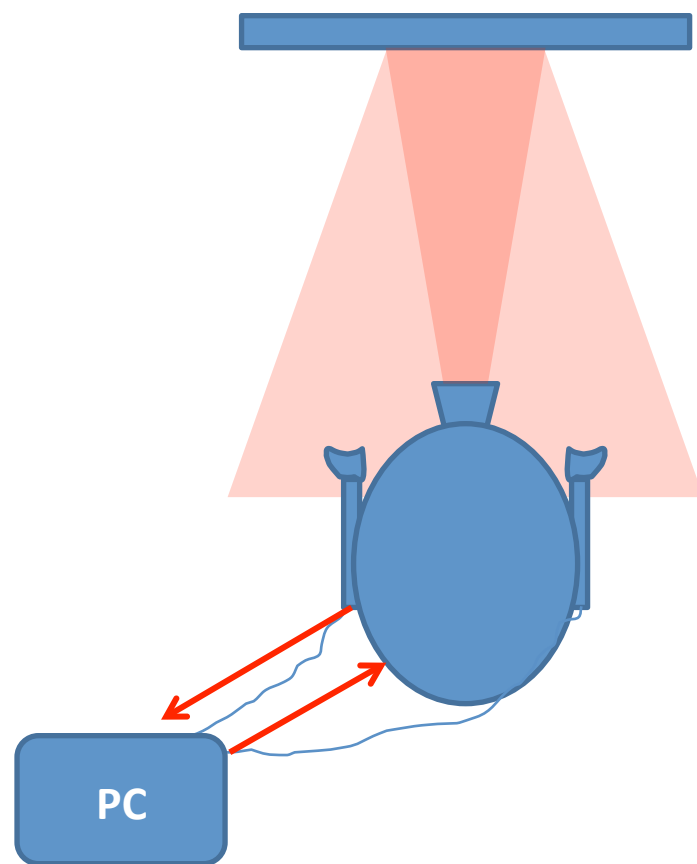


b) Example data, 2s of data in 20ms bins



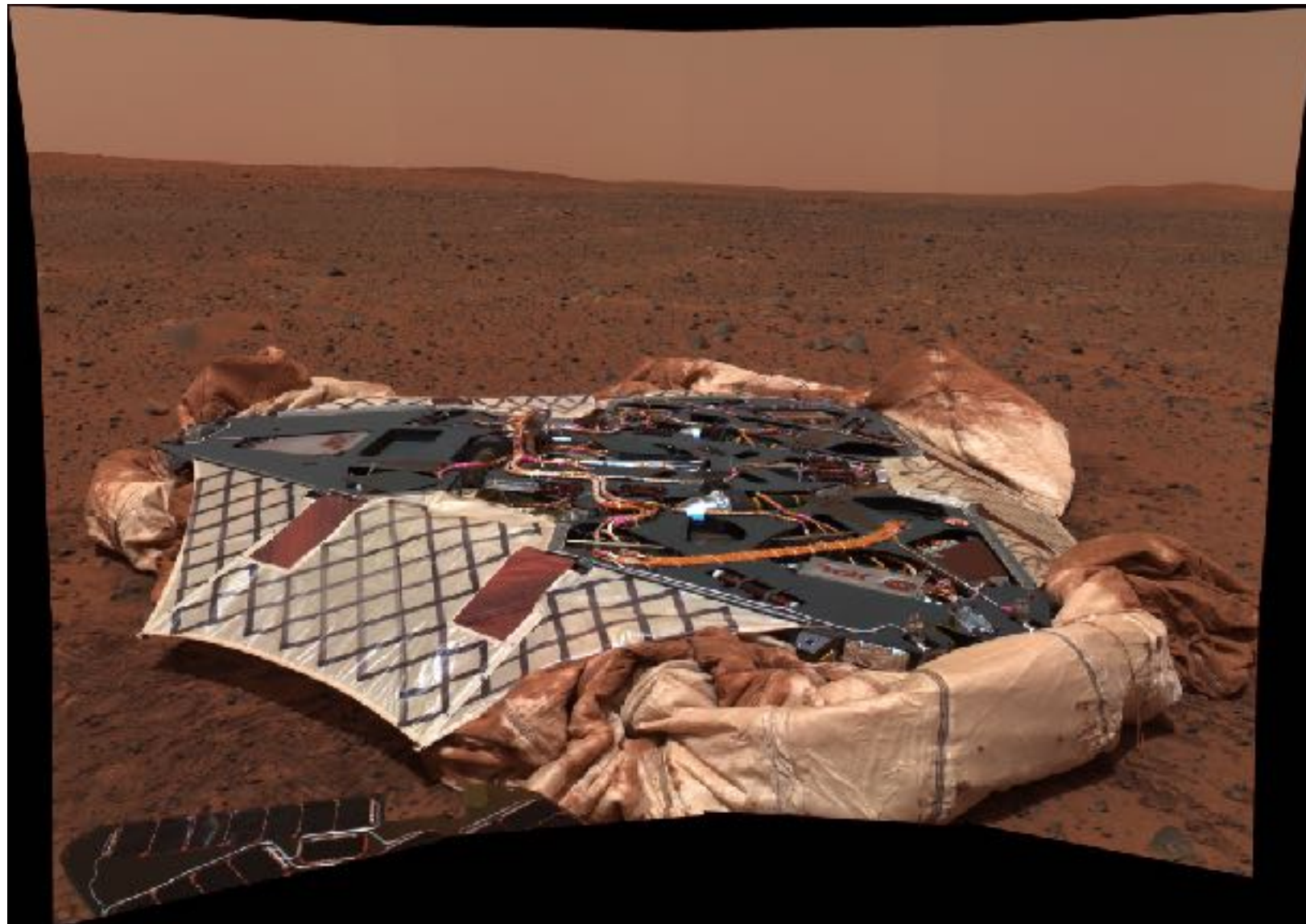
Understanding Real Data

- [Human ultrasonic echolocation](#): Blind assistive device [TBME 2015]



Understanding Real Data

- **Planetary science:** multispectral observations
[Science 2004a] [Science 2004b]



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

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

Continuous Time

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

$$p(\mathbf{x}^t | \mathbf{x}^{t+dt}) = \mathcal{N}\left(\mathbf{x}^t; \mathbf{x}^{t+dt} - \mathbf{x}^{t+dt} \frac{\exp(-\beta t)}{1 - \exp(-\beta t)} \beta dt - \frac{1}{2} \mathbf{x}^{t+dt} \beta dt + \frac{1}{2} f_0(\mathbf{x}^{t+dt}, t) \operatorname{csch}\left(\frac{\beta t}{2}\right) \beta dt, \beta dt\right)$$

$$\begin{aligned} D_{KL}(q(\mathbf{x}^t | \mathbf{x}^0, \mathbf{x}^{t+dt}) || p(\mathbf{x}^t | \mathbf{x}^{t+dt})) &= \frac{1}{2} \frac{\Sigma_q}{\Sigma_p} + \frac{1}{2} \log \frac{\Sigma_p}{\Sigma_q} + \frac{1}{2} \frac{(\mu_p - \mu_q)^2}{\Sigma_p} - \frac{1}{2} \\ &= \frac{1}{8} (f_0(\mathbf{x}^{t+dt}, t) - \mathbf{x}^0)^2 \operatorname{csch}^2\left(\frac{\beta t}{2}\right) \beta dt \end{aligned}$$

Denoising autoencoder penalty

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 — \rightarrow 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)