

Deep Unsupervised Learning using Nonequilibrium Thermodynamics

Jascha Sohl-Dickstein¹, Eric Weiss²,
Niru Maheswaranathan¹, Surya Ganguli¹

¹ Stanford University, ² University of California at Berkeley

Outline

- **Motivation:** The promise of deep unsupervised learning
- **Physical intuition:** Diffusion processes and time reversal
- **Diffusion probabilistic model:** Derivation and experimental results

Outline

- **Motivation: The promise of deep unsupervised learning**
- **Physical intuition:** Diffusion processes and time reversal
- **Diffusion probabilistic model:** Derivation and experimental results

The Promise of Deep Unsupervised Learning

The Promise of Deep Unsupervised Learning

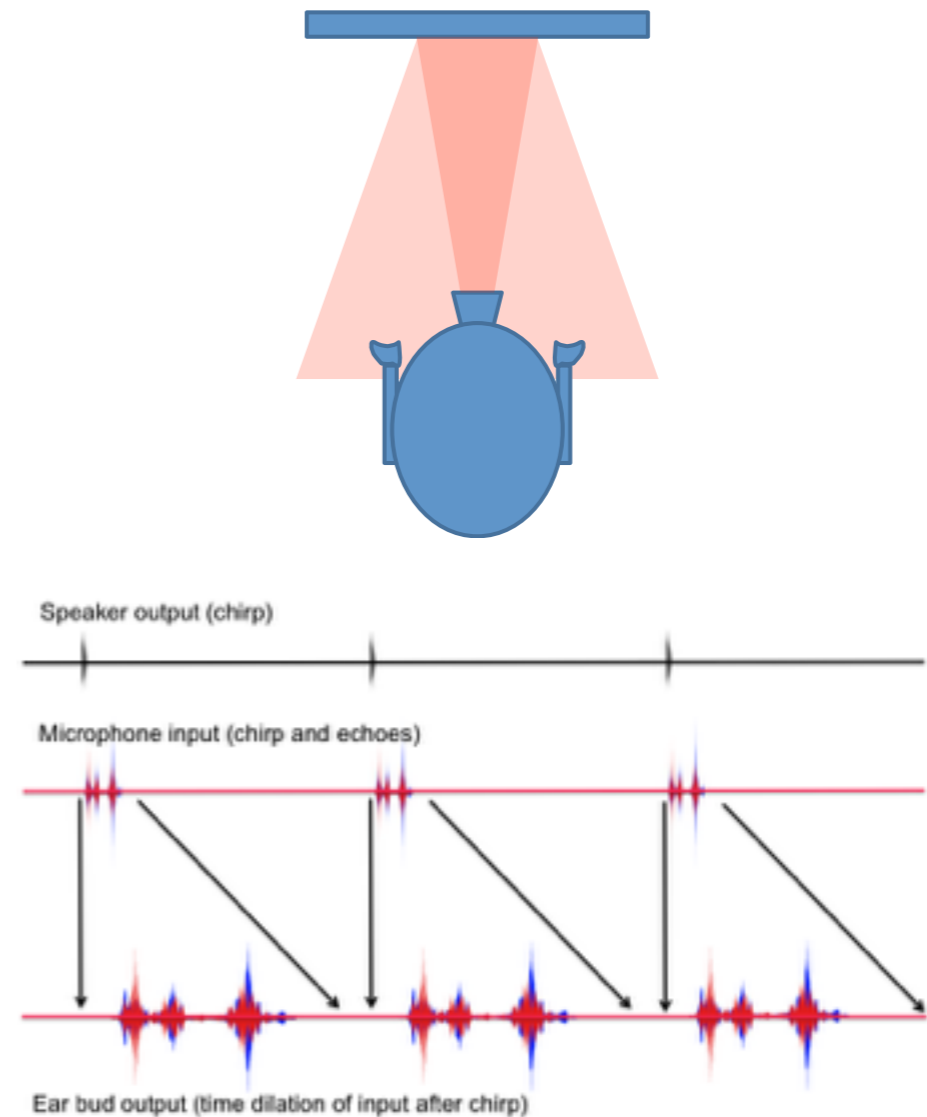
- Unknown features/labels

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

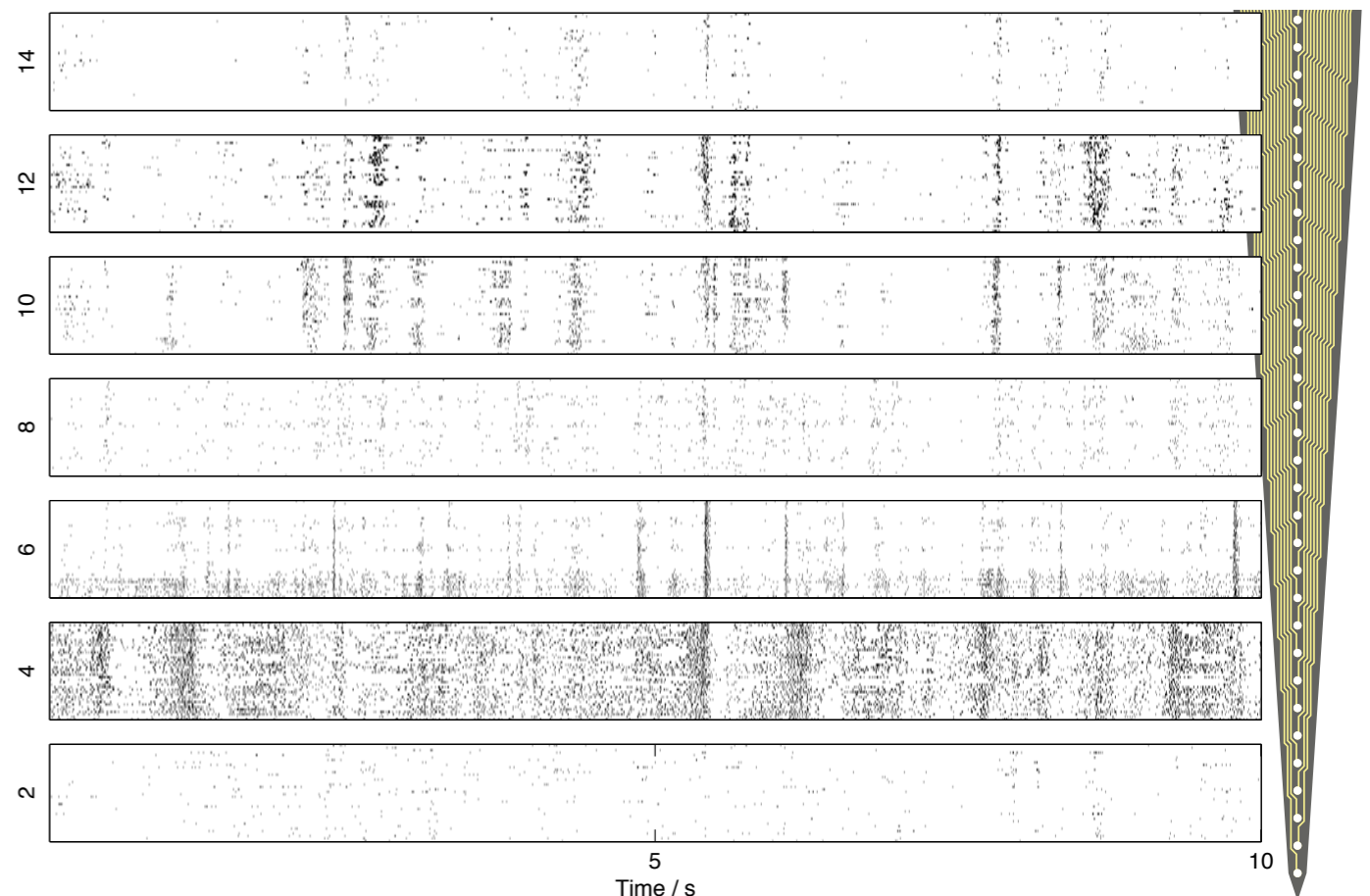
The Promise of Deep Unsupervised Learning

- 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

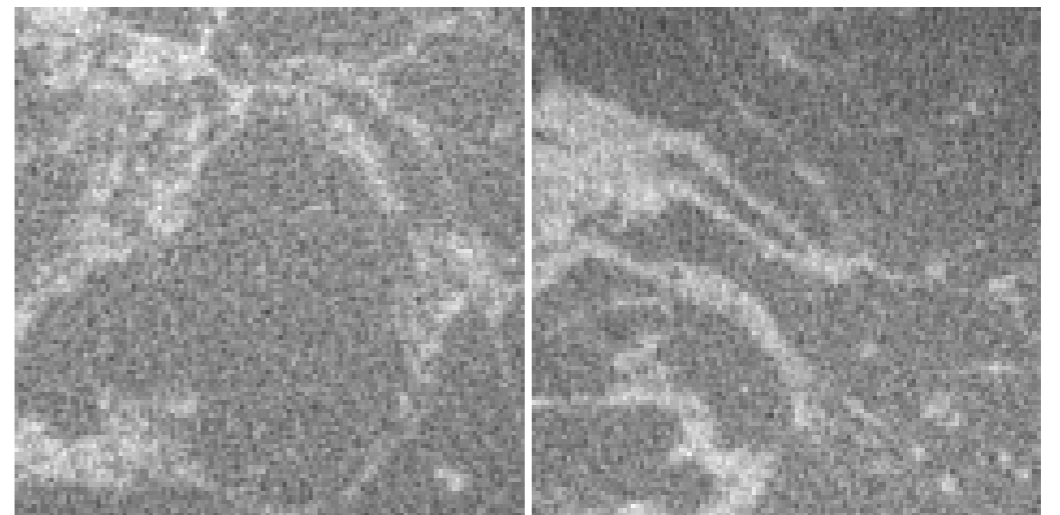
The Promise of Deep Unsupervised Learning

- 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

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

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



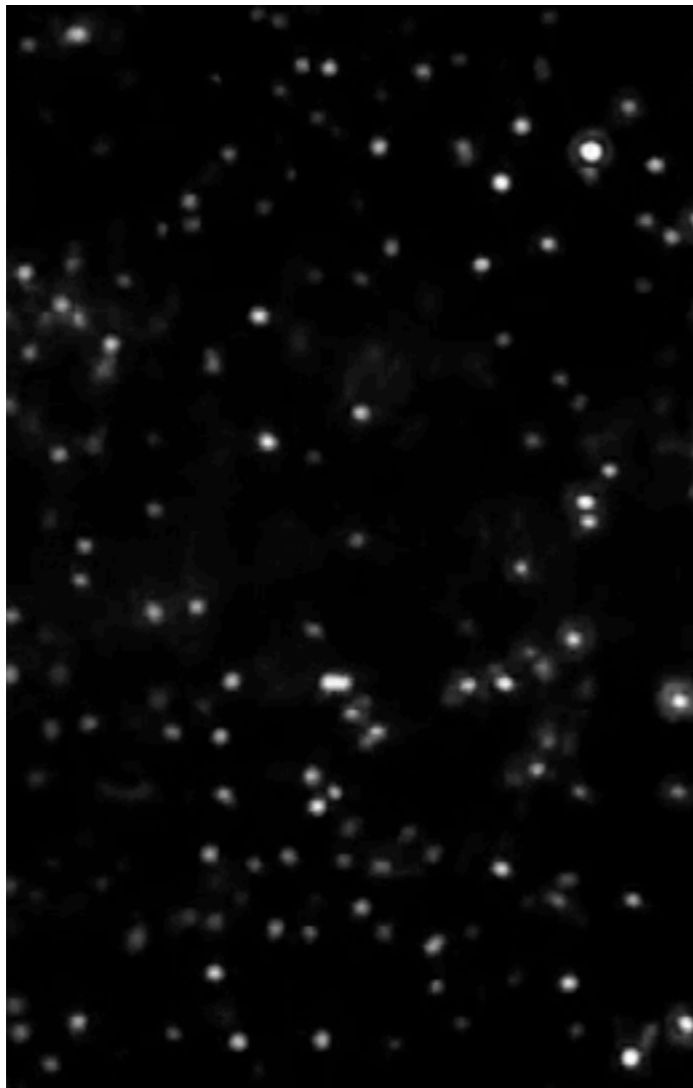
Uniform distribution

Core Idea: Recover Structure by Reversing Time



Core Idea: Recover Structure by Reversing Time

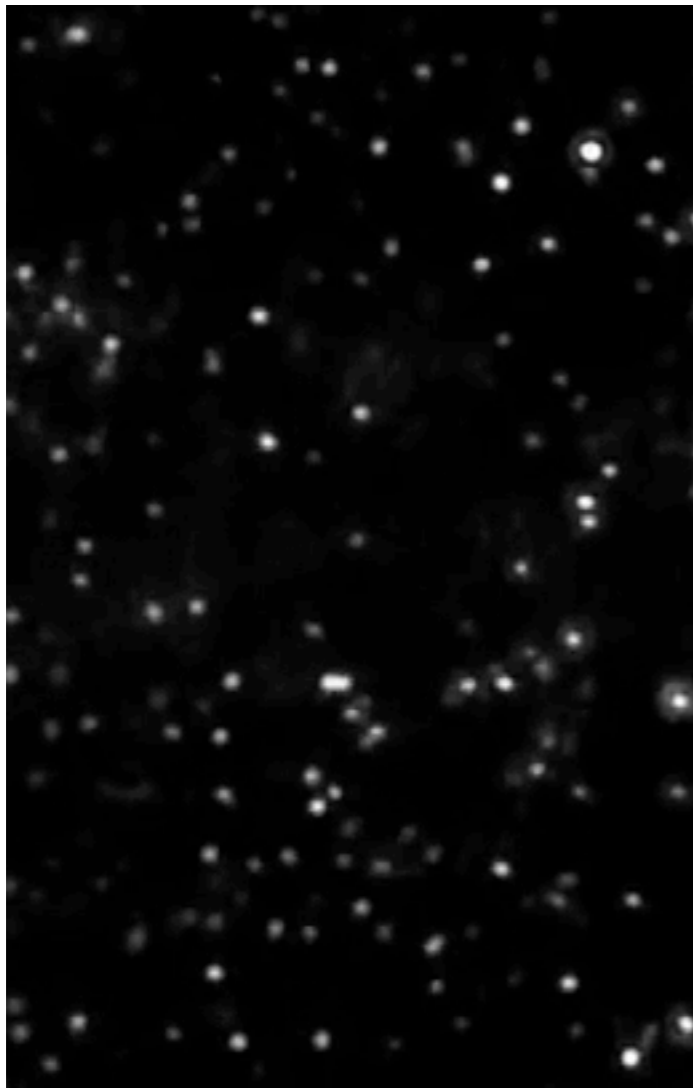
Observation 2: Microscopic Diffusion is Time Reversible



© Rutger Saly

- Microscopic view
- Brownian motion

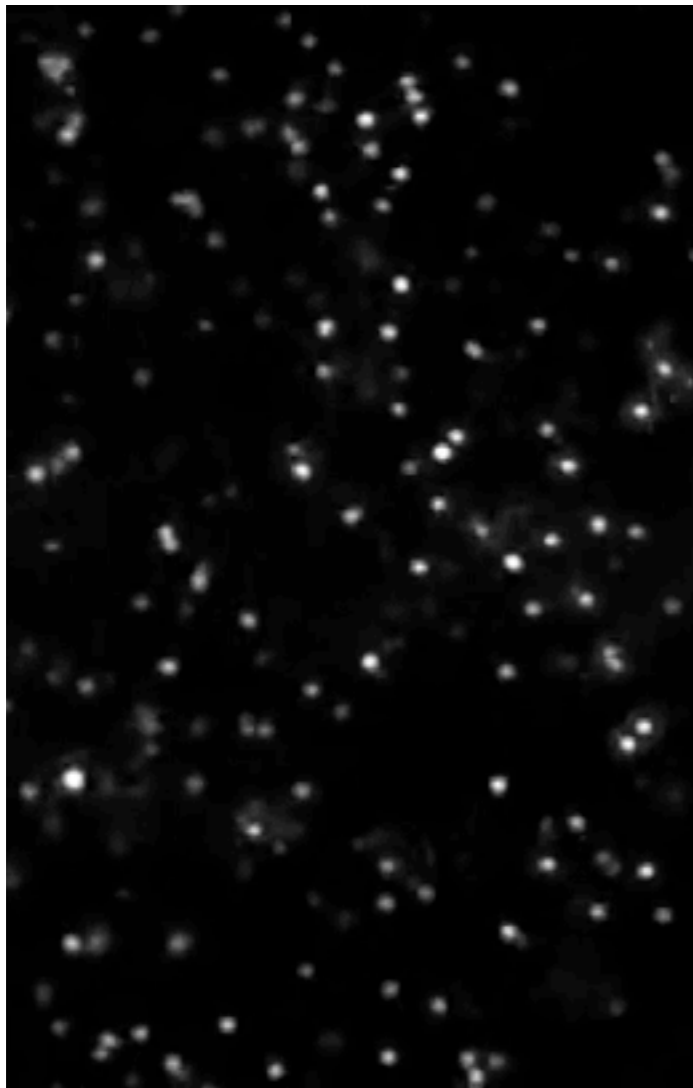
Observation 2: Microscopic Diffusion is Time Reversible



© Rutger Saly

- Microscopic view
- Brownian motion

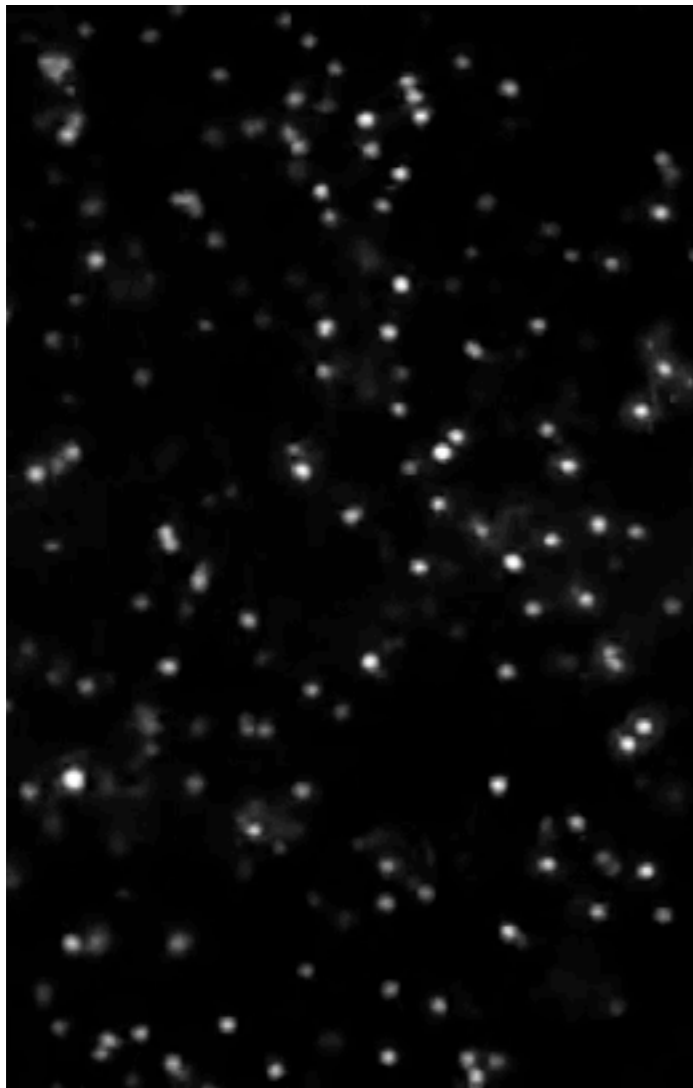
Observation 2: Microscopic Diffusion is Time Reversible



© Rutger Saly

- Microscopic view
- Brownian motion

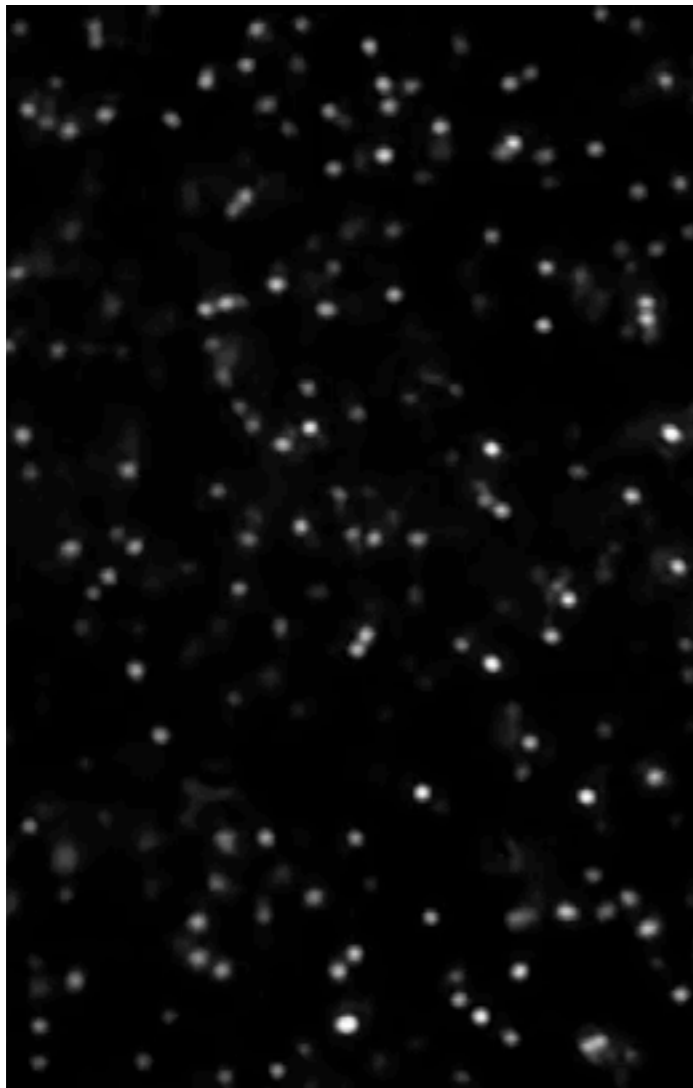
Observation 2: Microscopic Diffusion is Time Reversible



© Rutger Saly

- Microscopic view
- Brownian motion

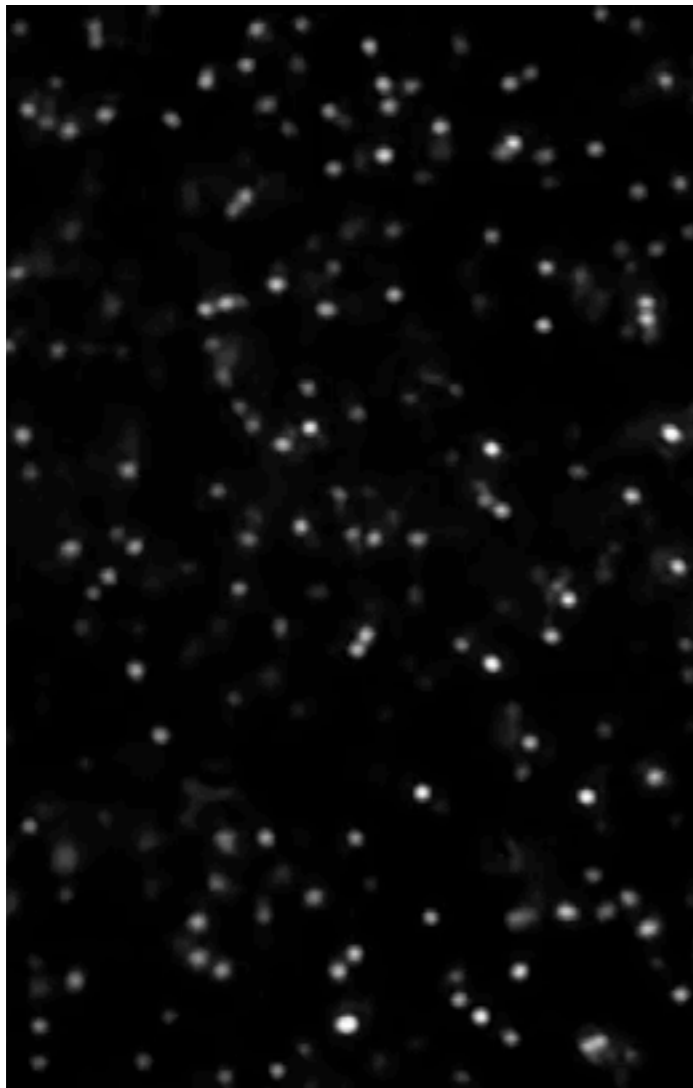
Observation 2: Microscopic Diffusion is Time Reversible



© Rutger Saly

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

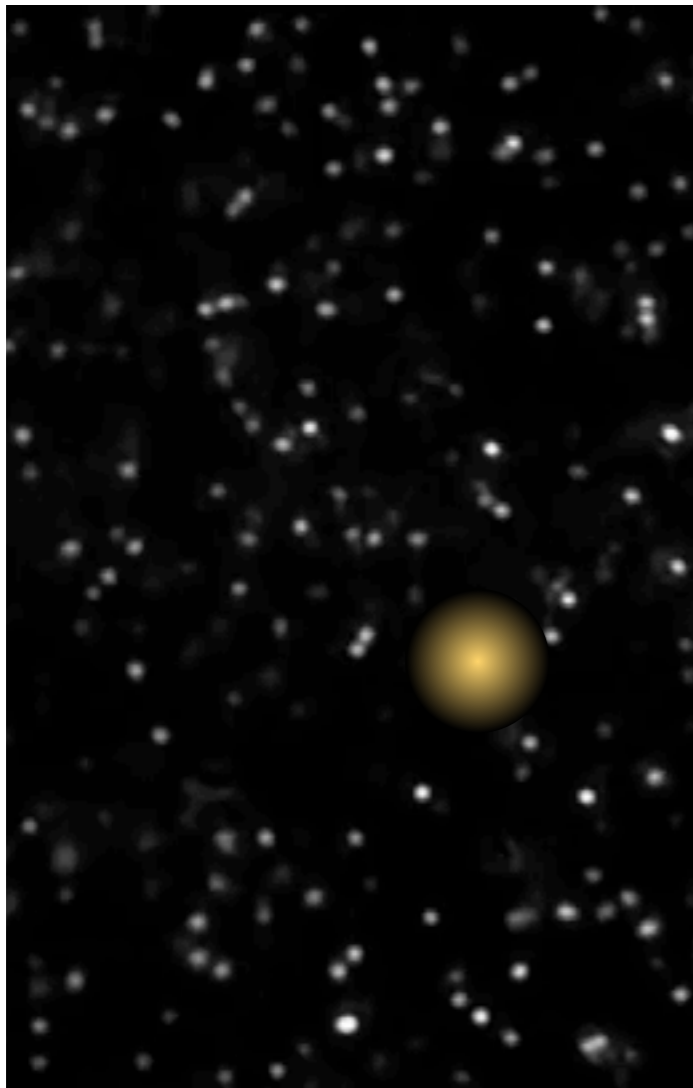
Observation 2: Microscopic Diffusion is Time Reversible



© Rutger Saly

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

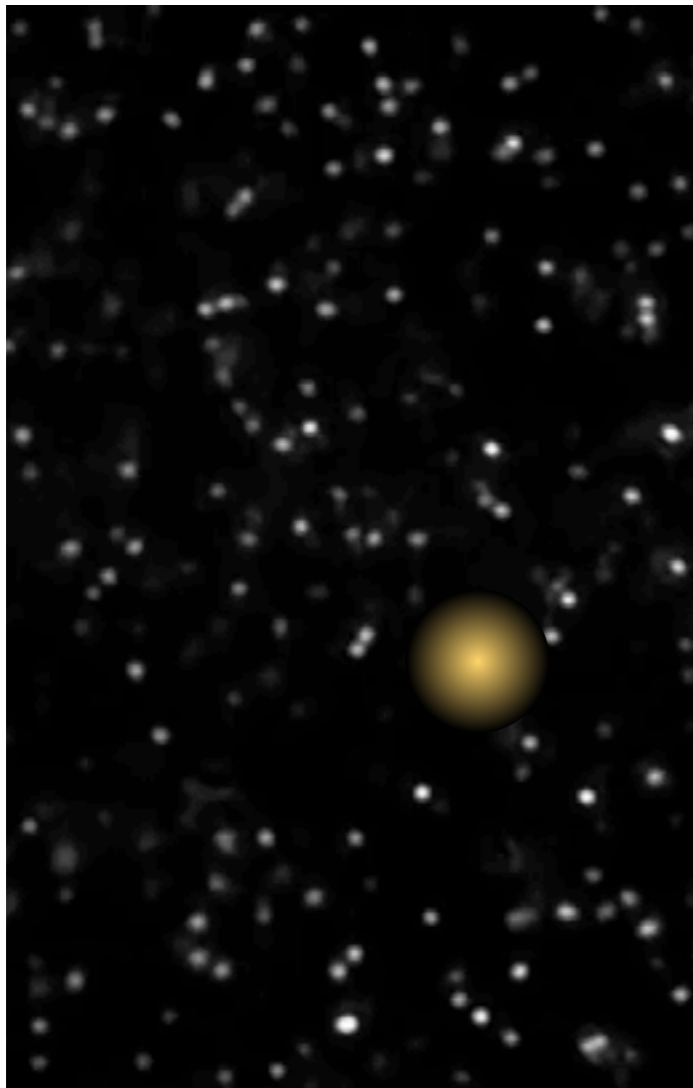
Observation 2: Microscopic Diffusion is Time Reversible



© 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

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

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

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



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

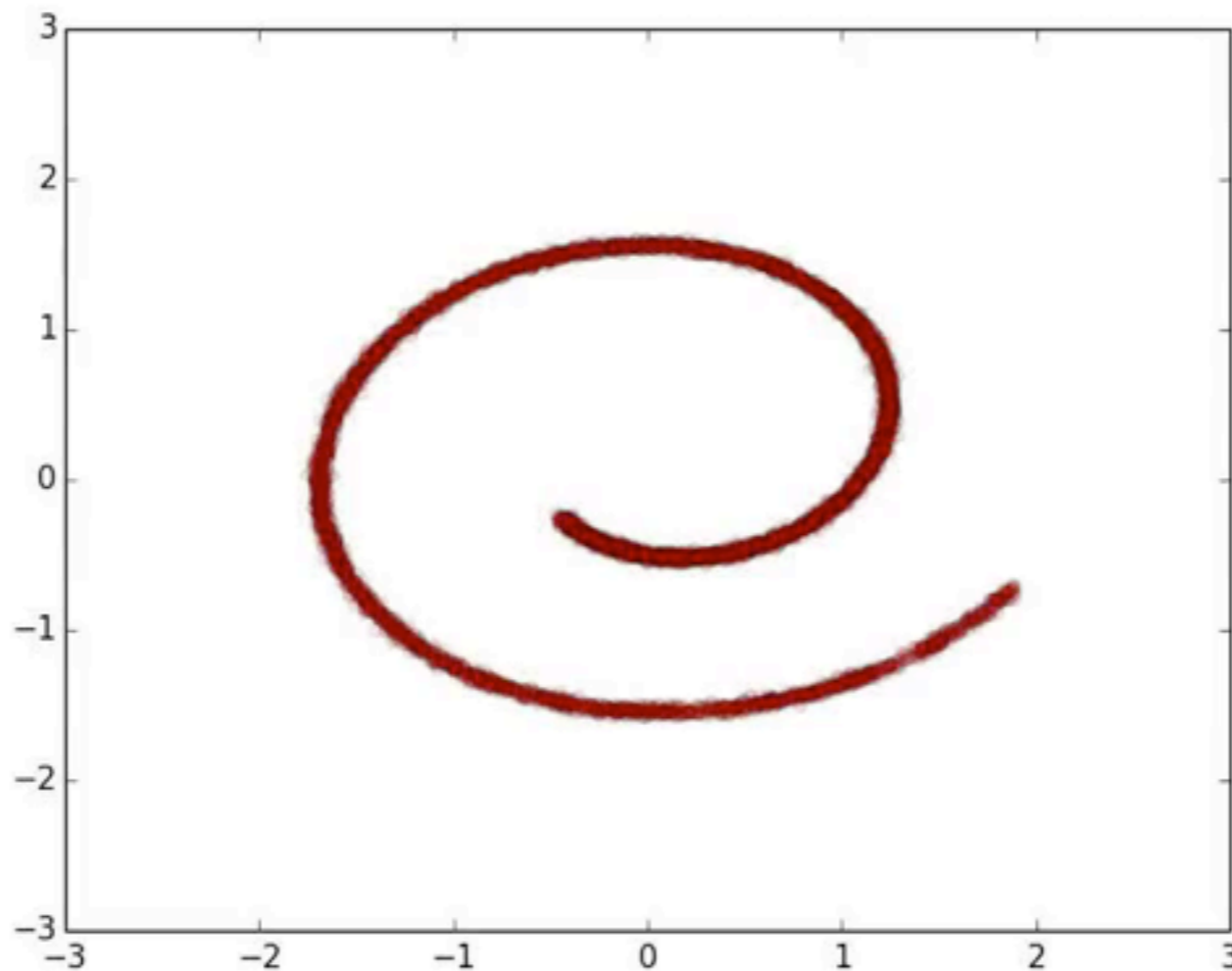
$$q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)}) = \mathcal{N}(\mathbf{x}^{(t)}; \underbrace{\mathbf{x}^{(t-1)} \sqrt{1 - \beta_t}}_{\text{Decay towards origin}}, \underbrace{\mathbf{I} \beta_t}_{\text{Add small noise}})$$

Decay towards origin

Add small noise

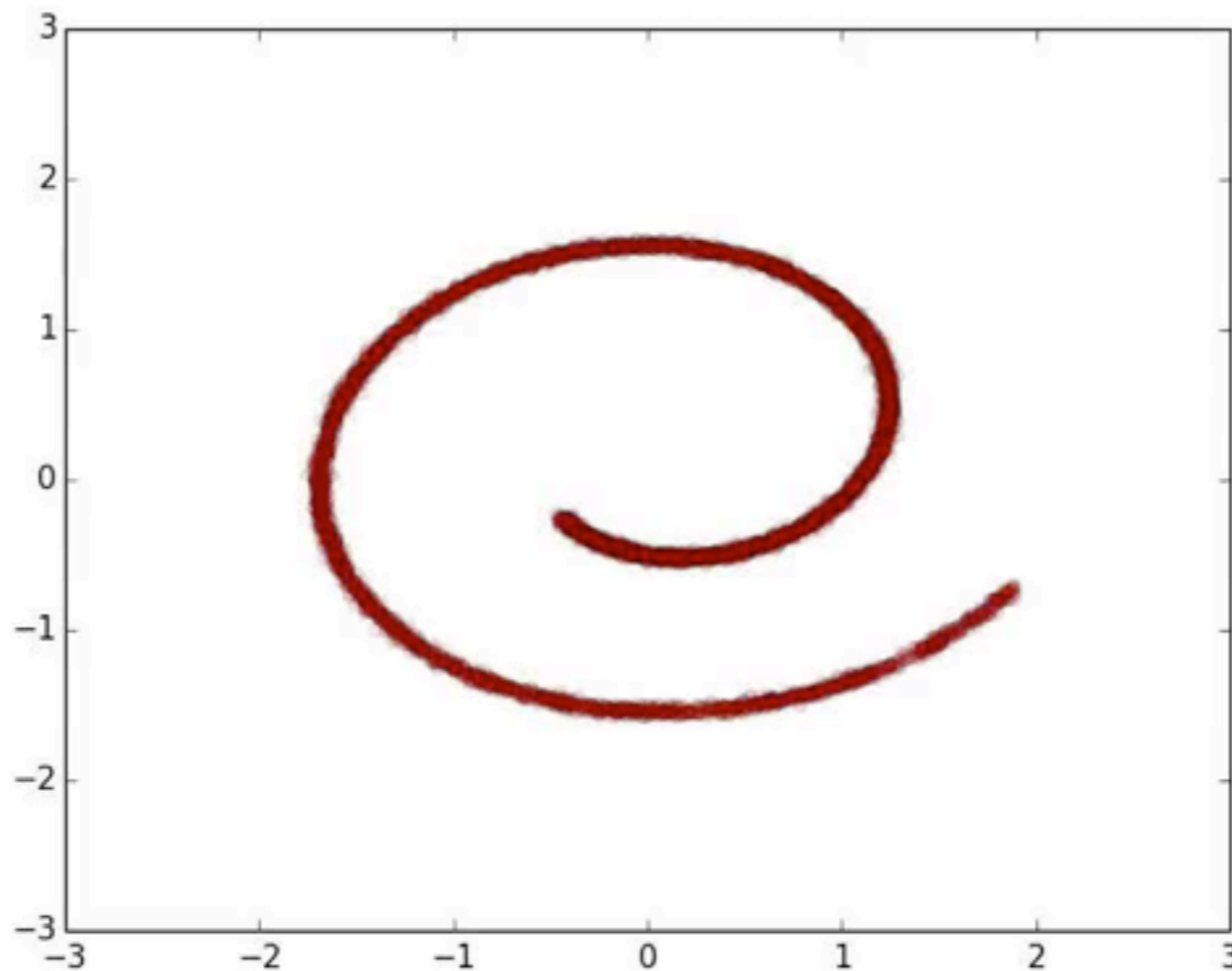
Forward Diffusion Process on Swiss Roll

- Start at data
- Run Gaussian diffusion until samples become Gaussian blob



Forward Diffusion Process on Swiss Roll

- Start at data
- Run Gaussian diffusion until samples become Gaussian blob



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



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

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

Learned drift and covariance functions

Recover Structure in Data using Reversal of Diffusion Process

Data
distribution

Reverse
diffusion

Noise
distribution

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



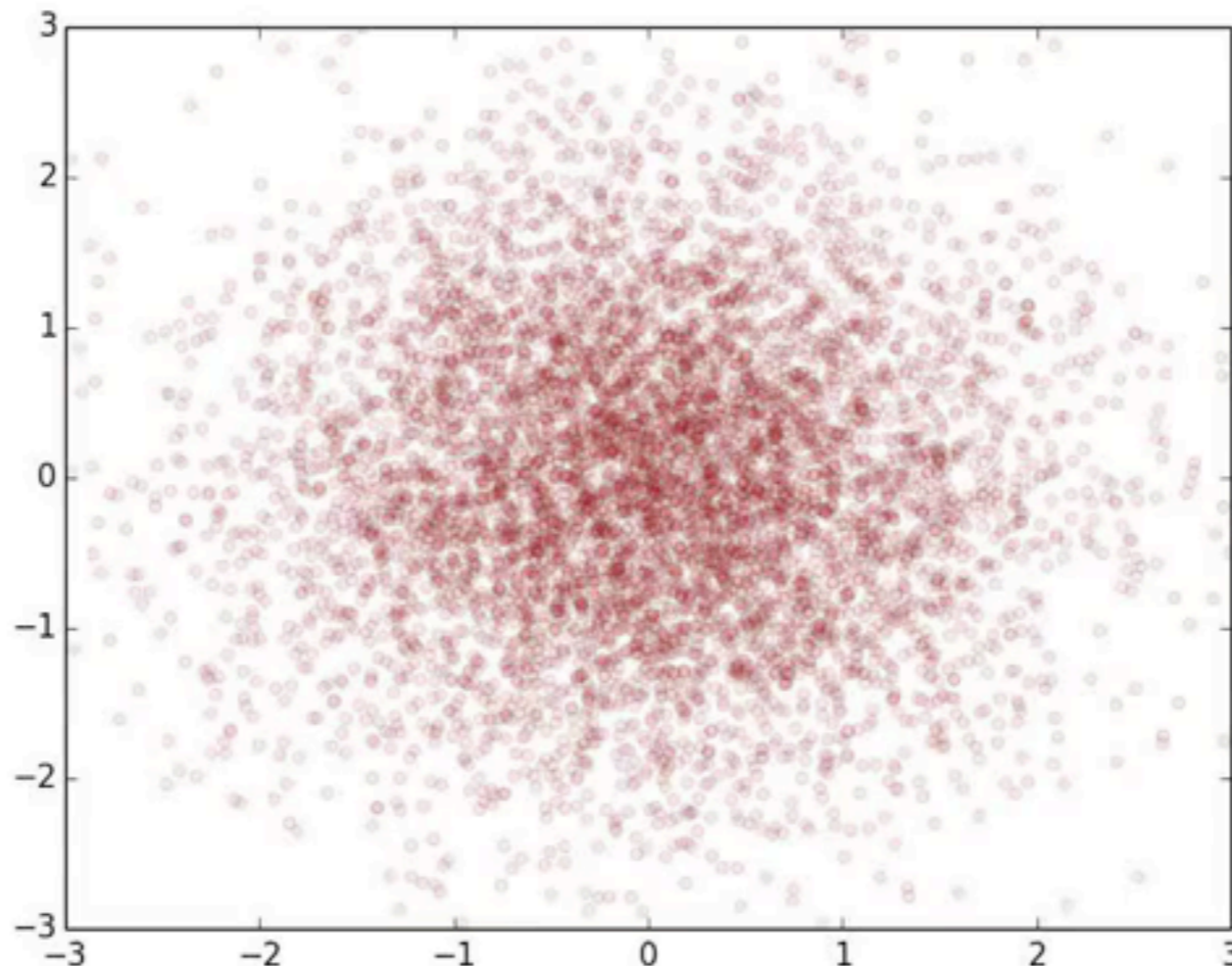
$$p(\mathbf{x}^{(T)}) = \mathcal{N}(\mathbf{x}^{(T)}; 0, \mathbf{I})$$

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

Learned drift and covariance functions

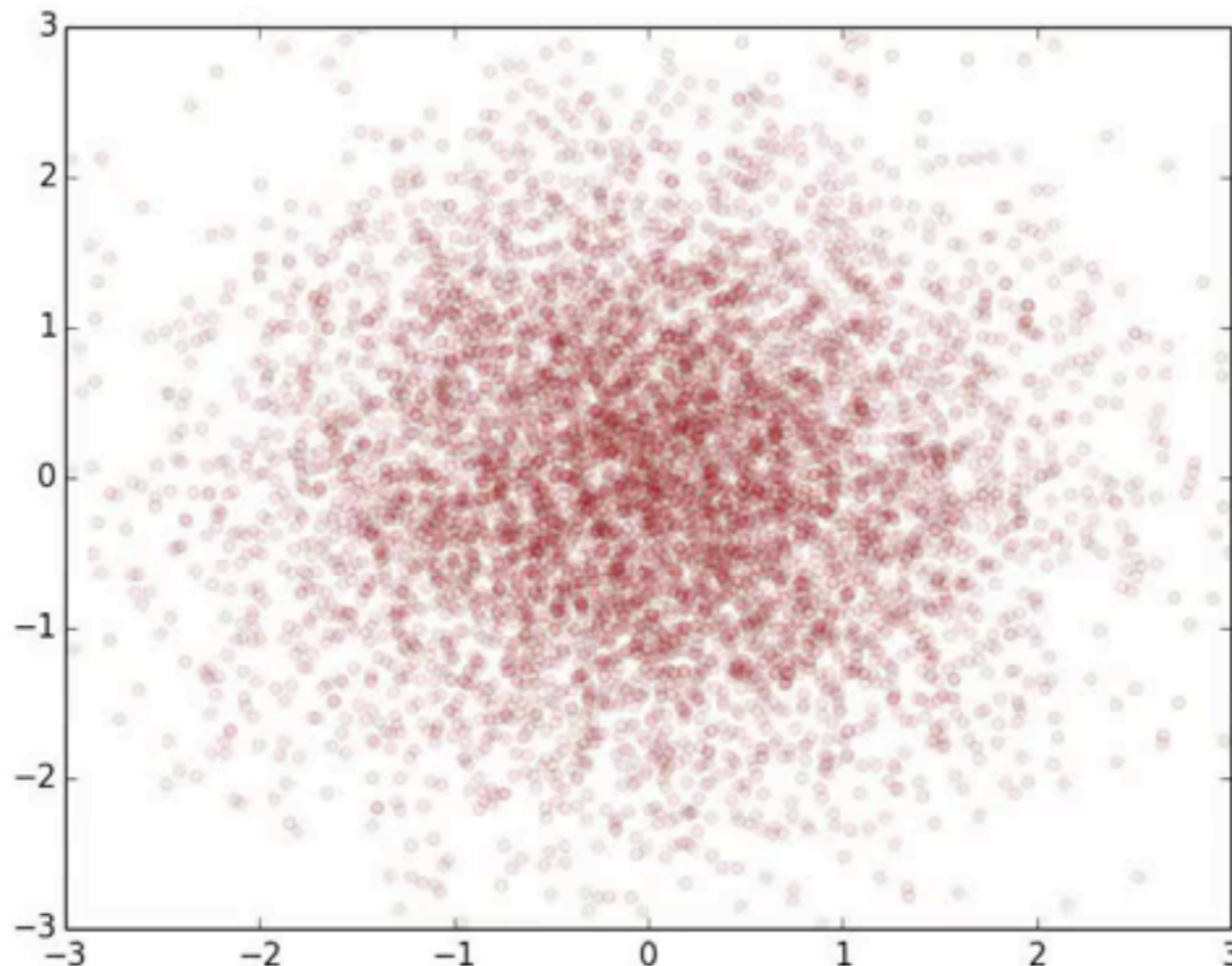
Learned Reverse Diffusion Process on Swiss Roll

- Start at Gaussian blob
- Run Gaussian diffusion until samples become data distribution

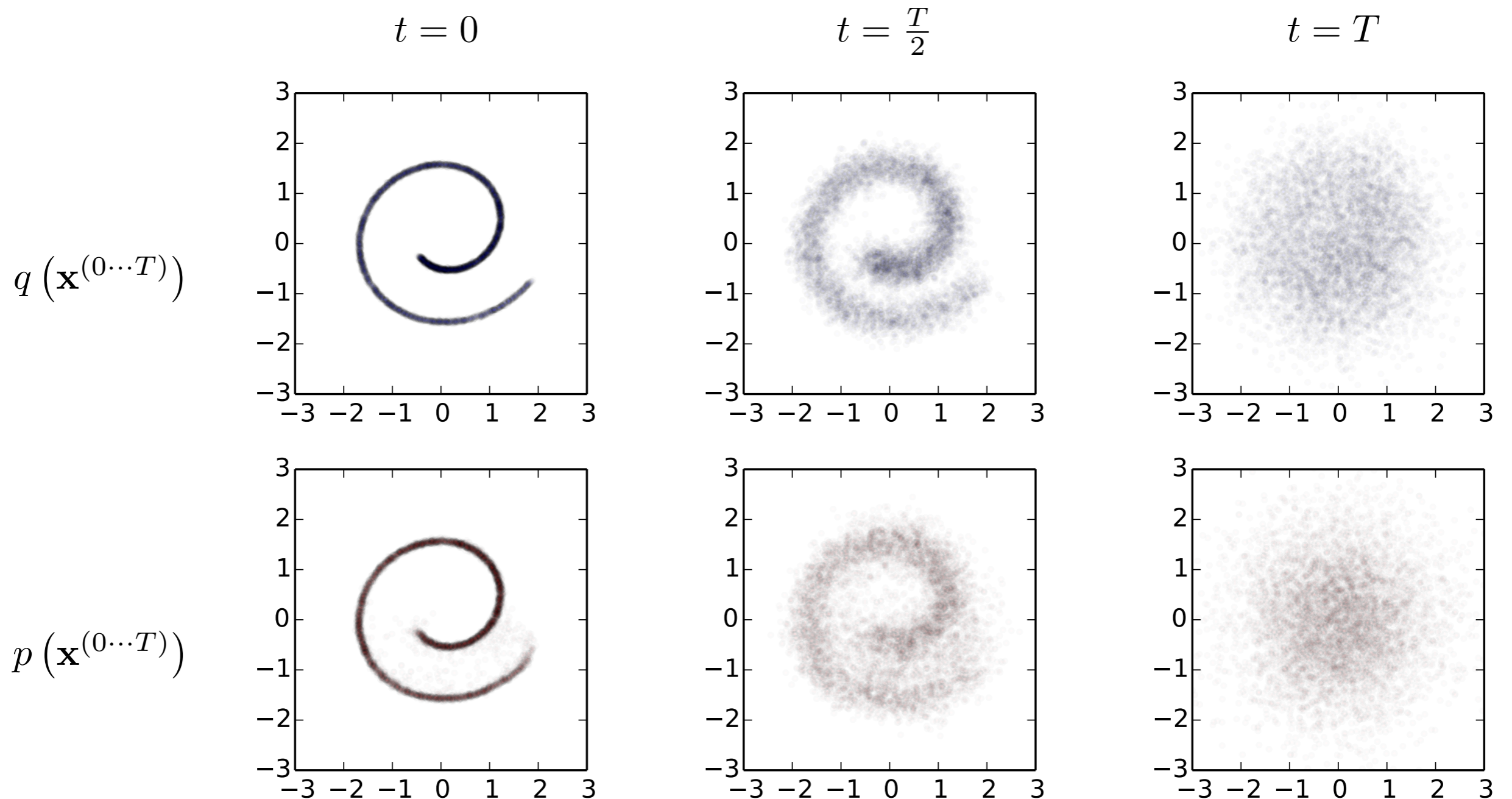


Learned Reverse Diffusion Process on Swiss Roll

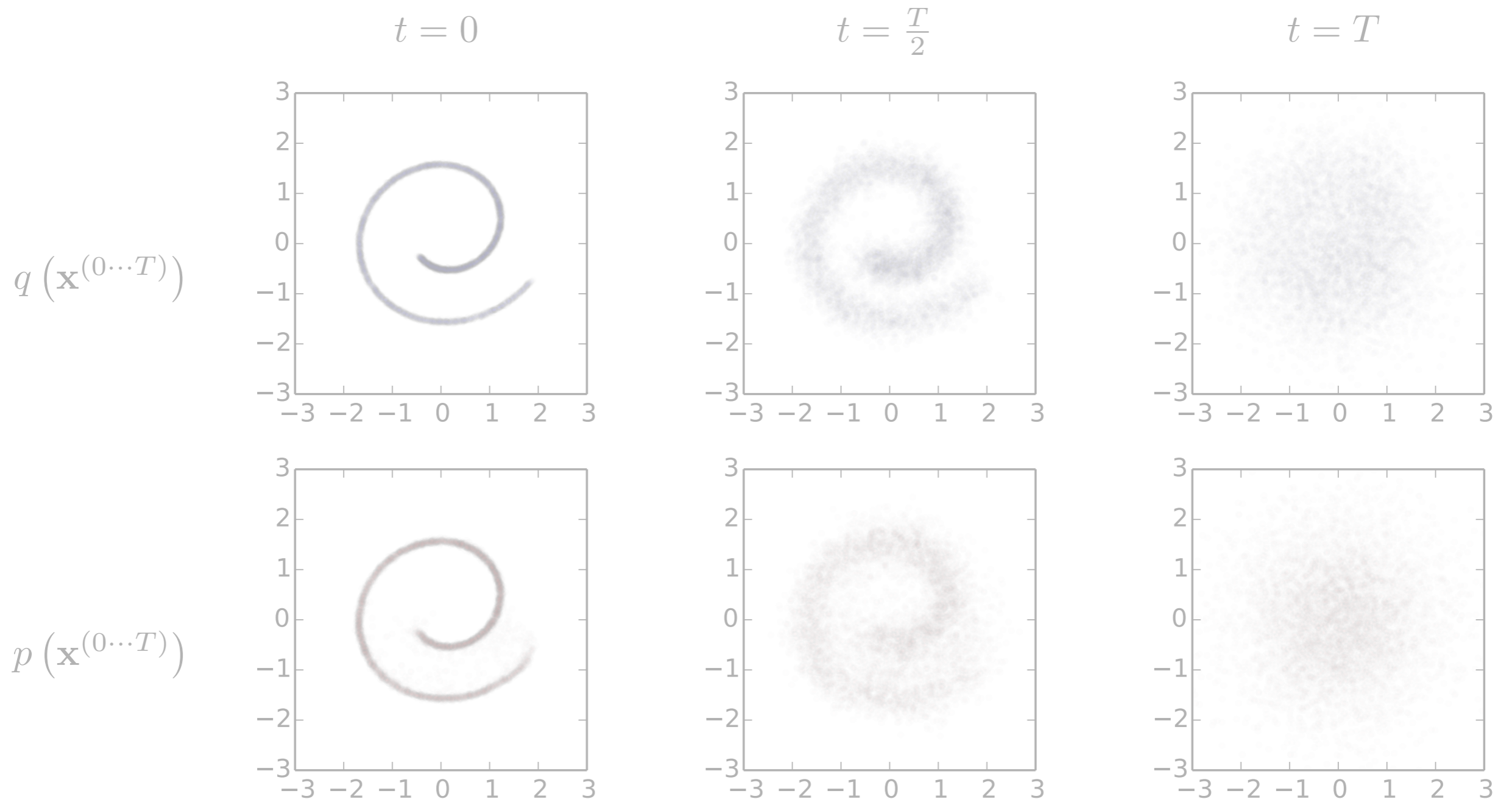
- Start at Gaussian blob
- Run Gaussian diffusion until samples become data distribution



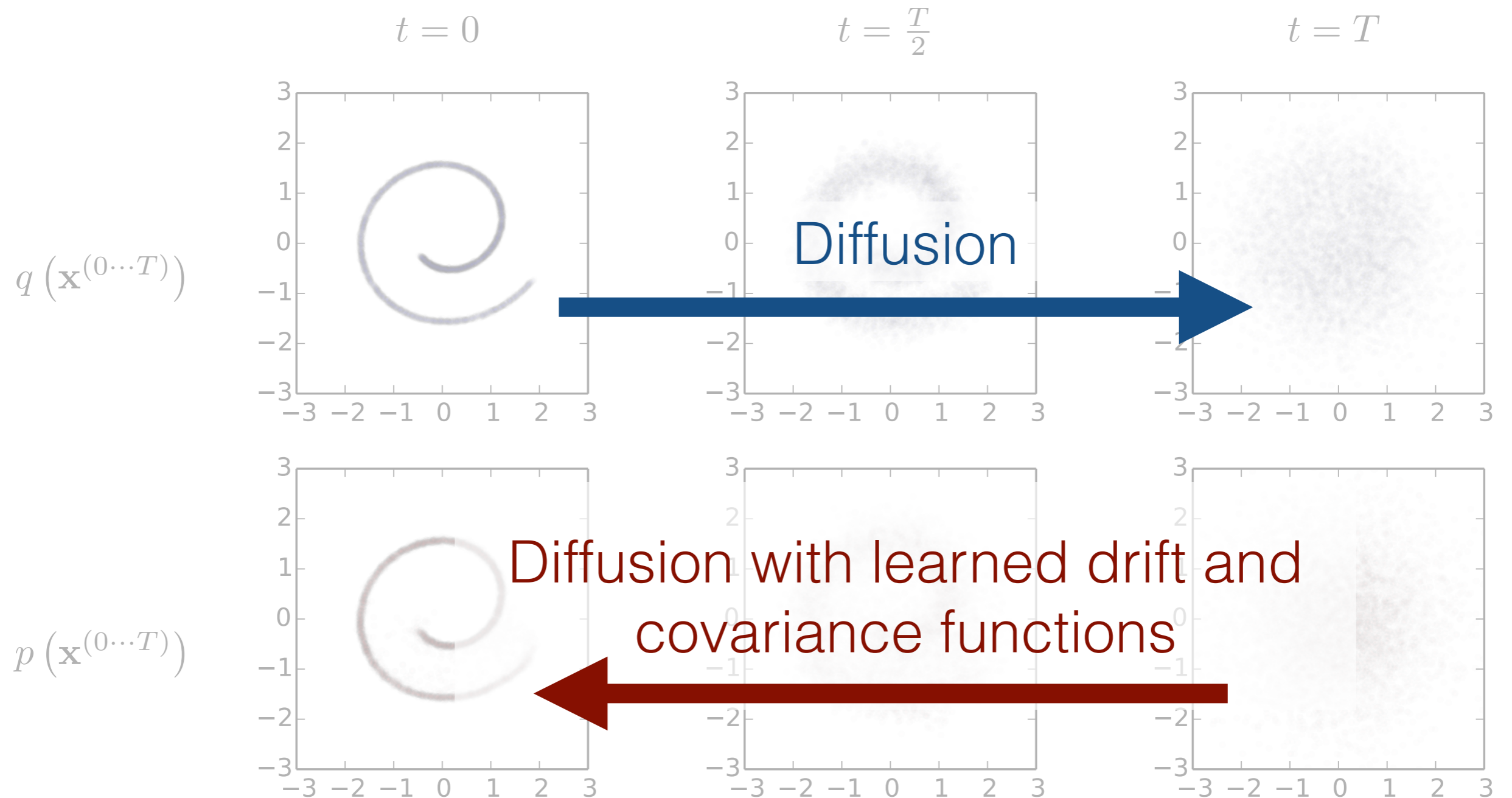
Summary of Forward and Reverse Diffusion on Swiss Roll



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

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




Gaussian

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

Gaussian

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

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


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

Training

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

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

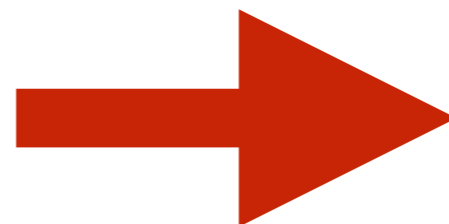


$$p\left(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}\right) = \mathcal{N}\left(\mathbf{x}^{(t-1)}; f_{\mu}\left(\mathbf{x}^{(t)}, t\right), f_{\Sigma}\left(\mathbf{x}^{(t)}, t\right)\right)$$

Training

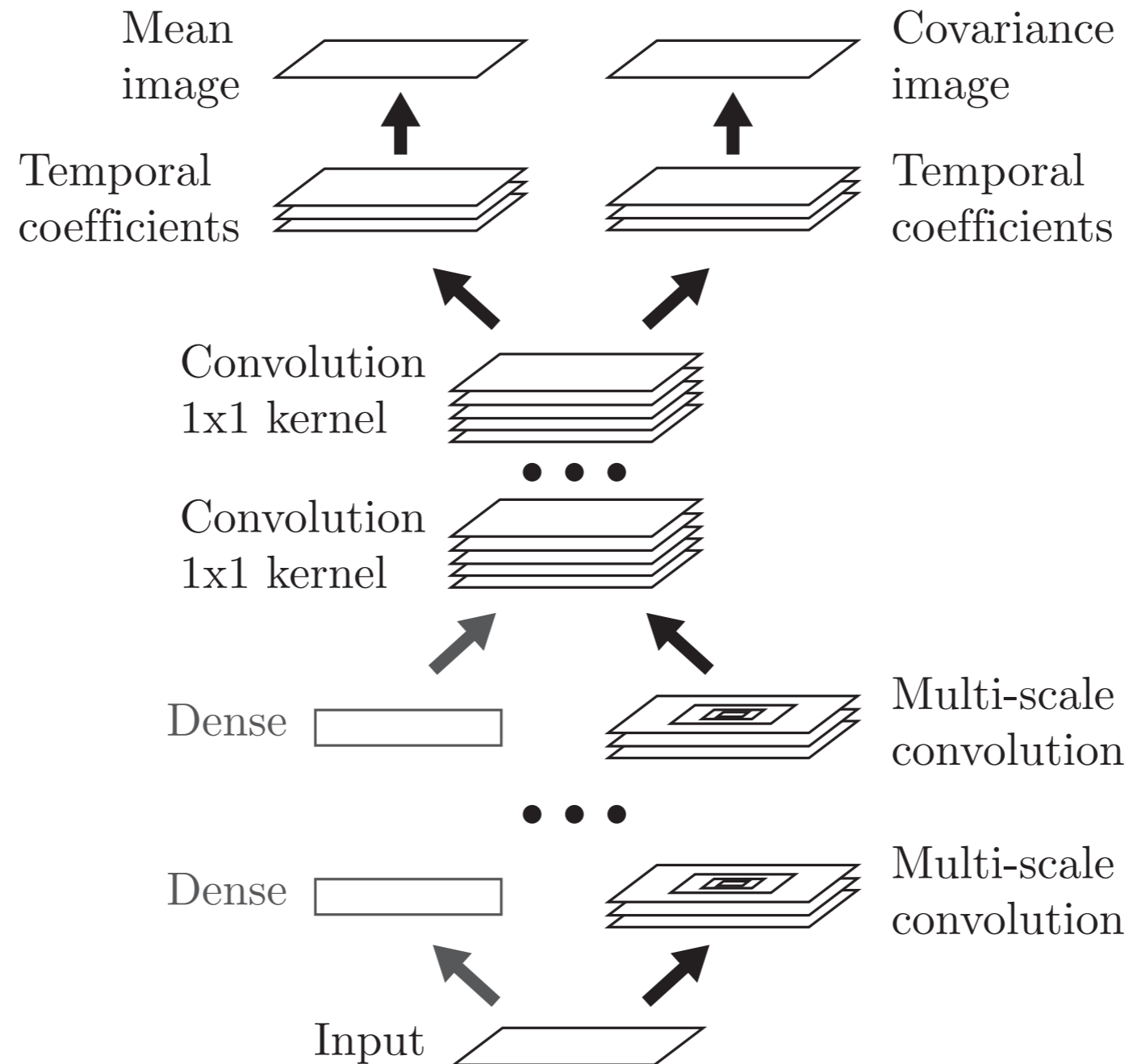
$$\underset{f_{\mu}(\mathbf{x}^{(t)}, t), f_{\Sigma}(\mathbf{x}^{(t)}, t)}{\text{argmin}} \mathbb{E} \left[D_{KL} \left(q\left(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}\right) \parallel p\left(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}\right) \right) \right]$$

Unsupervised
learning

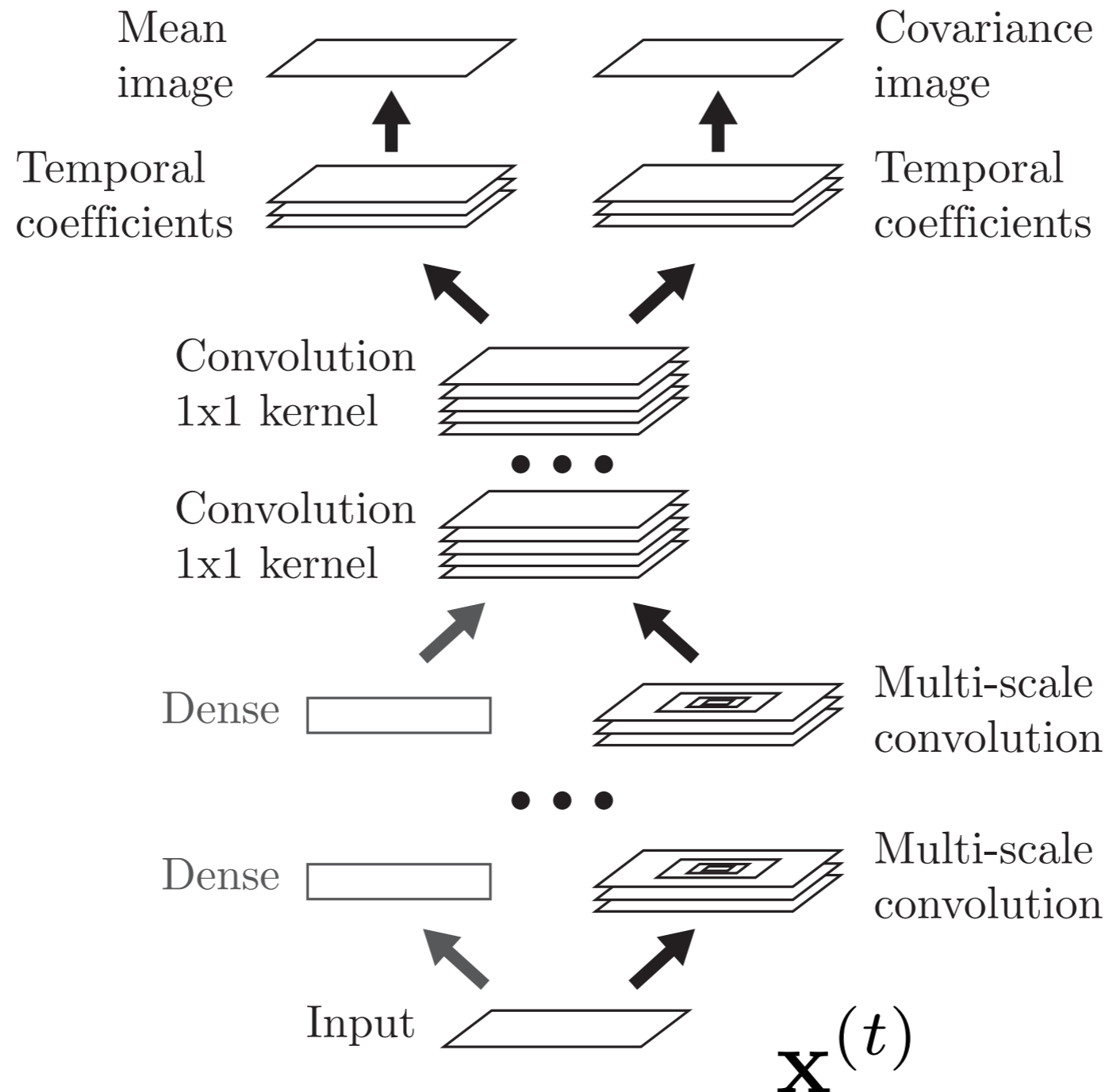


Regression

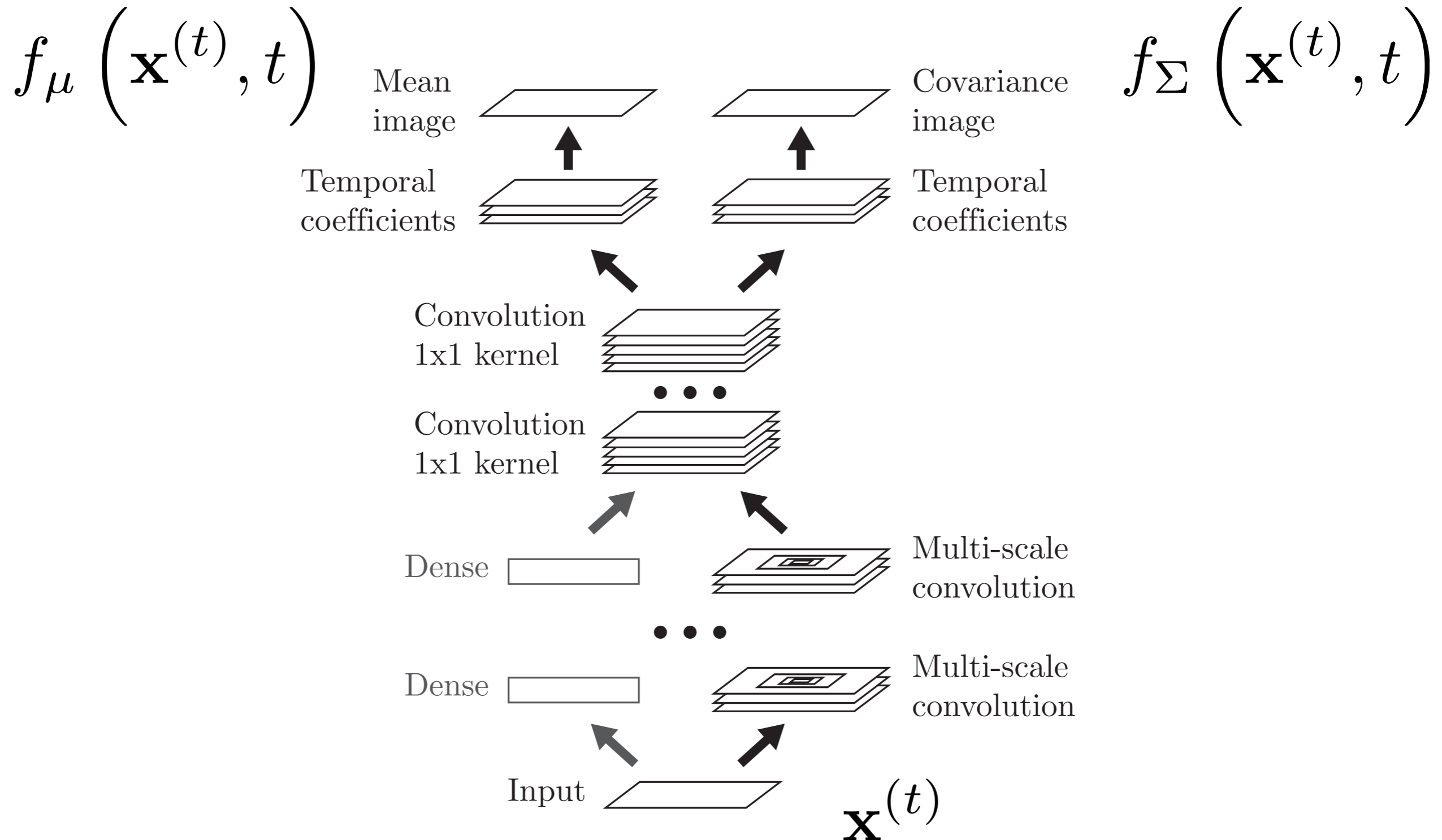
Use Deep Network as Function Approximator for Images



Use Deep Network as Function Approximator for Images



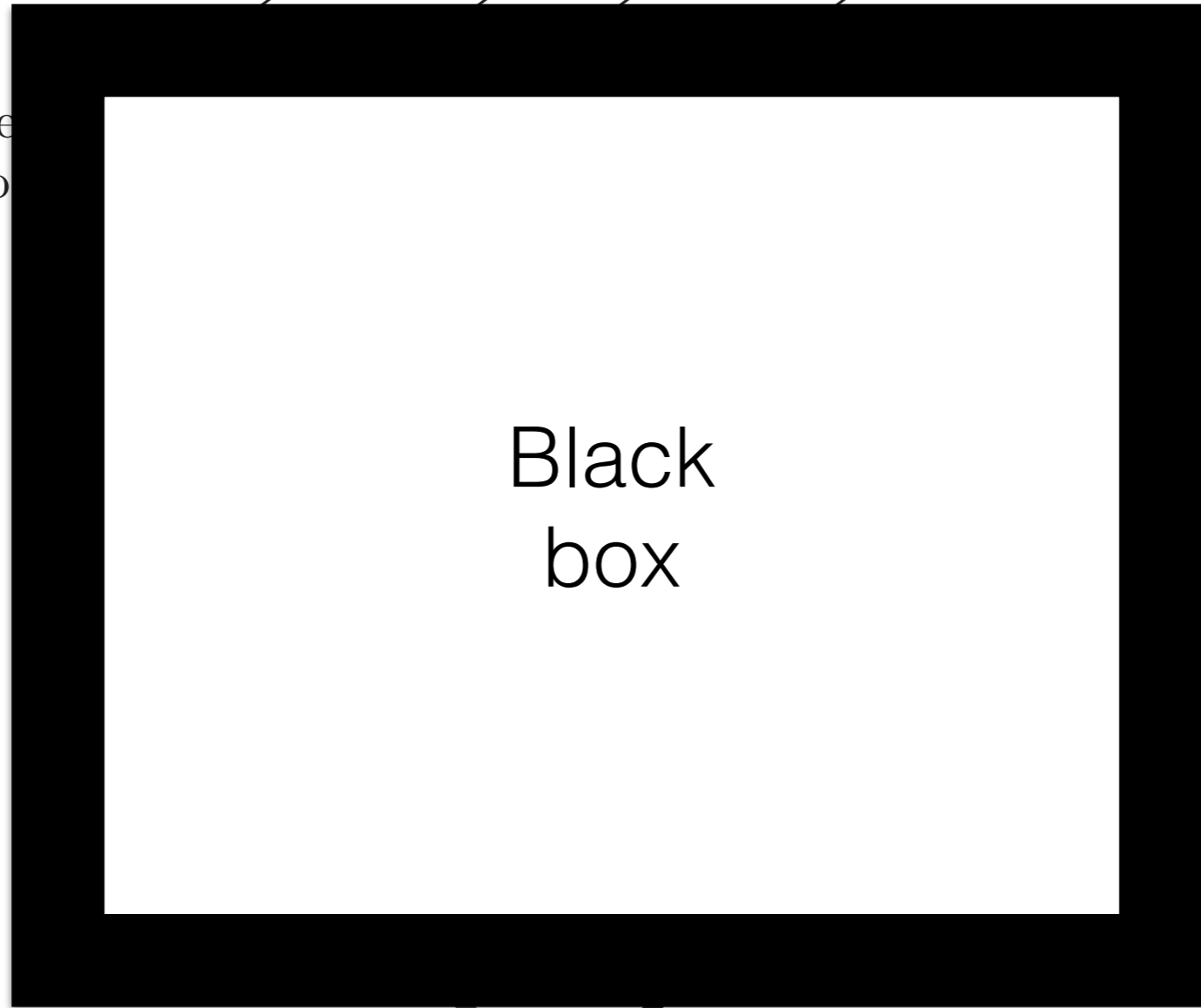
Use Deep Network as Function Approximator for Images



Use Deep Network as Function Approximator for Images

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

Te
co



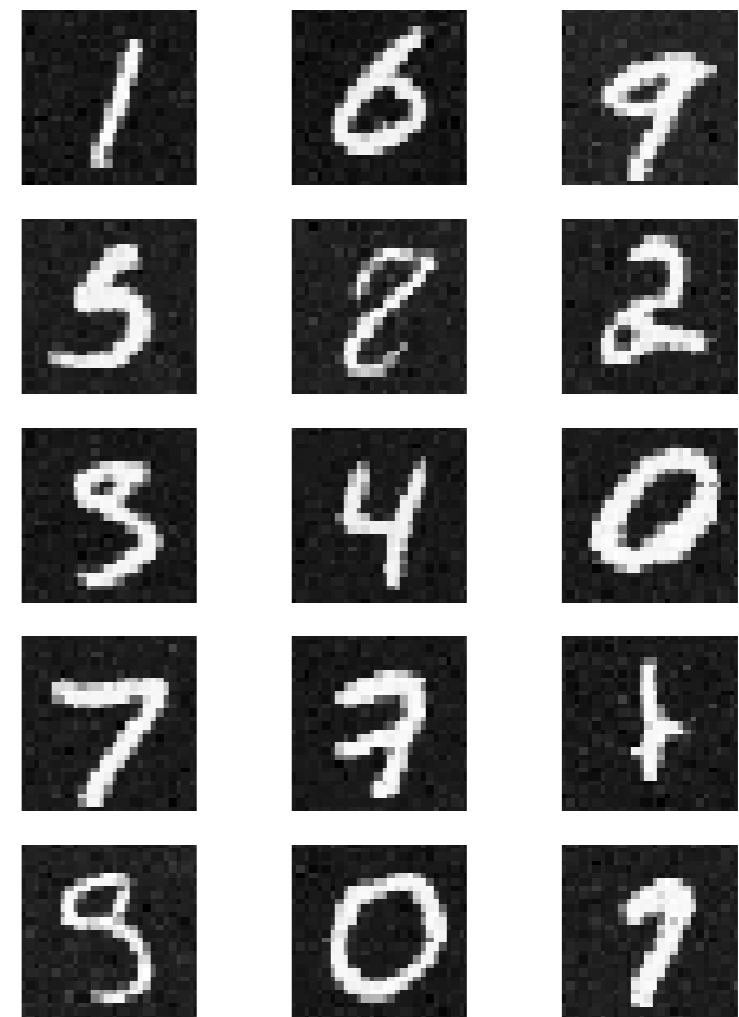
Black
box

Input

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

Diffusion Probabilistic Model Applied to MNIST

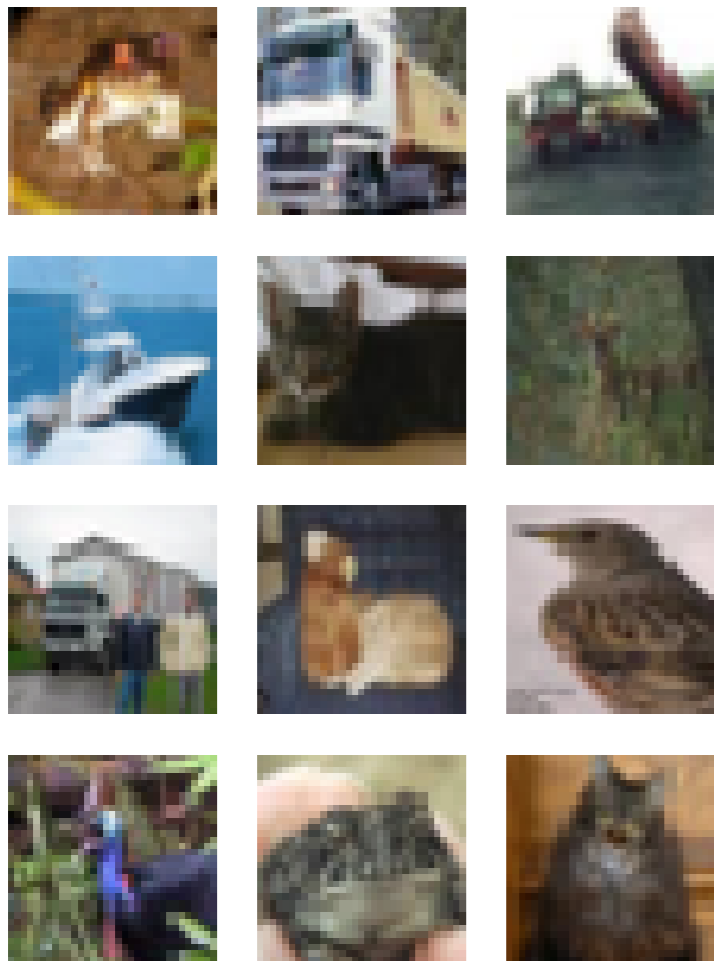
Model	Log likelihood estimate*
Stacked CAE	121 \pm 1.6 bits
DBN	138 \pm 2 bits
Deep GSN	214 \pm 1.1 bits
Diffusion	220 \pm 1.9 bits
Adversarial net	225 \pm 2 bits



Samples from
diffusion model

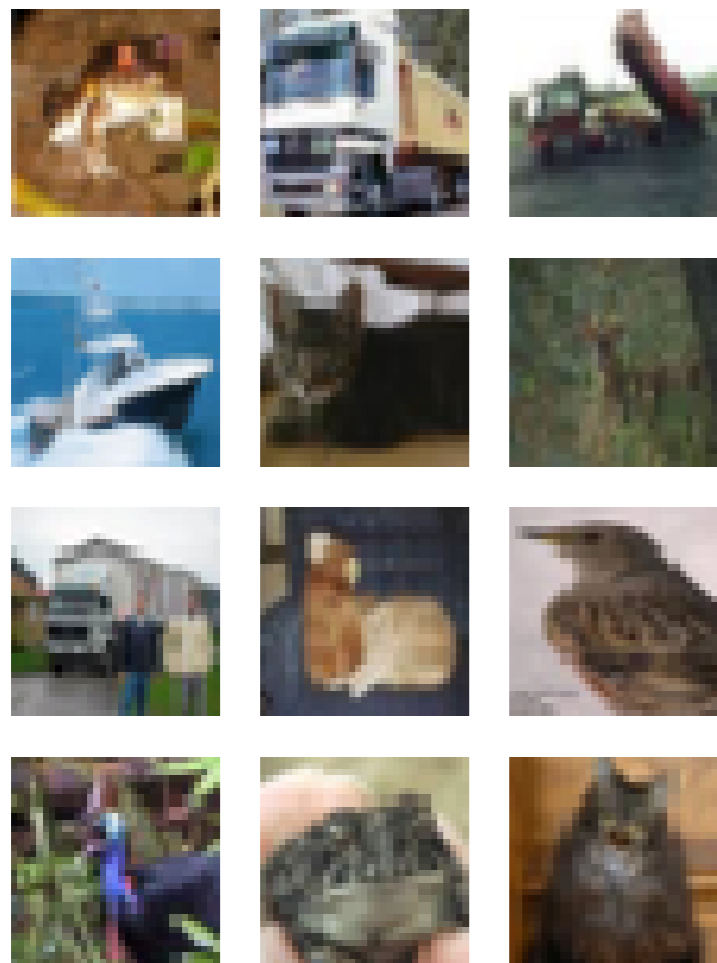
* via Parzen window code from [Goodfellow *et al*, 2014]

Diffusion Probabilistic Model Applied to CIFAR-10

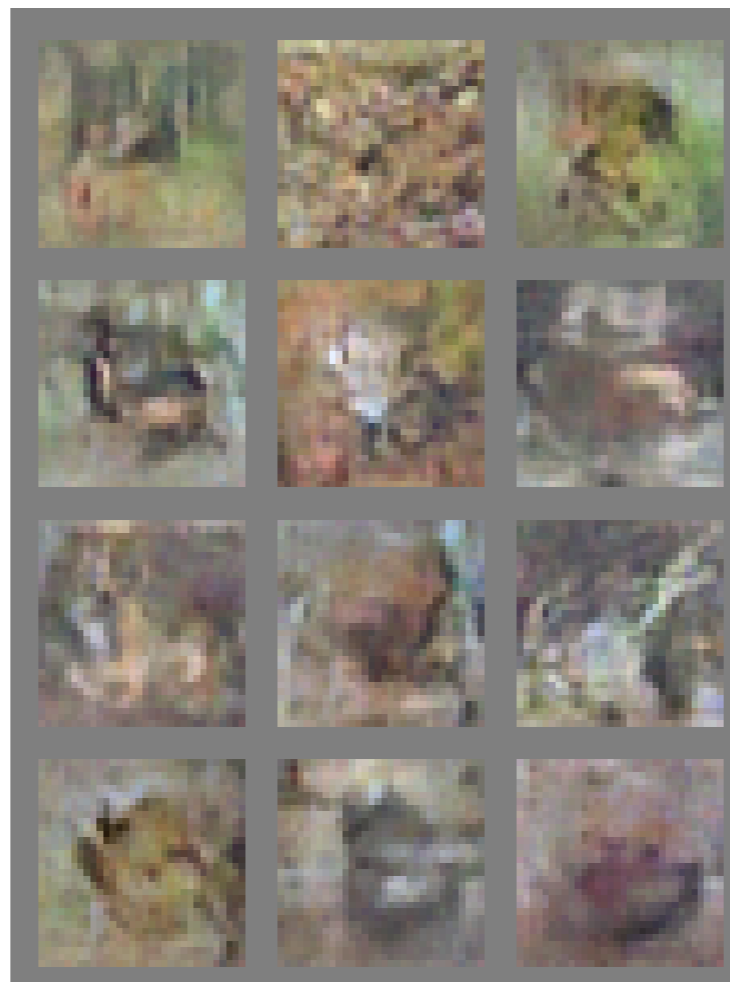


Training Data

Diffusion Probabilistic Model Applied to CIFAR-10

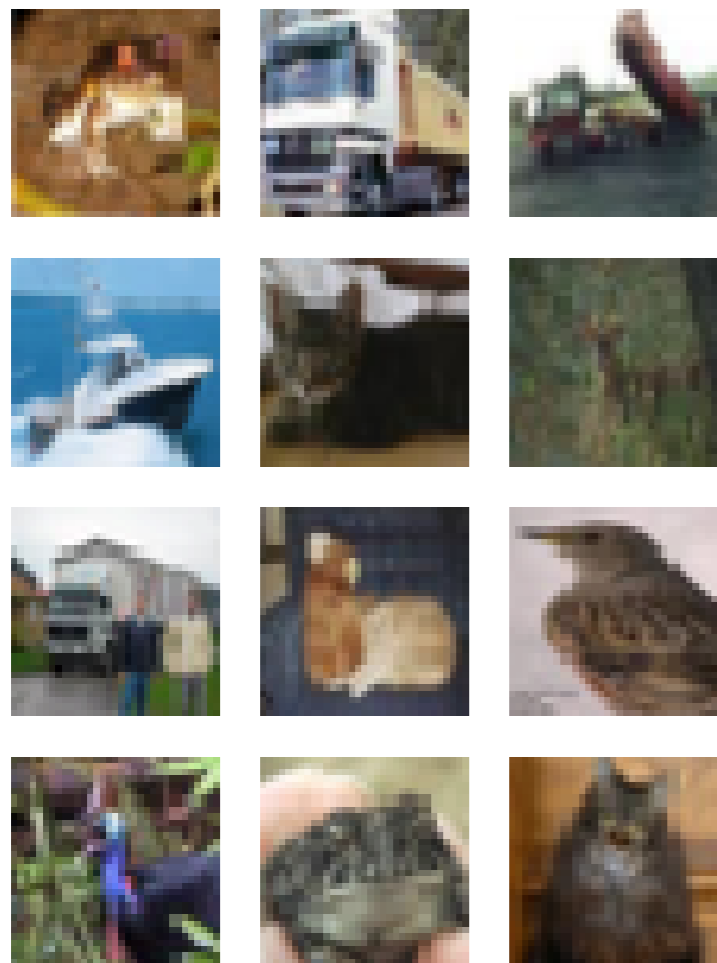


Training Data

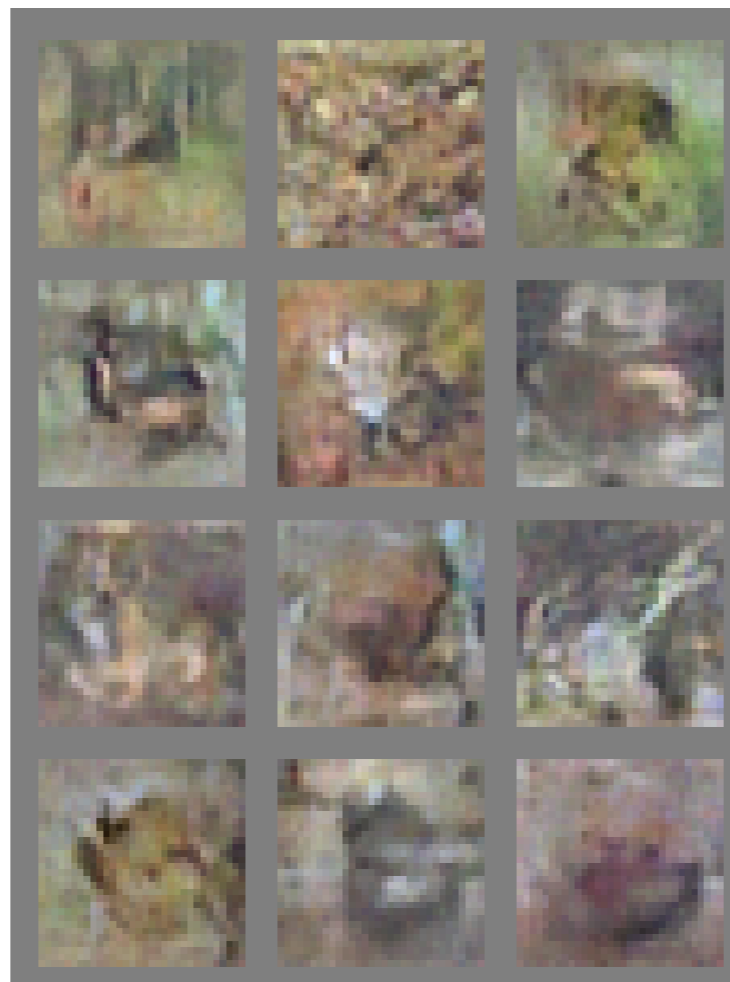


Samples from
Generative Adversarial
[Goodfellow *et al*, 2014]

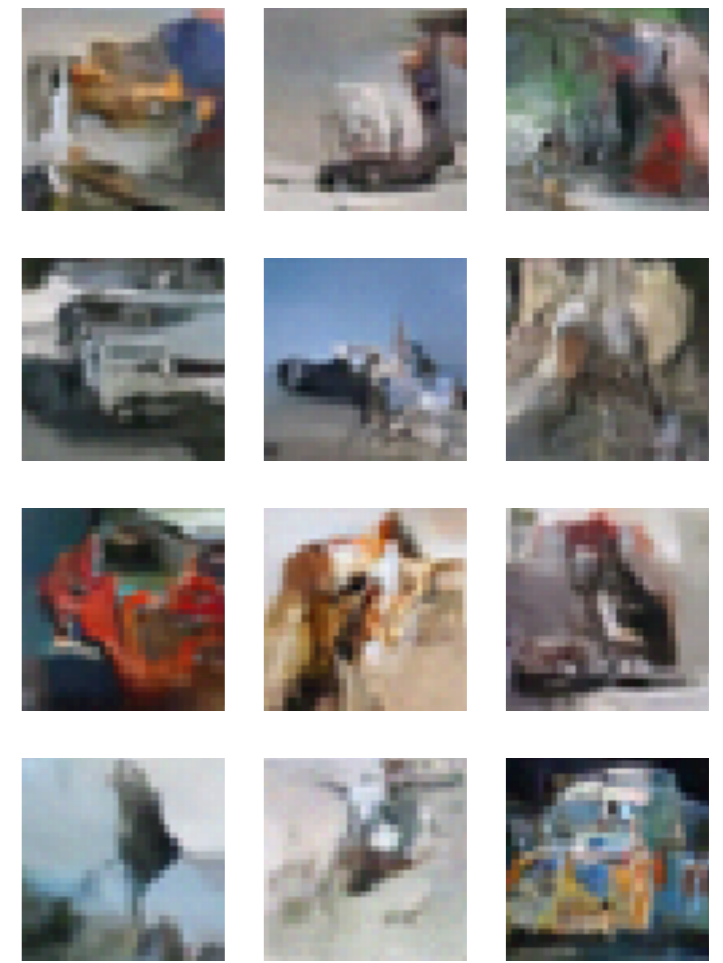
Diffusion Probabilistic Model Applied to CIFAR-10



Training Data

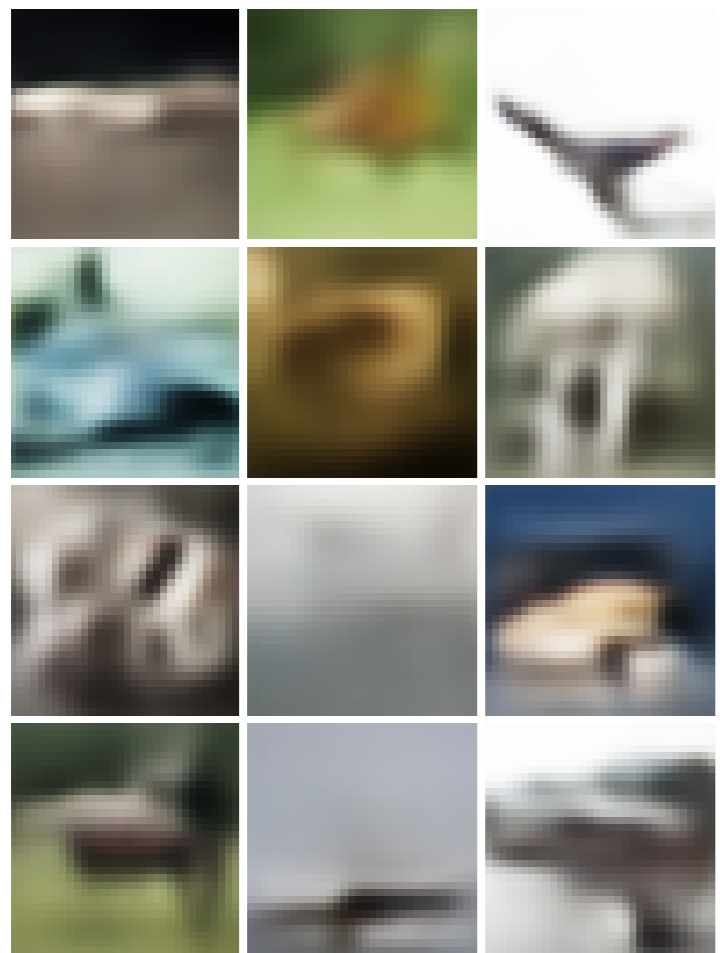


Samples from
Generative Adversarial
[Goodfellow *et al*, 2014]



Samples from
diffusion model

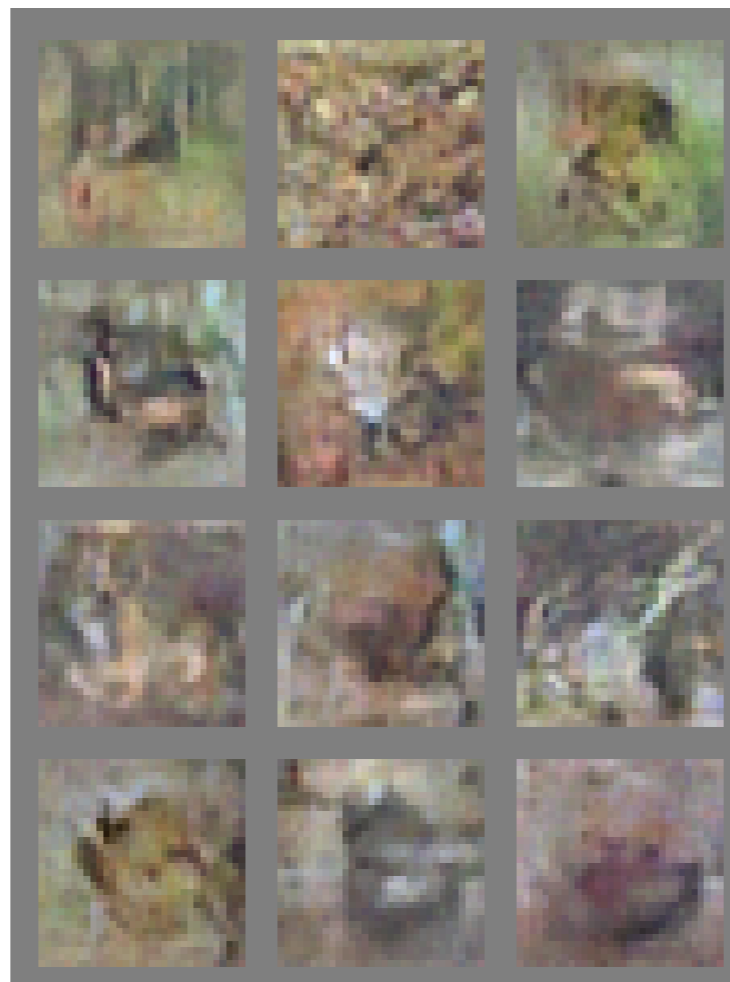
Diffusion Probabilistic Model Applied to CIFAR-10



Samples from
DRAW

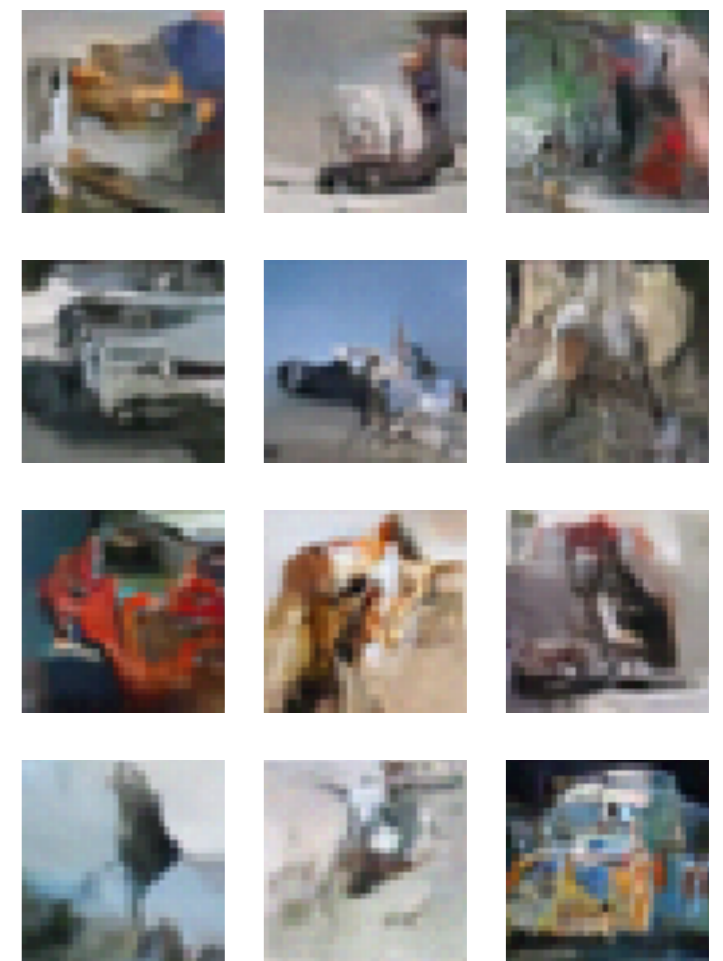
[Gregor *et al*, 2015]

Jascha Sohl-Dickstein



Samples from

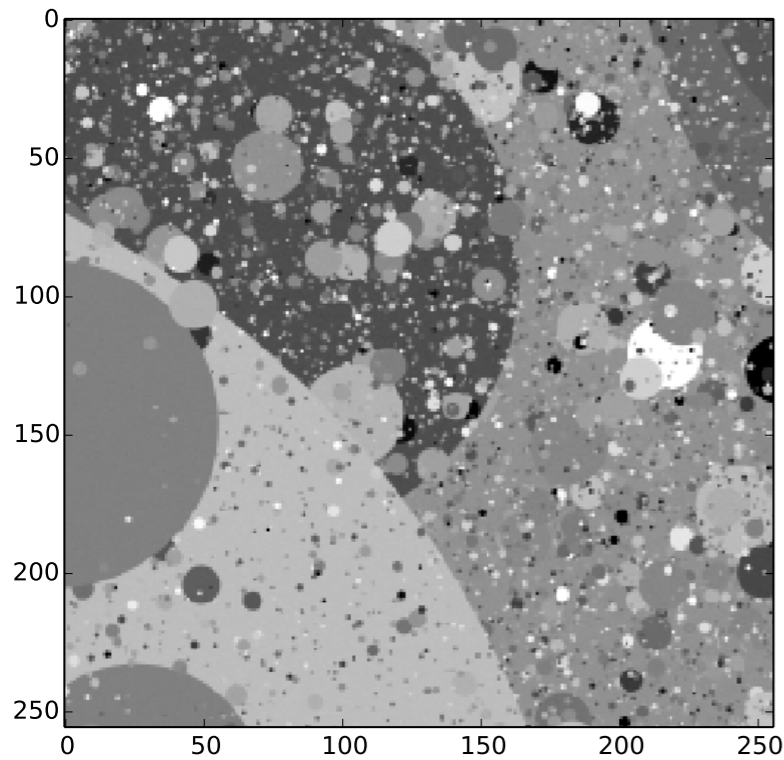
Generative Adversarial
[Goodfellow *et al*, 2014]



Samples from
diffusion model

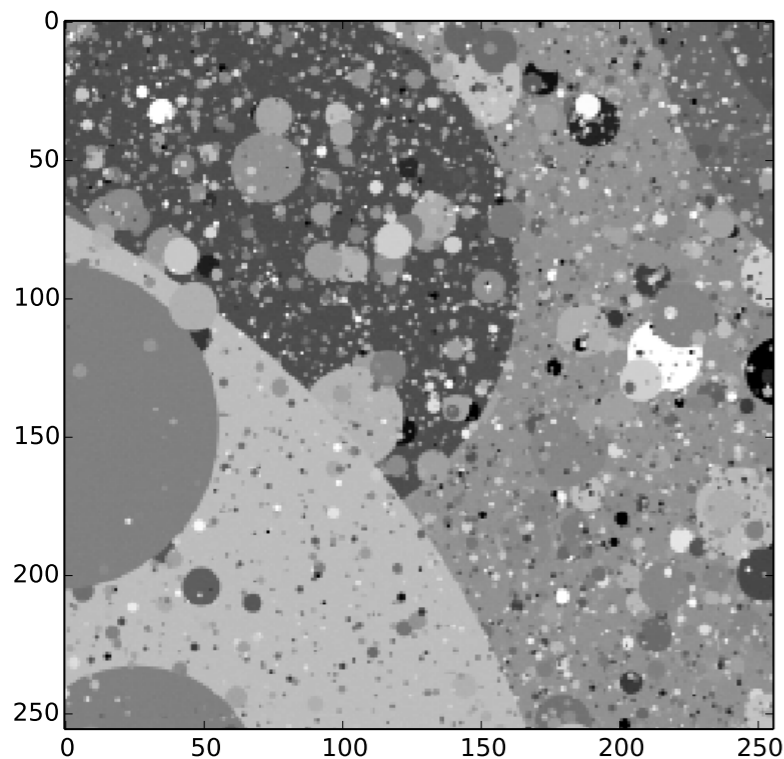
Diffusion Probabilistic Models

Diffusion Probabilistic Model Applied to Dead Leaves

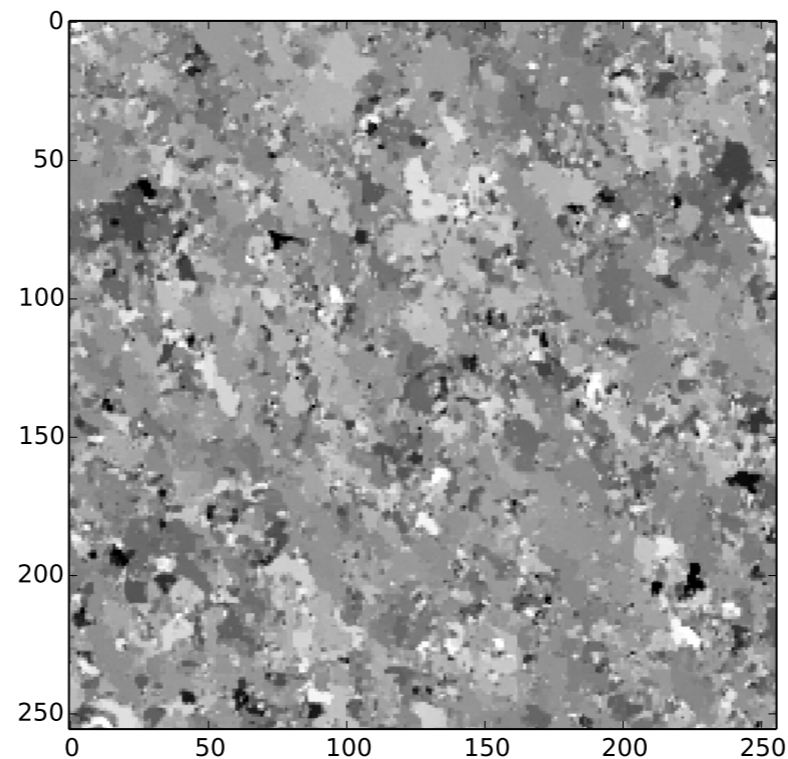


Training Data

Diffusion Probabilistic Model Applied to Dead Leaves

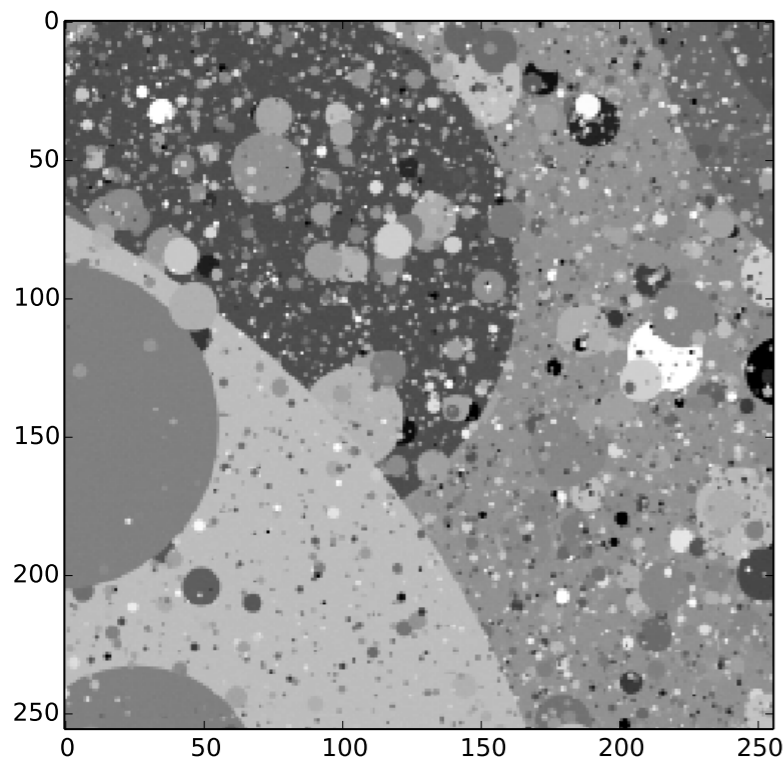


Training Data

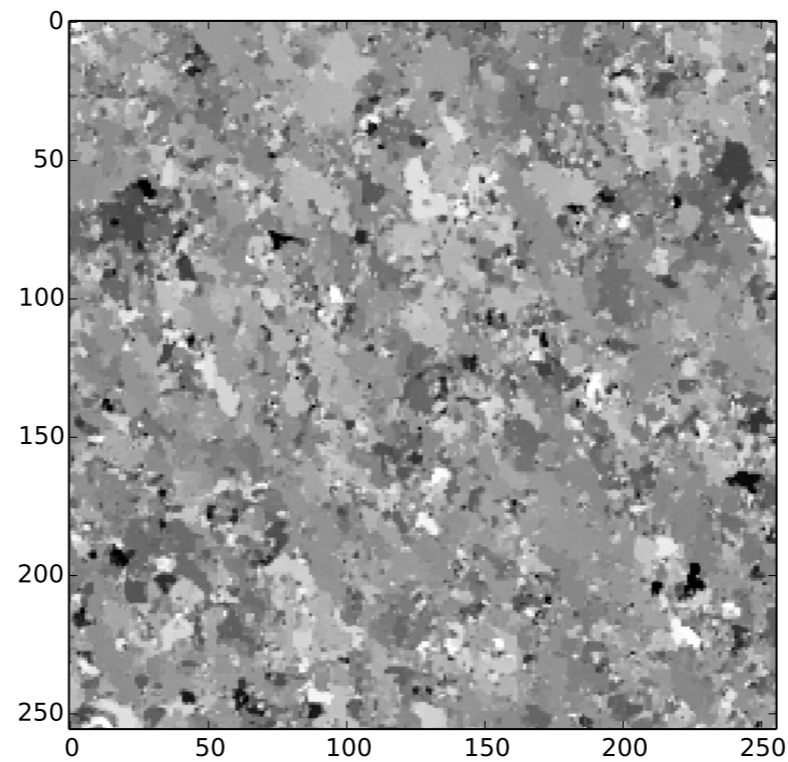


Sample from
[Theis *et al*, 2012]

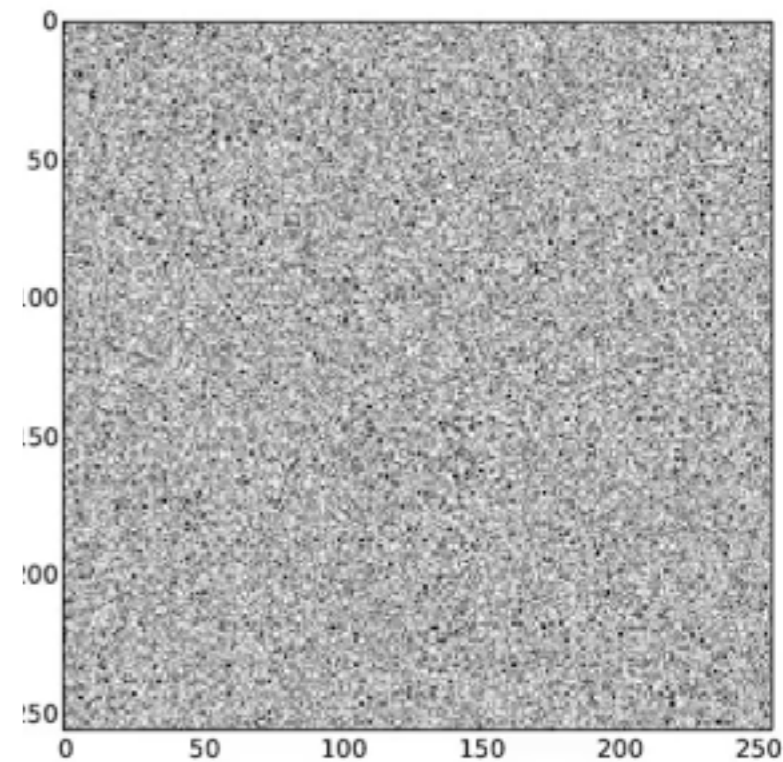
Diffusion Probabilistic Model Applied to Dead Leaves



Training Data

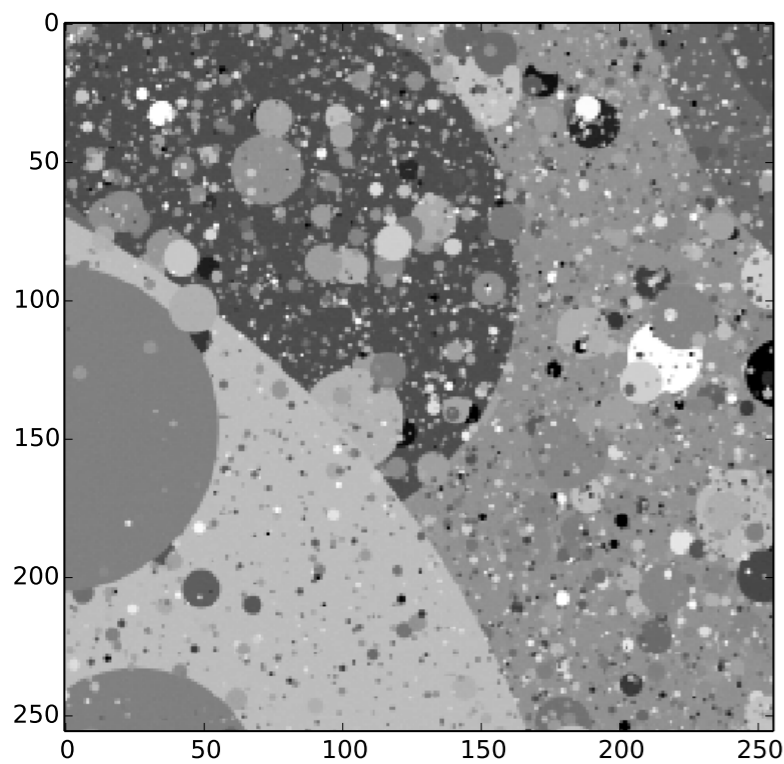


Sample from
[Theis *et al*, 2012]

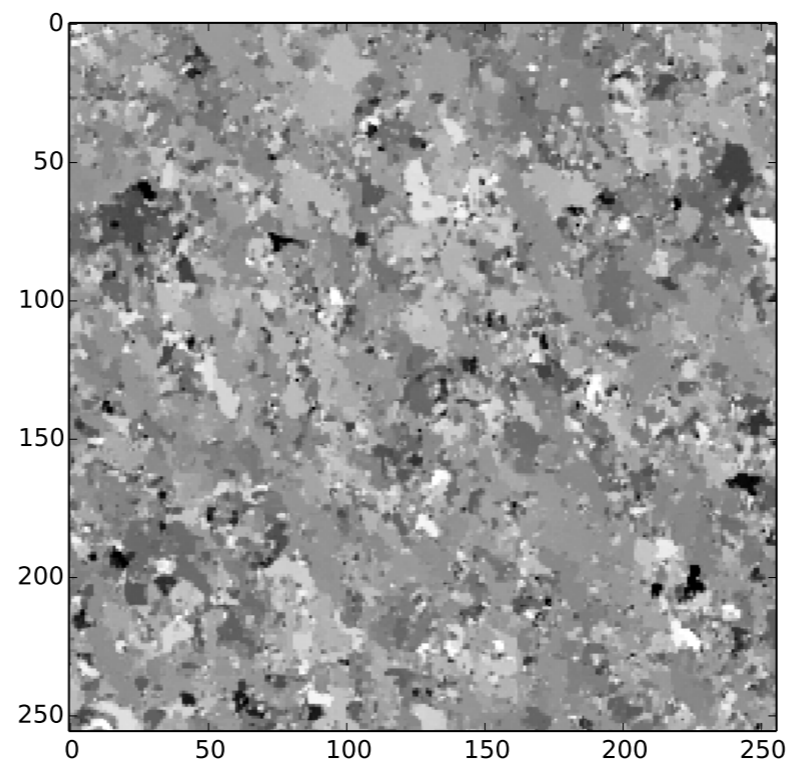


Diffusion Probabilistic Models

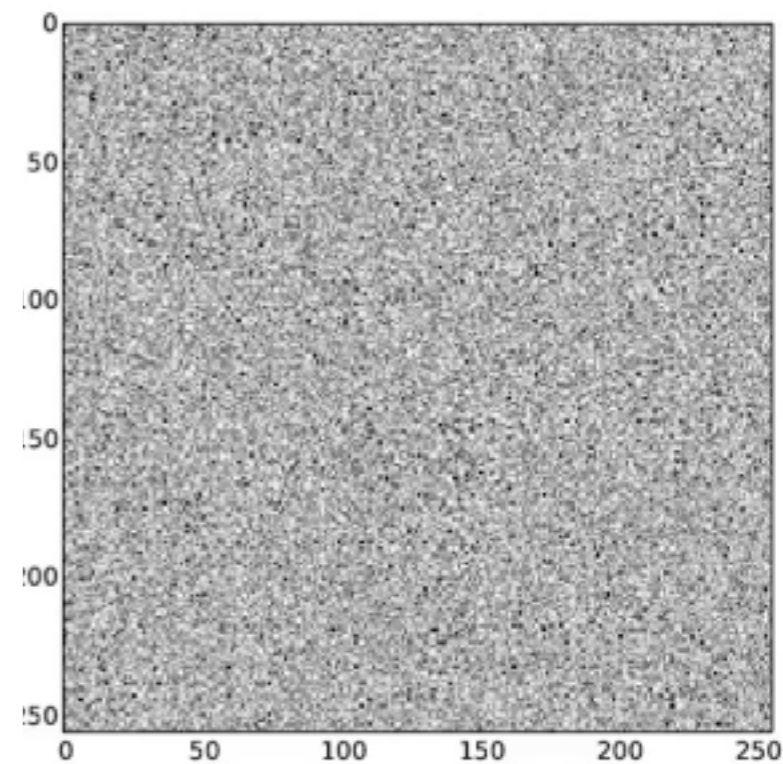
Diffusion Probabilistic Model Applied to Dead Leaves



Training Data

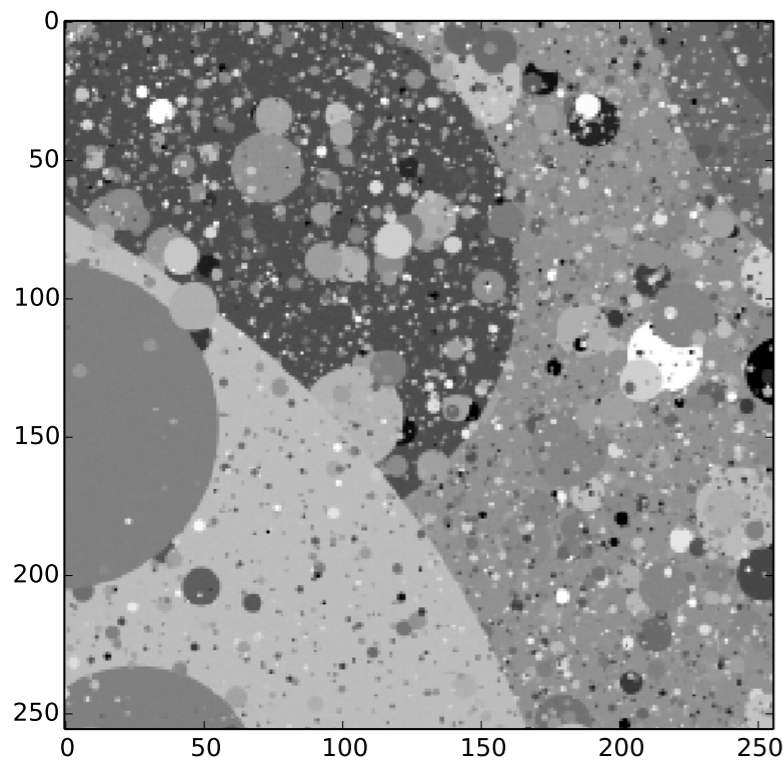


Sample from
[Theis *et al*, 2012]



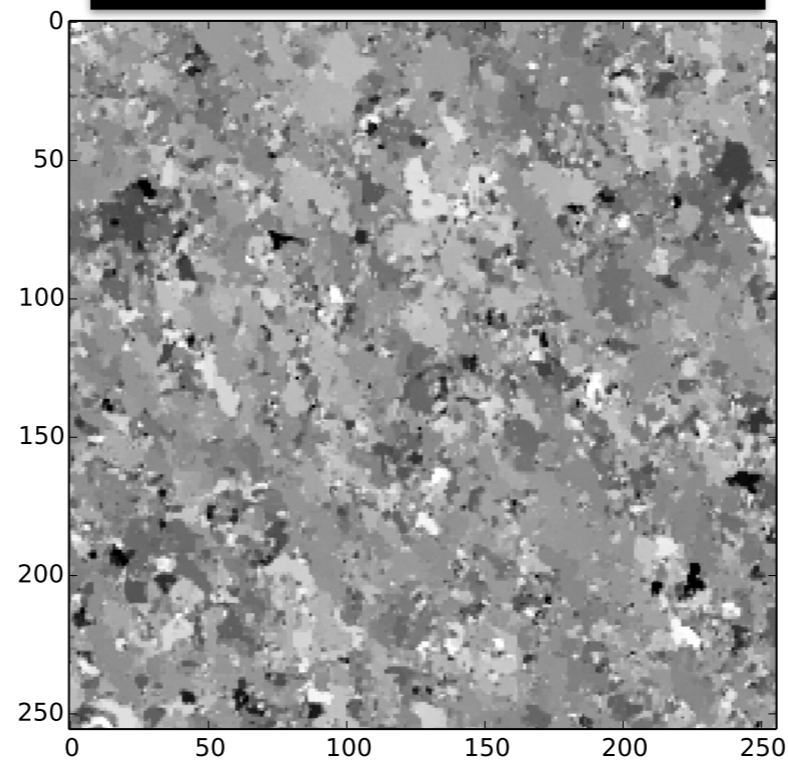
Sample from
diffusion model

Diffusion Probabilistic Model Applied to Dead Leaves



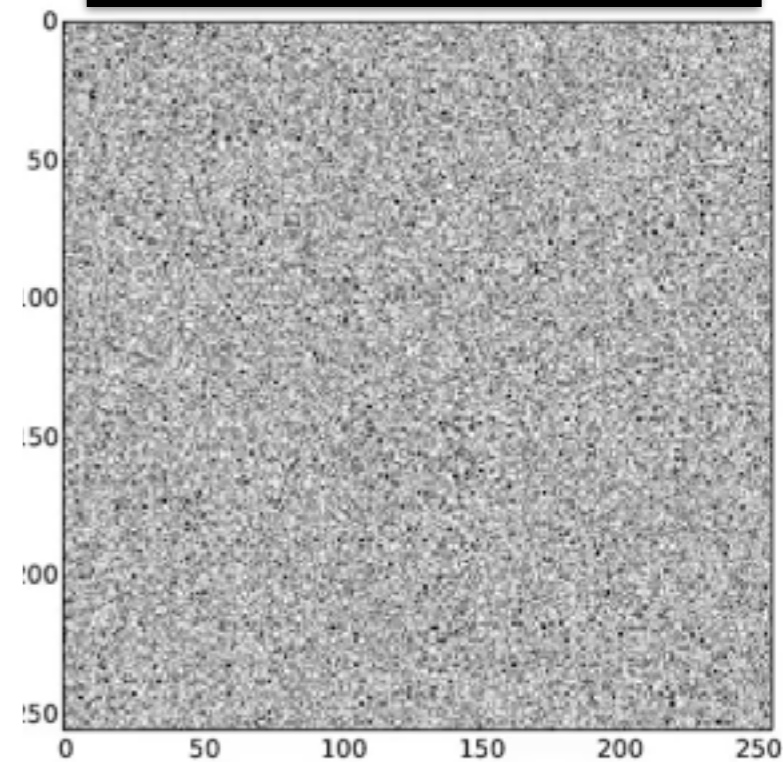
Training Data

Log likelihood
1.24 bits/pixel



Sample from
[Theis *et al*, 2012]

Log likelihood
1.49 bits/pixel



Sample from
diffusion model

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

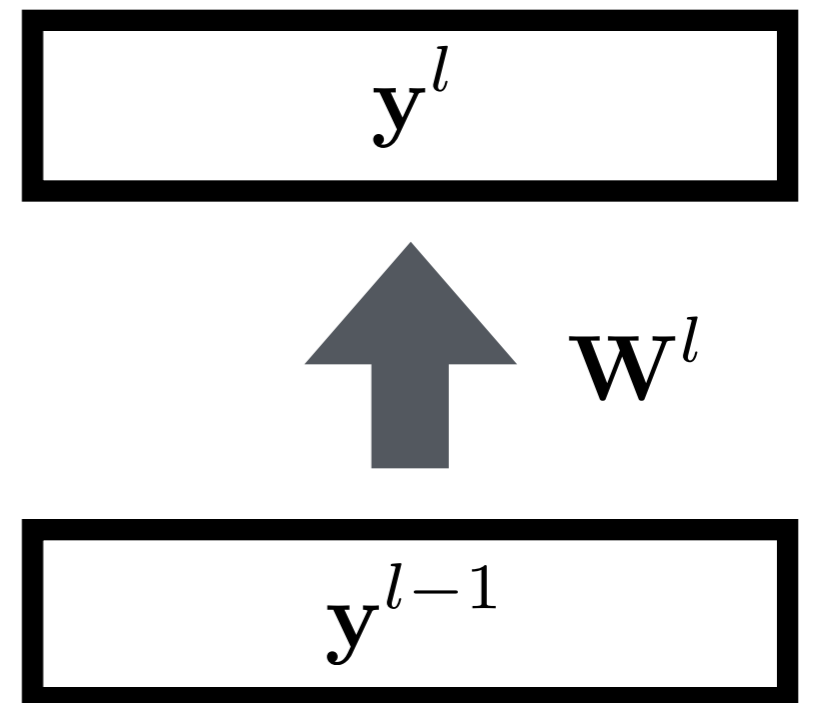
- Extremely flexible, parametric, function approximation
- **Single layer:** linear transformation, pointwise nonlinearity

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

Deep Networks

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

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



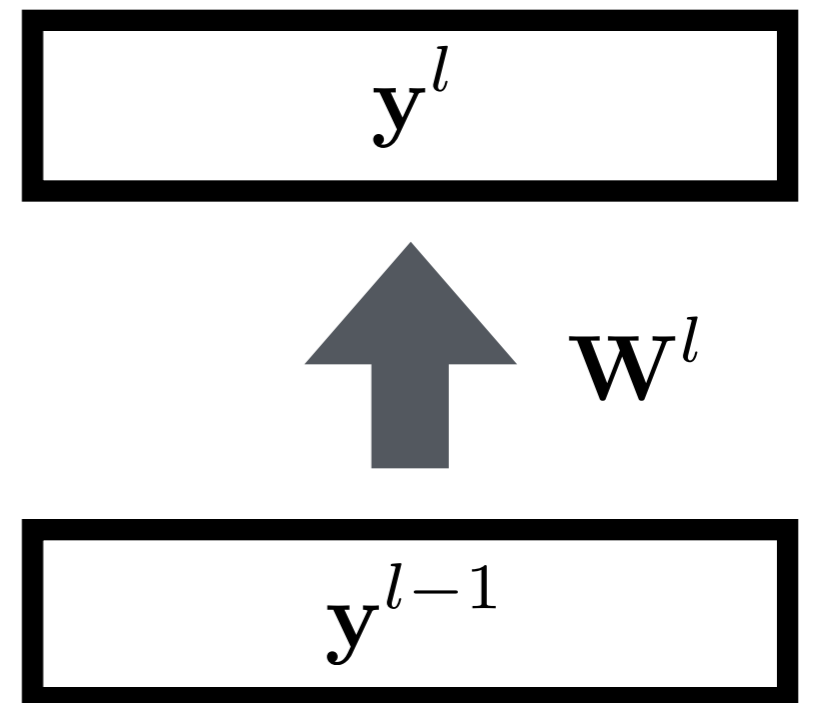
Deep Networks

- 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

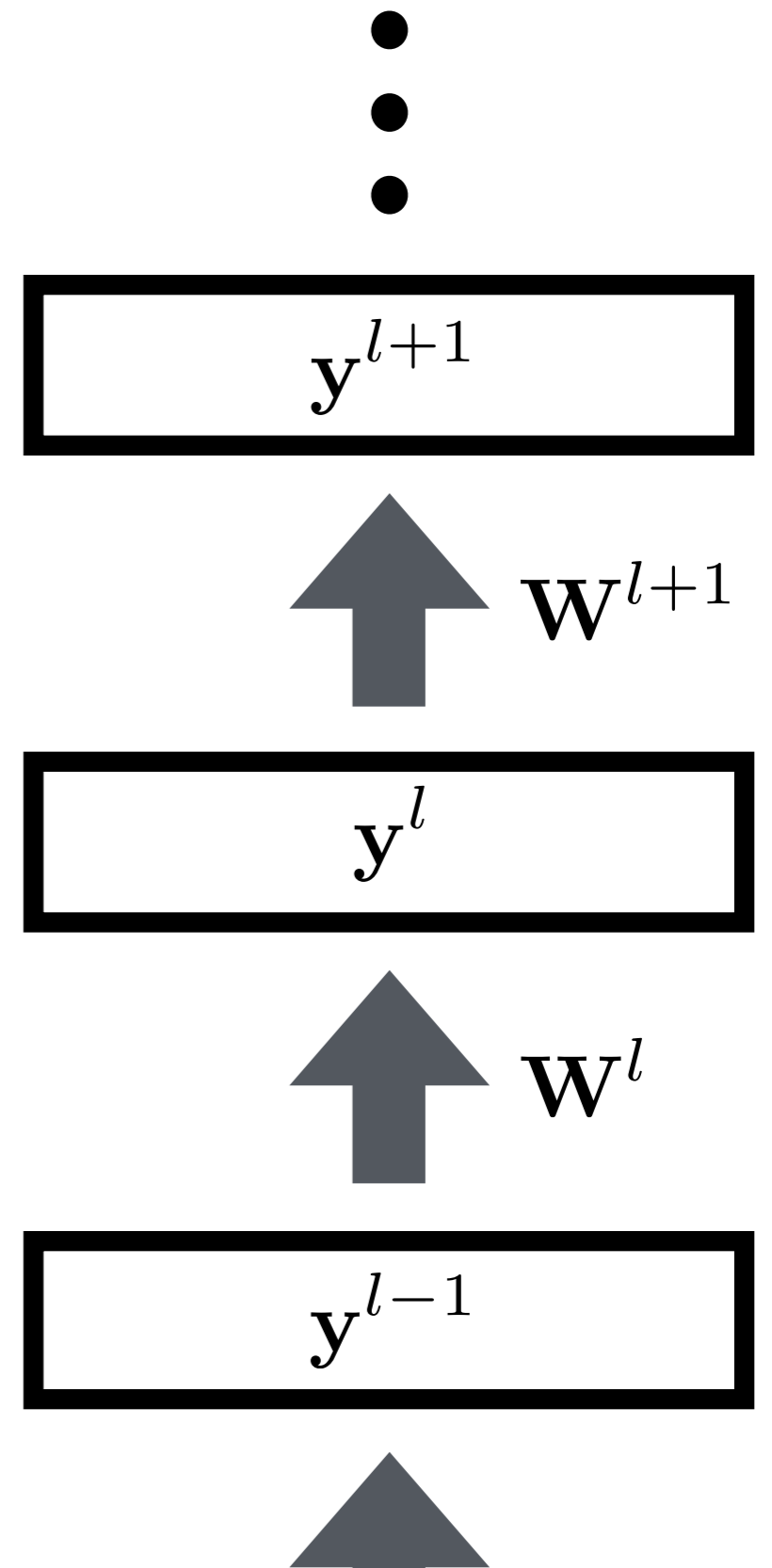
- Extremely flexible, parametric, function approximation
- **Single layer:** linear transformation, pointwise nonlinearity

Deep Networks

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

Deep Networks

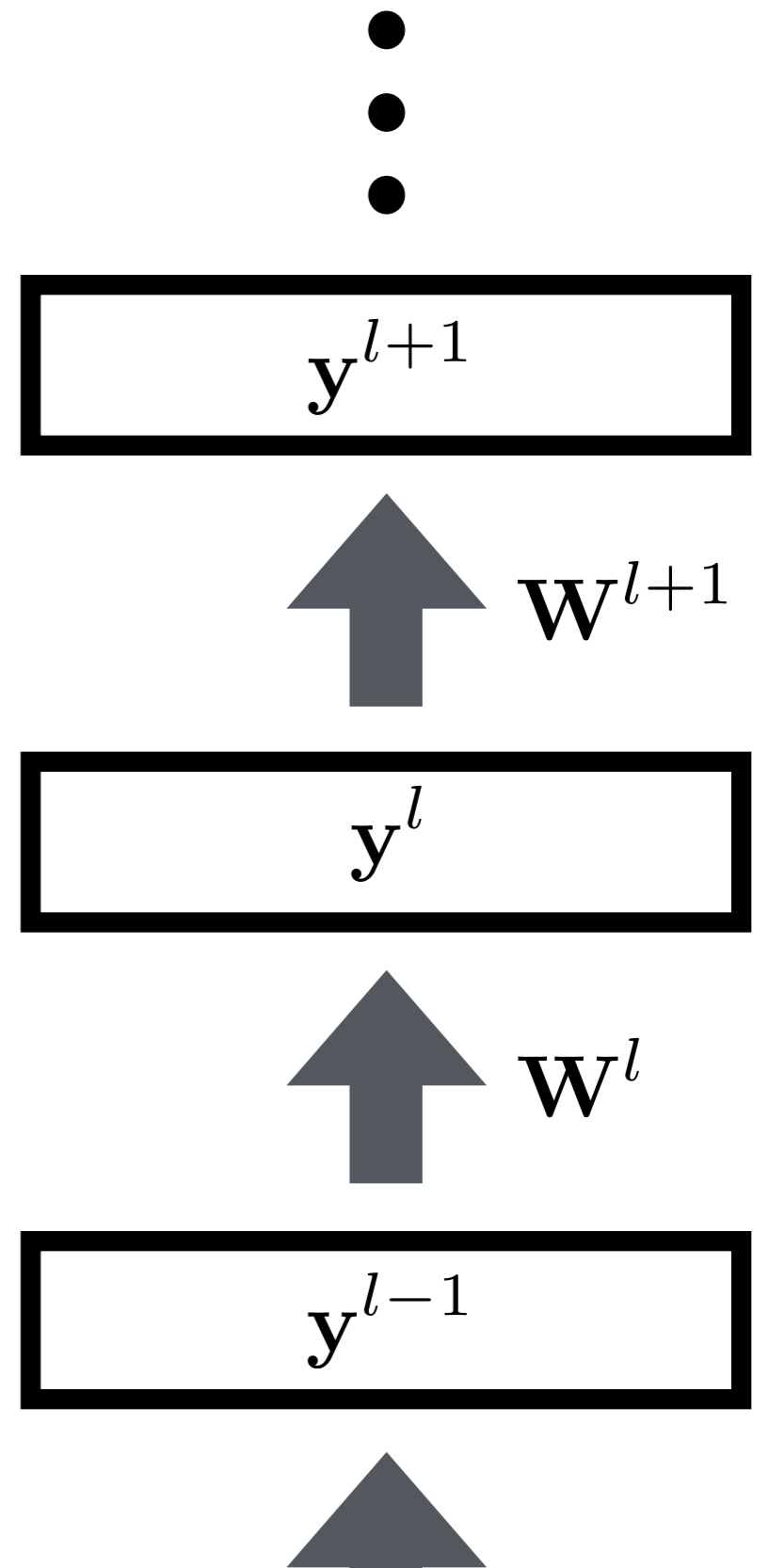
- Extremely flexible, parametric, function approximation
- **Single layer:** linear transformation, pointwise nonlinearity
- **Deep network:** stack single layers



Deep Networks

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

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

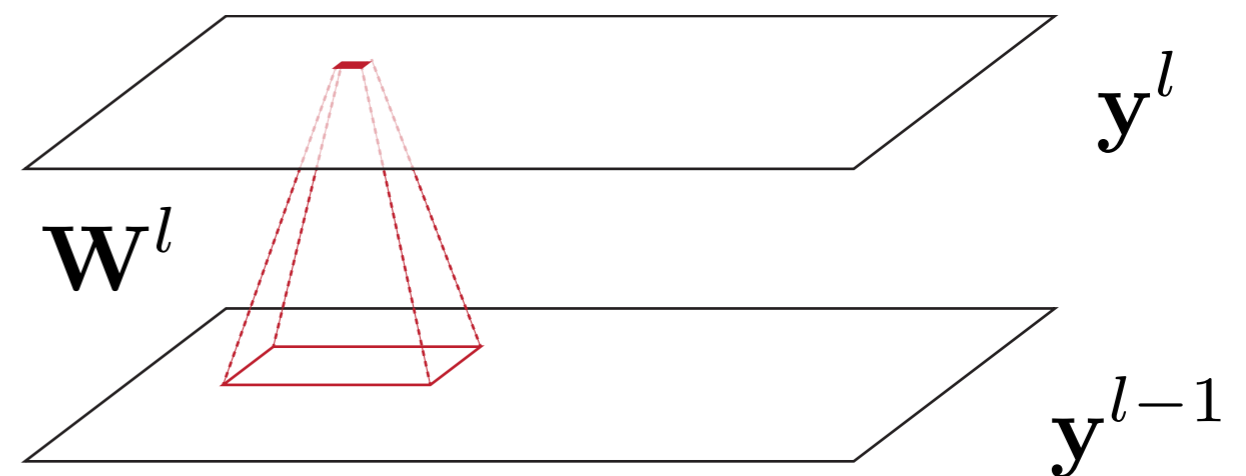


Convolutional Neural Network

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

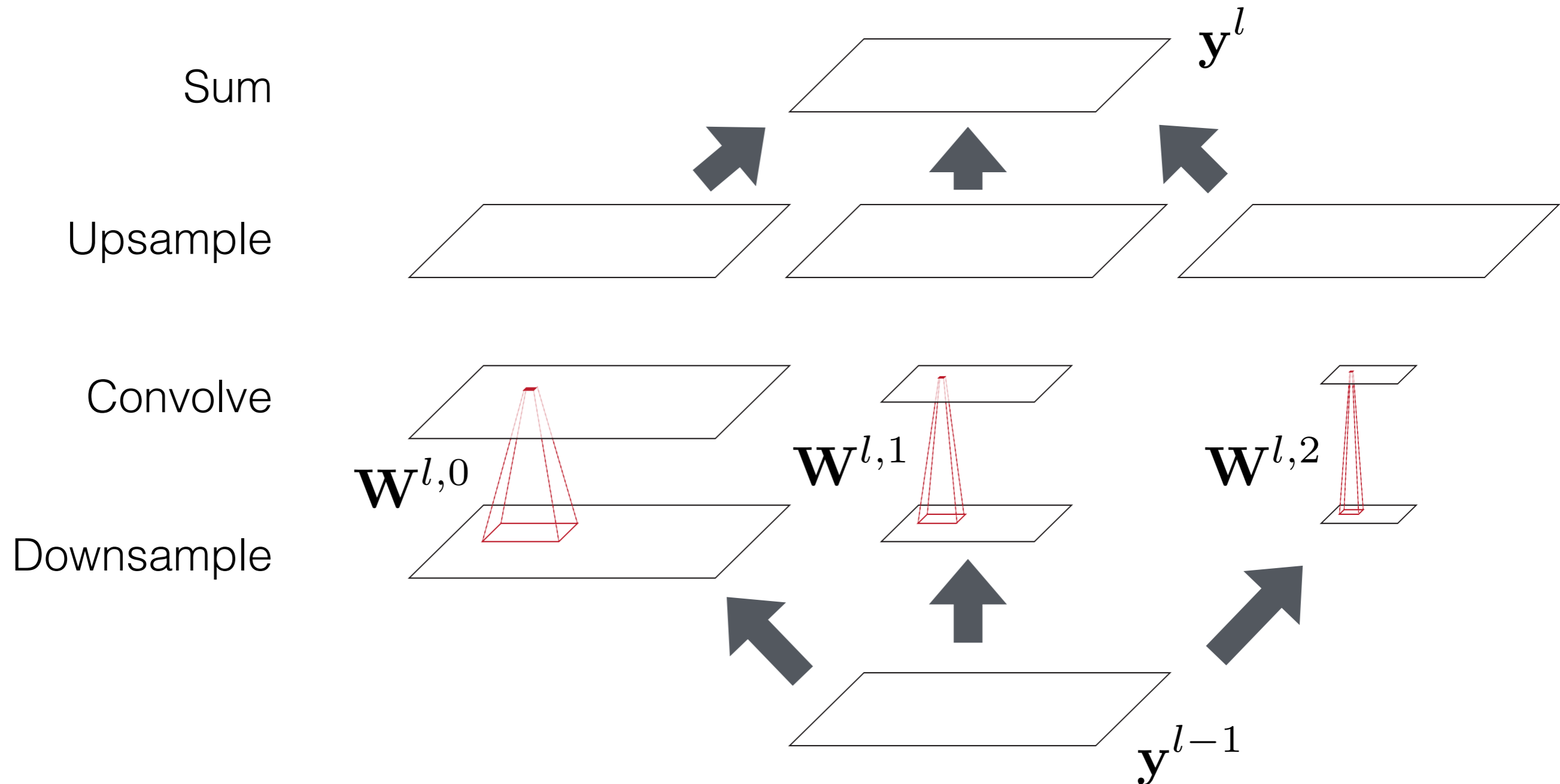
Convolutional Neural Network

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



Multiscale Convolution

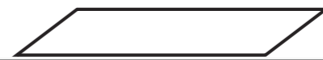
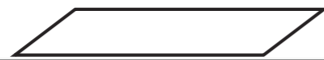
- Single multi-scale convolutional layer:



Deep Network Architecture for Diffusion

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

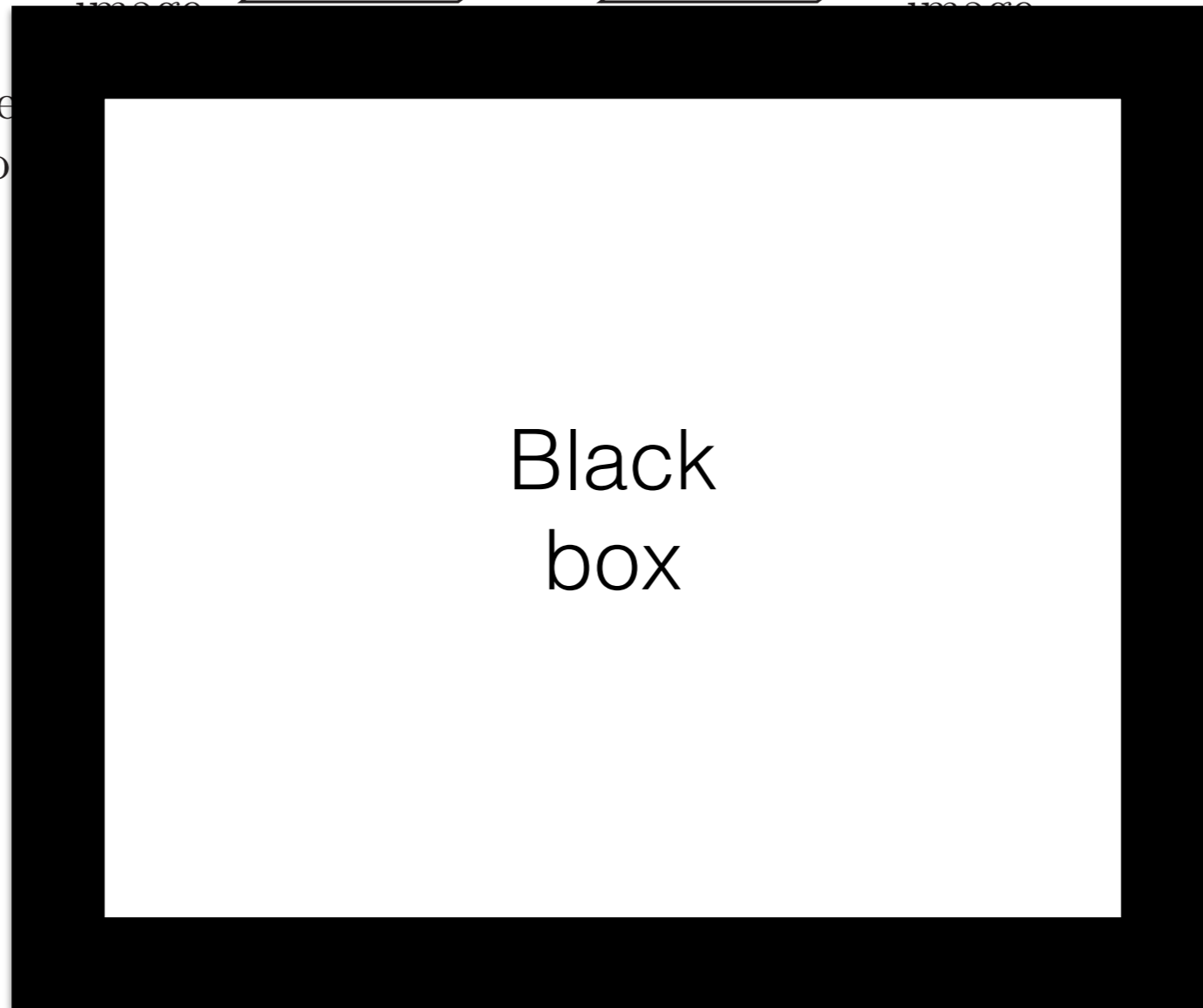
Mean



Covariance

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

Te
co



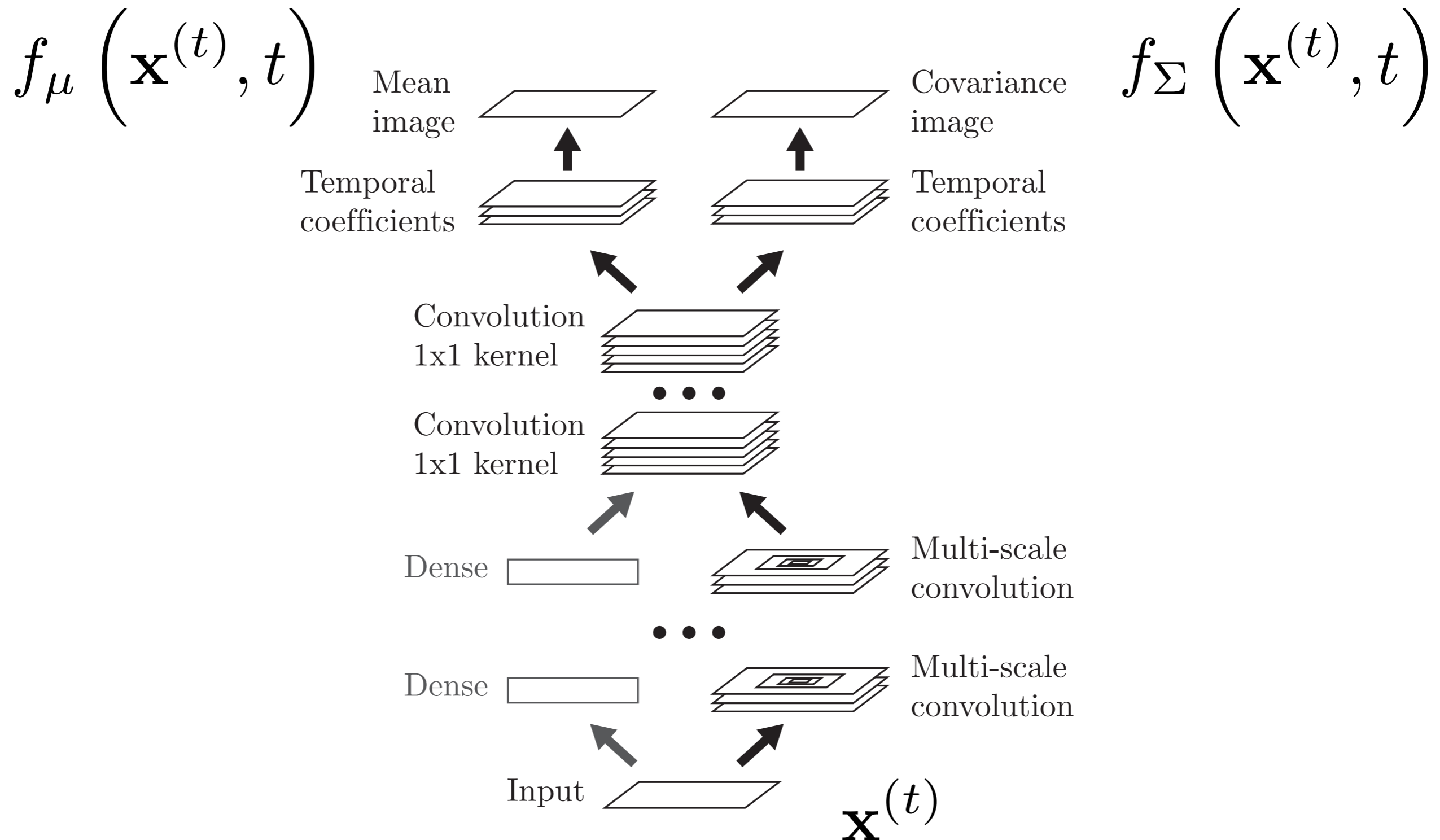
Black
box

Input

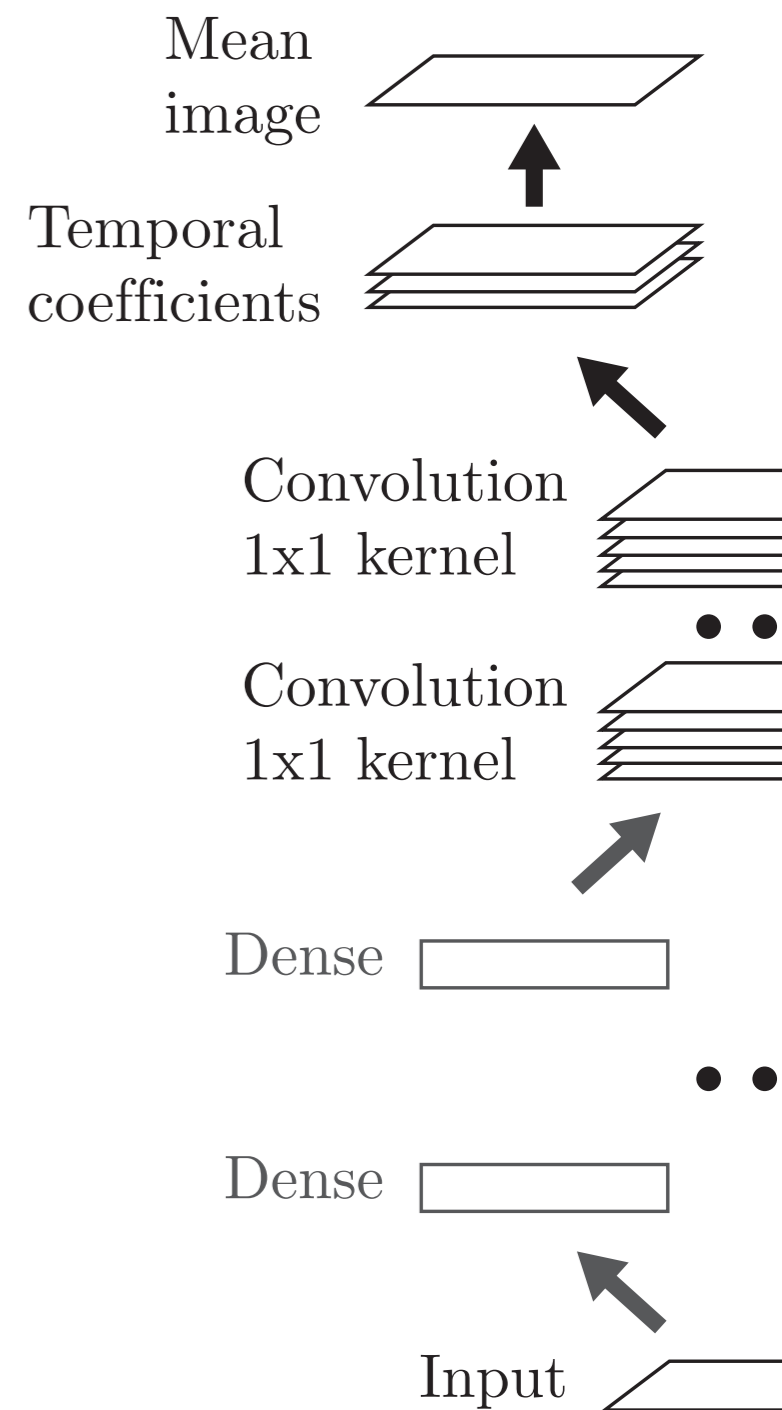


$\mathbf{x}^{(t)}$

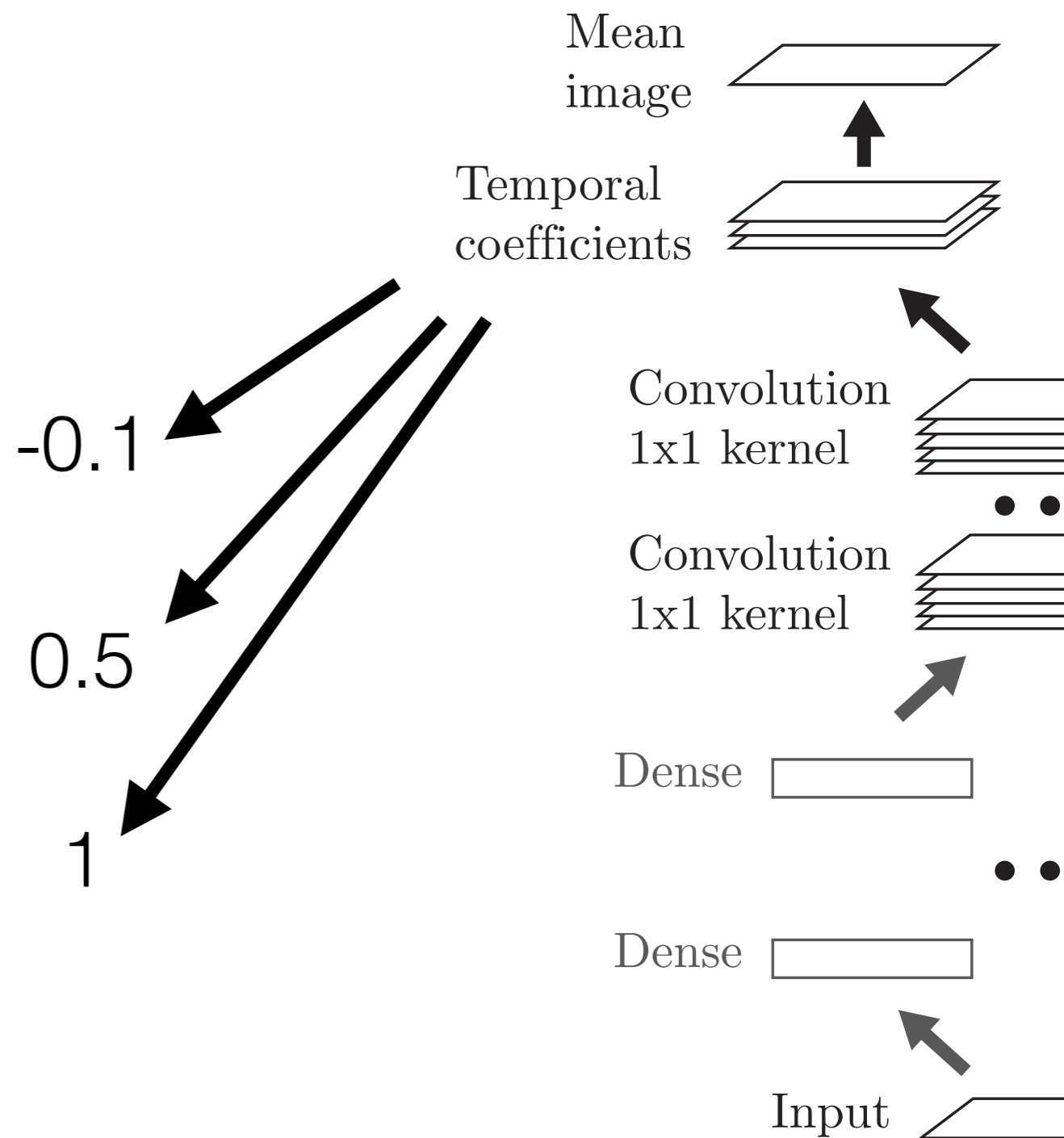
Deep Network Architecture for Diffusion



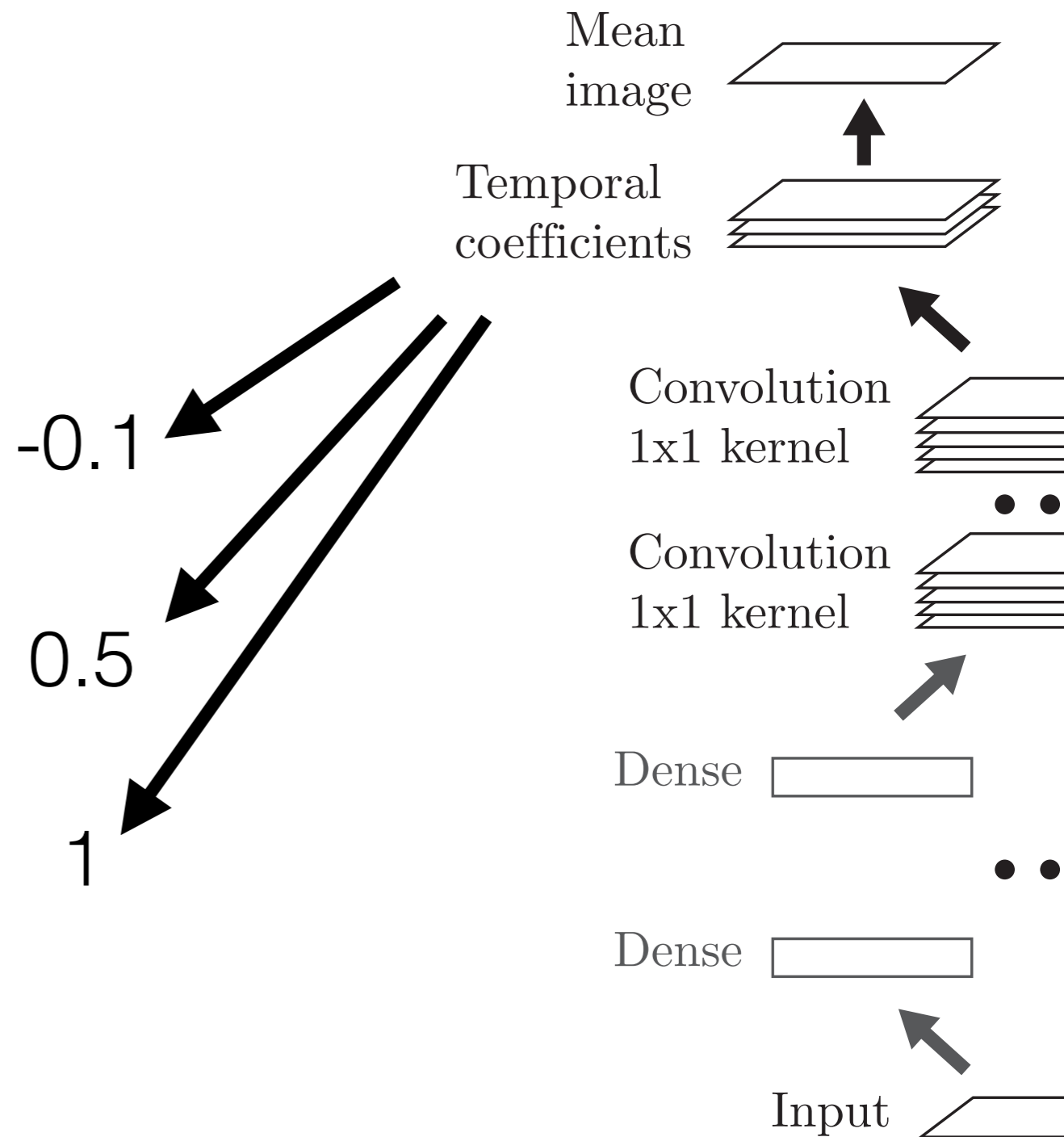
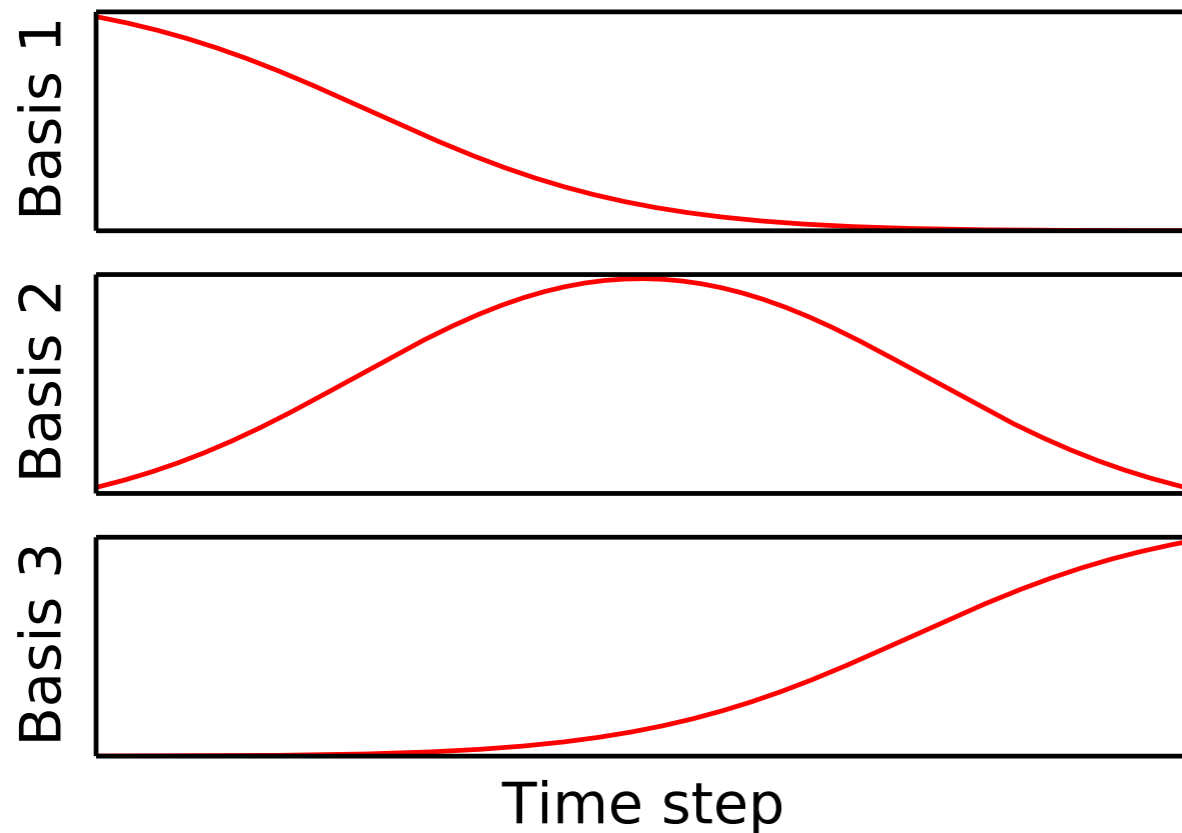
Time Dependence using Temporal Basis



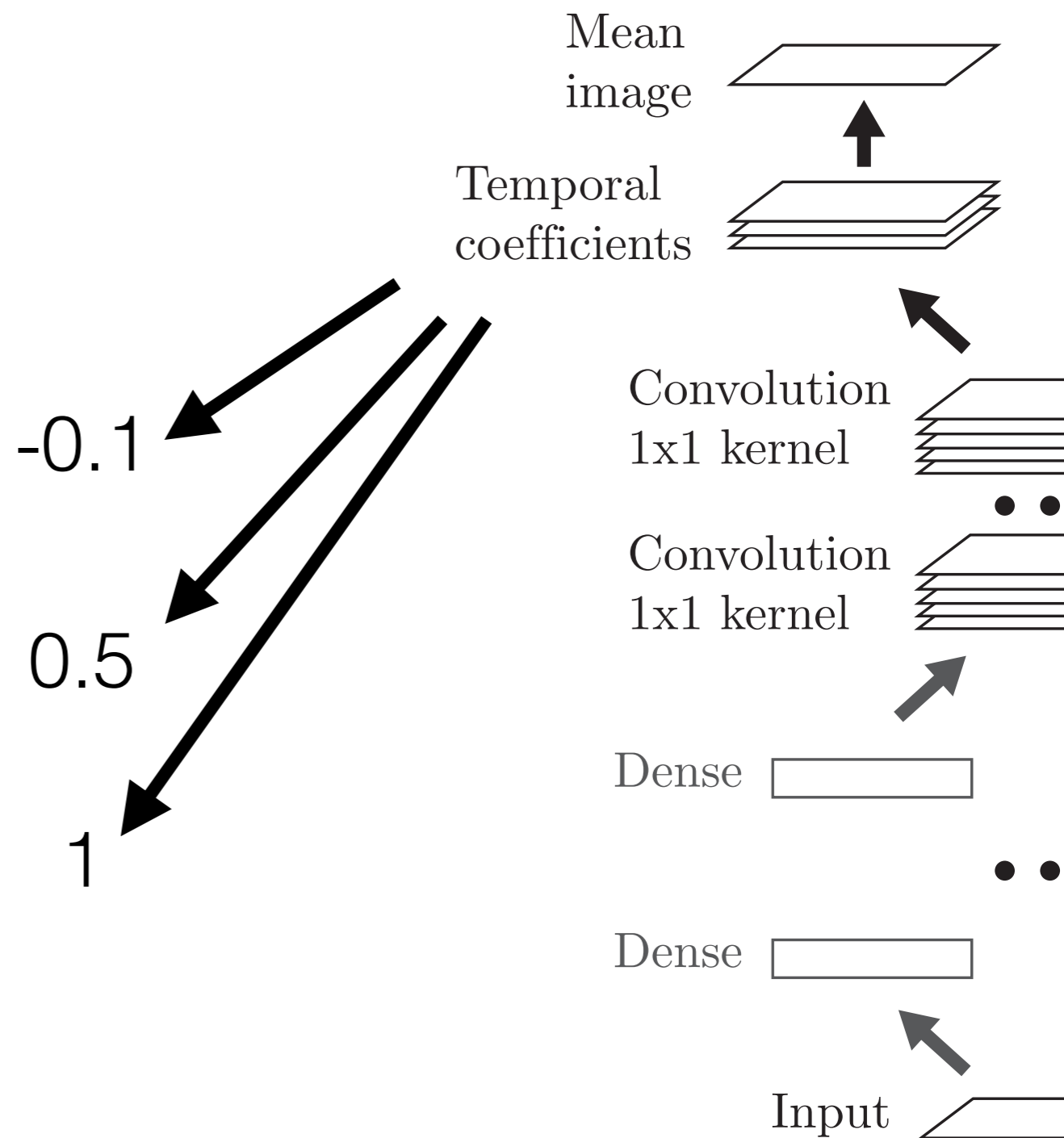
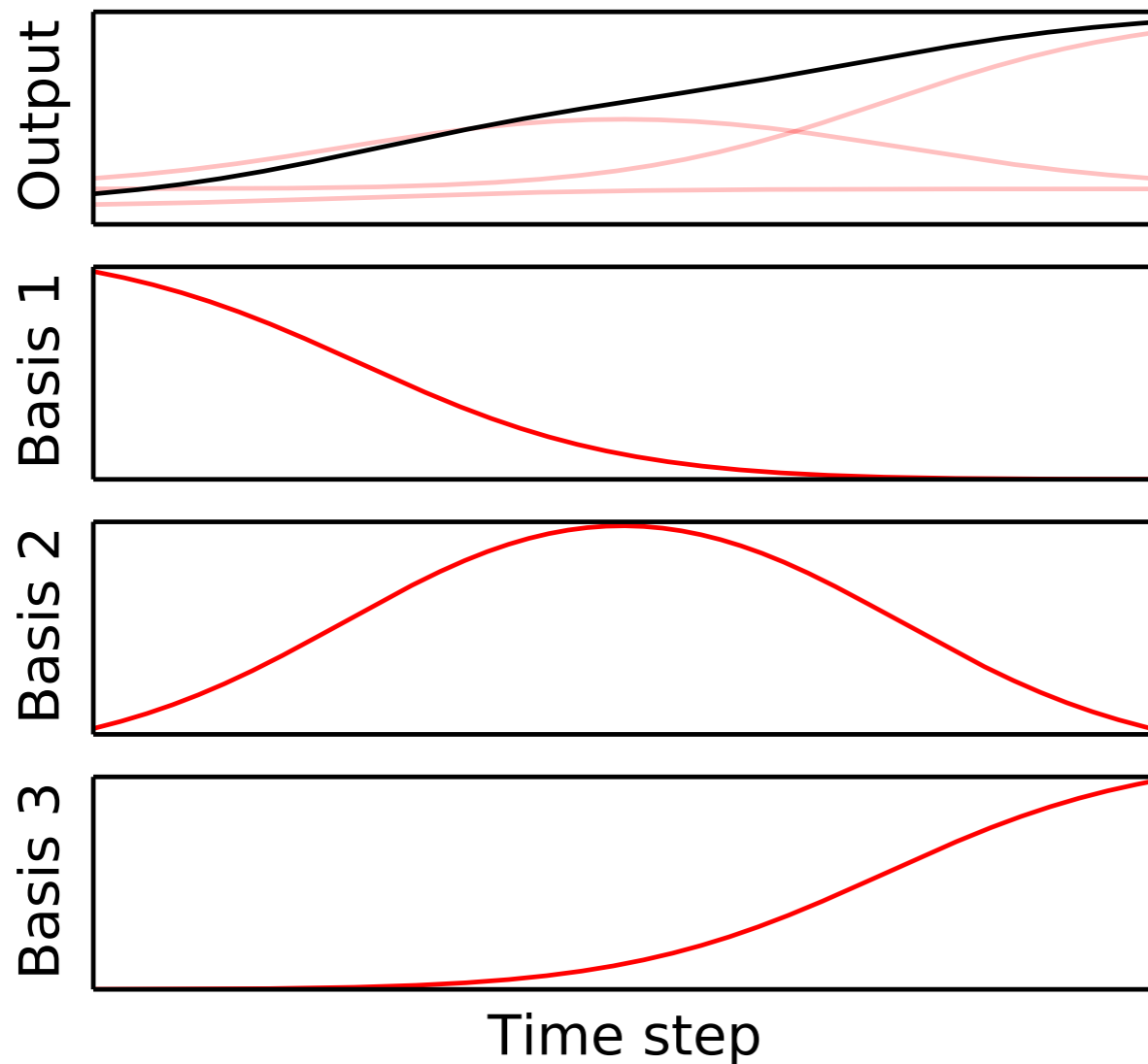
Time Dependence using Temporal Basis



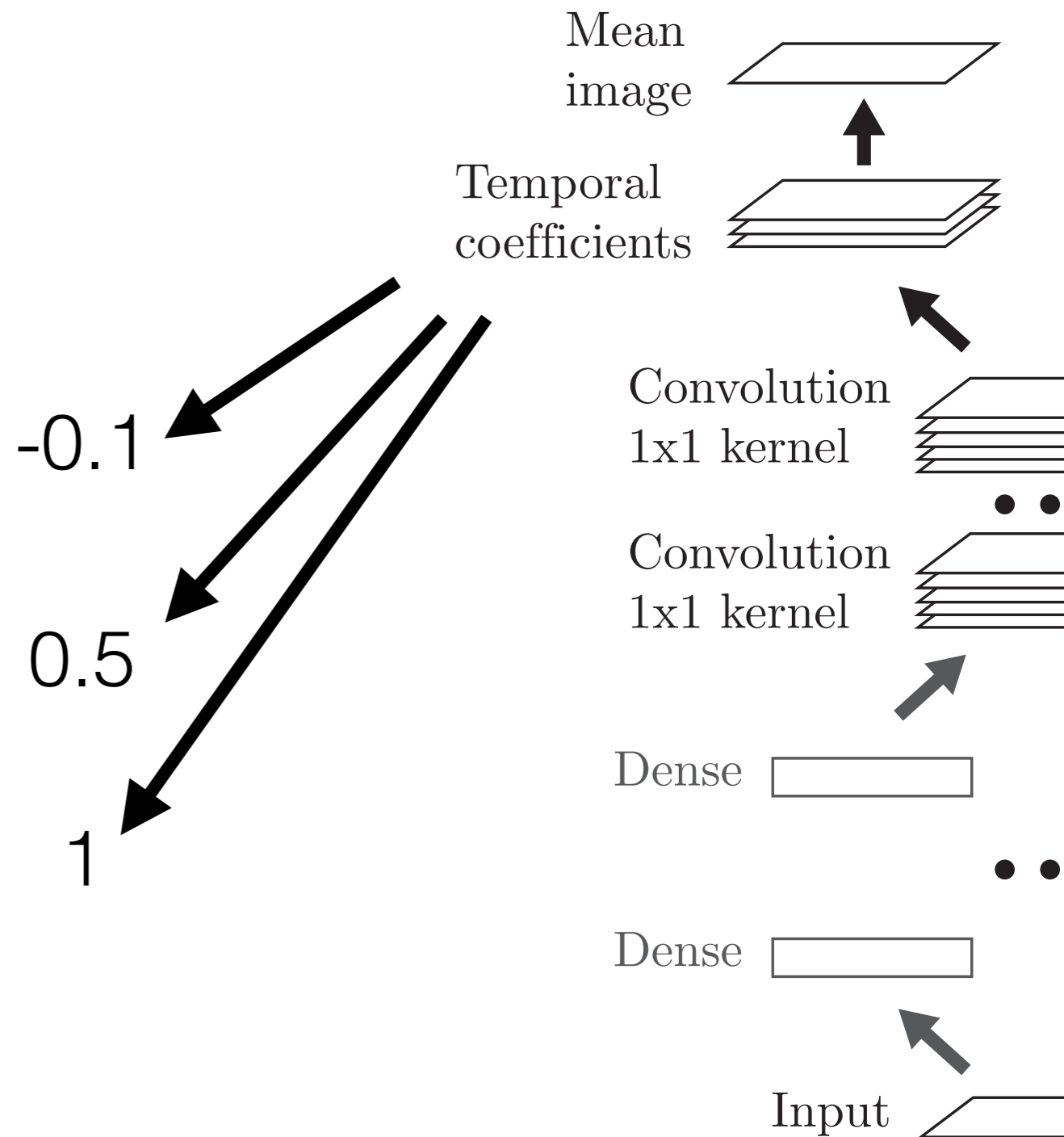
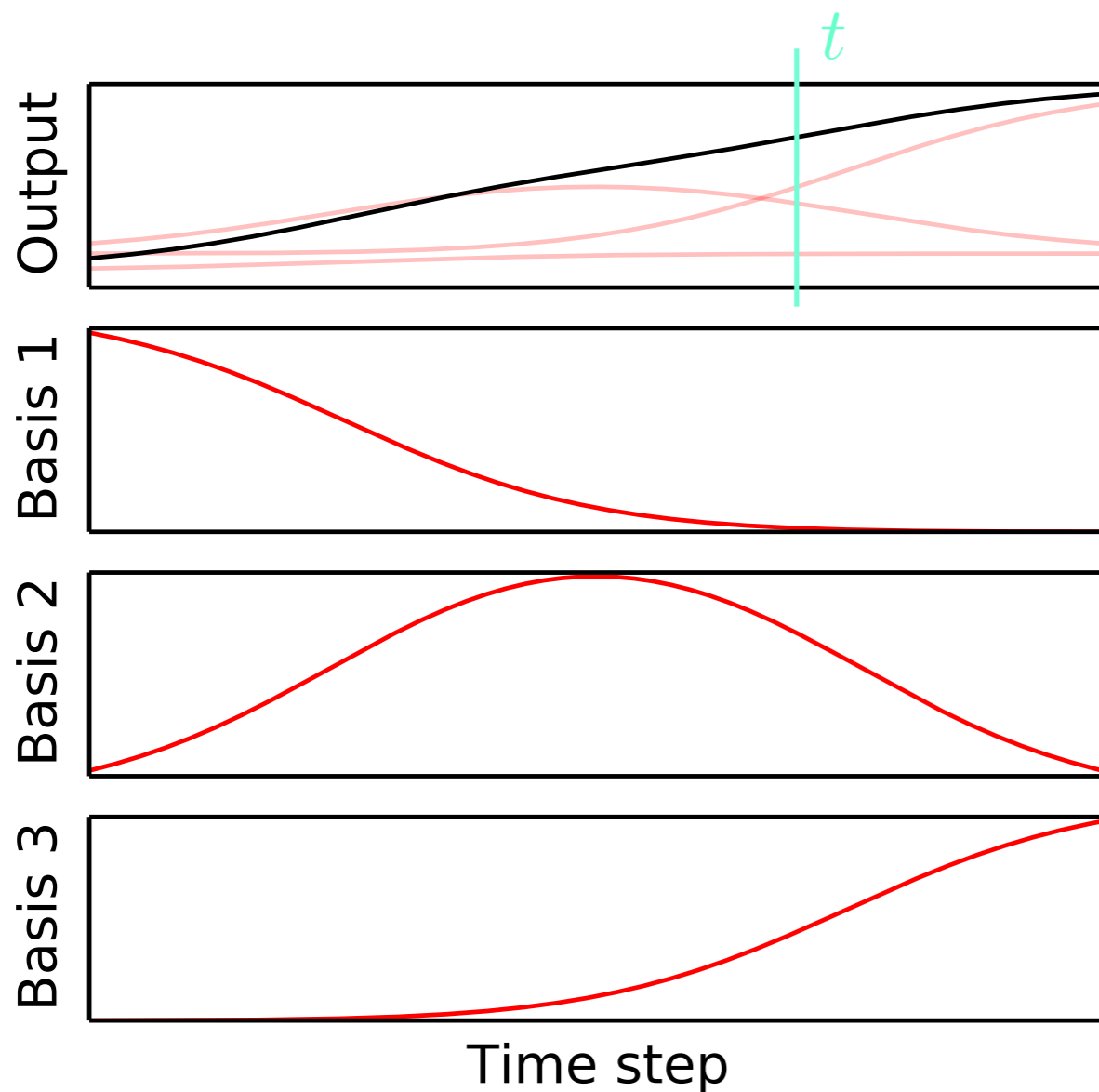
Time Dependence using Temporal Basis



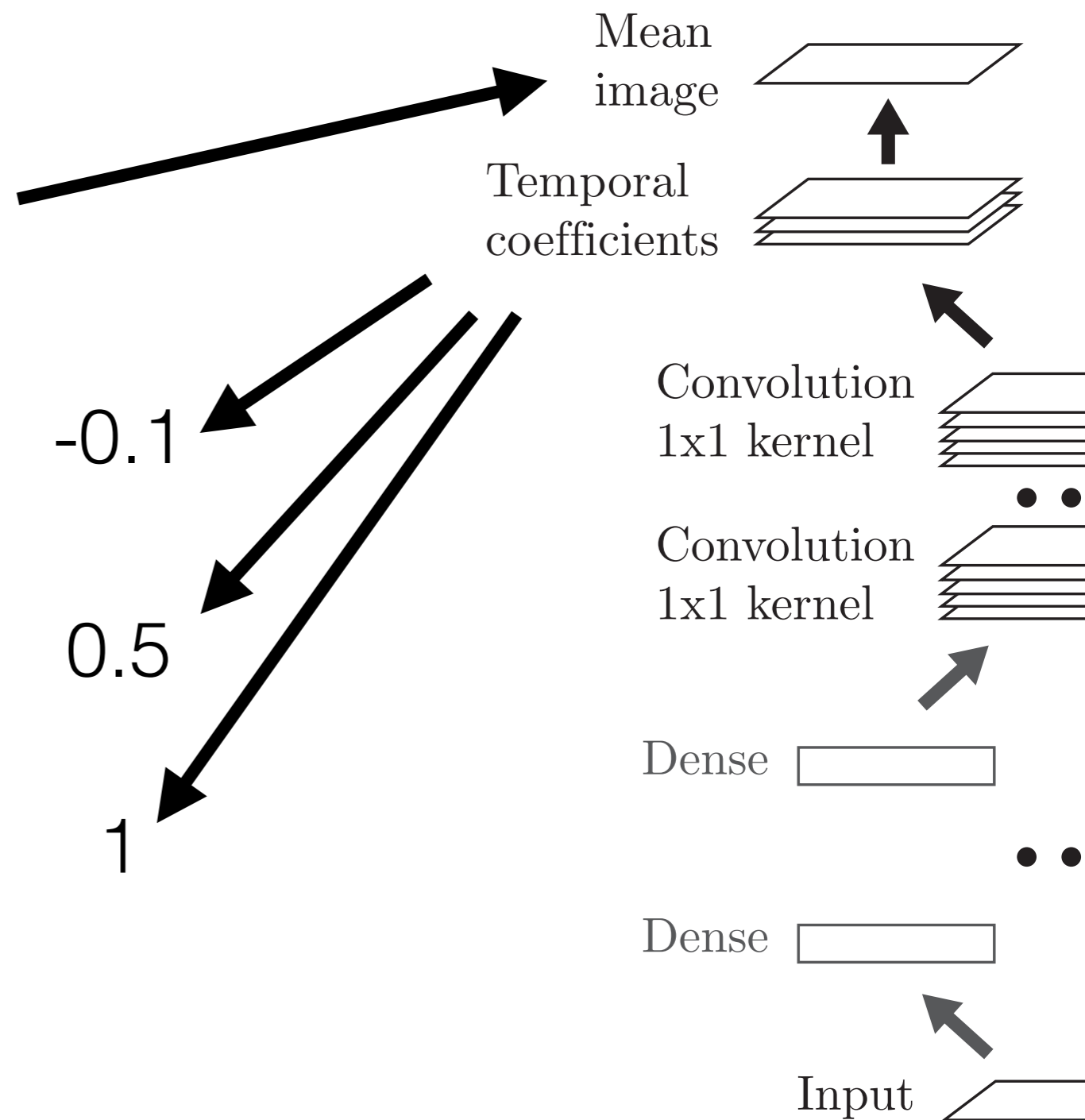
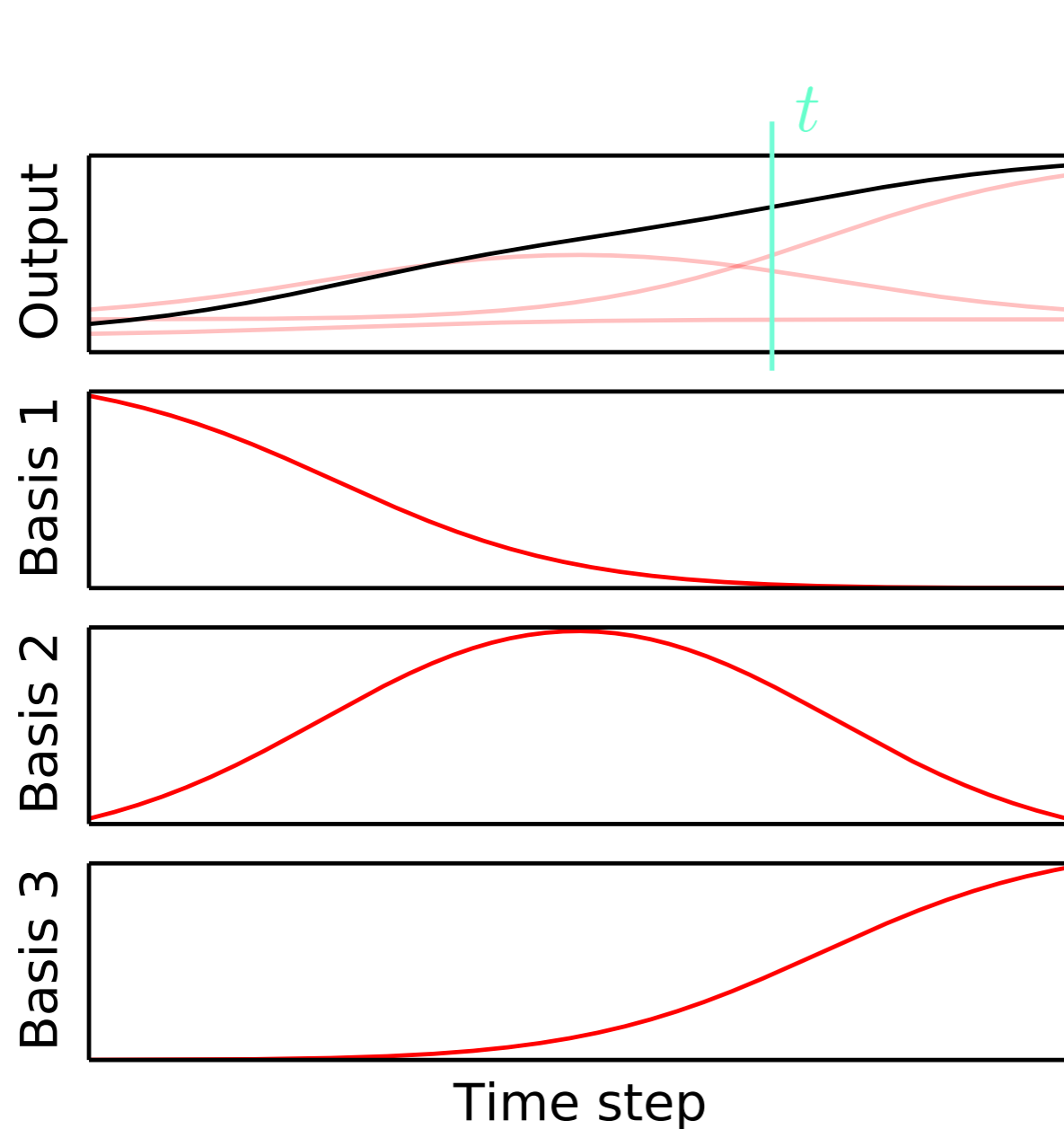
Time Dependence using Temporal Basis



Time Dependence using Temporal Basis



Time Dependence using Temporal Basis



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**

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. variational autoencoders, GSNs, NADEs, 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

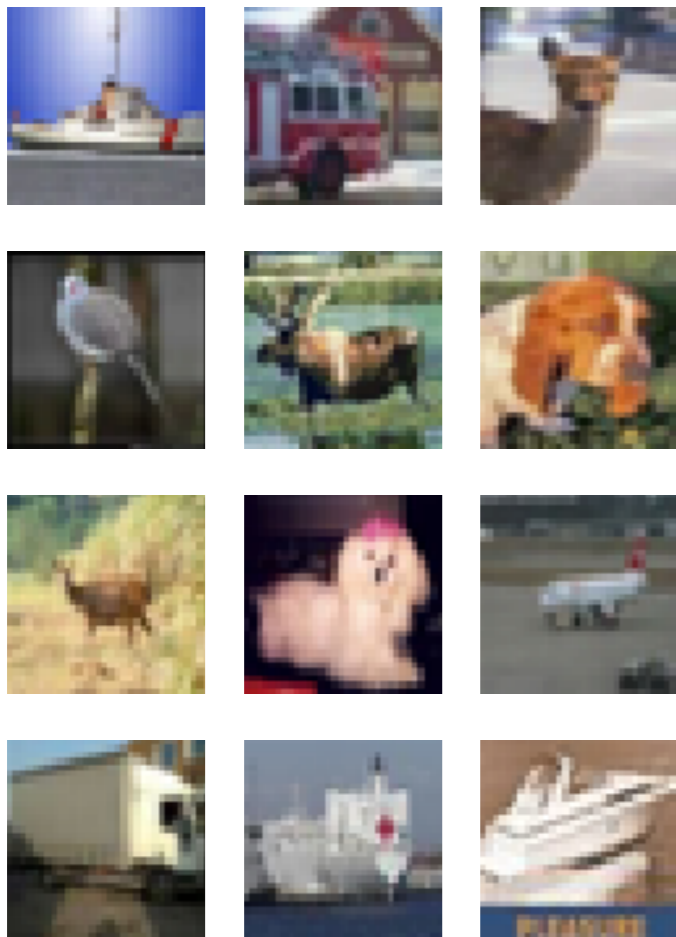
Interested in $\tilde{p}(\mathbf{x}^{(0)}) \propto p(\mathbf{x}^{(0)}) \underline{r(\mathbf{x}^{(0)})}$

Acts as small perturbation to diffusion process

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

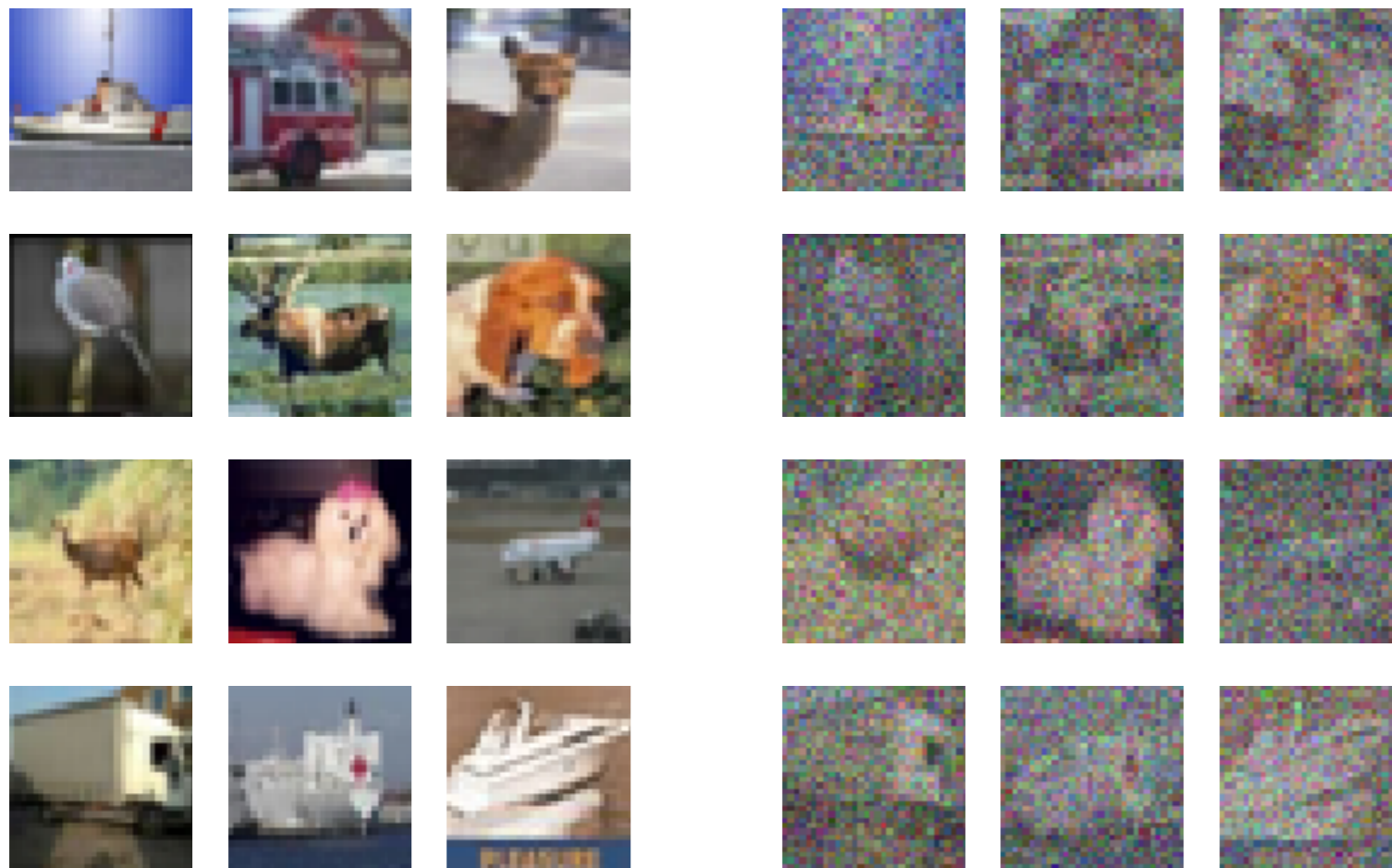
$$\tilde{p}(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}) \approx \mathcal{N}\left(x^{(t-1)}; \mathbf{f}_{\mu}(\mathbf{x}^{(t)}, t) + \mathbf{f}_{\Sigma}(\mathbf{x}^{(t)}, t) \frac{\partial \log r(\mathbf{x}^{(t-1)'})}{\partial \mathbf{x}^{(t-1)'}} \Big|_{\mathbf{x}^{(t-1)' = f_{\mu}(\mathbf{x}^{(t)}, t)},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

Image Denoising by Sampling from Posterior



Holdout Data

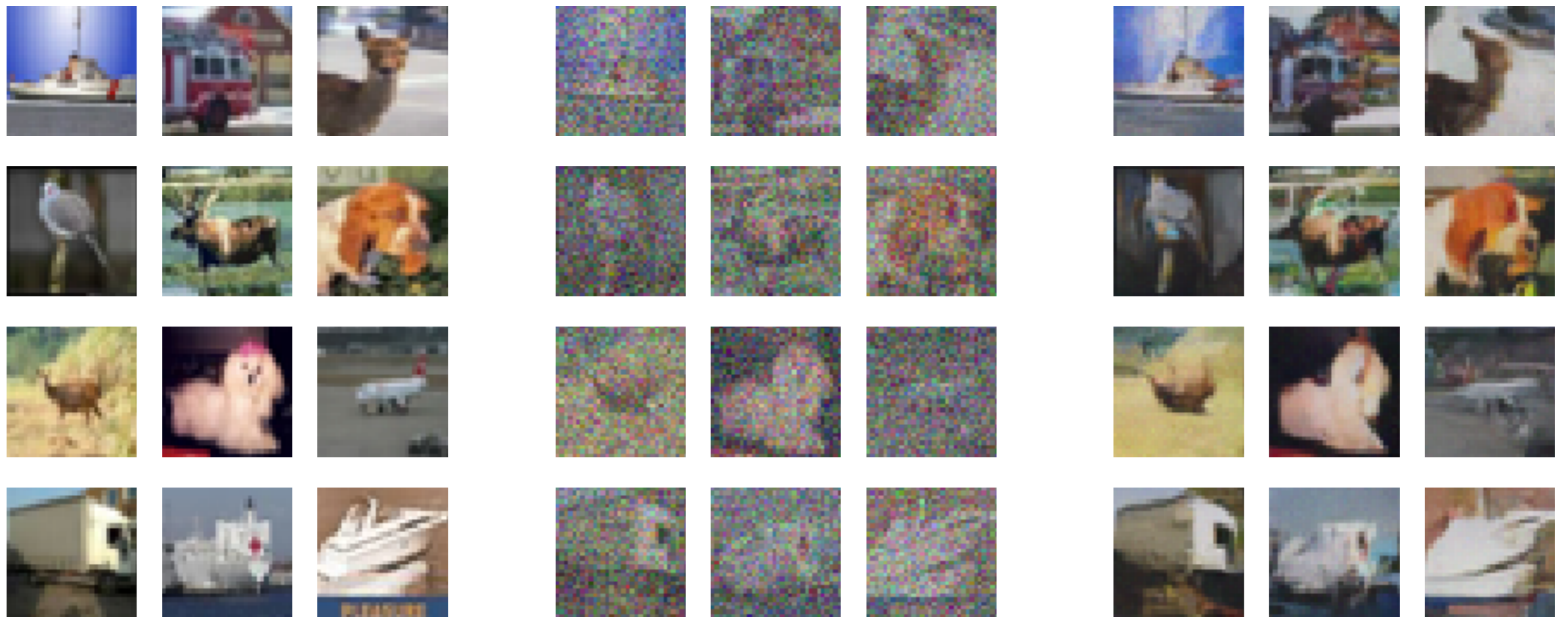
Image Denoising by Sampling from Posterior



Holdout Data

Corrupted
(SNR = 1)

Image Denoising by Sampling from Posterior



Holdout Data

Corrupted
(SNR = 1)

Denoised

Image Inpainting by Sampling from Posterior

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

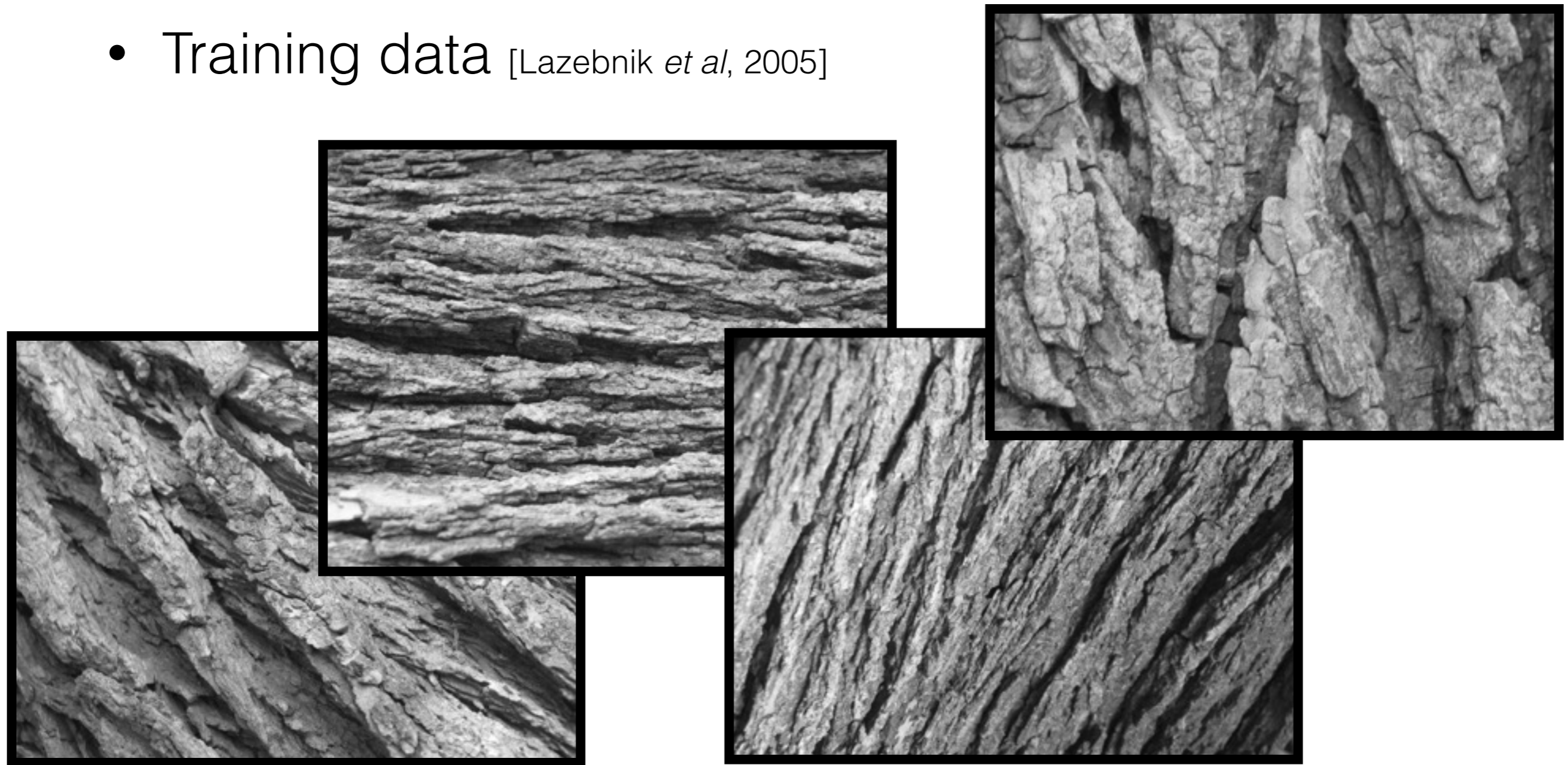
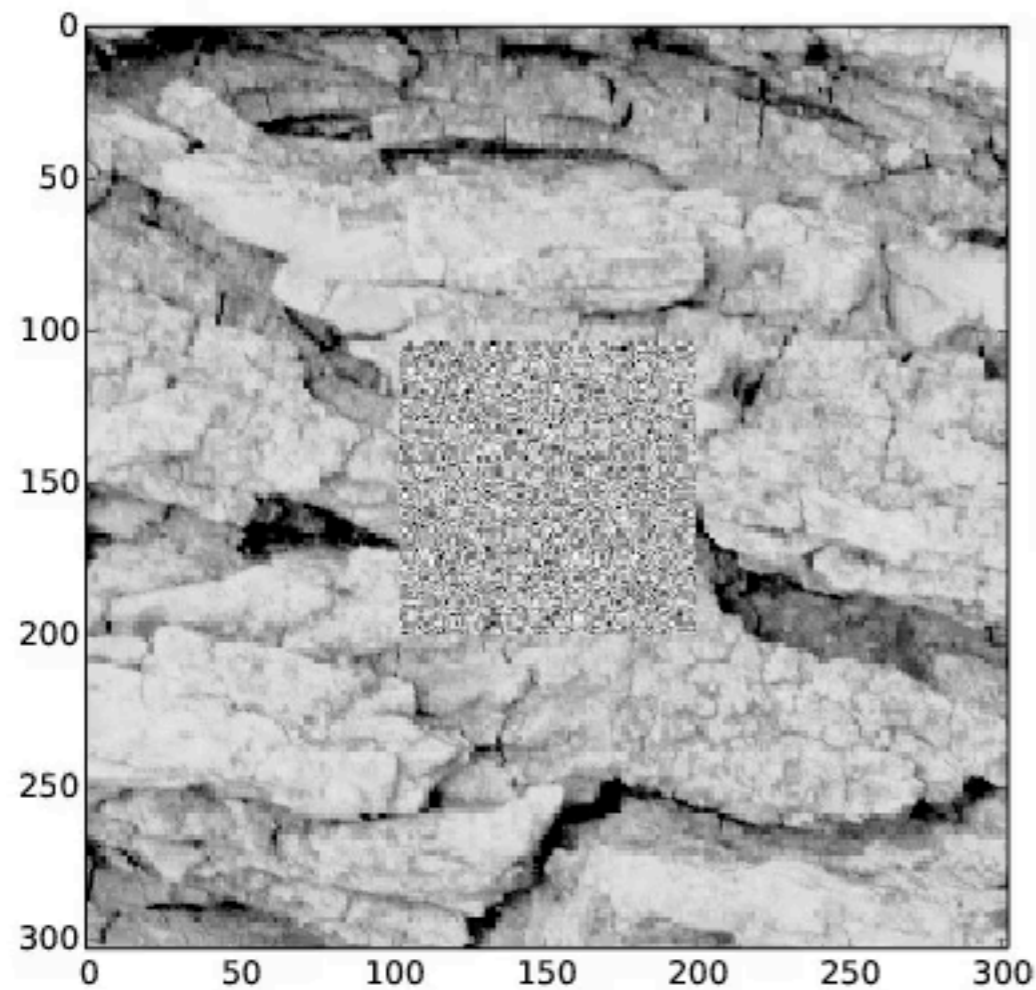
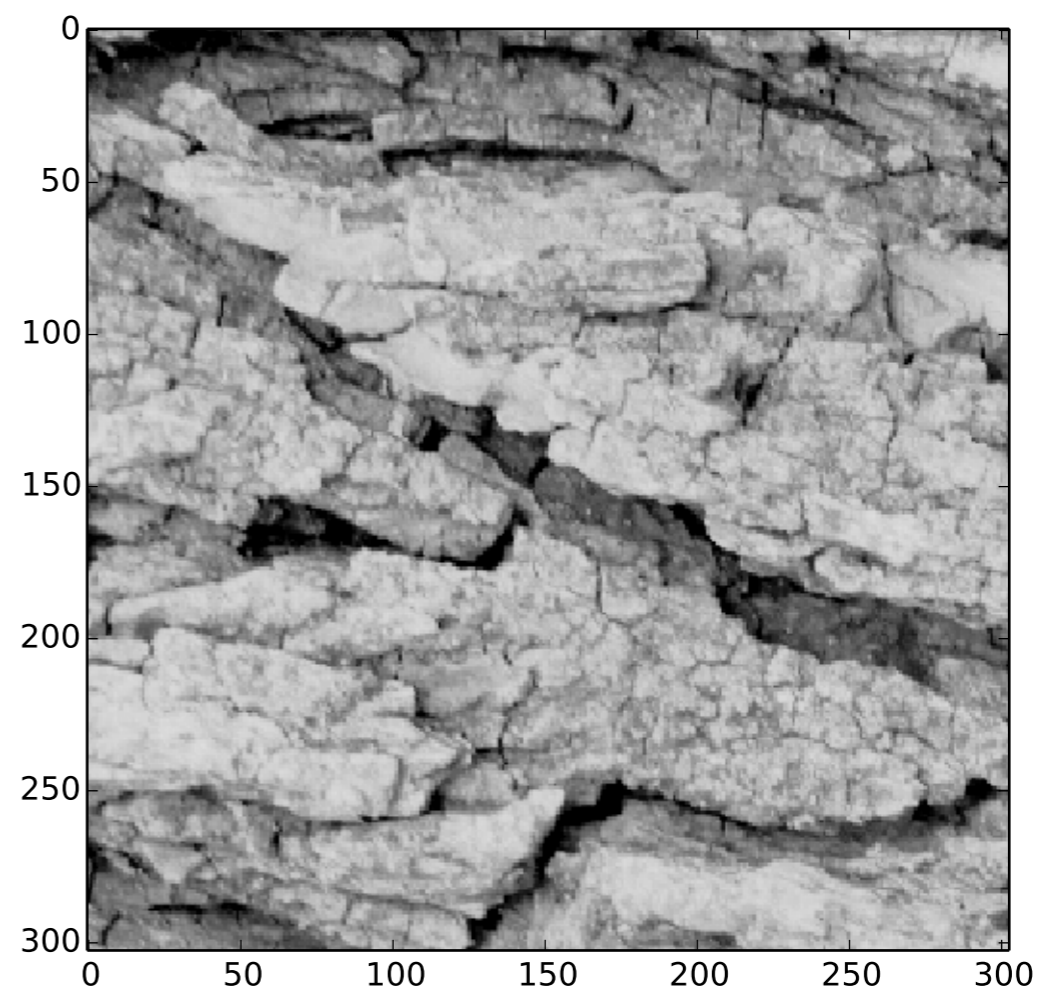


Image Inpainting by Sampling from Posterior

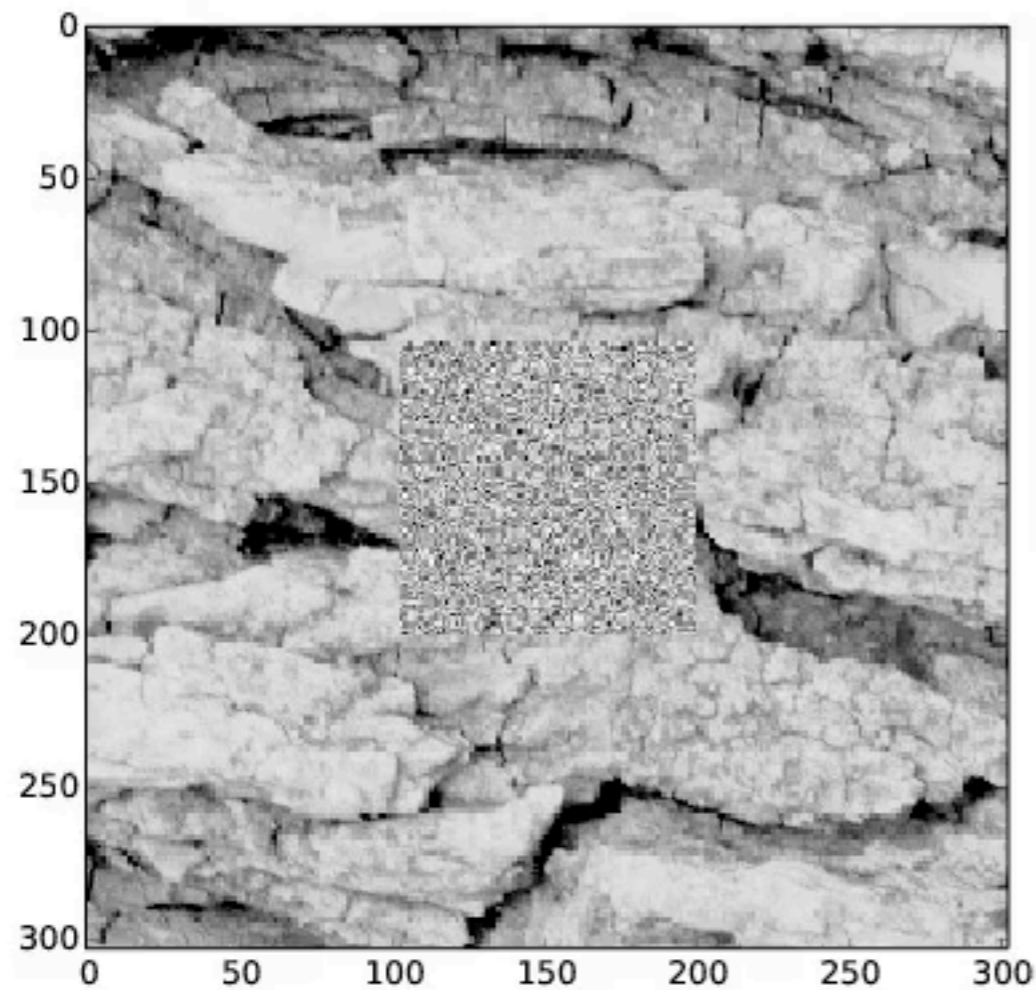


Inpainted image

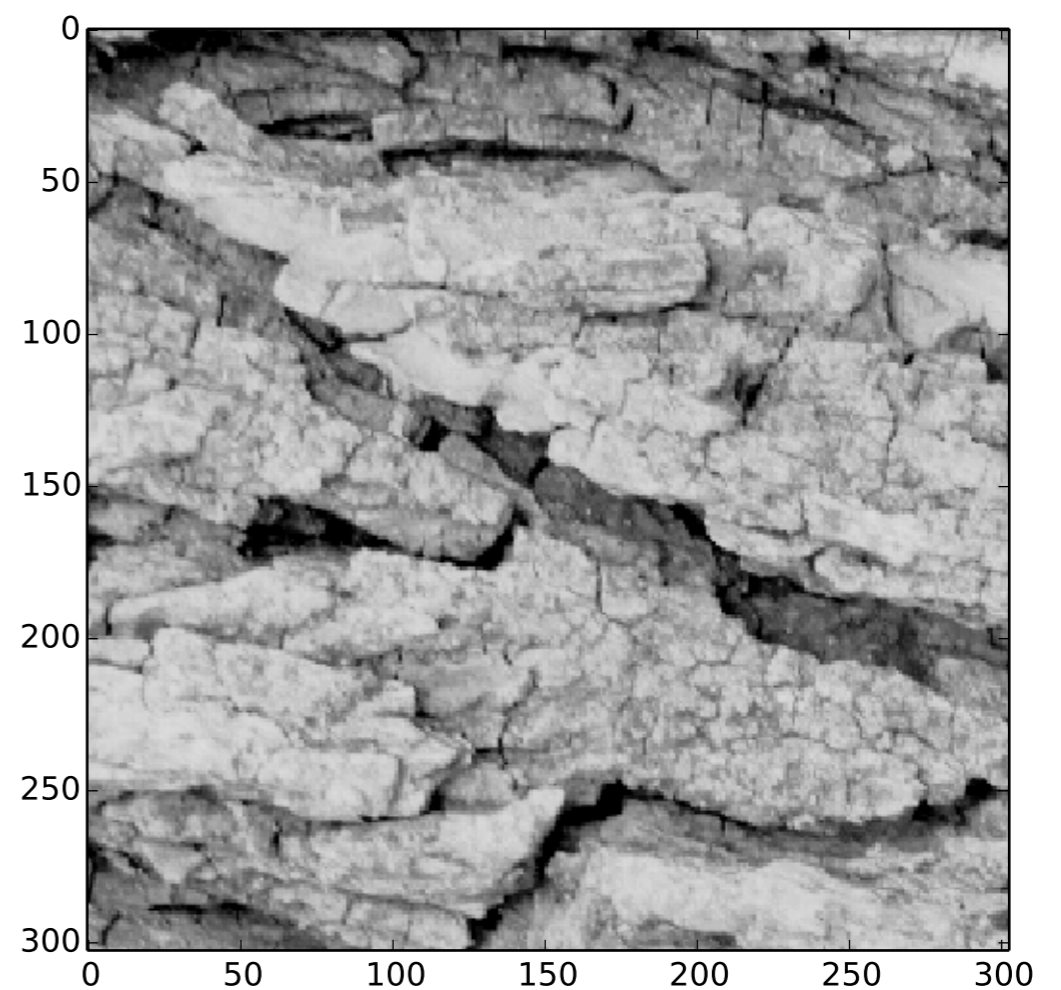


True image

Image Inpainting by Sampling from Posterior



Inpainted image



True image

Ongoing Work

Ongoing Work

- Binomial diffusion to neural spike trains

Ongoing Work

- Binomial diffusion to neural spike trains
- Full resolution color natural images

Ongoing Work

- Binomial diffusion to neural spike trains
- Full resolution color natural images
- Continuous time formulation

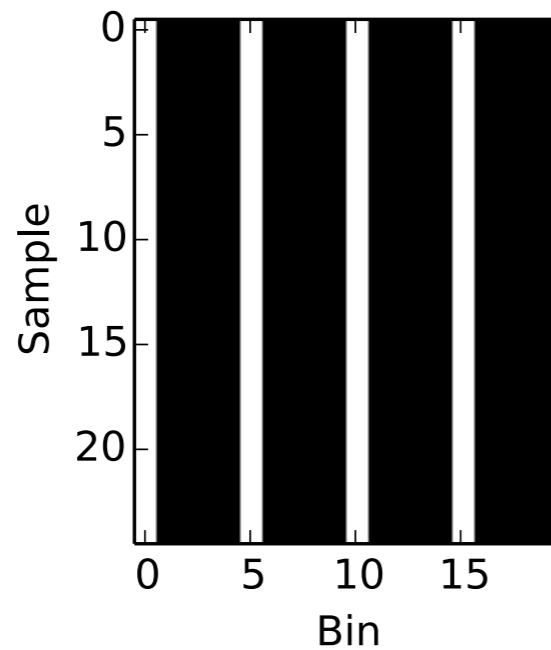
Ongoing Work

- Binomial diffusion to neural spike trains
- Full resolution color natural images
- Continuous time formulation
- Perturbation around energy based model

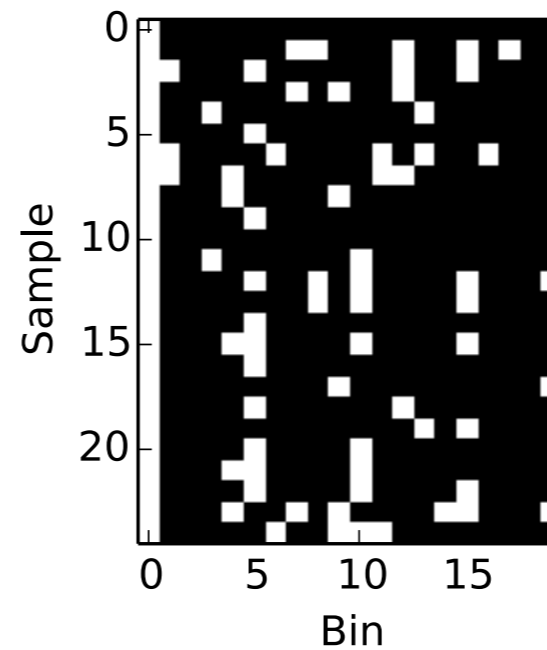
Toy Binary Sequence Learning

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

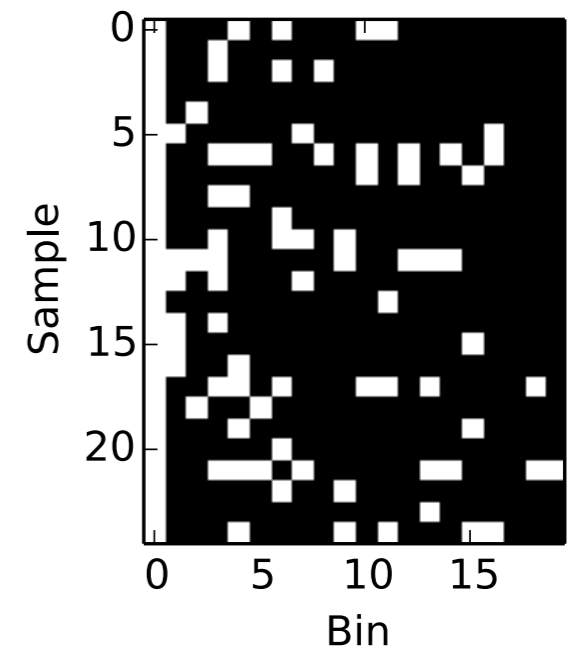
$t = 0$



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



$t = T$



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

Thanks!

Collaborators



Eric
Weiss



Niru
Maheswaranathan



Surya
Ganguli

Endless discussion

- The Ganguli Gang
- The Redwood Center for Theoretical Neuroscience

Setting Diffusion Rate

- Erase constant fraction of stimulus variance each step

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

- Can also train β_t