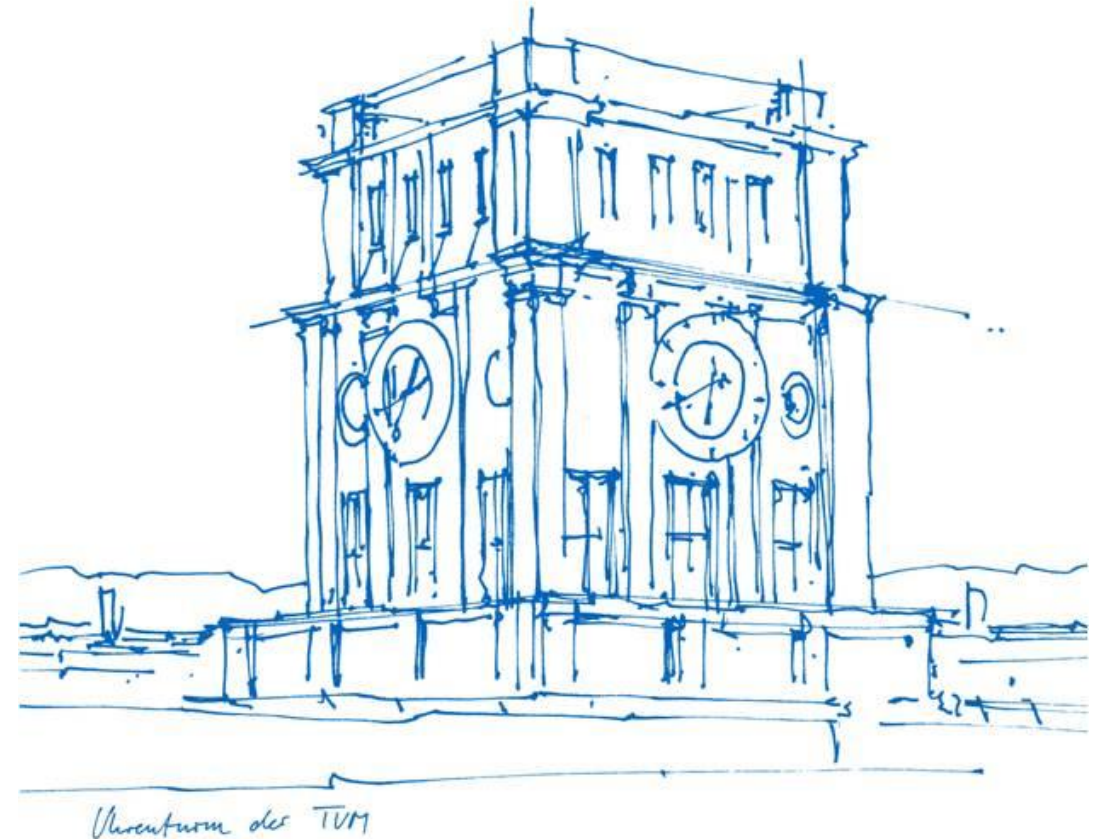


# Bayesian Flow Networks

Data Analytics and Machine Learning

Jed Guzelkabaagac

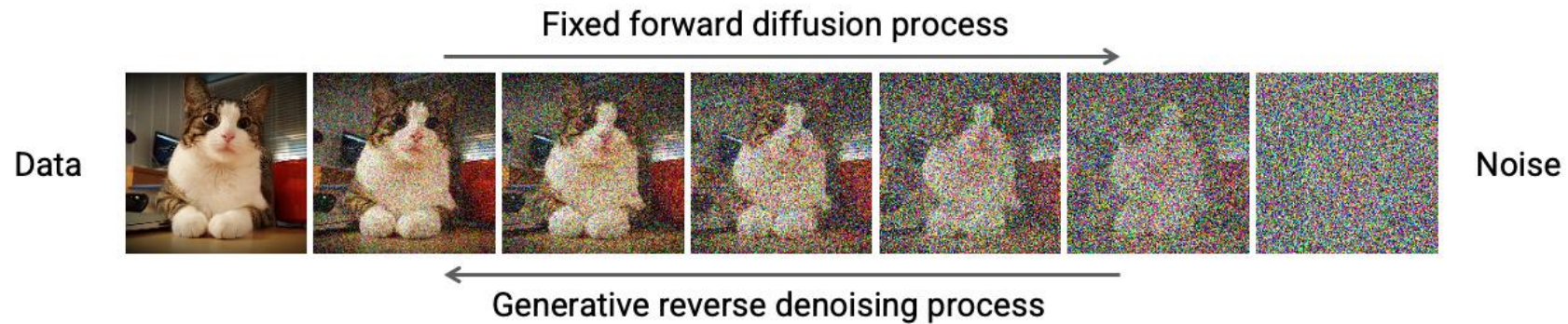
Supervised by Marcel Kolloviah



# Diffusion Models

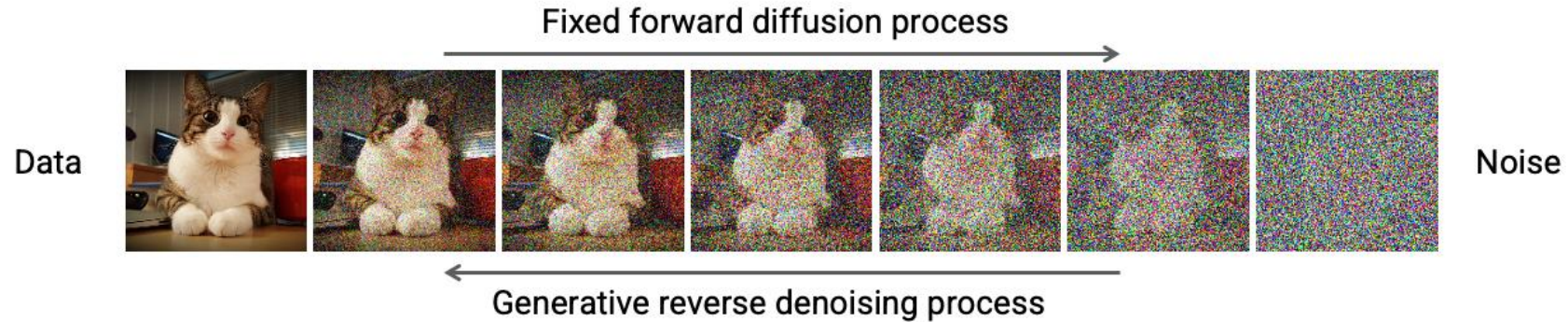
# Diffusion Models

Diffusion models typically operate directly on data samples  $x$



# Diffusion Models

Diffusion models typically operate directly on data samples  $x$



For discrete data, this leads to an **uninformative loss function**, impeding optimization.

$$\mathcal{L}_{\text{diffusion}} = \mathbb{E}_{x, \epsilon} [\|\epsilon - \epsilon_{\theta}(x_t, t)\|]$$

# Bayesian Flow Networks



**A unified differentiable generative framework**

# Bayesian Flow Networks



**A unified differentiable generative framework**

✓ Natively handles discrete data

# Bayesian Flow Networks



**A unified differentiable generative framework**

- ✓ Natively handles discrete data
- ✓ Great performance on 'hybrid data' tasks

# Bayesian Flow Networks

## **A unified differentiable generative framework**

- ✓ Natively handles discrete data
- ✓ Great performance on ‘hybrid data’ tasks
- ✓ No corruption process needed – start with a simple prior



# Bayesian Flow Networks

## **A unified differentiable generative framework**

- ✓ Natively handles discrete data
- ✓ Great performance on ‘hybrid data’ tasks
- ✓ No corruption process needed – start with a simple prior
- ✓ ‘Bayesian perspective’ on diffusion model

# Operating in Parameter Space



**Goal:** Operate on the parameter space to maintain continuity

# Operating in Parameter Space

**Goal:** Operate on the parameter space to maintain continuity

	Data Space	Parameter Space
Continuous	$\mathbf{x} \in \mathbb{R}^D$	$\mu \in \mathbb{R}^D$
Discrete	$\mathbf{x} \in \{0, 1\}^{KD}$	$\theta \in [0, 1]^{KD}$

# Operating in Parameter Space

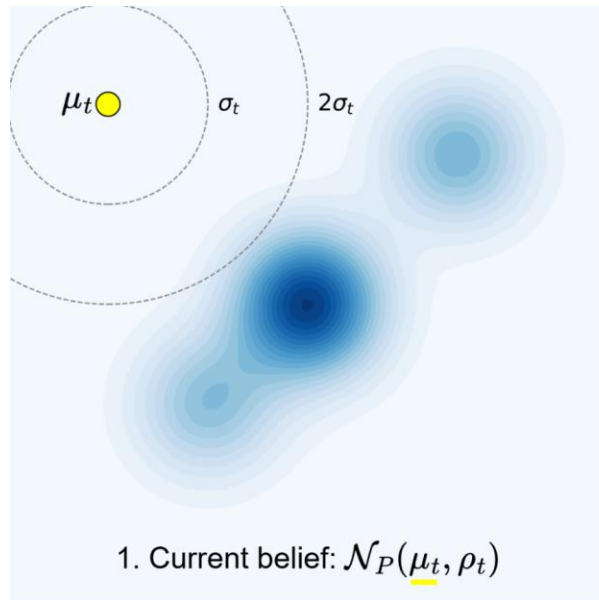
**Goal:** Operate on the parameter space to maintain continuity

	Data Space	Parameter Space
Continuous	$\mathbf{x} \in \mathbb{R}^D$	$\mu \in \mathbb{R}^D$
Discrete	$\mathbf{x} \in \{0, 1\}^{KD}$	$\theta \in [0, 1]^{KD}$

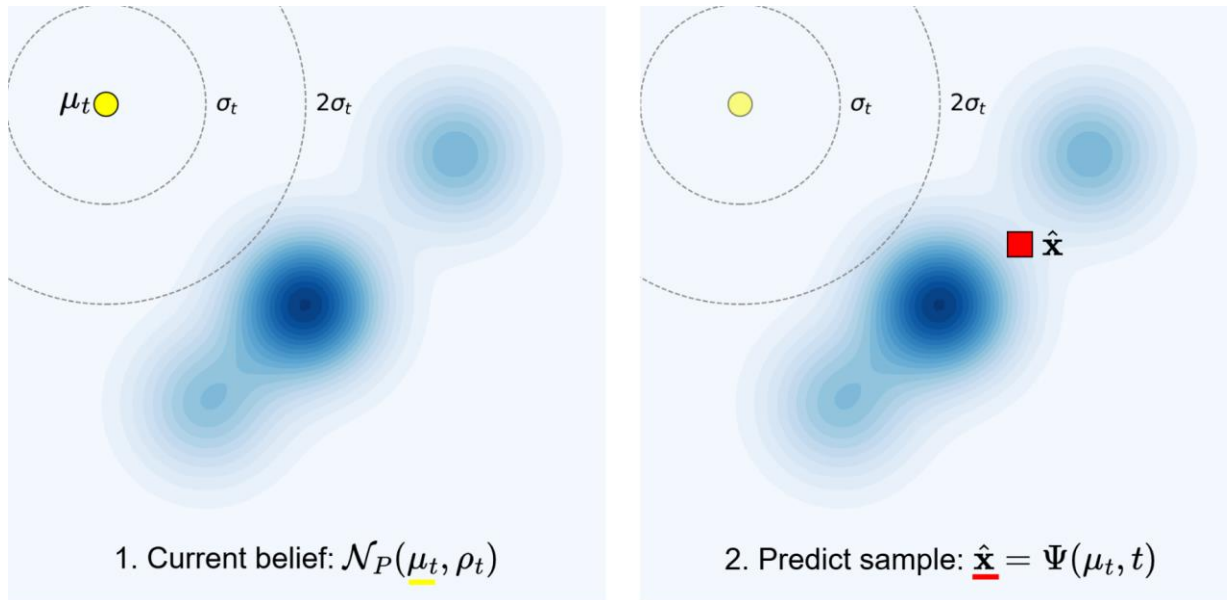
**Discrete data:** Relax each dimension's one-hot vector as a categorical distribution, which lives in a continuous simplex.

⇒ Maintain a **fully differentiable pipeline** even for discrete data

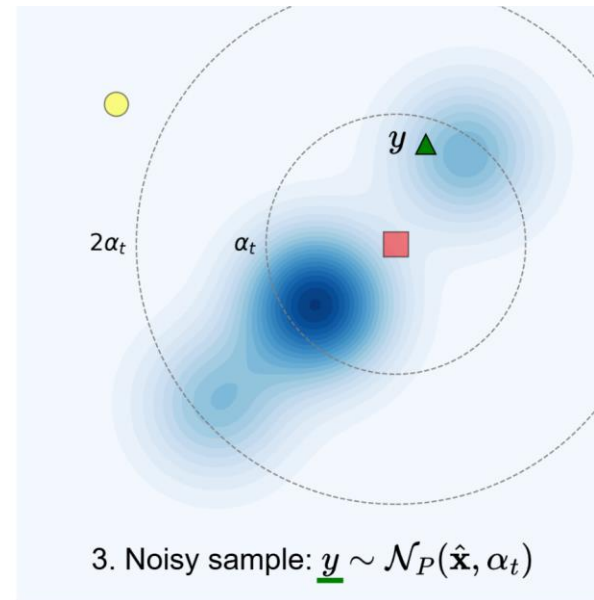
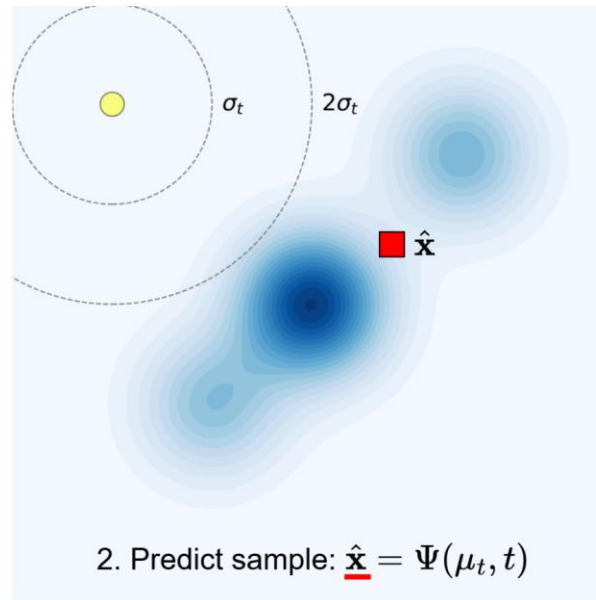
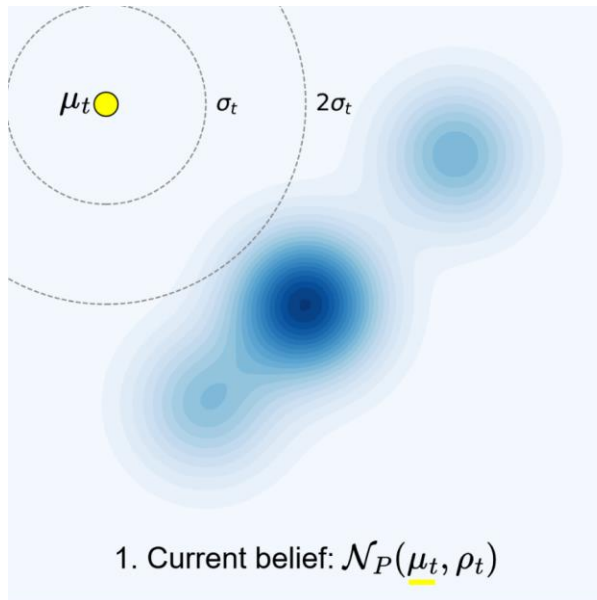
# Example: Continuous Generation



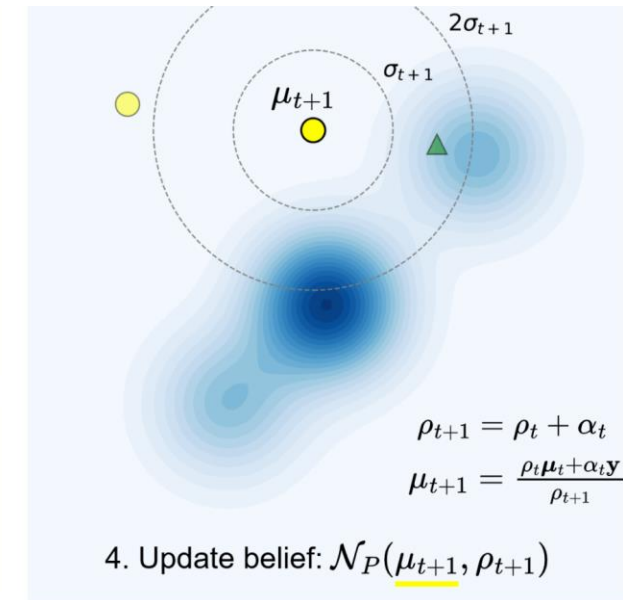
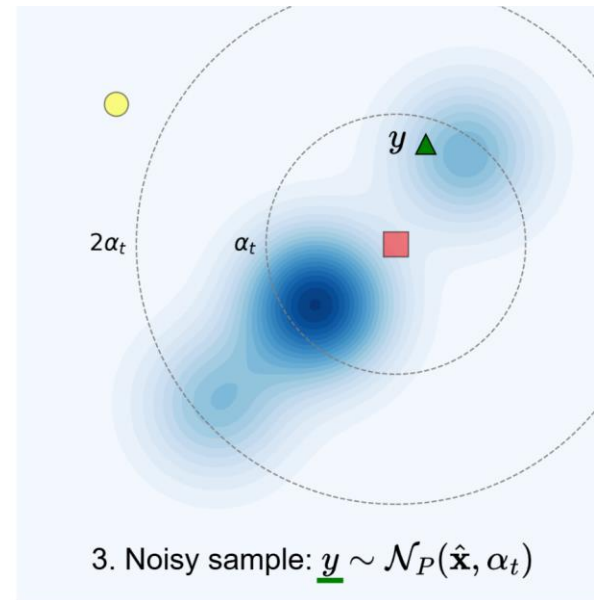
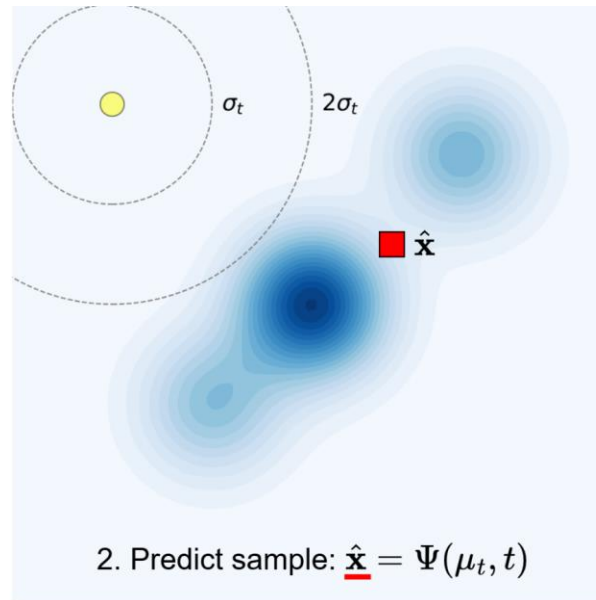
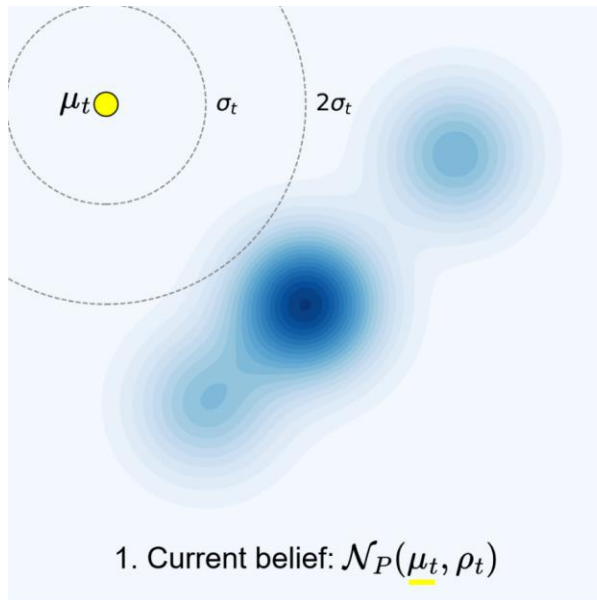
# Example: Continuous Generation



# Example: Continuous Generation

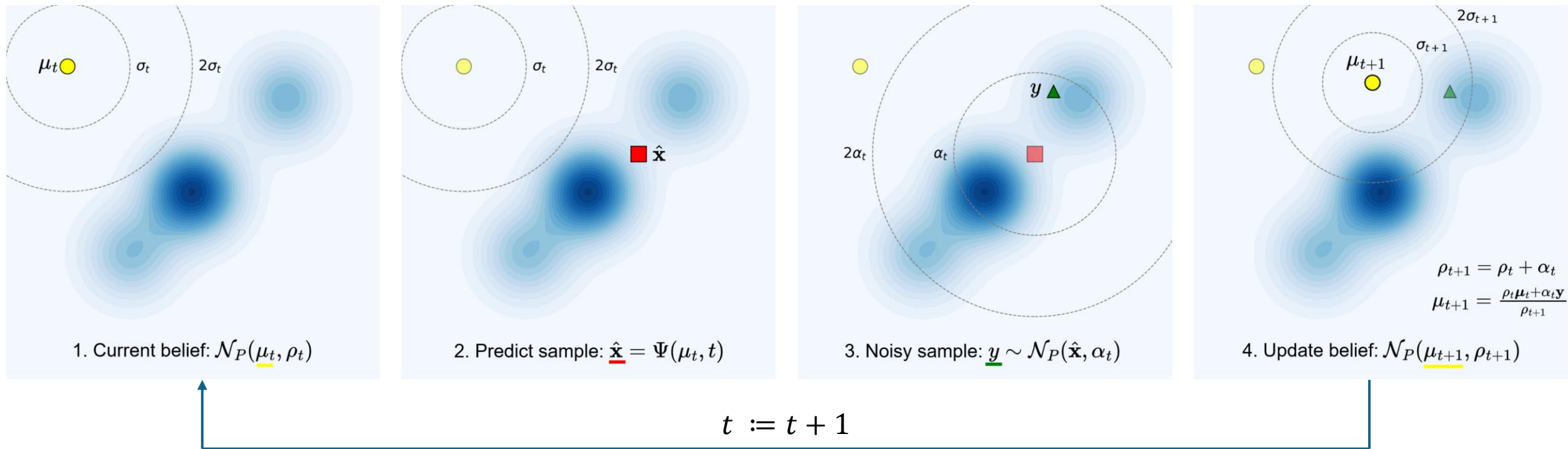


# Example: Continuous Generation

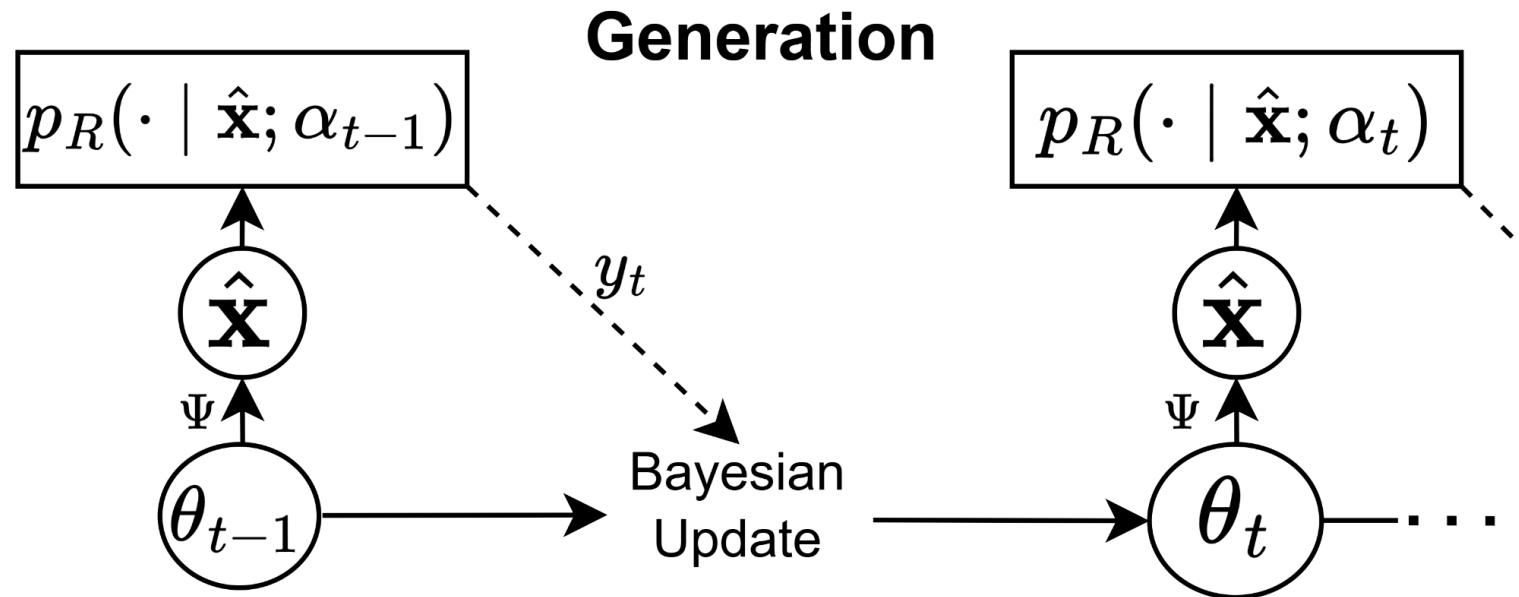




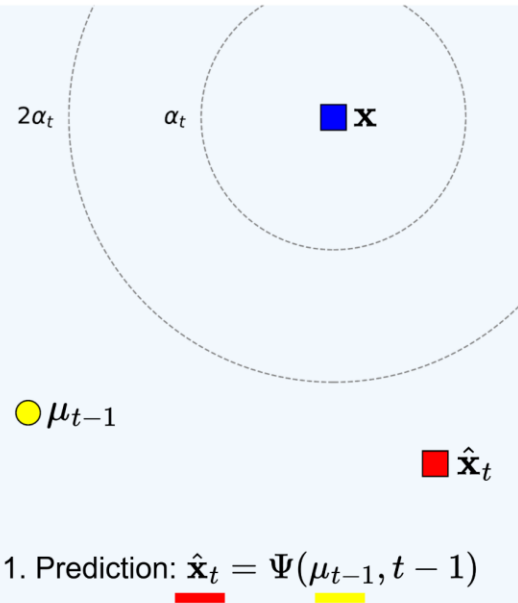
# Example: Continuous Generation



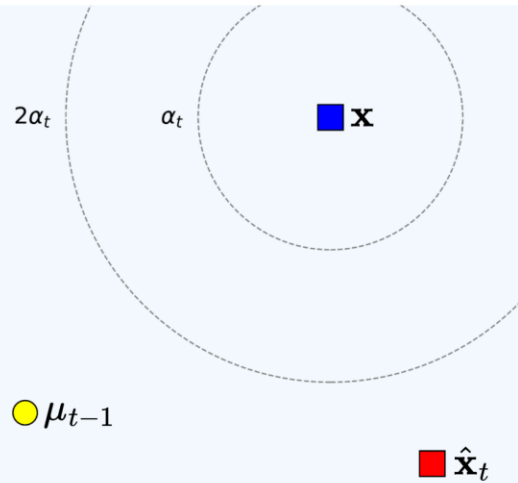
# Generative Process



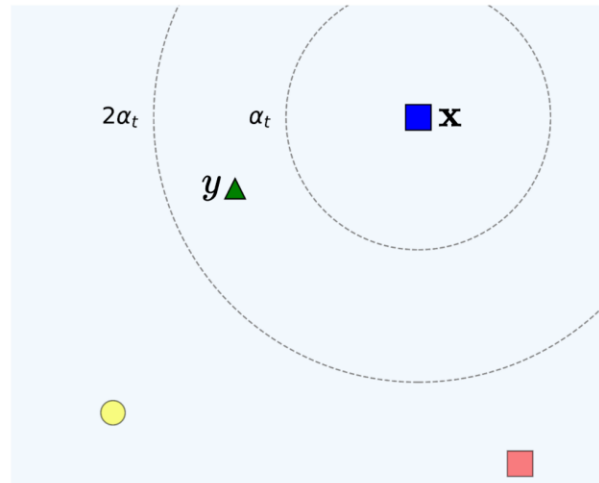
# Example: Continuous Training



# Example: Continuous Training

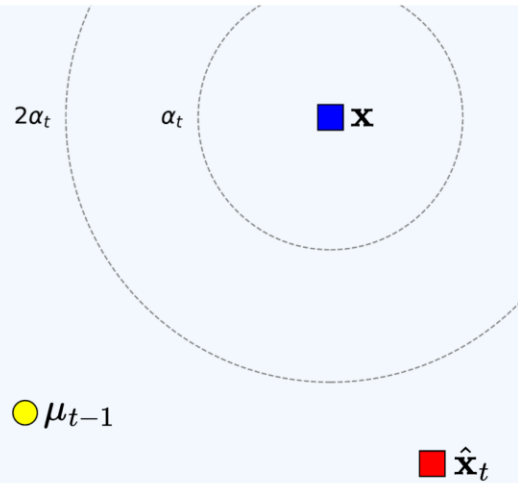


1. Prediction:  $\hat{\mathbf{x}}_t = \Psi(\mu_{t-1}, t-1)$

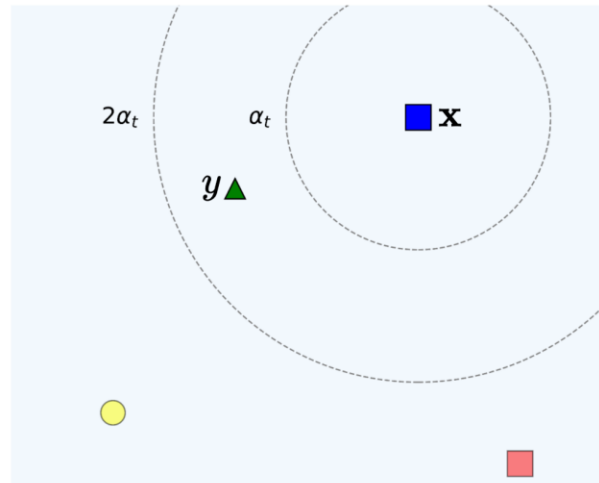


2. Noisy sample:  $y \sim \mathcal{N}_P(\mathbf{x}, \alpha_t)$

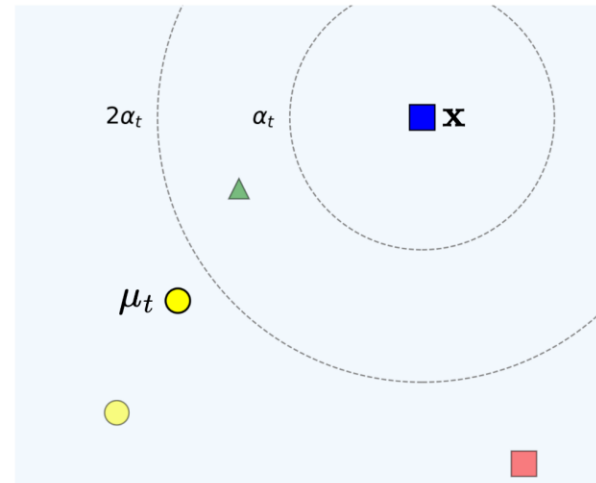
# Example: Continuous Training



1. Prediction:  $\hat{\mathbf{x}}_t = \Psi(\mu_{t-1}, t-1)$

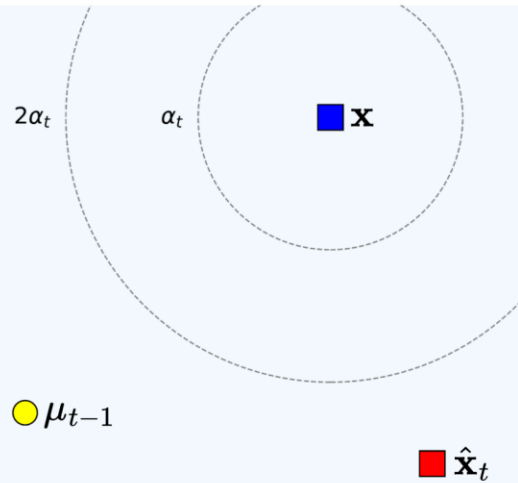


2. Noisy sample:  $y \sim \mathcal{N}_P(\mathbf{x}, \alpha_t)$

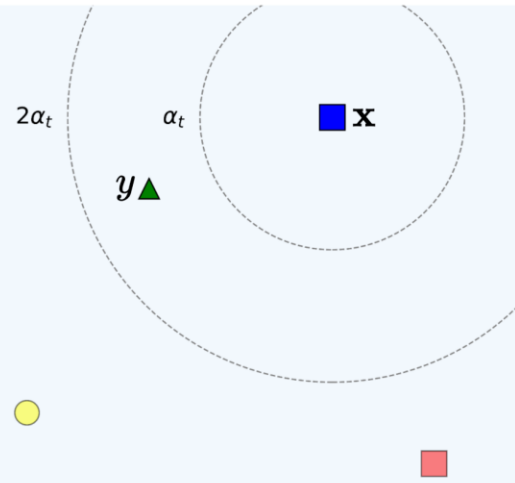


3. Bayesian update:  $\mu_{t-1} \rightarrow \mu_t$

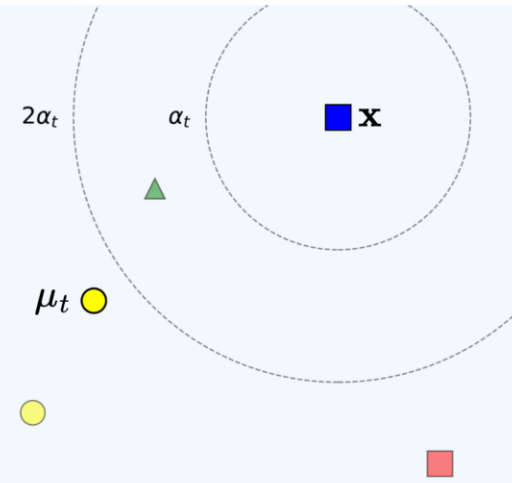
# Example: Continuous Training



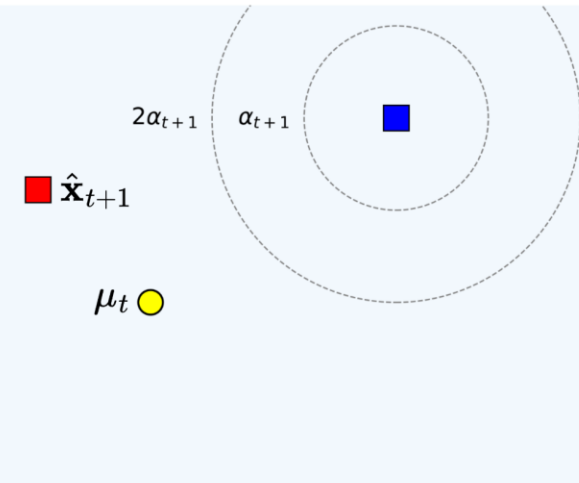
1. Prediction:  $\hat{\mathbf{x}}_t = \Psi(\mu_{t-1}, t-1)$



2. Noisy sample:  $y \sim \mathcal{N}_P(\mathbf{x}, \alpha_t)$

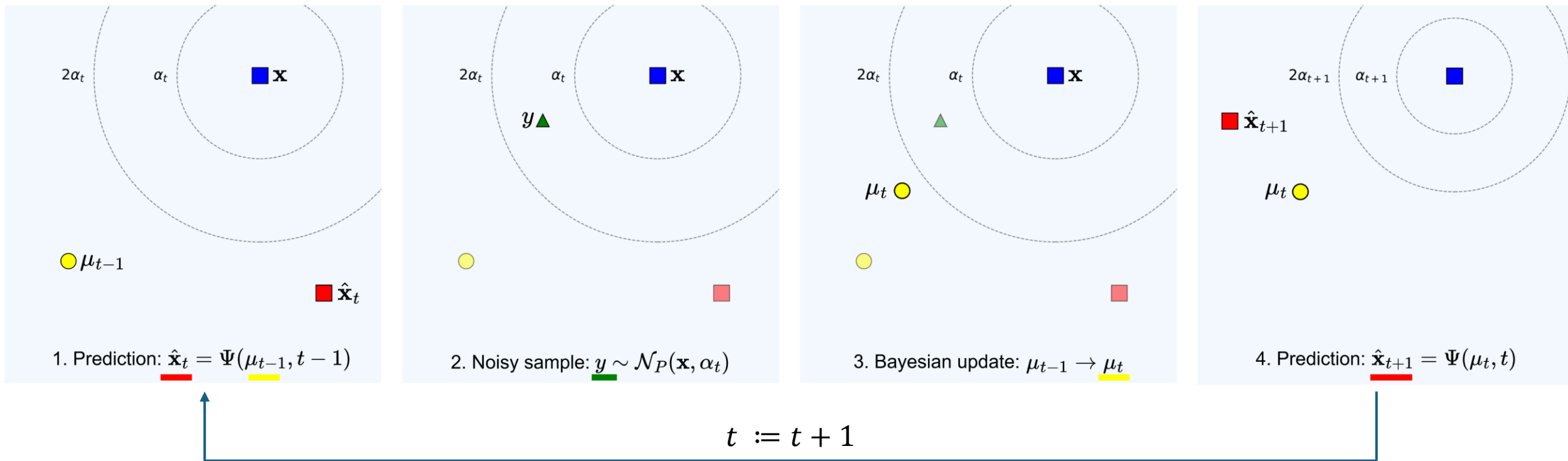


3. Bayesian update:  $\mu_{t-1} \rightarrow \mu_t$

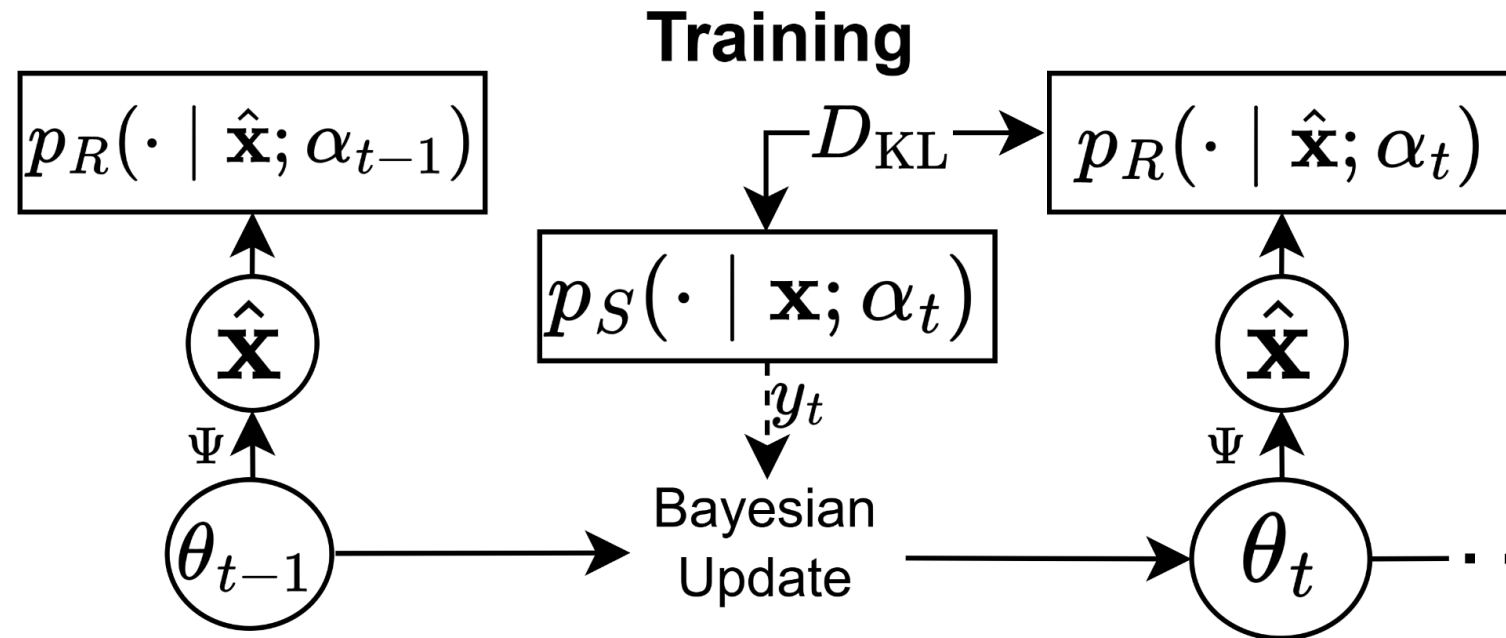


4. Prediction:  $\hat{\mathbf{x}}_{t+1} = \Psi(\mu_t, t)$

# Example: Continuous Training



# Training Process





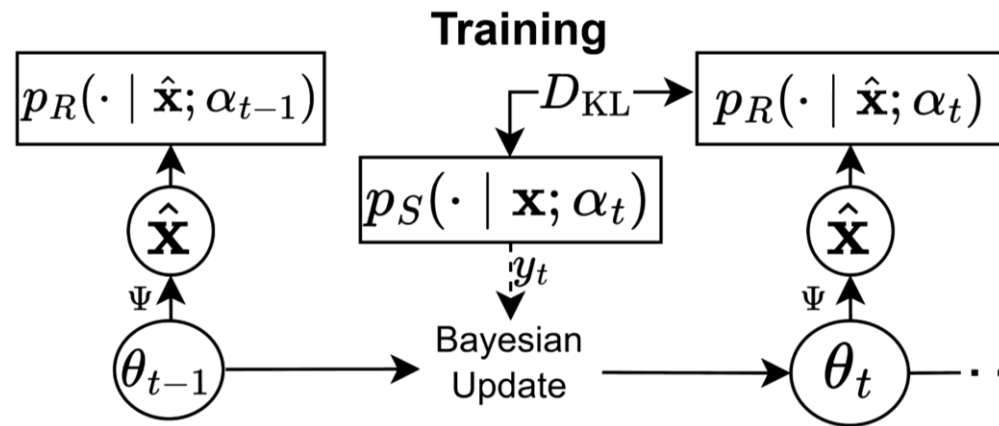
# Summary: Parameters $\theta$

## Purpose:

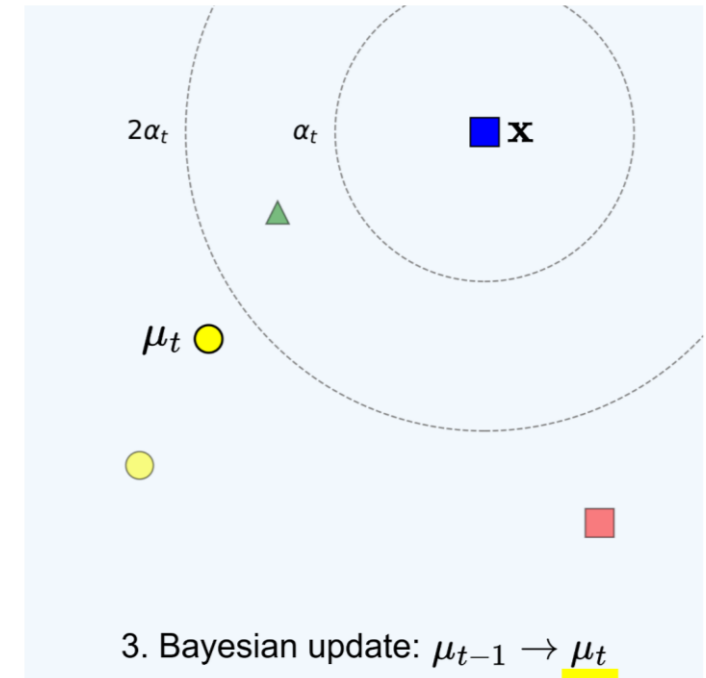
Parameters  $\theta$  are inputs to the neural network  $\Psi$

## Evolution:

During training, we perform Bayesian updates on  $\theta_t$ , using samples from the sender distribution (around  $\mathbf{x}$ )



$$\hat{\mathbf{x}}_t \triangleq \Psi(\theta_{t-1}, t-1)$$



# Sender and Receiver Distributions

# Sender and Receiver Distributions

Fix  $n$  steps of the network, choose an accuracy schedule  $\alpha_1 < \dots < \alpha_n$

$$p_S(\mathbf{y} \mid \mathbf{x}; \alpha) = \mathcal{N}(\mathbf{y} \mid \mathbf{x}, \alpha^{-1} \mathbf{I})$$

*Continuous case*

# Sender and Receiver Distributions

Fix  $n$  steps of the network, choose an accuracy schedule  $\alpha_1 < \dots < \alpha_n$

$$p_S(\mathbf{y} \mid \mathbf{x}; \alpha) = \mathcal{N}(\mathbf{y} \mid \mathbf{x}, \alpha^{-1} \mathbf{I})$$

*Continuous case*

- Depends on data type
- Parameterised by some point  $x$  and predetermined noise  $\alpha$
- Required to be factorizable (independent across dimensions)

# Sender and Receiver Distributions

Fix  $n$  steps of the network, choose an accuracy schedule  $\alpha_1 < \dots < \alpha_n$

$$p_S(\mathbf{y} \mid \mathbf{x}; \alpha) = \mathcal{N}(\mathbf{y} \mid \mathbf{x}, \alpha^{-1} \mathbf{I})$$

*Continuous case*

- Depends on data type
- Parameterised by some point  $x$  and predetermined noise  $\alpha$
- Required to be factorizable (independent across dimensions)



$$\sim p_S(\mathbf{y} \mid \text{cat}; \alpha_t)$$

# Sender and Receiver Distributions



Similarly, receiver distribution is parameterized by ‘current understanding’  $\hat{\mathbf{x}}$  and noise  $\alpha$

# Sender and Receiver Distributions



Similarly, receiver distribution is parameterized by ‘current understanding’  $\hat{\mathbf{x}}$  and noise  $\alpha$

$$p_R(\mathbf{y} \mid \hat{\mathbf{x}}; \alpha) \triangleq \mathbb{E}_{p_O(\mathbf{x}' \mid \hat{\mathbf{x}})} p_S(\mathbf{y} \mid \mathbf{x}'; \alpha)$$

The receiver **is a mixture of sender distributions** weighted by its belief about  $\mathbf{x}$

# Sender and Receiver Distributions

Similarly, receiver distribution is parameterized by ‘current understanding’  $\hat{\mathbf{x}}$  and noise  $\alpha$

$$p_R(\mathbf{y} \mid \hat{\mathbf{x}}; \alpha) \triangleq \mathbb{E}_{p_O(\mathbf{x}' \mid \hat{\mathbf{x}})} p_S(\mathbf{y} \mid \mathbf{x}'; \alpha)$$

The receiver **is a mixture of sender distributions** weighted by its belief about  $\mathbf{x}$

Now consider a loss function defined by the KL divergence

$$\text{KL}[p_S(\mathbf{y} \mid \mathbf{x}, \alpha) \parallel p_R(\mathbf{y} \mid \hat{\mathbf{x}}, \alpha)]$$



# Sender and Receiver Distributions

Similarly, receiver distribution is parameterized by ‘current understanding’  $\hat{\mathbf{x}}$  and noise  $\alpha$

$$p_R(\mathbf{y} \mid \hat{\mathbf{x}}; \alpha) \triangleq \mathbb{E}_{p_O(\mathbf{x}' \mid \hat{\mathbf{x}})} p_S(\mathbf{y} \mid \mathbf{x}'; \alpha)$$

The receiver **is a mixture of sender distributions** weighted by its belief about  $\mathbf{x}$

Now consider a loss function defined by the KL divergence

$$\text{KL}[p_S(\mathbf{y} \mid \mathbf{x}, \alpha) \parallel p_R(\mathbf{y} \mid \hat{\mathbf{x}}, \alpha)]$$

- ✓ Minimizing the loss pushes  $\hat{\mathbf{x}}$  towards  $\mathbf{x}$
- ✓ KL gives a ‘smooth training signal’ with respect to parameters of NN  $\Psi$

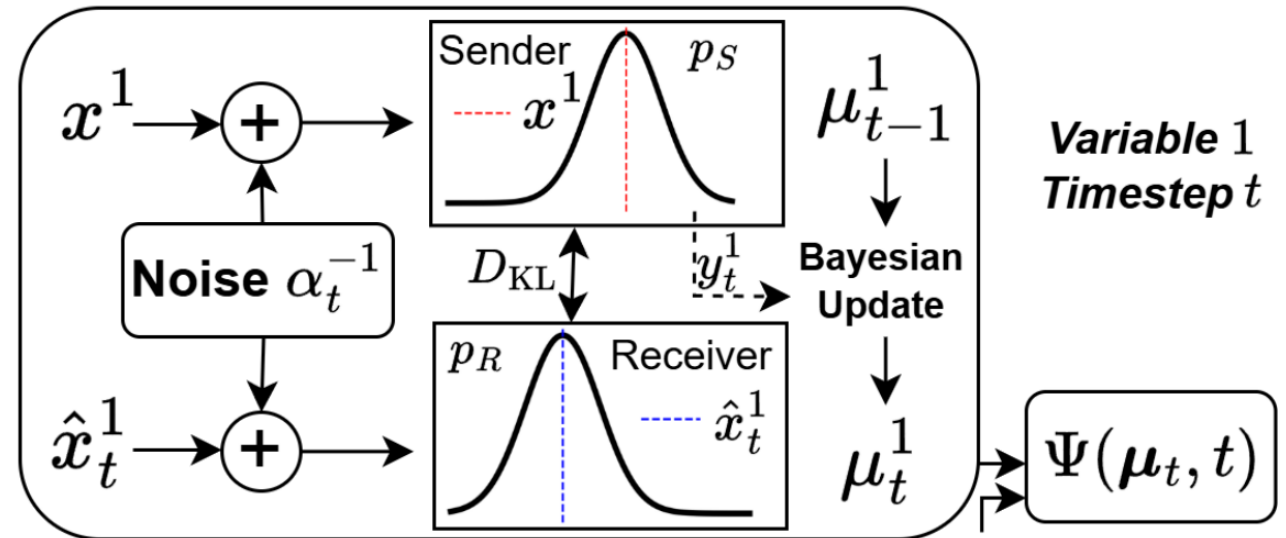
# Internal Distributions: Continuous

*Internal Distributions of BFN*

$$p_S(\cdot \mid \mathbf{x}; \alpha_t) = \mathcal{N}(\mathbf{x}, \alpha_t^{-1} \mathbf{I}).$$

$$p_O(\mathbf{x}' \mid \hat{\mathbf{x}}_t) = \delta(\mathbf{x}' - \hat{\mathbf{x}}_t).$$

$$\begin{aligned} p_R(\cdot \mid \hat{\mathbf{x}}_t; \alpha_t) &= \mathbb{E}_{p_O(\mathbf{x}' \mid \hat{\mathbf{x}}_t)} p_S(\cdot \mid \mathbf{x}'; \alpha_t) \\ &= \mathcal{N}(\hat{\mathbf{x}}_t, \alpha_t^{-1} \mathbf{I}) \end{aligned}$$



# Loss Function

Equivalent to a VAE loss where we have a **sequence** of latent samples  $\mathbf{y}_1, \dots, \mathbf{y}_n$

$$L^n(\mathbf{x}) = \mathbb{E}_{p(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{n-1})} \sum_{i=1}^n D_{\text{KL}}(p_S(\cdot \mid \mathbf{x}; \alpha_i) \parallel p_R(\cdot \mid \hat{\mathbf{x}}_i; \alpha_i))$$

# Loss Function

Equivalent to a VAE loss where we have a **sequence** of latent samples  $\mathbf{y}_1, \dots, \mathbf{y}_n$

$$L^n(\mathbf{x}) = \mathbb{E}_{p(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{n-1})} \sum_{i=1}^n D_{\text{KL}}(p_S(\cdot \mid \mathbf{x}; \alpha_i) \parallel p_R(\cdot \mid \hat{\mathbf{x}}_i; \alpha_i))$$

## Transmission interpretation:

*What is the cost of transmitting sample  $\mathbf{y}_i \sim p_S(\cdot \mid \mathbf{x}; \alpha_i)$  from the sender to a receiver that believes  $\mathbf{y}_i \sim p_R(\cdot \mid \hat{\mathbf{x}}_i; \alpha_i)$*

# Loss Function

Equivalent to a VAE loss where we have a **sequence** of latent samples  $\mathbf{y}_1, \dots, \mathbf{y}_n$

$$L^n(\mathbf{x}) = \mathbb{E}_{p(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{n-1})} \sum_{i=1}^n D_{\text{KL}}(p_S(\cdot \mid \mathbf{x}; \alpha_i) \parallel p_R(\cdot \mid \hat{\mathbf{x}}_i; \alpha_i))$$

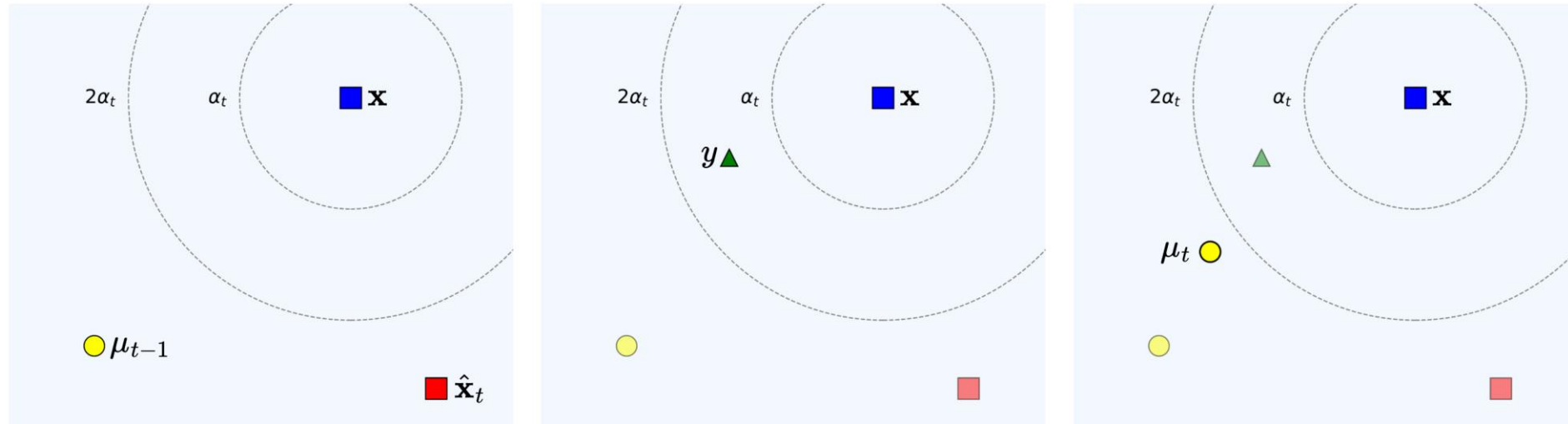
## Transmission interpretation:

*What is the cost of transmitting sample  $\mathbf{y}_i \sim p_S(\cdot \mid \mathbf{x}; \alpha_i)$  from the sender to a receiver that believes  $\mathbf{y}_i \sim p_R(\cdot \mid \hat{\mathbf{x}}_i; \alpha_i)$*

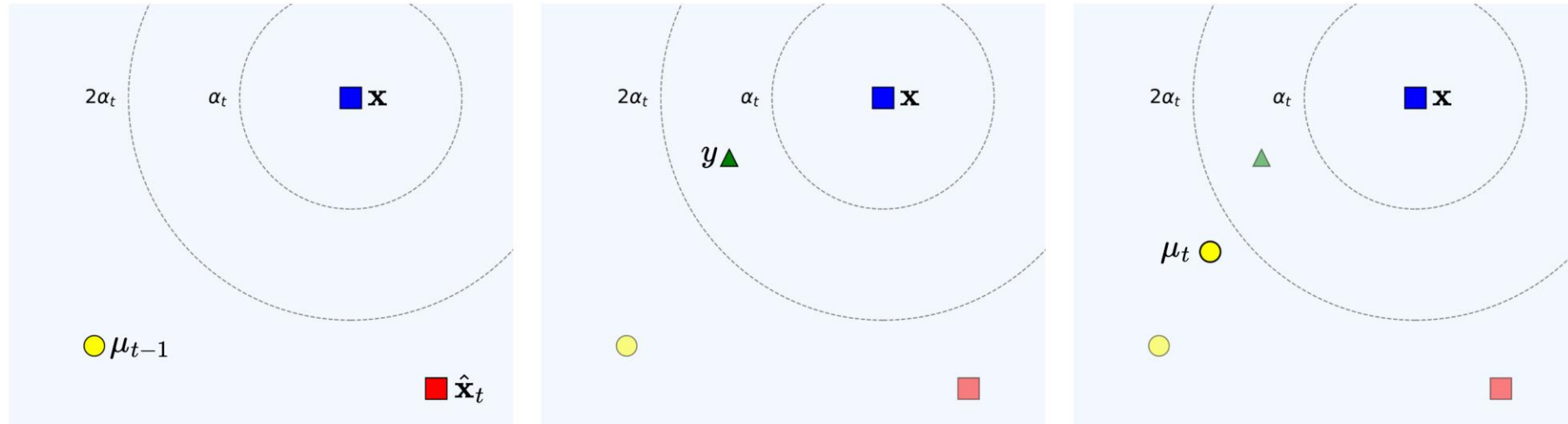
**Sampling Issue:**  $\boldsymbol{\theta}_1 \rightarrow \boldsymbol{\theta}_2 \rightarrow \dots \rightarrow \boldsymbol{\theta}_{n-1}$  from  $p(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{n-1}) = \prod_{i=1}^{n-1} p_U(\boldsymbol{\theta}_i \mid \boldsymbol{\theta}_{i-1}, \mathbf{x}; \alpha_i)$

✗ Computationally expensive ✗ Not parallelizable

# Update Distribution

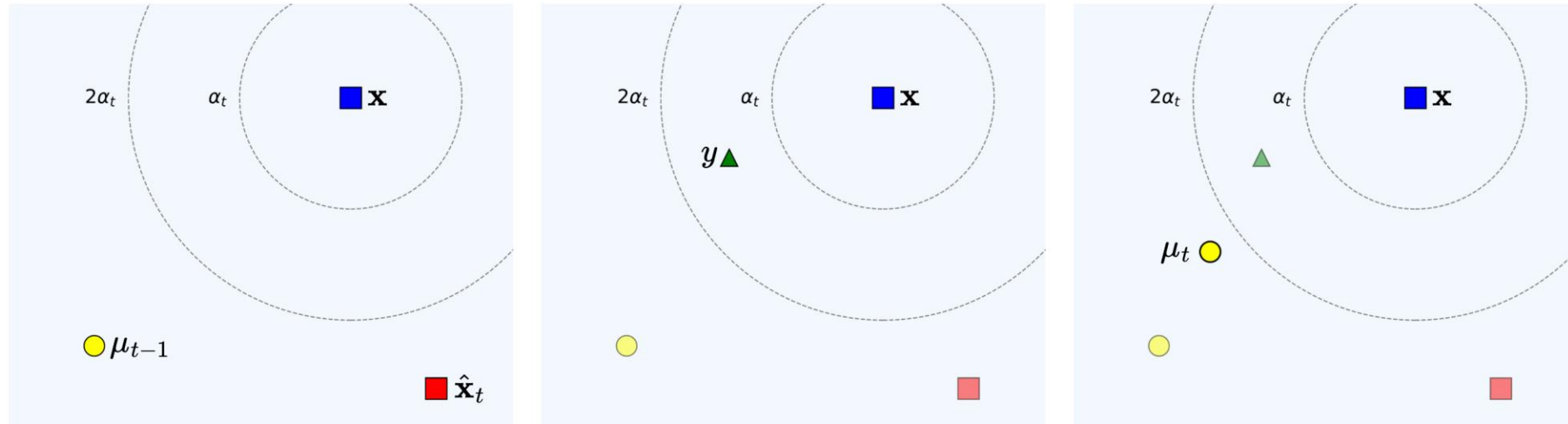


# Update Distribution



$$\mu_t = h(\mu_{t-1}, \mathbf{y}_t, \alpha_t) = \frac{\alpha_t \mathbf{y}_t + \rho_{t-1} \mu_{t-1}}{\rho_t}$$

# Update Distribution



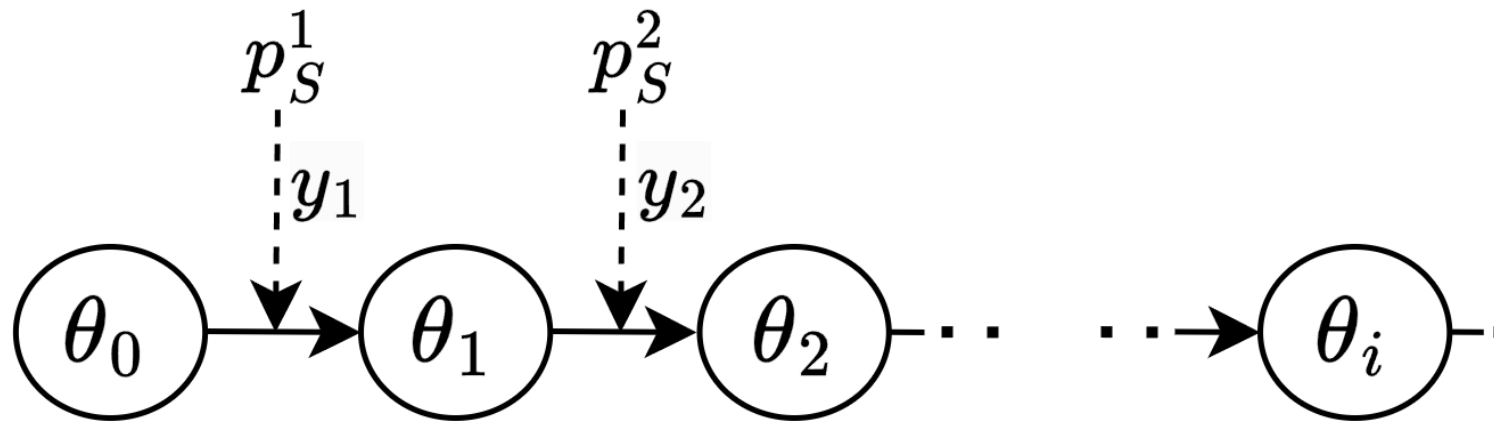
$$\mu_t = h(\mu_{t-1}, \mathbf{y}_t, \alpha_t) = \frac{\alpha_t \mathbf{y}_t + \rho_{t-1} \mu_{t-1}}{\rho_t}$$

$$p_U(\mu_t \mid \mu_{t-1}, \mathbf{x}; \alpha_t) = \mathcal{N}\left(\mu_t \mid \frac{\alpha_t \mathbf{x} + \mu_{t-1} \rho_{t-1}}{\rho_t}, \frac{\alpha_t}{\rho_t^2} \mathbf{I}\right)$$



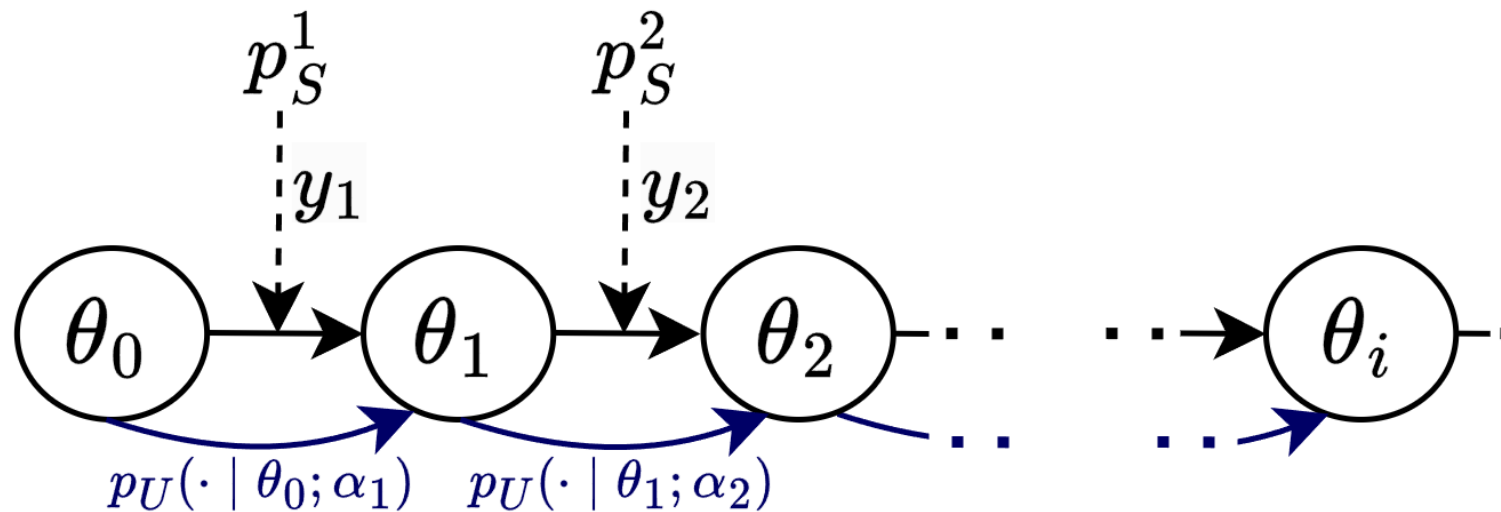
# Update Distribution

**Goal:** Sample  $\theta_i$  directly (without sampling  $\theta_1, \dots, \theta_{i-1}$ ) to simplify  $p(\theta_1, \dots, \theta_{n-1})$



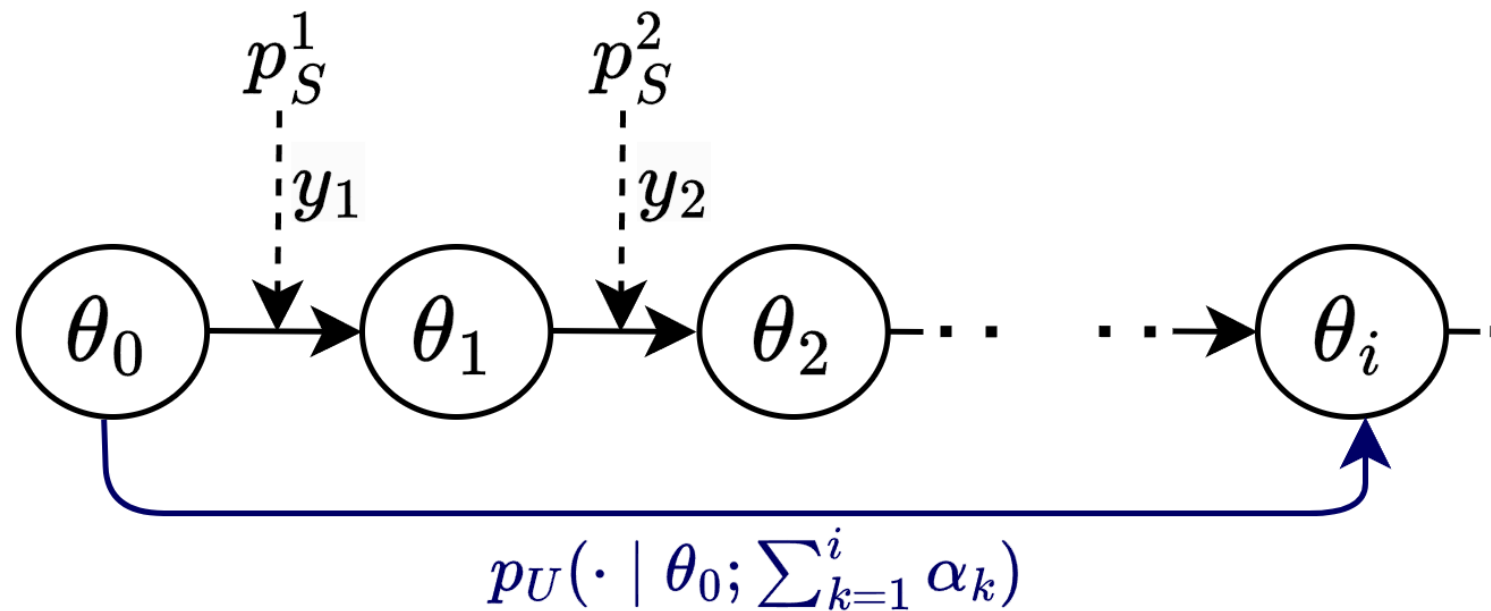
# Update Distribution

**Goal:** Sample  $\theta_i$  directly (without sampling  $\theta_1, \dots, \theta_{i-1}$ ) to simplify  $p(\theta_1, \dots, \theta_{n-1})$



# Update Distribution

**Goal:** Sample  $\theta_i$  directly (without sampling  $\theta_1, \dots, \theta_{i-1}$ ) to simplify  $p(\theta_1, \dots, \theta_{n-1})$



**Additive Accuracies:**  $\mathbb{E}_{p_U(\theta' \mid \theta, \mathbf{x}; \alpha_a)} p_U(\theta'' \mid \theta', \mathbf{x}; \alpha_b) \stackrel{!}{=} p_U(\theta'' \mid \theta, \mathbf{x}; \alpha_a + \alpha_b)$

# Loss Function

$$L^n(\mathbf{x}) = \mathbb{E}_{p(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{n-1})} \sum_{i=1}^n D_{\text{KL}} \left( p_S(\cdot \mid \mathbf{x}; \alpha_i) \parallel p_R(\cdot \mid \hat{\mathbf{x}}_i; \alpha_i) \right)$$

# Loss Function

$$\begin{aligned} L^n(\mathbf{x}) &= \mathbb{E}_{p(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{n-1})} \sum_{i=1}^n D_{\text{KL}} \left( p_S(\cdot \mid \mathbf{x}; \alpha_i) \parallel p_R(\cdot \mid \hat{\mathbf{x}}_i; \alpha_i) \right) \\ &= n \mathbb{E}_{i \sim U\{1, n\}} \underbrace{\mathbb{E}_{p(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{n-1})} D_{\text{KL}} \left( p_S(\cdot \mid \mathbf{x}; \alpha_i) \parallel p_R(\cdot \mid \hat{\mathbf{x}}_i; \alpha_i) \right)}_{\text{Depends on } \boldsymbol{\theta}_{i-1} \text{ only}} \end{aligned}$$

# Loss Function

$$\begin{aligned} L^n(\mathbf{x}) &= \mathbb{E}_{p(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{n-1})} \sum_{i=1}^n D_{\text{KL}} \left( p_S(\cdot \mid \mathbf{x}; \alpha_i) \parallel p_R(\cdot \mid \hat{\mathbf{x}}_i; \alpha_i) \right) \\ &= n \mathbb{E}_{i \sim U\{1, n\}} \underbrace{\mathbb{E}_{p(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{n-1})} D_{\text{KL}} \left( p_S(\cdot \mid \mathbf{x}; \alpha_i) \parallel p_R(\cdot \mid \hat{\mathbf{x}}_i; \alpha_i) \right)}_{\text{Depends on } \boldsymbol{\theta}_{i-1} \text{ only}} \\ &= n \mathbb{E}_{i \sim U\{1, \dots, n\}} \mathbb{E}_{p_U(\boldsymbol{\theta}_{i-1} \mid \boldsymbol{\theta}_0, \mathbf{x}; \sum_{j=1}^{i-1} \alpha_j)} D_{\text{KL}} \left( p_S(\cdot \mid \mathbf{x}; \alpha_i) \parallel p_R(\cdot \mid \hat{\mathbf{x}}_i; \alpha_i) \right) \end{aligned}$$

# Loss Function

$$\begin{aligned} L^n(\mathbf{x}) &= \mathbb{E}_{p(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{n-1})} \sum_{i=1}^n D_{\text{KL}} \left( p_S(\cdot \mid \mathbf{x}; \alpha_i) \parallel p_R(\cdot \mid \hat{\mathbf{x}}_i; \alpha_i) \right) \\ &= n \mathbb{E}_{i \sim U\{1, n\}} \underbrace{\mathbb{E}_{p(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{n-1})} D_{\text{KL}} \left( p_S(\cdot \mid \mathbf{x}; \alpha_i) \parallel p_R(\cdot \mid \hat{\mathbf{x}}_i; \alpha_i) \right)}_{\text{Depends on } \boldsymbol{\theta}_{i-1} \text{ only}} \\ &= n \mathbb{E}_{i \sim U\{1, \dots, n\}} \mathbb{E}_{p_U(\boldsymbol{\theta}_{i-1} \mid \boldsymbol{\theta}_0, \mathbf{x}; \sum_{j=1}^{i-1} \alpha_j)} D_{\text{KL}} \left( p_S(\cdot \mid \mathbf{x}; \alpha_i) \parallel p_R(\cdot \mid \hat{\mathbf{x}}_i; \alpha_i) \right) \end{aligned}$$

- ✓ Compute Monte Carlo without the  $n$ -step evolution of parameters
- ✓ Natural extension to continuous time loss

# Loss Function



$$L^n(\mathbf{x}) = n \mathbb{E}_{i \sim U\{1, \dots, n\}} \mathbb{E}_{p_U(\boldsymbol{\theta}_{i-1} | \boldsymbol{\theta}_0, \mathbf{x}; \sum_{j=1}^{i-1} \alpha_j)} D_{\text{KL}} \left( p_S(\cdot | \mathbf{x}; \alpha_i) \parallel p_R(\cdot | \hat{\mathbf{x}}_i; \alpha_i) \right)$$



# Loss Function

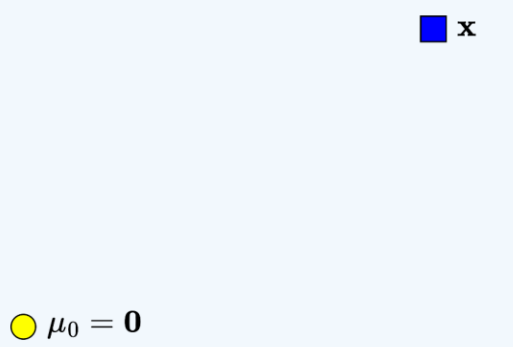
$$L^n(\mathbf{x}) = n \mathbb{E}_{i \sim U\{1, \dots, n\}} \mathbb{E}_{p_U(\boldsymbol{\theta}_{i-1} | \boldsymbol{\theta}_0, \mathbf{x}; \sum_{j=1}^{i-1} \alpha_j)} D_{\text{KL}} \left( p_S(\cdot | \mathbf{x}; \alpha_i) \parallel p_R(\cdot | \hat{\mathbf{x}}_i; \alpha_i) \right)$$

●  $\mu_0 = \mathbf{0}$

1. Sample:  $i \sim U\{1, \dots, n\}$

# Loss Function

$$L^n(\mathbf{x}) = n \mathbb{E}_{i \sim U\{1, \dots, n\}} \mathbb{E}_{p_U(\boldsymbol{\theta}_{i-1} | \boldsymbol{\theta}_0, \mathbf{x}; \sum_{j=1}^{i-1} \alpha_j)} D_{\text{KL}} \left( p_S(\cdot | \mathbf{x}; \alpha_i) \parallel p_R(\cdot | \hat{\mathbf{x}}_i; \alpha_i) \right)$$



$\mu_0 = 0$

1. Sample:  $i \sim U\{1, \dots, n\}$



$\mu_{i-1}$

2. Sample:  $\mu_{i-1} \sim p_U(\cdot | \mu_0, \mathbf{x}; \rho_{i-1})$

# Loss Function

$$L^n(\mathbf{x}) = n \mathbb{E}_{i \sim U\{1, \dots, n\}} \mathbb{E}_{p_U(\boldsymbol{\theta}_{i-1} | \boldsymbol{\theta}_0, \mathbf{x}; \sum_{j=1}^{i-1} \alpha_j)} D_{\text{KL}} \left( p_S(\cdot | \mathbf{x}; \alpha_i) \parallel p_R(\cdot | \hat{\mathbf{x}}_i; \alpha_i) \right)$$

●  $\mu_0 = \mathbf{0}$

■  $\mathbf{x}$

1. Sample:  $i \sim U\{1, \dots, n\}$

●  $\mu_{i-1}$

■  $\mathbf{x}$

2. Sample:  $\mu_{i-1} \sim p_U(\cdot | \mu_0, \mathbf{x}; \rho_{i-1})$

■  $\mathbf{x}$

■  $\hat{\mathbf{x}}_i$

3. Prediction:  $\hat{\mathbf{x}}_i = \Psi(\mu_{i-1}, i-1)$

# Loss Function

$$L^n(\mathbf{x}) = n \mathbb{E}_{i \sim U\{1, \dots, n\}} \mathbb{E}_{p_U(\boldsymbol{\theta}_{i-1} | \boldsymbol{\theta}_0, \mathbf{x}; \sum_{j=1}^{i-1} \alpha_j)} D_{\text{KL}} \left( p_S(\cdot | \mathbf{x}; \alpha_i) \parallel p_R(\cdot | \hat{\mathbf{x}}_i; \alpha_i) \right)$$

●  $\mu_0 = \mathbf{0}$

1. Sample:  $i \sim U\{1, \dots, n\}$

■  $\mathbf{x}$

●  $\mu_{i-1}$

■  $\mathbf{x}$

2. Sample:  $\mu_{i-1} \sim p_U(\cdot | \mu_0, \mathbf{x}; \rho_{i-1})$

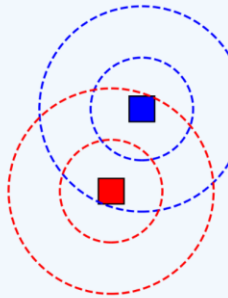
■  $\mathbf{x}$

■  $\hat{\mathbf{x}}_i$

3. Prediction:  $\hat{\mathbf{x}}_i = \Psi(\mu_{i-1}, i-1)$

4. Loss:

$$D_{\text{KL}} \left( \mathcal{N}_P(\mathbf{x}, \alpha_t) \parallel \mathcal{N}_P(\mathbf{x}, \alpha_t) \right)$$



# Flow Distribution



**Goal:** Extend to continuous time. First define the accuracy schedule  $\beta(t) = \int_0^t \alpha(t') dt'$

$$p_F(\boldsymbol{\theta} \mid \mathbf{x}; t) \triangleq p_U(\boldsymbol{\theta} \mid \boldsymbol{\theta}_0, \mathbf{x}; \beta(t))$$

# Flow Distribution



**Goal:** Extend to continuous time. First define the accuracy schedule  $\beta(t) = \int_0^t \alpha(t') dt'$

$$\begin{aligned} p_F(\boldsymbol{\theta} \mid \mathbf{x}; t) &\triangleq p_U(\boldsymbol{\theta} \mid \boldsymbol{\theta}_0, \mathbf{x}; \beta(t)) \\ &= \mathcal{N}\left(\boldsymbol{\mu} \mid \frac{\beta(t)}{1 + \beta(t)} \mathbf{x}, \frac{\beta(t)}{(1 + \beta(t))^2} I\right) \text{ Continuous case} \end{aligned}$$

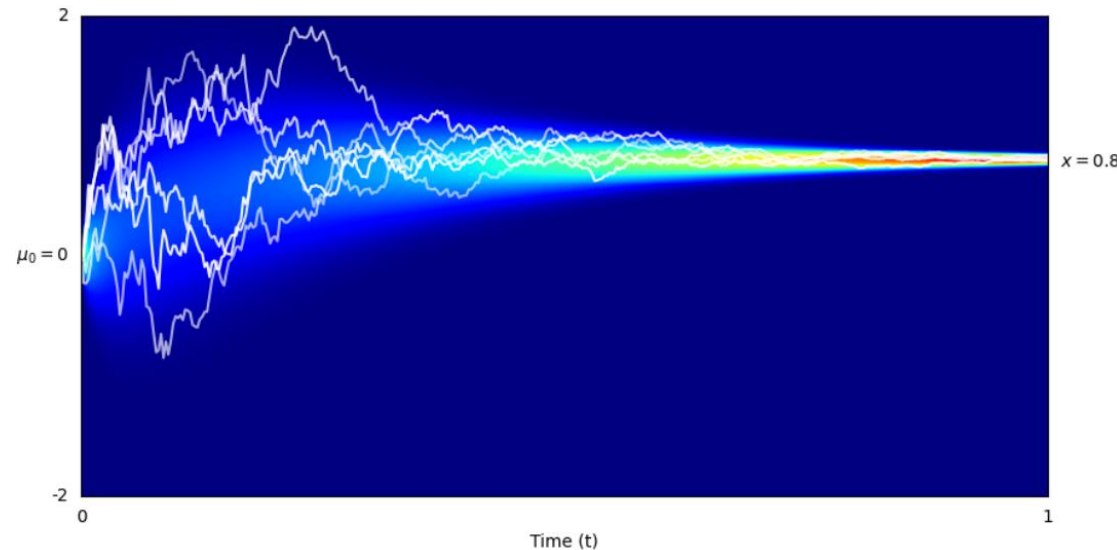
# Flow Distribution

**Goal:** Extend to continuous time. First define the accuracy schedule  $\beta(t) = \int_0^t \alpha(t') dt'$

$$\begin{aligned} p_F(\boldsymbol{\theta} \mid \mathbf{x}; t) &\triangleq p_U(\boldsymbol{\theta} \mid \boldsymbol{\theta}_0, \mathbf{x}; \beta(t)) \\ &= \mathcal{N}\left(\boldsymbol{\mu} \mid \frac{\beta(t)}{1 + \beta(t)} \mathbf{x}, \frac{\beta(t)}{(1 + \beta(t))^2} I\right) \text{ Continuous case} \end{aligned}$$

## Benefits:

- ✓ Permits continuous time loss  $L^\infty$
- ✓ Decoupled number of samples from  $t$



## Differences to Diffusion:

- Occurring in parameter space, not sample space
- Start with an initialisation instead of noisy sample

# Continuous Time Loss



Similarly to Variational Diffusion Models, we derive a continuous time loss

$$L^{\infty}(\mathbf{x}) \stackrel{\text{def}}{=} \lim_{n \rightarrow \infty} L^n(\mathbf{x})$$



# Continuous Time Loss

Similarly to Variational Diffusion Models, we derive a continuous time loss

$$\begin{aligned} L^\infty(\mathbf{x}) &\stackrel{\text{def}}{=} \lim_{n \rightarrow \infty} L^n(\mathbf{x}) \\ &= \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \mathbb{E}_{\substack{t \sim U(\epsilon, 1) \\ p_F(\boldsymbol{\theta} | \mathbf{x}, t - \epsilon)}} D_{KL} \left( p_S(\cdot | \mathbf{x}; \alpha(t, \epsilon)) \parallel p_R(\cdot | \hat{\mathbf{x}}; \alpha(t, \epsilon)) \right) \end{aligned}$$

# Continuous Time Loss

Similarly to Variational Diffusion Models, we derive a continuous time loss

$$\begin{aligned} L^\infty(\mathbf{x}) &\stackrel{\text{def}}{=} \lim_{n \rightarrow \infty} L^n(\mathbf{x}) \\ &= \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \mathbb{E}_{\substack{t \sim U(\epsilon, 1) \\ p_F(\boldsymbol{\theta} | \mathbf{x}, t - \epsilon)}} D_{KL} \left( p_S(\cdot | \mathbf{x}; \alpha(t, \epsilon)) \parallel p_R(\cdot | \hat{\mathbf{x}}; \alpha(t, \epsilon)) \right) \\ &= -\ln \sigma_1 \mathbb{E}_{t \sim U(0, 1), p_F(\boldsymbol{\theta} | \mathbf{x}; t)} \frac{\|\mathbf{x} - \hat{\mathbf{x}}(\boldsymbol{\theta}, t)\|^2}{\sigma_1^{2t}} \quad \text{Continuous case} \end{aligned}$$

# Continuous Time Loss

Similarly to Variational Diffusion Models, we derive a continuous time loss

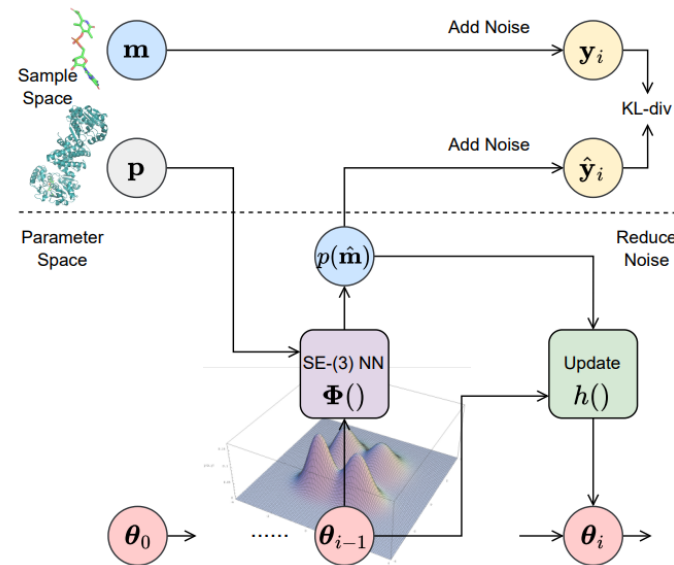
$$\begin{aligned} L^\infty(\mathbf{x}) &\stackrel{\text{def}}{=} \lim_{n \rightarrow \infty} L^n(\mathbf{x}) \\ &= \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \mathbb{E}_{\substack{t \sim U(\epsilon, 1) \\ p_F(\boldsymbol{\theta} | \mathbf{x}, t - \epsilon)}} D_{KL} \left( p_S(\cdot | \mathbf{x}; \alpha(t, \epsilon)) \parallel p_R(\cdot | \hat{\mathbf{x}}; \alpha(t, \epsilon)) \right) \\ &= -\ln \sigma_1 \mathbb{E}_{t \sim U(0, 1), p_F(\boldsymbol{\theta} | \mathbf{x}; t)} \frac{\|\mathbf{x} - \hat{\mathbf{x}}(\boldsymbol{\theta}, t)\|^2}{\sigma_1^{2t}} \quad \text{Continuous case} \end{aligned}$$

- ✓ No need to choose discrete number of steps
- ✓ Simplifies to a cleaner integral
- ✓ Flexible steps at inference time when trained on continuous loss

# Application: Hybrid Data

BFNs excel at **hybrid data tasks** due to their **unified framework**. For example, with structure-based drug design (SBDD) where we model (continuous) atom coordinates and (discrete) atom types

- ✓ Robust against mode collapse and noise
- ✓ Low variance updates and faster convergence
- ✓ Fewer sampling errors (incomplete or distorted molecules)



MolCRAFT: Structure-Based Drug Design in Continuous Parameter Space (Qu et al., 2024, arXiv:2404.12141)

Thank you for listening

# Supplementary Material



# ‘Full’ Loss



Define the  $n$ -step *discrete-time loss*  $L^n(\mathbf{x})$  as the expected number of nats required to first transmit  $\mathbf{y}_1, \dots, \mathbf{y}_n$ , and the *reconstruction loss*  $L^r(\mathbf{x})$  as the expected number of nats required to then transmit  $\mathbf{x}$ . Since — using a bits-back coding scheme [7, 11] — it requires  $D_{KL}(p_S \parallel p_R)$  nats to transmit a sample from  $p_S$  to a receiver with  $p_R$ ,

$$L^n(\mathbf{x}) = \mathbb{E}_{p(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{n-1})} \sum_{i=1}^n D_{KL}(p_S(\cdot \mid \mathbf{x}; \alpha_i) \parallel p_R(\cdot \mid \boldsymbol{\theta}_{i-1}; t_{i-1}, \alpha_i)), \quad (13)$$

$$L^r(\mathbf{x}) = - \mathbb{E}_{p_F(\boldsymbol{\theta} \mid \mathbf{x}, 1)} \ln p_O(\mathbf{x} \mid \boldsymbol{\theta}; 1).$$

Note that  $L^r(\mathbf{x})$  is not directly optimised in this paper; however it is indirectly trained by optimising  $L^n(\mathbf{x})$  since both are minimised by matching the output distribution to the data. Furthermore, as long as  $\beta(1)$  is high enough, the input distribution at  $t = 1$  will be very close to  $\mathbf{x}$ , making it trivial for the network to fit  $p_O(\mathbf{x} \mid \boldsymbol{\theta}; 1)$ .

The loss function  $L(\mathbf{x})$  is defined as the total number of nats required to transmit the data, which is the sum of the  $n$ -step and reconstruction losses:

$$L(\mathbf{x}) = L^n(\mathbf{x}) + L^r(\mathbf{x}) \quad (16)$$

Alternatively  $L(\mathbf{x})$  can be derived as the loss function of a variational autoencoder (VAE; [18]). Consider the sequence  $\mathbf{y}_1, \dots, \mathbf{y}_n$  as a latent code with posterior probability given by

$$q(\mathbf{y}_1, \dots, \mathbf{y}_n) = \prod_{i=1}^n p_S(\mathbf{y}_i \mid \mathbf{x}; \alpha_i), \quad (17)$$

and autoregressive prior probability given by

$$p(\mathbf{y}_1, \dots, \mathbf{y}_n) = \prod_{i=1}^n p_R(\mathbf{y}_i \mid \boldsymbol{\theta}_{i-1}; t_{i-1}, \alpha_i). \quad (18)$$

Then, noting that the decoder probability  $p(\mathbf{x} \mid \mathbf{y}_1, \dots, \mathbf{y}_n) = p_O(\mathbf{x} \mid \boldsymbol{\theta}_n; 1)$ , the complete transmission process defines a VAE with loss function given by the negative variational lower bound (VLB)

$$L(\mathbf{x}) = -\text{VLB}(\mathbf{x}) = D_{KL}(q \parallel p) - \mathbb{E}_{\mathbf{y}_1, \dots, \mathbf{y}_n \sim q} \ln p(\mathbf{x} \mid \mathbf{y}_1, \dots, \mathbf{y}_n) \quad (19)$$

$$= L^n(\mathbf{x}) + L^r(\mathbf{x}). \quad (20)$$



Given parameters  $\boldsymbol{\theta}$  and sender sample  $\mathbf{y}$  drawn with accuracy  $\alpha$  the *Bayesian update function*  $h$  is derived by applying the rules of Bayesian inference to compute the updated parameters  $\boldsymbol{\theta}'$ :

$$\boldsymbol{\theta}' \leftarrow h(\boldsymbol{\theta}, \mathbf{y}, \alpha). \quad (5)$$

The *Bayesian update distribution*  $p_U(\cdot \mid \boldsymbol{\theta}, \mathbf{x}; \alpha)$  is then defined by marginalizing out  $\mathbf{y}$ :

$$p_U(\boldsymbol{\theta}' \mid \boldsymbol{\theta}, \mathbf{x}; \alpha) = \mathbb{E}_{p_S(\mathbf{y} \mid \mathbf{x}; \alpha)} \delta(\boldsymbol{\theta}' - h(\boldsymbol{\theta}, \mathbf{y}, \alpha)), \quad (6)$$

where  $\delta(\cdot - \mathbf{a})$  is the multivariate Dirac delta distribution centred on the vector  $\mathbf{a}$ . In Sections 4.4 and 6.7 we will prove that both forms of  $p_U(\cdot \mid \boldsymbol{\theta}, \mathbf{x}; \alpha)$  considered in this paper have the following property: the accuracies are additive in the sense that if  $\alpha = \alpha_a + \alpha_b$  then

$$p_U(\boldsymbol{\theta}'' \mid \boldsymbol{\theta}, \mathbf{x}; \alpha) = \mathbb{E}_{p_U(\boldsymbol{\theta}' \mid \boldsymbol{\theta}, \mathbf{x}; \alpha_a)} p_U(\boldsymbol{\theta}'' \mid \boldsymbol{\theta}', \mathbf{x}; \alpha_b). \quad (7)$$

# Bayesian Updates: Continuous

Given a univariate Gaussian prior  $\mathcal{N}(\mu_a, \rho_a^{-1})$  over some unknown data  $x$  it can be shown [27] that the Bayesian posterior after observing a noisy sample  $y$  from a normal distribution  $\mathcal{N}(x, \alpha^{-1})$  with known precision  $\alpha$  is  $\mathcal{N}(\mu_b, \rho_b^{-1})$ , where

$$\rho_b = \rho_a + \alpha, \quad (46)$$

$$\mu_b = \frac{\mu_a \rho_a + y \alpha}{\rho_b}. \quad (47)$$

Since both  $p_I(\mathbf{x} \mid \boldsymbol{\theta})$  and  $p_S(\mathbf{y} \mid \mathbf{x}; \alpha)$  distributions are normal with diagonal covariance, Eqs. 46 and 47 can be applied to obtain the following Bayesian update function for parameters  $\boldsymbol{\theta}_{i-1} = \{\boldsymbol{\mu}_{i-1}, \rho_{i-1}\}$  and sender sample  $\mathbf{y}$  drawn from  $p_S(\cdot \mid \mathbf{x}; \alpha \mathbf{I}) = \mathcal{N}(\mathbf{x}, \alpha^{-1} \mathbf{I})$ :

$$h(\{\boldsymbol{\mu}_{i-1}, \rho_{i-1}\}, \mathbf{y}, \alpha) = \{\boldsymbol{\mu}_i, \rho_i\}, \quad (48)$$

with

$$\rho_i = \rho_{i-1} + \alpha, \quad (49)$$

$$\boldsymbol{\mu}_i = \frac{\boldsymbol{\mu}_{i-1} \rho_{i-1} + \mathbf{y} \alpha}{\rho_i}. \quad (50)$$