

# Forward diffusion process - improvements

Since  $q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})$  is Gaussian, we can directly compute  $q(\mathbf{x}^{(t)} | \mathbf{x}^{(0)})$  in closed-form!

$$\mathbf{x}^{(t)} = \prod_{i=1}^t \sqrt{1 - \beta_i} \mathbf{x}^{(0)} + \sqrt{1 - \prod_{i=1}^t (1 - \beta_i)} \epsilon$$

Let us define  $\alpha_t = 1 - \beta_t$ , and  $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$ .

$$\mathbf{x}^{(t)} = \sqrt{\bar{\alpha}_t} \mathbf{x}^{(0)} + \sqrt{1 - \bar{\alpha}_t} \epsilon; \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

# Reverse diffusion process

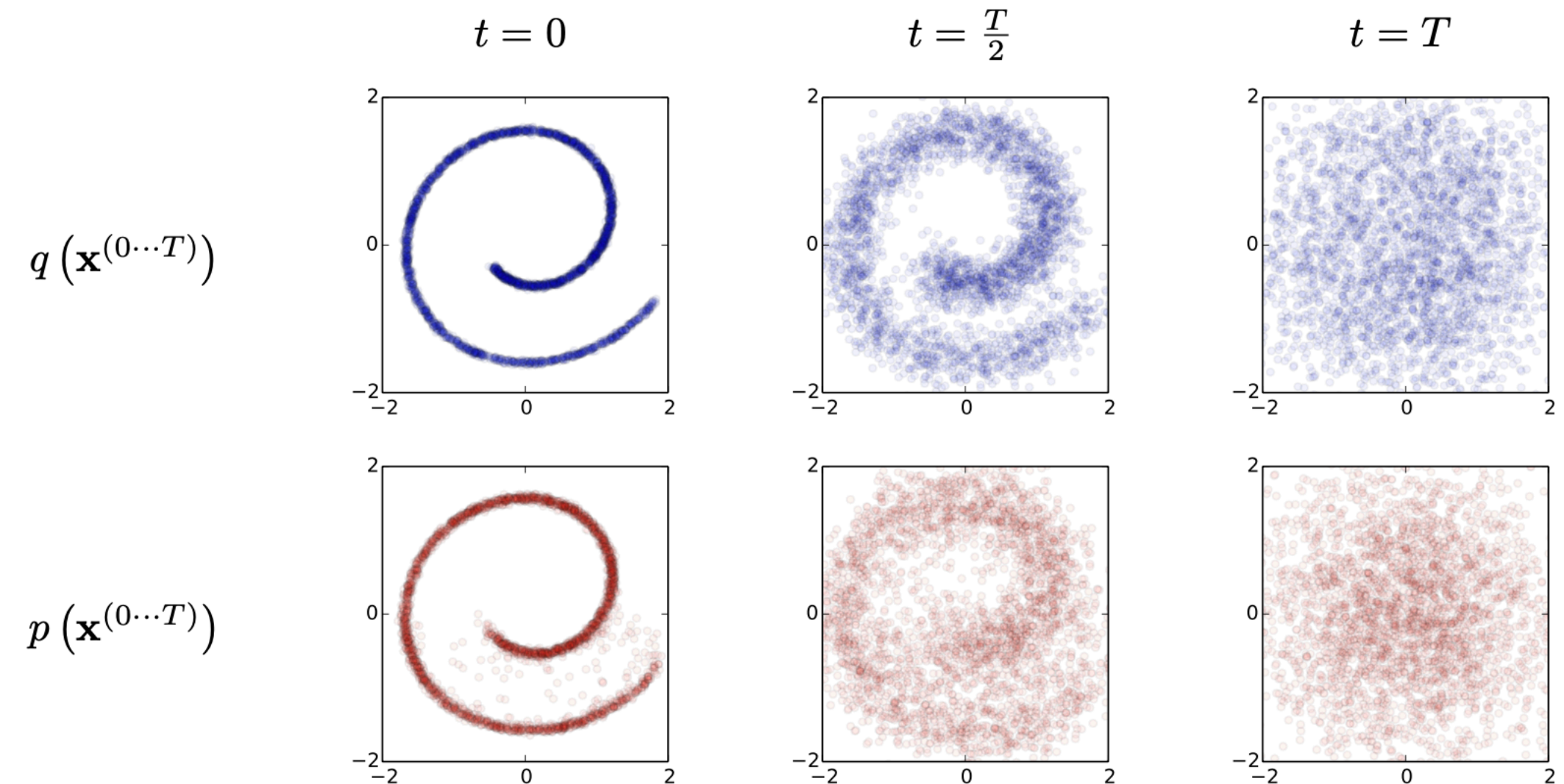
## 2.2. Reverse Trajectory

The generative distribution will be trained to describe the same trajectory, but in reverse,

$$p(\mathbf{x}^{(T)}) = \pi(\mathbf{x}^{(T)}) \quad (4)$$

$$p(\mathbf{x}^{(0:T)}) = p(\mathbf{x}^{(T)}) \prod_{t=1}^T p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}) \quad (5)$$

For both Gaussian and binomial diffusion, for continuous diffusion (limit of small step size  $\beta$ ) the reversal of the diffusion process has the identical functional form as the forward process (Feller, 1949). Since  $q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})$  is a Gaussian (binomial) distribution, and if  $\beta_t$  is small, then  $q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})$  will also be a Gaussian (binomial) distribution. The longer the trajectory the smaller the diffusion rate  $\beta$  can be made.



$$q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)}, \mathbf{x}^{(t-2)}) = q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})$$

$$q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(t+1)}) \neq q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})$$

$$p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(t+1)}) = p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})$$

$$p(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)}, \mathbf{x}^{(t-2)}) \neq p(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})$$

# Reverse diffusion process

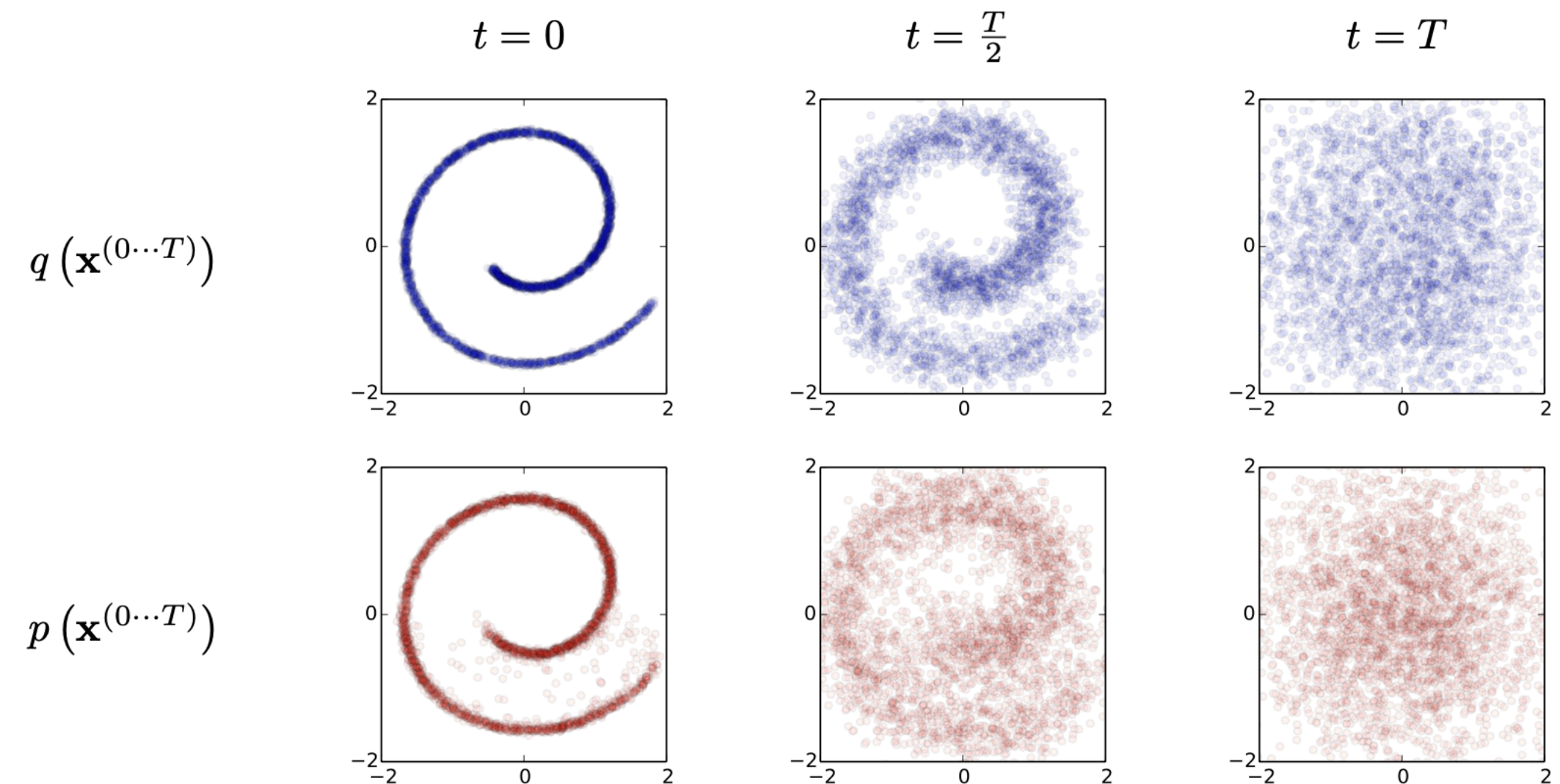
## 2.2. Reverse Trajectory

The generative distribution will be trained to describe the same trajectory, but in reverse,

$$p(\mathbf{x}^{(T)}) = \pi(\mathbf{x}^{(T)}) \quad (4)$$

$$p(\mathbf{x}^{(0 \dots T)}) = p(\mathbf{x}^{(T)}) \prod_{t=1}^T p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}) \quad (5)$$

For both Gaussian and binomial diffusion, for continuous diffusion (limit of small step size  $\beta$ ) the reversal of the diffusion process has the identical functional form as the forward process (Feller, 1949). Since  $q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})$  is a Gaussian (binomial) distribution, and if  $\beta_t$  is small, then  $q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})$  will also be a Gaussian (binomial) distribution. The longer the trajectory the smaller the diffusion rate  $\beta$  can be made.



1.  $\mathbf{x}^{(T)} \sim \pi(\mathbf{x}^{(T)})$
2.  $\mathbf{x}^{(T-1)} \sim p(\mathbf{x}^{(T-1)} | \mathbf{x}^{(T)})$
3.  $\mathbf{x}^{(T-2)} \sim p(\mathbf{x}^{(T-2)} | \mathbf{x}^{(T-1)})$
4.  $\dots$
5.  $\mathbf{x}^{(0)} \sim p(\mathbf{x}^{(0)} | \mathbf{x}^{(1)})$



# Reverse diffusion process

## 2.2. Reverse Trajectory

The generative distribution will be trained to describe the same trajectory, but in reverse,

$$p(\mathbf{x}^{(T)}) = \pi(\mathbf{x}^{(T)}) \quad (4)$$

$$p(\mathbf{x}^{(0:T)}) = p(\mathbf{x}^{(T)}) \prod_{t=1}^T p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}) . \quad (5)$$

For both Gaussian and binomial diffusion, for continuous diffusion (limit of small step size  $\beta$ ) the reversal of the diffusion process has the identical functional form as the forward process (Feller, 1949). Since  $q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})$  is a Gaussian (binomial) distribution, and if  $\beta_t$  is small, then  $q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})$  will also be a Gaussian (binomial) distribution. The longer the trajectory the smaller the diffusion rate  $\beta$  can be made.

During learning only the mean and covariance for a Gaussian diffusion kernel, or the bit flip probability for a binomial kernel, need be estimated. As shown in Table App.1,  $\mathbf{f}_\mu(\mathbf{x}^{(t)}, t)$  and  $\mathbf{f}_\Sigma(\mathbf{x}^{(t)}, t)$  are functions defining the mean and covariance of the reverse Markov transitions for a Gaussian, and  $\mathbf{f}_b(\mathbf{x}^{(t)}, t)$  is a function providing the bit flip probability for a binomial distribution. The computational cost of running this algorithm is the cost of these functions, times the number of time-steps. For all results in this paper, multi-layer perceptrons are used to define these functions. A wide range of regression or function fitting techniques would be applicable however, including nonparameteric methods.

1.  $\mathbf{x}^{(T)} \sim \pi(\mathbf{x}^{(T)})$
2.  $\mathbf{x}^{(T-1)} \sim p(\mathbf{x}^{(T-1)} | \mathbf{x}^{(T)})$
3.  $\mathbf{x}^{(T-2)} \sim p(\mathbf{x}^{(T-2)} | \mathbf{x}^{(T-1)})$
4.  $\dots$
5.  $\mathbf{x}^{(0)} \sim p(\mathbf{x}^{(0)} | \mathbf{x}^{(1)})$

|  |  | <i>Gaussian</i>  |
|--|--|--|
| Well behaved (analytically tractable) distribution | $\pi(\mathbf{x}^{(T)}) =$                    | $\mathcal{N}(\mathbf{x}^{(T)}; \mathbf{0}, \mathbf{I})$  |
| Forward diffusion kernel                           | $q(\mathbf{x}^{(t)}   \mathbf{x}^{(t-1)}) =$ | $\mathcal{N}(\mathbf{x}^{(t)}; \mathbf{x}^{(t-1)} \sqrt{1 - \beta_t}, \mathbf{I} \beta_t)$                     |
| Reverse diffusion kernel                           | $p(\mathbf{x}^{(t-1)}   \mathbf{x}^{(t)}) =$ | $\mathcal{N}(\mathbf{x}^{(t-1)}; \mathbf{f}_\mu(\mathbf{x}^{(t)}, t), \mathbf{f}_\Sigma(\mathbf{x}^{(t)}, t))$ |

# Training

$$\begin{aligned} & \mathbb{E}_{\mathbf{x}^{(0)} \sim \mathcal{D}_{train}} \log p(\mathbf{x}^{(0)}) \\ &= \mathbb{E}_{\mathbf{x}^{(0)} \sim q(\mathbf{x}^{(0)})} \log p(\mathbf{x}^{(0)}) \\ &= \int q(\mathbf{x}^{(0)}) \log p(\mathbf{x}^{(0)}) d\mathbf{x}^{(0)} \end{aligned}$$

## 2.4. Training

Training amounts to maximizing the model log likelihood,

$$L = \int d\mathbf{x}^{(0)} q(\mathbf{x}^{(0)}) \log p(\mathbf{x}^{(0)}) \quad (10)$$

# Training

## 2.3. Model Probability

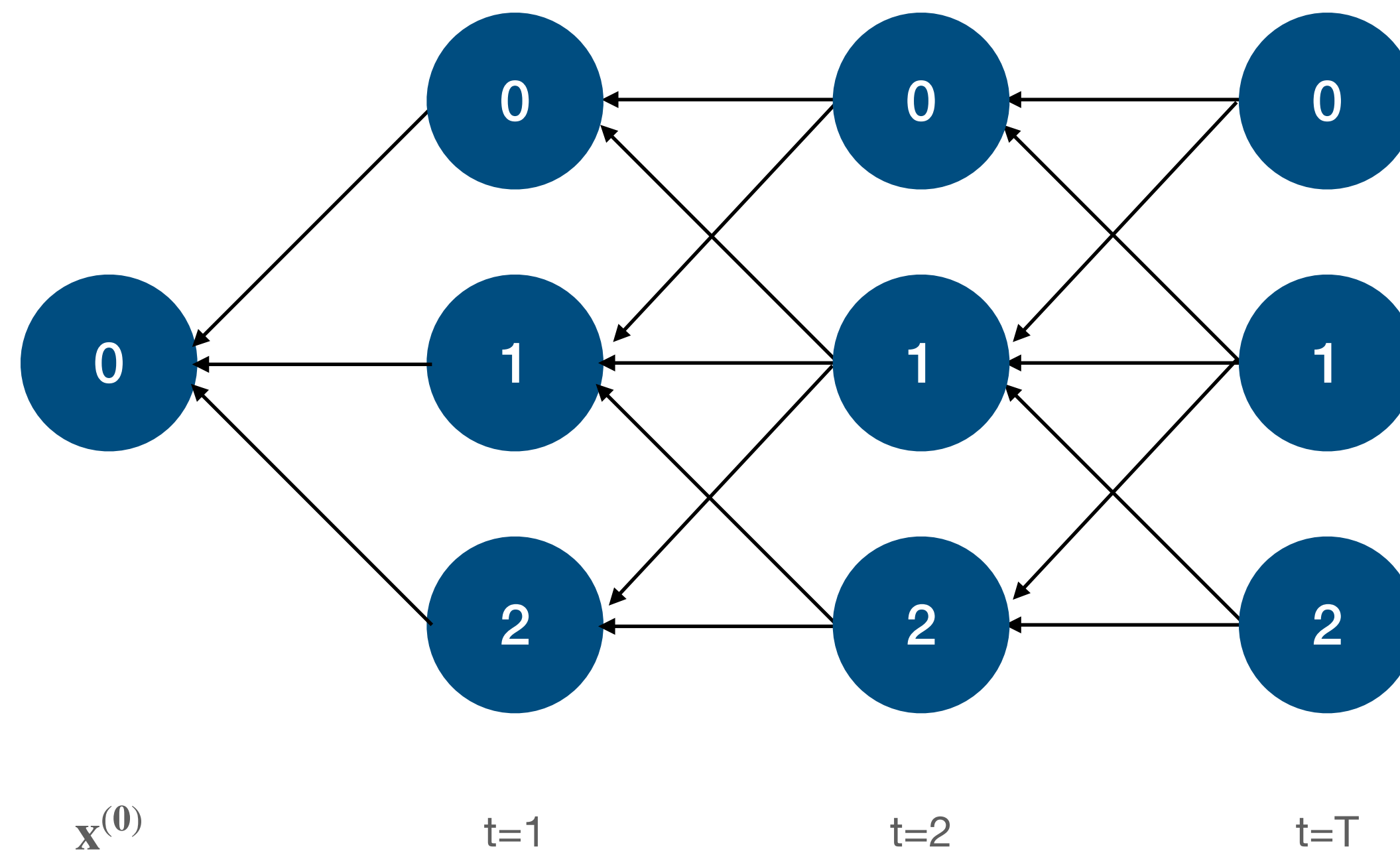
The probability the generative model assigns to the data is

$$p(\mathbf{x}^{(0)}) = \int d\mathbf{x}^{(1 \dots T)} p(\mathbf{x}^{(0 \dots T)}). \quad (6)$$

## 2.4. Training

Training amounts to maximizing the model log likelihood,

$$L = \int d\mathbf{x}^{(0)} q(\mathbf{x}^{(0)}) \log p(\mathbf{x}^{(0)}) \quad (10)$$



# Training

## 2.3. Model Probability

The probability the generative model assigns to the data is

$$p(\mathbf{x}^{(0)}) = \int d\mathbf{x}^{(1\cdots T)} p(\mathbf{x}^{(0\cdots T)}). \quad (6)$$

Naively this integral is intractable – but taking a cue from annealed importance sampling and the Jarzynski equality, we instead evaluate the relative probability of the forward and reverse trajectories, averaged over forward trajectories,

$$\begin{aligned} p(\mathbf{x}^{(0)}) &= \int d\mathbf{x}^{(1\cdots T)} p(\mathbf{x}^{(0\cdots T)}) \frac{q(\mathbf{x}^{(1\cdots T)}|\mathbf{x}^{(0)})}{q(\mathbf{x}^{(1\cdots T)}|\mathbf{x}^{(0)})} \quad (7) \\ &= \int d\mathbf{x}^{(1\cdots T)} q(\mathbf{x}^{(1\cdots T)}|\mathbf{x}^{(0)}) \frac{p(\mathbf{x}^{(0\cdots T)})}{q(\mathbf{x}^{(1\cdots T)}|\mathbf{x}^{(0)})} \quad (8) \end{aligned}$$

$$\begin{aligned} &= \int d\mathbf{x}^{(1\cdots T)} q(\mathbf{x}^{(1\cdots T)}|\mathbf{x}^{(0)}) \cdot \\ &\quad p(\mathbf{x}^{(T)}) \prod_{t=1}^T \frac{p(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)})}. \quad (9) \end{aligned}$$

## 2.4. Training

Training amounts to maximizing the model log likelihood,

$$L = \int d\mathbf{x}^{(0)} q(\mathbf{x}^{(0)}) \log p(\mathbf{x}^{(0)}) \quad (10)$$



# Training

## 2.3. Model Probability

The probability the generative model assigns to the data is

$$p(\mathbf{x}^{(0)}) = \int d\mathbf{x}^{(1\cdots T)} p(\mathbf{x}^{(0\cdots T)}). \quad (6)$$

Naively this integral is intractable – but taking a cue from annealed importance sampling and the Jarzynski equality, we instead evaluate the relative probability of the forward and reverse trajectories, averaged over forward trajectories,

$$p(\mathbf{x}^{(0)}) = \int d\mathbf{x}^{(1\cdots T)} p(\mathbf{x}^{(0\cdots T)}) \frac{q(\mathbf{x}^{(1\cdots T)}|\mathbf{x}^{(0)})}{q(\mathbf{x}^{(1\cdots T)}|\mathbf{x}^{(0)})} \quad (7)$$

$$= \int d\mathbf{x}^{(1\cdots T)} q(\mathbf{x}^{(1\cdots T)}|\mathbf{x}^{(0)}) \frac{p(\mathbf{x}^{(0\cdots T)})}{q(\mathbf{x}^{(1\cdots T)}|\mathbf{x}^{(0)})} \quad (8)$$

$$= \int d\mathbf{x}^{(1\cdots T)} q(\mathbf{x}^{(1\cdots T)}|\mathbf{x}^{(0)}) \cdot p(\mathbf{x}^{(T)}) \prod_{t=1}^T \frac{p(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)})}. \quad (9)$$

## 2.4. Training

Training amounts to maximizing the model log likelihood,

$$L = \int d\mathbf{x}^{(0)} q(\mathbf{x}^{(0)}) \log p(\mathbf{x}^{(0)}) \quad (10)$$

$$= \int d\mathbf{x}^{(0)} q(\mathbf{x}^{(0)}) \cdot \log \left[ \frac{\int d\mathbf{x}^{(1\cdots T)} q(\mathbf{x}^{(1\cdots T)}|\mathbf{x}^{(0)}) \cdot p(\mathbf{x}^{(T)}) \prod_{t=1}^T \frac{p(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)})}}{p(\mathbf{x}^{(0)})} \right], \quad (11)$$

which has a lower bound provided by Jensen's inequality,

$$L \geq \int d\mathbf{x}^{(0\cdots T)} q(\mathbf{x}^{(0\cdots T)}) \cdot \log \left[ p(\mathbf{x}^{(T)}) \prod_{t=1}^T \frac{p(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)})} \right]. \quad (12)$$

$$\log(\mathbb{E}[\mathbb{X}]) \leq (\mathbb{E}[\log \mathbb{X}])$$



# Training

$$L \geq \mathbb{E}_{\mathbf{x}^{(0..T)} \sim q(\mathbf{x}^{(0..T)})} \log[p(\mathbf{x}^{(T)}) \prod_{t=1}^T \frac{p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})}]$$

## 2.4. Training

Training amounts to maximizing the model log likelihood,

$$L = \int d\mathbf{x}^{(0)} q(\mathbf{x}^{(0)}) \log p(\mathbf{x}^{(0)}) \quad (10)$$

$$= \int d\mathbf{x}^{(0)} q(\mathbf{x}^{(0)}) \cdot \log \left[ \frac{\int d\mathbf{x}^{(1..T)} q(\mathbf{x}^{(1..T)} | \mathbf{x}^{(0)}) \cdot p(\mathbf{x}^{(T)}) \prod_{t=1}^T \frac{p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})}}{p(\mathbf{x}^{(T)}) \prod_{t=1}^T \frac{p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})}} \right], \quad (11)$$

which has a lower bound provided by Jensen's inequality,

$$L \geq \int d\mathbf{x}^{(0..T)} q(\mathbf{x}^{(0..T)}) \cdot \log \left[ p(\mathbf{x}^{(T)}) \prod_{t=1}^T \frac{p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})} \right]. \quad (12)$$

# Training

$$L \geq \mathbb{E}_{\mathbf{x}^{(0..T)} \sim q(\mathbf{x}^{(0..T)})} \log[p(\mathbf{x}^{(T)}) \prod_{t=1}^T \frac{p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})}]$$

## 2.4. Training

Training amounts to maximizing the model log likelihood,

$$L = \int d\mathbf{x}^{(0)} q(\mathbf{x}^{(0)}) \log p(\mathbf{x}^{(0)}) \quad (10)$$

$$= \int d\mathbf{x}^{(0)} q(\mathbf{x}^{(0)}) \cdot$$

$$\log \left[ \frac{\int d\mathbf{x}^{(1..T)} q(\mathbf{x}^{(1..T)} | \mathbf{x}^{(0)}) \cdot p(\mathbf{x}^{(T)}) \prod_{t=1}^T \frac{p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})}}{p(\mathbf{x}^{(T)}) \prod_{t=1}^T \frac{p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})}} \right], \quad (11)$$

which has a lower bound provided by Jensen's inequality,

$$L \geq \int d\mathbf{x}^{(0..T)} q(\mathbf{x}^{(0..T)}) \cdot \log \left[ p(\mathbf{x}^{(T)}) \prod_{t=1}^T \frac{p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})} \right]. \quad (12)$$

As described in Appendix B, for our diffusion trajectories this reduces to,

$$L \geq K \quad (13)$$

$$K = - \sum_{t=2}^T \int d\mathbf{x}^{(0)} d\mathbf{x}^{(t)} q(\mathbf{x}^{(0)}, \mathbf{x}^{(t)}) \cdot$$

$$D_{KL} \left( q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}) \parallel p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}) \right) + H_q(\mathbf{X}^{(T)} | \mathbf{X}^{(0)}) - H_q(\mathbf{X}^{(1)} | \mathbf{X}^{(0)}) - H_p(\mathbf{X}^{(T)}). \quad (14)$$

# Training

$$H(x) = \frac{1}{2} + \frac{1}{2} \log(2\pi\sigma^2)$$

$$H(\mathbf{x}^{(T)}) = \frac{1}{2} + \frac{1}{2} \log(2\pi)$$

$$H(\mathbf{x}^{(1)} | \mathbf{x}^{(0)}) = \frac{1}{2} + \frac{1}{2} \log(2\pi\beta_0)$$

## 2.4. Training

Training amounts to maximizing the model log likelihood,

$$L = \int d\mathbf{x}^{(0)} q(\mathbf{x}^{(0)}) \log p(\mathbf{x}^{(0)}) \quad (10)$$

$$= \int d\mathbf{x}^{(0)} q(\mathbf{x}^{(0)}) \cdot$$

$$\log \left[ \frac{\int d\mathbf{x}^{(1 \dots T)} q(\mathbf{x}^{(1 \dots T)} | \mathbf{x}^{(0)}) \cdot p(\mathbf{x}^{(T)}) \prod_{t=1}^T \frac{p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})}}{p(\mathbf{x}^{(T)}) \prod_{t=1}^T \frac{p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})}} \right], \quad (11)$$

which has a lower bound provided by Jensen's inequality,

$$L \geq \int d\mathbf{x}^{(0 \dots T)} q(\mathbf{x}^{(0 \dots T)}) \cdot \log \left[ p(\mathbf{x}^{(T)}) \prod_{t=1}^T \frac{p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})} \right]. \quad (12)$$

As described in Appendix B, for our diffusion trajectories this reduces to,

$$L \geq K \quad (13)$$

$$K = - \sum_{t=2}^T \int d\mathbf{x}^{(0)} d\mathbf{x}^{(t)} q(\mathbf{x}^{(0)}, \mathbf{x}^{(t)}) \cdot$$

$$D_{KL} \left( q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}) \parallel p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}) \right)$$

$$+ H_q(\mathbf{x}^{(T)} | \mathbf{x}^{(0)}) - H_q(\mathbf{x}^{(1)} | \mathbf{x}^{(0)}) - H_p(\mathbf{x}^{(T)}). \quad (14)$$



# Training

$$K = - \sum_{t=2}^T \mathbb{E}_{\mathbf{x}^{(0)} \sim \mathcal{D}_{train} \mathbf{x}^{(t)} \sim q(\mathbf{x}^{(t)} | \mathbf{x}^{(0)})} D_{KL}(q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}) || p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}))$$

$$D_{KL}(p || q) = \log \frac{\sigma_q}{\sigma_p} + \frac{\sigma_p^2 + (\mu_p - \mu_q)^2}{2\sigma_q^2} - \frac{1}{2}$$

## 2.4. Training

Training amounts to maximizing the model log likelihood,

$$L = \int d\mathbf{x}^{(0)} q(\mathbf{x}^{(0)}) \log p(\mathbf{x}^{(0)}) \quad (10)$$

$$= \int d\mathbf{x}^{(0)} q(\mathbf{x}^{(0)}) \cdot$$

$$\log \left[ \frac{\int d\mathbf{x}^{(1 \dots T)} q(\mathbf{x}^{(1 \dots T)} | \mathbf{x}^{(0)}) \cdot p(\mathbf{x}^{(T)}) \prod_{t=1}^T \frac{p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})}}{p(\mathbf{x}^{(T)}) \prod_{t=1}^T \frac{p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})}} \right], \quad (11)$$

which has a lower bound provided by Jensen's inequality,

$$L \geq \int d\mathbf{x}^{(0 \dots T)} q(\mathbf{x}^{(0 \dots T)}) \cdot$$

$$\log \left[ p(\mathbf{x}^{(T)}) \prod_{t=1}^T \frac{p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})} \right]. \quad (12)$$

As described in Appendix B, for our diffusion trajectories this reduces to,

$$L \geq K \quad (13)$$

$$K = - \sum_{t=2}^T \int d\mathbf{x}^{(0)} d\mathbf{x}^{(t)} q(\mathbf{x}^{(0)}, \mathbf{x}^{(t)}) \cdot$$

$$D_{KL} \left( q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}) || p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}) \right) \\ - H_q(\mathbf{x}^{(T)} | \mathbf{x}^{(0)}) - H_q(\mathbf{x}^{(1)} | \mathbf{x}^{(0)}) - H_p(\mathbf{x}^{(T)}). \quad (14)$$

# Training

$$K = - \sum_{t=2}^T \mathbb{E}_{\mathbf{x}^{(0)} \sim \mathcal{D}_{train} \mathbf{x}^{(t)} \sim q(\mathbf{x}^{(t)} | \mathbf{x}^{(0)})} D_{KL}(q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}) || p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}))$$

$$D_{KL}(p || q) = \log \frac{\sigma_q}{\sigma_p} + \frac{\sigma_p^2 + (\mu_p - \mu_q)^2}{2\sigma_q^2} - \frac{1}{2}$$

We don't know the mean  $\mu_q$  and standard deviation  $\sigma_q$  for  $(q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}))$ !

## 2.4. Training

Training amounts to maximizing the model log likelihood,

$$L = \int d\mathbf{x}^{(0)} q(\mathbf{x}^{(0)}) \log p(\mathbf{x}^{(0)}) \quad (10)$$

$$= \int d\mathbf{x}^{(0)} q(\mathbf{x}^{(0)}) \cdot$$

$$\log \left[ \frac{\int d\mathbf{x}^{(1 \dots T)} q(\mathbf{x}^{(1 \dots T)} | \mathbf{x}^{(0)}) \cdot p(\mathbf{x}^{(T)}) \prod_{t=1}^T \frac{p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})}}{p(\mathbf{x}^{(T)}) \prod_{t=1}^T \frac{p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})}} \right], \quad (11)$$

which has a lower bound provided by Jensen's inequality,

$$L \geq \int d\mathbf{x}^{(0 \dots T)} q(\mathbf{x}^{(0 \dots T)}) \cdot$$

$$\log \left[ p(\mathbf{x}^{(T)}) \prod_{t=1}^T \frac{p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})} \right]. \quad (12)$$

As described in Appendix B, for our diffusion trajectories this reduces to,

$$L \geq K \quad (13)$$

$$K = - \sum_{t=2}^T \int d\mathbf{x}^{(0)} d\mathbf{x}^{(t)} q(\mathbf{x}^{(0)}, \mathbf{x}^{(t)}) \cdot$$

$$D_{KL} \left( q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}) || p(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}) \right)$$

~~$$+ H_q(\mathbf{x}^{(T)} | \mathbf{x}^{(0)}) - H_q(\mathbf{x}^{(1)} | \mathbf{x}^{(0)}) - H_p(\mathbf{x}^{(T)}).$$~~

(14)

# Training

We don't know the mean  $\mu_q$  and standard deviation  $\sigma_q$  for  $q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)})$ !

$$p(x | y) = \frac{p(y | x)p(x)}{p(y)}$$

$$\frac{1}{\sigma\sqrt{(2\pi)}} \exp^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

$$q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}) = \frac{q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)}, \mathbf{x}^{(0)})q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(0)})}{q(\mathbf{x}^{(t)} | \mathbf{x}^{(0)})}$$

$$= \frac{1}{\sqrt{\frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}} \beta_t \sqrt{(2\pi)}} \exp \left[ -\frac{1}{2} \left( \frac{\mathbf{x}^{(t-1)} - \left( \frac{\sqrt{\bar{\alpha}_{t-1}} \beta_t}{1 - \bar{\alpha}_t} \mathbf{x}^{(0)} + \frac{\sqrt{\bar{\alpha}_t} (1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{x}^{(t)} \right)}{\sqrt{\frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}} \beta_t} \right)^2 \right]$$