

# Appendix

## A. Conditional Entropy Bounds Derivation

The conditional entropy  $H_q(\mathbf{X}^{(t-1)}|\mathbf{X}^{(t)})$  of a step in the reverse trajectory is

$$H_q(\mathbf{X}^{(t-1)}, \mathbf{X}^{(t)}) = H_q(\mathbf{X}^{(t)}, \mathbf{X}^{(t-1)}) \quad (24)$$

$$H_q(\mathbf{X}^{(t-1)}|\mathbf{X}^{(t)}) + H_q(\mathbf{X}^{(t)}) = H_q(\mathbf{X}^{(t)}|\mathbf{X}^{(t-1)}) + H_q(\mathbf{X}^{(t-1)}) \quad (25)$$

$$H_q(\mathbf{X}^{(t-1)}|\mathbf{X}^{(t)}) = H_q(\mathbf{X}^{(t)}|\mathbf{X}^{(t-1)}) + H_q(\mathbf{X}^{(t-1)}) - H_q(\mathbf{X}^{(t)}) \quad (26)$$

An upper bound on the entropy change can be constructed by observing that  $\pi(\mathbf{y})$  is the maximum entropy distribution. This holds without qualification for the binomial distribution, and holds for variance 1 training data for the Gaussian case. For the Gaussian case, training data must therefore be scaled to have unit norm for the following equalities to hold. It need not be whitened. The upper bound is derived as follows,

$$H_q(\mathbf{X}^{(t)}) \geq H_q(\mathbf{X}^{(t-1)}) \quad (27)$$

$$H_q(\mathbf{X}^{(t-1)}) - H_q(\mathbf{X}^{(t)}) \leq 0 \quad (28)$$

$$H_q(\mathbf{X}^{(t-1)}|\mathbf{X}^{(t)}) \leq H_q(\mathbf{X}^{(t)}|\mathbf{X}^{(t-1)}). \quad (29)$$

A lower bound on the entropy difference can be established by observing that additional steps in a Markov chain do not increase the information available about the initial state in the chain, and thus do not decrease the conditional entropy of the initial state,

$$H_q(\mathbf{X}^{(0)}|\mathbf{X}^{(t)}) \geq H_q(\mathbf{X}^{(0)}|\mathbf{X}^{(t-1)}) \quad (30)$$

$$H_q(\mathbf{X}^{(t-1)}) - H_q(\mathbf{X}^{(t)}) \geq H_q(\mathbf{X}^{(0)}|\mathbf{X}^{(t-1)}) + H_q(\mathbf{X}^{(t-1)}) - H_q(\mathbf{X}^{(0)}|\mathbf{X}^{(t)}) - H_q(\mathbf{X}^{(t)}) \quad (31)$$

$$H_q(\mathbf{X}^{(t-1)}) - H_q(\mathbf{X}^{(t)}) \geq H_q(\mathbf{X}^{(0)}, \mathbf{X}^{(t-1)}) - H_q(\mathbf{X}^{(0)}, \mathbf{X}^{(t)}) \quad (32)$$

$$H_q(\mathbf{X}^{(t-1)}) - H_q(\mathbf{X}^{(t)}) \geq H_q(\mathbf{X}^{(t-1)}|\mathbf{X}^{(0)}) - H_q(\mathbf{X}^{(t)}|\mathbf{X}^{(0)}) \quad (33)$$

$$H_q(\mathbf{X}^{(t-1)}|\mathbf{X}^{(t)}) \geq H_q(\mathbf{X}^{(t)}|\mathbf{X}^{(t-1)}) + H_q(\mathbf{X}^{(t-1)}|\mathbf{X}^{(0)}) - H_q(\mathbf{X}^{(t)}|\mathbf{X}^{(0)}). \quad (34)$$

Combining these expressions, we bound the conditional entropy for a single step,

$$H_q(\mathbf{X}^{(t)}|\mathbf{X}^{(t-1)}) \geq H_q(\mathbf{X}^{(t-1)}|\mathbf{X}^{(t)}) \geq H_q(\mathbf{X}^{(t)}|\mathbf{X}^{(t-1)}) + H_q(\mathbf{X}^{(t-1)}|\mathbf{X}^{(0)}) - H_q(\mathbf{X}^{(t)}|\mathbf{X}^{(0)}), \quad (35)$$

where both the upper and lower bounds depend only on the conditional forward trajectory  $q(\mathbf{x}^{(1 \dots T)}|\mathbf{x}^{(0)})$ , and can be analytically computed.

## B. Log Likelihood Lower Bound

The lower bound on the log likelihood is

$$L \geq K \quad (36)$$

$$K = \int d\mathbf{x}^{(0 \dots T)} q(\mathbf{x}^{(0 \dots T)}) \log \left[ p(\mathbf{x}^{(T)}) \prod_{t=1}^T \frac{p(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)})} \right] \quad (37)$$

$$(38)$$

**B.1. Entropy of  $p(\mathbf{X}^{(T)})$** 

We can peel off the contribution from  $p(\mathbf{X}^{(T)})$ , and rewrite it as an entropy,

$$K = \int d\mathbf{x}^{(0\dots T)} q(\mathbf{x}^{(0\dots T)}) \sum_{t=1}^T \log \left[ \frac{p(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)})} \right] + \int d\mathbf{x}^{(T)} q(\mathbf{x}^{(T)}) \log p(\mathbf{x}^{(T)}) \quad (39)$$

$$= \int d\mathbf{x}^{(0\dots T)} q(\mathbf{x}^{(0\dots T)}) \sum_{t=1}^T \log \left[ \frac{p(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)})} \right] + \int d\mathbf{x}^{(T)} q(\mathbf{x}^{(T)}) \log \pi(\mathbf{x}^{(T)}) \quad (40)$$

$$(41)$$

By design, the cross entropy to  $\pi(\mathbf{x}^{(t)})$  is constant under our diffusion kernels, and equal to the entropy of  $p(\mathbf{x}^{(T)})$ . Therefore,

$$K = \sum_{t=1}^T \int d\mathbf{x}^{(0\dots T)} q(\mathbf{x}^{(0\dots T)}) \log \left[ \frac{p(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)})} \right] - H_p(\mathbf{X}^{(T)}). \quad (42)$$

**B.2. Remove the edge effect at  $t = 0$** 

In order to avoid edge effects, we set the final step of the reverse trajectory to be identical to the corresponding forward diffusion step,

$$p(\mathbf{x}^{(0)}|\mathbf{x}^{(1)}) = q(\mathbf{x}^{(1)}|\mathbf{x}^{(0)}) \frac{\pi(\mathbf{x}^{(0)})}{\pi(\mathbf{x}^{(1)})} = T_\pi(\mathbf{x}^{(0)}|\mathbf{x}^{(1)}; \beta_1). \quad (43)$$

We then use this equivalence to remove the contribution of the first time-step in the sum,

$$K = \sum_{t=2}^T \int d\mathbf{x}^{(0\dots T)} q(\mathbf{x}^{(0\dots T)}) \log \left[ \frac{p(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)})} \right] + \int d\mathbf{x}^{(0)} d\mathbf{x}^{(1)} q(\mathbf{x}^{(0)}, \mathbf{x}^{(1)}) \log \left[ \frac{q(\mathbf{x}^{(1)}|\mathbf{x}^{(0)}) \pi(\mathbf{x}^{(0)})}{q(\mathbf{x}^{(1)}|\mathbf{x}^{(0)}) \pi(\mathbf{x}^{(1)})} \right] - H_p(\mathbf{X}^{(T)}) \quad (44)$$

$$= \sum_{t=2}^T \int d\mathbf{x}^{(0\dots T)} q(\mathbf{x}^{(0\dots T)}) \log \left[ \frac{p(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)})} \right] - H_p(\mathbf{X}^{(T)}), \quad (45)$$

where we again used the fact that by design  $-\int d\mathbf{x}^{(t)} q(\mathbf{x}^{(t)}) \log \pi(\mathbf{x}^{(t)}) = H_p(\mathbf{X}^{(T)})$  is a constant for all  $t$ .

**B.3. Rewrite in terms of posterior  $q(\mathbf{x}^{(t-1)}|\mathbf{x}^{(0)})$** 

Because the forward trajectory is a Markov process,

$$K = \sum_{t=2}^T \int d\mathbf{x}^{(0\dots T)} q(\mathbf{x}^{(0\dots T)}) \log \left[ \frac{p(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)})}{q(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)}, \mathbf{x}^{(0)})} \right] - H_p(\mathbf{X}^{(T)}). \quad (46)$$

Using Bayes' rule we can rewrite this in terms of a posterior and marginals from the forward trajectory,

$$K = \sum_{t=2}^T \int d\mathbf{x}^{(0\dots T)} q(\mathbf{x}^{(0\dots T)}) \log \left[ \frac{p(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)})}{q(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}, \mathbf{x}^{(0)})} \frac{q(\mathbf{x}^{(t-1)}|\mathbf{x}^{(0)})}{q(\mathbf{x}^{(t)}|\mathbf{x}^{(0)})} \right] - H_p(\mathbf{X}^{(T)}). \quad (47)$$

#### B.4. Rewrite in terms of KL divergences and entropies

We then recognize that several terms are conditional entropies,

$$K = \sum_{t=2}^T \int d\mathbf{x}^{(0\dots T)} q(\mathbf{x}^{(0\dots T)}) \log \left[ \frac{p(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)})}{q(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}, \mathbf{x}^{(0)})} \right] + \sum_{t=2}^T \left[ H_q(\mathbf{X}^{(t)}|\mathbf{X}^{(0)}) - H_q(\mathbf{X}^{(t-1)}|\mathbf{X}^{(0)}) \right] - H_p(\mathbf{X}^{(T)}) \quad (48)$$

$$= \sum_{t=2}^T \int d\mathbf{x}^{(0\dots T)} q(\mathbf{x}^{(0\dots T)}) \log \left[ \frac{p(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)})}{q(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}, \mathbf{x}^{(0)})} \right] + H_q(\mathbf{X}^{(T)}|\mathbf{X}^{(0)}) - H_q(\mathbf{X}^{(1)}|\mathbf{X}^{(0)}) - H_p(\mathbf{X}^{(T)}). \quad (49)$$

Finally we transform the log ratio of probability distributions into a KL divergence,

$$K = - \sum_{t=2}^T \int d\mathbf{x}^{(0)} d\mathbf{x}^{(t)} q(\mathbf{x}^{(0)}, \mathbf{x}^{(t)}) D_{KL} \left( q(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}, \mathbf{x}^{(0)}) \parallel p(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}) \right) + H_q(\mathbf{X}^{(T)}|\mathbf{X}^{(0)}) - H_q(\mathbf{X}^{(1)}|\mathbf{X}^{(0)}) - H_p(\mathbf{X}^{(T)}). \quad (50)$$

Note that the entropies can be analytically computed, and the KL divergence can be analytically computed given  $\mathbf{x}^{(0)}$  and  $\mathbf{x}^{(t)}$ .

#### C. Markov Kernel of Perturbed Distribution

In Equations 19 and 20, the perturbed diffusion kernels are set as follows (unlike in the text body, we include the normalization constant)

$$\tilde{q}(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)}) = \frac{q(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)}) r(\mathbf{x}^{(t)})}{\int d\mathbf{x}^{(t)} q(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)}) r(\mathbf{x}^{(t)})}, \quad (51)$$

$$\tilde{q}(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}) = \frac{q(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}) r(\mathbf{x}^{(t-1)})}{\int d\mathbf{x}^{(t-1)} q(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}) r(\mathbf{x}^{(t-1)})}, \quad (52)$$

or writing them instead in terms of the original transitions,

$$q(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)}) = \frac{\tilde{q}(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)}) \int d\mathbf{x}^{(t)} q(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)}) r(\mathbf{x}^{(t)})}{r(\mathbf{x}^{(t)})}, \quad (53)$$

$$q(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}) = \frac{\tilde{q}(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}) \int d\mathbf{x}^{(t-1)} q(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}) r(\mathbf{x}^{(t-1)})}{r(\mathbf{x}^{(t-1)})}. \quad (54)$$

Similarly, we write Equation 16 in terms of the original forward distributions,

$$q(\mathbf{x}^{(t)}) = \frac{\tilde{q}(\mathbf{x}^{(t)}) \tilde{Z}_t}{r(\mathbf{x}^{(t)})}. \quad (55)$$

We substitute into Equation 17,

$$q(\mathbf{x}^{(t-1)}) q(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)}) = q(\mathbf{x}^{(t)}) q(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}), \quad (56)$$

$$\frac{\tilde{q}(\mathbf{x}^{(t-1)}) \tilde{Z}_{t-1}}{r(\mathbf{x}^{(t-1)})} \frac{\tilde{q}(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)}) \int d\mathbf{x}^{(t)} q(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)}) r(\mathbf{x}^{(t)})}{r(\mathbf{x}^{(t)})} = \frac{\tilde{q}(\mathbf{x}^{(t)}) \tilde{Z}_t}{r(\mathbf{x}^{(t)})} \frac{\tilde{q}(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}) \int d\mathbf{x}^{(t-1)} q(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}) r(\mathbf{x}^{(t-1)})}{r(\mathbf{x}^{(t-1)})}, \quad (57)$$

$$\tilde{q}(\mathbf{x}^{(t-1)}) \tilde{q}(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)}) \tilde{Z}_{t-1} \int d\mathbf{x}^{(t)} q(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)}) r(\mathbf{x}^{(t)}) = \tilde{q}(\mathbf{x}^{(t)}) \tilde{q}(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}) \tilde{Z}_t \int d\mathbf{x}^{(t-1)} q(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}) r(\mathbf{x}^{(t-1)}). \quad (58)$$

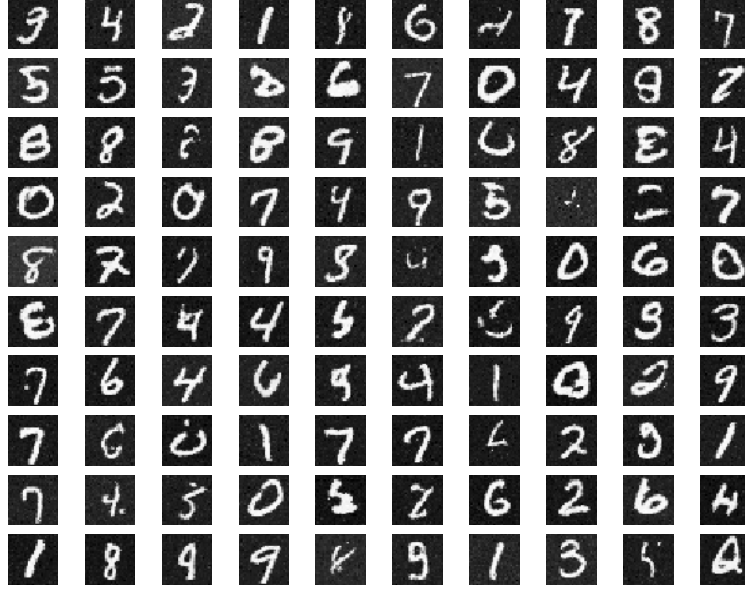


Figure App.1. Samples from a diffusion probabilistic model trained on MNIST digits. Note that unlike many MNIST sample figures, these are true samples rather than the mean of the Gaussian or binomial distribution from which samples would be drawn.

We then substitute  $\tilde{Z}_t = \int d\mathbf{x}^{(t)} q(\mathbf{x}^{(t)}) r(\mathbf{x}^{(t)})$ ,

$$\begin{aligned} \tilde{q}(\mathbf{x}^{(t-1)}) \tilde{q}(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)}) \int d\mathbf{x}^{(t-1)} q(\mathbf{x}^{(t-1)}) r(\mathbf{x}^{(t-1)}) \int d\mathbf{x}^{(t)} q(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)}) r(\mathbf{x}^{(t)}) \\ = \tilde{q}(\mathbf{x}^{(t)}) \tilde{q}(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}) \int d\mathbf{x}^{(t)} q(\mathbf{x}^{(t)}) r(\mathbf{x}^{(t)}) \int d\mathbf{x}^{(t-1)} q(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}) r(\mathbf{x}^{(t-1)}), \end{aligned} \quad (59)$$

$$\begin{aligned} \tilde{q}(\mathbf{x}^{(t-1)}) \tilde{q}(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)}) \int d\mathbf{x}^{(t-1)} d\mathbf{x}^{(t)} q(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)}) q(\mathbf{x}^{(t-1)}) r(\mathbf{x}^{(t)}) r(\mathbf{x}^{(t-1)}) \\ = \tilde{q}(\mathbf{x}^{(t)}) \tilde{q}(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}) \int d\mathbf{x}^{(t-1)} d\mathbf{x}^{(t)} q(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}) q(\mathbf{x}^{(t)}) r(\mathbf{x}^{(t-1)}) r(\mathbf{x}^{(t)}), \end{aligned} \quad (60)$$

$$\begin{aligned} \tilde{q}(\mathbf{x}^{(t-1)}) \tilde{q}(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)}) \int d\mathbf{x}^{(t-1)} d\mathbf{x}^{(t)} q(\mathbf{x}^{(t-1)}, \mathbf{x}^{(t)}) r(\mathbf{x}^{(t)}) r(\mathbf{x}^{(t-1)}) \\ = \tilde{q}(\mathbf{x}^{(t)}) \tilde{q}(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}) \int d\mathbf{x}^{(t-1)} d\mathbf{x}^{(t)} q(\mathbf{x}^{(t-1)}, \mathbf{x}^{(t)}) r(\mathbf{x}^{(t-1)}) r(\mathbf{x}^{(t)}). \end{aligned} \quad (61)$$

We can now cancel the identical integrals on each side, achieving our goal of showing that the choice of perturbed Markov transitions in Equations 19 and 20 satisfy Equation 18,

$$\tilde{q}(\mathbf{x}^{(t-1)}) \tilde{q}(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)}) = \tilde{q}(\mathbf{x}^{(t)}) \tilde{q}(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}). \quad (62)$$

|  | <i>Gaussian</i>  | <i>Binomial</i>   |
|--|--|---|
| Well behaved (analytically tractable) distribution | $\pi(\mathbf{x}^{(T)}) = \mathcal{N}(\mathbf{x}^{(T)}; \mathbf{0}, \mathbf{I})$  | $\mathcal{B}(\mathbf{x}^{(T)}; 0.5)$  |
| Forward diffusion kernel                           | $q(\mathbf{x}^{(t)}   \mathbf{x}^{(t-1)}) = \mathcal{N}(\mathbf{x}^{(t)}; \mathbf{x}^{(t-1)}, \sqrt{1 - \beta_t}, \mathbf{I}\beta_t)$  | $\mathcal{B}(\mathbf{x}^{(t)}; \mathbf{x}^{(t-1)}(1 - \beta_t) + 0.5\beta_t)$   |
| Reverse diffusion kernel                           | $p(\mathbf{x}^{(t-1)}   \mathbf{x}^{(t)}) = \mathcal{N}(\mathbf{x}^{(t-1)}; \mathbf{f}_\mu(\mathbf{x}^{(t)}, t), \mathbf{f}_\Sigma(\mathbf{x}^{(t)}, t))$  | $\mathcal{B}(\mathbf{x}^{(t-1)}; \mathbf{f}_b(\mathbf{x}^{(t)}, t))$  |
| Training targets                                   | $\mathbf{f}_\mu(\mathbf{x}^{(t)}, t), \mathbf{f}_\Sigma(\mathbf{x}^{(t)}, t), \beta_{1 \dots T}$   | $\mathbf{f}_b(\mathbf{x}^{(t)}, t)$   |
| Forward distribution                               | $q(\mathbf{x}^{(0 \dots T)}) = q(\mathbf{x}^{(0)}) \prod_{t=1}^T q(\mathbf{x}^{(t)}   \mathbf{x}^{(t-1)})$   |   |
| Reverse distribution                               | $p(\mathbf{x}^{(0 \dots T)}) = \pi(\mathbf{x}^{(T)}) \prod_{t=1}^T p(\mathbf{x}^{(t-1)}   \mathbf{x}^{(t)})$   |   |
| Log likelihood                                     | $L = \int d\mathbf{x}^{(0)} q(\mathbf{x}^{(0)}) \log p(\mathbf{x}^{(0)})$  |   |
| Lower bound on log likelihood                      | $K = -\sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}^{(0)}, \mathbf{x}^{(t)})} [D_{KL}(q(\mathbf{x}^{(t-1)}   \mathbf{x}^{(t)}, \mathbf{x}^{(0)})    p(\mathbf{x}^{(t-1)}   \mathbf{x}^{(t)}))] + H_q(\mathbf{X}^{(T)}   \mathbf{X}^{(0)}) - H_q(\mathbf{X}^{(1)}   \mathbf{X}^{(0)}) - H_p(\mathbf{X}^{(T)})$                            |   |
| Perturbed forward diffusion kernel                 | $\tilde{q}(\mathbf{x}^{(t)}   \mathbf{x}^{(t-1)}) = \mathcal{N}\left(\mathbf{x}^{(t)}; \mathbf{x}^{(t-1)}, \sqrt{1 - \beta_t} + \left(\frac{\beta_t}{2}\right)^{\frac{1}{2}} \frac{\partial \log r(\mathbf{x}^{(t)})}{\partial \mathbf{x}^{(t)}}, \mathbf{I}\beta_t\right)$  | $\mathcal{B}\left(x_i^{(t)}; \frac{b_i^t d_i^t}{x_i^t d_i^t + (1 - d_i^t)(1 - d_i^t)}\right)$                           |
| Perturbed reverse diffusion kernel                 | $\tilde{p}(\mathbf{x}^{(t-1)}   \mathbf{x}^{(t)}) = \mathcal{N}\left(\mathbf{x}^{(t-1)}; \mathbf{f}_\mu(\mathbf{x}^{(t)}, t) + \left(\frac{\mathbf{f}_\Sigma(\mathbf{x}^{(t)}, t)}{2}\right)^{\frac{1}{2}} \frac{\partial \log r(\mathbf{x}^{(t-1)})}{\partial \mathbf{x}^{(t)}}, \mathbf{f}_\Sigma(\mathbf{x}^{(t)}, t)\right)$ | $\mathcal{B}\left(x_i^{(t-1)}; \frac{c_i^{t-1} d_i^{t-1}}{x_i^{t-1} d_i^{t-1} + (1 - c_i^{t-1})(1 - d_i^{t-1})}\right)$ |

Table C.1. The key equations in this paper for the specific cases of Gaussian and binomial diffusion processes.  $\mathcal{N}(u; \mu, \Sigma)$  is a Gaussian distribution with mean  $\mu$  and covariance  $\Sigma$ .  $\mathcal{B}(u; r)$  is the distribution for a single Bernoulli trial, with  $u = 1$  occurring with probability  $r$ , and  $u = 0$  occurring with probability  $1 - r$ . Finally, for the perturbed Bernoulli trials  $b_i^t = \mathbf{x}^{(t-1)}(1 - \beta_t) + 0.5\beta_t$ ,  $c_i^t = \lceil \mathbf{f}_b(\mathbf{x}^{(t+1)}, t) \rceil_t$ , and  $d_i^t = r(x_i^{(t)} = 1)$ , and the distribution is given for a single bit  $i$ .

## D. Experimental Details

### D.1. Toy Problems

#### D.1.1. SWISS ROLL

A probabilistic model was built of a two dimensional swiss roll distribution. The generative model  $p(\mathbf{x}^{(0:T)})$  consisted of 40 time steps of Gaussian diffusion initialized at an identity-covariance Gaussian distribution. A (normalized) radial basis function network with a single hidden layer and 16 hidden units was trained to generate the mean and covariance functions  $\mathbf{f}_\mu(\mathbf{x}^{(t)}, t)$  and a diagonal  $\mathbf{f}_\Sigma(\mathbf{x}^{(t)}, t)$  for the reverse trajectory. The top, readout, layer for each function was learned independently for each time step, but for all other layers weights were shared across all time steps and both functions. The top layer output of  $\mathbf{f}_\Sigma(\mathbf{x}^{(t)}, t)$  was passed through a sigmoid to restrict it between 0 and 1. As can be seen in Figure 1, the swiss roll distribution was successfully learned.

#### D.1.2. BINARY HEARTBEAT DISTRIBUTION

A probabilistic model was trained on simple binary sequences of length 20, where a 1 occurs every 5th time bin, and the remainder of the bins are 0. The generative model consisted of 2000 time steps of binomial diffusion initialized at an independent binomial distribution with the same mean activity as the data ( $p(x_i^{(T)} = 1) = 0.2$ ). A multilayer perceptron with sigmoid nonlinearities, 20 input units and three hidden layers with 50 units each was trained to generate the Bernoulli rates  $\mathbf{f}_b(\mathbf{x}^{(t)}, t)$  of the reverse trajectory. The top, readout, layer was learned independently for each time step, but for all other layers weights were shared across all time steps. The top layer output was passed through a sigmoid to restrict it between 0 and 1. As can be seen in Figure 2, the heartbeat distribution was successfully learned. The log likelihood under the true generating process is  $\log_2(\frac{1}{5}) = -2.322$  bits per sequence. As can be seen in Figure 2 and Table 1 learning was nearly perfect.

### D.2. Images

#### D.2.1. ARCHITECTURE

**Readout** In all cases, a convolutional network was used to produce a vector of outputs  $\mathbf{y}_i \in \mathcal{R}^{2J}$  for each image pixel  $i$ . The entries in  $\mathbf{y}_i$  are divided into two equal sized subsets,  $\mathbf{y}^\mu$  and  $\mathbf{y}^\Sigma$ .

**Temporal Dependence** The convolution output  $\mathbf{y}^\mu$  is used as per-pixel weighting coefficients in a sum over time-dependent ‘‘bump’’ functions, generating an output  $\mathbf{z}_i^\mu \in \mathcal{R}$

for each pixel  $i$ ,

$$\mathbf{z}_i^\mu = \sum_{j=1}^J \mathbf{y}_{ij}^\mu g_j(t). \quad (63)$$

The bump functions consist of

$$g_j(t) = \frac{\exp\left(-\frac{1}{2w^2}(t - \tau_j)^2\right)}{\sum_{k=1}^J \exp\left(-\frac{1}{2w^2}(t - \tau_k)^2\right)}, \quad (64)$$

where  $\tau_j \in (0, T)$  is the bump center, and  $w$  is the spacing between bump centers.  $\mathbf{z}^\Sigma$  is generated in an identical way, but using  $\mathbf{y}^\Sigma$ .

For all image experiments a number of timesteps  $T = 1000$  was used, except for the bark dataset which used  $T = 500$ .

**Mean and Variance** Finally, these outputs are combined to produce a diffusion mean and variance prediction for each pixel  $i$ ,

$$\Sigma_{ii} = \sigma\left(z_i^\Sigma + \sigma^{-1}(\beta_t)\right), \quad (65)$$

$$\mu_i = (x_i - z_i^\mu)(1 - \Sigma_{ii}) + z_i^\mu. \quad (66)$$

where both  $\Sigma$  and  $\mu$  are parameterized as a perturbation around the forward diffusion kernel  $T_\pi(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)}; \beta_t)$ , and  $z_i^\mu$  is the mean of the equilibrium distribution that would result from applying  $p(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)})$  many times.  $\Sigma$  is restricted to be a diagonal matrix.

**Multi-Scale Convolution** We wish to accomplish goals that are often achieved with pooling networks – specifically, we wish to discover and make use of long-range and multi-scale dependencies in the training data. However, since the network output is a vector of coefficients for every pixel it is important to generate a full resolution rather than down-sampled feature map. We therefore define multi-scale-convolution layers that consist of the following steps:

1. Perform mean pooling to downsample the image to multiple scales. Downsampling is performed in powers of two.
2. Performing convolution at each scale.
3. Upsample all scales to full resolution, and sum the resulting images.
4. Perform a pointwise nonlinear transformation, consisting of a soft relu ( $\log[1 + \exp(\cdot)]$ ).

The composition of the first three linear operations resembles convolution by a multiscale convolution kernel, up to blocking artifacts introduced by upsampling. This method of achieving multiscale convolution was described in (Baron et al., 2013).

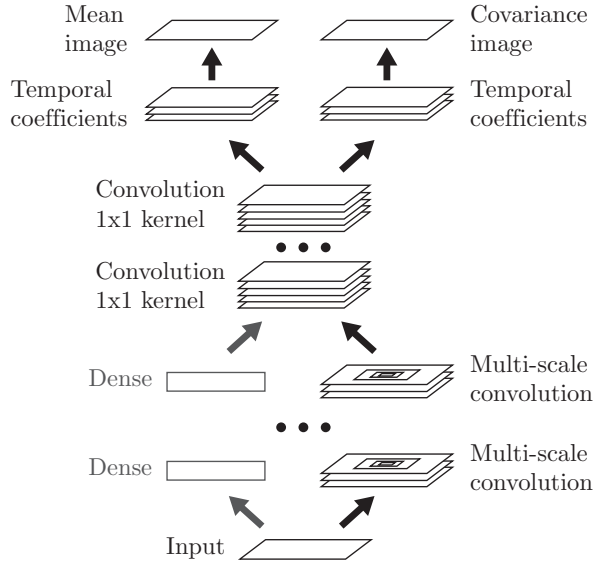


Figure D.1. Network architecture for mean function  $f_\mu(\mathbf{x}^{(t)}, t)$  and covariance function  $f_\Sigma(\mathbf{x}^{(t)}, t)$ , for experiments in Section 3.2. The input image  $\mathbf{x}^{(t)}$  passes through several layers of multi-scale convolution (Section D.2.1). It then passes through several convolutional layers with  $1 \times 1$  kernels. This is equivalent to a dense transformation performed on each pixel. A linear transformation generates coefficients for readout of both mean  $\mu^{(t)}$  and covariance  $\Sigma^{(t)}$  for each pixel. Finally, a time dependent readout function converts those coefficients into mean and covariance images, as described in Section D.2.1. For CIFAR-10 a dense (or fully connected) pathway was used in parallel to the multi-scale convolutional pathway. For MNIST, the dense pathway was used to the exclusion of the multi-scale convolutional pathway.

**Dense Layers** Dense (acting on the full image vector) and kernel-width-1 convolutional (acting separately on the feature vector for each pixel) layers share the same form. They consist of a linear transformation, followed by a tanh nonlinearity.