

Density Ratio Estimation with Conditional Probability Paths



Hanlin Yu¹, Arto Klami¹, Aapo Hyvärinen¹, Anna Korba², Omar Chehab²

1. University of Helsinki, Finland 2. ENSAE, CREST, IP Paris, France



Problem statement

Given samples from two distributions, $X_0 \sim p_0$ and $X_1 \sim p_1$, estimate the ratio $\frac{p_1(\mathbf{x})}{p_0(\mathbf{x})}$.

Algorithm [Choi et al., AISTATS 2022]

1. Interpolate samples: $X_t = \sqrt{1-t^2}X_0(x) + \sqrt{t^2}X_1(x)$. The law $p_t(x)$ is implicit.
2. Estimate the time score $\partial_t \log p_t(x)$.
3. Obtain the log ratio through numerical integration: $\log \frac{p_1(x)}{p_0(x)} = \int_0^1 \partial_t \log p_t(x) dt$

Learning objectives for the time score

Original regression $\mathcal{L}(\theta) = \mathbb{E}_{p(x,t)} [\lambda(t) (\partial_t \log p_t(x) - s_\theta(x, t))^2]$ not explicit

Integrate by parts
TSM $\mathcal{L}(\theta) = 2\mathbb{E}_{p_0(x)}[s_\theta(x, 0)] - 2\mathbb{E}_{p_1(x)}[s_\theta(x, 1)] + \mathbb{E}_{p(t,x)} [2\partial_t s_\theta(x, t) + 2\dot{\lambda}(t)s_\theta(x, t) + \lambda(t)s_\theta(x, t)^2]$ slow to differentiate

Condition (**ours**)
CTSM $\mathcal{L}(\theta) = \mathbb{E}_{p(x,z,t)} [\lambda(t) (\partial_t \log p_t(x|z) - s_\theta(x, t))^2]$ explicit

Factorize (**ours**)
CTSM-v $\mathcal{L}(\theta) = \mathbb{E}_{p(x,z,t)} \left[\lambda(t) \sum_{i=1}^D (\partial_t \log p_t(x^i | x^{<i}, z) - s_\theta^i(x, t))^2 \right]$

We also introduce the weighting function $\lambda(t) \propto 1/|\partial_t \log p_t(x|z)|$.

Theoretical guarantees (modified): for K integration steps and N samples,

$$\mathbb{E}_{\hat{p}_1} \left\| \log \frac{p_1}{p_0} - \widehat{\log \frac{p_1}{p_0}} \right\|_{L^2(p_1)}^2 \leq \underbrace{\frac{1}{2K^2} \mathbb{E}_{p_1(x)} [L(x)^2]}_{\text{integral discretization error}} + \underbrace{\frac{2}{N} e(\theta^*, \lambda, p_t)}_{\text{score estimation error}} + o\left(\frac{1}{N}\right)$$

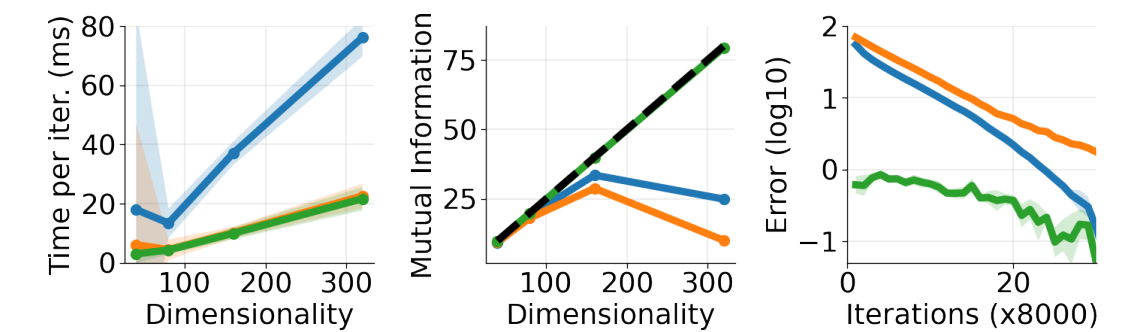
null if $t \rightarrow \partial_t \log p_t(x)$ constant, i.e. Lipschitz constant $L(x)$ is null

null if $\partial_t \log p_t(x|z) = \partial_t \log p_t(x)$

Applications of density-ratio estimation

Mutual information estimation

CTSM-v is faster and outperforms others especially in high dimensions.



— TSM — CTSM — CTSM-v

Space	Methods	Approx. BPD	Time per step
Latent space	TSM	1.30	347 ms
	Ours	1.26	58 ms
Pixel space	TSM	unstable	1103 ms
	Ours	1.03	142 ms

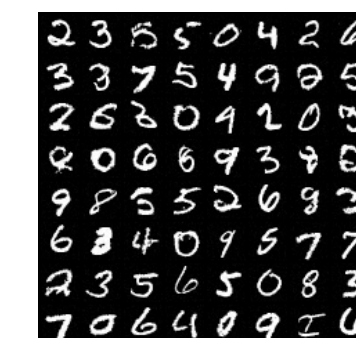
Likelihood estimation (in bits per dimension, BPD). We use

$$\log p_1(x) = \underbrace{\log p_0(x)}_{\text{Known}} + \underbrace{\int_0^1 \partial_t \log p_t(x) dt}_{\text{Estimated}}$$

Sample generation. We convert the estimated time scores into space scores and plug them into popular score-based samplers.

$$\nabla \log p_t(x) = \nabla \left(\underbrace{\log p_0(x)}_{\text{Known}} + \underbrace{\int_0^t \partial_s \log p_s(x) ds}_{\text{Estimated}} \right)$$

Annealed Langevin sampler



Probability flow ODE sampler



Theoretical guarantees: for K integration steps and N samples,

$$\mathbb{E}_{\hat{p}_1}[\text{KL}(p_1, \hat{p}_1)^2] \leq \underbrace{\frac{1}{\gamma K^2} \mathbb{E}_{p_1(x)}[L(x)^2]}_{\substack{\text{integral discretization error} \\ \text{null if } \partial_t \log p_t(\mathbf{x}) \text{ constant wrt } t}} + \underbrace{\frac{2}{N} e(\theta^*, \lambda, p)}_{\substack{\text{score estimation error} \\ \text{null if } \partial_t \log p_t(\mathbf{x} | \mathbf{z}) = \partial_t \log p_t(\mathbf{x})}} +$$

NCE

	Binary	Multi-class	Continuous
NCE	Yes	No	No
TRE	No	Consecutive	No
DRE-infinity	No	Consecutive	Yes
Michael	No	Full	No
Hanlin	No	Full	Yes

The literature suggests that multi-class (full) and continuous would be the thing to aim for.

CNCE

	Binary	Multi-class	Continuous
CNCE	Yes	No	No
No-one	No	Consecutive	No
No-one	No	Consecutive	Yes
Us	No	Full	No
Us	No	Full	Yes

$$p(Y = 1 \mid x) = \frac{p_1(x)}{p_0(x) + p_1(x)} = \frac{1}{1 + \frac{p_0}{p_1}(x)}$$

$$p(Y = k | x) = \frac{p_k(x)}{\sum_{j=1}^M p_j(x)} = \frac{1}{1 + \sum_{j \neq k}^M \frac{p_j}{p_k}(x)}$$

Similar to mean flow

$$= \frac{1}{1 + \sum_{j \neq k}^M \frac{\frac{p_j}{p_1}(x)}{\frac{p_k}{p_1}(x)}}$$

Similar to flow map

$$\log \frac{p(x', t')}{p(x, t)} = \log \frac{p(x' | t')}{p(x | t)} = \log \frac{p(x_{t'})}{p(x_t)} = - \int_t^{t'} \nabla \cdot v(s) ds$$

OLDER VERSIONS

$$\mathcal{L}_{\text{TSM}}(\boldsymbol{\theta}) = \mathbb{E}_{p(\mathbf{x}, t)} \left[\lambda(t) \|\underbrace{\partial_t \log p_t(\mathbf{x})}_{\text{Hard to compute.}} - s_{\boldsymbol{\theta}}(\mathbf{x}, t) \|^2 \right]$$

}

$$\mathcal{L}_{\text{CTSM}}(\boldsymbol{\theta}) = \mathbb{E}_{p(\mathbf{x}, \mathbf{z}, t)} \left[\lambda(t) \|\underbrace{\partial_t \log p_t(\mathbf{x} | \mathbf{z})}_{\text{Easy to compute.}} - s_{\boldsymbol{\theta}}(\mathbf{x}, t) \|^2 \right]$$

Possible sellings points / research directions

1. Link between multi-class NCE variants and modern papers (mean flows, flow map matching) that parameterize the jump between any two times t, t' . Link with ST-NCE.

2. Focusing the “spatial” part of ST-NCE.

NCE has in very recent years benefitted from a lot of modernization (TRE, DRE-infinity, Srivastava). Applying this modernization to CNCE is very unexplored yet!!

Presumably, binary CNCE estimates decently the EBM on the data manifold, but maybe the multi-class estimates the EBM much better *outside* of the data manifold. This OOD is of interest to many (Yilun, Florentin).

Psychologically, people love scores. The link between CNCE and **space score** matching seems very under-explored yet.

More generally, solving OOD estimating of the space score is very relevant now for composing EBMs.

3. ST-NCE.

Density Ratio Estimation with Conditional Probability Paths



Hanlin Yu¹, Arto Klami¹, Aapo Hyvärinen¹, Anna Korba², Omar Chehab²

1 University of Helsinki, Finland 2. ENSAE, CREST, IP Paris, France



Problem statement

We have samples from two distributions, $X_0 \sim p_0$ and $X_1 \sim p_1$ and want to estimate their density ratio $\frac{p_1(x)}{p_0(x)}$ [1][2][3].

Time Score Matching (TSM) [3]

1. Interpolate samples: $X_t = \sqrt{1-t^2}X_0(x) + \sqrt{t^2}X_1(x)$. Its probability law $p_t(x)$ is implicit.
2. Estimate the time score $\partial_t \log p_t(x)$.
3. Obtain the log ratio through numerical integration: $\log \frac{p_1(x)}{p_0(x)} = \int_0^1 \partial_t \log p_t(x) dt$

Step 2 is currently **slow**: it involves minimizing a loss with higher order gradients (against time and parameters). Can we obtain a faster variant?

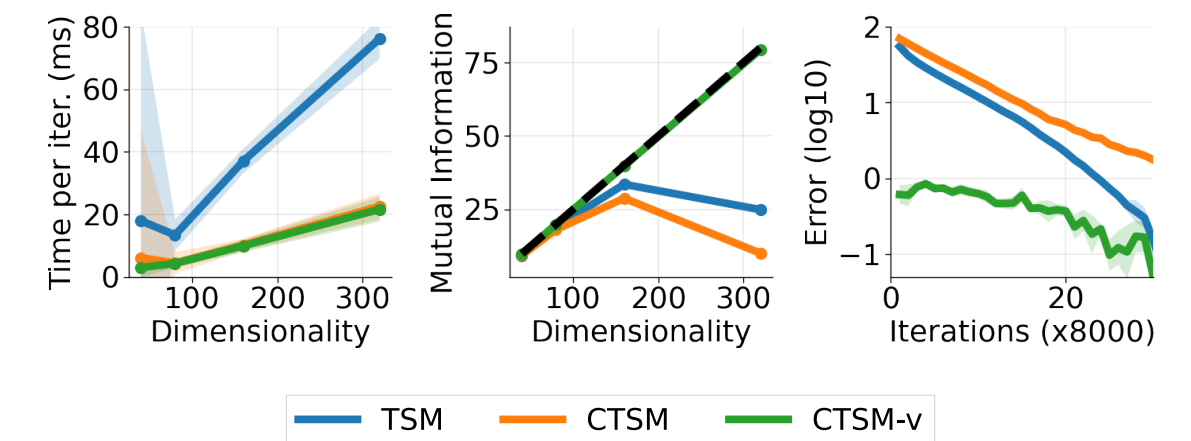
Conditional Time Score Matching (CTSM)

We propose a cheaper version of step 2, using a conditioning variable z to make the intermediate $p_t(x|z)$ known in closed-form, similar to diffusion [4] and flow matching [5]: $\mathcal{L}(\theta) = \mathbb{E}_{p(x,z,t)} [\lambda(t) (\partial_t \log p_t(x|z) - s_\theta(x, t))^2]$, $\lambda(t)$ are positive weights.

We also propose an efficient, **vectorized** version of this loss (CTSM-v), using instead vector of $\partial_t \log p_t(x^i | x^{<i}, z)$, as well as **theoretical guarantees** on the estimation errors.

Mutual information estimation

Can be reframed as a density ratio estimation problem. CTSM-v is faster and outperforms others especially in high dimensions.



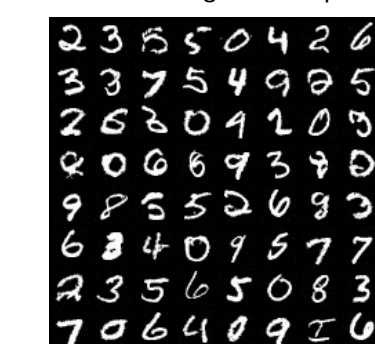
MNIST

Likelihood estimation (in bits per dimension, BPD). Our method is faster and more accurate, achieving good results without relying on pre-trained normalizing flows.

Sample generation. Our trained score network generates reasonable samples using popular diffusion-based samplers with ambient space scores induced by time scores..

Space	Methods	Approx. BPD	Time per step
Latent	TSM	1.30	347 ms
	Ours	1.26	58 ms
Ambient	TSM	unstable	1103 ms
	Ours	1.03	142 ms

Annealed Langevin sampler



PF-ODE sampler



References

1. Gutmann, M.U. and Hyvärinen, A. Noise-contrastive estimation of unnormalized statistical models, with applications to natural image statistics, JMLR 2012
2. Rhodes, B., Xu., K. and Gutmann, M.U., Telescoping Density-Ratio Estimation, NeurIPS 2020
3. Choi, K., Meng, C., Song, Y., and Ermon, S., Density ratio estimation via infinitesimal classification, AISTATS 2022
4. Song, Y., Sohl-Dickstein, J., Kingma, D.P., Kumar, A., Eamon, S., and Poole, B., Score-based generative modeling through stochastic differential equations, ICLR 2021
5. Lipman, Y., Chen, R. T.Q., Ben-Hamu, H., Nickel, M., and Le, M., Flow matching for generative modeling, ICLR 2023

Density Ratio Estimation with Conditional Probability Paths



Hanlin Yu¹, Arto Klami¹, Aapo Hyvärinen¹, Anna Korba², Omar Chehab²

1 University of Helsinki, Finland 2. ENSAE, CREST, IP Paris, France



Problem statement

We have samples from two distributions, $X_0 \sim p_0$ and $X_1 \sim p_1$ and want to estimate their density ratio $\frac{p_1(\mathbf{x})}{p_0(\mathbf{x})}$.

Can be solved using probabilistic classification [1][2][3].

Time Score Matching (TSM) [3]

1. Interpolate samples: $X_t = \sqrt{1-t^2}X_0(\mathbf{x}) + tX_1(\mathbf{x})$. Its probability law $p_t(\mathbf{x})$ is *implicit*.
2. Estimate the time score $\partial_t \log p_t(\mathbf{x}) \approx s_\theta(\mathbf{x}, t)$ using the TSM loss:

$$\mathcal{L}(\theta) = 2\mathbb{E}_{p_0(\mathbf{x})}[s_\theta(\mathbf{x}, 0)] - 2\mathbb{E}_{p_1(\mathbf{x})}[s_\theta(\mathbf{x}, 1)] + \mathbb{E}_{p(t, \mathbf{x})}[2\partial_t s_\theta(\mathbf{x}, t) + 2\dot{\lambda}(t)s_\theta(\mathbf{x}, t) + \lambda(t)s_\theta(\mathbf{x}, t)^2],$$
 Where $\lambda(t)$ are positive weights.

3. Obtain the ratio through numerical integration: $\frac{p_1(\mathbf{x})}{p_0(\mathbf{x})} = \exp\left(\int_0^1 s_\theta(\mathbf{x}, t)dt\right)$

Step 2 involves **higher order gradients** and thus is expensive. Can we obtain a cheaper variant with higher accuracies?

Conditional Time Score Matching (CTSM)

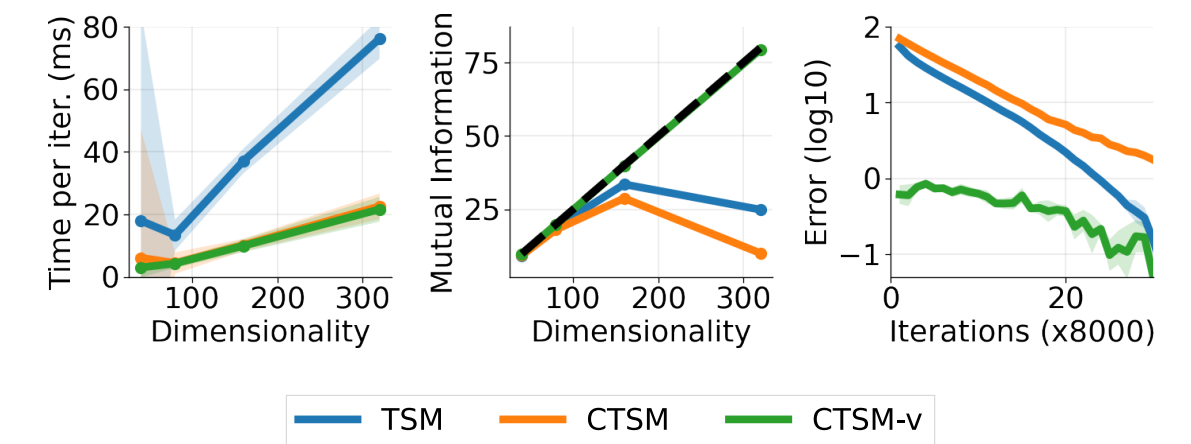
We propose a cheaper version of step 2. The idea is to introduce a conditioning variable \mathbf{z} to make the intermediate $p_t(\mathbf{x} | \mathbf{z})$ *known in closed-form*, similar to diffusion [4] and flow matching [5]:

$$\mathcal{L}(\theta) = \mathbb{E}_{p(\mathbf{x}, \mathbf{z}, t)}[\lambda(t)(\partial_t \log p_t(\mathbf{x} | \mathbf{z}) - s_\theta(\mathbf{x}, t))^2].$$

We also propose an efficient, **vectorized** version of this loss (CTSM-v), using instead vector of $\partial_t \log p_t(x^i | \mathbf{x}^{<i}, \mathbf{z})$, as well as **theoretical guarantees**.

Mutual information estimation

Can be reframed as a density ratio estimation problem. CTSM-v is faster and outperforms others especially in high dimensions.



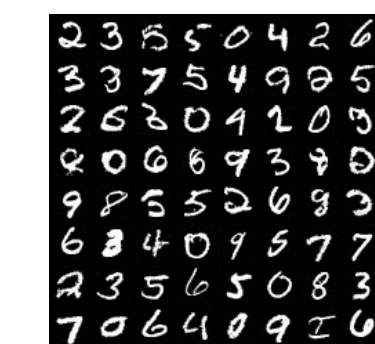
MNIST

Likelihood estimation (BPD, lower is better). Our method is faster and more accurate, achieving good results without relying on pre-trained normalizing flows.

Sample generation. Our trained score network generates reasonable samples using popular diffusion-based samplers with ambient space scores.

Space	Methods	Approx. BPD	Time per step
Latent	TSM	1.30	347 ms
	Ours	1.26	58 ms
Ambient	TSM	unstable	1103 ms
	Ours	1.03	142 ms

Annealed Langevin sampler



PF-ODE sampler



References

1. Gutmann, M.U. and Hyvärinen, A. Noise-contrastive estimation of unnormalized statistical models, with applications to natural image statistics, JMLR 2012
2. Rhodes, B., Xu., K. and Gutmann, M.U., Telescoping Density-Ratio Estimation, NeurIPS 2020
3. Choi, K., Meng, C., Song, Y., and Ermon, S., Density ratio estimation via infinitesimal classification, AISTATS 2022
4. Song, Y., Sohl-Dickstein, J., Kingma, D.P., Kumar, A., Eamon, S., and Poole, B., Score-based generative modeling through stochastic differential equations, ICLR 2021
5. Lipman, Y., Chen, R. T.Q., Ben-Hamu, H., Nickel, M., and Le, M., Flow matching for generative modeling, ICLR 2023

Density Ratio Estimation with Conditional Probability Paths



Hanlin Yu¹, Arto Klami¹, Aapo Hyvärinen¹, Anna Korba², Omar Chehab²

1. University of Helsinki, Finland 2. ENSAE, CREST, IP Paris, France



Problem statement

Given samples from two distributions, $X_0 \sim p_0$ and $X_1 \sim p_1$, estimate their density ratio $\frac{p_1(x)}{p_0(x)}$.

Density Ratio Estimation using a time score [Choi et al., AISTATS 2022]

1. Interpolate samples: $X_t = \sqrt{1-t^2}X_0(x) + \sqrt{t^2}X_1(x)$. Its probability law $p_t(x)$ is implicit.
2. Estimate the time score $\partial_t \log p_t(x)$.
3. Obtain the log ratio through numerical integration: $\log \frac{p_1(x)}{p_0(x)} = \int_0^1 \partial_t \log p_t(x) dt$

Our loss for learning the time score using positive weights $\lambda(t)$

Classical regression loss (Choi et al.)

$$\mathcal{L}(\theta) = \mathbb{E}_{p(x,t)} [\lambda(t) \|\partial_t \log p_t(x) - s_t(x; \theta)\|^2]$$

No explicit formula for $p_t(x)$.

Integration by parts (Choi et al.)

$$\mathcal{L}(\theta) = 2\mathbb{E}_{p_0(x)}[s_\theta(x,0)] - 2\mathbb{E}_{p_1(x)}[s_\theta(x,1)] + \mathbb{E}_{p(t,x)} [2\partial_t s_\theta(x,t) + 2\dot{\lambda}(t)s_\theta(x,t) + \lambda(t)s_\theta(x,t)^2]$$

*The loss gradient computes $\partial_\theta \partial_t s_\theta(x,t)$ which is **slow**.*

Conditioning variable (ours): like diffusion / flow

$$\mathcal{L}(\theta) = \mathbb{E}_{p(x,z,t)} [\lambda(t) \|\partial_t \log p_t(x|z) - s_\theta(x,t)\|^2]$$

*Conditioning variable z makes $p_t(x|z)$ **known in closed-form**.*

Efficient implementation of our loss

Reweigh the loss: $\lambda(t) \propto 1/\|\partial_t \log p_t(x|z)\|$

Vectorize the loss: replace $\partial_t \log p_t(x|z) \in \mathbb{R}$ by the vector of $\partial_t \log p_t(x^i|x^{<i},z)$

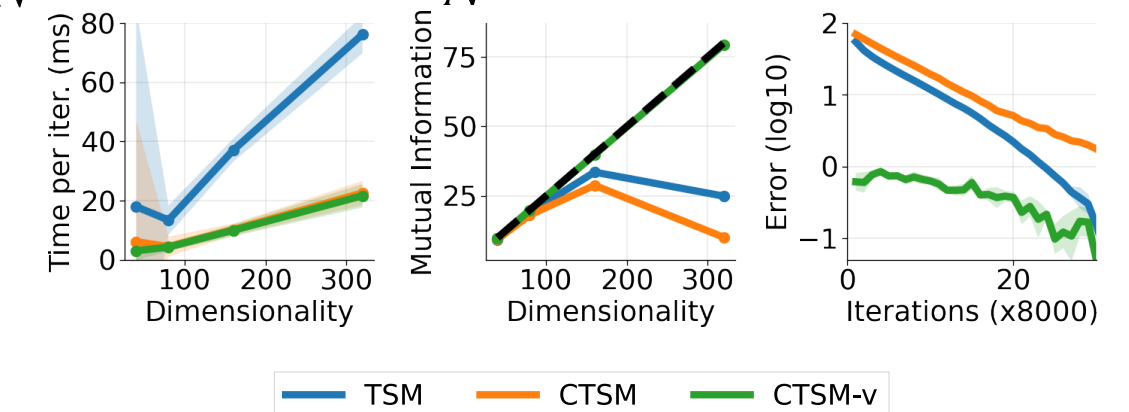
Our theoretical guarantees control the expected error as time discretization steps and number of samples increase.

$$\mathbb{E}_{\hat{p}_1}[\text{KL}(p_1, \hat{p}_1)^2] \leq \frac{1}{2K^2} \mathbb{E}_{p_1(x)}[L(x)^2] + \frac{2}{N} e(\theta^*, \lambda, p) + o\left(\frac{1}{N}\right)$$

Applications of density-ratio estimation

Mutual information estimation

CTSM-v is faster and outperforms others especially in high dimensions.

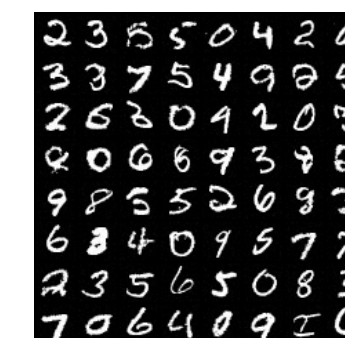


Likelihood estimation (in bits per dimension, BPD). Higher likelihoods and faster training directly in pixel space.

Sample generation. We convert the estimated time scores $\partial_t \log p_t(x)$ into space scores $\nabla \log p_t(x)$ and plug them into popular score-based samplers.

Space	Methods	Approx. BPD	Time per step
Latent	TSM	1.30	347 ms
	Ours	1.26	58 ms
Ambient	TSM	unstable	1103 ms
	Ours	1.03	142 ms

Annealed Langevin sampler



PF-ODE sampler

