

---

# Can We Estimate The Entropy Of Distributions Known Up To A Normalization Constant?

---

## Abstract

1 Computing the differential entropy for distributions known up to a normaliza-  
2 tion constant is a challenging problem with significant theoretical and practical  
3 applications. Variational inference is widely used for scalable approxima-  
4 tion of densities from samples, but is under-explored when *only unnormalized*  
5 *densities are available*. Messaoud et al. [2024] introduced **P-SVGD**, a particle-  
6 based variational method using Stein Variational Gradient Descent. However,  
7 we show that **P-SVGD** scales poorly to high-dimensional spaces. We propose  
8 **MET-SVGD**, an extension of **P-SVGD** that scales efficiently with convergence  
9 guarantees. Our method achieves SOTA results on scaling SVGD. We significantly  
10 outperform **P-SVGD** on entropy estimation, Maximum Entropy Reinforcement  
11 Learning, and image generation with Energy-Based Models benchmarks. Code:  
12 <https://tinyurl.com/2esyfx8j>.

## 1 Introduction

14 The differential entropy [Cover, 1999, Shannon, 2001] of a  $d$ -dimensional random variable  $X$  with  
15 a probability density function  $p(x) = \bar{p}(x)/Z$  is  $\mathcal{H}(p) = -\mathbb{E}_{x \sim p(x)}[\log p(x)] = -\int p(x) \log p(x) dx$ ,  
16 with  $Z = \int \bar{p}(x) dx$  being the normalization constant. The entropy plays a central role in machine  
17 learning [Rubinstein and Kroese, 2004]. While significant progress has been made on estimating  
18 entropies from samples [Beirlant et al., 1997, Paninski, 2003], settings where only the unnormalized  
19 density is available remain under-explored. Recently, Messaoud et al. [2024] introduced Parametrized  
20 Stein Variational Gradient Descent (**P-SVGD**), which leverages Stein Variational Gradient Descent  
21 (SVGD) sampler [Liu and Wang, 2016] to construct a variational distribution  $q^L$  from  $\bar{p}$ . SVGD  
22 updates a set of interacting particles  $\{x_i\}_{i=1}^M$  using a deterministic velocity field  $\phi(\cdot)$  that balances a  
23 gradient force and a repulsive one:

$$\phi(x_i^l) = \mathbb{E}_{x_j^l \sim q^l} \left[ \kappa(x_i^l, x_j^l) \nabla_{x_j^l} \log p(x_j^l) + \nabla_{x_j^l} \kappa(x_i^l, x_j^l) \right], \quad (1)$$

24 following the update rule  $x_i^{l+1} = x_i^l + \epsilon \phi(x_i^l)$ .  $\epsilon$  is the step-size,  $q^l$  is the particles distribution at  
25 step  $l \in [1, L]$  and  $\kappa(\cdot, \cdot)$  is typically an RBF kernel, *i.e.*,  $\kappa(x_i, x_j) = \exp(-\|x_i - x_j\|^2 / 2\sigma^2)$ , with  
26 bandwidth  $\sigma^2$  set heuristically as the median of the particles squared differences. **P-SVGD** derives a  
27 closed form expression of  $q^l$  at every step  $l$  including the last step  $L$ :

$$\log q^L(x^L) = \log q^0(x^0) - \epsilon \sum_{l=0}^{L-1} \mathbb{E}_{x_j^l \sim q^l} \left[ \text{Tr} \left( \partial \bar{\phi}(x^l, x_j^l) / \partial x^l \right) \right] + \mathcal{O}(\epsilon^2), \quad (2)$$

28 where  $\bar{\phi}(x^l, x_j^l)$  is the contribution of  $x_j^l$  to the update of particle  $x^l$  and  $\phi(x^l) = \mathbb{E}_{x_j^l} [\bar{\phi}(x^l, x_j^l)]$  is  
29 defined in Eq.1. The trace term is approximated using only first-order derivatives and, for efficiency,  
30 the authors omit a trace-of-Hessian term  $\text{Tr}(\nabla_{x^l}^2 \log p(x^l))$  by sampling  $x_j^l \neq x^l$ :

$$\log q^L(x_i^L) = \log q^0(x_i^0) - \epsilon \sum_{l=0}^{L-1} \sum_{i \neq j=0}^{M-1} \frac{\kappa(x_j^l, x_i^l)}{M\sigma^2} \left( d - \frac{\|x_i^l - x_j^l\|^2}{\sigma^2} - (x_i^l - x_j^l)^\top \nabla_{x_j^l} \log p(x_j^l) \right) + \mathcal{O}(\epsilon^2). \quad (3)$$

31  $q^L$  can approximate a broad class of densities under mild assumptions [Villani et al., 2009], enabling  
32 accurate estimation of  $\mathcal{H}(p)$  with compelling results in MaxEntr RL. Despite its promise, we show  
33 that **P-SVGD** has several limitations including sensitivity to SVGD hyperparameters, mode collapse,  
34 poor convergence to non-smooth targets and limited scalability in high-dimensional spaces (Fig. 1).

35 To address these challenges, we introduce **MET-SVGD**, an extension of **P-SVGD** that scales to  
36 high-dimensional distributions with improved accuracy and convergence guarantees. **MET-SVGD**  
37 learns step-wise kernel bandwidths and step-sizes, replacing the non-robust heuristics of **P-SVGD**;  
38 upgrades its local invertibility condition to a principled global one that satisfies the assumptions in  
39 the entropy derivation; consolidates two independent, informal step-size constraints into a unified,  
40 principled rule; efficiently restores the missing Hessian term via Hutchinson’s estimator; incorporates  
41 a Metropolis–Hastings correction step to ensure asymptotic convergence; and adaptively tunes the  
42 number of sampling steps to better handle complex, high-dimensional distributions.

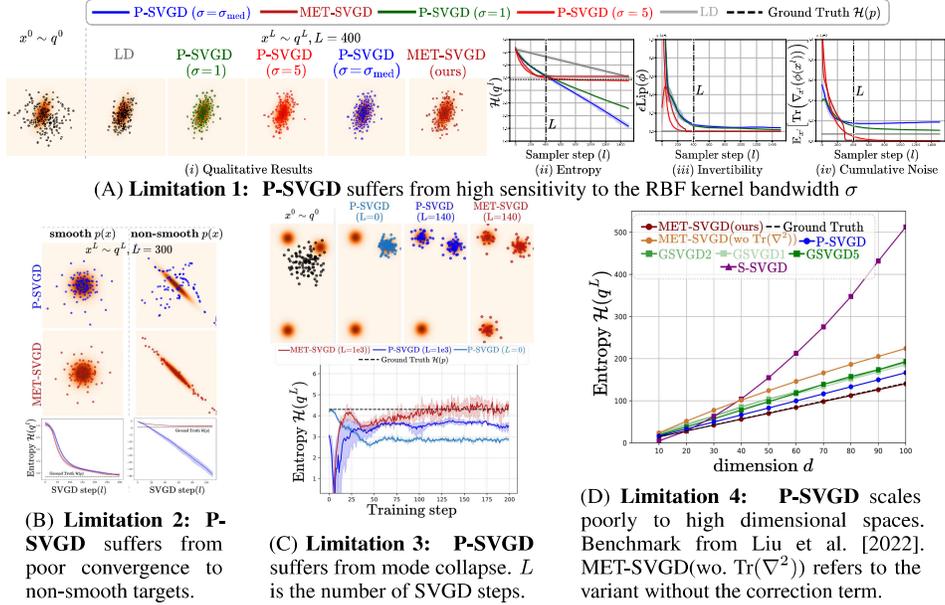


Figure 1: **P-SVGD** limitations.

43 **MET-SVGD** achieves SOTA performance on SVGD scalability benchmarks. Additionally, it significantly outperforms **P-SVGD** on image generation (20 vs 88 FID) and MaxEnr RL tasks (2.3% and 44 1.5% better final returns on Walker2d-v2 and Humanoid-v2 [Brockman et al., 2016], respectively). 45 By bridging variational inference, SVGD, and MH methods, **MET-SVGD** sets a novel framework for 46 entropy estimation from unnormalized densities and scalable sampling. 47

## 48 2 Approach

49 Similarly to **P-SVGD**, **MET-SVGD** is a VI-based method for computing the entropy of distributions 50  $p$  known up to a normalization constant, *i.e.*, it approximates  $p$  with a tractable, sample-efficient 51 distribution and estimates  $\mathcal{H}(p)$  using the entropy of this distribution. **MET-SVGD** introduces a 52 series of optimizations to address **P-SVGD**'s key limitations as illustrated in Fig. 2B: [L<sub>1</sub>-F<sub>1</sub>] **P-** 53 **SVGD** introduces two informal independent constraints on the step-size including a local invertibility 54 one, although CVF requires global invertibility; **MET-SVGD** unifies these constraints into a single 55 principled one satisfying global invertibility (Sec. 2.1). [L<sub>2</sub>-F<sub>2</sub>] **P-SVGD** exhibits high sensitivity to 56 hyperparameters (Fig. 1A-iii) with no tuning guidelines; we show that this is due to the accumulation 57 of noise in the trace term of Eq. 2, leading to entropy divergence and mitigates this by learning the 58 SVGD hyperparameters via reverse KL minimization (Sec. 2.2). [L<sub>3</sub>-F<sub>3</sub>] **P-SVGD** suffers from 59 poor convergence to non-smooth targets and sampling mode collapse (Fig. 1B and Fig. 1C), due 60 to its divergence control heuristic and the absence of convergence guarantees in the finite particle 61 regime; **MET-SVGD** replaces this heuristic with a MH step, guaranteeing asymptotic convergence 62 independently from the number of particles (Sec. 2.4). [L<sub>4</sub>-F<sub>4</sub>] **P-SVGD**'s omission of the trace- 63 of-Hessian term limits its scalability to high-dimensions (Fig. 1D) which **MET-SVGD** efficiently 64 restores as explained in Sec. 2.3. [L<sub>5</sub>-F<sub>5</sub>] **P-SVGD** uses a fixed number of steps ( $L$ ), which may be 65 insufficient for convergence; **MET-SVGD** determines  $L_c$  using the Stein Identity (Sec. ??).

### 66 2.1 Conditions On The SVGD Step-Size For Invertibility and log-det Approximation

67 In **P-SVGD**, Eq. 3 was derived by (1) leveraging the CVF (App. 5.2) assuming invertibility, and 68 (2) approximating the log-det term in the CVF with an efficient trace one. These steps introduce 69 two conditions on the SVGD step-size: (1)  $\epsilon \ll \sigma$  and (2)  $\epsilon \|\nabla_{x^l} \phi(x^l)\|_\infty \ll 1$ . However, these 70 conditions present two major issues: (1) Both are informal (use of  $\ll$ ); in practice,  $\epsilon$  is simply set to 71 an arbitrarily small value, hoping that both constraints hold, which may not be true and often results 72 in more steps than necessary. (2) The step-size condition only guarantees *local* invertibility, whereas 73 the CVF requires *global* invertibility. To address this, we extend a sufficient condition for invertible 74 residual networks (Behrmann et al. [2019]) to SVGD. We also derive a precise condition on  $\epsilon$  for the 75 log-det approximation (Proposition 2.2) and unify both into a single efficient bound (Corollary 2.3). 76 **Proposition 2.1** (Sufficient condition for global SVGD invertibility). *Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$  with  $f =$  77  $(f^1 \circ \dots \circ f^L)$  denote a sequence of SVGD updates with  $f^l = I + \epsilon \phi^l$ . We denote by  $\text{Lip}(\phi^l)$  the 78 Lipschitz constant of the velocity  $\phi^l$  at step  $l$ .  $f$  is invertible if:  $\epsilon \text{Lip}(\phi^l) < 1$ , for all  $l \in [0, L - 1]$ .* 79

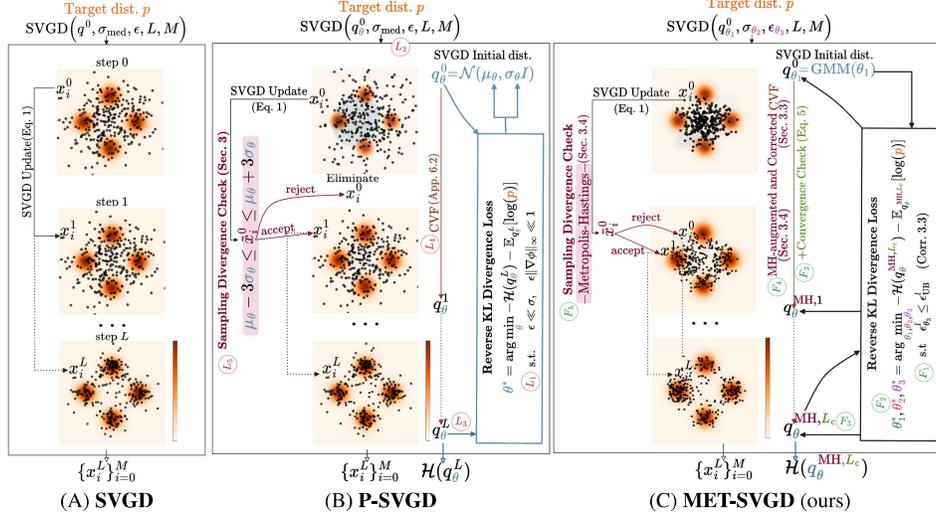


Figure 2: Novelty over P-SVGD. Limitations (L) and corresponding Fixes (F).

80 The proof is in App. 7.1. Using the mean value theorem (App.5.6), we estimate this Lipschitz constant  
 81 as  $\|\nabla\phi^l\|_2 = \max_{x^l} \|\nabla_{x^l}\phi(x^l)\|_2$  with  $\|\cdot\|_2$  being the spectral norm (largest singular value).

82 **Proposition 2.2** (Sufficient condition for log-det Approximation). *Let  $\phi^l: \mathbb{R}^d \rightarrow \mathbb{R}^d$ ,  $\log|\det(I +$   
 83  $\epsilon\nabla\phi^l)| = \epsilon\text{Tr}(\nabla\phi^l)$  if  $\epsilon|\lambda_{\max}(\nabla\phi^l)| < 1$  for all  $l \in [0, L-1]$ , with  $\lambda_{\max}$  being the largest eigenvalue  
 84 value and  $\nabla$  the gradient operator w.r.t the input.*

85 **Corollary 2.3.** *The distribution induced by the SVGD update (Eq. 1) using an RBF kernel is given  
 86 by Eq. 3 if  $\epsilon < \epsilon_{UB}^l = 1/\sup_x \sqrt{\text{Tr}(\nabla\phi^l(x)\nabla\phi^{l,T}(x))}$  for all  $l \in [0, L-1]$ .*

87 *Proof Sketch:* Given  $A \in \mathbb{R}^{d \times d}$ , the following always holds:  $|\lambda_{\max}(A)| \leq \|A\|_2 \leq \sqrt{\text{Tr}(AA^T)}$ .  
 88 Proof is in App. 7.3, where we also show that  $\text{Tr}(AA^T)$  can be efficiently computed using only  
 89 first-order derivatives and vector dot products, making this condition practical.

## 90 2.2 Optimized SVGD Parameters

91 A major drawback of P-SVGD is its high sensitivity to the RBF kernel bandwidth  $\sigma$ , which Messaoud  
 92 et al. [2024] attribute to violation of the invertibility of the SVGD update rule (Eq. 1): In a 2- $d$   
 93 Gaussian target setup (reproduced in Fig. 1A), they show that, paradoxically, although SVGD and  
 94 Langevin Dynamics (LD) (update rule in App. 5.1) *qualitatively* converge to the target, *i.e.*, the  
 95 particles reach high-density regions (Fig. 1A-i), the entropy estimate only converges for specific  
 96  $\sigma$  values, *e.g.*,  $\sigma = 5$  (Fig. 1A-ii). The authors hypothesise that this is due to LD being inherently  
 97 non-invertible and SVGD being invertible only for certain  $\sigma$  values. This is incorrect: in Fig. 1A-  
 98 iii, we show that the step-size condition from Corollary 2.3 is always satisfied. Instead we show  
 99 that the poor quantitative convergence of  $\mathcal{H}(q^L)$  to  $\mathcal{H}(p)$  arises from the cumulative residual noise  
 100 in the trace term of Eq. 2, *i.e.*,  $\mathbb{E}_{x^l \sim q^l}[\text{Tr}(\nabla_{x^l}\phi(x^l))] \rightarrow 0$  as  $l \rightarrow \infty$  (Fig. 1A-iv), resulting in a  
 101 quasi-linear growth in the entropy with the number of steps (Fig. 1A-ii). To address this, we propose  
 102 leveraging the closed-form expression of  $q_{\theta_2}^{L_c}$  to learn a step-wise kernel bandwidth  $\sigma_{\theta_2}^l$  and step-size  
 103  $\epsilon_{\theta_3}^l$  alongside the initial distribution  $q_{\theta_1}^0$  by minimizing the reverse KL-divergence:

$$\theta^* = \arg \min_{\theta} \mathbb{E}_{x^{L_c} \sim q_{\theta}^{L_c}} [\log q_{\theta}^{L_c}(x^{L_c}) - \log p(x^{L_c})], \quad \text{s.t.} \quad \epsilon_{\theta_3}^l \leq \epsilon_{UB}^l \quad \forall l \in [0, L_c - 1],$$

104 with  $\epsilon_{UB}^l$  being the upper-bound from Corollary 2.3 and  $\theta = \{\theta_i\}_{i=1}^3$ . Besides, we derive an efficient  
 105 convergence check based on the Stein Identity (SI) enabling an adaptive number of steps  $L_c$ , rather  
 106 than fixing it a priori:  $\text{SI}(q_{\theta}^l, p) = \sqrt{\mathbb{E}_{x^l} [\phi_{\theta}(x^l)^T \nabla_{x^l} \log p(x^l) + \text{Tr}(\nabla_{x^l} \phi_{\theta}(x^l))]}$ . Note that the  
 107 expression above only depends on  $\text{Tr}(\nabla_{x^l} \phi(x^l))$  which can be computed efficiently (App. 8.3).

## 108 2.3 Corrected Derivation of $q_{\theta}^l$

109 As explained in Sec. 2, P-SVGD approximate the expectation over particles in Eq. 2 by excluding  
 110 the updated particle itself ( $x^l \neq x_j^l$ ), so that the trace term can be estimated using only first-order  
 111 derivatives and thereby avoiding explicit Hessian computation. While this approximation is valid  
 112 asymptotically, it breaks in the finite-particle regime and translates into inconsistent updates across  
 113 particles ( $i=j$  term is missing), leading to particles following different distributions, and making  
 114 the entropy ill-defined. This is a key source of P-SVGD's poor scalability as shown in Fig. 1D  
 115 (orange vs brown). To address this, MET-SVGD adds the missing term to the entropy using (1) the  
 116 Hutchinson estimator [Hutchinson, 1989] and (2) the double differentiation trick [Song et al., 2020]:

117  $Tr(\nabla_{x_i}^2 \log p(x_i)) \stackrel{(1)}{=} \mathbb{E}_{v \sim p_v} [v^T \nabla_{x_i}^2 \log p(x_i) v] \stackrel{(2)}{=} \mathbb{E}_{v \sim p_v} [\nabla_{x_i} (v^T \nabla_{x_i} \log p(x_i)) v]$ , where  $p_v$  is chosen such  
 118 that  $\mathbb{E}[v] = 0$  and  $\mathbb{E}[vv^T] = I$  (e.g.,  $p_v$  is the Radamacher distr). Importantly, SVGD is less sensitive  
 119 to trace approximation errors compared to other MCMC methods (e.g., LD) as shown in Fig. 12.  
 120 Notably, the trace term in SVGD is scaled by the number of particles  $M$ :

$$\log q_\theta^L(x^L) = \log q_{\theta_1}^0(x^0) + \epsilon_{\theta_3} \sum_{l=0}^{L-1} \sum_{x^l \neq x^l} Tr \left( \frac{\partial \bar{\phi}_\theta(x^l, x_j^l)}{\partial x^l} \right) + \frac{\epsilon_{\theta_3}^l}{MV} \sum_{v=0}^{V-1} \nabla_{x^l} (v^T \nabla_{x^l} \log p(x^l)) v,$$

121 unlike LD:  $\log q_\theta^L(x^L) = \log q_{\theta_1}^0(x^0) + (\epsilon_{\theta_3}^l/V) \sum_{l=0}^{L-1} \sum_{v=0}^{V-1} \nabla_{x^l} (v^T \nabla_{x^l} \log p(x^l)) v$  (App. 7.8).

## 122 2.4 Divergence Control via Metropolis Hastings

123 To prevent divergence during sampling due to steepness in the target, **P-SVGD** introduces a heuristic  
 124 that removes particles deviating beyond a fixed number of standard deviations from the mean of  
 125 the initial Gaussian distribution  $q_{\theta_1}^0$  (Fig. 2B). This heuristic, however, exacerbates mode collapse  
 126 by discouraging exploration of distant modes (Fig. 1C). Instead, we propose a more principled  
 127 MH-based divergence control Robert et al. [2004]. After each update, the proposed position  $\tilde{x}^l =$   
 128  $x^{l-1} + \epsilon_{\theta_3} \phi_\theta(x^{l-1})$  is accepted with probability  $\alpha_\theta^l$  (i.e.,  $x^l = \tilde{x}^l$ ), otherwise the old position is  
 129 retained (i.e.,  $x^l = x^{l-1}$ ). We compute  $\alpha_\theta^l$  efficiently by leveraging  $Tr(\nabla_{x^l} \phi_\theta(x^l))$ :  $\log \alpha_\theta^l =$   
 130  $\min[0, \log \bar{p}(\tilde{x}^l) - \log \bar{p}(x^{l-1}) + \epsilon_{\theta_3} Tr(\nabla_{x^l} \phi_\theta(x^l))]$  (proof in App. 8). **MET-SVGD** is an MH  
 131 with an efficient SVGD-based proposal distribution. It therefore inherits asymptotic **convergence**  
 132 **guarantees** from MH (details in App. 6.8).

## 133 3 Experiment

134 **Entropy Estimation on Gaussian (Fig.1A, Fig.1D, Fig.13) and**

135 **GMM (Fig.15, Fig.20) Targets. MET-SVGD** consistently out-

136 performs **P-SVGD**, notably, Fig.1D and Fig.20 show that, while

137 **P-SVGD** and projection-based baselines, e.g., S-SVGD [Gong et al.,

138 2021] struggle to scale beyond 20-d spaces, **MET-SVGD** achieves

139 high accuracy in up to 100-d. Note that **MET-SVGD** mitigates the

140 vanishing repulsive force (Fig.13c) identified as the root cause of

141 SVGD’s poor scalability in Ba et al. [2022]. **Learning EBMs.**

142 Training EBMs  $p_\phi(x) = \bar{p}_\phi(x)/Z$  via maximum likelihood is in-

143 tractable due to the partition function  $Z$ . When the sampler has

144 a tractable distribution  $q_\theta$ , a tight lower bound can be computed:

145  $\mathcal{L}_{ELBO}(\phi, \theta) = \mathbb{E}_{x \sim q_\theta} [\log \bar{p}_\phi(x)] - \mathbb{E}_{x \sim p_d} [\log \bar{p}_\phi(x)] + \mathcal{H}(q_\theta)$ , as

146 detailed in App. 10. The entropy is often omitted due to its compu-

147 tational complexity, yielding the commonly used contrastive divergence loss  $\mathcal{L}_{CD}(\phi)$ . We optimize

148  $\mathcal{L}_{ELBO}(\phi, \theta)$  using both **P-SVGD** and **MET-SVGD**, and train with  $\mathcal{L}_{CD}(\phi)$  using LD. Experiments

149 are conducted on the **Moon dataset** [Rezende and Mohamed, 2015b] (Fig. 21) and CIFAR10 (Fig. 3).

150 For CIFAR10, we report the Frechet Inception Distance (FID) over 5 seeds. In Fig. 3, we show that

151 not including the trace of Hessian in **MET-SVGD** (purple) leads to divergence. Using an adaptive

152 number of steps  $L_c$  stabilizes the training (green). Replacing  $\sigma_{med}$  with the learnable one (red) im-

153 proves stability and yields significantly better FID scores relative to **P-SVGD** (orange). Additionally,

154 learning the step-size (brown) enables faster convergence to the target ( $\epsilon_{\theta_3}^l \gg \epsilon$  in Fig. 25) and results

155 in smoother landscapes (Fig. 24). Yet, experiments with MH diverge. In App. 10, we show that

156 MH-augmented updates lead to a high rejection rate due to landscape complexity. This results in poor

157 sampling and eventually divergence. To mitigate this, in future work, we plan to explore controlling

158 the Lipschitz constant of the target. Also, learning only the kernel bandwidth does improve over

159 **P-SVGD**. We attribute this to vanishing gradients in high-dimensions, i.e., the kernel collapses to zero.

160 We plan to explore dimension-wise decomposable kernels. Qualitative results and implementation

161 details are in App. 10. **Max-Entropy RL** additional results are in App. 11.

## 162 4 Conclusion

164 **MET-SVGD** is a novel VI approach for entropy estimation from unnormalized density. It bridges the

165 gap between VI, particle-based inference and MCMC; sets a new SOTA for scaling SVGD sampling

166 to high-dimensional and non-smooth densities; and introduces a framework for learning sampler

167 parameters end-to-end and is also a new residual flow model with full rank Jacoian and adaptive

168 number of layers.

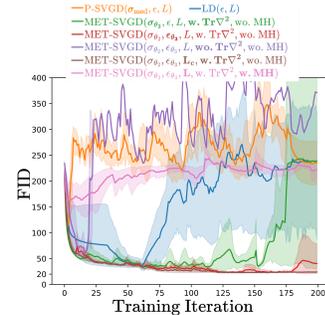


Figure 3: FID on CIFAR10.

The modification between two consecutive configs is **bolded**.

## References

- 169 John Wiley Sons, Ltd, 1992.
- 170 I. Ahmad and P.-E. Lin. A nonparametric estimation of the entropy for absolutely continuous  
172 distributions. *IEEE Trans. Inf. Theory*, 1976.
- 173 et al. Ahmed, Zafarali. Understanding the impact of entropy on policy optimization. *ICML*, 2019.
- 174 et al. Alemi, A. A. Deep variational information bottleneck. *ICLR*, 2016.
- 175 Gil Ariel and Yoram Louzoun. Estimating differential entropy using recursive copula splitting.  
176 *Entropy*, 2020.
- 177 Jimmy Ba, Murat A. Erdogdu, Marzyeh Ghassemi, Shengyang Sun, Taiji Suzuki, Denny Wu, and  
178 Tianzong Zhang. Understanding the variance collapse of svgd in high dimensions. In *ICML*, 2022.
- 179 Jens Behrmann, Will Grathwohl, Ricky TQ Chen, David Duvenaud, and Jörn-Henrik Jacobsen.  
180 Invertible residual networks. In *ICML*, 2019.
- 181 Jan Beirlant, Edward J Dudewicz, László Györfi, Edward C Van der Meulen, et al. Nonparametric  
182 entropy estimation: An overview. *Int. J. Math. Stat. Sci.*, 1997.
- 183 et al. Belghazi, M. I. Mutual information neural estimation. *ICML*, 2018.
- 184 Jose M Bernardo. Reference posterior distributions for bayesian inference. *Journal of the Royal*  
185 *Statistical Society Series B: Statistical Methodology*, 1979.
- 186 Vladimir Igorevich Bogachev, Aleksandr Viktorovich Kolesnikov, and Kirill Vladimirovich  
187 Medvedev. Triangular transformations of measures. *Sb. Math.*, 2005.
- 188 Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and  
189 Wojciech Zaremba. Openai gym, 2016.
- 190 Markov Chain Monte Carlo. Honest exploration of intractable probability distributions via. *Stat. Sci.*,  
191 2001.
- 192 Anthony L Caterini, Arnaud Doucet, and Dino Sejdinovic. Hamiltonian variational auto-encoder.  
193 *NeurIPS*, 2018.
- 194 Yogendra P. Chaubey and Pranab K. Sen. On nonparametric estimation of the density of a non-  
195 negative function of observations. *Calcutta Stat. Assoc. Bull.*, 2013.
- 196 Ricky TQ Chen, Jens Behrmann, David K Duvenaud, and Jörn-Henrik Jacobsen. Residual flows for  
197 invertible generative modeling. *NeurIPS*, 2019.
- 198 Wei-Chia Chen, Ammar Tareen, and Justin B Kinney. Density estimation on small data sets. *Phys.*  
199 *Rev. Lett.*, 2018.
- 200 Siddhartha Chib. Markov chain monte carlo methods: computation and inference. *Handbook of*  
201 *econometrics*, 2001.
- 202 Thomas M Cover. *Elements of information theory*. John Wiley & Sons, 1999.
- 203 Bo Dai, Hanjun Dai, Arthur Gretton, Le Song, Dale Schuurmans, and Niao He. Kernel exponential  
204 family estimation via doubly dual embedding. In *AISTATS*, 2019a.
- 205 Bo Dai, Zhen Liu, Hanjun Dai, Niao He, Arthur Gretton, Le Song, and Dale Schuurmans. Exponential  
206 family estimation via adversarial dynamics embedding. *NeurIPS*, 32, 2019b.
- 207 Jay L Devore, Kenneth N Berk, Matthew A Carlton, et al. *Modern mathematical statistics with*  
208 *applications*. Springer, 2012.
- 209 Yu G. Dmitriev and F. P. Tarasenko. On estimation of functionals of the probability density function  
210 and its derivatives. *Theory Probab. Its Appl.*, 1973.

- 211 M. D. Donsker and S. R. S. Varadhan. Asymptotic evaluation of certain markov process expectations  
212 for large time—ii. *Commun. Pure Appl. Math.*, 1975.
- 213 Andrew Duncan, Nikolas Nüsken, and Lukasz Szpruch. On the geometry of stein variational gradient  
214 descent. *J. Mach. Learn. Res.*, 2023.
- 215 Charles W Fox and Stephen J Roberts. A tutorial on variational bayesian inference. *Artif. Intell. Rev.*,  
216 2012.
- 217 Tomas Geffner and Justin Domke. Langevin diffusion variational inference. In *AISTATS*, 2023.
- 218 Aidan N Gomez, Mengye Ren, Raquel Urtasun, and Roger B Grosse. The reversible residual network:  
219 Backpropagation without storing activations. *NeurIPS*, 2017.
- 220 Wenbo Gong, Yingzhen Li, and José Miguel Hernández-Lobato. Sliced kernelized stein discrepancy.  
221 In *ICLR*, 2021.
- 222 László Györfi and Edward C Van der Meulen. Density-free convergence properties of various  
223 estimators of entropy. *Comput. Stat. Data Anal.*, 1987.
- 224 Eldad Haber, Lars Ruthotto, Elliot Holtham, and Seong-Hwan Jun. Learning across scales - a  
225 multiscale method for convolution neural networks. 2017.
- 226 Peter Hall and Sally C Morton. On the estimation of entropy. *Ann. Inst. Stat. Math.*, 1993.
- 227 Kakade S. M. Singh K. Hazan, E. and A. Van Soest. Provably efficient maximum entropy exploration.  
228 *NeurIPS*, 2019.
- 229 Matthew D Hoffman. Learning deep latent gaussian models with markov chain monte carlo. In  
230 *ICML*, 2017.
- 231 Michael F Hutchinson. A stochastic estimator of the trace of the influence matrix for laplacian  
232 smoothing splines. *Commun. Stat. Simul. Comput.*, 1989.
- 233 Jörn-Henrik Jacobsen, Arnold Smeulders, and Edouard Oyallon. i-revnet: Deep invertible networks.  
234 *ICLR*, 2018.
- 235 Harry Joe. Estimation of entropy and other functionals of a multivariate density. *Ann. Inst. Stat.*  
236 *Math.*, 1989.
- 237 Galin L Jones and James P Hobert. Sufficient burn-in for gibbs samplers for a hierarchical random  
238 effects model. *Ann. Stat.*, 2004.
- 239 et al. Kandasamy, Kirthevasan. Nonparametric von mises estimators for entropies, divergences and  
240 mutual informations. *NeurIPS*, 2015.
- 241 Sudheesh Kumar Kattumannil. On stein’s identity and its applications. *Stat. Probab. Lett.*, 2009.
- 242 D. P. Kingma and M. Welling. Auto-encoding variational bayes. *ICLR*, 2013a.
- 243 D. P. Kingma and M. Welling. Auto-encoding variational bayes. *ICLR*, 2013b.
- 244 Ivan Kobyzev, Simon JD Prince, and Marcus A Brubaker. Normalizing flows: An introduction and  
245 review of current methods. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2020.
- 246 Anna Korba, Adil Salim, Michael Arbel, Giulia Luise, and Arthur Gretton. A non-asymptotic analysis  
247 for stein variational gradient descent. *NeurIPS*, 2020.
- 248 A. Kraskov. Estimating mutual information. *Phys. Rev. E*, 2004.
- 249 Ullrich Köthe. A review of change of variable formulas for generative modeling. 2023.
- 250 Qiang Liu. A short introduction to kernelized stein discrepancy. 2016. URL [https://api.  
251 semanticscholar.org/CorpusID:16209224](https://api.semanticscholar.org/CorpusID:16209224).
- 252 Qiang Liu. Stein variational gradient descent as gradient flow. *NeurIPS*, 2017.

- 253 Qiang Liu and Dilin Wang. Stein variational gradient descent: A general purpose bayesian inference  
254 algorithm. *NeurIPS*, 2016.
- 255 Qiang Liu and Dilin Wang. Stein variational gradient descent as moment matching. *NeurIPS*, 2018.
- 256 Qiang Liu, Jason D. Lee, and Michael I. Jordan. A kernelized stein discrepancy for goodness-of-fit  
257 tests and model evaluation. *ICML*, 2016.
- 258 Tianle Liu, Promit Ghosal, Krishnakumar Balasubramanian, and Natesh Pillai. Towards understanding  
259 the dynamics of gaussian-stein variational gradient descent. *NeurIPS*, 2024.
- 260 Xing Liu, Harrison Zhu, Jean-François Ton, George Wynne, and Andrew Duncan. Grassmann stein  
261 variational gradient descent. *AISTATS*, 2022.
- 262 Safa Messaoud, Billel Mokeddem, Zhenghai Xue, Linsey Pang, Bo An, Haipeng Chen, and Sanjay  
263 Chawla.  $\mathcal{S}^2$  ac: Energy-based reinforcement learning with stein soft actor critic. *ICLR*, 2024.
- 264 Kevin R Moon, Kumar Sricharan, Kristjan Greenewald, and Alfred O Hero III. Ensemble estimation  
265 of information divergence. *Entropy*, 2018.
- 266 Liam Paninski. Estimation of entropy and mutual information. *Neural Comput.*, 2003.
- 267 George Papamakarios, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji  
268 Lakshminarayanan. Normalizing flows for probabilistic modeling and inference. *JMLR*, 2021.
- 269 E. Parzen. On estimation of a probability density function and mode. *Ann. Math. Stat.*, 1962.
- 270 Colombo P. Boudiaf Pichler, G. Knife: Kernelized-neural differential entropy estimation. *ICLR*,  
271 2022.
- 272 D. J. Rezende and S. Mohamed. Variational inference with normalizing flows. *ICML*, 2015a.
- 273 Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *ICML*,  
274 2015b.
- 275 Christian P Robert, George Casella, Christian P Robert, and George Casella. The metropo-  
276 lis—hastings algorithm. *Monte Carlo statistical methods*, 2004.
- 277 CP Robert. Monte carlo statistical methods, 1999.
- 278 M. Rosenblatt. Remarks on some nonparametric estimates of a density function. *Ann. Math. Statist.*,  
279 1956.
- 280 Jeffrey S Rosenthal. Minorization conditions and convergence rates for markov chain monte carlo. *J.*  
281 *Am. Stat. Assoc.*, 1995.
- 282 Reuven Y. Rubinstein and Dirk P. Kroese. *The Cross Entropy Method: A Unified Approach To Com-*  
283 *binatorial Optimization, Monte-Carlo Simulation (Information Science and Statistics)*. Springer-  
284 Verlag, 2004.
- 285 Adil Salim, Lukang Sun, and Peter Richtarik. A convergence theory for svgd in the population limit  
286 under talagrand’s inequality t1. In *ICML*. PMLR, 2022.
- 287 Tim Salimans, Diederik Kingma, and Max Welling. Markov chain monte carlo and variational  
288 inference: Bridging the gap. In *ICML*. PMLR, 2015.
- 289 N. N. Schraudolph. Gradient-based manipulation of nonparametric entropy estimates. *IEEE Trans.*  
290 *Neural Netw. Learn. Syst.*, 2004.
- 291 Claude Elwood Shannon. A mathematical theory of communication. *Bell Syst. Tech. J.*, 1948.
- 292 Claude Elwood Shannon. A mathematical theory of communication. *ACM SIGMOBILE mobile*  
293 *computing and communications review*, 2001.
- 294 Jiaxin Shi and Lester Mackey. A finite-particle convergence rate for stein variational gradient descent.  
295 *NeurIPS*, 2022.

- 296 Jiaxin Shi and Lester Mackey. A finite-particle convergence rate for stein variational gradient descent.  
297 *NeurIPS*, 2024.
- 298 P Shyam. Model-based active exploration. *ICLR*, 2019.
- 299 Jure Sokolić, Raja Giryes, Guillermo Sapiro, and Miguel RD Rodrigues. Robust large margin deep  
300 neural networks. *IEEE Trans. Signal Process*, 2017.
- 301 Yang Song, Sahaj Garg, Jiaxin Shi, and Stefano Ermon. Sliced score matching: A scalable approach  
302 to density and score estimation. In *UAI*, 2020.
- 303 Kumar Sricharan, Dennis Wei, and Alfred O Hero. Ensemble estimators for multivariate entropy  
304 estimation. *IEEE Trans. Inf. Theory*, 2013.
- 305 Lukang Sun, Avetik Karagulyan, and Peter Richtarik. Convergence of stein variational gradient  
306 descent under a weaker smoothness condition. In *AISTATS*. PMLR, 2023.
- 307 Achille Thin, Nikita Kotelevskii, Jean-Stanislas Denain, Leo Grinsztajn, Alain Durmus, Maxim  
308 Panov, and Eric Moulines. Metflow: a new efficient method for bridging the gap between markov  
309 chain monte carlo and variational inference. *arXiv preprint arXiv:2002.12253*, 2020.
- 310 Luke Tierney. Markov chains for exploring posterior distributions. *Ann. Stat.*, 1994.
- 311 O. Vasicek. A test for normality based on sample entropy. *J R Stat Soc Series B Stat Methodol*, 1976.
- 312 Cédric Villani et al. *Optimal transport: old and new*. Springer, 2009.
- 313 Nicol Schraudolph Viola, Paul and Terrence J. Sejnowski. Empirical entropy manipulation for  
314 real-world problems. *NeurIPS*, 1995.
- 315 Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics. In *ICML*,  
316 2011.
- 317 Christopher S. Withers and Saralees Nadarajah.  $\log \det a = \text{tr} \log a$ . *Int. J. Math. Educ. Sci. Technol.*,  
318 2010.
- 319 Henry Wolkowicz and George PH Styan. Bounds for eigenvalues using traces. *Linear Algebra Appl.*,  
320 1980.
- 321 Jiaxi Wu, Jiaxin Chen, and Di Huang. Entropy-based active learning for object detection with  
322 progressive diversity constraint. In *CVPR*, 2022.
- 323 A. Zellner. An introduction to bayesian inference in econometrics. 1971.
- 324 Jiankui Zhou and Yue Qiu. Augmented message passing stein variational gradient descent. *arXiv*  
325 *preprint arXiv:2305.10636*, 2023.
- 326 Jingwei Zhuo, Chang Liu, Jiaxin Shi, Jun Zhu, Ning Chen, and Bo Zhang. Message passing stein  
327 variational gradient descent. In *ICML*. PMLR, 2018.

## 328 **NeurIPS Paper Checklist**

### 329 **1. Claims**

330 Question: Do the main claims made in the abstract and introduction accurately reflect the  
331 paper's contributions and scope?

332 Answer: [\[Yes\]](#)

333 Justification: The abstract and introduction accurately describe **MET-SVGD** contributions.  
334 Empirical results (GMM entropy estimation, EBMs, MaxEnt RL tasks) align closely with  
335 claims made.

336 Guidelines:

### 337 **2. Limitations**

338 Question: Does the paper discuss the limitations of the work performed by the authors?

339 Answer: [\[Yes\]](#)

340 Justification: We discuss the limitations at the end of Sec. 4.2 and Sec. 4.3.

341 Guidelines:

- 342 • The answer NA means that the paper has no limitation while the answer No means that  
343 the paper has limitations, but those are not discussed in the paper.
- 344 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 345 • The paper should point out any strong assumptions and how robust the results are to  
346 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
347 model well-specification, asymptotic approximations only holding locally). The authors  
348 should reflect on how these assumptions might be violated in practice and what the  
349 implications would be.
- 350 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
351 only tested on a few datasets or with a few runs. In general, empirical results often  
352 depend on implicit assumptions, which should be articulated.
- 353 • The authors should reflect on the factors that influence the performance of the approach.  
354 For example, a facial recognition algorithm may perform poorly when image resolution  
355 is low or images are taken in low lighting. Or a speech-to-text system might not be  
356 used reliably to provide closed captions for online lectures because it fails to handle  
357 technical jargon.
- 358 • The authors should discuss the computational efficiency of the proposed algorithms  
359 and how they scale with dataset size.
- 360 • If applicable, the authors should discuss possible limitations of their approach to  
361 address problems of privacy and fairness.
- 362 • While the authors might fear that complete honesty about limitations might be used by  
363 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
364 limitations that aren't acknowledged in the paper. The authors should use their best  
365 judgment and recognize that individual actions in favor of transparency play an impor-  
366 tant role in developing norms that preserve the integrity of the community. Reviewers  
367 will be specifically instructed to not penalize honesty concerning limitations.

### 368 **3. Theory Assumptions and Proofs**

369 Question: For each theoretical result, does the paper provide the full set of assumptions and  
370 a complete (and correct) proof?

371 Answer: [\[Yes\]](#)

372 Justification: We add assumptions and proofs in supplementary.

373 Guidelines:

- 374 • The answer NA means that the paper does not include theoretical results.
- 375 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
376 referenced.
- 377 • All assumptions should be clearly stated or referenced in the statement of any theorems.

- 378 • The proofs can either appear in the main paper or the supplemental material, but if
- 379 they appear in the supplemental material, the authors are encouraged to provide a short
- 380 proof sketch to provide intuition.
- 381 • Inversely, any informal proof provided in the core of the paper should be complemented
- 382 by formal proofs provided in appendix or supplemental material.
- 383 • Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 384 4. Experimental Result Reproducibility

385 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
 386 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
 387 of the paper (regardless of whether the code and data are provided or not)?

388 Answer: [Yes]

389 Justification: We provide experimental details in Appendix and share a link to the code.

390 Guidelines:

- 391 • The answer NA means that the paper does not include experiments.
- 392 • If the paper includes experiments, a No answer to this question will not be perceived
- 393 well by the reviewers: Making the paper reproducible is important, regardless of
- 394 whether the code and data are provided or not.
- 395 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 396 to make their results reproducible or verifiable.
- 397 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 398 For example, if the contribution is a novel architecture, describing the architecture fully
- 399 might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 400 be necessary to either make it possible for others to replicate the model with the same
- 401 dataset, or provide access to the model. In general, releasing code and data is often
- 402 one good way to accomplish this, but reproducibility can also be provided via detailed
- 403 instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 404 of a large language model), releasing of a model checkpoint, or other means that are
- 405 appropriate to the research performed.
- 406 • While NeurIPS does not require releasing code, the conference does require all submis-
- 407 sions to provide some reasonable avenue for reproducibility, which may depend on the
- 408 nature of the contribution. For example
  - 409 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
  - 410 to reproduce that algorithm.
  - 411 (b) If the contribution is primarily a new model architecture, the paper should describe
  - 412 the architecture clearly and fully.
  - 413 (c) If the contribution is a new model (e.g., a large language model), then there should
  - 414 either be a way to access this model for reproducing the results or a way to reproduce
  - 415 the model (e.g., with an open-source dataset or instructions for how to construct
  - 416 the dataset).
  - 417 (d) We recognize that reproducibility may be tricky in some cases, in which case
  - 418 authors are welcome to describe the particular way they provide for reproducibility.
  - 419 In the case of closed-source models, it may be that access to the model is limited in
  - 420 some way (e.g., to registered users), but it should be possible for other researchers
  - 421 to have some path to reproducing or verifying the results.

#### 422 5. Open access to data and code

423 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
 424 tions to faithfully reproduce the main experimental results, as described in supplemental  
 425 material?

426 Answer: [Yes]

427 Justification: provide code for reproducing all the figures in the paper.

428 Guidelines:

- 429 • The answer NA means that paper does not include experiments requiring code.
- 430 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
 431 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.

- 432 • While we encourage the release of code and data, we understand that this might not be  
433 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
434 including code, unless this is central to the contribution (e.g., for a new open-source  
435 benchmark).
- 436 • The instructions should contain the exact command and environment needed to run to  
437 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 438 • The authors should provide instructions on data access and preparation, including how  
439 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 440 • The authors should provide scripts to reproduce all experimental results for the new  
441 proposed method and baselines. If only a subset of experiments are reproducible, they  
442 should state which ones are omitted from the script and why.
- 443 • At submission time, to preserve anonymity, the authors should release anonymized  
444 versions (if applicable).
- 445 • Providing as much information as possible in supplemental material (appended to the  
446 paper) is recommended, but including URLs to data and code is permitted.

## 448 6. Experimental Setting/Details

449 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
450 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
451 results?

452 Answer: [Yes]

453 Justification: We share important details in the paper and the rest in supplementary.

454 Guidelines:

- 455 • The answer NA means that the paper does not include experiments.
- 456 • The experimental setting should be presented in the core of the paper to a level of detail  
457 that is necessary to appreciate the results and make sense of them.
- 458 • The full details can be provided either with the code, in appendix, or as supplemental  
459 material.

## 460 7. Experiment Statistical Significance

461 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
462 information about the statistical significance of the experiments?

463 Answer: [Yes]

464 Justification: Yes, across all our experiments, we run for several seeds (exact number  
465 depends on the experiment and is mentioned in the paper). We report mean and variance.  
466 For the RL experiments, we report the IQM (inter-quantile mean).

467 Guidelines:

- 468 • The answer NA means that the paper does not include experiments.
- 469 • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
470 dence intervals, or statistical significance tests, at least for the experiments that support  
471 the main claims of the paper.
- 472 • The factors of variability that the error bars are capturing should be clearly stated (for  
473 example, train/test split, initialization, random drawing of some parameter, or overall  
474 run with given experimental conditions).
- 475 • The method for calculating the error bars should be explained (closed form formula,  
476 call to a library function, bootstrap, etc.)
- 477 • The assumptions made should be given (e.g., Normally distributed errors).
- 478 • It should be clear whether the error bar is the standard deviation or the standard error  
479 of the mean.
- 480 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
481 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
482 of Normality of errors is not verified.

- 483
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- 484
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.
- 485
- 486
- 487

## 488 8. Experiments Compute Resources

489 Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

492 Answer: [Yes]

493 Justification: We provide the required details in supplementary.

494 Guidelines:

- The answer NA means that the paper does not include experiments.
  - The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
  - The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
  - The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).
- 495
- 496
- 497
- 498
- 499
- 500
- 501
- 502

## 503 9. Code Of Ethics

504 Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

506 Answer: [Yes]

507 Justification: We took the required measures to ensure that our submission is anonymous.

508 Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
  - If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
  - The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).
- 509
- 510
- 511
- 512
- 513

## 514 10. Broader Impacts

515 Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

517 Answer: [NA]

518 Justification: Our paper proposes a new sampling technique that can be leveraged in standard machine learning applications.

520 Guidelines:

- The answer NA means that there is no societal impact of the work performed.
  - If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
  - Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
  - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- 521
- 522
- 523
- 524
- 525
- 526
- 527
- 528
- 529
- 530
- 531
- 532
- 533
- 534

- 535
- 536
- 537
- 538
- 539
- 540
- 541
- 542
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
  - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 543 11. Safeguards

544 Question: Does the paper describe safeguards that have been put in place for responsible  
545 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
546 image generators, or scraped datasets)?

547 Answer: [NA]

548 Justification: we propose a new algorithm for sampling.

549 Guidelines:

- 550
- 551
- 552
- 553
- 554
- 555
- 556
- 557
- 558
- 559
- The answer NA means that the paper poses no such risks.
  - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
  - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
  - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 560 12. Licenses for existing assets

561 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
562 the paper, properly credited and are the license and terms of use explicitly mentioned and  
563 properly respected?

564 Answer: [NA]

565 Justification: We are the original owners.

566 Guidelines:

- 567
- 568
- 569
- 570
- 571
- 572
- 573
- 574
- 575
- 576
- 577
- 578
- 579
- 580
- 581
- The answer NA means that the paper does not use existing assets.
  - The authors should cite the original paper that produced the code package or dataset.
  - The authors should state which version of the asset is used and, if possible, include a URL.
  - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
  - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
  - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
  - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
  - If this information is not available online, the authors are encouraged to reach out to the asset's creators.

## 582 13. New Assets

583 Question: Are new assets introduced in the paper well documented and is the documentation  
584 provided alongside the assets?

585 Answer: [Yes]

586 Justification: We share the code

587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

**14. Crowdsourcing and Research with Human Subjects**

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

**15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA] .

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

## Supplementary Material

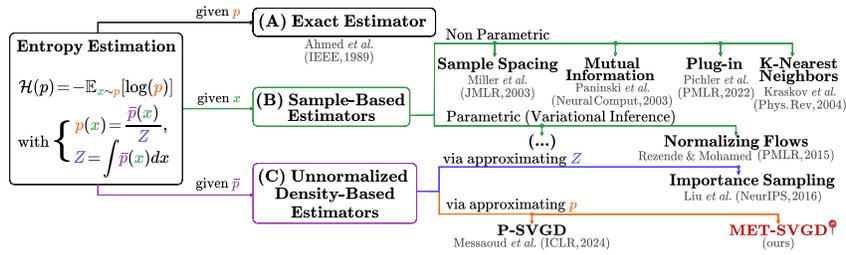


Figure 4: **MET-SVGD** is a new variational inference approach for entropy estimation of distributions known up to a normalization constant. It extends **P-SVGD** [Messaoud et al., 2024] to high-dimensional spaces by addressing its key limitations (see Fig. 1).

631 **MET-SVGD** is a novel variational inference approach for entropy estimation that overcomes key  
 632 limitations of **P-SVGD** Messaoud et al. [2024], particularly poor convergence and scalability in  
 633 high-dimensional spaces (Fig. 1). To achieve this, it introduces: (1) Sufficient condition for global  
 634 invertibility. (2) Optimized parameter search for improved stability (Sec. 2.2). (3) Metropolis-  
 635 Hastings augmented SVGD updates to ensure asymptotic convergence (Sec. 2.4). (4) A correction  
 636 term to the density estimation in **P-SVGD** (Sec. 2.3). **MET-SVGD** maintains computational efficiency,  
 637 requiring no significant additional memory or runtime overhead. Its full workflow is illustrated  
 638 in Algorithm 1. Beyond entropy estimation, **MET-SVGD** can be valuable to different research  
 639 communities:

- 640 • **MET-SVGD** bridges the gap between Metropolis-Hastings algorithms (MH) Robert et al.  
 641 [2004], particle-based sampling techniques (SVGD) Liu and Wang [2016], and parametrized  
 642 variational inference (P-VI) Fox and Roberts [2012], leveraging the strengths of each (Tab. 1):  
 643 (1) scalability from P-VI, (2) expressivity, convergence detection, and particle efficiency  
 644 from SVGD, as well as (3) convergence guarantees from MH. See Fig. 3
- 645 • **MET-SVGD** is a new approach for unprecedentedly scaling SVGD to high-dimensional  
 646 spaces while being computationally more efficient than all proposed approaches in the  
 647 literature Gong et al. [2021], Liu et al. [2022]
- 648 • **MET-SVGD** is a new approach for end-to-end learning of sampler parameters. It enables  
 649 training samplers via KL-divergence minimization, achieving compelling results for both  
 650 LD (Fig. 6B) and SVGD (Fig. 6A).
- 651 • **MET-SVGD** is a new normalizing flow model with (1) an adaptive number of updates  
 652 controlled by a convergence check and (2) a full-rank Jacobian for improved flexibility  
 653 and expressivity (Fig. ??). We plan to extend **MET-SVGD** to image generation using  
 654 flow-matching in future work.

655 The detailed algorithm is in Alg.1. We build a library for **MET-SVGD**. Our code is available  
 656 at: [https://anonymous.4open.science/r/Variational-Inference-with-SVGD--3F81/](https://anonymous.4open.science/r/Variational-Inference-with-SVGD--3F81/README.md)  
 657 [README.md](https://anonymous.4open.science/r/Variational-Inference-with-SVGD--3F81/README.md).

Criterion	P-VI	MCMC	SVGD	P-SVGD	MET-SVGD
Expressivity	✗	✓	✓	✓	✓✓
Convergence Detection	✓	✗	✓	✓	✓
Convergence Guarantees	✗	✓	✗	✗	✓
Sampling Efficiency	✓	✗	✓	✓	✓
Tractable Entropy	✓	✗	✗	✓	✓
Parameter Efficiency	✓	-	-	✓✓	✓✓

Table 1: **MET-SVGD** inherits advantages of different approximate inference methods: VI, SVGD, and MCMC.

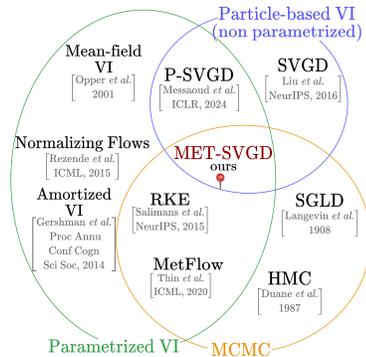


Figure 5: Bridges the gap ...

Table 2: P-SVGD vs MET-SVGD

Category	P-SVGD	MET-SVGD
Invertibility Condition	Local (Implicit Function Theorem); imprecise: $\epsilon \ll \sigma$ (Proposition 3.2, P-SVGD paper)	Global (Banach Theorem); precise: $\epsilon < \sqrt{\text{Tr}(\nabla \phi^l \nabla \phi^{l, \top})}$ (Corollary 2.3)
Entropy Trace Approximation	Imprecise: $\epsilon \ \nabla \phi^l\ _\infty \ll 1$ (Theorem 3.1, P-SVGD paper)	Automatically implied by invertibility condition (Corollary 2.3)
Divergence Control	Heuristic: particles truncation beyond 3 std from $q_{\theta_1}^{(0)}$ mean (Eq. 9, P-SVGD paper)	Metropolis-Hastings correction (Section 2.4)
Tr(Hessian) in Entropy	Omitted; invalid for finite particles (Theorem 3.3, P-SVGD paper)	Restored via Hutchinson estimator (Section 2.3)
Kernel Bandwidth $\sigma$	Median heuristic: $O(M^2)$	Learned via lightweight GNN (Section 2.2)
Step Size $\epsilon$	Fixed	Learned via lightweight GNN (Section 2.2)
Number of Steps $L$	Fixed	Adaptive via Stein Identity (Section 2.2)
Computation	Grid search for $\epsilon$ , median heuristic for $\sigma^2$ ( $O(M^2)$ )	Efficient reuse of $\text{Tr}(\nabla \phi^l)$ for the invertibility bound (Corollary 2.3), MH correction (Proposition 2.4), and convergence check; GNN inference adds minor overhead (Section 2.2)
Memory	-	Two small GNNs for $\sigma, \epsilon$ (Section 2.2)
Convergence Guarantee	$L, M \rightarrow \infty$	$L \rightarrow \infty$
Empirical Performance	Sensitive to hyperparameters (Fig. ??); mode collapse (Fig. ??); poor scalability to non-smooth and high-dimensional targets (Fig. ?? and Fig. ??)	SoTA entropy on G/GMM (Fig. 11 and Fig. 15); better FID, stability in EBMs for image generation (Fig. 3); improved Max-Ent RL returns (Fig. ??)

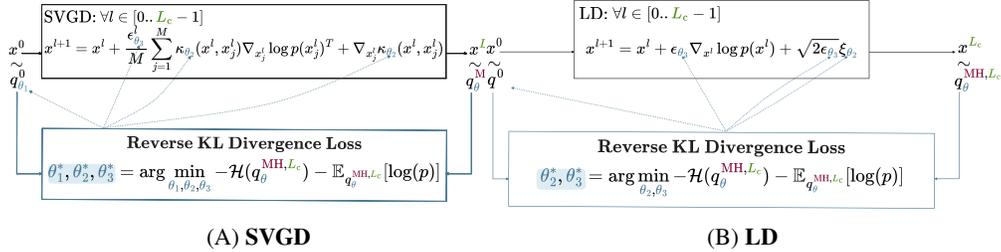


Figure 6: MET-SVGD provides a principled approach to learn sampler parameters via first computing the particles induced density, then Learning the parameters through KLD minimization.

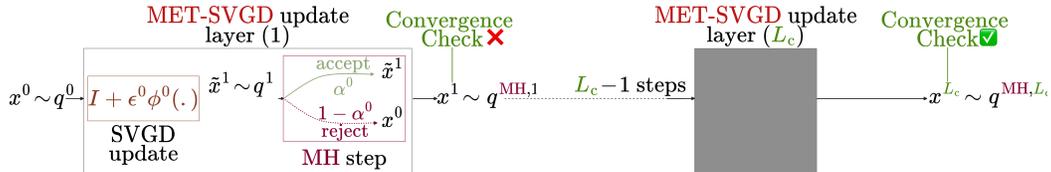


Figure 7: MET-SVGD is a normalizing flow model with a full rank Jacobian and an adaptive number of layers.

658 The rest of the appendix is organized as follows:

- 659 • Appendix 5: **Preliminaries**, including the Change of Variable formula for probability  
660 densities, Jacobi's formula corollary, the Stein Identity, the Banach Theorem and the implicit  
661 function theorem.

- 662 • Appendix 6: **Related work and Background** on entropy estimation, variational inference,
- 663 sampling-based variational inference, Normalizing Flows, Metropolis-Hastings, SVGD and
- 664 the Stein Identity.
- 665 • Appendix 7: **Derivation of closed-form density expressions for LD and SVGD samplers**
- 666 using RBF, Bilinear, and DKEF kernels. This section also includes derivation of the sufficient
- 667 condition on the step-size.
- 668 • Appendix 8 **Derivation of the Metropolis-Hastings augmented entropy**
- 669 • Appendix 9: **Additional results on entropy estimation**
- 670 • Appendix 10: **Additional results on learning EBMs for image generation**
- 671 • Appendix 11: **Additional results on MaxEnt RL**

---



---

Algorithm 1: MET-SVGD (Training)

**input** : Unnormalized density  $\bar{p}$ . SVGD parameters: (i) initial distr.  $q_{\theta_1}^0$ , (ii) number of particles  $M$ ,  
 (iii) maximum number of steps  $L$ , (vi) RBF kernel variance deepnet  $\sigma_{\theta_2}$  and (v) learning rate  
 deepnet  $\epsilon_{\theta_3}$ .

**output** :  $\theta^* = \{\theta_1^*, \theta_2^*, \theta_3^*\}$ .

- 1: **for** Each training iteration **do**
- 2:    $l = 0$  % initialize the number of SVGD steps
- 3:    $\{x_i^0\}_{i=0}^{M-1} \sim q_{\theta_1}^0$  % sample initial particles from  $q_{\theta_1}^0$
- 4:    $q_{\text{MH}}^0 = q_{\theta_1}^0$  % Initialize  $q_{\text{MH}}^0$
- 5:   % Run SVGD chain to convergence of si(Eq. 2.3)
- 6:   **while**  $(l \leq L)$  and  $(\Delta \text{SI}(q_{\theta}^{\text{MH},l}, p) \leq 0)$  **do**
- 7:      $\epsilon_{\theta_3}^l = \text{GNN}(\{x_i^l\}_{i=0}^{M-1}; \theta_3)$  % Compute learning rate
- 8:      $\epsilon_{\theta_3}^l = \min(\epsilon_{\theta_3}^l, \epsilon_{\text{UB}}^l)$  % Learning rate truncation (Corr.2.3)
- 9:      $\sigma_{\theta_2}^l = \text{GNN}(\{x_i^l\}_{i=0}^{M-1}; \theta_2)$  % Compute kernel variance
- 10:      $\tilde{x}_i^{l+1} \leftarrow x_i^l + \epsilon_{\theta_3}^l \phi(x_i^l), \forall i \in [0, M-1]$  % SVGD update (Eq. 1)
- 11:     % Metropolis Hastings Step (Sec. 2.4)
- 12:      $\alpha_{i,\theta_{2,3}}^l = (\alpha_{i,\theta_{2,3}}^l)^{a_i}, a_i \in \{0, 1\}, \forall i \in [0, M-1]$  % MH acceptance probability
- 13:      $u_i^l \sim \mathcal{N}(0, I)$  % Generate uniform random number
- 14:      $x_i^{l+1} = \tilde{x}_i^{l+1}$  If  $u_i^l \geq \alpha_i$ , Else  $x_i^{l+1} = \tilde{x}_i^l, \forall i \in [0, M-1]$  % Update
- 15:     % Update  $q_{\theta}^{\text{MH},l}$  (Eq. 2.4)
- 16:     
$$\log q_{\theta}^{\text{MH},l}(x^l) = \log q_{\theta}^{\text{MH},l-1}(x^{l-1}) + \log \left[ \exp(\log(\alpha_{\theta_{2,3}}^l) - \epsilon_{\theta_3} \text{Tr}(\nabla_{x^l} \phi(x^l))) + \right.$$

$$\left. \exp(\log(1 - \alpha_{\theta_{2,3}}^l)) \right]$$
- 17:      $l \leftarrow l + 1$  % Update number of steps
- 18:   **end while**
- 19:    $L_c \leftarrow l$
- 20:    $\mathcal{H}(q_{\theta}^{L_c}) = \frac{-1}{M} \sum_{i=0}^{M-1} \log q^{\text{MH},L_c}(x_i^{L_c})$  % Compute entropy
- 21:   % Update  $\theta$
- 22:    $\{\theta_1^*, \theta_2^*, \theta_3^*\} = \arg \max_{\theta_1, \theta_2, \theta_3} \mathbb{E}_{x^{L_c} \sim q_{\theta}^{\text{MH},L_c}} [\log p(x)] + \mathcal{H}(q_{\theta}^{\text{MH},L_c})$
- 23: **end for**
- 24: **Return**  $\theta^* = \{\theta_1^*, \theta_2^*, \theta_3^*\}$

---

672 **5 Preliminaries**

673 In the following, we review preliminaries about Langevin Dynamics, the Change of Variable formula  
 674 for pdfs, the corollary of the Jacobi formula, the Banach theorem, the Mean Value theorem, a  
 675 Sufficient Condition for residual flows invertibility and the Stein Identity.

676 **5.1 Langevin Dynamics**

677 **SGLD** [Welling and Teh, 2011] is a popular Markov chain Monte Carlo (MCMC) method for  
 678 sampling from a distribution. It first initializes a sample  $x^0$  from a random initial distribution. Then  
 679 at every step, it adds the gradient of the current proposal distribution  $p(x)$  to the previous sample  $x^l$ ,  
 680 together with a Brownian motion  $\xi \sim \mathcal{N}(0, I)$ . We denote with  $\epsilon$  the step size. The iterative update  
 681 for SGLD is:

$$x^{l+1} = x^l + \epsilon \nabla_{x^l} \log p(x^l) + \sqrt{2\epsilon} \xi. \quad (4)$$

682 **5.2 Change of Variable Formula (CVF)**

683 We first introduce the concept of an Invertible Function.

684 According to [Köthe, 2023], the following holds: if  $F : Z \rightarrow X$  is an invertible function then:

$$p_X(x) = p_Z(z) \left| \det \frac{\partial F^{-1}(x)}{\partial x} \right| = p_Z(z) \left| \det \frac{\partial F(z)}{\partial z} \right|^{-1}$$

685 **5.3 Implicit Function Theorem**

686 Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$  be continuously differentiable on some open set containing  $a$ , and suppose  
 687  $\det(\nabla_x f(x)) \neq 0$ . Then, there is some open set  $V$  containing  $x$  and an open  $W$  containing  $f(x)$   
 688 such that  $f : V \rightarrow W$  has a continuous inverse  $f^{-1} : W \rightarrow V$  which is differentiable  $\forall y \in W$ .

689 **5.4 Corollary of Jacobi's Formula**

690 Given an invertible matrix  $A$ , the following equality holds:

$$\log(\det A) = \text{Tr}(\log A) = \text{Tr}\left(\sum_{k=1}^{\infty} (-1)^{k+1} \frac{(A - I)^k}{k}\right). \quad (5)$$

691 The second equation is obtained by taking the power series of  $\log A$ . Hence, under the assumption  
 692  $\|A - I\|_{\infty} \ll 1$ , we obtain:  $\log(\det A) \approx \text{tr}(A - I)$ , where  $\|\cdot\|_{\infty}$  is the infinity norm.

693 **5.5 Banach Theorem**

694 We begin by introducing the concepts of a Cauchy sequence and a contractive mapping. Next, we  
 695 discuss the Banach Fixed Point theorem.

696 **Theorem 5.1** (Cauchy Sequence). *If a sequence  $\{x_n\}_{n \in \mathbb{N}}$  satisfy **either** of the following conditions:*

- 697 1.  $|x_{n+1} - x_n| \leq \alpha^n, \quad \forall n \in \mathbb{N}$
- 698 2.  $|x_{n+2} - x_{n+1}| \leq \alpha |x_{n+1} - x_n|, \quad \forall n \in \mathbb{N}$ ,

699 *where  $0 < \alpha < 1$ , then  $\{x_n\}$  is a Cauchy sequence.*

700 **Theorem 5.2** (Contractive Mapping). *Let  $(\mathcal{X}, d)$  be a metric space with  $d$  a distance function and let  
 701  $\phi : \mathcal{X} \rightarrow \mathcal{X}$  be a mapping on  $\mathcal{X}$ .  $\phi$  is called a contraction if and only if:*

$$\exists K \in [0, 1[ \quad \text{s.t.} \quad d(\phi(x), \phi(\tilde{x})) \leq K d(x, \tilde{x}), \quad \forall x, \tilde{x} \in \mathcal{X} \quad (6)$$

**Theorem 5.3** (Banach Fixed Point). *Let  $(X, d)$  be a complete metric space (i.e., all **Cauchy Sequences** are convergent) with  $d$  a distance function. If  $\phi$  is a contraction, then it has a **unique fixed point**  $x^* \in X$ , i.e.,  $\phi(x^*) = x^*$  and*

$$\forall x_0 \in \mathbf{X}, \quad \lim_{n \rightarrow \infty} \phi^n(x_0) = x^*, \quad \text{with} \quad \phi^n(x_0) = \underbrace{\phi \circ \phi \circ \dots \circ \phi}_{n \text{ times}}(x_0) = x_n.$$

702 *Proof.* The proof is structured in two main parts: we first establish the *existence* of a fixed point by  
 703 showing that  $(x_n)_{n \in \mathbb{N}}$  is a Cauchy sequence. Then prove *uniqueness* of the fixed point using a proof  
 704 by contradiction.

705 **Step 1: Existence of a fixed point.**  $(x_n)_{n \in \mathbb{N}}$  is a Cauchy sequence, we distinguish two cases:  
 706 consecutive samples and non-consecutive samples.

707 • *consecutive samples:*

$$d(x_{n+1}, x_n) = d(\phi(x_n), \phi(x_{n-1})) \leq K d(x_n, x_{n-1}) \leq K^2 d(x_{n-1}, x_{n-2}) \leq \dots \leq K^n d(x_1, x_0)$$

708 • *non-consecutive samples  $x_n$  and  $x_m$  with  $n < m$*

$$\begin{aligned} d(x_n, x_m) &\leq d(x_n, x_{n-1}) + d(x_{n-1}, x_{n-2}) + \dots + d(x_{m+1}, x_m) \\ &\leq (K^{n-1} + K^{n-2} + \dots + K^m) d(x_1, x_0) \\ &\leq K^m \underbrace{\sum_{k=0}^{n-1-m} K^k}_{\leq \sum_{k=0}^{\infty} K^k} d(x_1, x_0) \\ &\leq K^m \left( \sum_{k=0}^{\infty} K^k \right) d(x_1, x_0) = \frac{K^m}{1-q} d(x_1, x_0) \end{aligned}$$

709 It follows that  $\{x_n\}_{n \in \mathbb{N}}$  is a Cauchy sequence since  $d(x_n, x_m) \rightarrow 0$  as  $n, m \rightarrow \infty$ . Because the  
 710 metric space is complete, this implies convergence to a limit  $x^* \in \mathcal{X}$ : *i.e.*,  $x^* = \lim_{n \rightarrow \infty} x_n$ .  
 711 Additionally, since  $\phi$  is continuous,

$$\phi(x^*) = \phi\left(\lim_{n \rightarrow \infty} x_n\right) = \lim_{n \rightarrow \infty} \phi(x_n) = \lim_{n \rightarrow \infty} x_{n+1} = x^*.$$

712 Hence,  $x^*$  is a fixed point of  $\phi$ .

713 **Step 2: Uniqueness of the fixed point.** Assume that there exist two distinct fixed points  $x^*$  and  $\hat{x}$  such  
 714 that  $\phi(x^*) = x^*$  and  $\phi(\hat{x}) = \hat{x}$ . Then, If  $x^* \neq \hat{x} \Rightarrow d(x^*, \hat{x}) = d(\phi(x^*), \phi(\hat{x})) \leq K d(x^*, \hat{x})$   
 715 Which implies  $\Rightarrow \frac{d(x^*, \hat{x})}{d(x^*, \hat{x})} \leq K \Rightarrow 1 \leq K$ . which contradicts the assumption that  $K < 1$ . Hence, the  
 716 fixed point exists and is unique. We can compute it using the following algorithm:

---

Algorithm 2: Inverse of  $g(x)$  via fixed point iteration

---

```

input  $y^0 = g(x)$ , number of fixed-point iterations  $n$ 
  1: for  $i = 0 \dots n - 1$  do
  2:    $y^{i+1} = y^0 - g(y^i)$ 
  3: end for
  4: Return  $y^n = g(x)^{-1}$ 

```

---

717

□

## 718 5.6 The Mean Value Theorem

719 **Theorem 5.4.** Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be differentiable on  $\mathbb{R}^n$  with a Lipschitz continuous gradient  $\nabla f$ .  
 720 Then for given  $x$  and  $\bar{x}$  in  $\mathbb{R}^n$ , there is  $y = x + t(x - \bar{x})$  with  $t \in [0, 1]$ , such that

$$f(x) - f(\bar{x}) = \nabla f(y) \cdot (x - \bar{x}).$$

721 **5.7 Stein Identity ([Liu, 2016])**

722 Let  $p(x)$  be a continuously differentiable density supported on  $\mathcal{X} \subseteq \mathbb{R}^d$ , and let  $\phi(x) =$   
 723  $[\phi_1(x), \dots, \phi_d(x)]^T$  be a vector-valued function. **Stein's identity** states that for sufficiently regular  
 724  $\phi$ , we have:

$$\mathbb{E}_{x \sim p}[\mathcal{A}_p \phi(x)] = 0,$$

725 where the Stein operator  $\mathcal{A}_p$  is defined as:  $\mathcal{A}_p \phi(x) = \phi(x) \nabla_x \log p(x) + \nabla_x \phi(x)$ .

726 *Proof.* We can verify this identity using integration by parts under mild boundary assumptions: either  
 727  $p(x)\phi(x) = 0, \quad \forall x \in \partial\mathcal{X}$  when  $\mathcal{X}$  is compact, or  $\lim_{\|x\| \rightarrow \infty} \phi(x)p(x) = 0$  when  $\mathcal{X} = \mathbb{R}^d$ .

728 In the following we assume  $\mathcal{X} = [a, b]$ :

$$\begin{aligned} \mathbb{E}_{x \sim p}[\mathcal{A}_p \phi(x)] &= \int_a^b p(x) \phi(x) \nabla_x \log p(x) + p(x) \nabla_x \phi(x) dx \\ &\stackrel{(i)}{=} \int_a^b \phi(x) \nabla_x p(x) + p(x) \nabla_x \phi(x) dx \stackrel{(ii)}{=} [\phi(x)p(x)]_a^b \stackrel{(iii)}{=} 0 \end{aligned}$$

729 (i) Uses the identity  $\nabla_x \log p(x) = \frac{\nabla_x p(x)}{p(x)}$ .

730 (ii) Applies integration by parts:  $\int_a^b f(x)g'(x) + f'(x)g(x) dx = [f(x)g(x)]_a^b$ .

731 (iii) Boundary term vanishes under the stated assumptions.

732

□

733 **6 Related Work**

734 In the following, we review work on the differential entropy, variational inference, sampling-based  
 735 variational inference, Normalizing Flows, Stein Variational Gradient Descent (SVGD), Parametrized-  
 736 SVGD (P-SVGD), and Metropolis Hastings (MH) convergence.

737 **6.1 Differential Entropy**

Differential entropy, first introduced by Shannon in his foundational work on information theory Shannon [1948], has been widely studied in statistics Box [1992], Zellner [1971], Bernardo [1979]. For a continuous random variable  $z$  with density  $p(z)$ , the entropy is defined as:

$$\mathcal{H}(x) = - \int_{-\infty}^{\infty} p(x) \log(p(x)) dx.$$

738 **Applications of Entropy:** Entropy plays a crucial role in machine learning, Bayesian inference (BI),  
 739 reinforcement learning (RL), and variational inference (VI): (i) In classification & calibration, the  
 740 entropy measures model confidence Shyam [2019], used in active learning Wu et al. [2022]. (ii) In  
 741 Bayesian Inference, the Maximum Entropy principle ensures the least informative prior Bernardo  
 742 [1979]. (iii) In Reinforcement learning, it prevents overly deterministic policies by incorporating  
 743 entropy into the reward function Hazan and Van Soest [2019], Ahmed [2019]. (iv) Variational  
 744 inference & generative Models: The entropy appears in ELBO Kingma and Welling [2013b] for  
 745 posterior approximation and mitigates mode collapse in GANs and VAEs Alemi [2016], Belghazi  
 746 [2018].

748 **Challenges in Entropy Estimation:** Despite its simple definition, entropy is analytically tractable  
 749 only for limited distributions. For instance, for a uniform  $p(x) = \frac{1}{b-a}$  for  $x \in [a, b]$  and  $p(x) = 0$  for  
 750  $x \notin [a, b]$  the entropy is  $\mathcal{H}(p) = \frac{1}{2} [1 + \log(2\pi\sigma^2)]$ . for a Gaussian  $p(x) = \mathcal{N}(\mu, \sigma^2)$ , the entropy  
 751 is  $\mathcal{H}(p(y|\mu, \sigma^2)) = \frac{1}{2} (1 + \log(2\pi\sigma^2))$ . For general distributions, numerical integration (*e.g.*, Monte  
 752 Carlo) is required as direct computation is often infeasible. Different methods have been developed  
 753 for entropy estimation from samples.

754 **Entropy estimation methods from samples** can be classified into:

- 755 • *Plug-in Estimators:* Estimate density from data, then apply entropy formula. Given a sample  
 756  $x = \{x_i\}_{i=1}^M$ , the plug-in method estimates the pdf  $\hat{p}(x)$  from the data and then substitutes  
 757 this estimate into the entropy formula:  $\mathcal{H}^{\text{PLUGIN}}(p) \approx -\frac{1}{M} \sum_{i=1}^M \log \hat{p}(x_i)$ . This approach  
 758 was first proposed by Dmitriev et al. Dmitriev and Tarasenko [1973] and later investigated by  
 759 others using kernel density estimator Joe [1989], Hall and Morton [1993], Moon et al. [2018],  
 760 Pichler [2022], histogram estimator Györfi and Van der Meulen [1987], Hall and Morton  
 761 [1993] and field-theoretic approaches Chen et al. [2018]. Early approaches leverage kernels  
 762 that capture pairwise distances between the particles. For instance, Parzen-Rosenblatt  
 763 estimator Rosenblatt [1956], Parzen [1962]:  $\hat{p}(x) = \frac{1}{w^n} \sum_{i=1}^M \kappa\left(\frac{x-x_i}{w}\right)$ , where  $w$  denotes  
 764 the bandwidth and  $\kappa$  is a kernel density. The resulting entropy estimator was analyzed  
 765 by Ahmad and Lin [1976]. Schraudolph [2004] extended this approach using a kernel  
 766 estimator:  $\hat{p}(x) = \frac{1}{M} \sum_{i=1}^M \kappa_{\Sigma_i}(x - x_i)$ , where  $\Sigma = (\Sigma_1, \dots, \Sigma_n)$  are distinct diagonal  
 767 covariance matrices and  $\kappa_{\Sigma}(x) \sim \mathcal{N}(0, \Sigma)$  is a centered Gaussian density with covariance  
 768 matrix  $\Sigma$ . Pichler [2022] introduced KNIFE, a kernel-based estimator for density estimation  
 769 (DE) defined as:  $\hat{p}^{\text{KNIFE}}(x; \theta) = \sum_{i=1}^M \mu_i \kappa_{\Sigma_i}(x - b_i)$ , where  $\Sigma = (\Sigma_1, \Sigma_2, \dots, \Sigma_n)$ , and  
 770  $\theta = (\Sigma, b, \mu)$ , with the constraints  $\sum_{i=1}^M \mu_i = 1$ . The covariance matrices  $\Sigma_i$  are symmetric  
 771 and positive definite but not necessarily diagonal. Despite its advantages, the method  
 772 has a significant limitation in its simple structure, being restricted to either individual  
 773 Gaussian kernels or Gaussian Mixture Models (GMMs) with a fixed number of components  
 774  $n$ . This can limit its flexibility in modeling complex data distributions. Traditional off-  
 775 the-shelf density estimators often suffer from key drawbacks, such as non-differentiability,  
 776 computational intractability, or an inability to adapt to changes in the underlying data  
 777 distribution. These limitations make them unsuitable for applications requiring integration  
 778 into neural network training pipelines as regularizers. To improve density estimation for  
 779

780 non-negative random variables, recent studies have suggested replacing Gaussian kernels  
781 with Poisson weight-based estimators to fit counts or rate-based data Chaubey and Sen  
782 [2013] defined as:  $\hat{p}^{\text{POIS}}(x) = k \sum_{i=0}^{\infty} \left( F_n\left(\frac{i+1}{k}\right) - F_n\left(\frac{i}{k}\right) \right) e^{-kx} \frac{(kx)^i}{i!}$ , where  $F_n(\cdot)$  is the  
783 empirical distribution function, and  $k$  is a smoothing parameter. Additionally, the concept  
784 of learning kernel parameters end-to-end has been explored, providing a foundation for  
785 modern differentiable approaches. The idea of learning kernel parameters end-to-end has  
786 also been explored previously Viola and Sejnowski [1995], Schraudolph [2004], providing a  
787 foundation for modern differentiable approaches.

- 788 • *Sample-spacing Estimates* use distances between ordered samples (e.g., Vasicek estimator  
789 Vasicek [1976]). Sample spacing methods rely on the spacing of sorted samples and was  
790 initiated by Vasicek Vasicek [1976]:  $H^{\text{Vasicek}}(p) \approx -\frac{1}{M} \sum_{i=1}^M \log\left(\frac{n}{2m} (x_{i+1} - x_i)\right)$ ,  
791 where  $x_i$  are the order statistics and  $m$  is a positive integer smaller than  $\frac{n}{2}$ . One of the  
792 greatest weakness of sample-spacing-based estimator is the choice of spacing parameter  $m$ ,  
793 which does not have the optimal form.
- 794 • *Nearest-Neighbor Methods*: leverage distances to  $k$ -th nearest neighbor Kraskov [2004].  
795 This method estimates entropy using distances to the  $k$ -th nearest neighbor in the sample  
796 space Kraskov [2004], i.e.,  $\mathcal{H}(p) \approx \psi(n) - \psi(k) + \log(c_d) + \frac{d}{n} \sum_{i=1}^M \log \epsilon_i$ , where  $\psi$  is the  
797 digamma function defined as the logarithmic derivative of the gamma function  $\frac{d}{dx} \ln(\Gamma(x))$ ,  
798  $c_d$  is the volume of the unit  $d$ -dimensional ball, and  $\epsilon_i$  is the distance to the  $k$ -th nearest  
799 neighbor.
- 800 • *Variational Inference*: Optimizes a surrogate distribution  $q(x)$  to approximate  $p(x)$  Kingma  
801 and Welling [2013a]. The entropy is computed as Kingma and Welling [2013a]:  $\mathcal{H}(p) \approx$   
802  $-\mathbb{E}_{q(x)}[\log q(x)]$ , where  $q(y)$  is optimized to approximate  $p(x)$ .  $q$  is chosen to be easy  
803 to sample from, e.g., Gaussians, GMMs and Normalizing Flows Rezende and Mohamed  
804 [2015a].
- 805 • *Mutual Information (MI) Estimators*: Approximate entropy indirectly via MI relation-  
806 ships, i.e.,  $I(x, y) = \mathcal{H}(p_x) + \mathcal{H}(p_y) - \mathcal{H}(p_{x,y})$ , Belghazi [2018], where  $p_{x,y}$  is the joint  
807 distribution and  $p_x \cdot p_y$  is the product of the marginal distributions  $p_x$  and  $p_y$ . Neural  
808 networks were used to approximate the mutual information between two variables using  
809 the Donsker-Varadhan representation of the KL-Divergence Donsker and Varadhan [1975]:  
810  $D_{\text{KL}}(p||q) = \sup_{T \in \mathcal{T}} (\mathbb{E}_p[T(x)] - \log \mathbb{E}_q[e^{T(x)}])$ , where  $\mathcal{T}$  is a class of functions where  
811  $P_{x,y}$  is the joint distribution and  $p_x \cdot p_y$  is the product of the marginal distributions. The  
812 MI lower bound is expressed as:  $I_{\theta}(x; y) = \sup_{\theta} (\mathbb{E}_{p_{x,y}}[T_{\theta}(x, y)] - \log \mathbb{E}_{p_x \cdot p_y}[e^{T_{\theta}(x, y)}])$ ,  
813 where:  $T_{\theta}(x, y)$  is the output of a neural network parameterized by  $\theta$ ,  $\mathbb{E}_{p_{x,y}}[T_{\theta}(x, y)]$  is  
814 the expectation over samples from the joint distribution  $p_{x,y}$ , and  $\mathbb{E}_{p_x \cdot p_y}[e^{T_{\theta}(x, y)}]$  is the  
815 expectation over samples from the product of the marginals. The neural network is trained to  
816 maximize this bound, providing an approximation of  $I(x; y)$ . If two of these three entropies  
817  $\mathcal{H}(p_x)$ ,  $\mathcal{H}(p_y)$  or  $\mathcal{H}(p_{x,y})$  are available, the third one can be computed.
- 818 • *Ensemble Methods*: Weight different entropy estimators adaptively Sricharan et al. [2013].  
819 The estimators in the ensemble are assigned different weights, and the overall entropy  
820 estimate is calculated as a weighted combination of the individual estimators where optimal  
821 weights are determined by solving a convex optimization problem. Ariel and Louzoun  
822 [2020] proposed an innovative approach to estimating the entropy of high-dimensional data  
823 by decomposing the target entropy into two components:  $\mathcal{H}^{\text{CADEE}}(x) = \sum_{i=1}^d x_i + \mathcal{H}_{\text{copula}}$ ,  
824 where  $\mathcal{H}(y)$  is the total entropy of the multivariate distribution,  $\mathcal{H}(x_i)$  is the marginal entropy  
825 of each variable, and  $H_{\text{copula}}$  represents the entropy of the copula, capturing the dependencies  
826 between variables. The idea comes from the fact that any density distribution  $p(x)$  can  
827 be decomposed as the following:  $p(x) = p_1(x_1) \dots p_d(x_d) c(F_1(x_1), \dots, F_d(x_d))$ , where  
828  $c(u_1, \dots, u_d)$  is the density of copula. The copula entropy is estimated *recursively* by  
829 splitting the data into subgroups based on statistically dependent dimensions. This recursive  
830 process (1) identifies pairs or groups of dimensions with high statistical dependence, (2)  
831 splits the data along these dimensions and (3) repeats the process within each subgroup  
832 until the dependencies are resolved. [Kandasamy, 2015] proposed a leave-one-out technique  
833 to improve the robustness of entropy estimation using the von Mises expansion-based

834 estimator. The key idea is to iteratively remove one data point from the sample and compute  
835 the entropy estimate using the remaining data points. This procedure helps reduce bias  
836 and ensures that the estimator is not overly influenced by any single data point. The leave-  
837 one-out entropy is given by:  $\mathcal{H}^{\text{LOO}}(x) = \frac{1}{M} \sum_{i=1}^M \mathcal{H}(x_{-i})$  where  $\mathcal{H}(x_{-i})$ , is calculated for  
838  $x_{-i} = \{x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n\}$ . This approach provides a more robust estimate of the  
839 entropy by mitigating the influence of outliers or anomalous data points.

840 A summary of these methods is provided in Tab. 3.

Method	Formula	Key Idea
Analytical	$\mathcal{H}(x)$	Closed-form expressions
Plugin	$-\frac{1}{M} \sum_{i=1}^M \log \hat{p}(x_i)$	Sampling-based estimation
KDE	$-\frac{1}{M} \sum_{i=1}^M \log \left( \frac{1}{nh} \sum_{j=1}^M \kappa \left( \frac{x_i - x_j}{h} \right) \right)$	Density smoothing
KNIFE	$\sum_{i=1}^M \mu_i \kappa_{\Sigma_i}(x - b_i)$	Kernel-based estimator
Nearest-Neighbor	$\psi(n) - \psi(k) + \log(c_d) + \frac{d}{n} \sum_{i=1}^M \log \epsilon_i$	Distance-based estimation
Vasicek	$-\frac{1}{M} \sum_{i=1}^M \log \left\{ \frac{n}{2m} (x_{i+1} - x_i) \right\}$	Sorted sample spacing
Variational Inference	$-\mathbb{E}_{q(x)}[\log q(x)]$	Surrogate distribution
MINE	$\sup_{\theta} \left( \mathbb{E}_{p_{x,y}}[T_{\theta}(x,y)] - \log \mathbb{E}_{p_x, p_y}[e^{T_{\theta}(x,y)}] \right)$	Calculate it via Infrmion
CADEE	$= \sum_{i=1}^d \mathcal{H}(y_i) + \mathcal{H}_{\text{copula}}$	Marginal via copula
LOO	$\frac{1}{M} \sum_{i=1}^M \mathcal{H}(x_{-i})$	Data driven approach

Table 3: Summary of Differential Entropy Approximations

## 841 6.2 Variational Inference

842 Variational Inference (VI) [Fox and Roberts, 2012] approximates a target distribution  $p(x) = \bar{p}(x)/Z$ ,  
843 known up to the normalizing constant  $Z$ , via a simpler-to-sample-from distribution  $q^*(x)$  from a  
844 predefined family  $\mathcal{Q} = \{q\}$ , by maximizing the KL-divergence *i.e.*,  $q^* = \arg \max_{q \in \mathcal{Q}} D_{\text{KL}}(p||q)$ . The  
845 choice of  $\mathcal{Q}$  significantly impacts performance; more expressive families yield better approximations.  
846 At convergence, the target entropy can be estimated as  $\mathcal{H}(p) = -\mathbb{E}_{q^*}[\log q^*]$ . While VI is scalable, it  
847 may not achieve the optimal  $q^*$  due to, either the limited expressivity of  $\mathcal{Q}$ , *i.e.*, easy to sample from  
848 distributions are usually over-simplistic (*e.g.*, Gaussians), or optimization challenges (*e.g.*, mode  
849 collapse in Normalizing Flows [Papamakarios et al., 2021]).

## 850 6.3 Sampling-based Variational Inference.

851 Bridging the gap between parametric variational inference (VI) and Markov Chain Monte Carlo  
852 (MCMC) has been a key research focus to achieve both expressivity and scalability in inference.  
853 A central challenge is deriving an analytical expression for the marginal distribution of the last  
854 sample in an MCMC chain, which is often intractable. To address this, prior work [Salimans et al.,  
855 2015, Geffner and Domke, 2023] introduced auxiliary variables to construct augmented variational  
856 distributions that include all samples from the chain. However, this approach requires optimizing  
857 a looser ELBO and estimating the reverse Markov kernel, which introduces additional parameters  
858 and complex design choices. Several extensions have been proposed to avoid estimating the reverse  
859 kernel: (i) Hoffman [2017] optimize ELBO with respect to the initial distribution and only uses the  
860 MCMC steps to produce “better” samples to the target distribution. However, this method lacks direct  
861 feedback between the final marginal distribution and variational parameters, limiting full unification  
862 of VI and MCMC, (ii) Caterini et al. [2018] propose a deterministic Hamiltonian MCMC by removing  
863 resampling and the accept-reject step. However, this sacrifices MCMC guarantees, (iii) Thin et al.  
864 [2020] introduce MetFlow, a Metropolis-Hastings method that models the proposal distribution  
865 as a normalizing flow, removing the need for inverse kernel estimation. **MET-SVGD** has several  
866 advantages compared with the aforementioned approaches: It computes the exact loglikelihood,  
867 *i.e.*, via using the change of variable formula (Sec. 5.2). Hence, there is no need in the variational  
868 approximation on the joint distribution of the samples of the Markov chain, to estimate the reverse  
869 dynamics. Besides, it leverages knowledge of the unnormalized density unlike classical flow models.  
870 This makes our approach very easy to integrate in modern day deep learning pipelines. The idea

871 of approximating log-likelihoods for distributions known up to a normalization constant using  
872 MCMC and the change-of-variable formula was first explored by [Dai et al., 2019b], applying it to  
873 Hamiltonian Monte Carlo (HMC) and Langevin Dynamics (LD). Since, they augment the input with  
874 noise or velocity variable for LD and HMC, respectively, the derived log-likelihood of the sampling  
875 distribution turns out to be –counter-intuitively– independent of the sampler’s dynamics and equal to  
876 the initial distribution, which is then parameterized using a normalizing flow model [Kobyzev et al.,  
877 2020]. Our derived log-likelihood is more intuitive as it depends on the SVGD dynamics.

## 878 6.4 Normalizing Flows, Residual Flows and Neural ODEs

879 We review Normalizing Flows in general and focus on residual flows as **MET-SVGD** is one. We  
880 also draw the connection to neural ODEs. **Normalizing Flows** are generative models that produce  
881 tractable distributions where both sampling and density evaluation can be efficient and exact. This  
882 is achieved by transforming a simple probability distribution (*e.g.*, a standard normal) into a more  
883 complex distribution by a sequence of invertible and differentiable mappings. The density of a sample  
884 can be evaluated by transforming it back to the original simple distribution and then computing the  
885 product of the density of the inverse-transformed sample under this distribution and the associated  
886 change in volume induced by the sequence of inverse transformations. The change in volume is  
887 the product of the absolute values of the determinants of the Jacobians for each transformation, as  
888 required by the change of variables formula (See App.5.2). Formally, Let  $x = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d$   
889 be a random variable with a known and tractable probability density function  $p_x : \mathbb{R}^d \rightarrow \mathbb{R}$ . Let  $g$  be  
890 an invertible function and  $x = F(z)$ . Then using the change of variables formula, one can compute  
891 the probability density function of the random variable  $y$ :

$$p_x(x) = p_z(F^{-1}(x)) \left| \det \nabla_x g^{-1}(x) \right| \quad (7)$$

892 Intuitively, if the transformation  $F$  can be arbitrarily complex, one can generate any distribution  $p_x$   
893 from any base distribution  $p_z$  under reasonable assumptions on the two distributions. This has been  
894 formally proven [Bogachev et al., 2005]. However, constructing arbitrarily complicated non-linear  
895 invertible functions can be difficult. Additionally,  $F$  should be sufficiently expressive to model the  
896 distribution of interest and computationally efficient, both in terms of computing  $F$ , its inverse and  
897 the determinant of the Jacobian  $\nabla_x F^{-1}(x)$ .

898 Different types of flows have been constructed: (1) **Elementwise Flows**, (2) **Linear Flows**, (3)  
899 **Planar Flows**, (4) **Radial Flows**, (5) **Coupling Flows**, (6) **Autoregressive Flows**, and (7) **Residual**  
900 **Flows**, which we focus on due to relevance to **MET-SVGD**.

901 **Residual Flows** are compositions of the function of the form  $g(x) = x + \phi(x)$ . The first attempts  
902 to build a reversible network architecture based on residual connections was motivated by saving  
903 memory (each layer activation can be reconstructed from the previous layer) Gomez et al. [2017],  
904 Jacobsen et al. [2018] and was achieved via partitioning units in each layer into two groups and  
905 defining coupling functions as:

$$y^A = x^A + F(x^B)y^B = x^B + G(y^A), \quad (8)$$

906 where  $x = (x^A, x^B)$  and  $y = (y^A, y^B)$  are respectively the input and output activations,  
907  $F : \mathbb{R}^{D-d} \rightarrow \mathbb{R}^d$  and  $G : \mathbb{R}^d \rightarrow \mathbb{R}^{D-d}$  are residual blocks. The Jacobian of such a transfor-  
908 mation is, however inefficient to compute and constrains the architecture. To address this, to enable  
909 unconstrained architectures for each residual block, Behrmann et al. [2019] proved the following  
910 statement:

911 **Proposition 6.1.** *A residual connection is invertible if the Lipschitz constant of the residual block is*  
912  *$Lip(\phi) < 1$ , where  $Lip(\phi) = \sup_{x \neq y} \frac{|\phi(x) - \phi(y)|}{|x - y|}$ . By the mean value theorem 5.6, if  $\phi$  is differentiable*  
913  *$\forall x$ , then  $Lip(\phi) = \sup_x \|\nabla_x \phi(x)\|_2$  with  $\|\cdot\|$  being the spectral norm.*

914 The detailed proof is in (App. 7.1) . Controlling the Lipschitz constant of a neural network is not  
915 trivial. Note, that regularizing the spectral norm of the Jacobian of  $\phi$  Sokolić et al. [2017] only  
916 reduces it locally and does not guarantee the above condition. Instead, Jacobsen et al. [2018] proposes  
917 constraining the spectral radius of each convolutional layer in this network to be less than one.

918 In residual flows, the density is also derived using the change of variable formula (App. 5.2). A  
919 different approach is proposed to approximate the log-det term:

$$\log |\det(I + \nabla_x \phi(x))| \stackrel{(i)}{=} \text{Tr}(\log(I + \nabla_x \phi(x))) \stackrel{(ii)}{=} \sum_{k=1}^{\infty} (-1)^{k+1} \frac{\text{Tr}(\nabla_x \phi(x))^k}{k}$$

920 Where (i) is obtained using the matrix identity result  $\log \det(A) = \text{Tr}(\log(A))$  for non-singular  
 921  $A \in \mathbb{R}^{d \times d}$  Withers and Nadarajah [2010] and (ii) follows from replacing the trace of the matrix by  
 922 its power series. By truncating this series one can calculate an approximation to the log Jacobian  
 923 determinant. To efficiently compute each member of the truncated series, the Hutchinson trick is used.  
 924 However, this resulted in a biased estimate of the log Jacobian determinant. An unbiased stochastic  
 925 estimator was proposed by [Chen et al., 2019]. In a model they called a Residual flow [Chen et al.,  
 926 2019], the authors used a Russian roulette estimator instead of truncation. Informally, the next term  
 927 is added to the partial sum while calculating the series, one flips a coin to decide if the calculation  
 928 should be continued or stopped.

929 **Neural ODEs.** Due to the similarity of ResNets and Euler discretizations, there are many connections  
 930 between the i-ResNet and ODEs. Residual connections can be viewed as discretizations of a first  
 931 order ordinary differential equation (ODE) [Haber et al., 2017]:

$$\frac{d}{dt} \mathbf{x}(t) = F(\mathbf{x}(t), \theta(t)), \quad (9)$$

932 where  $F : \mathbb{R}^D \times \Theta \rightarrow \mathbb{R}^D$  is a function which determines the dynamic (the *evolution function*),  $\Theta$  is  
 933 a set of parameters and  $\theta : \mathbb{R} \rightarrow \Theta$  is a parameterization. The discretization of this equation (Euler’s  
 934 method) is

$$\mathbf{x}_{n+1} - \mathbf{x}_n = \varepsilon F(\mathbf{x}_n, \theta_n), \quad (10)$$

935 and this is equivalent to a residual connection with a residual block  $\varepsilon F(\cdot, \theta_n)$ .

## 936 6.5 Stein Variational Gradient Descent

937 In the following, we provide an explanation of SVGD, the RBF kernel variance and its effect on the  
 938 SVGD dynamics, followed by the formal derivation of SVGD and related work on its convergence  
 939 rate.

940 **Stein Variational Gradient Descent (SVGD)** [Liu and Wang, 2016] is a sampling algorithm with  
 941 update rule given by Eq. 1. Traditionally, an RBF kernel is used, with its bandwidth  $\sigma$  set via the  
 942 median heuristic:  $\sigma_{\text{med}} = \text{median}\{\|x_i^l - x_j^l\|\}_{i,j=1}^M / \log M$ . Bandwidth  $\sigma$  determines the influence  
 943 of neighboring particles  $\{x_j^l\}$  on the update of each particle  $x_i^l$ : larger values lead to broader  
 944 neighborhoods, while setting  $\sigma = 0$  decouples the particles, making their updates independent (Fig. 8  
 945 in App. 6.5). SVGD has several advantages compared to other *approximate inference* approaches:  
 946 unlike classical variational inference (VI) methods, SVGD can sample from arbitrary complex  
 947 distributions under smoothness assumptions [Villani et al., 2009]. Compared to Markov Chain Monte  
 948 Carlo (MCMC) methods [Chib, 2001], SVGD is more particle efficient and its convergence can be  
 949 easily checked using the Stein Identity [Kattumannil, 2009]. However, *SVGD convergence* is only  
 950 proved under certain conditions, such as sub-Gaussian targets [Shi and Mackey, 2022] or infinite  
 951 particles [Salim et al., 2022, Liu and Wang, 2018]. Additionally, SVGD suffers from *poor scalability*  
 952 *in high-dimensional spaces* due to diminishing repulsive forces [Zhuo et al., 2018]. To address  
 953 this, existing solutions are based on dimensionality reduction via projections into low-dimensional  
 954 manifolds [Gong et al., 2021, Liu et al., 2022], this however leads to inflated variance (*e.g.*, S-SVGD  
 955 [Gong et al., 2021]) or impractical hyperparameters (*e.g.*, GSVGD [Liu et al., 2022]). Alternatively,  
 956 [Zhuo et al., 2018, Zhou and Qiu, 2023] propose using local kernels based on the Markov blanket, but  
 957 this requires prior knowledge of the target’s probabilistic graphical model. In this paper, we extend  
 958 Messaoud et al. [2024] work on deriving a closed-form expression of  $q^l$  enabling better scalability.

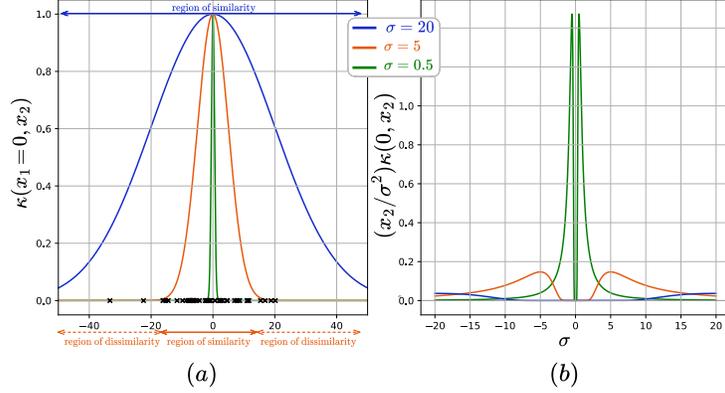


Figure 8: (a) Regions of similarity/dissimilarity for an RBF  $\kappa(x_1, x_2)$  evaluated at  $x_1 = 0$ . (b) Repulsion term in the SVGD update as a function of  $\sigma$ .

959 **RBF Kernel Variance Interpretation.** RBF kernels are the most generalized form of kernelization  
 960 and is one of the most widely used kernels due to its similarity to the Gaussian distribution. The RBF  
 961 kernel function for two points  $x_1$  and  $x_2$  computes the similarity or how close they are to each other.  
 962 This kernel can be mathematically represented as follows:  $\kappa(x_1, x_2) = \exp(-\frac{\|x_1 - x_2\|^2}{\sigma^2})$ , where  $\sigma$  is  
 963 the kernel variance and  $\|x_1 - x_2\|$  is the  $L_2$  distance between  $x_1$  and  $x_2$ . The maximum value that  
 964 the RBF kernel can reach is 1 when  $x_1 = x_2$ . When a large distance separates the points, the kernel  
 965 value is less than 1 and close to 0 indicating dissimilarity between  $x_1$  and  $x_2$ . This also means that  
 966 the particles are independent, *i.e.*, they follow their own gradients (the expectation in the SVGD  
 967 update is reduced to one term corresponding to  $x_i = x_j$ ). The width of the region of similarity is  
 968 controlled by  $\sigma$ , *i.e.*, a larger sigma results in a larger region of similarity with  $\kappa(x_1, x_2) \neq 0$  (Fig. 8  
 969 which also means that the particle update is impacted by its neighbors' gradients (b)).

970 Setting  $\sigma$  in the SVGD update rule;

$$x^{l+1} = x^l + \underbrace{\epsilon \mathbb{E}_{x_j^l} \left[ \kappa(x^l, x_j^l) \nabla_{x_j^l} \log p(x_j^l) \right]}_{\text{drift term}} + \underbrace{\frac{(x^l - x_j^l)}{\sigma^2} \kappa(x^l, x_j^l)}_{\text{repulsion term}} \quad (11)$$

971 is not obvious.  $\sigma$  in the drift term determines the neighboring samples  $x_j^l$  that will contribute with their  
 972 scores to the update. A larger  $\sigma$  implies, more influence from the neighbors. For the repulsion term,  
 973 both a very small or a very large  $\sigma$  value can result in setting the repulsion term to 0 as shown in Fig. 8  
 974 (a). Classically, the median trick is used to set the  $\sigma$ , *i.e.*,  $\sigma_{\text{med}} = \text{median}\{\|x_i^l - x_j^l\|\}_{i,j=1}^M / \log M$   
 975 with  $M$  being the number of particles. In our experiments, we show that this is suboptimal and that  
 976 that a more optimal  $\sigma$  can be learnt end-to-end via minimizing the KL-divergence (Eq. 2.2).

**SVGD Derivation.** Liu and Wang [2016] The goal is to approximate a target via a variational  
 distribution  $q \in \mathcal{Q}$  *i.e.*,

$$q^* = \arg \min_{q \in \mathcal{Q}} D_{KL}(q||p).$$

$\mathcal{Q}$  is obtained by transforming a reference density  $q^0$  via an invertible map  $F : \mathcal{X} \rightarrow \mathcal{X}$ , where for  
 any particle  $x \sim q^0$ , we define  $y = F(x)$ . The distributions of  $y$  and  $x$  are related by CVF (App. 5.2):

$$q_{[F]}(y) = q(F^{-1}(y)) \cdot |\det(\nabla_y F^{-1}(y))|$$

In this setup,  $F(x)$  is chosen to have a specific form:  $F(x) = x + \epsilon \phi(x)$ , where  $\epsilon$  is a stepsize and  $\phi$   
 is a perturbation direction chosen to maximally decrease the KL divergence:

$$\phi^* = \arg \max_{\phi} \{D_{KL}(q||p) - D_{KL}(q_{[F]}||p)\} = \arg \max_{\phi \in \mathcal{F}} \nabla_{\epsilon} D_{KL}(q_{[F]}||p)$$

This maximization has a closed form expression if we constrain the space of perturbations  $\mathcal{F}$  to be a  
 reproducing kernel Hilbert space (RKHS) with a positive kernel  $\kappa(\cdot, \cdot)$ , and  $\|\phi\|_{\mathcal{F}} \leq 1$ . In this case  
 $\arg \max_{\phi \in \mathcal{F}} \nabla_{\epsilon} D_{KL}(q_{[F]}||p) = \mathbb{E}_q[\text{Tr}(\mathcal{A}_p \phi)]$ . The optimal perturbation direction  $\phi^*$  is, hence, the  
 one that maximizes the Stein Discrepancy [Liu et al., 2016]:

$$\mathbb{S}(q, p) = \max_{\phi \in \mathcal{F}} \{\mathbb{E}_q[\text{Tr}(\mathcal{A}_p \phi)] \quad \text{s.t.} \quad \|\phi\|_{\mathcal{F}} \leq 1\}$$

and given by:

$$\phi_{p,q}^*(\cdot) = \mathbb{E}_q \left[ \kappa(x, \cdot) \nabla_x \log p + \nabla_x \kappa(x, \cdot) \right].$$

977 **SVG D Convergence Rate.** SVG D is difficult to analyze theoretically because it involves a system  
 978 of particles that interact with each other in a complex way. In the infinite particles case, [Liu, 2017]  
 979 proved that SVG D converges (weakly) to  $p$  in KSD. [Korba et al., 2020, Salim et al., 2022, Sun et al.,  
 980 2023] refined these results with path-independent constants, weaker smoothness conditions, and  
 981 explicit rates of convergence. [Duncan et al., 2023] provides conditions for exponential convergence.  
 982 For the finite particles case, [Liu, 2017] shows that finite particles SVG D converges to infinite  
 983 particles SVG D in bounded-Lipschitz distance but only under boundedness assumptions violated by  
 984 most applications of SVG D. [Korba et al., 2020] explicitly bounded the expected squared Wasserstein  
 985 distance between  $n$ -particle and continuous SVG D but only under the assumption of bounded  $\log p$ .  
 986 Also they do not provide convergence rates. [Liu et al., 2024] show that SVG D with finite particles  
 987 achieves linear convergence in KL divergence under a very limited setting where the target distribution  
 988 is Gaussian. [Shi and Mackey, 2024] shows that SVG D convergence rate is  $\mathcal{O}(1/\sqrt{\log \log n})$  under  
 989 the assumption that the target is sub-Gaussian with a Lipschitz score.

## 990 6.6 Parametrized-SVG D

991 **Parametrized SVG D (P-SVG D)** Messaoud et al. [2024] is a VI approach for entropy estimation  
 992 from unnormalized densities. Under *invertibility assumption* of the SVG D update rule (Eq. 1), it  
 993 computes the density of the SVG D particles  $q^L(x^L)$  by sequentially applying the Change of Variable  
 994 formula (CVF) [Devore et al., 2012] over  $L$  steps under an invertibility condition derived from the  
 995 implicit function theorem (App. 5.3):  $\log q^{l+1}(x^{l+1}) = \log q^l(x^l) - \log |\det(I + \epsilon \nabla_{x^l} \phi(x^l))|$ . To  
 996 avoid computing the full Jacobian, two approximations are used: (1) If  $\epsilon \|\nabla_{x^l} \phi(x^l)\|_\infty \ll 1$ , the  
 997 Jacobian determinant is reduced to its trace following Jacobi’s formula (App. 5.4) and leading to  
 998 Eq. 3.

## 999 6.7 Metropolis–Hastings

The Metropolis–Hastings algorithm’s goal is to generate a Markov Chain  $\{x^{(l)}\}_{l=0}^\infty$  that simulates  
 samples from a given probability distribution  $p$  Robert et al. [2004]. The chain starts with samples  
 from an initial distribution  $q^{(0)}$  and updates its state by leveraging a proposal distribution  $q(\tilde{x}|x^{(l)})$  as

$$x^{(l+1)}|x^{(l)} = \begin{cases} \tilde{x}, & \text{if } \alpha^l \leq \frac{p(\tilde{x})q(x^{(l)}|\tilde{x})}{p(x^{(l)})q(\tilde{x}|x^{(l)})} \\ x^{(l)}, & \text{otherwise} \end{cases}$$

1000 where  $\alpha^{(l)} \sim \mathcal{U}(0, 1)$ .

1001 Importantly, because the update only involves ratios of  $p$ , its normalization constant is not required.  
 1002 Furthermore, by construction, a chain that is constructed using the Metropolis-Hastings algorithm is  
 1003 reversible [Tierney, 1994], which means that if  $x^{(0)} \sim p$ , then  $x^l \sim p$  for all iterations  $l$ .

As an example, the Metropolis-Adjusted Langevin Algorithm employs the following proposal  
 distribution

$$q(\tilde{x}^{(l+1)}|x^{(l)}) = \mathcal{N}_d \left( x^{(l)} + \epsilon \nabla \log p(x^{(l)}), 2\epsilon I_d \right)$$

## 1004 6.8 Convergence of Metropolis Hastings

1005 Under relatively weak conditions, generating samples from an MCMC algorithm such as Metropolis-  
 1006 Hastings asymptotically draws samples from the target distribution [Robert, 1999]. The finite number  
 1007 of steps required for the marginal distribution of the Markov chain to reach the target under a  
 1008 discrepancy measure, has been heavily studied for both the total variation and Wasserstein distances  
 1009 [Carlo, 2001, Jones and Hobert, 2004, Rosenthal, 1995, Villani et al., 2009]. A popular approach is  
 1010 to show geometric ergodicity and provide an exponential convergence rate to the target distribution  
 1011 from any point of initialization in total variation. Explicit convergence rates have been rare with  
 1012 the exception of some Metropolis-Hastings independence samplers [Tierney, 1994]. To quantify  
 1013 said convergence, discrepancy measures are used. Notably, the total variation distance between two  
 1014 densities  $p$  and  $q$  defined as:  $d_{TV}(p, q) = \frac{1}{2} \int_{\mathcal{X}} |p(x) - q(x)| dx$ .

1015 An upper bound on the convergence rate can be computed as:

$$d_{\text{TV}}(q^L, p) \leq \left(1 - \frac{1}{\beta}\right)^L \quad \text{with } \beta = \sup_{x \in \mathcal{X}} \frac{p(x)}{q(x)}$$

1016 A lower bound can be computed as:

$$d_{\text{TV}}(q^l, p) \geq (1 - \alpha(x))^l \quad \text{with } \alpha(x) = \mathbb{E} \left[ \min \left( \frac{p(\tilde{x})q(x | \tilde{x})}{p(x)q(\tilde{x} | x)}, 1 \right) \right]$$

1017 In our case computing the lower bounds for **MET-SVGD** is possible as we have a closed-form  
1018 expression for the acceptance probability.

1019 **Notation:** We start by introducing the notation for this section. We compute the first and second  
 1020 order derivatives of the kernel as follows:

$$\begin{aligned} \forall i, j \in \{1..M\}^2 \quad \gamma &= \frac{1}{2\sigma^2} \quad \text{and} \quad \delta_{i,j} = (x_i^l - x_j^l) \quad \text{hence we express } \kappa, \nabla_{x_i} \kappa, \nabla_{x_i} \nabla_{x_j} \kappa \text{ as follows:} \\ \kappa(x_i^l, x_j^l) &= \exp(-\gamma \|x_i^l - x_j^l\|^2), \\ \nabla_{x_j^l} \kappa(x_i^l, x_j^l) &= 2\gamma \delta_{i,j} \kappa(x_i^l, x_j^l) \\ \nabla_{x_i^l} \kappa(x_i^l, x_j^l) &= -2\gamma \delta_{i,j} \kappa(x_i^l, x_j^l) = -\nabla_{x_j^l} \kappa(x_i^l, x_j^l) \\ \nabla_{x_i^l} \nabla_{x_j^l} \kappa(x_i^l, x_j^l) &= \nabla_{x_i^l} (2\gamma \delta_{i,j} \kappa(x_i^l, x_j^l)) = 2\gamma (I - 2\gamma \delta_{i,j} \delta_{i,j}^T) \kappa(x_i^l, x_j^l) \end{aligned}$$

## 1021 7 SVGD Density Derivation

1022 **Theorem 7.1.** Let  $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$  be an invertible transformation of the form  $F(x) = x + \epsilon\phi(x)$ . We  
 1023 denote by  $q^L(x^L)$  the distribution obtained from repeatedly ( $L$  times) applying  $F$  to a set of action  
 1024 samples (called ‘‘particles’’)  $\{x^0\}_{i=1}^M$  from an initial distribution  $q^0(x^0)$ , i.e.,  $x^L = F \circ F \circ \dots \circ F(x^0)$ .  
 1025 Under the condition  $\epsilon < \epsilon_{UB}^l = 1/\sup_x \sqrt{\text{Tr}(\nabla\phi^l(x)\nabla\phi^{l,T}(x))}$ ,  $\forall l \in [0..L]$ , the closed-form  
 1026 expression of  $\log q^L(x^L)$  is:

$$\log q^L(x^L) = \log q^0(x^0) - \epsilon \sum_{l=0}^{L-1} \text{Tr}(\nabla_{x^l} \phi(x^l)) + \mathcal{O}(\epsilon^2) \quad (12)$$

1027 *Proof.* Based on the change of variable formula (5.2), when for every iteration  $l \in [1, L]$ , the  
 1028 transformation  $x^l = F(x^{l-1})$  is invertible and we have:

$$q^l(x^l) = q^{l-1}(x^{l-1}) |\det \nabla_{x^l} \phi(x^l)|^{-1}, \forall l \in [1, L].$$

1029 By induction, we derive the probability distribution of sample  $x^L$ :

$$q^L(x^L) = q^0(x^0) \prod_{l=0}^{L-1} |\det (I + \epsilon \nabla_{x^l} \phi(x^l))|^{-1}$$

1030 By taking the log for both sides, we obtain:

$$\log q^L(x^L) = \log q^0(x^0) - \sum_{l=0}^{L-1} \log |\det (I + \epsilon \nabla_{x^l} \phi(x^l))|.$$

1031 This, however, requires computing the Jacobian  $\nabla_{x^l} \phi(x^l)$ . Next, we show that  $\log |\det (I +$   
 1032  $\epsilon \nabla_{x^l} \phi(x^l))|$  can be approximated efficiently via  $\epsilon \text{Tr}(\nabla_{x^l} \phi(x^l)) + \mathcal{O}(\epsilon^2)$  under an assumption  
 1033 on the learning rate in section 7.3, that’s satisfied by the invertibility assumption (Sec.7.1) and derive  
 1034 the expression of  $\text{Tr}(\nabla\phi)$  for the RBF, Bilinear and DKEF Kernels (Sec 7.5).  $\square$

### 1035 7.1 Sufficient Condition For $x + \epsilon\phi(x)$ Invertibility (Prop. 2.1)

1036 **Proposition 3.1** (Sufficient condition for invertible SVGD).

1037 Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$  with  $f = (f^1 \circ \dots \circ f^L)$  denote a sequence of SVGD updates with  $f^l = I + \epsilon\phi^l$ . We  
 1038 denote by  $\text{Lip}(\phi^l)$  the Lipschitz constant of the velocity  $\phi^l$  at step  $l$ .  $f$  is invertible if  $\epsilon \text{Lip}(\phi^l) < 1$ ,  
 1039 for all  $l \in [0, L - 1]$ .

1040 *Proof.* Given  $x^{l+1}$ , the goal is to find  $x^l$ . We denote by  $c$   $x^{l+1}$  resulting in  $x^{l+1} = c - \epsilon\phi(x^l)$ . Hence,  
 1041 we are interested in the invertibility of the function  $g(x) = c - \epsilon\phi(x)$ . for this, we show that  $g$  is a  
 1042 contractive mapping:

$$\begin{aligned}
d(g(x), g(\tilde{x})) &= d(c - \epsilon\phi(x), c - \epsilon\phi(\tilde{x})) \\
&\stackrel{(i)}{=} d(-\epsilon\phi(x), -\epsilon\phi(\tilde{x})) \\
&\stackrel{(ii)}{=} |\epsilon|d(\phi(x), \phi(\tilde{x})) \\
&\stackrel{(iii)}{\leq} |\epsilon|K \cdot d(x, \tilde{x}), \quad \text{with } |\epsilon|K < 1
\end{aligned}$$

1043 (i) The distance is translation invariant.

1044 (ii) The distance is absolutely homogeneous.

1045 (iii)  $\epsilon\phi$  is a contractive mapping, i.e.,  $d(\epsilon\phi(x), \epsilon\phi(\tilde{x})) \leq \epsilon K d(x, \tilde{x})$  with  $\epsilon K \leq 1$ . Note that  
1046  $\text{Lip}(\epsilon\phi) = \sup_x \frac{d(\epsilon\phi(x), \epsilon\phi(\tilde{x}))}{d(x, \tilde{x})} = \epsilon K$

1047 Therefore,  $g(x)$  is a contractive mapping, and by the **Banach fixed point theorem 2**, it has a unique  
1048 fixed point. This implies that the inverse of the mapping  $x^{l+1} = x^l + \epsilon\phi(x^l)$  exists and is unique.

1049 Hence, we demonstrate that  $f^l = I + \epsilon\phi^l$  is invertible if  $\epsilon \text{Lip}(\phi^l) < 1$ . Since  $f$  is a composition of  
1050  $f^l$  ( $l \in [1 \cdots L]$ ), we conclude that  $f$  is invertible.  $\square$

## 1051 7.2 A sufficient condition for invertibility check - an upper bound (Corr. 2.3)

1052 **Corollary 3.3.** *The distribution induced by the SVGD update (Eq. 1) using an RBF kernel is given by*  
1053 *Eq. 3 if  $\epsilon < \epsilon_{UB}^l = 1 / \sup_x \sqrt{\text{Tr}(\nabla\phi^l(x)\nabla\phi^{l,T}(x))} \quad \forall l \in [0, L-1]$*

1054 *Proof.*  $x^{l+1} = x^l + \epsilon\phi(x^l)$  is invertible if  $\text{Lip}(\phi(x)) < 1$  as we demonstrate in (App. 7.1)

1055 We compute the Lipschitz constant:

$$\text{Lip}(\phi) = \sup_x \frac{\|\phi(x) - \phi(y)\|}{\|x - y\|} \stackrel{(i)}{=} \sup_x \|\nabla_x \phi(x)\|_2 \stackrel{(ii)}{=} \sup_x \sigma_{\max}\{\nabla_x \phi(x)\}$$

1056 (i) We consider the  $\ell_2$  norm in computing the operator norm.

1057 (ii) Using the definition of the Lipschitz constant via Jacobian norm:  $\|\phi(x) - \phi(y)\| \leq$   
1058  $\sup_x \|\nabla_x \phi(x)\| \cdot \|x - y\|$ .

1059 The following always holds:  $\lambda_{\max}\{\nabla\phi\}$  is upper bounded by  $\|\nabla\phi\|_2 \leq \sqrt{\text{Tr}(\nabla\phi\nabla\phi^T)}$ .  $\square$

## 1060 7.3 A sufficient condition for log-det approximation (Prop. 2.2)

1061 **Proposition 3.2**(Condition for log-det Approximation) *Let  $\phi^l : \mathbb{R}^d \rightarrow \mathbb{R}^d$ ,  $\log|\det(I + \epsilon\nabla\phi^l)| =$*   
1062  *$\epsilon \text{Tr}(\nabla\phi^l)$  if  $\epsilon|\lambda_{\max}(\nabla\phi^l)| < 1$  for all  $l \in [0, L-1]$ , with  $\lambda_{\max}$  being the largest eigenvalue value*  
1063 *and  $\nabla$  is the gradient operator w.r.t the input.*

1064 *Proof.* We discuss two approaches leveraging the corollary of Jacobi's formula and the bounds on  
1065 the eigenvalues of  $\nabla_{x^l}\phi(x^l)$ :

1066 **Method 1 (P-SVD): Leveraging the Corollary of the Jacobi's formula.** Let  $A = I + \epsilon\nabla_{x^l}\phi(x^l)$ ,  
1067 under the assumption  $\epsilon\|\nabla_{x_i}\phi(x_i)\|_\infty \ll 1$ , i.e.,  $\|A - I\|_\infty \ll 1$ , we apply the collorary of Jacobi's  
1068 formula (App. 5.4) and get

$$\begin{aligned}
\log q^L(x^L) &= \log q^0(x^0) - \sum_{l=0}^{L-1} \text{Tr}(\log(I + \epsilon\nabla_{x^l}\phi(x^l))) + \mathcal{O}(\epsilon^2) \\
&= \log q^0(x^0) - \epsilon \sum_{l=0}^{L-1} \text{Tr}((I + \epsilon\nabla_{x^l}\phi(x^l) - I)) + \mathcal{O}(\epsilon^2)
\end{aligned}$$

$$= \log q^0(x^0) - \epsilon \sum_{l=0}^{L-1} \text{Tr}(\nabla_{x^l} \phi(x^l)) + \mathcal{O}(\epsilon^2)$$

1069 In practice, since this bound is informal, [Messaoud et al., 2024] recommend choosing a small enough  
1070 learning rate.

1071 **Method 2 (MET-SVGD): Leveraging bounds on the eigenvalues of  $\nabla_{x^l} \phi(x^l)$ .** In the following  
1072 we denote by  $\lambda_i\{A\}$  the eigenvalue of matrix  $A$

$$\begin{aligned} |\det(I + \epsilon \nabla_{x^l} \phi(x^l))| &\stackrel{(i)}{=} \left| \prod_{i=1}^d \lambda_i\{I + \epsilon \nabla_{x^l} \phi(x^l)\} \right| = \prod_{i=1}^d |\lambda_i\{I + \epsilon \nabla_{x^l} \phi(x^l)\}| \\ &\stackrel{(ii)}{=} \prod_{j=1}^d |1 + \epsilon \lambda_j\{\nabla_{x^l} \phi(x^l)\}| = \exp\left(\sum_{j=1}^d \ln |1 + \epsilon \lambda_j\{\nabla_{x^l} \phi(x^l)\}|\right) \\ &= \exp\left(\sum_{j=1}^d \ln(1 + \epsilon \lambda_j\{\nabla_{x^l} \phi(x^l)\})\right) \quad \text{if } \lambda_j\{\nabla_{x^l} \phi(x^l)\} > \frac{-1}{\epsilon} \\ &\stackrel{(iii)}{=} \exp\left(\sum_{j=1}^d \epsilon \lambda_j\{\nabla_{x^l} \phi(x^l)\} + \mathcal{O}(\epsilon^2)\right) \\ &= \exp\left(\epsilon \text{Tr}(\nabla_{x^l} \phi(x^l)) + \mathcal{O}(\epsilon^2)\right) \end{aligned}$$

1073 (i) By definition of the determinant.

1074 (ii) Let  $\lambda_i$  be the eigenvalue of  $\{I + \epsilon \nabla_{x^l} \phi(x^l)\}$  associated with the eigenvector  $v_i$ . We show  
1075 that  $\lambda_i - 1$  is the eigenvalue associated with  $\epsilon \nabla_{x^l} \phi(x^l)$ :

$$\begin{aligned} \Leftrightarrow (I + \epsilon \nabla_{x^l} \phi(x^l))v_i &= \lambda_i v_i \\ \Rightarrow \epsilon \nabla_{x^l} \phi(x^l)v_i &= (\lambda_i - 1)v_i \\ \Rightarrow \lambda_j &= (\lambda_i - 1) \text{ is an eigenvalue of } \epsilon \nabla_{x^l} \phi(x^l) \end{aligned}$$

1076 (iii) We use Taylor expansion of  $\ln(1 + \epsilon a) = \sum_i \frac{(-1)^{i-1}(\epsilon a)^i}{i} = \epsilon a + \mathcal{O}(\epsilon^2)$  around  $\epsilon a \rightarrow 0$ .

1077 Hence, under the condition  $\lambda_i\{\nabla_{x^l} \phi(x^l)\} > \frac{-1}{\epsilon}$ , the approximation  $\log |\det(I + \epsilon \nabla_{x^l} \phi(x^l))| =$   
1078  $\epsilon \text{Tr}(\nabla_{x^l} \phi(x^l))$  holds exactly.

$$\begin{aligned} \lambda_i\{\nabla_{x^l} \phi(x^l)\} &> \frac{-1}{\epsilon} \quad \forall i \in [1..d] \\ \Leftrightarrow |\lambda_i \epsilon| &< 1 \\ \Leftrightarrow |\lambda_i \epsilon| &< |\lambda_{\max} \epsilon| < 1 \quad \text{s.t. } \forall i \quad \lambda_i < \lambda_{\max} \\ \Leftrightarrow \epsilon &< \underbrace{1}_{\epsilon_{\text{UB}}} / |\lambda_{\max}| \end{aligned}$$

1079 Even though the condition  $\epsilon < \frac{\alpha}{|\lambda_{\max}|}$  is more exact than the one derived by [Messaoud et al., 2024],  
1080 it's still impractical as it requires computing the Jaccobian.

1081 **7.4 Unifying the sufficient conditions for invertibility and  $\log |\det(I + \epsilon A)| = \epsilon \text{Tr}(A) + \mathcal{O}(\epsilon^2)$**

**Corollary 7.2.** *Following [Wolkowicz and Styan, 1980], Let  $A$  be an  $d \times d$  complex matrix, and let  $A^*$  be the Hermitian of  $A$ :*

$$|\lambda_i| \leq \sigma_i \leq (\text{Tr}(A^*A))^{1/2} \quad \forall i \in [1..d]$$

1082 Where  $\sigma_i$  is the  $i$ -th singular value of  $A$ .

1083 *Proof.* In our setup  $A = \nabla_{x_i^l} \phi(x_i^l)$ , which we can easily compute as illustrated in the following:

$$\nabla_{x_i^l} \phi(x_i^l) = \underbrace{\frac{1}{M} \sum_{j=1, j \neq i}^M \nabla_{x_i^l} \kappa(x_i^l, x_j^l) s_p(x_j^l)^T}_{A_i} + \underbrace{\frac{1}{M} \underbrace{\kappa(x_i^l, x_i^l)}_{=1} \nabla_{x_i^l} s_p(x_i^l)^T}_{B_i}$$

1084 Next we compute  $\text{Tr}(\nabla_{x_i^l} \phi(x_i^l))$ . We denote by  $A_i$  and  $B_i$  the two terms of  $\nabla_{x_i^l} \phi(x_i^l)$ :

$$\begin{aligned} \text{Tr}((\nabla_{x_i^l} \phi(x_i^l))^T \nabla_{x_i^l} \phi(x_i^l)) &= \text{Tr}(A^T A) = \text{Tr}((A_i + B_i)^T (A_i + B_i)) \\ &= \text{Tr}(A_i^T A_i + B_i^T B_i + A_i^T B_i + A_i B_i^T) \\ &= \text{Tr}(A_i^T A_i) + \text{Tr}(B_i^T B_i) + 2\text{Tr}(A_i B_i^T) \\ &= \underbrace{\text{Tr}(A_i^T A_i)}_{(1)} + \underbrace{\text{Tr}(B_i^T B_i)}_{(2)} + \underbrace{2\text{Tr}(A_i B_i^T)}_{(3)} \end{aligned}$$

1085 For a term by term breakdown:

$$\text{Term (1)} = \text{Tr}(A_i^T A_i)$$

$$\begin{aligned} &= \text{Tr} \left( \left( \frac{1}{M} \sum_{\substack{j=1 \\ j \neq i}}^M \overbrace{\nabla_{x_i^l} \kappa(x_i^l, x_j^l) s_p(x_j^l)^T}^{C_{i,j}} + \overbrace{\nabla_{x_i^l} \nabla_{x_j^l} \kappa(x_i^l, x_j^l)}^{D_{i,j}} \right)^T \left( \frac{1}{M} \sum_{\substack{r=1 \\ r \neq i}}^M \underbrace{\nabla_{x_i^l} \kappa(x_i^l, x_r^l) s_p(x_r^l)^T}_{C_{i,r}} + \underbrace{\nabla_{x_i^l} \nabla_{x_r^l} \kappa(x_i^l, x_r^l)}_{D_{i,r}} \right) \right) \\ &= \text{Tr} \left( \frac{1}{M^2} \sum_{\substack{j=1 \\ j \neq i}}^M \sum_{\substack{r=1 \\ r \neq i}}^M (C_{i,j}^T + D_{i,j}^T) (C_{i,r} + D_{i,r}) \right) \\ &= \frac{1}{M^2} \sum_{\substack{j=1 \\ j \neq i}}^M \sum_{\substack{r=1 \\ r \neq i}}^M \underbrace{\text{Tr}(C_{i,j}^T C_{i,r} + D_{i,j}^T C_{i,r})}_{(1a)} + \underbrace{\text{Tr}(D_{i,j}^T D_{i,r} + C_{i,j}^T D_{i,r})}_{(1b)} \end{aligned}$$

$$\text{Term (1a)} = \text{Tr}(C_{i,r}^T C_{i,j} + D_{i,r}^T C_{i,j})$$

$$\begin{aligned} &= \text{Tr} \left( (\nabla_{x_i^l} \kappa(x_i^l, x_r^l) s_p(x_r^l)^T)^T (\nabla_{x_i^l} \kappa(x_i^l, x_j^l) s_p(x_j^l)^T) + (\nabla_{x_i^l} \nabla_{x_r^l} \kappa(x_i^l, x_r^l))^T (\nabla_{x_i^l} \kappa(x_i^l, x_j^l) s_p(x_j^l)^T) \right) \\ &= \text{Tr} \left( s_p(x_r^l) \nabla_{x_i^l} \kappa(x_i^l, x_r^l)^T \nabla_{x_i^l} \kappa(x_i^l, x_j^l) s_p(x_j^l)^T + (\nabla_{x_i^l} \nabla_{x_r^l} \kappa(x_i^l, x_r^l))^T (\nabla_{x_i^l} \kappa(x_i^l, x_j^l) s_p(x_j^l)^T) \right) \\ &= \text{Tr} \left( 4\gamma^2 \kappa(x_i^l, x_r^l) \kappa(x_i^l, x_j^l) s_p(x_r^l) \delta_{i,r}^T \delta_{i,j} s_p(x_j^l)^T - 4\gamma^2 \kappa(x_i^l, x_r^l) \kappa(x_i^l, x_j^l) (I - 2\gamma \delta_{i,r} \delta_{i,r}^T) \delta_{i,j} s_p(x_j^l)^T \right) \\ &= \text{Tr} \left( 4\gamma^2 \kappa(x_i^l, x_r^l) \kappa(x_i^l, x_j^l) (s_p(x_r^l) \delta_{i,r}^T - I + 2\gamma \delta_{i,r} \delta_{i,r}^T) \delta_{i,j} s_p(x_j^l)^T \right) \\ &= 4\gamma^2 \kappa(x_i^l, x_r^l) \kappa(x_i^l, x_j^l) (\delta_{i,r}^T s_p(x_r^l) - d + 2\gamma \|\delta_{i,r}\|^2) s_p(x_j^l)^T \delta_{i,j} \end{aligned}$$

$$\text{Term (1b)} = \text{Tr}(D_{i,r}^T D_{i,j} + C_{i,r}^T D_{i,j})$$

$$\begin{aligned} &= \text{Tr} \left( (\nabla_{x_i^l} \nabla_{x_r^l} \kappa(x_i^l, x_r^l))^T (\nabla_{x_i^l} \nabla_{x_j^l} \kappa(x_i^l, x_j^l)) + (\nabla_{x_i^l} \kappa(x_i^l, x_r^l) s_p(x_r^l)^T)^T (\nabla_{x_i^l} \nabla_{x_j^l} \kappa(x_i^l, x_j^l)) \right) \\ &= \text{Tr} \left( 4\gamma^2 \kappa(x_i^l, x_r^l) \kappa(x_i^l, x_j^l) (I - 2\gamma \delta_{i,r} \delta_{i,r}^T) (I - 2\gamma \delta_{i,j} \delta_{i,j}^T) - 4\gamma^2 \kappa(x_i^l, x_r^l) \kappa(x_i^l, x_j^l) s_p(x_r^l) \delta_{i,r}^T (I - 2\gamma \delta_{i,j} \delta_{i,j}^T) \right) \\ &= \text{Tr} \left( 4\gamma^2 \kappa(x_i^l, x_r^l) \kappa(x_i^l, x_j^l) (I - 2\gamma \delta_{i,r} \delta_{i,r}^T - s_p(x_r^l) \delta_{i,r}^T) (I - 2\gamma \delta_{i,j} \delta_{i,j}^T) \right) \end{aligned}$$

$$= 4\gamma^2 \kappa(x_i^l, x_r^l) \kappa(x_i^l, x_j^l) \left( d - \delta_{i,r}^T s_p(x_r^l) - 2\gamma |\delta_{i,r}|^2 \right) \left( d - 2\gamma |\delta_{i,j}|^2 \right)$$

1086 Adding these sub-terms together

$$\begin{aligned} \text{Term ①} &= \frac{1}{M^2} \sum_{\substack{j=1 \\ j \neq i}}^M \sum_{\substack{r=1 \\ r \neq i}}^M 4\gamma^2 \kappa(x_i^l, x_r^l) \kappa(x_i^l, x_j^l) \left( \delta_{i,r}^T s_p(x_r^l) - d + 2\gamma \|\delta_{i,r}\|^2 \right) s_p(x_j^l)^T \delta_{i,j} \\ &\quad + 4\gamma^2 \kappa(x_i^l, x_r^l) \kappa(x_i^l, x_j^l) \left( d - \delta_{i,r}^T s_p(x_r^l) - 2\gamma \|\delta_{i,r}\|^2 \right) \left( d - 2\gamma \|\delta_{i,j}\|^2 \right) \\ &= \frac{1}{M^2} \sum_{\substack{j=1 \\ j \neq i}}^M \sum_{\substack{r=1 \\ r \neq i}}^M 4\gamma^2 \kappa(x_i^l, x_r^l) \kappa(x_i^l, x_j^l) \left( \delta_{i,r}^T s_p(x_r^l) - d + 2\gamma \|\delta_{i,r}\|^2 \right) \left( s_p(x_j^l)^T \delta_{i,j} - d + 2\gamma \|\delta_{i,j}\|^2 \right) \end{aligned}$$

$$\begin{aligned} \text{Term ②} &= \text{Tr}(B_i^T B_i) \\ &= \text{Tr} \left( \frac{1}{M^2} \nabla_{x_i^l} s_p(x_i^l) \left( \nabla_{x_i^l} s_p(x_i^l) \right)^T \right) \\ &= \frac{1}{VM^2} \sum_{t=1}^V v_t^T \nabla_{x_i^l} s_p(x_i^l) \left( \nabla_{x_i^l} s_p(x_i^l) \right)^T v_t \\ &= \frac{1}{VM^2} \sum_{t=1}^V \left\| \nabla_{x_i^l} \left( v_t^T s_p(x_i^l) \right) \right\|^2 \end{aligned}$$

$$\begin{aligned} \text{Term ③} &= \text{Tr}(B_i^T A_i) \\ &\approx \frac{1}{V} \sum_{t=1}^V v_t^T \left( \frac{1}{M^2} \sum_{\substack{j=1 \\ j \neq i}}^M \left[ \underbrace{-2\gamma \nabla_{x_i^l} s_p(x_i^l) \delta_{i,j} s_p(x_j^l)^T}_{E_{i,j}} + \underbrace{2\gamma \nabla_{x_i^l} s_p(x_i^l) (I - 2\gamma \delta_{i,j} \delta_{i,j}^T)}_{F_{i,j}} \right] \kappa(x_i^l, x_j^l) \right) v_t \end{aligned}$$

Using Hutchinson Trace Estimation Hutchinson [1989]

$$\begin{aligned} &\approx \frac{1}{V} \sum_{t=1}^V \frac{1}{M^2} \sum_{\substack{j=1 \\ j \neq i}}^M \kappa(x_i^l, x_j^l) \left[ v_t^T E_{i,j} v_t + v_t^T F_{i,j} v_t \right] \\ &\approx \frac{1}{VM^2} \sum_{t=1}^V \sum_{\substack{j=1 \\ j \neq i}}^M \kappa(x_i^l, x_j^l) \left[ -2\gamma (v_t^T \nabla_{x_i^l} s_p(x_i^l)) (\delta_{i,j}^T v_t) s_p(x_j^l)^T + 2\gamma (v_t^T \nabla_{x_i^l} s_p(x_i^l)) (v_t - 2\gamma (\delta_{i,j}^T v_t) \delta_{i,j}) \right] \end{aligned}$$

1087 By combining **Terms ①, ② and ③**, we obtain:

$$\begin{aligned} (1) + (2) + (3) &= \frac{1}{M^2} \sum_{j=1}^M \sum_{\substack{r=1 \\ r \neq i}}^M 4\gamma^2 \kappa(x_i^l, x_r^l) \kappa(x_i^l, x_j^l) \left( \delta_{i,r}^T s_p(x_r^l) - d + 2\gamma |\delta_{i,r}|^2 \right) \left( s_p(x_j^l)^T \delta_{i,j} - d + 2\gamma |\delta_{i,j}|^2 \right) \\ &\quad + \frac{2}{M^2} \left| \nabla_{x_i^l} s_p(x_i^l) \right|^2 \\ &\quad + \frac{1}{VM^2} \sum_{t=1}^V \sum_{\substack{j=1 \\ j \neq i}}^M \kappa(x_i^l, x_j^l) \left[ -2\gamma (v_t^T \nabla_{x_i^l} s_p(x_i^l)) (\delta_{i,j}^T v_t) s_p(x_j^l)^T + 2\gamma (v_t^T \nabla_{x_i^l} s_p(x_i^l)) (v_t - 2\gamma (\delta_{i,j}^T v_t) \delta_{i,j}) \right] \end{aligned}$$

1088

□

1089 **7.5 Computing  $\text{Tr}(\nabla_{x^l} \phi(x^l))$  with RBF kernel**

1090 We show that the closed-form estimate of the log-likelihood  $\log q^L(x^L)$  for the SVGD-based sampler  
 1091 with an RBF kernel  $\kappa(\cdot, \cdot)$  is

$$\log q^L(x^L) \approx \log q^0(x^0) - \frac{\epsilon}{M\sigma^2} \sum_{l=0}^{L-1} \sum_{\substack{j=1 \\ x_j^l \neq x_i^l}}^M \left( \kappa(x_j^l, x^l) \left( -(x^l - x_j^l)^\top \nabla_{x_j^l} s_p(x_j^l) - \frac{\alpha}{\sigma^2} \|x^l - x_j^l\|^2 + d\alpha \right) \right) - \frac{\epsilon}{M} \text{Tr} \left( \nabla_{x_i^l}^2 \log p(x_i^l) \right)$$

1092 *Proof.* We explicitly compute  $\text{Tr}(\nabla_{x^l} \phi(x^l))$  as follows:

$$\log q^L(a_i^L) = \log q^0(x_i^0) - \frac{\epsilon}{m} \sum_{l=0}^{L-1} \left[ \sum_{\substack{j=1 \\ x_j^l \neq x_i^l}}^{m-1} \underbrace{\left[ \text{Tr} \left( \nabla_{x_i^l} (\kappa(x_i^l, x_j^l) \nabla_{x_j^l} \log p(x_j^l)) \right) \right]}_{\textcircled{1}} + \underbrace{\text{Tr} \left( \nabla_{x_i^l} \nabla_{x_j^l} \kappa(x_i^l, x_j^l) \right)}_{\textcircled{2}} \right. \\ \left. + \underbrace{\text{Tr} \left( \nabla_{x_i^l}^2 \log p(x_i^l) \right)}_{\textcircled{3}} \right]$$

1093 Next we compute simplifications for all subterms  $\textcircled{1}$  and  $\textcircled{2}$  respectively. In the following, we denote  
 1094 by  $(\cdot)^{(k)}$  the  $k$ -th dimension of the vector.

1095 **Term  $\textcircled{1}$ :**  
 1096

$$\begin{aligned} \text{Tr} \left( \nabla_{x_i^l} (\kappa(x_j^l, x_j^l) \nabla_{x_j^l} s_p(x_j^l)^T) \right) &= \text{Tr} \left( \nabla_{x_i^l} \kappa(x_j^l, x_j^l) (\nabla_{x_j^l} s_p(x_j^l))^\top + \kappa(x_j^l, x_j^l) \nabla_{x_i^l} \nabla_{x_j^l} s_p(x_j^l) \right) \\ &= \sum_{t=1}^d \frac{\partial \kappa(x_j^l, x_j^l)}{\partial (x_i^l)^{(t)}} \frac{\partial s_p(x_j^l)}{\partial (x_j^l)^{(t)}} + 0 \\ &= (\nabla_{x_i^l} \kappa(x_j^l, x_j^l))^\top \nabla_{x_j^l} s_p(x_j^l) \\ &= -\frac{1}{2\sigma^2} \kappa(x_j^l, x_j^l) (x_i^l - x_j^l)^\top \nabla_{x_j^l} s_p(x_j^l) \end{aligned}$$

1097 **Term  $\textcircled{2}$ :**

$$\begin{aligned} \text{Tr} \left( \nabla_{x_i^l} \nabla_{x_j^l} \kappa(x_i^l, x_j^l) \right) &= \text{Tr} \left( \nabla_{x_i^l} \left( \frac{1}{\sigma^2} \kappa(x_i^l, x_j^l) (x_i^l - x_j^l) \right) \right) \\ &= \frac{1}{\sigma^2} \sum_{k=1}^d \left( \frac{\partial \kappa(x_i^l, x_j^l)}{\partial (x_i^l)^{(k)}} (x_i^l - x_j^l)^{(k)} + \kappa(x_i^l, x_j^l) \right) \\ &= \frac{1}{\sigma^2} \left( \nabla_{x_i^l} \kappa(x_i^l, x_j^l)^\top (x_i^l - x_j^l) + d \times \kappa(x_i^l, x_j^l) \right) \\ &= \frac{1}{\sigma^2} \left( \nabla_{x_i^l} \kappa(x_i^l, x_j^l)^\top (x_i^l - x_j^l) + d \times \kappa(x_i^l, x_j^l) \right) \\ &= -\frac{1}{2\sigma^4} \times \kappa(x_i^l, x_j^l) \|x_i^l - x_j^l\|^2 + \frac{1}{2\sigma^2} \times d \times \kappa(x_i^l, x_j^l) \\ &= \kappa(x_i^l, x_j^l) \left( -\frac{1}{2\sigma^4} \|x_i^l - x_j^l\|^2 + \frac{d}{2\sigma^2} \right) \end{aligned}$$

1098 **Term  $\textcircled{3}$ : Using Hutchinson Trace Estimation Hutchinson [1989]**

$$\text{Tr} \left( \nabla_{x_i^l}^2 \log p(x_i^l) \right) \approx \frac{1}{V} \sum_{t=1}^V v_t^T \nabla_{x_i^l}^2 \log p(x_i^l) v_t$$

1099 By combining **Terms ①**, **②** and **③**, we obtain:

$$\begin{aligned} \log q^L(x_i^L) &= \log q^0(x_i^0) - \frac{\epsilon}{M\sigma^2} \sum_{l=0}^{L-1} \sum_{j=1}^M \kappa(x_j^l, x_i^l) \left( -(x_i^l - x_j^l)^\top \nabla_{x_j^l} s_p(x_j^l) - \frac{\alpha}{\sigma^2} \|x_i^l - x_j^l\|^2 + d\alpha \right) \\ &\quad - \frac{\epsilon}{MV} \sum_{t=1}^V v_t^\top \nabla_{x_i^l}^2 \log p(x_i^l) v_t \end{aligned}$$

1100 Proof done if we take a generic action particle  $x_i$  in place of  $x$ .

1101

□

## 1102 7.6 Computing $\text{Tr}(\nabla_{x^l} \phi(x^l))$ with Bilinear kernel

1103 [Liu et al., 2024] show that, for a Gaussian initial distribution,  $q^0(x) = \mathcal{N}(\mu_0, \Sigma_0)$ , and a tar-  
1104 get distribution  $p(x) = \mathcal{N}(b, Q)$  such that  $Q \in \mathbb{R}^{d \times d}$ . Applying SVGD with a Bilinear kernel  
1105  $\kappa(x_i, x_j) = \frac{x_j^\top x_i}{C} + 1$  and explicitly showing  $\log p(x^l) = -V(x^l) = -\frac{1}{2}(x^l - b)^\top Q^{-1}(x^l - b)$   
1106 produces a Gaussian density  $q^l$  at every step with mean  $\mu^l$  and covariance matrix  $\Sigma^l$ , satisfying the  
1107 following system of equations:

$$\begin{cases} \mu^{l+1} = \mu^l + \epsilon^l \left[ (I - (\Sigma^l + \mu^l \mu^{lT}) Q^{-1} + \mu^l b^\top Q^{-1}) \frac{\mu^l}{C} + (b - \mu^l)^\top Q^{-1} \right], \\ \Sigma^{l+1} = \Sigma^l + \epsilon^l \left[ 2 \frac{\Sigma^l}{C} - \frac{\Sigma^l}{C} Q^{-1} (\Sigma^l + (\mu^l - b) \mu^{lT}) - (\Sigma^l + \mu^l (\mu^l - b)^\top) Q^{-1} \frac{\Sigma^l}{C} \right]. \end{cases} \quad (13)$$

1108 We use this property to verify that the intermediate distributions  $q^l$  with the bilinear kernel are also  
1109 Gaussian. We make use of the bilinear kernel's expression in the proof

1110 *Proof.* For  $\kappa(x_i^l, x_j^l) = \frac{x_j^{lT} x_i^l}{C} + 1$ ,  $\nabla_{x_j^l} \kappa(x_i^l, x_j^l) = \frac{x_i^l}{C}$ , and  $\nabla_{x_i^l} \nabla_{x_j^l} \kappa(x_i^l, x_j^l) = \frac{I}{C}$  we have  
1111 the following SVGD dynamics at every step:

$$(x_i^{l+1} - x_i^l) / \epsilon^l = \frac{1}{M} \sum_{j=1}^M \nabla_{x_j^l} \kappa(x_i^l, x_j^l) + \frac{1}{M} \sum_{j=1}^M \kappa(x_i^l, x_j^l) \nabla_{x_i^l} \log p(x^l).$$

1112 Hence, substituting these into the dynamics we obtain:

$$\begin{aligned} (x_i^{l+1} - x_i^l) / \epsilon^l &= \frac{1}{M} \sum_{j=1}^M \frac{x_i^l}{C} - \frac{1}{M} \sum_{j=1}^M \left( \frac{x_j^{lT} x_i^l}{C} + 1 \right) \nabla V(x_j^l) \\ &= \frac{x_i^l}{C} - \frac{1}{M} \sum_{j=1}^M Q^{-1}(x_j^l - b) \left( \frac{x_j^{lT} x_i^l}{C} + 1 \right) \\ &= \frac{x_i^l}{C} - \frac{1}{M} \sum_{j=1}^M Q^{-1}(x_j^l - b) - \frac{1}{M} \sum_{j=1}^M Q^{-1}(x_j^l - b) \frac{x_j^{lT}}{C} x_i^l \\ &= \frac{x_i^l}{C} - \frac{Q^{-1}}{M} \sum_{j=1}^M (x_j^l - b) - \frac{Q^{-1}}{M} \sum_{j=1}^M x_j^l \frac{x_j^{lT}}{C} x_i^l + \frac{Q^{-1} b}{M} \sum_{j=1}^M \frac{x_j^{lT}}{C} x_i^l \\ &= \frac{x_i^l}{C} - Q^{-1}(\mu^l - b) - \frac{Q^{-1}}{C} (\Sigma^l + \mu^l \mu^{lT}) x_i^l + \frac{Q^{-1}}{C} b \mu^{lT} x_i^l \\ &= \left( \frac{I}{C} - \frac{Q^{-1}}{C} (\Sigma^l + \mu^l \mu^{lT}) + \frac{Q^{-1}}{C} b \mu^{lT} \right) x_i^l - Q^{-1}(\mu^l - b) \quad (i) \end{aligned}$$

We substitute  $\mu^l = \frac{1}{M} \sum_{j=1}^M x_j^l$  and that  $\Sigma^l + \mu^l \mu^{lT} = \frac{1}{M} \sum_{j=1}^M x_j^l x_j^{lT}$ , and obtain

$$\begin{aligned}
\frac{1}{M} \sum_{i=1}^M (x_i^{l+1} - x_i^l) / \epsilon^l &= \left( \frac{I}{C} - \frac{Q^{-1}}{C} (\Sigma^l + \mu^l \mu^{lT}) + \frac{Q^{-1}}{C} b \mu^{lT} \right) \mu^l - Q^{-1} (\mu^l - b) \\
&= \left( \frac{I}{C} - \frac{Q^{-1}}{C} \Sigma^l - \frac{Q^{-1}}{C} (\mu^l - b) \mu^{lT} \right) \mu^l - Q^{-1} (\mu^l - b) \\
&= \left( I - Q^{-1} \Sigma^l \right) \frac{\mu^l}{C} - \frac{Q^{-1}}{C} (\mu^l - b) \mu^{lT} \mu^l - Q^{-1} (\mu^l - b) \\
&= \left( I - Q^{-1} \Sigma^l \right) \frac{\mu^l}{C} - Q^{-1} (\mu^l - b) \left( \frac{\mu^{lT} \mu^l}{C} + 1 \right) \\
\text{Hence } (\mu^{l+1} - \mu^l) / \epsilon^l &= \left( I - Q^{-1} \Sigma^l \right) \frac{\mu^l}{C} - Q^{-1} (\mu^l - b) \left( \frac{\mu^{lT} \mu^l}{C} + 1 \right) \quad (\text{ii})
\end{aligned}$$

Knowing that

$$\begin{aligned}
(\Sigma^{l+1} - \Sigma^l) / \epsilon^l &= \frac{\partial \Sigma^l}{\partial l} \\
&= \frac{\partial}{\partial l} \frac{1}{M} \sum_{j=1}^M x_j^l x_j^{lT} - \mu^l \mu^{lT}
\end{aligned}$$

Taking into consideration (i) and (ii)

$$(\Sigma^{l+1} - \Sigma^l) / \epsilon^l = 2 \frac{\Sigma^l}{C} - \frac{\Sigma^l}{C} Q^{-1} (\Sigma^l + (\mu^l - b) \mu^l) - (\Sigma^l + \mu^l (\mu^l - b)^T) Q^{-1} \frac{\Sigma^l}{C} \quad (\text{iii})$$

1113

□

1114 In Fig. 9, we empirically verify that SVGD intermediate distributions coincide with the derived  
1115 Gaussians.

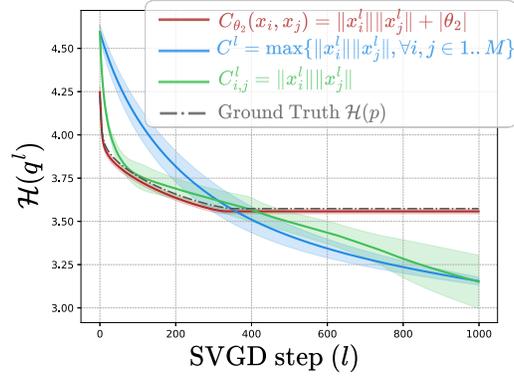


Figure 9: Bilinear kernel.  $q_\theta^l$  coincides with theoretically derived intermediate distributions by Liu et al. [2024].

## 1116 7.7 Computing $\text{Tr}(\nabla_{x^l} \phi(x^l))$ with DKEF kernel

1117 In the following, we show that using the deep exponential kernel (DKEF)  $\kappa(x_i^l, x_j^l) =$   
1118  $\exp(-\|\psi(x_i^l) - \psi(x_j^l)\|^2)$  where we denote by  $\psi(x) \in \mathbb{R}^m$  an  $m$  dimensional vector, is com-  
1119 putationally inefficient due to the requirement of computing  $\nabla_x \psi(x)$  in our entropy derivation, *i.e.*,

1120 terms 1 and 2 in the density below:

$$\log q^L(a_i^L) = \log q^0(x_i^0) - \frac{\epsilon}{m} \sum_{l=0}^{L-1} \left[ \sum_{\substack{j=1 \\ x_i^l \neq x_j^l}}^{m-1} \left[ \underbrace{\text{Tr}(\nabla_{x_i^l}(\kappa(x_i^l, x_j^l)) \nabla_{x_j^l} \log p(x_j^l))}_{\textcircled{1}} + \underbrace{\text{Tr}(\nabla_{x_i^l} \nabla_{x_j^l} \kappa(x_i^l, x_j^l))}_{\textcircled{2}} \right] + \underbrace{\text{Tr}(\nabla_{x_i^l}^2 \log p(x_i^l))}_{\textcircled{3}} \right]$$

1121 *Proof. Term ①:*

$$\begin{aligned} \text{Tr}(\nabla_{x_i^l}(\kappa(x_j^l, x_j^l)) \nabla_{x_j^l} \log p(x_j^l)) &= \text{Tr}(\nabla_{x_i^l} \kappa(x_j^l, x_j^l) (\nabla_{x_j^l} \log p(x_j^l))^\top + \kappa(x_j^l, x_j^l) \nabla_{x_i^l} \nabla_{x_j^l} \log p(x_j^l)) \\ &= \sum_{t=1}^d \frac{\partial \kappa(x_j^l, x_j^l)}{\partial (x_i^l)^{(t)}} \frac{\partial \log p(x_j^l)}{\partial (x_j^l)^{(t)}} + 0 \\ &= (\nabla_{x_i^l} \kappa(x_j^l, x_j^l))^\top \nabla_{x_j^l} \log p(x_j^l) \\ &= -\frac{1}{\sigma^2} \kappa(x_j^l, x_j^l) \nabla_{x_i} \psi(x_i) (\psi(x_i^l) - \psi(x_j^l))^\top \nabla_{x_j^l} \log p(x_j^l) \end{aligned}$$

1122 **Term ②:**

$$\begin{aligned} \text{Tr}(\nabla_{x_i^l} \nabla_{x_j^l} \kappa(x_i^l, x_j^l)) &= \text{Tr}\left(\nabla_{x_i^l} \left(\frac{1}{\sigma^2} \kappa(x_i^l, x_j^l) (\psi(x_i^l) - \psi(x_j^l))\right)\right) \\ &= \frac{1}{\sigma^2} \text{Tr}\left(\nabla_{x_i^l} \kappa(x_i^l, x_j^l) (\psi(x_i^l) - \psi(x_j^l))^\top + \kappa(x_i^l, x_j^l) \cdot I\right) \\ &= \frac{1}{\sigma^2} \text{Tr}\left(-\frac{1}{\sigma^2} \kappa(x_i^l, x_j^l) \nabla_{x_i^l} \psi(x_i^l) (\psi(x_i^l) - \psi(x_j^l)) (\psi(x_i^l) - \psi(x_j^l))^\top + \kappa(x_i^l, x_j^l) \cdot I\right) \end{aligned}$$

1123 **Term ③:**

$$\text{Tr}(\nabla_{x_i^l}^2 \log p(x_i^l)) \approx \frac{1}{V} \sum_{t=1}^V v_t^\top \nabla_{x_i^l}^2 \log p(x_i^l) v_t$$

1124  $\nabla_{x_i^l} \psi(x_i^l) \in \mathbb{R}^{m \times d}$  is required for both **Term ①** and **Term ②**, which introduces significant  
1125 computational overhead and renders the density intractable.  $\square$

## 1126 7.8 Derivation of the LD density

1127 Similarly to the derivation above the LD induced density can be derived as:

*Proof.*

$$\begin{aligned} \log q^{l+1}(x^{l+1}) &= \log q^l(x^l) + \log |\det \nabla_{x^l} \phi(x^l)|^{-1} \quad \text{where } \phi(x^l) = x^l - \epsilon \nabla_{x^l} \log p(x^l) + \sqrt{2\epsilon} \xi \\ &\stackrel{(i)}{=} \log q^l(x^l) - \epsilon \text{Tr}(\nabla_{x^l} \phi(x^l)), \quad \text{if } \lambda_i \{\nabla_{x^l} \phi(x^l)\} > -\frac{1}{\epsilon} \quad \forall i \\ &= \log q^l(x^l) - \epsilon \text{Tr}(\nabla_{x^l}^2 \log p(x^l)) \\ &\stackrel{(ii)}{=} \log q^l(x^l) - \frac{\epsilon}{V} \sum_{t=1}^V v_t^\top \nabla_{x_i^l}^2 \log p(x_i^l) v_t \end{aligned}$$

1128 (i) Using CVF (Eq 5.2)

1129 (ii) Using the Hutchinson Estimator, where  $p_v$  is chosen such that  $\mathbb{E}[vv^T] = I$  (eg.  $p_v$  is  
1130 Radamacher distribution.)

1131 In conclusion:

$$\log q^L(x^L) = \log q^0(x^0) - \frac{\epsilon}{V} \sum_{l=0}^L \sum_{t=1}^V v_t^T \nabla_{x_i}^2 \log p(x_i^l) v_t$$

1132

□

1133 **8 Metropolis Hastings augmented entropy**

1134 In the following, we provide derivations for (1) the acceptance probability (Prop.??), (2) convergence  
1135 check (Eq. ??) and (3) MH-augmented SVGD density (Eq. ??).

1136 **8.1 Acceptance Probability (Prop.??)**

1137 **Proposition 3.4** Given a target  $p = \bar{p}/Z$ , the log-likelihood of the MH acceptance probability for an  
1138 SVGD update of a particle  $x^{l-1}$  at step  $l$  is

$$\log \alpha(x^{l-1}, \tilde{x}^l) = \min \left[ 0, \log \bar{p}(\tilde{x}^l) - \log \bar{p}(x^{l-1}) + \epsilon \text{Tr}(\nabla_{x^l} \phi(x^l)) \right].$$

1139 *Proof.* By leveraging Bayes' rule, we compute:

$$\begin{aligned} \frac{q^l(x^{l-1} | \tilde{x}^l)}{q^l(\tilde{x}^l | x^{l-1})} &= \frac{q^l(x^{l-1}, \tilde{x}^l)}{q^l(\tilde{x}^l, x^{l-1})} \cdot \frac{q(x^{l-1})}{q(\tilde{x}^l)} \\ &= \frac{q(x^{l-1})}{q(\tilde{x}^l)} \\ &= \frac{q(x^{l-1})}{q(x^{l-1}) |\det(I + \epsilon \nabla_{x^{l-1}} \phi(x^{l-1}))|^{-1}} \\ &= |\det(I + \epsilon \nabla_{x^{l-1}} \phi(x^{l-1}))| \end{aligned}$$

1140 Thus, the Metropolis-Hastings ratio becomes:

$$\frac{p(x^l)}{p(x^{l-1})} \cdot \frac{q^l(x^{l-1} | \tilde{x}^l)}{q^l(\tilde{x}^l | x^{l-1})} = \frac{p(x^l)}{p(x^{l-1})} \cdot |\det(I + \epsilon \nabla_{x^{l-1}} \phi(x^{l-1}))|$$

1141 Taking logs:

$$\log \left( \frac{p(x^l)}{p(x^{l-1})} \cdot \frac{q^l(x^{l-1} | \tilde{x}^l)}{q^l(\tilde{x}^l | x^{l-1})} \right) = \log \left( \frac{p(x^l)}{p(x^{l-1})} \right) + \log |\det(I + \epsilon \nabla_{x^{l-1}} \phi(x^{l-1}))|$$

1142 Finally, using the first-order approximation  $\log |\det(I + A)| \approx \text{Tr}(A)$  for small  $\epsilon$ , we obtain:

$$= \log \left( \frac{p(x^l)}{p(x^{l-1})} \right) + \epsilon \text{Tr}(\nabla_{x^{l-1}} \phi(x^{l-1}))$$

1143 □

1144 **8.2 Motivation: learning the SVGD learning rate (Sec. ??-Step-Size)**

1145 Learning the kernel bandwidth alone is generally insufficient to ensure convergence of the entropy  
1146 term. Specifically, the expectation  $\mathbb{E}_{x^l \sim q^l} [\epsilon \text{Tr}(\nabla_{x^l} \phi(x^l))]$  does not necessarily vanish as  $l \rightarrow \infty$ . We  
1147 show, via a Taylor expansion around 0, that this cumulative trace term corresponds to a 4<sup>th</sup>-degree  
1148 polynomial in whose convergence to zero requires the existence of at least one real root. However,  
1149 the coefficients of this polynomial depend on the particle positions and are not guaranteed to yield a  
1150 real root during training, making this condition both non-trivial and fragile.

*Proof.*

$$\begin{aligned} \text{Tr}(\nabla_{x^{L_c}} \phi(x^{L_c})) &= \frac{\epsilon}{M\sigma^2} \sum_{\substack{j=1 \\ x^{L_c} \neq x_j^{L_c}}}^M \kappa(x_j^{L_c}, x_j^{L_c}) \left( (x^{L_c} - x_j^{L_c})^T \nabla_{x_j^{L_c}} \log p(x_j^{L_c}) + \frac{1}{\sigma^2} \|x^{L_c} - x_j^{L_c}\| - d \right) \\ &\quad + \frac{\epsilon}{M} \text{Tr} \left( \nabla_{x_j^{L_c}}^2 \log p(x_j^{L_c}) \right) \end{aligned}$$

1151 We approximate the RBF kernel using a Taylor expansion:

$$\exp\left(-\frac{\|x^{L_c} - x_j^{L_c}\|^2}{\sigma^2}\right) \approx 1 - \frac{\|x^{L_c} - x_j^{L_c}\|^2}{\sigma^2} + \frac{\|x^{L_c} - x_j^{L_c}\|^4}{2\sigma^4}$$

1152 We substitute in the formula above and obtain:

$$\begin{aligned} \text{Tr}(\nabla_{x^{L_c}} \phi(x^{L_c})) &= \frac{\epsilon}{M\sigma^2} \sum_{\substack{j=1 \\ x^{L_c} \neq x_j^{L_c}}}^M \left( 1 - \frac{\|x^{L_c} - x_j^{L_c}\|^2}{\sigma^2} + \frac{\|x^{L_c} - x_j^{L_c}\|^4}{2\sigma^4} \right) \times \left( (x^{L_c} - x_j^{L_c})^T \nabla_{x_j^{L_c}} \log p(x_j^{L_c}) \right) \\ &\quad + \frac{1}{\sigma^2} \|x^{L_c} - x_j^{L_c}\| - d + \frac{\epsilon}{M} \text{Tr} \left( \nabla_{x_j^{L_c}}^2 \log p(x_j^{L_c}) \right) \\ &= \frac{\epsilon\sigma^8}{M\sigma^2} \sum_{\substack{j=1 \\ x^{L_c} \neq x_j^{L_c}}}^M \left( \sigma^4 - \sigma^2 \|x^{L_c} - x_j^{L_c}\|^2 + \frac{\|x^{L_c} - x_j^{L_c}\|^4}{2} \right) \times \left( \sigma^2 (x^{L_c} - x_j^{L_c})^T \nabla_{x_j^{L_c}} \log p(x_j^{L_c}) \right) \\ &\quad + \|x^{L_c} - x_j^{L_c}\| - d\sigma^2 + \frac{\epsilon\sigma^8}{M} \text{Tr} \left( \nabla_{x_j^{L_c}}^2 \log p(x_j^{L_c}) \right) \end{aligned}$$

1153 Since we have a polynomial of degree 8, we have 8 roots, not all of which are guaranteed to be real  
 1154 for  $\sigma$ . In fact, we need the number of real roots to be computed as the signature of the Hermitian  
 1155 matrix. ie the number of real roots is equal to the number of positive values. In Fig. 10, we show that  
 1156 SI converges to 0 under different sampler configs of the experiments.

### 1157 8.3 Convergence Check (Eq. ??)

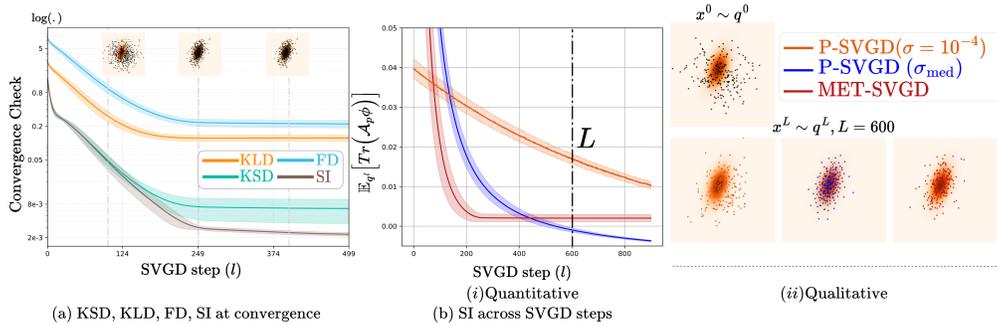


Figure 10: (a) SI shows the same convergence trend as other convergence metrics such as Fisher Divergence (FD) and Kernelized Stein Discrepancy (KSD). With SI being the tractable and computationally efficient metric. (b) SI can be used to check SVGD convergence across steps, for **MET-SVGD** and **P-SVGD**( $\sigma_{\text{med}}$ ,  $\sigma = 10^{-4}$ ).

#### 1158 8.3.1 Derivation of the Stein Identity

1159 The following is a proof of proposition ??:

1160 *Proof.* The Stein Identity is the square of the Kernelized Stein Discrepancy [Liu et al., 2016]:

$$\mathbb{S}(q^l, p) = \max_{\phi \in \mathcal{H}^d} \left[ \mathbb{E}_{x^l} [\text{Tr}(\mathcal{A}_p \phi(x^l))] \right]^2, \quad \text{s.t } \|\phi\|_{\mathcal{H}^d} \leq 1$$

1161 The optimal perturbation  $\phi$  is given by:

$$\phi(x^l) = \frac{\phi_{q,p}^*(x^l)}{\|\phi_{q,p}^*\|_{\mathcal{H}^d}}, \quad \text{with } \phi_{q,p}^*(\cdot) = \mathbb{E}_{x^l} [\mathcal{A}_p k(x^l, \cdot)] \quad \text{and } \mathbb{S}(q^l, p) = \|\phi_{q,p}^*\|_{\mathcal{H}^d}$$

1162 We compute:

$$\begin{aligned} \|\phi_{q,p}^*\|_{\mathcal{H}^d}^2 &= \langle \phi_{q,p}^*, \phi_{q,p}^* \rangle_{\mathcal{H}^d} \\ &= \mathbb{E}_{x_i^l} \mathbb{E}_{x_j^l} \left[ (k(x_i^l, x_j^l) \nabla_{x_j^l} \log p(x_j^l) + \nabla_{x_j^l} k(x_i^l, x_j^l)) \cdot (k(x_j^l, x_i^l) \nabla_{x_i^l} \log p(x_i^l) + \nabla_{x_i^l} k(x_j^l, x_i^l)) \right] \end{aligned}$$

1163 And obtain

$$\mathbb{S}(q^l, p) = \mathbb{E}_{x^l \sim q^l} \left[ \text{Tr} \left( \frac{\phi^*(x^l)}{\|\phi^*\|} \nabla_{x^l} \log p(x^l)^T + \nabla_{x^l} \frac{\phi^*(x^l)}{\|\phi^*\|} \right) \right] = \mathbb{E}_{x^l \sim q^l} \left[ \frac{\phi(x^l)^T \nabla_{x^l} \log p(x^l) + \text{Tr}(\nabla_{x^l} \phi(x^l))}{\|\phi^*\|} \right]$$

1164

□

### 1165 8.3.2 Intractability of KSD, FD

1166 Even though Fisher Divergence  $\mathbb{F}(q^l, p)$  and Kernelized Stein Discrepancy  $\mathbb{S}(q^l, p)$  follow the same  
1167 trend with the Stein Identity, they cannot be used as convergence metrics as they require computing  
1168  $\nabla_{x^l} \log q^l(x^l)$  which will require computing a jacobian on  $\nabla_{x^{l-1}} \phi(x^{l-1})$ .

1169 *Proof.* Note that  $\forall x, x'$  being i.i.d. draws from  $q^l$ ,  $\mathbb{F}(q^l, p) = \mathbb{E}_{x^l} [\nabla_x \log q^l(x) - \nabla_x \log p(x)]$   
1170 and  $\mathbb{S}(q^l, p) = \mathbb{E}_{x, x' \sim q^l} \left[ (\nabla_x \log q^l(x) - \nabla_x \log p(x))^T k(x, x') (\nabla_{x'} \log q^l(x') - \nabla_{x'} \log p(x')) \right]$ .  
1171 Computing either is possible thanks to the closed form expression of the log density 12, which  
1172 circumvents the intractability issue encountered due to the score of  $q^l, \forall l \in [0 \dots L]$ :

$$\nabla_{x^l} \log q^l(x^l) = \frac{\partial \log q^l(x^l)}{\partial x^l} = \left( \frac{\partial x^l}{\partial x^{l-1}} \right)^{-1} \cdot \frac{\partial \log q^l(x^l)}{\partial x^{l-1}} = \left( \nabla_{x^{l-1}} \phi(x^{l-1}) \right)^{-1} \cdot \nabla_{x^{l-1}} \log q^l(x^l)$$

$$\text{Such that } \log q^l(x^l) = \log q^{l-1}(x^{l-1}) - \epsilon \log \left| \det \nabla_{x^{l-1}} \phi(x^{l-1}) \right|$$

$$\text{So } \nabla_{x^l} \log q^l(x^l) = \underbrace{\left( \nabla_{x^{l-1}} \phi(x^{l-1}) \right)^{-1}}_{\text{Intractable}} \cdot \left( \nabla_{x^{l-1}} \log q^{l-1}(x^{l-1}) - \epsilon \nabla_{x^{l-1}} \text{Tr}(\nabla_{x^{l-1}} \phi(x^{l-1})) \right)$$

1173

□

### 1174 8.4 MH-augmented SVGD Density (Eq ??)

**Proposition 8.1.** *The MH-augmented density over particles after incorporating MH correction is:*

$$q_{\theta}^{\text{MH},l}(x^l) = \alpha_{\theta_{2,3}}^l q_{\theta}^{\text{MH},l-1}(x^{l-1}) |\det \nabla_{x^l} \phi_{\theta_2}(x^l)|^{-1} + (1 - \alpha_{\theta_{2,3}}^l) q_{\theta}^{\text{MH},l-1}(x^{l-1}), \quad \text{with } q_{\theta_1}^{\text{MH},0} = q_{\theta_1}^0$$

1175 *Proof.* We prove the statement above by induction. We leverage:

$$q_{\theta}^{\text{MH},l} = \alpha_{\theta_{2,3}}^{l-1} q_{\theta}^{\text{MH},l-1}(x^{l-1}) \left| \det(I + \epsilon \nabla_{x^{l-1}} \phi(x^{l-1})) \right|^{-1} + (1 - \alpha_{\theta_{2,3}}^{l-1}) q_{\theta}^{\text{MH},l-1} \quad (14)$$

$$\begin{cases} \text{Case } l = 0: & q_\theta^{\text{MH},0} = q_\theta^0 \\ \text{Case } l = 1: & q_\theta^{\text{MH},1}(x^1) = \alpha_{\theta_{2,3}}^0 q_\theta^0(x^0) \left| \det(I + \epsilon \nabla_{x^0} \phi(x^0)) \right| + (1 - \alpha_{\theta_{2,3}}^0) q_\theta^0(x^0) \end{cases}, \text{with } a_{1:L} = a_1$$

1176 On the other hand, evaluating  $q_\theta^{\text{MH},1}(x^1, a_1)$  by marginalizing over  $a_1 \in \{0, 1\}$ :

$$\begin{cases} q_\theta^{\text{MH},1}(x^1, a_1 = 0) = q_\theta^0(x_0)(1 - \alpha_{\theta_{2,3}}^1) \\ q_\theta^{\text{MH},1}(x^1, a_1 = 1) = q_\theta^0(x_0)(\alpha_{\theta_{2,3}}^1) \left| \det(I + \epsilon \nabla_{x_0} \phi(x_0)) \right| \end{cases}$$

1177 Therefore:

$$q_\theta^{\text{MH},1}(x^1) = q_\theta^0(x_0) \alpha_{\theta_{2,3}}^1 \left| \det(I + \epsilon \nabla_{x_0} \phi(x_0)) \right| + q_\theta^0(x_0)(1 - \alpha_{\theta_{2,3}}^1)$$

1178 Case  $l > 1$ : We assume

$$q_\theta^{\text{MH},l}(x^l) = q_\theta^l(x^l) \prod_{l=1}^l \alpha_{\theta_{2,3}}^l + \sum_{a_{1:l} \neq 1} q_\theta^{\text{MH},l}(x^l, a_{1:l}), \quad \forall l \in ]1, L_c] \quad (15)$$

1179 and show that it holds for  $q_\theta^{\text{MH},l+1}$ .

1180 We know that:

$$q_\theta^{\text{MH},L_c+1}(x^{L_c+1}) = \alpha_{\theta_{2,3}}^{L_c+1} q_\theta^{\text{MH},L_c}(x^{L_c}) \left| \det(I + \epsilon \nabla_{x^{L_c}} \phi(x^{L_c})) \right| + (1 - \alpha_{\theta_{2,3}}^{L_c}) q_\theta^{\text{MH},L_c}(x^{L_c}) \quad (16)$$

1181 We plug Eq. 15 into Eq. 16:

$$\begin{aligned} q_\theta^{\text{MH},L_c+1}(x^{L_c+1}) &= q_\theta^{\text{MH},L_c}(x^{L_c}) \times \left[ \alpha_{\theta_{2,3}}^{L_c+1} \left| \det(I + \epsilon \nabla_{x^{L_c}} \phi(x^{L_c})) \right| + (1 - \alpha_{\theta_{2,3}}^{L_c}) \right] \\ &= \left[ q_\theta^{L_c}(x^{L_c}) \prod_{l=1}^{L_c} \alpha_{\theta_{2,3}}^l + \sum_{a_{1..L_c} \neq 1} q_\theta^{\text{MH},L_c}(x^{L_c}, a_{1..L_c}) \right] \times \left[ \alpha_{\theta_{2,3}}^{L_c+1} \left| \det(I + \epsilon \nabla_{x^{L_c}} \phi(x^{L_c})) \right| + (1 - \alpha_{\theta_{2,3}}^{L_c}) \right] \\ &= \underbrace{q_\theta^{L_c}(x^{L_c}) \prod_{l=1}^{L_c+1} \alpha_{\theta_{2,3}}^l \left| \det(I + \epsilon \nabla_{x^{L_c}} \phi(x^{L_c})) \right|}_{q_\theta^{L_c+1}(x^{L_c+1}) \prod_{l=1}^{L_c+1} \alpha_{\theta_{2,3}}^l} + \underbrace{q_\theta^{L_c}(x^{L_c}) \prod_{l=1}^{L_c} \alpha_{\theta_{2,3}}^l (1 - \alpha_{\theta_{2,3}}^{L_c})}_{q_\theta^{\text{MH},L_c+1}(x^{L_c+1}, a_{1..L_c}=1, a^{L_c+1}=1)} \\ &\quad + \underbrace{\sum_{a_{1..L_c} \neq 1} q_\theta^{\text{MH},L_c}(x^{L_c}, a_{1..L_c}) \alpha_{\theta_{2,3}}^{L_c+1} \left| \det(I + \epsilon \nabla_{x^{L_c}} \phi(x^{L_c})) \right|}_{q_\theta^{\text{MH},L_c+1}(x^{L_c+1}, a_{1..L_c} \neq 0, a^{L_c+1}=1)} + \underbrace{\sum_{a_{1..L_c} \neq 1} q_\theta^{\text{MH},L_c}(x^{L_c}, a_{1..L_c}) (1 - \alpha_{\theta_{2,3}}^{L_c})}_{q_\theta^{\text{MH},L_c+1}(x^{L_c+1}, a_{1..L_c} \neq 0, a^{L_c+1}=0)} \end{aligned}$$

1182  $\square$

1183 **9 Additional Results on Entropy Estimation**

1184 In the following, we provide implementation details and additional experiments for the toy experi-  
 1185 ments.

1186 **9.1 Optimized SVGD parameters**

1187 **9.1.1 Kernel Bandwidth**

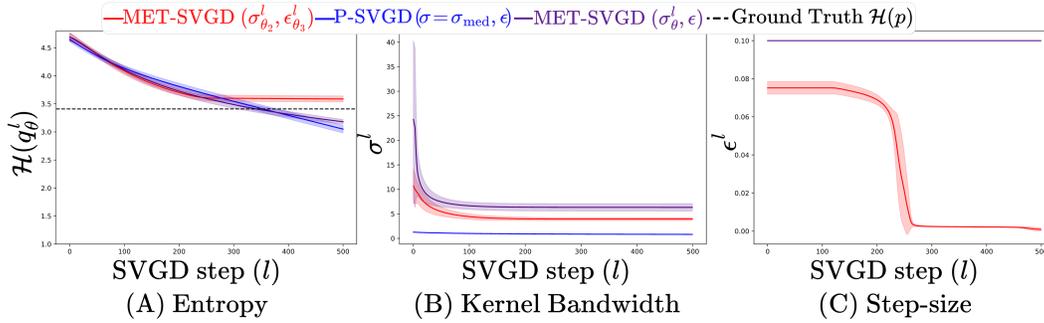


Figure 11: (A) Entropy, (B) RBF kernel bandwidth, and (C) step-size across SVGD steps. Target is the Gaussian target from Fig. ??.

1188 **9.2 Gaussian Targets**

1189 **Implementation Details** for Fig. ?? are reported in Tab. 4.

Algorithm	Parameter	Value
LD	Target $p$	$p = \mathcal{N}([-0.69, 0.8], [[1.13, 0.82], [0.82, 3.39]])$
P-SVGD	Number of Particles $M$	$M = 200$
MET-SVGD	Number of Steps $L$	$L = 1500$
	Initial Distribution $q^0$	$\mathcal{N}(0, 6I)$
LD	Learning Rate $\epsilon$	$\epsilon = 0.1$
P-SVGD	Kernel Bandwidth $\sigma$	$\sigma \in \{1, 5, \sigma_{\text{med}}\}$
	<b>Kernel Architecture</b>	
	Architecture	$\sigma_{\theta_2} = \text{GNN}(\{x_i^l\}_{i=1}^M; \theta_2)$
	# Layers	3
	Activation	{ReLU, Exponential, Truncate}
	<b>Learning Rate Architecture ??</b>	
MET-SVGD	Architecture	$\epsilon_{\theta_3}^l = \min(\epsilon_{\theta_3}^0, \epsilon_{\theta_3}^0 d^{l/s_{\theta_3}})$
	Initial Learning Rate $\epsilon_{\theta_3}^0$	0.1
	Decay Factor $d_{\theta_3}$	$5 \times 10^{-3}$
	Decay Scale $s_{\theta_3}$	$L$
	<b>Training Parameters</b>	
	Optimizer	Adam
	Learning Rate	$5 \cdot 10^{-3}$
	Epochs	300
	Loss	{KL Divergence}
Resources	GPU	Tesla V100-SXM2-32GB
	RAM	2 GB
	Per-epoch runtime	2.6 seconds

Table 4: Experimental setup for Fig. ??

1190 **Implementation Details** (Fig. 9) & (Fig. 10)

Parameter	Figure 9	Figure 10
Target $p$	$p = \mathcal{N}([-0.69, 0.8], [[1.13, 0.82], [0.82, 3.39]])$	
Number of Particles $M$	$M = 500$	$M = 500$
Number of Steps $L$	$L = 1000$	$L = 500$
Initial Distribution $q^0$	$\mathcal{N}(0, 6I)$	
Kernel Architecture		
Architecture	$C_{\theta_2} = \text{GNN}(\{x_i\}_{i=1}^M; \theta_2)$	$\sigma_{\theta_2} = \text{GNN}(\{x_i\}_{i=1}^M; \theta_2)$
Kernel Type	Bilinear: $\frac{x_i^T x_j}{C_{\theta_2}} + 1$	RBF: $\exp(-\ x_i - x_j\ ^2 / 2\sigma^2)$
Learning Rate Architecture ??		
Architecture	$\epsilon_{\theta_3}^l = \min(\epsilon_{\theta_3}^0, \epsilon_{\theta_3}^0 d^{l/s_{\theta_3}})$	
Initial Learning Rate $\epsilon_{\theta_3}^0$	$\epsilon_{\theta_3}^0 = 0.1$	
Decay Factor $d_{\theta_3}$	$d_{\theta_3} = 5 \times 10^{-3}$	
Decay Scale $s_{\theta_3}$	$s_{\theta_3} = L$	
Training Parameters		
Optimizer	Adam	
Learning Rate	$5 \cdot 10^{-3}$	
Epochs	300	

Table 5: Experimental setup for Figs. 9 and 10.

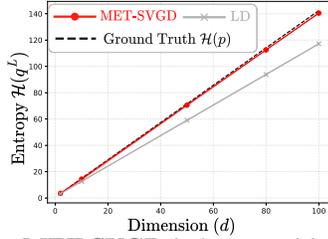


Figure 12: **MET-SVGD** is less sensitive to the  $Tr$  approximation than LD. Target is a slanted Gaussian (details in App. 6).

Parameter	LD	MET-SVGD
Number of particles	$M = 100$	
Number of iterations	$L = 1000$	
Target distribution	$p = \mathcal{N}(0, I_d), \quad d \in \{2, 10, 50, 80, 100\}$	
Initial distribution	$\mathcal{N}(0, 6I_d)$ (augmented)	
Algorithm-Specific Learned Parameters		
Learned parameter	Gaussian noise	Kernel variance $\sigma$
Learned LR	$\epsilon_{\theta_2}^l = \min(\epsilon_{\theta_2}^0, \epsilon_{\theta_2}^0 d^{l/s_{\theta_2}})$	
Training Parameters		
Optimizer	Adam	
Learning rate	$5 \cdot 10^{-3}$	
Epochs	300	

Table 6: Experimental setup for Fig. 12

1191 **Implementation Details (Fig. 12) Langevin Dynamics.** In Fig. 12, we show that we can learn the  
1192 Langevin dynamics' parameters  $\theta = \{\epsilon_{\theta_2}, \mu_{\theta_3}\}$  such that  $x^{l+1} = x^l + \epsilon_{\theta_2} \nabla_{x^l} \log p(x^l) + \sqrt{2\epsilon_{\theta_2}} \xi_{\theta_3}$   
1193 end-to-end by minimizing the reverse KL-Divergence:

$$\theta^* = \arg \min_{\theta} \mathbb{E}_{x^L \sim q_{\theta}^L} [\log q_{\theta}^L(x^L) - \log p(x^L)] \quad \text{s.t.} \quad \epsilon_{\theta_3}^l \leq \epsilon_{\text{UB}}^l, \quad \forall l \in [0, L-1]$$

1194 Here  $q^L = q^0 + \epsilon \sum_{l=0}^{L-1} \text{Tr}(\nabla_{x^l}^2 p(x^l))$ , where the trace is approximated via the Hutchinson estimator  
1195 (Eq. 7.5). We show SVGD is less sensitive to this approximation than LD in high dimensions.  
1196 **Scalability.** (a) Multivariate Gaussian. In Fig. 13, we visualize the learnt SVGD learning rate  
1197 for the setup in Fig. 8 in the main paper. The target is a  $d$ -dimensional multivariate Gaussian  
1198  $p(x) = \mathcal{N}(x; 0, I_d)$ . For each method, 100 particles are initialized from  $\mathcal{N}(x; 2\mathbf{1}, 2I_d)$ , where  
1199  $\mathbf{1} \in \mathbb{R}^d$  denotes the vector of ones. This is a standard benchmark illustrating the diminishing variance  
1200 issue of SVGD. We show that MET-SVGD outperforms other baselines in high dimensional spaces  
1201 as measured by the entropy and the variance across dimensions (Fig. 13-a,b). The SVGD repulsive  
1202 force is large during the first updates, preventing the particles from collapsing, unlike other baselines  
1203 (Fig. 13-c). P-SVGD is not scalable due to a missing term in the entropy (see Sec. 8.3). This is  
1204 subsequently fixed in the MET-SVGD update.

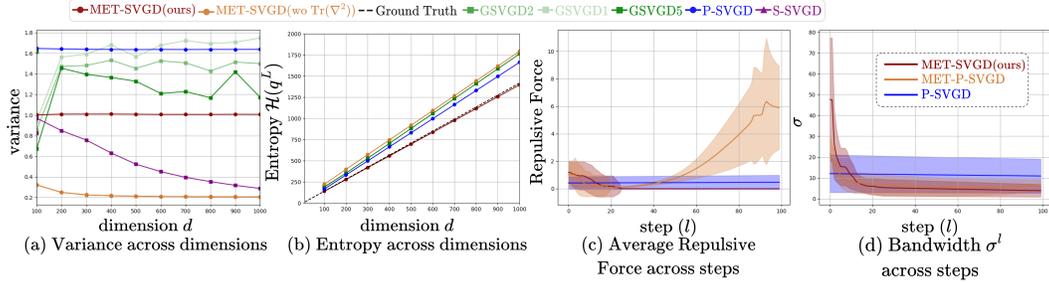


Figure 13: Scalability Results: **MET-SVGD** achieves higher accuracy on both (a) Entropy Estimation and (b) Variance across dimensions

1205 **Accelerating Convergence.** In Fig. 14, we show that, for the setup of Fig. 10, convergence can be  
 1206 accelerated either by (1) adding a regularization on the decay rate in the optimization objective (see  
 1207 Sec.2.2) or (2) randomizing the maximum number of steps during training (Eq. ??). We observe a  
 1208 quicker drop in the kernel variance (particles are less correlated initially).

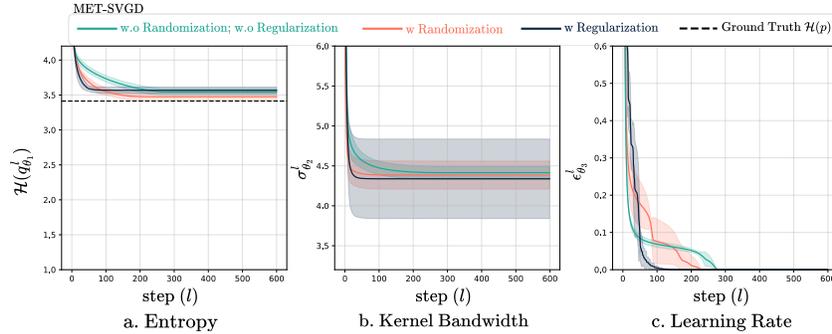


Figure 14: Accelerated convergence via regularization and number of SVGD steps randomization for the same MET-SVGD setup in Tab.6

1209 **9.3 Gaussian Mixture Model**

1210 **9.3.1 Experiment on 2D GMM with moving component**

1211 In the following, we provide implementation details and additional experiments for the toy experi-  
 1212 ments where the target is a 2D GMM.

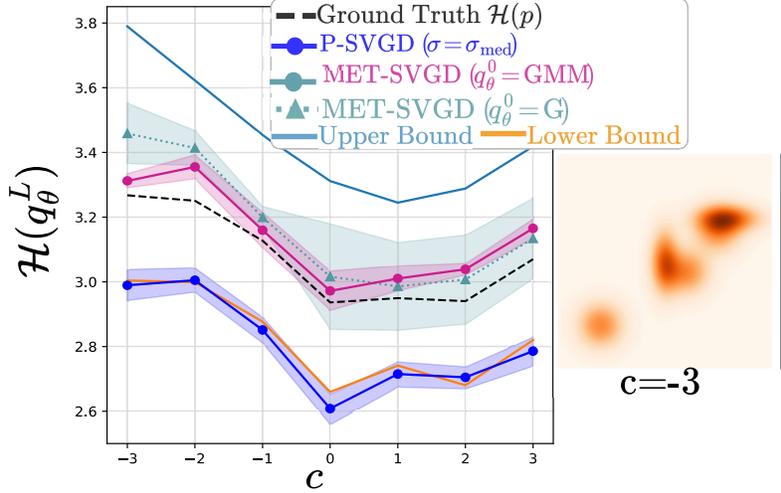


Figure 15: **MET-SVGD** outperforms **P-SVGD** on GMM entropy estimation.

1213 **Implementation Details** for (Fig. 15) are reported in Tab.7.

Parameter	P-SVGD	MET-SVGD
Distribution type	GMM with 5 components	
Fixed components	$\mu_1 = (0.0, 0.0), \Sigma_1 = 0.16I_2$ $\mu_2 = (3.0, 2.0), \Sigma_2 = I_2$ $\mu_3 = (1.0, -0.5), \Sigma_3 = 0.5I_2$ $\mu_4 = (2.5, 1.5), \Sigma_4 = 0.5I_2$	
Mobile component	$\mu_5 = (c, c), \Sigma_5 = 0.5I_2$ where $c \in [-3, 3]$	
Dimension	$d = 2$	
Number of particles	$M \in \{50, 100, 500\}$	
Number of iterations	$L = 1500$	
<b>Initial Distribution Settings</b>		
Setting 1	$\mathcal{N}(0, I_2)$	
Setting 2	GMM( $K = 10$ ) with $\mu_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{U}([-4, 4]^2)$ and $\Sigma_k = I_2 \forall k$	
<b>Algorithm-Specific Parameters</b>		
Kernel variance	$\sigma \in \{1, 5, \text{median}\}$	$\sigma = \text{GNN}(\{x_i\}_{i=1}^M; \theta_2)$
Learning rate	$\epsilon = 0.1$	$\epsilon_{\theta_3}^l = \min(\epsilon_{\theta_3}^0, \epsilon_{\theta_3}^0 d^{l/s_{\theta_3}})$
Base learning rate	-	$\epsilon_{\theta_3}^0 = 0.1$
Decay factor	-	$d_{\theta_3} = 5 \times 10^{-3}$
Decay scale	-	$s_{\theta_3} = L$
<b>Training Parameters</b>		
Optimizer	Adam	
Learning rate	$5 \cdot 10^{-3}$	
Epochs	300	

Table 7: Experimental setup for Fig. 15

1214 **Qualitative results** for different  $c$  values, as well as the KL-divergence, entropy, kernel Bandwidth  
1215 and step-size are reported in Fig. 16. We observe that **P-SVGD** with  $\sigma_{\text{med}}$  has poor convergence. In  
1216 fact,  $\sigma_{\theta_2}$  shows a different trend entirely. Whereas  $\epsilon_{\theta_3}$  converges consistently across all configurations.

1217 Additionally, we show in Fig. 17 that the number of particles is key to an accurate, low variance  
1218 estimation. Adding more particles significantly helps improve the accuracy. **P-SVGD** performs  
1219 poorly. In fact, the estimation is worse than a trivial lower bound that we derive as follows:

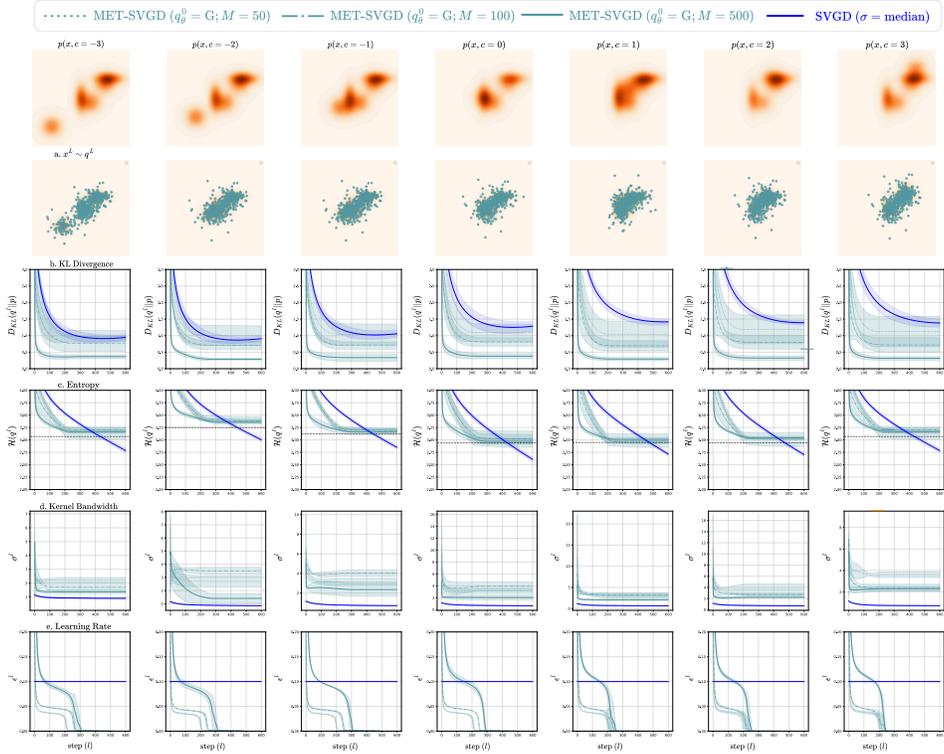


Figure 16: Results on entropy estimation for different  $c$  configurations. (a) KLD, (b) Entropy, (c) Learnt kernel bandwidth and (d) SVGD step-size.

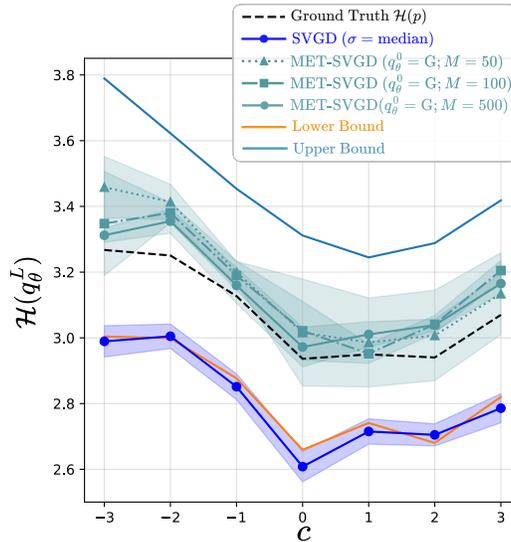


Figure 17: Entropy results on a 2D GMM with a moving component(Fig. 16). **MET-SVGD** significantly outperforms **P-SVGD**. Increasing the number of particles reduces the variances. Results are reported on 5 different seeds.

1220 Next we explain the derivation of the Upper & Lower Bounds for a GMM target as seen in Fig. 17 as  
 1221 follows:

The pdf of A GMM with  $K$  components is given by the formula:

$$p(x) = \sum_{i=1}^K \omega_i \mathcal{N}(x; \mu_i, C_i),$$

1222 where  $\{\omega_i\}_{i=1}^K$  are non-negative weighting coefficients such that  $\sum_i \omega_i = 1$  and  $\mathcal{N}(x; \mu_i, C_i)$   
 1223 is a gaussian density with mean  $\mu_i$  and covariance  $C_i$ . Note that the entropy generally can-  
 1224 not be calculated in closed form for GMM due to the logarithm of a sum of exponential func-  
 1225 tions (except for the special case of a single Gaussian density). We derive two trivial bounds  
 1226 to assess the performance of the different baselines. We start by deriving the lower bound  
 1227  $\text{LB}(x) = -\sum_{i=1}^K \omega_i \log(\sum_{j=1}^K \omega_j \mathcal{N}(\mu_i; \mu_j, C_i + C_j))$ .

*Proof.*

$$\begin{aligned} \mathcal{H}(x) &= -\sum_{i=1}^K \omega_i \int_{\mathbb{R}^N} \mathcal{N}(x; \mu_i, C_i) \log p(x) dx \\ &\geq -\sum_{i=1}^K \omega_i \log \left[ \int_{\mathbb{R}^N} \mathcal{N}(x; \mu_i, C_i) p(x) dx \right] \quad \text{by Jensen inequality} \\ &\geq -\sum_{i=1}^K \omega_i \log \left[ \int_{\mathbb{R}^N} \mathcal{N}(x; \mu_i, C_i) \sum_{j=1}^K \omega_j \mathcal{N}(x; \mu_j, C_j) dx \right] \geq \underbrace{-\sum_{i=1}^K \omega_i \log \left[ \sum_{j=1}^K \omega_j \mathcal{N}(\mu_i; \mu_j, C_i + C_j) \right]}_{\text{LB}(x)} \end{aligned}$$

1228 **Upper-bound UB:** We prove that  $\mathcal{H}(p(x)) \leq \text{UB}$  by deriving the entropy of the Gaussian enveloping  
 1229 the target:

$$\text{UB} = \mathcal{H}(\mathcal{N}(\mu_G, C_G)), \quad \text{with} \begin{cases} \mu_G = \sum_{i=1}^L \omega_i \mu_i \\ C_G = \sum_{i=1}^L \omega_i (C_i + \mu_i \mu_i^T) - \sum_{i=1}^L \sum_{j=1}^L \omega_i \omega_j \mu_i \mu_j^T \end{cases}$$

1230 *Proof.* Consider a Gaussian mixture model  $p(x) = \sum_{i=1}^L \omega_i \mathcal{N}(x; \mu_i, C_i)$ . The mean is  $\mu =$   
 1231  $\mathbb{E}[p(x)] = \sum_{i=1}^L \omega_i \mu_i$ . The covariance can be computed as:

$$\begin{aligned} C &= \mathbb{E}[(x - \mu)(x - \mu)^T] = \mathbb{E}[xx^T] - \mu\mu^T = \int_x \left( \sum_{i=1}^L \omega_i \mathcal{N}(x; \mu_i, C_i) \right) xx^T dx - \mu\mu^T \\ &= \sum_{i=1}^L \omega_i \int_x xx^T \mathcal{N}(x; \mu_i, C_i) dx - \mu\mu^T \end{aligned}$$

1232 For a single Gaussian component,  $C_i = \mathbb{E}[(x - \mu_i)(x - \mu_i)^T] = \mathbb{E}[xx^T] - \mu_i \mu_i^T$ , which means  
 1233  $\mathbb{E}_{x \sim \mathcal{N}(x; \mu_i, C_i)}[xx^T] = C_i + \mu_i \mu_i^T$ . Substituting this in the above:

$$\begin{aligned} C &= \sum_{i=1}^L \omega_i (C_i + \mu_i \mu_i^T) - \mu\mu^T \\ &= \sum_{i=1}^L \omega_i (C_i + \mu_i \mu_i^T) - \left( \sum_{i=1}^L \omega_i \mu_i \right) \left( \sum_{j=1}^L \omega_j \mu_j^T \right) = \sum_{i=1}^L \omega_i (C_i + \mu_i \mu_i^T) - \sum_{i,j=1}^L \omega_i \omega_j \mu_i \mu_j^T = C_G \end{aligned}$$

1234 Since a Gaussian distribution maximizes entropy among all distributions with the same mean and  
 1235 covariance, we have  $\mathcal{H}(p(x)) \leq \mathcal{H}(\mathcal{N}(\mu_G, C_G)) = \text{UB}(x)$ .

1236 **9.3.2 Targets with distant modes**

1237 **Implementation Details (Fig. 18)** are reported in Tab.8. We show that the divergence control  
 1238 heuristic based on eliminating particles further than 3 standard deviations of the initial distribution  
 1239 mean exacerbates mode collapse is already an issue when using the reverse KL-divergence.

1240 **Qualitative results** particles from the initial distribution are visualized in (i) a. We apply the  
 1241 truncation heuristic to different setups (**P-SVGD** with 0 steps, ie with only a learnable initial  
 distribution and **P-SVGD** with  $L = 140$  steps).

Algorithm	Parameter	Value
No SVGD P-SVGD MET-SVGD	Target $p$	GMM with 3 components $\mu_1 = (-4.0, 2.0), \Sigma_1 = I_2$ $\mu_2 = (4.0, 2.0), \Sigma_2 = I_2$ $\mu_3 = (0.0, -12.0), \Sigma_3 = 2I_2$
	Number of Particles	$M = 200$
	Number of Steps	$L = 1000$
	Initial Distribution	$\mathcal{N}(0, 6I)$
P-SVGD	Learning Rate $\epsilon$	$\epsilon = 0.1$
	Kernel	$\sigma \in \{1, 5, \text{median}\}$
	Bandwidth $\sigma$	
	<b>Kernel Architecture</b>	
	Architecture	$\sigma_{\theta_2} = \text{GNN}(\{x_i^l\}_{i=1}^M; \theta_2)$
	<b>Learning Rate Architecture</b>	
	Architecture	$\epsilon_{\theta_3}^l = \min(\epsilon_{\theta_3}^0, \epsilon_{\theta_3}^0 d^{l/s_{\theta_3}})$
	Initial Learning Rate $\epsilon_{\theta_3}^0$	$\epsilon_{\theta_3}^0 = 0.1$
	Decay Factor $d_{\theta_3}$	$d_{\theta_3} = 5 \times 10^{-3}$
	Decay Scale $s_{\theta_3}$	$s_{\theta_3} = L$
	<b>Training Parameters</b>	
	Optimizer	Adam
	Learning Rate	$5 \cdot 10^{-3}$
	Epochs	300
	Activation Functions	{Exponential(), Truncate(min, max)}

Table 8: Experimental setup for Fig. 18

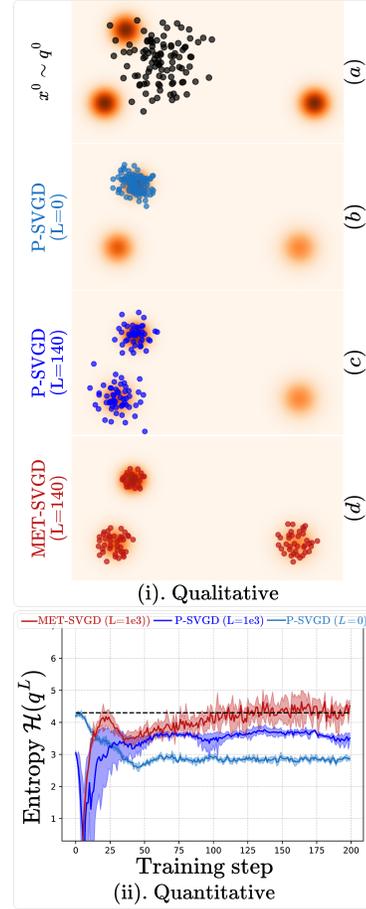


Figure 18: P-SVGD struggles with distributions with distant modes.

1242

1243 **Effect of the Initial Distribution.** In Fig. 19, we compare the results of using a Gaussian and a  
 1244 GMM with 10 components. We show that when using a GMM, less steps and particles are needed  
 1245 to learn the target. All experiments with a Gaussian initial distribution resulted in mode collapse.  
 1246 Using a GMM mitigates the mode collapse issue when learning the parameters using the reverse  
 1247 KL-divergence. In the future, we will explore training with the forward KLD while leveraging  
 1248 importance sampling.

1249 **Implementation details are in Tab.9.**

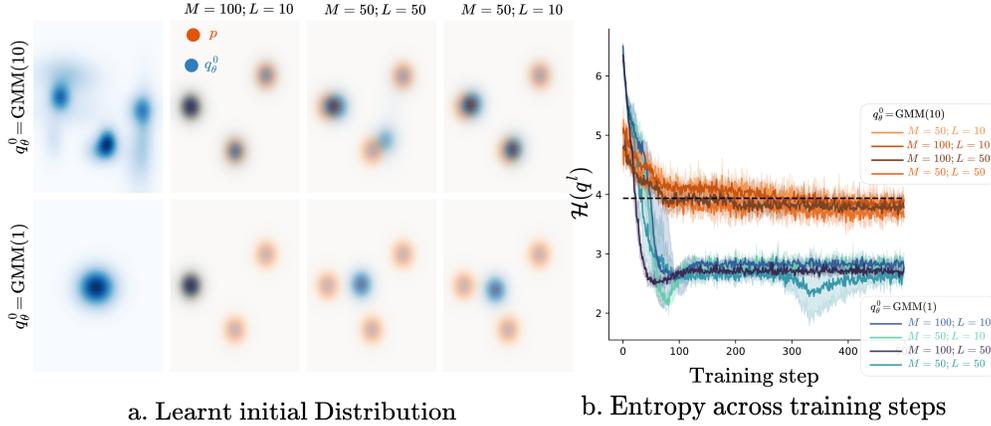


Figure 19: Effect of the initial distribution on a GMM setup.

Algorithm	Parameter	Value
MET-SVGD	Target $p$	GMM with 3 components (orange spots in figure)
	Number of Particles $M$ Number of Steps $L$	$M \in \{50, 100\}$ $L \in \{10, 50\}$
	Initial Distribution $q_\theta^0$	Two settings: GMM(1): Single component (left column) GMM(10): 10 components (right column)
	<b>Kernel Architecture</b> Architecture	$\sigma_{\theta_2} = \text{GNN}(\{x_i^l\}_{i=1}^M; \theta_2)$
	<b>Learning Rate Architecture</b> Architecture Initial Learning Rate $\epsilon_{\theta_3}^0$ Decay Factor $d_{\theta_3}$ Decay Scale $s_{\theta_3}$	$\epsilon_{\theta_3}^l = \min(\epsilon_{\theta_3}^0, \epsilon_{\theta_3}^0 d^l / s_{\theta_3})$ $\epsilon_{\theta_3}^0 = 0.1$ $d_{\theta_3} = 5 \times 10^{-3}$ $s_{\theta_3} = L$
	<b>Training Parameters</b> Optimizer Learning Rate Epochs Activation Functions	Adam $5 \cdot 10^{-3}$ 300 {Exponential(), Truncate(min, max)}

Table 9: Experimental setup for Fig. 19

## 1250 9.4 High Dimensional GMMs

1251 The goal of this experiment is to further assess the scalability of **MET-SVGD**.

1252 **Implementation Details** The target distribution is a mixture of 4 d-dimensional Gaussian distributions  
 1253  $p(x) = \sum_{k=1}^4 0.25 \mathcal{N}(x, \mu_k, I_d)$  with uniform mixture ratios. The first two coordinates of the mean  
 1254 vectors are equally spaced on a circle, while the other coordinates are set to 0 (Fig. 20.A-D). Particles  
 1255 are initialized from  $c\mathcal{N}(0, I_d)$  and only the first two dimensions need to be learned. In Fig. 20A, we  
 1256 show that **MET-SVGD** efficiently recovers the low-dimensional structure.

	Parameter	Value
	Target distribution Initial distribution	$p(x) = \sum_{k=1}^4 0.25\mathcal{N}(x, \mu_k, I_d)$ $q^0 = \mathcal{N}(0, I)$
Default SVGD parameters	Learning rate	$\epsilon_{\text{P-SVGD}} = 0.1$ $\epsilon_{\text{MET-SVGD}} = \text{nn.Parameter}(0.1)$
	Number of steps Number of particles	$L = 100$ $M = 100$
	Kernel variance	$\sigma_{\text{P-SVGD}} = \sqrt{\frac{\text{med}(\ x_i - x_j\ ^2)}{2 \ln M}}$ $\sigma_{\text{MET-SVGD}} = \text{nn.Parameter}(1.0)$
Training	Optimizer	Adam
	Learning rate	$10^{-2}$
	Epochs	500

Table 10: Experimental configuration for high-dimensional GMM results.

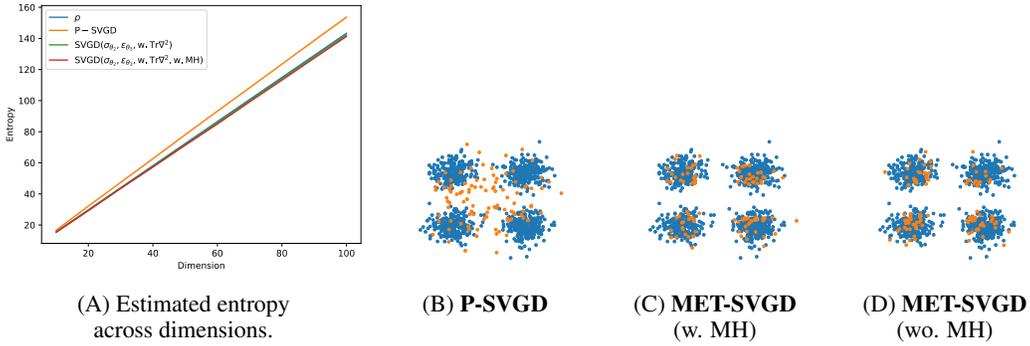


Figure 20: Scalability results. Target is a high-dimensional GMM. **MET-SVGD** successfully recovers the entropy of the low-dimensional GMM.

## 1257 10 Additional Results: Energy Based Models

1258 **Proposition 10.1** (Sec.3). *Training EBMs  $p_\theta(x) = \bar{p}_\phi(x)/Z$  via maximum likelihood ( $\mathcal{L}_{\text{ebm}}(\phi, \theta) =$*   
1259  *$-\mathbb{E}_{x \sim p_d}[\log p_\theta(x)]$ ) is intractable due to the partition function  $Z$ . When the sampler has a tractable*  
1260 *distribution  $q_\phi$ , a tight lower bound can be computed in return:  $\mathcal{L}_{\text{ELBO}}(\phi) = \mathbb{E}_{x \sim q}[\log \bar{p}_\phi(x)] -$*   
1261  *$\mathbb{E}_{x \sim p_d}[\log \bar{p}_\phi(x)] + \mathcal{H}(q_\phi)$  with  $p_d$  being the data distribution,*

1262 *Proof.* Given:

$$\mathcal{L}_{\text{ebm}}(\phi, \theta) = -\mathbb{E}_{x \sim p_d}[\log p_\theta(x)] = -\mathbb{E}_{x \sim p_d}[\log \bar{p}_\phi(x)] + \log Z(\theta).$$

1263 We bound the partition function using the KL-divergence:

$$\begin{aligned} \log Z(\theta) &\geq \log Z(\theta) - D_{\text{KL}}(q_\phi(x) \| p_\theta(x)) \\ &\geq \log Z(\theta) + \int_x q_\phi(x) \log \frac{p_\theta(x)}{q_\phi(x)} dx \\ &\geq \log Z(\theta) + \int_x q_\phi(x) \log \frac{\bar{p}_\phi(x)}{Z(\theta)} dx \\ &\geq \log Z(\theta) + \int_x q_\phi(x) \log \bar{p}_\phi(x) dx - \int_x q_\phi(x) \log Z(\theta) dx - \int_x q_\phi(x) \log q_\phi(x) dx \\ &\geq \log Z(\theta) + \mathbb{E}_{x \sim q_\phi}[\log \bar{p}_\phi(x)] - \log Z(\theta) + \mathcal{H}(q_\phi) \end{aligned}$$

$$\geq \mathbb{E}_{x \sim q_\phi} [\log \bar{p}_\phi(x)] + \mathcal{H}(q_\phi).$$

1264 Substituting back into the MLE objective:

$$\begin{aligned} \mathcal{L}_{\text{ebm}}(\phi, \theta) &= -\mathbb{E}_{x \sim p_d} [\log \bar{p}_\phi(x)] + \log Z(\theta) \\ &\geq -\mathbb{E}_{x \sim p_d} [\log \bar{p}_\phi(x)] + \mathbb{E}_{x \sim q_\phi} [\log \bar{p}_\phi(x)] + \mathcal{H}(q_\phi) \\ &\geq \mathcal{L}_{\text{ELBO}}(\phi). \end{aligned}$$

1265

□

## 1266 10.1 Synthetic Experiment: Moon Distribution

1267 We evaluate on the Moon dataset [Rezende and Mohamed, 2015b] with varying smoothness. **Implementation Details (Fig. 21)** are described in Tab. 11. **Performance:** As smoothness decreases,

	Parameter	Value
RED <sub>B</sub> O <sub>X</sub>	Target distribution	$p_\theta(x) = \frac{\exp f_\theta(x)}{Z_\theta}$
	Initial distribution	$f_\theta(x) = \text{MLP}_\theta(128, \text{Swish}, 128, \text{Swish}, 128, \text{Swish}, 1)$ $q^0 = \mathcal{N}([0, 0], 7I)$
RED <sub>B</sub> O <sub>X</sub> Default SVGD parameters	Learning rate	$\epsilon = \epsilon_{\theta_2}^l$ ( $l$ free learnable parameters)
	Number of steps	$L = 100$
	Number of particles	$m = 129$
	Kernel variance	$\sigma = \sigma_{\theta_2}^l$ ( $l$ free learnable parameters)
Training	Optimizer	Adam
	$\theta$ Learning rate	$10^{-3}$
	$\phi$ Learning rate	$10^{-2}$
	Epochs	1250
Resources	GPU	Tesla V100-SXM2-32GB
	RAM	2 GB
	Per-epoch runtime	2.6 seconds

Table 11: Experimental configuration for EBM results.

1268

1269 **MET-SVGD** consistently outperforms all baselines in terms of the MMD score [Dai et al., 2019a],  
 1270 where  $\text{MMD}(p, q) = \mathbb{E}_{x, x' \sim p} [\kappa(x, x')] + \mathbb{E}_{y, y' \sim q} [\kappa(y, y')] - 2\mathbb{E}_{x \sim p, y \sim q} [\kappa(x, y)]$ , s.t.  $\kappa(x, y) =$   
 1271  $\exp(-\|x - y\|^2 / 2\sigma^2)$

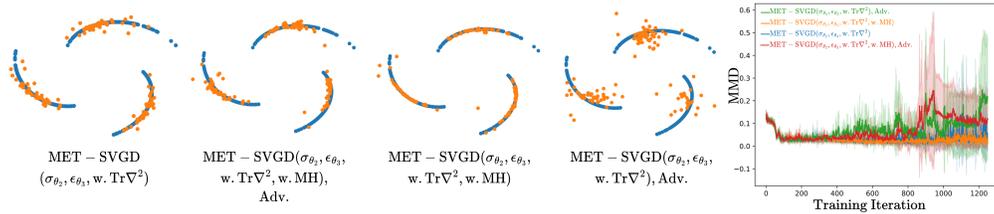


Figure 21: EBM Results. **MET-SVGD** outperforms **P-SVGD** and **LD** on learning EBMs to fit non-smooth data distributions.

## 1272 10.2 Image generation

1273 **Implementation Details (Fig. 23)** are reported in Tab. 12. All experiments were conducted on a  
 1274 single NVIDIA A100 80GB with 8GB allocated memory; average runtime was around 6 seconds per  
 1275 iteration (processing one batch of data composed of 64 particles). **Qualitative Results.** In Fig. 23,  
 1276 we visualize generated images sampled from the different models. **FID and Inception Score** In  
 1277 Fig. 23, we report the FID and IS scores for baselines: (1) LD trained with contrastive divergence:

$$\min_{\theta} \mathcal{L}_{\text{CD}}(\theta) = \min_{\theta} -\mathbb{E}_{x \sim p_d} [f_\theta(x)] + \mathbb{E}_{x \sim q} [f_\theta(x)], \quad (17)$$



Figure 22: Image generation using EBMs across different configurations.

	Parameter	Value
RED <sub>BOX</sub>	Target distribution	$p_\theta(x) = \frac{\exp f_\theta(x)}{Z_\theta}$
	Initial distribution	$f_\theta(x)$ is a WideResnet(28,10) network $q^0$ is a replay buffer initialized using a GMM whose modes are based on the class-conditional means and covariances
Default SVGD parameters	Learning rate	$\epsilon_{LD} = \epsilon_{P-SVGD} = 64$ (this number is divided by $m$ in the SVGD update formula)
	Number of steps	$\epsilon_{MET-SVGD} = \text{GNN}(\{x_i^l\}, \{\nabla_{x_i} \log p_\theta(x_i^l)\}; \theta_3)$
	Number of particles	$L = 5$
	Kernel variance	$L_c = 10$ $M = 64$ $\sigma_{LD} = 0$ $\sigma_{P-SVGD} = \sqrt{\frac{\text{med}(\ x_i - x_j\ ^2)}{2 \ln M}}$ $\sigma_{MET-SVGD} = \text{GNN}(\{x_i^l\}; \theta_2)$
Training	Optimizer	SGD
	$\theta$ Learning rate	$10^{-1}$ with 1000 iterations warm-up and decay at epochs 60, 120, and 180.
	$\phi$ Learning rate	$10^{-4}$
Resources	Epochs	200
	GPU	NVIDIA A100 80GB
	RAM	8 GB
	Per-iteration runtime	6 seconds

Table 12: Experimental configuration for EBM results.

1278 with  $q$  being the empirical distribution induced by the LD particles; (2) **P-SVGD** and **MET-SVGD**  
1279 variants trained adversarially by alternating between learning the sampler parameters

$$\min_{\phi} -\mathcal{L}_{\text{ELBO}}(\phi) = \max_{\phi} -\mathbb{E}_{x \sim p_d} [f_\theta(x)] + \mathbb{E}_{x \sim q} [f_\theta(x)] + \mathcal{H}(q_\phi), \quad (18)$$

1280 and minimizing the contrastive divergence

$$\min_{\theta} \mathcal{L}_{\text{CD}}(\theta) = \min_{\theta} \mathbb{E}_{x \sim q} [\log \bar{p}_\phi(x)] - \mathbb{E}_{x \sim p_d} [\log \bar{p}_\phi(x)] \quad (19)$$

1281 in Fig. 25A and Fig. 23C. (3) **MET-SVGD** variants trained adversarially using the ELBO loss for  
1282 both learning the sampler and the energy:

$$\min_{\theta} \max_{\phi} \mathcal{L}_{\text{ELBO}}(\theta, \phi) = \min_{\theta} \max_{\phi} \mathbb{E}_{x \sim q} [\log \bar{p}_\phi(x)] - \mathbb{E}_{x \sim p_d} [\log \bar{p}_\phi(x)] + \mathcal{H}(q_\phi), \quad (20)$$

1283 in Fig. 25B and Fig. 23D. The second setup (Fig. 25A, Fig. 23C) led to the best result. **Performance:**  
1284 The best FID (lowest) was obtained when both the kernel bandwidth and step-size are learnt. Compar-  
1285 ative results with better scalability are obtained with an adaptive number of steps  $L_c$ . Removing  
1286 the trace led to early divergence showcasing the importance of this correction. Both LD and **P-SVGD**  
1287 resulted in early divergence. This was also the case for the setup with MH due to high rejection rates  
1288 and hence limited number of steps (constrained by the GPU memory) leading to poor convergence to  
1289 the target. In the future, we will explore optimizing the memory usage to afford more steps. Also,  
1290 we find that the performance improvement obtained from learning the step-size on top of the kernel  
1291 bandwidth, *i.e.*, **MET-SVGD**( $\sigma_{\theta_2}, \epsilon_{\theta_3}$ ) vs. **MET-SVGD**( $\sigma_{\theta_2}$ ), is due to the learning the step-size  
1292 resulting in smoother energy landscapes (Fig. 24). The third setup (Fig. 25B, Fig. 23D), where both  
1293 the sampler and the energy are learnt using the ELBO (with the entropy term in learning the energy)  
1294 didn't work as well. In Fig. 24, we show that this is due to the score exploding frequently leading  
1295 to almost zero learning rates. We plan to address this by constraining the Lipschitz constant of the  
1296 deepnet in the future.

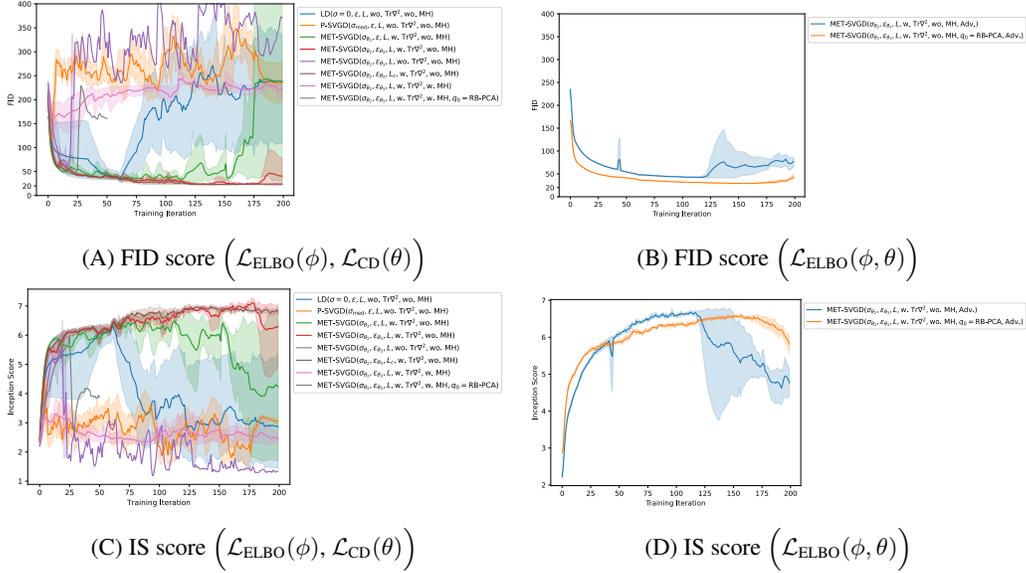


Figure 23: EBM Results. We report the FID and IS scores across training iterations for both (A-B) the set-up where the sampler ( $\phi$ ) is trained using  $\mathcal{L}_{\text{ELBO}}(\phi)$  and energy ( $\theta$ ) using  $\mathcal{L}_{\text{CD}}(\theta)$ , and (B-C) the setup where the sampler and the energy are learnt adversarially using  $\mathcal{L}_{\text{ELBO}}(\theta, \phi)$ .  $Tr\nabla^2$  stands for including the trace of Hessian term and  $q_0 = \text{RB-PCA}$  for initializing the replay buffer with samples obtained via a linear combination of the principal components of the data samples.

1297 **Smoothness.** In Fig. 24, we visualize the scores of the learnt distribution  $\nabla_x f_\theta(x)$  across training  
 1298 iterations. The setups with the lowest FID and highest IS score are associated with the smoothest  
 1299 landscapes, *i.e.*, lowest scores. This is the case of **MET-SVGD** $(\sigma_{\theta_2}, \epsilon_{\theta_3}, w, Tr\nabla^2)$  and **MET-**  
 1300 **SVGD** $(\sigma_{\theta_2}, \epsilon_{\theta_3}, L_C, w, Tr\nabla^2)$ . The diverging setups are associated with frequently exploding  
 1301 scores.

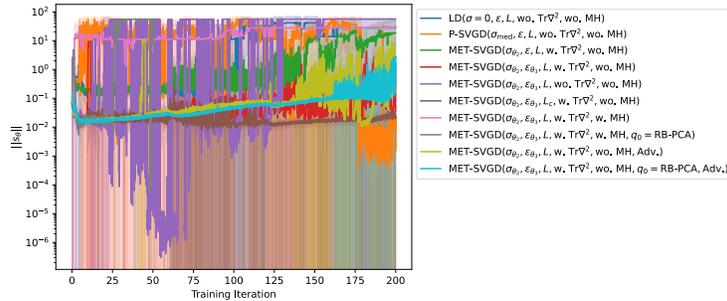


Figure 24: EBM results. The L2-norm of the learnt EBM score  $\nabla_x f_\theta(x)$ .

1302 **Trainable SVGD Hyperparameters.** In Fig. 25, we visualize the SVGD step-size and kernel-  
 1303 bandwidth across training iterations. We observe that  $\sigma_{\text{med}}$  is frequently higher than the learned ones.  
 1304 The learned step-size is also higher than the **P-SVGD** one in setups that led to the best FID. In setups  
 1305 with high FID, we observe that  $\epsilon_{\theta_3}$  is associated with high variance: it frequently becomes very small.  
 1306 This is mostly driven by the score of the energy.

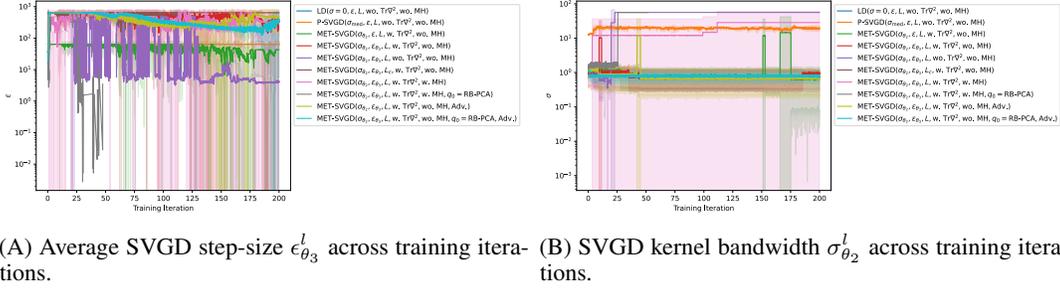


Figure 25: EBM Results. Visualization of the learnt kernel bandwidth and step-size across training iterations.

1307 **11 Additional Results: MaxEntr RL**

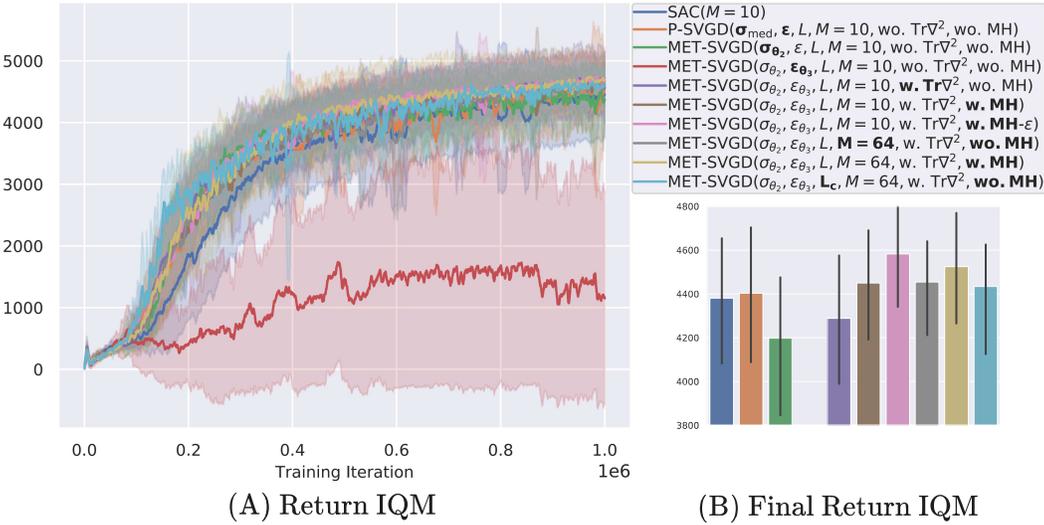


Figure 26: Return IQM for Walker2d-v2. MH variants yield the best returns.

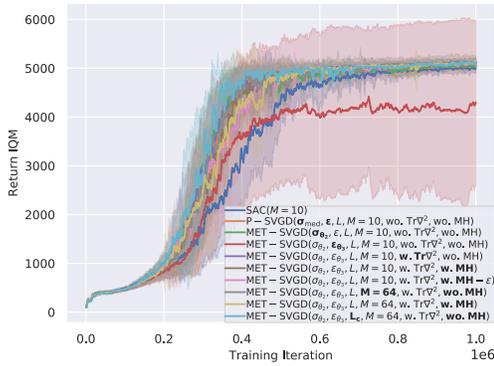
1308 **Implementation Details** for (Fig. 26) are reported in Tab. 13

1309 **Performance.** In Fig. 27A and Fig. 27B, we report the Inter Quantile Mean (IQM) return values  
 1310 averaged over 5 runs, where every run is the average of 10 evaluations of the policy. We refer to the  
 1311 approach leveraging **P-SVG**D as proposed by Messaoud et al. [2024] as  $S^2AC$ , and our approach as  
 1312  $S^2AC^+$ . We follow the same convention form EBM experiments for naming the different methods:  
 1313 for the different  $S^2AC^+$  variants, we only include arguments that are different from the  $S^2AC$  setup.  
 1314 The default parameters for  $S^2AC$  are  $(d_{\theta_1}^0, \epsilon = 1e^{-4}, M = 10, L = 3)$ , divergence control w.  
 1315 particles truncation, wo.  $Tr\nabla^2$ ). In the Humanoid environment (Fig. 27A), adding the missing trace  
 1316 of Hessian term resulted in faster convergence. Increasing the number of particles and incorporating  
 1317 MH for divergence control led to improved performance: better exploration through more particles  
 1318 and exploitation through the MH step.

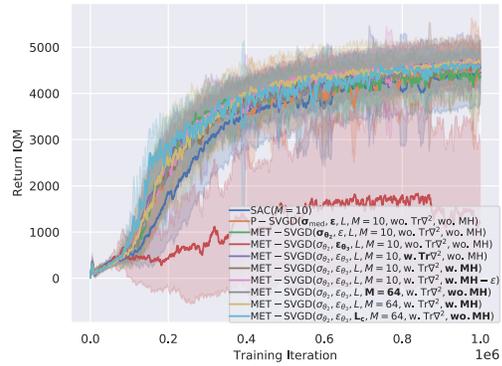
1319 In Walker environment (Fig. 27B), MH helped achieve higher return. We also see that  $S^2AC(\sigma_{\theta_2}, \epsilon_{\theta_3})$   
 1320 failed because the learning rate became very small for many iterations (Fig. 33A), which is visible in  
 1321 the high variance of the score norm (Fig. 34). This is the same phenomenon observed in the other  
 1322 EBM experiments.

Table 13: Hyperparameters

	Hyperparameter	Value
Training	Optimizer	Adam
	Actor and Critic Learning rate	$10^{-4}$ for Humanoid and $10^{-3}$ for all other environments
	Batch size	100
Deepnet	Number of hidden layers	2 Critic and 3 Actor
	Number of hidden units per layer	256
	Nonlinearity	ELU
RL	Target smoothing coefficient	0.005
	Discount $\gamma$	0.99
	Target update interval	1
	Entropy weight $\alpha$	0.2
	Replay buffer size $ \mathcal{D} $	$10^6$
SVGD	Initial distribution	$q_0 = \mathcal{N}(\mu_\theta, \text{diag}(\sigma_\theta))$
	Number of steps	$L = 3$
	Number of particles	$M = 10$
	Kernel variance	$\sigma^2 = \frac{\sum_{i,j} \ a_i - a_j\ ^2}{4(2 \log m + 1)}$ $\sigma = \text{GNN}(s_t, \{x_i^t\}, \{\nabla_{x_i^t} \log p(x_i^t)\}; \theta_2)$
	Learning rate	$\epsilon = 0.1$ $\epsilon = \text{GNN}(s_t, \{x_i^t\}, \{\nabla_{x_i^t} \log p(x_i^t)\}; \theta_3)$

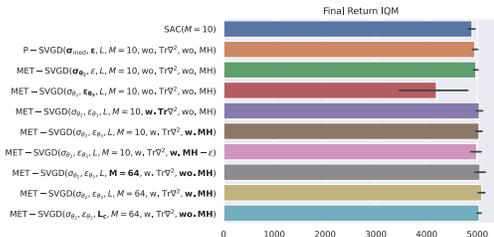


(A) Humanoid's IQM Return.

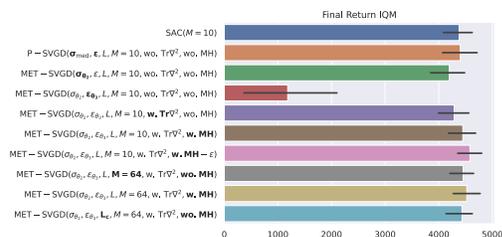


(B) Walker's IQM Return.

Figure 27: IQM return scores across environments.



(A) Humanoid's Final IQM Return.



(B) Walker's Final IQM Return.

Figure 28: Final IQM return scores across environments.

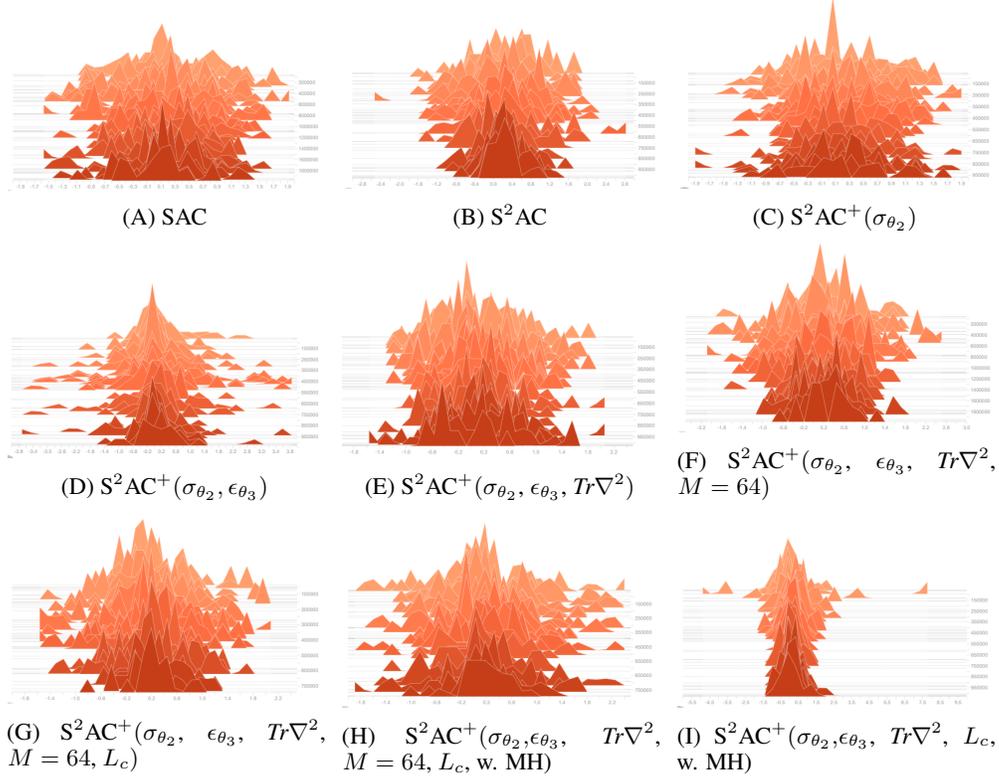
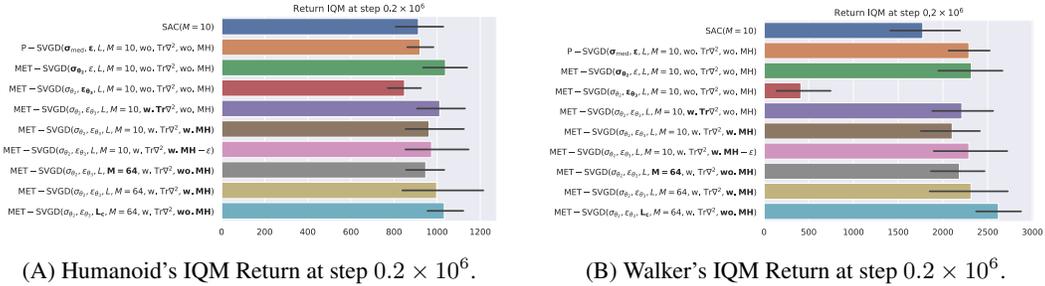


Figure 30: Histogram of the mean of  $q^0$  across training iterations (Humanoid env).



(A) Humanoid’s IQM Return at step  $0.2 \times 10^6$ .

(B) Walker’s IQM Return at step  $0.2 \times 10^6$ .

Figure 29: IQM return scores at step  $0.2 \times 10^6$  across environments.

1323 **Trainable SVGD Parameters (Humanoid Env)** In Fig. 30, we visualize a histogram of the mean of  
 1324 **initial distribution**  $q_{\theta_1}^0$  across training iterations. We observe that the mean has several components  
 1325 outside the  $[-1,1]$  range of valid actions. While the actions are truncated to satisfy the constraints,  
 1326 this still limits the exploration as many particles would end-up having  $-1/1$  as values. This trend  
 1327 is exacerbated across  $S^2AC^+$  variants, especially for the cases of learnable step-size  $\epsilon_{\theta_3}$ , adaptive  
 1328 number of steps  $L_c$  and larger number of particles  $M = 64$ . In the future, we will explore mechanisms  
 1329 for constraining the support of the policy distribution to the valid range. This is not a trivial problem  
 1330 as the obvious solution of truncating the mean leads to vanishing gradients and modeling  $q_{\theta_1}^0$  as a  
 1331 distribution with a limited support (e.g., beta distribution) is not obvious as such a distribution is  
 1332 highly sensitive to parameters with big ranges. Also, enforcing the constraints in the Q-value through  
 1333 reward is not trivial as it can lead to non-smooth hard-to-learn landscapes. This poor exploration  
 1334 limits the effect of our contribution, as our approach helps explore better in the local neighborhood of  
 1335 the modes identified through exploration.

1336 In Fig. 33, we present a histogram of the learned **kernel bandwidth** across training iterations. Note  
 1337 that in these experiments  $\sigma_{\theta_2} \in \mathbb{R}^d$ . We observe that for certain dimensions the bandwidth was small

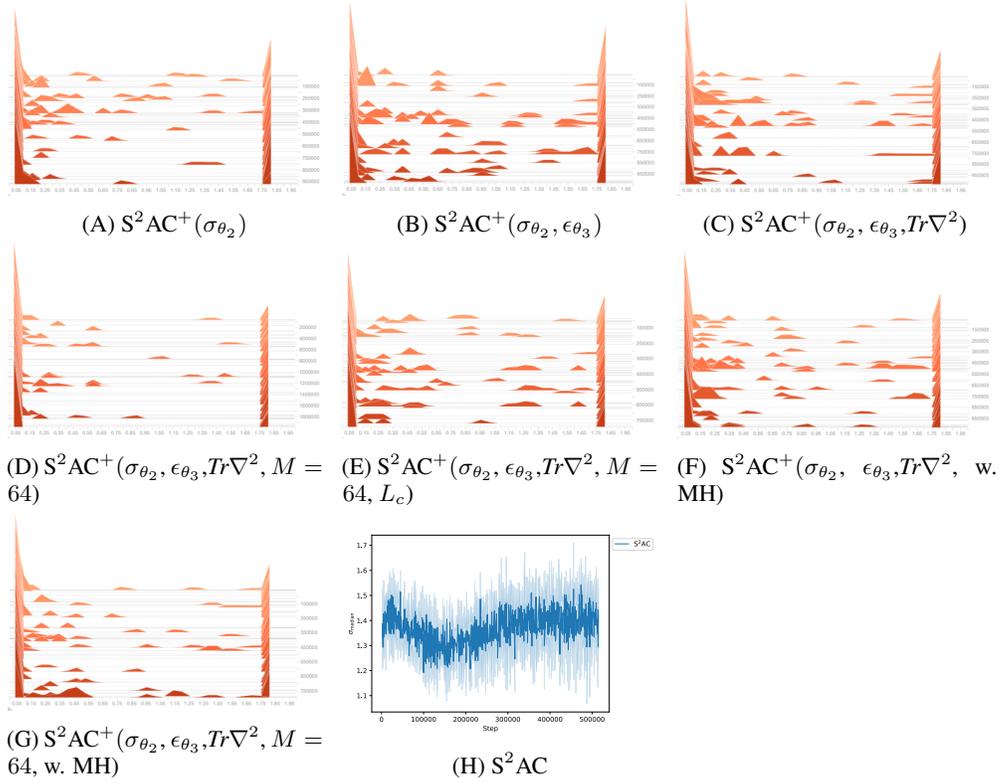


Figure 31: Histogram of the kernel bandwidth  $\sigma_{\theta_2}^l$

- 1338 indicating independant particles while for other states the particles were more interdependent (large  
 1339  $\sigma_{\theta_2}$  values). Also, note that the kernel bandwidth values for S2AC are consistently large.  
 1340 The SVGD step-size  $\epsilon_{\theta_3}^l$  is visualized in Fig. 32.

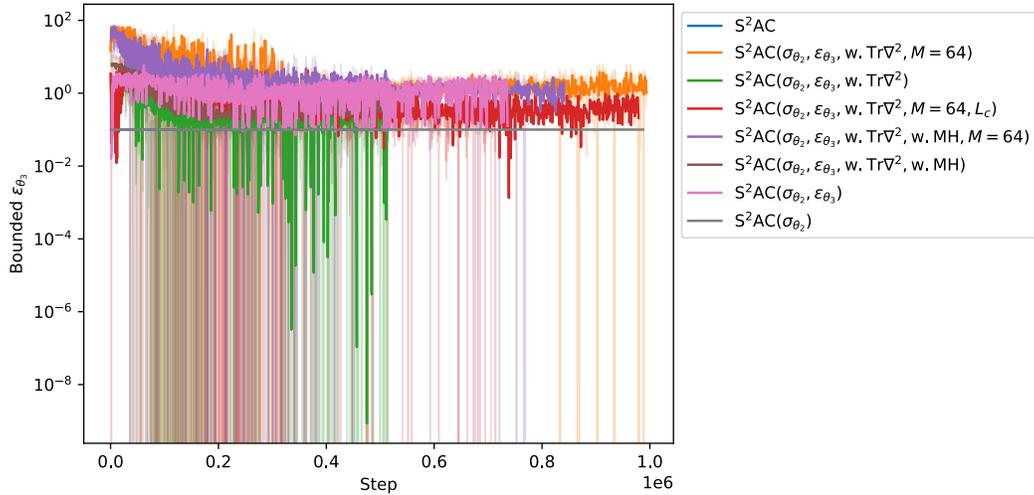


Figure 32: Learned step-size  $\epsilon_{\theta_3}^l$  for Humanoid env.

- 1341 **Trainable SVGD Parameters (Walker Env).** We visualize the scores and step-size in Fig. 33A and  
 1342 Fig. 34.

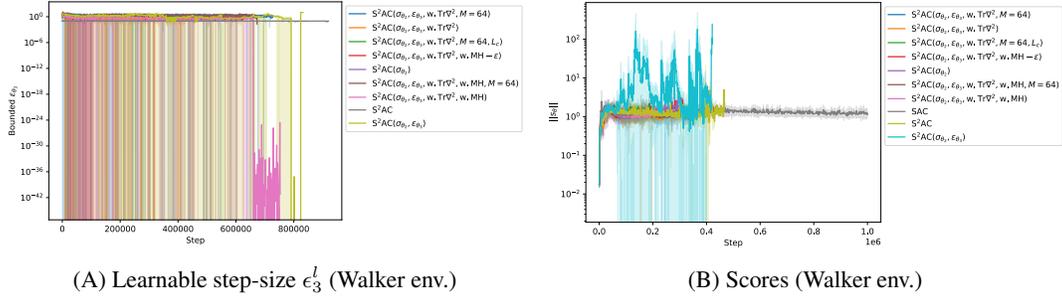


Figure 33: Learned step-size and scores in Walker env.

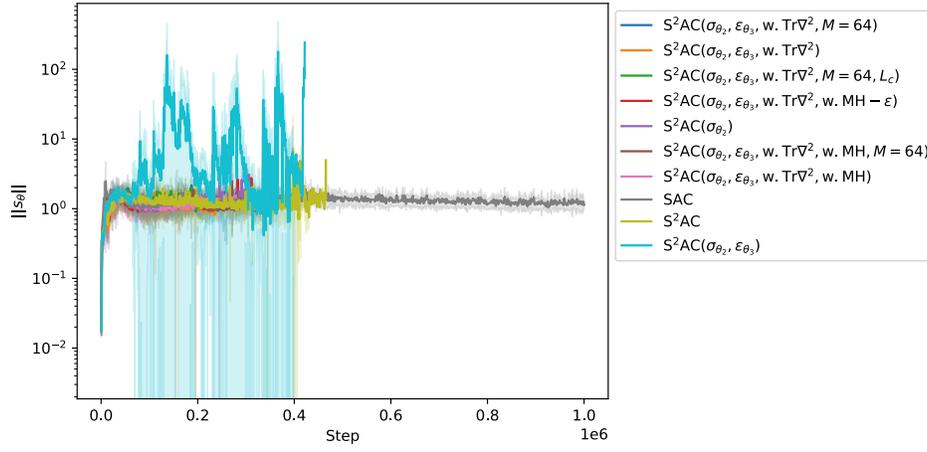


Figure 34: Scores (Walker env.)