
Provably Fast Finite Particle Variants of SVGD via Virtual Particle Stochastic Approximation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 SVGD is a popular particle-based variational inference algorithm with well studied
2 mean-field dynamics. However, its finite-particle behavior is far less understood.
3 Our work introduces the notion of *virtual particles* to develop novel stochastic ap-
4 proximations of mean-field SVGD dynamics in the space of probability measures,
5 that are exactly realizable using finite particles. As a result, we design two compu-
6 tationally efficient variants of SVGD (VP-SVGD and GB-SVGD) with provably
7 fast finite-particle convergence rates. Our algorithms are specific random-batch
8 approximations of SVGD which are computationally more efficient than ordinary
9 SVGD. We show that the n output particles of VP-SVGD and GB-SVGD, run
10 for T steps with batchsize K , are as good as i.i.d samples from a measure whose
11 Kernel Stein Discrepancy to the target is at most $O(d^{1/3}/(KT)^{1/6})$ under standard
12 assumptions. We prove similar results under a mild growth condition on the score
13 function, which is weaker than the assumptions of prior works. Our convergence
14 rates for the empirical measure (of the particles output by VP-SVGD and GB-
15 SVGD) to the target distribution enjoys a *double exponential improvement* over
16 the best known finite-particle analysis of SVGD. Furthermore, our results give the
17 *first known polynomial oracle complexity in dimension*, completely eliminating the
18 curse of dimensionality exhibited by previously known finite-particle rates.

19 1 Introduction

20 Sampling from a distribution over \mathbb{R}^d whose density $\pi^*(\mathbf{x}) \propto \exp(-F(\mathbf{x}))$ is known only upto
21 a normalizing constant, is a fundamental problem in machine learning [44, 19, 25] and statistics
22 [35, 31, 15]. Stein Variational Gradient Descent (SVGD) by Liu and Wang [27] is a popular algorithm
23 for this problem. It uses a positive definite kernel k to evolve n interacting particles $(\mathbf{x}_t^{(i)})_{i \in [n], t \in \mathbb{N}}$:

$$\mathbf{x}_{t+1}^{(i)} \leftarrow \mathbf{x}_t^{(i)} - \frac{\gamma}{n} \sum_{j=1}^n \left[k(\mathbf{x}_t^{(i)}, \mathbf{x}_t^{(j)}) \nabla F(\mathbf{x}_t^{(j)}) - \nabla_2 k(\mathbf{x}_t^{(i)}, \mathbf{x}_t^{(j)}) \right] \quad (1)$$

24 SVGD exhibits remarkable empirical performance in various Bayesian inference, generative mod-
25 eling and reinforcement learning tasks [27, 43, 21, 29] and usually converges rapidly to the target
26 density while using only a few particles, often outperforming Markov Chain Monte Carlo methods.
27 However, in contrast to its wide practical applicability, theoretical analysis of its behavior is relatively
28 unexplored. Prior works on the analysis of SVGD [23, 14, 26, 36, 7] mainly consider the mean-field
29 limit (or population limit), where the number of particles $n \rightarrow \infty$. These works assume that the
30 initial distribution of the (infinite number of) particles has a finite KL divergence to the target π^*
31 and subsequently, interpret mean-field SVGD dynamics as ‘Projected’ Gradient Descent (GD) of KL
32 divergence on the space of probability measures, equipped with the Wasserstein geometry. Under
33 suitable assumptions on the target density, these works use the theory of Wasserstein Gradient Flows

Result	Algorithm	Assumption	Rate	Oracle Complexity
Korba et al. [23]	Population Limit SVGD	Uniformly Bounded $\text{KSD}_{\pi^*}(\bar{\mu}_t \pi^*)$	$\frac{\text{poly}(d)}{\sqrt{T}}$	Not Implementable
Salim et al. [36]	Population Limit SVGD	Sub-gaussian π^*	$\frac{d^{3/2}}{\sqrt{T}}$	Not Implementable
Shi and Mackey [37]	SVGD	Sub-gaussian π^*	$\frac{\text{poly}(d)}{\sqrt{\log \log n^{\Theta(1/d)}}}$	$\frac{\text{poly}(d)}{\epsilon^2} e^{\Theta(d e^{\text{poly}(d)/\epsilon^2})}$
Ours, Corollary 1	VP-SVGD	Sub-gaussian π^*	$(d/n)^{1/4} + (d/n)^{1/2}$	d^4/ϵ^{12}
Ours, Corollary 1	GB-SVGD	Sub-gaussian π^*	$d^{1/3}/n^{1/12} + (d/n)^{1/2}$	d^6/ϵ^{18}
Ours, Corollary 1	VP-SVGD	Sub-exponential π^*	$\frac{d^{1/3}}{n^{1/6}} + \frac{d}{n^{1/2}}$	d^6/ϵ^{16}
Ours, Corollary 1	GB-SVGD	Sub-exponential π^*	$\frac{d^{3/8}}{n^{1/16}} + \frac{d}{n^{1/2}}$	d^9/ϵ^{24}

Table 1: Comparison of our results with prior works. d , T , and n denote the dimension, no. of iterations and no. of output particles respectively. Oracle Complexity denotes number of evaluations of ∇F needed to achieve $\text{KSD}_{\pi^*}(\cdot || \pi^*) \leq \epsilon$ (with n and T appropriately optimized), and Rate denotes convergence rate w.r.t KSD metric. Note that: 1. Population Limit SVGD is not implementable as it requires infinite particles 2. The uniformly bounded $\text{KSD}_{\pi^*}(\bar{\mu}_t || \pi^*)$ assumption is much stronger than subgaussianity and cannot be verified a priori (see Salim et al. [36] Section 1.2.1)

[1] to establish non-asymptotic (in time) convergence of mean-field SVGD to π^* in the Kernel Stein Discrepancy (KSD) metric. While this framework suffices to explain the behavior of SVGD in the mean-field limit, the same techniques are insufficient for analyzing finite-particle regime. This is mainly due to the fact that the empirical measure $\hat{\mu}^{(n)}$ of a finite number of particles does not admit a density (w.r.t Lebesgue Measure), and thus, its KL divergence to the target is always infinite. Moreover, a direct analysis of the dynamics of finite-particle SVGD becomes prohibitively difficult due to complex inter-particle dependencies. To the best of our knowledge, Shi and Mackey [37] is the only result that obtains an explicit convergence rate for finite-particle SVGD by tracking the deviation between the law of n -particle SVGD and mean-field SVGD. The authors show that for subgaussian π^* , the empirical measure of n -particle SVGD converges to π^* at $O(\sqrt{\frac{\text{poly}(d)}{\log \log n^{\Theta(1/d)}}})$ rate in KSD (we explicate the d dependence in Shi and Mackey [37] by closely following their analysis). The obtained rate (which suffers from curse of dimensionality) is quite slow and fails to adequately explain the practical performance of SVGD.

Our work deliberately avoids computing the deviation between mean-field SVGD and finite-particle SVGD. Instead, we directly analyze the dynamics of KL divergence along a carefully constructed trajectory in the space of distributions. Our proposed algorithm, Virtual Particle SVGD (VP-SVGD) devises an *unbiased stochastic approximation (in the space of measures) to mean-field SVGD*. We achieve this by considering additional particles called *virtual particles* which evolve in time but aren't part of the output (i.e. *real particles*). These virtual particles are used only to compute information about the current population-level distribution of the real particles, and enable exact implementation of our stochastic approximation to mean-field SVGD, while using only a finite number of particles. Our analysis is similar in spirit to non-asymptotic analyses of Stochastic Gradient Descent (SGD) that do not attempt to track GD (analogous to mean-field SVGD in this case), but instead track the evolution of the objective function along the SGD trajectory using appropriate descent lemmas [20, 11]. The key feature of our proposed stochastic approximation is the fact that it can be exactly implemented using only a finite number of particles. This allows us to design faster variants of SVGD with provably fast finite-particle convergence.

1.1 Contributions and Technical Challenges

1.2 Contributions

VP-SVGD and GB-SVGD We propose two variants of SVGD that enjoy provably fast finite-particle convergence guarantees: Virtual Particle SVGD (VP-SVGD, Algorithm 1) and Global Batch SVGD (GB-SVGD, Algorithm 2). VP-SVGD is a conceptually elegant stochastic approximation (in the space of probability measures) of mean-field SVGD, and GB-SVGD is a practically efficient version of SVGD which achieves good empirical performance. Our analysis of GB-SVGD builds upon that of VP-SVGD. When the potential F is smooth and satisfies a quadratic growth condition (which

69 holds under subgaussianity of π^* , a common assumption in prior works [36, 37]), we show that
 70 the n particles output by T steps of our algorithms, run with batch-size K , are at least as good as
 71 i.i.d draws from a distribution whose KSD to π^* is at most $O(d^{1/3}/(KT)^{1/6})$. Our results also hold
 72 under a mild subquadratic growth condition for F , which is much weaker than isoperimetric (e.g.
 73 Poincare Inequality) or information-transport (e.g. Talagrand’s Inequality T_1) assumptions generally
 74 considered in the sampling literature [41, 36, 37, 8, 2].

75 **State-of-the-art Finite Particle Guarantees** As corollaries of the above result, we establish that
 76 *VP-SVGD and GB-SVGD exhibit the best known finite-particle guarantees in the literature which*
 77 *significantly outperform that of prior works.* Our results are summarized in Table 1. In particular,
 78 under subgaussianity of π^* , we show that the empirical measure of the n particles output by VP-
 79 SVGD converges to π^* in KSD at a $O((d/n)^{1/4} + (d/n)^{1/2})$ rate. Similarly, the empirical measure of the
 80 n output particles of GB-SVGD converges to π^* at a KSD rate of $O(d^{1/3}/n^{1/12} + (d/n)^{1/2})$. Both these
 81 results are a *double exponential improvement* over the $O(\frac{\text{poly}(d)}{\sqrt{\log \log n^{\Theta(1/d)}}})$ KSD rate of n -particle
 82 SVGD obtained by Shi and Mackey [37], which, to our knowledge, is the best known finite-particle
 83 rate for SVGD so far. In terms of gradient oracle complexity (i.e., the number of ∇F evaluations
 84 required to achieve $\text{KSD}_{\pi^*}(\cdot|\pi^*) \leq \epsilon$), we show that for subgaussian π^* , the oracle complexity
 85 of VP-SVGD is $O(d^4/\epsilon^{12})$ while that of GB-SVGD is $O(d^6/\epsilon^{18})$. To the best of our knowledge, our
 86 result presents the *first known oracle complexity guarantee with polynomial dimension dependence*,
 87 and consequently, does not suffer from a curse of dimensionality unlike prior works. Furthermore,
 88 as discussed above, the conditions under which our result holds is far weaker than subgaussianity
 89 of π^* , and as such, includes sub-exponential targets and beyond. In particular, *our guarantees for*
 90 *sub-exponential target distributions are (to the best of our knowledge) the first of its kind.*

91 **Empirical Evaluation** Our experiments in Appendix 8 show that GB-SVGD obtains similar perfor-
 92 mance as SVGD *but requires fewer computations.*

93 Our analysis resolves the following important technical challenges of independent interest:

94 **Stochastic Approximation in the Space of Probability Measures** Stochastic approximations are
 95 widely used in optimization and and sampling [24, 44]. In sampling, such approximations are
 96 generally implemented in path space, e.g., Stochastic Gradient Langevin Dynamics [44] takes
 97 a stochastic approximation of the form $\mathbf{x}_{t+1} = \mathbf{x}_t - \frac{\eta}{K} \sum_{j=0}^{K-1} \nabla f(\mathbf{x}_t, \xi_j) + \sqrt{2\eta}\epsilon_t$, $\epsilon_t \sim$
 98 $\mathcal{N}(0, \mathbf{I})$; $\mathbb{E}[f(\mathbf{x}_t, \xi_j)|\mathbf{x}_t] = F(\mathbf{x}_t)$. Such stochastic approximations are analyzed using the theo-
 99 ry of stochastic processes over \mathbb{R}^d [12, 34, 22]. However, when viewed in the space of probability
 100 measures (i.e, $\mu_t = \text{Law}(\mathbf{x}_t)$), the time-evolution of these algorithms is deterministic. In contrast, our
 101 approach designs *stochastic approximations in the space of probability measures*. In particular, the
 102 time-evolution of the law of any particle in VP-SVGD and GB-SVGD are a stochastic approximation
 103 of the dynamics of mean-field SVGD. Careful design ensures that our stochastic approximation
 104 requires only a finite number of particles for exact implementation.

105 **Tracking KL Divergence in the Finite-Particle Regime** The population limit ($n \rightarrow \infty$) ensures that
 106 the initial empirical distribution (μ_0) of SVGD admits a density (w.r.t the Lebesgue measure). Prior
 107 works on population-limit SVGD analyze the time-evolution of the KL divergence to π^* . However,
 108 this approach cannot be directly used for finite-particle SVGD since the empirical distribution of a
 109 finite number of particles does not admit a density, and thus its KL divergence to π^* is infinite. Our
 110 analysis of VP-SVGD and GB-SVGD circumvents this obstacle by considering the dynamics of an
 111 infinite number of particles, whose empirical measure then admits a density. However, the careful
 112 design ensures that the dynamics of n of these particles can be computed exactly, using only a finite
 113 total number of (real + virtual) particles. When conditioned on the virtual particles, these particles
 114 are i.i.d. and their conditional law is close to the target distribution with high probability.

115 2 Notation and Problem Setup

116 We use $\|\cdot\|$, $\langle \cdot, \cdot \rangle$ to denote the Euclidean norm and inner product over \mathbb{R}^d respectively, while other
 117 norms and inner products are subscripted with their underlying space. $\mathcal{B}(R)$ denotes the ball of radius
 118 R in $(\mathbb{R}^d, \langle \cdot, \cdot \rangle)$. $\mathcal{P}_2(\mathbb{R}^d)$ denotes the space of probability measures on \mathbb{R}^d with finite second moment,
 119 with the Wasserstein-2 metric denoted as $\mathcal{W}_2(\mu, \nu)$ for $\mu, \nu \in \mathcal{P}_2(\mathbb{R}^d)$. For any two probability
 120 measures μ, ν , we denote their KL divergence as $\text{KL}(\mu||\nu)$. For any function $f : X \rightarrow Y$ and

121 any probability measure μ over X , we let $f_{\#}\mu$ denote the law of $f(\mathbf{x}) : \mathbf{x} \sim \mu$. Given a sigma
122 algebra \mathcal{F} over some space Ω , and a measurable space \mathcal{X} , $\mu(\cdot; \cdot) : \mathcal{F} \times \mathcal{X} \rightarrow \mathbb{R}^+$ is a probability
123 kernel if $\forall x \in \mathcal{X}$, $\mu(\cdot; x)$ is a measure over \mathcal{F} and $\forall A \in \mathcal{F}$, the map $x \rightarrow \mu(A; x)$ is measurable.
124 We use probability measures $\mu(\cdot; \mathbf{x})$, where \mathbf{x} is a random element of some appropriate space \mathcal{X} ,
125 resulting in random probability measures. We use $[m]$ and (m) to denote the sets $\{1, \dots, m\}$ and
126 $\{0, \dots, m-1\}$ respectively, and $S_{(m)}$ to denote the set of all permutations of (m) . We use the O
127 notation to characterize the dependence of our rates on the number of iterations T , dimension d and
128 batch-size K , suppressing numerical and problem-dependent constants. We use \lesssim to denote \leq upto
129 universal constants. We fix a symmetric positive definite kernel $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ and denote the
130 corresponding reproducing kernel Hilbert space (RKHS) [38] as \mathcal{H}_0 . We denote the product RKHS
131 as $\mathcal{H} = \prod_{i=1}^d \mathcal{H}_0$, equipped with the standard inner product for product spaces. We assume k is
132 differentiable in both its arguments and let $\nabla_2 k(\mathbf{x}, \mathbf{y})$ denote its gradient w.r.t the second argument.
133 For any $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, we assume $\mathcal{H} \subset L^2(\mu)$ and the inclusion map $i_\mu : \mathcal{H} \rightarrow L^2(\mu)$ is continuous.
134 We use $P_\mu : L^2(\mu) \rightarrow \mathcal{H}$ to denote the adjoint of i_μ , i.e., the unique operator which satisfies
135 $\langle f, i_\mu g \rangle_{L^2(\mu)} = \langle P_\mu f, g \rangle_{\mathcal{H}}$ for any $f \in L^2(\mu), g \in \mathcal{H}$. Carmeli et al. [6] shows that P_μ can be
136 expressed as a kernel convolution, i.e., $(P_\mu f)(\mathbf{x}) = \int k(\mathbf{x}, \mathbf{y}) f(\mathbf{y}) d\mu(\mathbf{y})$. We define the function $h :$
137 $\mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ as $h(\mathbf{x}, \mathbf{y}) = k(\mathbf{x}, \mathbf{y}) \nabla F(\mathbf{y}) - \nabla_2 k(\mathbf{x}, \mathbf{y})$ and $h_\mu \in \mathcal{H}$ as $h_\mu = P_\mu(\nabla_{\mathbf{x}} \log(\frac{d\mu}{d\pi^*}(\mathbf{x})))$
138 for any $\mu \in \mathcal{P}_2(\mathbb{R}^d)$. Integration by parts shows that $h_\mu(\mathbf{x}) = \int h(\mathbf{x}, \mathbf{y}) d\mu(\mathbf{y})$. Similar to prior
139 works [36, 23, 37] we use Kernel Stein Discrepancy (KSD) as a convergence metric.

140 **Definition 1** (Kernel Stein Discrepancy[28, 9]). *Define the Langevin Stein Operator of π^* acting*
141 *on any differentiable $g : \mathbb{R}^d \rightarrow \mathbb{R}^d$ as $(T_{\pi^*} g)(\mathbf{x}) = \nabla \cdot g(\mathbf{x}) - \langle \nabla F(\mathbf{x}), g(\mathbf{x}) \rangle$. Then, for any*
142 *two probability measures μ, ν , the Kernel Stein Discrepancy between μ and ν w.r.t π^* is defined as*
143 $\text{KSD}_{\pi^*}(\mu|\nu) = \sup_{\|g\|_{\mathcal{H}} \leq 1} \mathbb{E}_\mu[T_{\pi^*} g] - \mathbb{E}_\nu[T_{\pi^*} g] = \|h_\mu - h_\nu\|_{\mathcal{H}}$.

144 3 Background on Mean-Field SVGD

145 We briefly introduce the analysis of mean-field SVGD using the theory of Wasserstein Gradient Flows
146 and refer the readers to prior work [23, 36] for a detailed treatment. The metric space $(\mathcal{P}_2(\mathbb{R}^d), \mathcal{W}_2)$ is
147 called the Wasserstein space, which admits the following Riemannian structure : For any $\mu \in \mathcal{P}_2(\mathbb{R}^d)$,
148 the tangent space $T_\mu \mathcal{P}_2(\mathbb{R}^d)$ can be identified with the Hilbert space $L^2(\mu)$. We can then define
149 differentiable functionals $\mathcal{L} : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$ and compute their Wasserstein gradients $\nabla_{\mathcal{W}_2} \mathcal{L}$. Note
150 that the target π^* is the unique minimizer over of the functional $\mathcal{L}[\mu] = \text{KL}(\mu|\pi^*)$ over $\mathcal{P}_2(\mathbb{R}^d)$, and
151 its Wasserstein Gradient is $\nabla_{\mathcal{W}_2} \mathcal{L}[\mu] = \nabla_{\mathbf{x}} \log(\frac{d\mu}{d\pi^*}(\mathbf{x}))$ [1]. This powerful machinery has served as
152 a backbone for the analysis of algorithms such as LMC [45, 3] and mean-field SVGD [14, 23, 36].
153 In particular, mean-field SVGD can be viewed as ‘Projected’ Gradient Descent in $\mathcal{P}_2(\mathbb{R}^d)$. To infer
154 this, let $\hat{\mu}_t^n$ denote the empirical measures of the SVGD particles $(\mathbf{x}_t^{(i)})_{i \in [n]}$ at step t and recall that
155 $h_\mu(\mathbf{x}) = P_\mu(\nabla_{\mathbf{x}} \log(\frac{d\mu}{d\pi^*}))(\mathbf{x}) = \int h(\mathbf{x}, \mathbf{y}) d\mu(\mathbf{y})$ (Sec. 2). The SVGD updates in (1) can be recast
156 as $\hat{\mu}_{t+1}^n = (I - \gamma h_{\hat{\mu}_t^n})_{\#} \hat{\mu}_t^n$. In the limit of infinite particles $n \rightarrow \infty$, suppose the empirical measure
157 $\hat{\mu}_t^n$ converges to the population measure $\bar{\mu}_t$. In this mean-field limit, the updates can be expressed as,

$$\bar{\mu}_{t+1} = (I - h_{\bar{\mu}_t})_{\#} \bar{\mu}_t = (I - \gamma P_{\bar{\mu}_t}(\nabla \log(\frac{d\bar{\mu}_t}{d\pi^*})))_{\#} \bar{\mu}_t = (I - \gamma P_{\bar{\mu}_t}(\nabla_{\mathcal{W}_2} \text{KL}(\bar{\mu}_t|\pi^*)))_{\#} \bar{\mu}_t$$

158 Recall from Sec. 2 that $P_{\bar{\mu}_t} : L^2(\bar{\mu}_t) \rightarrow \mathcal{H}$ is the adjoint of $i_{\bar{\mu}_t}$. Since $\mathcal{H} \subset L^2(\bar{\mu}_t)$,
159 the updates of mean-field SVGD can be seen as ‘Projected’ Wasserstein Gradient Descent for
160 $\mathcal{L}[\mu] = \text{KL}(\mu|\pi^*)$, with the Wasserstein Gradient at each step being projected onto the RKHS \mathcal{H} .
161 Assuming $\text{KL}(\bar{\mu}_0|\pi^*) < \infty$, convergence of population limit SVGD is then established by tracking
162 the evolution of $\text{KL}(\bar{\mu}_t|\pi^*)$ under appropriate structural assumptions (such as subgaussianity) on π^* .

163 4 Algorithm and Intuition

164 In this section, we derive VP-SVGD (Algorithm 1), and build upon it to obtain GB-SVGD. Consider
165 a countably infinite collection of particles $\mathbf{x}_0^{(l)} \in \mathbb{R}^d$, $l \in \mathbb{N} \cup \{0\}$, sampled i.i.d from a measure
166 μ_0 , having a density w.r.t. the Lebesgue measure. By the strong law of large numbers, the empirical
167 measure of $\mathbf{x}_0^{(l)}$ is almost surely equal to μ_0 [13, Theorem 11.4.1]. Let $K \in \mathbb{N}$ denote the batch size
168 and define the filtration $\mathcal{F}_t = \sigma(\{\mathbf{x}_0^{(l)} \mid l \leq Kt - 1\})$, $\forall t \in \mathbb{N}$ with \mathcal{F}_0 being the trivial σ algebra.

Algorithm 1 Virtual Particle SVGD (VP-SVGD)

Input: Number of steps T , number of output particles n , batch size K , Initial positions $\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(n+KT-1)}$ *i.i.d.* μ_0 , Kernel k , step size γ .

- 1: **for** $t \in \{0, \dots, T-1\}$ **do**
- 2: **for** $s \in \{0, \dots, KT+n-1\}$ **do**
- 3: $\mathbf{x}_{t+1}^{(s)} = \mathbf{x}_t^{(s)} - \frac{\gamma}{K} \sum_{l=0}^{K-1} [k(\mathbf{x}_t^{(s)}, \mathbf{x}_t^{(tK+l)}) \nabla F(\mathbf{x}_t^{(tK+l)}) - \nabla_2 k(\mathbf{x}_t^{(s)}, \mathbf{x}_t^{(tK+l)})]$
- 4: **end for**
- 5: **end for**
- 6: Draw S uniformly at random from $\{0, \dots, T-1\}$
- 7: Output $(\mathbf{y}^{(0)}, \dots, \mathbf{y}^{(n-1)}) = (\mathbf{x}_S^{(TK)}, \dots, \mathbf{x}_S^{(TK+n-1)})$

169 For ease of exposition, we discuss the case of $K = 1$ below and present a complete derivation for
170 arbitrary $K \geq 1$ in Section C. Recall from Section 3 that the updates of mean-field SVGD in $\mathcal{P}_2(\mathbb{R}^d)$
171 is as follows:

$$\bar{\mu}_{t+1} = (I - \gamma h_{\bar{\mu}_t})_{\#} \bar{\mu}_t \quad (2)$$

172 We aim to design a stochastic approximation in $\mathcal{P}_2(\mathbb{R}^d)$ for the updates (2), such that it admits a
173 finite-particle realization. To this end, we propose the following dynamics in \mathbb{R}^d

$$\mathbf{x}_{t+1}^{(s)} = \mathbf{x}_t^{(s)} - \gamma h(\mathbf{x}_t^{(s)}, \mathbf{x}_t^{(t)}), \quad s \in \mathbb{N} \cup \{0\} \quad (3)$$

174 Now, for each time-step t , we focus on the time evolution of the particles $(\mathbf{x}_t^{(l)})_{l \geq t}$ (called the *lower*
175 *triangular evolution*). From (3), we observe that for any $t \in \mathbb{N}$ and $l \geq t$, $\mathbf{x}_t^{(l)}$ depends only on
176 $\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(t-1)}, \mathbf{x}_0^{(l)}$. Therefore there exists a deterministic, measurable function H_t such that:

$$\mathbf{x}_t^{(l)} = H_t(\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(t-1)}, \mathbf{x}_0^{(l)}); \quad \text{for every } l \geq t \quad (4)$$

177 Since $\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(t-1)}, \mathbf{x}_0^{(l)}$ *i.i.d.* μ_0 , we conclude from (4) that $(\mathbf{x}_t^{(l)})_{l \geq t}$ are i.i.d when conditioned
178 on $\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(t-1)}$. To this end, we define the random measure $\mu_t | \mathcal{F}_t$ as the law of $\mathbf{x}_t^{(t)}$ conditioned
179 on \mathcal{F}_t , i.e., $\mu_t | \mathcal{F}_t$ is a probability kernel $\mu_t(\cdot; \mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(t-1)})$, where $\mu_0 | \mathcal{F}_0 := \mu_0$. By the strong
180 law of large numbers, $\mu_t | \mathcal{F}_t$ is equal to the empirical measure of $(\mathbf{x}_t^{(l)})_{l \geq t}$ conditioned on \mathcal{F}_t . We
181 will use $\mu_t | \mathcal{F}_t$ and $\mu_t(\cdot; \mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(t-1)})$ interchangeably.

182 Define the random function $g_t : \mathbb{R}^d \rightarrow \mathbb{R}^d$ as $g_t(\mathbf{x}) := h(\mathbf{x}, \mathbf{x}_t^{(t)})$. From (4), we note that g_t is \mathcal{F}_{t+1}
183 measurable. From (3), we infer that the particles satisfy the following relation:

$$\mathbf{x}_{t+1}^{(s)} = (I - \gamma g_t)(\mathbf{x}_t^{(s)}), \quad s \geq t+1$$

184 Recall that $\mathbf{x}_{t+1}^{(s)} | \mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(t)} \sim \mu_{t+1} | \mathcal{F}_{t+1}$ for any $s \geq t+1$. Furthermore, from Equation (4), we
185 note that for $s \geq t+1$, $\mathbf{x}_t^{(s)}$ depends only on $\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(t-1)}$ and $\mathbf{x}_0^{(s)}$. Hence, we conclude that
186 $\text{Law}(\mathbf{x}_t^{(s)} | \mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(t)}) = \text{Law}(\mathbf{x}_t^{(s)} | \mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(t-1)}) = \mu_t | \mathcal{F}_t \forall s \geq t+1$. With this insight, the
187 dynamics of the lower-triangular evolution in $\mathcal{P}_2(\mathbb{R}^d)$ that the following holds almost surely:

$$\mu_{t+1} | \mathcal{F}_{t+1} = (I - \gamma g_t)_{\#} \mu_t | \mathcal{F}_t \quad (5)$$

188 $\mathbf{x}_t^{(t)} | \mathcal{F}_t \sim \mu_t | \mathcal{F}_t$ implies $\mathbb{E}[g_t(\mathbf{x}) | \mathcal{F}_t] = h_{\mu_t | \mathcal{F}_t}(\mathbf{x})$. Thus *lower triangular dynamics* (5) is a stochas-
189 tic approximation in $\mathcal{P}_2(\mathbb{R}^d)$ to the population limit of SVGD (2). Setting the batch size to general
190 K and tracking the evolution of the first $KT+n$ particles, we obtain VP-SVGD (Algorithm 1).

191 **Virtual Particles** In Algorithm 1, $(\mathbf{x}_t^{(l)})_{KT \leq l \leq KT+n-1}$ are the *real particles* which constitute the
192 output. $(\mathbf{x}_t^{(l)})_{l < KT}$ are *virtual particles* which propagate information about the probability measure
193 $\mu_t | \mathcal{F}_t$ to enable computation of g_t , an unbiased estimate of the projected Wasserstein gradient $h_{\mu_t | \mathcal{F}_t}$.

194 **Intuition Behind GB-SVGD** We note that VP-SVGD (Algorithm 1) is a without-replacement
195 random-batch approximation of SVGD (1), where a different batch is used across timesteps, but
196 the same batch is used across particles given a fixed timestep. With i.i.d. initialization, picking the
197 ‘virtual particles’ in a fixed order or from a random permutation does not change the evolution of the
198 real particles. With this insight, we design GB-SVGD (Algorithm 2) where we consider n particles
199 and output n particles (instead of wasting KT particles as ‘virtual particles’) via a random-batch

200 approximation of SVGD. In GB-SVGD, with replacement sampling means selecting a batch of K
 201 particles i.i.d. from $\text{Uniform}(\mathcal{S}(n))$. Without replacement sampling means fixing a random permutation
 202 $\sigma \sim \text{Uniform}(\mathcal{S}(n))$ and selecting the batches in the order specified by σ (essentially ensuring that no
 data point is repeated during an iteration).

Algorithm 2 Global Batch SVGD (GB-SVGD)

Input: # of time steps T , # of particles n , $\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(n-1)}$ i.i.d. μ_0 , Kernel k , step size γ , Batch size K ,
 Sampling method $\in \{\text{with replacement, without replacement}\}$

```

1: for  $t \in \{0, \dots, T-1\}$  do
2:    $\mathcal{K}_t \leftarrow$  random subset of  $[n]$  of size  $K$  (via. sampling method)
3:   for  $s \in \{0, \dots, n-1\}$  do
4:      $\mathbf{x}_{t+1}^{(s)} = \mathbf{x}_t^{(s)} - \frac{\gamma}{K} \sum_{r \in \mathcal{K}_t} [k(\mathbf{x}_t^{(s)}, \mathbf{x}_t^{(r)}) \nabla F(\mathbf{x}_t^{(r)}) - \nabla_2 k(\mathbf{x}_t^{(s)}, \mathbf{x}_t^{(r)})]$ 
5:   end for
6: end for
7: Draw  $S$  uniformly at random from  $\{0, 1, \dots, T-1\}$ 
8: Output  $(\bar{\mathbf{y}}^{(0)}, \dots, \bar{\mathbf{y}}^{(n-1)}) = (\mathbf{x}_S^{(0)}, \dots, \mathbf{x}_S^{(n-1)})$ 

```

203

204 5 Assumptions

205 We now discuss the key assumptions required for our analysis of VP-SVGD and GB-SVGD.

206 **Assumption 1** (L-Smoothness). ∇F exists and is L Lipschitz. Moreover $\|\nabla F(0)\| \leq \sqrt{L}$.

207 Lipschitzness of ∇F is standard in optimization and sampling. It is also easy find a point \mathbf{x}^* such
 208 that $\|\nabla F(\mathbf{x}^*)\| \leq \sqrt{L}$ (e.g., using $\Theta(1)$ steps of GD [32]) and center the initialization at \mathbf{x}^* . We
 209 take $\mathbf{x}^* = 0$ without loss of generality. We now impose the following growth condition on F .

210 **Assumption 2** (Growth Condition). There exist $\alpha, d_1, d_2 > 0$ such that $F(\mathbf{x}) \geq d_1 \|\mathbf{x}\|^\alpha - d_2$

211 Note that Assumption 1 ensures $\alpha \leq 2$. Assumption 2 is a tail decay assumption on $\pi^*(\mathbf{x}) \propto e^{-F(\mathbf{x})}$,
 212 ensuring that its tails decay as $\propto e^{-\|\mathbf{x}\|^\alpha}$. Thus, it holds with $\alpha = 2$ when π^* is subgaussian and with
 213 $\alpha = 1$ when π^* is subexponential (See Appendix B for proofs). Subgaussianity is equivalent to π^*
 214 satisfying the \mathbb{T}_1 inequality [5], commonly assumed in prior works on SVGD [36, 37]. Moreover,
 215 subexponentiality holds whenever π^* satisfies the Poincare Inequality [4], which is a mild condition in
 216 the sampling literature [41, 8, 2, 12, 7]. This makes Assumption 1 much weaker than the isoperimetric
 217 or information-transport assumptions considered in prior works. We also make the following mild
 218 assumptions on the k that appear in prior work [23, 17] and are satisfied by several standard kernels
 219 (e.g. RBF Kernels, Matérn kernels of order $\geq 3/2$)

220 **Assumption 3** (Kernel Regularity). For any $\mathbf{y} \in \mathbb{R}^d$, $k(\cdot, \mathbf{y})$ satisfies $\|k(\cdot, \mathbf{y})\|_{\mathcal{H}_0} \leq B$ and
 221 $\nabla_2 k(\cdot, \mathbf{y}) \in \mathcal{H}$ with $\|\nabla_2 k(\cdot, \mathbf{y})\|_{\mathcal{H}} \leq B$. Moreover, there exist $A_1, A_2, A_3 > 0$ such that
 222 $0 \leq k(\mathbf{x}, \mathbf{y}) \leq \frac{A_1}{1+\|\mathbf{x}-\mathbf{y}\|^2}$, $\|\nabla_2 k(\mathbf{x}, \mathbf{y})\| \leq A_2$, and $\|\nabla_2 k(\mathbf{x}, \mathbf{y})\|^2 \leq A_3 k(\mathbf{x}, \mathbf{y})$.

223 For ease of exposition, we make the following mild assumption on the initialization.

224 **Assumption 4** (Initialization). The initial density is $\mu_0 = \text{Uniform}(\mathcal{B}(R))$ with $\text{KL}(\mu_0 \|\pi^*) < \infty$.

225 Since $\mathcal{N}(0, \mathbf{I})$ and $\text{Uniform}(\mathcal{B}(R))$ are nearly indistinguishable with high probability when $R =$
 226 $\tilde{\Theta}(\sqrt{d})$, Assumption 4 can be easily replaced by the Gaussian initialization assumed in prior works.
 227 Furthermore, we show in Appendix B that $R = \sqrt{d/L}$ ensures $\text{KL}(\mu_0 \|\pi^*) = O(d)$

228 6 Results

229 6.1 VP-SVGD

230 Our first result, proved in Appendix C, shows that the law of the *real particles* of VP-SVGD, when
 231 conditioned on the virtual particles, is close to π^* in KSD. Consequently, it shows that the particles
 232 output by VP-SVGD are i.i.d. samples from a random probability measure $\bar{\mu}(\cdot; \mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(KT-1)}, S)$
 233 which is close to π^* in KSD. Appendix C also presents a high-probability version of Theorem 1.

234 **Theorem 1 (Convergence of VP-SVGD).** Let μ_t be as defined in Section 4. Let Assumptions 1, 2, 3,
 235 and 4 be satisfied and let $\gamma \leq \min\{1/2A_1L, 1/(4+L)B\}$. There exist $(\zeta_i)_{0 \leq i \leq 3}$ depending polynomially
 236 on $A_1, A_2, A_3, B, L, d_1, d_2$ for any fixed $\alpha \in (0, 2]$, such that whenever $\gamma\xi \leq \frac{1}{2B}$, with $\xi =$
 237 $\zeta_0 + \zeta_1(\gamma T)^{1/\alpha} + \zeta_2(\gamma^2 T)^{1/\alpha} + \zeta_3 R^{2/\alpha}$, the following holds:

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} [\text{KSD}_{\pi^*}^2(\mu_t | \mathcal{F}_t | \pi^*)] \leq \frac{2\text{KL}(\mu_0 | \pi^*)}{\gamma T} + \frac{\gamma B(4+L)\xi^2}{K}$$

238 Define the probability kernel $\bar{\mu}(\cdot; \cdot)$ as follows: For any $x_\tau \in \mathbb{R}^d$, $\tau \in (KT)$ and $s \in (T)$,
 239 $\bar{\mu}(\cdot; x_0, \dots, x_{KT-1}, s) := \mu_s(\cdot; x_0, \dots, x_{Ks-1})$ and $\bar{\mu}(\cdot; x_0, \dots, x_{KT-1}, s=0) := \mu_0(\cdot)$. Con-
 240 ditioned on $\mathbf{x}_\tau^{(0)} = x_\tau$, $S = s$ for every $\tau \in (KT)$, the outputs $\mathbf{y}^{(0)}, \dots, \mathbf{y}^{(n-1)}$ of VP-SVGD are
 241 i.i.d samples from $\bar{\mu}(\cdot; x_0, \dots, x_{KT-1}, s)$. Furthermore,

$$\mathbb{E}[\text{KSD}_{\pi^*}^2(\bar{\mu}(\cdot; \mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(KT-1)}, S) | \pi^*)] \leq \frac{2\text{KL}(\mu_0 | \pi^*)}{\gamma T} + \frac{\gamma B(4+L)\xi^2}{K}$$

242 **Convergence Rates** Setting $R = \sqrt{d/L}$ ensures $\text{KL}(\mu_0 | \pi^*) = O(d)$ (see Appendix B). Hence,
 243 choosing $\gamma = O(\frac{(Kd)^\eta}{T^{1-\eta}})$ ensures that $\mathbb{E}[\text{KSD}_{\pi^*}^2(\bar{\mu} | \pi^*)] = O(\frac{d^{1-\eta}}{(KT)^\eta})$ where $\eta = \frac{\alpha}{2(1+\alpha)}$. Thus, for
 244 $\alpha = 2$, (i.e, sub-Gaussian π^*), $\text{KSD}^2 = O(\frac{d^{2/3}}{(KT)^{1/3}})$. For $\alpha = 1$ (i.e, sub-Exponential π^*), the rate
 245 (in squared KSD) becomes $O(\frac{d^{3/4}}{(KT)^{1/4}})$. To the best of our knowledge, our convergence guarantee
 246 for sub-exponential π^* is the first of its kind.

247 **Comparison with Prior Works** Salim et al. [36] analyzes population-limit SVGD for subgaussian
 248 π^* , obtaining $\text{KSD}^2 = O(d^{3/2}/T)$ rate. We note that population-limit SVGD is not implementable
 249 whereas VP-SVGD is an implementable algorithm whose outputs are samples from a distribution
 250 with guaranteed convergence to π^* .

251 6.2 GB-SVGD

252 We now use VP-SVGD as the basis to analyze GB-SVGD. Assume $n > KT$. Then, with probability
 253 $\geq 1 - K^2 T^2/n$ (for with-replacement sampling) and 1 (for without-replacement sampling), the random
 254 batches \mathcal{K}_t in GB-SVGD (Algorithm 2) are disjoint and contain distinct elements. Conditioned on
 255 this event \mathcal{E} , we note that the $n - KT$ particles that were not included in any random batch \mathcal{K}_t evolve
 256 exactly like the n real particles of VP-SVGD. With this insight, we show that, conditioned on \mathcal{E} , the
 257 outputs of VP-SVGD and GB-SVGD can be coupled such that the first $n - KT$ particles output by
 258 both the algorithms are exactly equal. This can be used to derive the following squared KSD bound
 259 between their empirical measures. We prove this result in Appendix D

260 **Theorem 2 (KSD Bounds for GB-SVGD).** Let $n > KT$ and let $\mathbf{Y} = (\mathbf{y}^{(0)}, \dots, \mathbf{y}^{(n-1)})$ and
 261 $\bar{\mathbf{Y}} = (\bar{\mathbf{y}}^{(0)}, \dots, \bar{\mathbf{y}}^{(n-1)})$ denote the outputs of VP-SVGD and GB-SVGD respectively. Moreover, let
 262 $\hat{\mu}^{(n)} = \frac{1}{n} \sum_{i=0}^{n-1} \delta_{\mathbf{y}^{(i)}}$ and $\hat{\nu}^{(n)} = \frac{1}{n} \sum_{i=0}^{n-1} \delta_{\bar{\mathbf{y}}^{(i)}}$ denote their respective empirical measures. Under
 263 the assumptions and parameter settings of Theorem 1, there exists a coupling of \mathbf{Y} and $\bar{\mathbf{Y}}$ such that:

$$\mathbb{E}[\text{KSD}_{\pi^*}^2(\hat{\nu}^{(n)} | \hat{\mu}^{(n)})] \leq \begin{cases} \frac{2K^2 T^2 \xi^2}{n^2} & (\text{without replacement sampling}) \\ \frac{2K^2 T^2 \xi^2}{n^2} \left(1 - \frac{K^2 T^2}{n}\right) + \frac{2K^2 T^2 \xi^2}{n} & (\text{with replacement sampling}) \end{cases} \quad (6)$$

264 6.3 Convergence of the Empirical Measure to the Target

265 As a corollary of Theorem 1 and Theorem 2, we show that the empirical measure of the output of
 266 VP-SVGD and GB-SVGD rapidly converges to π^* in KSD. We refer to Appendix E for the full
 267 statement and proof.

268 **Corollary 1 (VP-SVGD and GB-SVGD: Fast Finite Particle Rates).** Let the assumptions and
 269 parameter settings of Theorem 1 be satisfied. Let $\hat{\mu}^{(n)}$ be the empirical measures of the n particles
 270 output by VP-SVGD, run with run with $KT = d^{2+\alpha}$, $R = \sqrt{d/L}$ and appropriately chosen γ . Then:

$$\mathbb{E}[\text{KSD}_{\pi^*}^2(\hat{\mu}^{(n)} | \pi^*)] \leq O\left(\frac{2}{n^{2+\alpha}} + \frac{d^{2/\alpha}}{n}\right)$$

271 Let $\hat{\nu}^{(n)}$ be the empirical measure of the output of GB-SVGD under without-replacement sampling,
 272 run with $KT = \sqrt{n}$, $R = \sqrt{d/L}$ and appropriately chosen γ . Then, the following holds:

$$\mathbb{E}[\text{KSD}_{\pi^*}^2(\hat{\nu}^{(n)}|\pi^*)] \leq O\left(\frac{d^{2/\alpha}}{n} + \frac{1}{n^{\frac{1+2\alpha}{2(1+\alpha)}}} + \frac{d^{\frac{2+\alpha}{2(1+\alpha)}}}{n^{\frac{\alpha}{4(1+\alpha)}}}\right)$$

273 **Comparison to Prior Work** For subgaussian π^* (i.e. $\alpha = 2$), VP-SVGD has a finite-particle rate
 274 of $\mathbb{E}[\text{KSD}_{\pi^*}(\hat{\mu}^{(n)}|\pi^*)] = O((d/n)^{1/4} + (d/n)^{1/2})$ while that of GB-SVGD is $\mathbb{E}[\text{KSD}_{\pi^*}(\hat{\nu}^{(n)}|\pi^*)] =$
 275 $O(d^{1/3}/n^{1/12} + (d/n)^{1/2})$. Both these rates are a *double exponential improvement* over the
 276 $\tilde{O}\left(\frac{\text{poly}(d)}{\sqrt{\log \log n}^{\Theta(1/d)}}\right)$ KSD rate obtained by Shi and Mackey [37] for SVGD with subgaussian π^* .

277 For subexponential π^* (i.e. $\alpha = 1$) the KSD rate of VP-SVGD is $O\left(\frac{d^{1/3}}{n^{1/6}} + \frac{d}{n^{1/2}}\right)$ while that of
 278 GB-SVGD is $O\left(\frac{d^{3/8}}{n^{1/16}} + \frac{d}{n^{1/2}}\right)$. To our knowledge, both these results are the first of their kind.

279 **Oracle Complexity** As illustrated in Section E.3, for subgaussian π^* , the oracle complexity of
 280 VP-SVGD to achieve ϵ -convergence in KSD is $O(d^4/\epsilon^{12})$ and that of GB-SVGD is $O(d^6/\epsilon^{18})$. To our
 281 knowledge, these results are the *first known oracle complexities for this problem with polynomial*
 282 *dimension dependence*, and significantly improve upon the $O\left(\frac{\text{poly}(d)}{\epsilon^2} e^{\Theta(de^{\text{poly}(d)}/\epsilon^2)}\right)$ oracle complexity
 283 of SVGD as implied by Shi and Mackey [37]. For subexponential π^* , the oracle complexity of
 284 VP-SVGD is $O(d^6/\epsilon^{16})$ and that of GB-SVGD is $O(d^9/\epsilon^{24})$.

285 7 Proof Sketch

286 We now present a sketch of our analysis. As shown in Section 4, the particles $(\mathbf{x}_t^{(l)})_{l \geq Kt}$ are i.i.d
 287 conditioned on the filtration \mathcal{F}_t , and the random measure $\mu_t|\mathcal{F}_t$ is the law of $(\mathbf{x}_t^{(Kt)})$ conditioned on
 288 $\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(Kt-1)}$. Moreover, from equation (5), we know that $\mu_t|\mathcal{F}_t$ is a stochastic approximation
 289 of population limit SVGD dynamics, i.e., $\mu_{t+1}|\mathcal{F}_{t+1} = (I - \gamma g_t)_\# \mu_t|\mathcal{F}_t$. Lemma 1 (similar to
 290 Salim et al. [36, Proposition 3.1] and Korba et al. [23, Proposition 5]) shows that under appropriate
 291 conditions, the KL between $\mu_t|\mathcal{F}_t$ and π^* satisfies a (stochastic) descent lemma. Hence $\mu_t|\mathcal{F}_t$ admits
 292 a density and $\text{KL}(\mu_t|\mathcal{F}_t|\pi^*)$ is almost surely finite.

293 **Lemma 1** (Descent Lemma for $\mu_t|\mathcal{F}_t$). *Let Assumptions 1, 3 and 4 be satisfied and let $\beta > 1$ be an*
 294 *arbitrary constant. On the event $\gamma \|g_t\|_{\mathcal{H}} \leq \frac{\beta-1}{\beta B}$, the following holds almost surely*

$$\text{KL}(\mu_{t+1}|\mathcal{F}_{t+1}|\pi^*) \leq \text{KL}(\mu_t|\mathcal{F}_t|\pi^*) - \gamma \langle h_{\mu_t|\mathcal{F}_t}, g_t \rangle_{\mathcal{H}} + \frac{\gamma^2(\beta^2 + L)B}{2} \|g_t\|_{\mathcal{H}}^2$$

295 Lemma 1 is analogous to the noisy descent lemma which is used in the analysis of SGD for smooth
 296 functions. Notice that $\mathbb{E}[g_t|\mathcal{F}_t] = h_{\mu_t|\mathcal{F}_t}$ (when interpreted as a Gelfand-Pettis integral [40], as
 297 discussed in Appendix B and Appendix C) and hence in expectation, the KL divergence decreases in
 298 time. In order to apply Lemma 1, we establish an almost-sure bound on $\|g_t\|_{\mathcal{H}}$ below.

299 **Lemma 2.** *Let Assumptions 1, 2, 3 and 4 hold. Then, for $\gamma \leq 1/2A_1L$, $\|g_t\|_{\mathcal{H}} \leq \xi$ holds almost surely,*
 300 *where ξ is as defined in Theorem 1*

301 Let $K = 1$ for clarity. To prove Lemma 2, we first note via smoothness of $F(\cdot)$ and Assumption 3 that
 302 $\|g_t\|_{\mathcal{H}} \leq C_0 \|\mathbf{x}_t^{(t)}\| + C_1$, and then bound $\|\mathbf{x}_t^{(t)}\|$. Now, $g_s(\mathbf{x}) = k(\mathbf{x}, \mathbf{x}_s^{(s)}) \nabla F(\mathbf{x}_s^{(s)}) - \nabla_2 k(\mathbf{x}, \mathbf{x}_s^{(s)})$.
 303 When $\|\mathbf{x}_s^{(s)} - \mathbf{x}\|$ is large, $\|g_s(\mathbf{x})\|$ is small due to decay assumptions on the kernel (Assumption 3)
 304 implying that the particle does not move much. When $\mathbf{x}_s^{(s)} \approx \mathbf{x}$, we have $g_s(\mathbf{x}) \approx k(\mathbf{x}, \mathbf{x}_s^{(s)}) \nabla F(\mathbf{x}) -$
 305 $\nabla_2 k(\mathbf{x}, \mathbf{x}_s^{(s)})$ and $k(\mathbf{x}, \mathbf{x}_s^{(s)}) \geq 0$. This is approximately a gradient descent update on $F(\cdot)$ along with
 306 a bounded term $\nabla_2 k(\mathbf{x}, \mathbf{x}_s^{(s)})$. Thus, the value of $F(\mathbf{x}_t^{(t)})$ cannot grow too large after T iterations.
 307 By Assumption 2, $F(\mathbf{x}_t^{(t)})$ being small implies that $\|\mathbf{x}_t^{(t)}\|$ is small.

308 Equipped with Lemma 2, we set the step-size γ to ensure that the descent lemma (Lemma 1) always
 309 holds. The remainder of the proof involves unrolling through Lemma 1 by taking iterated expectations
 310 on both sides. To this end we control $\langle h_{\mu_t|\mathcal{F}_t}, g_t \rangle_{\mathcal{H}}$ and $\|g_t\|_{\mathcal{H}}^2$ in expectation, in Lemma 3.

311 **Lemma 3.** *Let Assumptions 1, 2, 3 and 4 hold and ξ be as defined in Theorem 1. Then, for $\gamma \leq 1/2A_1L$,*
 312 $\mathbb{E}[\langle h_{\mu_t|\mathcal{F}_t}, g_t \rangle_{\mathcal{H}} | \mathcal{F}_t] = \|h_{\mu_t|\mathcal{F}_t}\|_{\mathcal{H}}^2$ and $\mathbb{E}[\|g_t\|_{\mathcal{H}}^2] \leq \xi^2/K + \|h_{\mu_t|\mathcal{F}_t}\|_{\mathcal{H}}^2$

313 **8 Experiments**

314 We compare the performance of GB-SVGD and SVGD on the standard baselines used by prior work
 315 [27]. We take $n = 100$ and use the Laplace kernel with $h = 1$ for both the algorithms. We pick the
 316 stepsize γ by a grid search independently for each algorithm. For both our experimental setups, we
 317 observe that while SVGD takes fewer iterations to converge, the compute time for GB-SVGD is
 318 considerably lower. This is similar to the typical behavior of stochastic optimization algorithms like
 SGD.

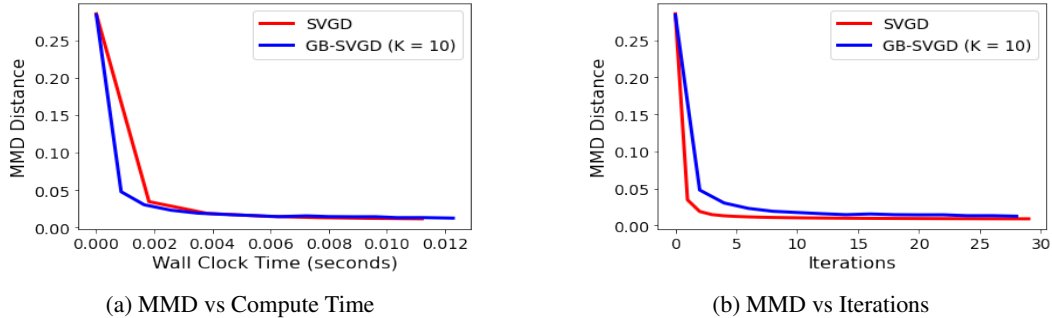


Figure 1: Gaussian Experiment Comparing SVGD and GB-SVGD averaged over 10 experiments.

319

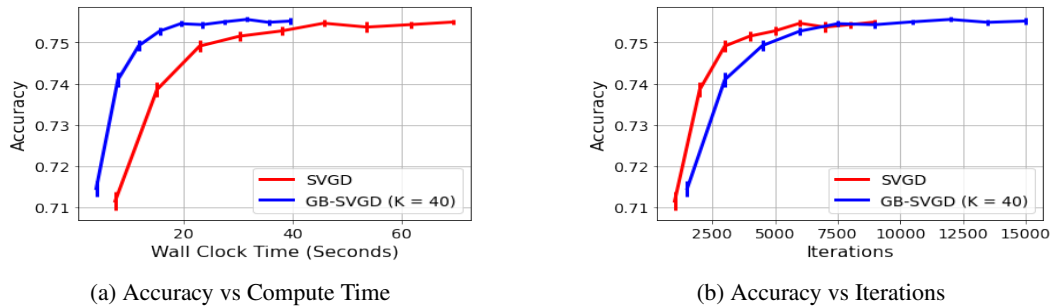


Figure 2: Covertypes Experiment, averaged over 50 runs. The error bars represent 95% CI.

320 **Sampling from Isotropic Gaussian (Figure 1):** As a sanity check, we set $\pi^* = \mathcal{N}(0, \mathbf{I})$ with $d = 5$.
 321 We pick $K = 10$ for GB-SVGD. The metric of convergence is MMD with respect to the empirical
 322 measure of 1000 i.i.d. sampled Gaussians.

323 **Bayesian Logistic Regression (Figure 2)** We consider the Covertypes dataset which contains \sim
 324 580,000 data points with $d = 54$. We consider the same priors suggested in Gershman et al. [16]
 325 and implemented in Liu and Wang [27]. We take $K = 40$ for GB-SVGD. For both VP-SVGD and
 326 GB-SVGD, we use AdaGrad with momentum to set the step-sizes as per Liu and Wang [27]

327 We ran our experiments using Python 3 on a 2.20 GHz Intel Xeon CPU with 13 GB of memory.

328 **9 Conclusion**

329 We develop two computationally efficient variants of SVGD with provably fast convergence guar-
 330 antees in the finite-particle regime, and present a wide range of improvements over prior work.
 331 A promising avenue of future work could be to establish convergence guarantees for SVGD with
 332 general non-logconcave targets, as was considered in recent works on LMC and SGLD [2, 12]. Other
 333 important avenues include establishing minimax lower bounds for SVGD and related particle-based
 334 variational inference algorithms. Beyond this, we also conjecture that the rates of GB-SVGD can be
 335 improved even in the regime $n \ll KT$. However, we believe this requires new analytic tools.

336 Acknowledgements

337 We thank Jiaxin Shi, Lester Mackey and the anonymous reviewers for their helpful feedback. We
338 are particularly grateful to Lester Mackey for providing insightful pointers on the properties of
339 Kernel Stein Discrepancy, which greatly helped us in removing the curse of dimensionality from our
340 finite-particle convergence guarantees.

341 References

- 342 [1] L. Ambrosio, N. Gigli, and G. Savaré. *Gradient flows: in metric spaces and in the space of*
343 *probability measures*. Springer Science & Business Media, 2005.
- 344 [2] K. Balasubramanian, S. Chewi, M. A. Erdogdu, A. Salim, and S. Zhang. Towards a theory of
345 non-log-concave sampling: first-order stationarity guarantees for langevin monte carlo. In P.-L.
346 Loh and M. Raginsky, editors, *Proceedings of Thirty Fifth Conference on Learning Theory*,
347 volume 178 of *Proceedings of Machine Learning Research*, pages 2896–2923. PMLR, 02–05
348 Jul 2022. URL <https://proceedings.mlr.press/v178/balasubramanian22a.html>.
- 349 [3] E. Bernton. Langevin monte carlo and jko splitting. In *Conference On Learning Theory*, pages
350 1777–1798. PMLR, 2018.
- 351 [4] S. Bobkov and M. Ledoux. Poincaré’s inequalities and talagrand’s concentration phenomenon
352 for the exponential distribution. *Probability Theory and Related Fields*, 107:383–400, 1997.
- 353 [5] S. G. Bobkov and F. Götze. Exponential integrability and transportation cost related to logarithmic
354 sobolev inequalities. *Journal of Functional Analysis*, 163(1):1–28, 1999.
- 355 [6] C. Carmeli, E. De Vito, A. Toigo, and V. Umanitá. Vector valued reproducing kernel hilbert
356 spaces and universality. *Analysis and Applications*, 8(01):19–61, 2010.
- 357 [7] S. Chewi, T. Le Gouic, C. Lu, T. Maunu, and P. Rigollet. Sgvd as a kernelized wasserstein
358 gradient flow of the chi-squared divergence. *Advances in Neural Information Processing*
359 *Systems*, 33:2098–2109, 2020.
- 360 [8] S. Chewi, M. A. Erdogdu, M. Li, R. Shen, and S. Zhang. Analysis of langevin monte carlo from
361 poincare to log-sobolev. In P.-L. Loh and M. Raginsky, editors, *Proceedings of Thirty Fifth Con-*
362 *ference on Learning Theory*, volume 178 of *Proceedings of Machine Learning Research*, pages
363 1–2. PMLR, 02–05 Jul 2022. URL [https://proceedings.mlr.press/v178/chewi22a.](https://proceedings.mlr.press/v178/chewi22a.html)
364 [html](https://proceedings.mlr.press/v178/chewi22a.html).
- 365 [9] K. Chwialkowski, H. Strathmann, and A. Gretton. A kernel test of goodness of fit. In
366 *International conference on machine learning*, pages 2606–2615. PMLR, 2016.
- 367 [10] J. B. Conway. *A course in functional analysis*, volume 96. Springer, 2019.
- 368 [11] A. Das, B. Schölkopf, and M. Muehlebach. Sampling without replacement leads to faster rates
369 in finite-sum minimax optimization. *Advances in Neural Information Processing Systems*, 2022.
- 370 [12] A. Das, D. Nagaraj, and A. Raj. Utilising the clt structure in stochastic gradient based sampling:
371 Improved analysis and faster algorithms. In *Conference on Learning Theory*, 2023.
- 372 [13] R. M. Dudley. *Real analysis and probability*. CRC Press, 2018.
- 373 [14] A. Duncan, N. Nüsken, and L. Szpruch. On the geometry of stein variational gradient descent.
374 *arXiv preprint arXiv:1912.00894*, 2019.
- 375 [15] A. El Alaoui, A. Montanari, and M. Sellke. Sampling from the sherrington-kirkpatrick gibbs
376 measure via algorithmic stochastic localization. In *2022 IEEE 63rd Annual Symposium on*
377 *Foundations of Computer Science (FOCS)*, pages 323–334. IEEE, 2022.
- 378 [16] S. Gershman, M. D. Hoffman, and D. M. Blei. Nonparametric variational inference. In
379 *Proceedings of the 29th International Conference on International Conference on Machine*
380 *Learning*, 2012.

- 381 [17] J. Gorham and L. Mackey. Measuring sample quality with kernels. In *International Conference*
382 *on Machine Learning*, pages 1292–1301. PMLR, 2017.
- 383 [18] J. Gorham, A. Raj, and L. Mackey. Stochastic stein discrepancies. *Advances in Neural*
384 *Information Processing Systems*, 33:17931–17942, 2020.
- 385 [19] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. *Advances in Neural*
386 *Information Processing Systems*, 33:6840–6851, 2020.
- 387 [20] P. Jain, D. M. Nagaraj, and P. Netrapalli. Making the last iterate of sgd information theoretically
388 optimal. *SIAM Journal on Optimization*, 31(2):1108–1130, 2021.
- 389 [21] P. Jains, L. Holdijk, and M. Welling. Learning equivariant energy based models with equivariant
390 stein variational gradient descent. *Advances in Neural Information Processing Systems*, 34:
391 16727–16737, 2021.
- 392 [22] Y. Kinoshita and T. Suzuki. Improved convergence rate of stochastic gradient langevin dynamics
393 with variance reduction and its application to optimization. *arXiv preprint arXiv:2203.16217*,
394 2022.
- 395 [23] A. Korba, A. Salim, M. Arbel, G. Luise, and A. Gretton. A non-asymptotic analysis for
396 stein variational gradient descent. *Advances in Neural Information Processing Systems*, 33:
397 4672–4682, 2020.
- 398 [24] H. J. Kushner and D. S. Clark. *Stochastic approximation methods for constrained and uncon-*
399 *strained systems*, volume 26. Springer Science & Business Media, 2012.
- 400 [25] Y. T. Lee and S. S. Vempala. The manifold joys of sampling (invited talk). In *49th International*
401 *Colloquium on Automata, Languages, and Programming (ICALP 2022)*. Schloss Dagstuhl-
402 Leibniz-Zentrum für Informatik, 2022.
- 403 [26] Q. Liu. Stein variational gradient descent as gradient flow. *Advances in neural information*
404 *processing systems*, 30, 2017.
- 405 [27] Q. Liu and D. Wang. Stein variational gradient descent: A general purpose bayesian inference
406 algorithm. *Advances in neural information processing systems*, 29, 2016.
- 407 [28] Q. Liu, J. Lee, and M. Jordan. A kernelized stein discrepancy for goodness-of-fit tests. In
408 *International conference on machine learning*, pages 276–284. PMLR, 2016.
- 409 [29] Y. Liu, P. Ramachandran, Q. Liu, and J. Peng. Stein variational policy gradient. *arXiv preprint*
410 *arXiv:1704.02399*, 2017.
- 411 [30] J. Lu, Y. Lu, and J. Nolen. Scaling limit of the stein variational gradient descent: The mean
412 field regime. *SIAM Journal on Mathematical Analysis*, 51(2):648–671, 2019.
- 413 [31] R. M. Neal et al. Mcmc using hamiltonian dynamics. *Handbook of markov chain monte carlo*,
414 2(11):2, 2011.
- 415 [32] Y. Nesterov. Introductory lectures on convex programming volume i: Basic course. *Lecture*
416 *notes*, 3(4):5, 1998.
- 417 [33] N. Nüsken and D. Renger. Stein variational gradient descent: many-particle and long-time
418 asymptotics. *arXiv preprint arXiv:2102.12956*, 2021.
- 419 [34] M. Raginsky, A. Rakhlin, and M. Telgarsky. Non-convex learning via stochastic gradient
420 langevin dynamics: a nonasymptotic analysis. In *Conference on Learning Theory*, pages
421 1674–1703. PMLR, 2017.
- 422 [35] G. O. Roberts and R. L. Tweedie. Exponential convergence of langevin distributions and their
423 discrete approximations. *Bernoulli*, pages 341–363, 1996.
- 424 [36] A. Salim, L. Sun, and P. Richtarik. A convergence theory for svgd in the population limit
425 under talagrand’s inequality t1. In *International Conference on Machine Learning*, pages
426 19139–19152. PMLR, 2022.

- 427 [37] J. Shi and L. Mackey. A finite-particle convergence rate for stein variational gradient descent.
428 *arXiv preprint arXiv:2211.09721*, 2022.
- 429 [38] I. Steinwart and A. Christmann. *Support vector machines*. Springer Science & Business Media,
430 2008.
- 431 [39] L. Sun, A. Karagulyan, and P. Richtarik. Convergence of stein variational gradient descent
432 under a weaker smoothness condition. In *International Conference on Artificial Intelligence*
433 *and Statistics*, pages 3693–3717. PMLR, 2023.
- 434 [40] M. Talagrand. *Pettis integral and measure theory*. American Mathematical Soc., 1984.
- 435 [41] S. Vempala and A. Wibisono. Rapid convergence of the unadjusted langevin algorithm:
436 Isoperimetry suffices. *Advances in neural information processing systems*, 32, 2019.
- 437 [42] R. Vershynin. *High-dimensional probability: An introduction with applications in data science*,
438 volume 47. Cambridge university press, 2018.
- 439 [43] D. Wang, Z. Zeng, and Q. Liu. Stein variational message passing for continuous graphical
440 models. In *International Conference on Machine Learning*, pages 5219–5227. PMLR, 2018.
- 441 [44] M. Welling and Y. W. Teh. Bayesian learning via stochastic gradient langevin dynamics.
442 In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages
443 681–688, 2011.
- 444 [45] A. Wibisono. Sampling as optimization in the space of measures: The langevin dynamics as a
445 composite optimization problem. In *Conference on Learning Theory*, pages 2093–3027. PMLR,
446 2018.

447	Contents	
448	1 Introduction	1
449	1.1 Contributions and Technical Challenges	2
450	1.2 Contributions	2
451	2 Notation and Problem Setup	3
452	3 Background on Mean-Field SVGD	4
453	4 Algorithm and Intuition	4
454	5 Assumptions	6
455	6 Results	6
456	6.1 VP-SVGD	6
457	6.2 GB-SVGD	7
458	6.3 Convergence of the Empirical Measure to the Target	7
459	7 Proof Sketch	8
460	8 Experiments	9
461	9 Conclusion	9
462	A Additional Notation and Organization	15
463	B Preliminaries	15
464	B.1 Gelfand-Pettis Integrals for Reproducing Kernel Hilbert Spaces	19
465	C Analysis of VP-SVGD	20
466	C.1 Population Level Dynamics : Proof of Lemma 1	20
467	C.2 Iterate Bounds : Proof of Lemma 2	21
468	C.3 Controlling g_t in Expectation : Proof of Lemma 3	22
469	C.4 Proof of Theorem 1	23
470	C.5 High-Probability Guarantees	24
471	D Analysis of GB-SVGD	27
472	D.1 Proof of Theorem 2	28
473	E Finite-Particle Convergence Guarantees for VP-SVGD and GB-SVGD	30
474	E.1 VP-SVGD	31
475	E.2 GB-SVGD	33
476	E.3 Oracle Complexity of SVGD, VP-SVGD and GB-SVGD	34
477	E.3.1 SVGD	34

478	E.3.2	VP-SVGD	34
479	E.3.3	GB-SVGD	35
480	F	Literature Review	35

481 A Additional Notation and Organization

482 We use Γ to denote the Gamma function $\Gamma(x) = \int_0^\infty t^{x-1} e^{-t} dt$, and recall that for any $n \in \mathbb{N}$,
 483 $\Gamma(n) = (n-1)!$. For any Lebesgue measurable $A \subseteq \mathbb{R}^d$, we use $\text{vol}(A)$ to denote it's Lebesgue
 484 Measure and $\text{Uniform}(A)$ to denote the uniform distribution supported on A . We use $\mathcal{B}(R)$ to denote
 485 the ball of radius R centered at the origin, and recall that $\text{vol}(\mathcal{B}(R)) = \frac{\pi^{d/2}}{\Gamma(d/2+1)} R^d$. For ease of
 486 exposition, we assume $d \geq 2$. We further assume $\pi^*(\mathbf{x}) = e^{-F(\mathbf{x})}$. We note that this can be easily
 487 ensured by absorbing the normalizing constant into $F(0)$, and does not affect the dynamics of SVGD,
 488 VP-SVGD or GB-SVGD (since they only use the gradient information of F). We highlight that
 489 both these assumptions are made purely for the sake of clarity and are very easily removable with
 490 negligible changes to our analysis.

491 We empirically benchmark SVGD and GB-SVGD in Appendix 8. In Appendix B, we discuss the
 492 technical lemmas used in our analysis, and present a short exposition to the Gelfand-Pettis integral
 493 in Appendix B.1, which we use to analyze VP-SVGD. We analyze VP-SVGD in Appendix C and
 494 GB-SVGD in Appendix D. Convergence guarantees for the empirical measure of VP-SVGD and
 495 GB-SVGD are presented in Appendix E. We give a brief review of the related work in Section F.

496 B Preliminaries

497 The following lemma shows that setting the initial distribution $\mu_0 = \text{Uniform}(\mathcal{B}(R))$ with $R = \sqrt{d/L}$
 498 suffices to ensure $\text{KL}(\mu_0 || \pi^*) = O(d)$. The proof of this result is similar to that of Vempala and
 499 Wibisono [41, Lemma 1] with the Gaussian initialization replaced by $\text{Uniform}(\mathcal{B}(R))$ initialization.

500 **Lemma 4 (KL Upper Bound for Uniform Initialization).** *Let Assumption 1 be satisfied and let*
 501 $\mu_0 = \text{Uniform}(\mathcal{B}(R))$ *with* $R = \sqrt{d/L}$. *Then, the following holds:*

$$\text{KL}(\mu_0 || \pi^*) \leq \frac{d}{2} \log(L/2\pi) + d + F(0) + 1/2 \leq O(d)$$

502 *Proof.* For any $\mathbf{x} \in \mathbb{R}^d$, the following holds by Assumption 1

$$\begin{aligned} F(\mathbf{x}) &\leq F(0) + \langle \nabla F(0), \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{x}\|^2 \\ &\leq F(0) + \sqrt{L} \|\mathbf{x}\| + \frac{L}{2} \|\mathbf{x}\|^2 \\ &\leq F(0) + 1/2 + L \|\mathbf{x}\|^2 \end{aligned}$$

503 where the second inequality uses $\|\nabla F(0)\| \leq \sqrt{L}$ and the Cauchy Schwarz inequality, and the last
 504 inequality uses the identity $ab \leq a^2 + b^2/4$. It follows that,

$$\mathbb{E}_{\mathbf{x} \sim \mu_0} [F(\mathbf{x})] \leq F(0) + 1/2 + LR^2$$

505 By a slight abuse of notation, let μ_0 denote the density of $\text{Uniform}(\mathcal{B}(R))$. Clearly, $\mu_0(\mathbf{x}) =$
 506 $\frac{1}{\text{vol}(\mathcal{B}(R))} \mathbb{1}_{\mathbf{x} \in \mathcal{B}(R)}$. It follows that,

$$\int_{\mathbb{R}^d} \mu_0(\mathbf{x}) \ln(\mu_0(\mathbf{x})) d\mathbf{x} = \int_{\mathcal{B}(R)} \frac{1}{\text{vol}(\mathcal{B}(R))} \log(1/\text{vol}(\mathcal{B}(R))) d\mathbf{x} = -\log(\text{vol}(\mathcal{B}(R)))$$

507 Now, $\text{vol}(\mathcal{B}(R)) = \frac{\pi^{d/2}}{\Gamma(d/2+1)} R^d$. Furthermore, by Stirling's Approximation, $(x/e)^{x-1} \leq \Gamma(x) \leq$
 508 $(x/2)^{x-1}$. Hence,

$$\frac{d}{2} \log\left(\frac{2\pi R^2}{d/2+1}\right) \leq \log(\text{vol}(\mathcal{B}(R))) \leq \frac{d}{2} \log\left(\frac{e\pi R^2}{d/2+1}\right)$$

509 Without loss of generality, assume $\pi^*(\mathbf{x}) = e^{-F(\mathbf{x})}$ (this can be easily ensured by appropriately
 510 adjusting $F(0)$ upto constant factors). It follows that,

$$\begin{aligned} \text{KL}(\mu_0 || \pi^*) &= \int_{\mathbb{R}^d} \mu_0(\mathbf{x}) \log\left(\frac{\mu_0(\mathbf{x})}{\pi^*(\mathbf{x})}\right) d\mathbf{x} = \int_{\mathbb{R}^d} \mu_0(\mathbf{x}) \ln(\mu_0(\mathbf{x})) d\mathbf{x} + \mathbb{E}_{\mathbf{x} \sim \mu_0} [F(\mathbf{x})] \\ &\leq -\log(\text{vol}(\mathcal{B}(R))) + F(0) + 1/2 + LR^2 \\ &\leq \frac{d}{2} \log\left(\frac{d/2+1}{2\pi R^2}\right) + F(0) + 1/2 + LR^2 \end{aligned}$$

511 Setting $R = \sqrt{d/L}$, we conclude that,

$$\text{KL}(\mu_0 \parallel \pi^*) \leq \frac{d}{2} \log(L/2\pi) + d + F(0) + 1/2 \leq O(d)$$

512

□

513 We now show that the growth condition on F , i.e. Assumption 2 is more general than specific
514 concentration assumptions on π^* (e.g. subgaussianity, subexponentiality etc.). To this end, we define
515 the notion of α -tail decay as follows:

516 **Definition 2** (α -Tail Decay). *A probability distribution ν on \mathbb{R}^d is said to satisfy α -tail decay for
517 some $\alpha > 0$ if there exists some $C > 0$ such that $\mathbb{E}_{\mathbf{x} \sim \nu} [\exp(\|\frac{\mathbf{x}}{C}\|^\alpha)] < \infty$*

518 The α -tail decay condition essentially implies that the tails of π^* decay as $\propto e^{-\|\mathbf{x}\|^\alpha}$. In particular,
519 Vershynin [42, Proposition 2.5.2 and Proposition 2.7.1] shows that π^* satisfying the tail decay
520 condition with $\alpha = 2$ is equivalent to π^* being subgaussian, whereas tail decay with $\alpha = 1$ is
521 equivalent to π^* being subexponential.

522 In the following lemma, we establish that, under smoothness of F , the α -tail decay condition is
523 equivalent to the growth condition on F with the same exponent α . Consequently, Assumption 2 is
524 much weaker than the standard isoperimetric and information transport assumptions generally used
525 in the literature.

526 **Lemma 5** (Growth Condition and Tail Decay). *Let Assumption 2 be satisfied for some $\alpha > 0$. Then,
527 π^* satisfies the α -tail decay condition. Conversely, let Assumption 1 be satisfied and suppose π^*
528 satisfies the α -tail decay condition. Then, F satisfies Assumption 2 with the same exponent α .*

529 *Proof.* **Growth Condition Implies Tail Decay** Since Assumption 2 is satisfied, $F(\mathbf{x}) \geq d_1 \|\mathbf{x}\|^\alpha - d_2$
530 for some $d_1, d_2, \alpha > 0$. Let $C = (2/d_1)^{1/\alpha}$. It follows that,

$$\begin{aligned} \mathbb{E}_{\mathbf{x} \sim \pi^*} [e^{\|\mathbf{x}/C\|^\alpha}] &= \int_{\mathbb{R}^d} e^{\frac{d_1}{2} \|\mathbf{x}\|^\alpha} \pi^*(\mathbf{x}) d\mathbf{x} \\ &\leq \int_{\mathbb{R}^d} e^{\frac{d_1}{2} \|\mathbf{x}\|^\alpha - d_1 \|\mathbf{x}\|^\alpha + d_2} d\mathbf{x} \\ &= e^{d_2} \int_{\mathbb{R}^d} e^{-\frac{d_1}{2} \|\mathbf{x}\|^\alpha} d\mathbf{x} < \infty \end{aligned}$$

531 From Definition 2, we conclude that π^* satisfies α -tail decay.

532 **Smoothness and Tail Decay Imply the Growth Condition** Since F is smooth, it suffices to consider
533 $\alpha \in (0, 2]$. By Assumption 1, the following inequalities hold,

$$F(\mathbf{y}) - F(\mathbf{x}) \leq \|\nabla F(\mathbf{x})\| \|\mathbf{y} - \mathbf{x}\| + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2 \leq (L\|\mathbf{x}\| + \sqrt{L}) \|\mathbf{y} - \mathbf{x}\| + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2 \quad (7)$$

534 We now prove this result by contradiction. Since π^* satisfies α -tail decay, there exists a constant
535 $C > 0$ such that $\mathbb{E}_{\mathbf{x} \sim \pi^*} [e^{\|\mathbf{x}/C\|^\alpha}] < \infty$. Now, suppose F does not satisfy the growth condition with
536 exponent α , i.e., *assume there does not exist* any $d_1, d_2 > 0$ such that $F(\mathbf{x}) \geq d_1 \|\mathbf{x}\|^\alpha - d_2 \forall \mathbf{x} \in \mathbb{R}^d$.
537 This implies that, $\liminf_{\|\mathbf{x}\| \rightarrow \infty} \frac{F(\mathbf{x})}{\|\mathbf{x}\|^\alpha} = 0$. Thus, without loss of generality, we can assume there
538 exists a diverging sequence $a_n \in \mathbb{R}$ and a diverging sequence $\mathbf{x}_n \in \mathbb{R}^d$ that satisfy the following for
539 every $n \in \mathbb{N}$:

$$\frac{F(\mathbf{x}_n)}{\|\mathbf{x}_n\|^\alpha} \leq \frac{1}{a_n}, \quad \|\mathbf{x}_n\| \geq 2n, \quad \|\mathbf{x}_{n+1} - \mathbf{x}_n\| \geq 1 \quad (8)$$

540 where, without loss of generality, we assume $a_n, \|\mathbf{x}_n\| > 0$. Now, let $r_n = \frac{1}{\|\mathbf{x}_n\|^2}$ and $B_n \subseteq \mathbb{R}^d$
541 denote the ball of radius r_n centered at \mathbf{x}_n . Since $r_n \leq 1/4n^2$ and $\|\mathbf{x}_{n+1} - \mathbf{x}_n\| \geq 1$, B_n is a family
542 of disjoint subsets of \mathbb{R}^d . We shall now prove that there exists some diverging sequence $b_n \in \mathbb{R}$ such
543 that $\frac{F(\mathbf{y})}{\|\mathbf{y}\|^\alpha} \leq \frac{1}{b_n}$ for every $\mathbf{y} \in B_n$.

544 Consider any arbitrary $n \in \mathbb{N}$ and let $\mathbf{y} \in B_n$. Applying (7) to \mathbf{y} and \mathbf{x}_n , we obtain,

$$\begin{aligned} \frac{F(\mathbf{y})}{\|\mathbf{x}_n\|^\alpha} &\leq \frac{F(\mathbf{x}_n)}{\|\mathbf{x}_n\|^\alpha} + \frac{L\|\mathbf{x}_n\|r_n}{\|\mathbf{x}_n\|^\alpha} + \frac{r_n\sqrt{L}}{\|\mathbf{x}_n\|^\alpha} + \frac{Lr_n^2}{2\|\mathbf{x}_n\|^\alpha} \\ &\leq \frac{1}{a_n} + \frac{L}{\|\mathbf{x}\|^{\alpha+1}} + \frac{\sqrt{L}}{\|\mathbf{x}_n\|^{\alpha+2}} + \frac{L}{2\|\mathbf{x}_n\|^{\alpha+4}} \end{aligned} \quad (9)$$

545 where we use (8) and $r_n = 1/\|\mathbf{x}_n\|^2$. Moreover, we note that

$$\|\mathbf{y}\| \geq \|\mathbf{x}_n\| - \|\mathbf{y} - \mathbf{x}_n\| \geq \|\mathbf{x}_n\| - r_n = \|\mathbf{x}_n\| - \frac{1}{\|\mathbf{x}_n\|^2} \geq \frac{\|\mathbf{x}_n\|}{2} \quad (10)$$

546 where we use the fact that $\|\mathbf{x}_n\| \geq 2n > 2^{1/3}$. It follows that,

$$\begin{aligned} \frac{F(\mathbf{y})}{\|\mathbf{y}\|^\alpha} &\leq \frac{2^\alpha F(\mathbf{y})}{\|\mathbf{x}_n\|^\alpha} \\ &\leq \frac{4}{a_n} + \frac{4L}{\|\mathbf{x}\|^{\alpha+1}} + \frac{4\sqrt{L}}{\|\mathbf{x}_n\|^{\alpha+2}} + \frac{2L}{\|\mathbf{x}_n\|^{\alpha+4}} \end{aligned} \quad (11)$$

547 where we use (9) and the fact that $\alpha \in (0, 2]$. We now define the sequence $b_n \in \mathbb{R}$ as follows:

$$b_n = \left(\frac{4}{a_n} + \frac{4L}{\|\mathbf{x}\|^{\alpha+1}} + \frac{4\sqrt{L}}{\|\mathbf{x}_n\|^{\alpha+2}} + \frac{2L}{\|\mathbf{x}_n\|^{\alpha+4}} \right)^{-1}$$

548 Since $\alpha > 0$, and $a_n, \|\mathbf{x}_n\| \rightarrow \infty$, it is clear that b_n is a diverging sequence. Furthermore, from (11),

549 we conclude that $\frac{F(\mathbf{y})}{\|\mathbf{y}\|^\alpha} \leq \frac{1}{b_n} \forall \mathbf{y} \in B_n$. Equipped with this construction, we note that

$$\begin{aligned} \mathbb{E}_{\mathbf{x} \sim \pi^*} \left[\exp \left(\frac{\|\mathbf{x}\|^\alpha}{C^\alpha} \right) \right] &= \int_{\mathbb{R}^d} \exp \left(\frac{\|\mathbf{y}\|^\alpha}{C^\alpha} \right) \exp(-F(\mathbf{y})) d\mathbf{y} \\ &\geq \sum_{n=1}^{\infty} \int_{B_n} \exp \left(\frac{\|\mathbf{y}\|^\alpha}{C^\alpha} \right) \exp(-F(\mathbf{y})) d\mathbf{y} \\ &\geq \sum_{n=1}^{\infty} \int_{B_n} \exp \left(\frac{\|\mathbf{y}\|^\alpha}{C^\alpha} - \frac{\|\mathbf{y}\|^\alpha}{b_n} \right) d\mathbf{y} \end{aligned}$$

550 where the second inequality use the fact that B_n is a disjoint family of subsets of \mathbb{R}^d and the third

551 inequality uses the fact that $\frac{F(\mathbf{y})}{\|\mathbf{y}\|^\alpha} \leq \frac{1}{b_n} \forall \mathbf{y} \in B_n$. Since b_n is a diverging sequence, there exists

552 some $N_0 \in \mathbb{N}$ such that $b_n \geq 2C^\alpha \forall n \geq N_0$. It follows that,

$$\begin{aligned} \mathbb{E}_{\mathbf{x} \sim \pi^*} \left[\exp \left(\frac{\|\mathbf{x}\|^\alpha}{C^\alpha} \right) \right] &\geq \sum_{n=1}^{\infty} \int_{B_n} \exp \left(\frac{\|\mathbf{y}\|^\alpha}{C^\alpha} - \frac{\|\mathbf{y}\|^\alpha}{b_n} \right) d\mathbf{y} \\ &\geq \sum_{n=N_0}^{\infty} \int_{B_n} \exp \left(\frac{\|\mathbf{y}\|^\alpha}{2C^\alpha} \right) d\mathbf{y} \\ &= \sum_{n=N_0}^{\infty} \text{vol}(B_n) \mathbb{E}_{\mathbf{y} \sim \text{Uniform}(B_n)} \left[\exp \left(\frac{\|\mathbf{y}\|^\alpha}{2C^\alpha} \right) \right] \end{aligned}$$

553 Consider the function $g : [0, \infty) \rightarrow [0, \infty)$ defined as $g(t) = e^{t^\alpha}$. We note that for $\alpha \geq 1$, g is a

554 convex function for every $t \geq 0$, and for $\alpha \in (0, 1)$, g is convex for every $t \geq (1/\alpha - 1)^{1/\alpha}$. From

555 (10), we note that $\|\mathbf{y}\| \geq \|\mathbf{x}\|/2 \geq n$ for every $\mathbf{y} \in B_n$. Hence, there exists an $N_1 \in \mathbb{N}$ such that e^{t^α}

556 is a convex function for all $t \geq \|\mathbf{y}\|/2$, $\forall \mathbf{y} \in B_n$, $n \geq N_1$. Let $N = \max\{N_0, N_1\} + 1$. Then,

$$\begin{aligned}
\mathbb{E}_{\mathbf{x} \sim \pi^*} \left[\exp \left(\frac{\|\mathbf{x}\|^\alpha}{C^\alpha} \right) \right] &\geq \sum_{n=N_0}^{\infty} \text{vol}(B_n) \mathbb{E}_{\mathbf{y} \sim \text{Uniform}(B_n)} \left[\exp \left(\frac{\|\mathbf{y}\|^\alpha}{2C^\alpha} \right) \right] \\
&\geq \sum_{n=N}^{\infty} \text{vol}(B_n) \exp \left(\frac{1}{2C^\alpha} \mathbb{E}_{\mathbf{y} \sim \text{Uniform}(B_n)} [\|\mathbf{y}\|^\alpha] \right) \\
&\geq \sum_{n=N}^{\infty} \text{vol}(B_n) \exp \left(\frac{1}{2C^\alpha} \|\mathbb{E}_{\mathbf{y} \sim \text{Uniform}(B_n)} [\mathbf{y}]\|^\alpha \right) \\
&\geq \sum_{n=N}^{\infty} \text{vol}(B_n) \exp \left(\frac{1}{2C^\alpha} \|\mathbf{x}_n\|^\alpha \right) \\
&= \sum_{n=N}^{\infty} C_d (r_n)^d \exp \left(\frac{1}{2C^\alpha} \|\mathbf{x}_n\|^\alpha \right) \\
&= \sum_{n=N}^{\infty} \frac{C_d \exp \left(\frac{1}{2C^\alpha} \|\mathbf{x}_n\|^\alpha \right)}{\|\mathbf{x}_n\|^{2d}}
\end{aligned}$$

557 where $C_d = \frac{\pi^{d/2}}{\Gamma(d/2+1)}$. Let k be any positive integer such that $\alpha k \geq 2d + 1$. It follows that,

$$\frac{\exp \left(\frac{1}{2C^\alpha} \|\mathbf{x}_n\|^\alpha \right)}{\|\mathbf{x}_n\|^{2d}} \geq \frac{C_d}{2^k k! C^{\alpha k}} \|\mathbf{x}_n\|^{\alpha k - 2d} \geq \frac{C_d n}{2^{k-1} k! C^{\alpha k}}$$

558 Thus, we infer that,

$$\mathbb{E}_{\mathbf{x} \sim \pi^*} \left[\exp \left(\frac{\|\mathbf{x}\|^\alpha}{C^\alpha} \right) \right] \geq \frac{C_d}{2^{k-1} k! C^{\alpha k}} \sum_{n=N_0}^{\infty} n = \infty$$

559 which is a contradiction. Thus, there exists some $d_1, d_2 > 0$ such that $F(\mathbf{x}) \geq d_1 \|\mathbf{x}\|^\alpha - d_2$, i.e., F
560 satisfies the growth condition with exponent α . \square

561 The following lemma establishes boundedness and contractivity properties of the function $h(\mathbf{x}, \mathbf{y}) =$
562 $k(\mathbf{x}, \mathbf{y}) \nabla F(\mathbf{y}) - \nabla_2 k(\cdot, \mathbf{y})$, that are vital for proving almost-sure bounds such as Lemma 2.

563 **Lemma 6** (Properties of h). *Let Assumptions 1 and 3 be satisfied. Then, the following holds,*

$$\begin{aligned}
\|h(\cdot, \mathbf{y})\|_{\mathcal{H}} &\leq BL\|\mathbf{y}\| + B\|\nabla F(0)\| + B \\
\|h(\mathbf{x}, \mathbf{y})\| &\leq \frac{A_1 L}{2} + A_2 + k(\mathbf{x}, \mathbf{y}) \|\nabla F(\mathbf{x})\| \\
-\langle \nabla F(\mathbf{x}), h(\mathbf{x}, \mathbf{y}) \rangle &\leq -\frac{1}{2} k(\mathbf{x}, \mathbf{y}) \|\nabla F(\mathbf{x})\|^2 + L^2 A_1 + A_3
\end{aligned}$$

564 *Proof.* Recalling the definition of h from Section 2, we observe that,

$$h(\cdot, \mathbf{y}) = k(\cdot, \mathbf{y}) \nabla F(\mathbf{y}) - \nabla_2 k(\cdot, \mathbf{y})$$

565 Thus, by triangle inequality of $\|\cdot\|_{\mathcal{H}}$, Assumptions 1 and 3, we obtain

$$\begin{aligned}
\|h(\cdot, \mathbf{y})\|_{\mathcal{H}} &\leq \|\nabla F(\mathbf{y})\| \|k(\cdot, \mathbf{y})\|_{\mathcal{H}_0} + \|\nabla_2 k(\cdot, \mathbf{y})\|_{\mathcal{H}} \\
&\leq BL\|\mathbf{y}\| + B\|\nabla F(0)\| + B
\end{aligned}$$

566 To prove the remaining inequalities, we first note that,

$$\begin{aligned}
h(\mathbf{x}, \mathbf{y}) &= k(\mathbf{x}, \mathbf{y}) \nabla F(\mathbf{y}) - \nabla_2 k(\mathbf{x}, \mathbf{y}) \\
&= k(\mathbf{x}, \mathbf{y}) \nabla F(\mathbf{x}) + k(\mathbf{x}, \mathbf{y}) [\nabla F(\mathbf{y}) - \nabla F(\mathbf{x})] - \nabla_2 k(\mathbf{x}, \mathbf{y})
\end{aligned} \tag{12}$$

567 Using Assumptions 1 and 3, we note that,

$$\begin{aligned}
\|h(\mathbf{x}, \mathbf{y})\| &\leq k(\mathbf{x}, \mathbf{y}) \|\nabla F(\mathbf{x})\| + \frac{LA_1 \|\mathbf{x} - \mathbf{y}\|}{1 + \|\mathbf{x} - \mathbf{y}\|^2} + A_2 \\
&\leq \frac{A_1 L}{2} + A_2 + k(\mathbf{x}, \mathbf{y}) \|\nabla F(\mathbf{x})\|
\end{aligned}$$

568 where the second inequality uses the fact $\frac{t}{1+t^2} \leq 1/2$

569 To prove the last inequality, we infer the following from (12)

$$\begin{aligned}
-\langle \nabla F(\mathbf{x}), h(\mathbf{x}, \mathbf{y}) \rangle &\leq -k(\mathbf{x}, \mathbf{y}) \|\nabla F(\mathbf{x})\|^2 + k(\mathbf{x}, \mathbf{y}) \|\nabla F(\mathbf{x}) - \nabla F(\mathbf{y})\| \|\nabla F(\mathbf{x})\| \\
&\quad + \|\nabla_2 k(\mathbf{x}, \mathbf{y})\| \|\nabla F(\mathbf{x})\| \\
&\leq -k(\mathbf{x}, \mathbf{y}) \|\nabla F(\mathbf{x})\|^2 + L \sqrt{k(\mathbf{x}, \mathbf{y})} \sqrt{\frac{A_1 \|\mathbf{x} - \mathbf{y}\|^2}{1 + \|\mathbf{x} - \mathbf{y}\|^2}} \|\nabla F(\mathbf{x})\| \\
&\quad + \sqrt{A_3 k(\mathbf{x}, \mathbf{y})} \|\nabla F(\mathbf{x})\| \\
&\leq -\frac{1}{2} k(\mathbf{x}, \mathbf{y}) \|\nabla F(\mathbf{x})\|^2 + L^2 A_1 + A_3
\end{aligned}$$

570 where the second inequality uses Assumptions 1 and 3, and the last inequality uses the identity
571 $ab \leq a^2 + b^2/4$ \square

572 To analyze the dynamics of VP-SVGD in the Wasserstein space, we use the following lemma
573 presented in Salim et al. [36]

574 **Lemma 7** (Salim et al. [36], Proposition 3.1). *Let Assumptions 1 and 3 be satisfied. Consider any*
575 *$\nu_0 \in \mathcal{P}_2(\mathbb{R}^d)$ with $\text{KL}(\nu_0 \|\pi^*) < \infty$, $f \in \mathcal{H}$ and let $\nu_1 = (I - \eta f)_\# \nu_0$ with $\eta \|f\|_{\mathcal{H}} \leq \frac{\beta-1}{\beta B}$ for*
576 *some $\beta > 1$. Then, the following holds,*

$$\text{KL}(\mu_1 \|\pi^*) \leq \text{KL}(\mu_0 \|\pi^*) - \eta \langle h_{\mu_0}, f \rangle + \frac{\eta^2 (\beta^2 + L) B}{2} \|f\|_{\mathcal{H}}^2$$

577 B.1 Gelfand-Pettis Integrals for Reproducing Kernel Hilbert Spaces

578 The Gelfand-Pettis integral is a generalization of the Lebesgue integral to functions that take values
579 in an arbitrary topological vector space. In this section, we describe the Gelfand-Pettis integral for
580 an arbitrary Hilbert space $(V, \langle \cdot, \cdot \rangle_V)$ and refer the readers to Talagrand [40] for a more general
581 treatment.

582 Let (X, Σ, λ) be a measure space and $(V, \langle \cdot, \cdot \rangle_V)$ be a Hilbert Space. A function $g : X \rightarrow V$
583 is said to be Gelfand-Pettis integrable if there exists a vector $w_g \in V$ such that $\langle u, w_g \rangle_V =$
584 $\int_X \langle u, g(x) \rangle_V d\lambda(x) \forall u \in V$. The vector w_g is called the Gelfand-Pettis integral of g

585 We now establish the following lemma for Gelfand-Pettis integrals with respect to the RKHS \mathcal{H} ,
586 which is a key component of our analysis of VP-SVGD.

587 **Lemma 8.** *Let μ be a probability measure on \mathbb{R}^d . Let $G : \mathbb{R}^d \times \mathbb{R}^d$ be a function such that*
588 *for every $\mathbf{y} \in \mathbb{R}^d$, $G(\cdot, \mathbf{y}) \in \mathcal{H}$ with $\|G(\cdot, \mathbf{y})\|_{\mathcal{H}} \leq C$ holding μ -almost surely. Let $G_\mu(\mathbf{x}) =$*
589 *$\mathbb{E}_{\mathbf{y} \sim \mu}[G(\mathbf{x}, \mathbf{y})]$. Then, the map $\psi : \mathbb{R}^d \rightarrow \mathcal{H}$ defined as $\psi(\mathbf{y}) = G(\cdot, \mathbf{y})$ is Gelfand-Pettis integrable*
590 *and G_μ is the Gelfand-Pettis integral of ψ with respect to μ , i.e. $G_\mu \in \mathcal{H}$ and for any $f \in \mathcal{H}$,*
591 *$\mathbb{E}_{\mathbf{y} \sim \mu}[\langle f, G(\cdot, \mathbf{y}) \rangle_{\mathcal{H}}] = \langle f, G_\mu \rangle_{\mathcal{H}}$*

592 *Proof.* Let $\Phi : \mathcal{H} \rightarrow \mathbb{R}$ denote the map $\Phi(f) = \mathbb{E}_{\mathbf{y} \sim \mu}[\langle f, G(\cdot, \mathbf{y}) \rangle_{\mathcal{H}}] \forall f \in \mathcal{H}$. By linear-
593 ity of expectations and inner products, we note that Φ is a linear functional on \mathcal{H} . Further-
594 more, since $\|G(\cdot, \mathbf{y})\|_{\mathcal{H}} \leq C$ holds μ -almost surely, we note that for any $f \in \mathcal{H}$, $|\Phi(f)| \leq$
595 $\mathbb{E}_{\mathbf{y} \sim \mu}[\langle f, G(\cdot, \mathbf{y}) \rangle_{\mathcal{H}}] \leq C \|f\|_{\mathcal{H}}$ by Jensen's inequality and Cauchy Schwarz inequality for \mathcal{H} . We
596 conclude that Φ is a bounded linear functional of \mathcal{H} . Thus, by Riesz Representation Theorem [10],
597 there exists $g \in \mathcal{H}$ such that for any $f \in \mathcal{H}$, the following holds

$$\mathbb{E}_{\mathbf{y} \in \mu}[\langle f, G(\cdot, \mathbf{y}) \rangle_{\mathcal{H}}] = \langle f, g \rangle_{\mathcal{H}}$$

598 Hence, we conclude that the map ψ is Gelfand-Pettis integrable. We now use the reproducing property
599 of \mathcal{H} to show that $g = G_\mu$, i.e., G_μ is the Gelfand-Pettis integral of ψ . To this end, let $\mathbf{x} \in \mathbb{R}^d$ be
600 arbitrary. Setting $f = k(\mathbf{x}, \cdot)$ and using the fact that $g \in \mathcal{H}$, $G(\cdot, \mathbf{y}) \in \mathcal{H}$ for any $\mathbf{y} \in \mathbb{R}^d$,

$$g(\mathbf{x}) = \mathbb{E}_{\mathbf{y} \in \mu}[G(\mathbf{x}, \mathbf{y})] = G_\mu(\mathbf{x})$$

601 Hence, $g = G_\mu$, i.e., $\mathbb{E}_{\mathbf{y} \sim \mu}[\langle f, G(\cdot, \mathbf{y}) \rangle_{\mathcal{H}}] = \langle f, G_\mu \rangle_{\mathcal{H}}$ \square

602 **C Analysis of VP-SVGD**

603 In this section, we present our analysis of VP-SVGD. Throughout this section, we define the random
 604 function $g_t : \mathbb{R}^d \times \mathbb{R}^d$ as $g_t(\mathbf{x}) = \frac{1}{K} \sum_{l=0}^{K-1} h(\mathbf{x}, \mathbf{x}_t^{(Kt+l)})$ where $t \in \mathbb{N} \cup \{0\}$, K is the batch-size of
 605 VP-SVGD, and $h : \mathbb{R}^d \times \mathbb{R}^d$ is as defined in Section 2, i.e., $h(\mathbf{x}, \mathbf{y}) = k(\mathbf{x}, \mathbf{y}) \nabla F(\mathbf{y}) - \nabla_2 k(\mathbf{x}, \mathbf{y})$.

606 After proving the key lemmas required for our analysis of VP-SVGD, we present the proof of
 607 Theorem 1 in Appendix C.4. We also present a high-probability version of Theorem 1 in Appendix
 608 C.5

609 **C.1 Population Level Dynamics : Proof of Lemma 1**

610 *Proof.* We now derive the population-limit dynamics of VP-SVGD for arbitrary batch-size K , and
 611 subsequently prove the descent lemma (i.e. Lemma 1) for VP-SVGD. The arguments of this section
 612 are a straightforward generalization of that used in Section 4.

613 To this end, we recall from Section 4 that the countably infinite number of particles $\mathbf{x}_0^{(l)}$, $l \in \mathbb{N} \cup \{0\}$
 614 are i.i.d samples from the measure μ_0 , which has a density w.r.t the Lebesgue measure. Thus, by
 615 the strong law of large numbers (Dudley [13, Theorem 11.4.1]), the empirical measure of $(\mathbf{x}_0^{(l)})_{l \geq 0}$
 616 is almost surely equal to μ_0 . Furthermore, we recall the filtration \mathcal{F}_t defined in Section 4 as
 617 $\mathcal{F}_t = \sigma(\mathbf{x}_0^{(l)} \mid l \leq Kt - 1)$, $t \in \mathbb{N}$ with \mathcal{F}_0 being the trivial σ algebra. We now consider the following
 618 dynamics in \mathbb{R}^d :

$$\mathbf{x}_{t+1}^{(s)} = \mathbf{x}_t^{(s)} - \frac{\gamma}{K} \sum_{l=0}^{K-1} h(\mathbf{x}_t^{(s)}, \mathbf{x}_t^{(tK+l)}), \quad s \in \mathbb{N} \cup \{0\} \quad (13)$$

619 We note that the above updates are the same as that of VP-SVGD for $s \in \{0, \dots, KT + n - 1\}$.
 620 Now, for each time-step t , we focus on the lower triangular evolution, i.e., the time evolution of
 621 the particles $(\mathbf{x}_t^{(l)})_{l \geq Kt}$. From (13), we infer that for any $t \in \mathbb{N}$ and $s \geq Kt$, $\mathbf{x}_t^{(s)}$ depends only on
 622 $(\mathbf{x}_0^{(l)})_{l \leq Kt-1}$ and $\mathbf{x}_0^{(s)}$. Hence, there exists a measurable function H_t for every $t \in \mathbb{N}$ such that the
 623 following holds almost surely:

$$\mathbf{x}_t^{(s)} = H_t(\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(Kt-1)}, \mathbf{x}_0^{(s)}); \quad \forall s \geq Kt \quad (14)$$

624 Since $\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(Kt-1)}, \mathbf{x}_0^{(s)}$ $\overset{i.i.d.}{\sim} \mu_0$, we conclude from (14) that $(\mathbf{x}_t^{(s)})_{s \geq Kt}$ are i.i.d when
 625 conditioned on $\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(Kt-1)}$. To this end, we define the random measure $\mu_t | \mathcal{F}_t$ as the law of
 626 $\mathbf{x}_t^{(Kt)}$ conditioned on \mathcal{F}_t , i.e. $\mu_t | \mathcal{F}_t$ is a probability kernel $\mu_t(\cdot; \mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(Kt-1)})$ with $\mu_0 | \mathcal{F}_0 := \mu_0$.
 627 By the strong law of large numbers, $\mu_t | \mathcal{F}_t$ is equal to the empirical measure of $(\mathbf{x}_t^{(l)})_{l \geq Kt}$ conditioned
 628 on \mathcal{F}_t . Furthermore, we infer from (13) that the particles satisfy the following:

$$\mathbf{x}_{t+1}^{(s)} = (I - \gamma g_t)(\mathbf{x}_t^{(s)}), \quad s \geq K(t+1)$$

629 Recall that $\mathbf{x}_{t+1}^{(s)} | \mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(K(t+1)-1)} \sim \mu_{t+1} | \mathcal{F}_{t+1}$ for any $s \geq K(t+1)$. Furthermore, from
 630 Equation (14), we note that for $s \geq K(t+1)$, $\mathbf{x}_t^{(s)}$ depends only on $\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(Kt-1)}$ and $\mathbf{x}_0^{(s)}$, which
 631 implies that $\text{Law}(\mathbf{x}_t^{(s)} | \mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(K(t+1)-1)}) = \text{Law}(\mathbf{x}_t^{(s)} | \mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(Kt-1)}) = \mu_t | \mathcal{F}_t$. Finally,
 632 we note that g_t is an \mathcal{F}_{t+1} -measurable random function. With these insights, we conclude that the
 633 population-level dynamics of the lower triangular evolution in $\mathcal{P}_2(\mathbb{R}^d)$ is almost surely described by
 634 the following update:

$$\mu_{t+1} | \mathcal{F}_{t+1} = (I - \gamma g_t) \# \mu_t | \mathcal{F}_t \quad (15)$$

635 Setting $\gamma \|g_t\|_{\mathcal{H}} \leq \frac{\beta-1}{\beta B}$ for some arbitrary $\beta > 1$ and applying Lemma 7 to the population-level
 636 update (15), we conclude that the following holds almost surely:

$$\text{KL}(\mu_{t+1} | \mathcal{F}_{t+1} \| \pi^*) \leq \text{KL}(\mu_t | \mathcal{F}_t \| \pi^*) - \gamma \langle h_{\mu_t | \mathcal{F}_t}, g_t \rangle_{\mathcal{H}} + \frac{\gamma^2(\beta^2 + L)B}{2} \|g_t\|_{\mathcal{H}}^2$$

637

□

638 **C.2 Iterate Bounds : Proof of Lemma 2**

639 To establish almost sure bounds on $\|g_t\|_{\mathcal{H}}$, we prove the following result which is stronger than
640 Lemma 2.

641 **Lemma 9** (Almost-Sure Iterate Bounds for VP-SVGD). *Let Assumptions 1, 2, 3 and 4 be satisfied.*
642 *Then, the following holds almost surely for any $s \in \mathbb{N} \cup \{0\}$ and $t \in (T + 1)$ whenever $\gamma \leq 1/2A_1L$*

$$\begin{aligned}\|\mathbf{x}_t^{(s)}\| &\leq \zeta_0 + \zeta_1(\gamma T)^{1/\alpha} + \zeta_2(\gamma^2 T)^{1/\alpha} + \zeta_3 R^{2/\alpha} \\ \|h(\cdot, \mathbf{x}_t^{(s)})\|_{\mathcal{H}} &\leq \zeta_0 + \zeta_1(\gamma T)^{1/\alpha} + \zeta_2(\gamma^2 T)^{1/\alpha} + \zeta_3 R^{2/\alpha} \\ \|g_t\|_{\mathcal{H}} &\leq \zeta_0 + \zeta_1(\gamma T)^{1/\alpha} + \zeta_2(\gamma^2 T)^{1/\alpha} + \zeta_3 R^{2/\alpha}\end{aligned}$$

643 where ζ_0, \dots, ζ_3 are problem-dependent constants that depend polynomially on
644 $A_1, A_2, A_3, B, d_1, d_2, L$ for any fixed α .

645 *Proof.* Let $c_t^{(s)} = \frac{1}{K} \sum_{l=0}^{K-1} k(\mathbf{x}_t^{(s)}, \mathbf{x}_t^{(Kt+l)})$. Note that by Assumption 3, $c_t^{(s)} \geq 0$ Since $\mathbf{x}_{t+1}^{(s)} =$
646 $\mathbf{x}_t^{(s)} - \gamma g_t(\mathbf{x}_t^{(s)})$, it follows from the smoothness of F that,

$$F(\mathbf{x}_{t+1}^{(s)}) - F(\mathbf{x}_t^{(s)}) \leq -\gamma \left\langle \nabla F(\mathbf{x}_t^{(s)}), g_t(\mathbf{x}_t^{(s)}) \right\rangle + \frac{\gamma^2 L}{2} \|g_t(\mathbf{x}_t^{(s)})\|^2 \quad (16)$$

647 By Lemma 6, we note that,

$$\begin{aligned}-\gamma \left\langle \nabla F(\mathbf{x}_t^{(s)}), g_t(\mathbf{x}_t^{(s)}) \right\rangle &= -\frac{\gamma}{K} \sum_{l=0}^{K-1} \left\langle \nabla F(\mathbf{x}_t^{(s)}), h(\mathbf{x}_t^{(s)}, \mathbf{x}_t^{(tK+l)}) \right\rangle \\ &\leq \frac{\gamma}{K} \sum_{l=0}^{L-1} \left[-\frac{1}{2} k(\mathbf{x}_t^{(s)}, \mathbf{x}_t^{(tK+l)}) \|\nabla F(\mathbf{x}_t^{(s)})\|^2 + L^2 A_1 + A_3 \right] \\ &\leq -\frac{\gamma c_t^{(s)}}{2} \|\nabla F(\mathbf{x}_t^{(s)})\|^2 + \gamma L^2 A_1 + \gamma A_3\end{aligned} \quad (17)$$

648 Moreover, by Jensen's Inequality and Lemma 6

$$\begin{aligned}\|g_t(\mathbf{x}_t^{(s)})\|^2 &\leq \frac{1}{K} \sum_{l=0}^{K-1} \|h(\mathbf{x}_t^{(s)}, \mathbf{x}_t^{(Kt+l)})\|^2 \\ &\leq \frac{1}{K} \sum_{l=0}^{K-1} 2(A_1 L/2 + A_2)^2 + 2k(\mathbf{x}_t^{(s)}, \mathbf{x}_t^{(tK+l)})^2 \|F(\mathbf{x}_t^{(s)})\|^2 \\ &\leq \frac{1}{K} \sum_{l=0}^{K-1} 2(A_1 L/2 + A_2)^2 + 2A_1 k(\mathbf{x}_t^{(s)}, \mathbf{x}_t^{(tK+l)}) \|F(\mathbf{x}_t^{(s)})\|^2 \\ &\leq 2(A_1 L/2 + A_2)^2 + 2A_1 c_t^{(s)} \|F(\mathbf{x}_t^{(s)})\|^2\end{aligned} \quad (18)$$

649 Substituting (17) and (18) into (16), we obtain,

$$\begin{aligned}F(\mathbf{x}_{t+1}^{(s)}) - F(\mathbf{x}_t^{(s)}) &\leq -\frac{\gamma c_t^{(s)}}{2} \|\nabla F(\mathbf{x}_t^{(s)})\|^2 + \gamma L^2 A_1 + \gamma A_3 \\ &\quad + \gamma^2 L(A_1 L/2 + A_2)^2 + \gamma^2 L A_1 c_t^{(s)} \|F(\mathbf{x}_t^{(s)})\|^2 \\ &\leq -\frac{\gamma c_t^{(s)}}{2} (1 - 2A_1 L \gamma) \|\nabla F(\mathbf{x}_t^{(s)})\|^2 + \gamma A_3 + \gamma L^2 A_1 + \gamma^2 L(A_1 L/2 + A_2)^2 \\ &\leq \gamma A_3 + \gamma L^2 A_1 + \gamma^2 L(A_1 L/2 + A_2)^2\end{aligned}$$

650 where the last inequality uses the fact that $c_t^{(s)} \geq 0$ and $\gamma \leq 1/2A_1L$. Now, iterating through the above
651 inequality, we obtain the following for any $t \in [T]$, $s \in \mathbb{N} \cup \{0\}$

$$F(\mathbf{x}_t^{(s)}) \leq F(\mathbf{x}_0^{(s)}) + \gamma T L^2 A_1 + \gamma T A_3 + \gamma^2 T L(A_1 L/2 + A_2)^2 \quad (19)$$

652 Furthermore, by Assumption 1

$$\begin{aligned} F(\mathbf{x}_0^{(s)}) &\leq F(0) + \|\nabla F(0)\| \|\mathbf{x}_0^{(s)}\| + \frac{L}{2} \|\mathbf{x}_0^{(s)}\|^2 \\ &\leq F(0) + 1/2 + L \|\mathbf{x}_0^{(s)}\|^2 \end{aligned}$$

653 Substituting the above inequality into (19), and using Assumption 2, we obtain the following for any
654 $t \in [T]$, $s \in \mathbb{N} \cup \{0\}$

$$\begin{aligned} d_1 \|\mathbf{x}_t^{(s)}\|^\alpha - d_2 &\leq F(\mathbf{x}_t^s) \leq F(0) + 1/2 + L \|\mathbf{x}_0^{(s)}\|^2 + \gamma T L^2 A_1 + \gamma T A_3 \\ &\quad + \gamma^2 T L (A_1 L/2 + A_2)^2 \end{aligned}$$

655 Rearranging and applying Assumption 4, we obtain

$$\begin{aligned} \|\mathbf{x}_t^{(s)}\| &\leq d_1^{-1/\alpha} [F(0) + 1/2 + L R^2 + \gamma T L^2 A_1 + \gamma T A_3 + \gamma^2 T L (A_1 L + A_2)^2]^{1/\alpha} \\ &\leq \tilde{\zeta}_0 + \tilde{\zeta}_1 (\gamma T)^{1/\alpha} + \tilde{\zeta}_2 (\gamma^2 T)^{1/\alpha} + \tilde{\zeta}_3 R^{2/\alpha} \end{aligned}$$

656 where $\tilde{\zeta}_0, \dots, \tilde{\zeta}_3$ are constants that depend polynomially on L, A_1, A_2, A_3, R . We note that, since
657 $0 < \alpha \leq 2$, the above inequality also holds for $t = 0$.

658 Using the above inequality Lemma 6 and Assumption 1, we conclude that the following holds almost
659 surely for any $t \in (T + 1)$, $s \in \mathbb{N} \cup \{0\}$

$$\begin{aligned} \|h(\cdot, \mathbf{x}_t^{(s)})\|_{\mathcal{H}} &\leq B L \|\mathbf{x}_t^{(s)}\| + B \sqrt{L} + B \\ &\leq \tilde{\eta}_0 + \tilde{\eta}_1 (\gamma T)^{1/\alpha} + \tilde{\eta}_2 (\gamma^2 T)^{1/\alpha} + \tilde{\eta}_3 R^{2/\alpha} \end{aligned}$$

660 where $\tilde{\eta}_0, \dots, \tilde{\eta}_3$ are constants that depend polynomially on L, B, A_1, A_2, A_3, R . Using the above
661 inequality, we conclude that the following also holds for any $t \in (T + 1)$.

$$\|g_t\|_{\mathcal{H}} \leq \tilde{\eta}_0 + \tilde{\eta}_1 (\gamma T)^{1/\alpha} + \tilde{\eta}_2 (\gamma^2 T)^{1/\alpha} + \tilde{\eta}_3 R^{2/\alpha}$$

662 Taking $\zeta_i = \max\{\tilde{\zeta}_i, \tilde{\eta}_i\}$, the proof is complete. \square

663 C.3 Controlling g_t in Expectation : Proof of Lemma 3

664 *Proof.* Let $\xi = \zeta_0 + \zeta_1 (\gamma T)^{1/\alpha} + \zeta_2 (\gamma^2 T)^{1/\alpha} + \zeta_3 R^{2/\alpha}$ where ζ_0, \dots, ζ_3 are as defined in Lemma
665 9. Recall that $g_t = \frac{1}{K} \sum_{l=0}^{K-1} h(\cdot, \mathbf{x}_t^{(Kt+l)})$. Since $\gamma \leq 1/2 A_1 L$, $\|h(\cdot, \mathbf{x}_t^{(Kt+l)})\|_{\mathcal{H}} \leq \xi$ holds almost
666 surely y Lemma 9.

667 Consider any $l \in (K)$. Conditioned on the filtration \mathcal{F}_t , $\text{Law}(\mathbf{x}_t^{(Kt+l)} | \mathcal{F}_t) = \mu_t | \mathcal{F}_t$. Moreover,
668 for any $\mathbf{x} \in \mathbb{R}^d$, $\mathbb{E}_{\mathbf{x}_t^{(Kt+l)}} [h(\mathbf{x}, \mathbf{x}_t^{(Kt+l)}) | \mathcal{F}_t] = h_{\mu_t | \mathcal{F}_t}(\mathbf{x})$. Thus, from Lemma 8, we conclude that
669 $h_{\mu_t | \mathcal{F}_t}$ is the Gelfand-Pettis Integral of the map $\mathbf{x} \rightarrow h(\mathbf{x}, \mathbf{x}_t^{(Kt+l)})$ with respect to $\mu_t | \mathcal{F}_t$. Hence, the
670 following holds

$$\mathbb{E}_{\mathbf{x}_t^{(Kt+l)}} \left[\left\langle h(\cdot, \mathbf{x}_t^{(Kt+l)}), f \right\rangle_{\mathcal{H}} \middle| \mathcal{F}_t \right] = \langle h_{\mu_t | \mathcal{F}_t}, f \rangle_{\mathcal{H}}$$

671 In particular, setting $f = h_{\mu_t | \mathcal{F}_t}$ and using linearity of expectation, we conclude,

$$\begin{aligned} \mathbb{E} \left[\langle g_t, h_{\mu_t | \mathcal{F}_t} \rangle_{\mathcal{H}} \middle| \mathcal{F}_t \right] &= \frac{1}{K} \sum_{l=0}^{K-1} \mathbb{E}_{\mathbf{x}_t^{(Kt+l)}} \left[\left\langle h(\cdot, \mathbf{x}_t^{(Kt+l)}), h_{\mu_t | \mathcal{F}_t} \right\rangle_{\mathcal{H}} \middle| \mathcal{F}_t \right] \\ &= \|h_{\mu_t | \mathcal{F}_t}\|_{\mathcal{H}}^2 \end{aligned}$$

672 To control $\mathbb{E}[\|g_t\|_{\mathcal{H}}^2 | \mathcal{F}_t]$, we note that,

$$\begin{aligned} \|g_t\|_{\mathcal{H}}^2 &= \frac{1}{K^2} \sum_{l_1, l_2=0}^{K-1} \left\langle h(\cdot, \mathbf{x}_t^{(Kt+l_1)}), h(\cdot, \mathbf{x}_t^{(Kt+l_2)}) \right\rangle_{\mathcal{H}} \\ &= \frac{1}{K^2} \sum_{l=0}^{K-1} \|h(\cdot, \mathbf{x}_t^{(Kt+l)})\|_{\mathcal{H}}^2 + \sum_{0 \leq l_1 \neq l_2 \leq K-1} \left\langle h(\cdot, \mathbf{x}_t^{(Kt+l_1)}), h(\cdot, \mathbf{x}_t^{(Kt+l_2)}) \right\rangle_{\mathcal{H}} \\ &\leq \frac{\xi^2}{K} + \sum_{0 \leq l_1 \neq l_2 \leq K-1} \left\langle h(\cdot, \mathbf{x}_t^{(Kt+l_1)}), h(\cdot, \mathbf{x}_t^{(Kt+l_2)}) \right\rangle_{\mathcal{H}} \end{aligned}$$

673 where the last inequality uses the fact that $\|h(\cdot, \mathbf{x}_t^{(Kt+l)})\|_{\mathcal{H}} \leq \xi$ almost surely as per Lemma 9.

674 To control the off-diagonal terms, let $i = Kt + l_1$ and $j = Kt + l_2$ for any arbitrary l_1, l_2 with
 675 $0 \leq l_1 \neq l_2 \leq K - 1$. Conditioned on \mathcal{F}_t , $\mathbf{x}_t^{(i)}$ and $\mathbf{x}_t^{(j)}$ are i.i.d samples from $\mu_t | \mathcal{F}_t$. Thus, by
 676 Lemma 8 and Fubini's Theorem,

$$\begin{aligned} \mathbb{E}_{\mathbf{x}_t^{(i)}, \mathbf{x}_t^{(j)}} \left[\left\langle h(\cdot, \mathbf{x}_t^{(i)}), h(\cdot, \mathbf{x}_t^{(j)}) \right\rangle_{\mathcal{H}} \middle| \mathcal{F}_t \right] &= \mathbb{E}_{\mathbf{x}_t^{(i)}} \left[\mathbb{E}_{\mathbf{x}_t^{(j)}} \left[\left\langle h(\cdot, \mathbf{x}_t^{(i)}), h(\cdot, \mathbf{x}_t^{(j)}) \right\rangle_{\mathcal{H}} \middle| \mathcal{F}_t \right] \right] \\ &= \mathbb{E}_{\mathbf{x}_t^{(i)}} \left[\left\langle h_{\mu_t | \mathcal{F}_t}, h(\cdot, \mathbf{x}_t^{(i)}) \right\rangle_{\mathcal{H}} \middle| \mathcal{F}_t \right] \\ &= \|h_{\mu_t | \mathcal{F}_t}\|_{\mathcal{H}}^2 \end{aligned}$$

677 Thus, we conclude that,

$$\mathbb{E} [\|g_t\|_{\mathcal{H}}^2 | \mathcal{F}_t] \leq \|h_{\mu_t | \mathcal{F}_t}\|_{\mathcal{H}}^2 + \frac{\xi^2}{K}$$

678 □

679 C.4 Proof of Theorem 1

680 *Proof.* Let $\xi = \zeta_0 + \zeta_1(\gamma T)^{1/\alpha} + \zeta_2(\gamma^2 T)^{1/\alpha} + \zeta_3 R^{2/\alpha}$ where ζ_0, \dots, ζ_3 are as defined in Lemma 9.
 681 Since $\gamma \leq 1/2A_1L$, $\|g_t\|_{\mathcal{H}} \leq \xi$ holds almost surely as per Lemma 9.

682 Since $\gamma\xi \leq 1/2B$, Lemma 1 ensures that the following holds almost surely

$$\text{KL}(\mu_{t+1} | \mathcal{F}_{t+1} | \pi^*) \leq \text{KL}(\mu_t | \mathcal{F}_t | \pi^*) - \gamma \langle h_{\mu_t | \mathcal{F}_t}, g_t \rangle_{\mathcal{H}} + \frac{\gamma^2(4+L)B}{2} \|g_t\|_{\mathcal{H}}^2$$

683 Taking conditional expectations w.r.t \mathcal{F}_t on both sides and applying Lemma 3, we obtain,

$$\begin{aligned} \mathbb{E} [\text{KL}(\mu_{t+1} | \mathcal{F}_{t+1} | \pi^*) | \mathcal{F}_t] &\leq \text{KL}(\mu_t | \mathcal{F}_t | \pi^*) - \gamma \left(1 - \frac{\gamma(4+L)B}{2} \right) \|h_{\mu_t | \mathcal{F}_t}\|_{\mathcal{H}}^2 + \frac{\gamma^2(4+L)B\xi^2}{2K} \\ &\leq \text{KL}(\mu_t | \mathcal{F}_t | \pi^*) - \frac{\gamma}{2} \|h_{\mu_t | \mathcal{F}_t}\|_{\mathcal{H}}^2 + \frac{\gamma^2(4+L)B\xi^2}{2K} \\ &= \text{KL}(\mu_t | \mathcal{F}_t | \pi^*) - \frac{\gamma}{2} \text{KSD}_{\pi^*}(\mu_t | \mathcal{F}_t | \pi^*)^2 + \frac{\gamma^2(4+L)B\xi^2}{2K} \end{aligned}$$

684 where the second inequality uses the fact that $\gamma \leq 1/(4+L)B$. Taking expectations on both sides and
 685 rearranging,

$$\frac{\gamma}{2} \mathbb{E} [\text{KSD}_{\pi^*}(\mu_t | \mathcal{F}_t | \pi^*)^2] \leq \mathbb{E} [\text{KL}(\mu_t | \mathcal{F}_t | \pi^*) - \text{KL}(\mu_{t+1} | \mathcal{F}_{t+1} | \pi^*)] + \frac{\gamma^2(4+L)B\xi^2}{2K}$$

686 Telescoping and averaging, we conclude,

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} [\text{KSD}_{\pi^*}(\mu_t | \mathcal{F}_t | \pi^*)^2] \leq \frac{2\text{KL}(\mu_0 | \mathcal{F}_0 | \pi^*)}{\gamma T} + \frac{\gamma(4+L)B\xi^2}{K} \quad (20)$$

687 Now, recall from the proof of Lemma 1 in Section C.1 that for any $t \in [T]$ and $l \geq Kt$, $\mathbf{x}_t^{(l)}$
 688 depends only on $\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(Kt-1)}, \mathbf{x}_0^{(l)}$, i.e., there exists a deterministic measurable function H_t
 689 such that $\mathbf{x}_t^{(l)} = H_t(\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(Kt-1)}, \mathbf{x}_0^{(l)})$ holds almost surely. We note that the output $\mathbf{Y} =$
 690 $(\mathbf{y}^{(0)}, \dots, \mathbf{y}^{(n-1)})$ satisfies $\mathbf{y}^{(l)} = \mathbf{x}_S^{(Kt+l)} \forall l \in (n)$, where $S \sim \text{Uniform}((T))$ is sampled
 691 independently of everything else.

692 Thus, we infer that $\mathbf{y}^{(l)}$ depends only on $\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(KT-1)}, S, \mathbf{x}_0^{(KT+l)}$, i.e., there exists a determin-
 693 istic measurable function G such that $\mathbf{y}^{(l)} = G(\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(KT-1)}, S, \mathbf{x}_0^{(KT+l)})$ for every $l \in (n)$.
 694 Since $\mathbf{x}_0^{(KT)}, \dots, \mathbf{x}_0^{(KT+n-1)} \stackrel{i.i.d.}{\sim} \mu_0$, we infer that $\mathbf{y}^{(0)}, \dots, \mathbf{y}^{(n-1)}$ are i.i.d when conditioned on
 695 $\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(KT-1)}, S$.

696 We now show that, when conditioned on $\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(KT-1)}, S$, $\mathbf{y}^{(l)}$ is distributed as $\bar{\mu}$, where $\bar{\mu}$ is
 697 the probability kernel defined as $\bar{\mu}(\cdot; \mathbf{x}_0^{(0)} = x_0, \dots, \mathbf{x}_0^{(KT-1)} = \mathbf{x}_{KT-1}, S = s) := \mu_s(\cdot, \mathbf{x}_0^{(0)} =$

698 $\mathbf{x}_0, \dots, \mathbf{x}_0^{(Ks-1)} = \mathbf{x}_{Ks-1}$. For any arbitrary fixed $s \in (T)$, note that, under the event $S = s$,
699 $\mathbf{y}^{(l)} = \mathbf{x}_s^{(KT+l)}$ for every $l \in (n)$. Thus, for any Borel measurable set $A \subseteq \mathbb{R}^d$, $\{\mathbf{y}^{(l)} \in A\} \cap$
700 $\{S = s\} = \{\mathbf{x}_s^{(KT+l)} \in A\} \cap \{S = s\}$. For the sake of clarity, we denote the conditioning
701 $\mathbf{x}_0^{(0)} = \mathbf{x}_0, \mathbf{x}_0^{(KT-1)} = \mathbf{x}_{KT-1}$ as \mathcal{C} , only in this proof. Since S is independent of $\mathbf{x}_t^{(l)}$ for every
702 $t \in (T+1), l \in (KT+n)$, we infer the following:

$$\begin{aligned} \mathbb{P}\left(\{\mathbf{y}^{(l)} \in A\} | \mathcal{C}, S = s\right) &= \frac{\mathbb{P}\left(\{\mathbf{y}^{(l)} \in A\} \cap \{S = s\} | \mathcal{C}\right)}{\mathbb{P}(S = s)} \\ &= T\mathbb{P}\left(\{\mathbf{x}_s^{(KT+l)} \in A\} \cap \{S = s\} | \mathcal{C}\right) \\ &= T\mathbb{P}(S = s) \mathbb{P}\left(\{\mathbf{x}^{(KT+l)} \in A\} | \mathcal{C}\right) \\ &= \mathbb{P}\left(\{\mathbf{x}_s^{(KT+l)} \in A\} | \mathcal{C}\right) \end{aligned}$$

703 As discussed above, $\mathbf{x}_s^{(KT+l)}$ depends only on $\mathbf{x}_0^{(0)}, \mathbf{x}_0^{(Ks-1)}, \mathbf{x}_0^{(KT+l)}$. It follows that
704 $\mathbb{P}\left(\{\mathbf{x}_s^{(KT+l)} \in A\} | \mathcal{C}\right) = \mu_s(A; \mathbf{x}_0^{(0)} = x_0, \dots, \mathbf{x}_0^{(Ks-1)} = x_{Ks-1})$ and,

$$\begin{aligned} \mathbb{P}\left(\{\mathbf{y}^{(l)} \in A\} | \mathcal{C}, S = s\right) &= \mu_s(A; \mathbf{x}_0^{(0)} = x_0, \dots, \mathbf{x}_0^{(Ks-1)} = x_{Ks-1}) \\ &= \bar{\mu}(A; \mathbf{x}_0^{(0)} = x_0, \dots, \mathbf{x}_0^{(KT-1)} = x_{KT-1}, S = s) \end{aligned}$$

705 Thus, $\mathbf{y}^{(0)}, \dots, \mathbf{y}^{(n-1)}$ are i.i.d samples from $\bar{\mu}$ when conditioned on $\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(KT-1)}, S$.

706 We now obtain an upper bound on the expected squared KSD between $\bar{\mu}$ and π^* . We recall from
707 the proof of Lemma 1 in Section C.1 that, for any $t \in (T+1)$, conditioned on $\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(Kt-1)}$,
708 $(\mathbf{x}_t^{(l)})_{l \geq t}$ are i.i.d samples from $\mu_t | \mathcal{F}_t$ where $\mu_t | \mathcal{F}_t := \mu_t(\cdot; \mathbf{x}_0^{(0)}, \mathbf{x}_0^{(Kt-1)})$. Hence, from (20), we
709 conclude that,

$$\begin{aligned} \mathbb{E}[\text{KSD}_{\pi^*}(\bar{\mu}(\cdot; (\mathbf{x}_0^{(l)})_{l \in (KT)}, S) | \pi^*)^2] &= \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[\mathbb{E} \left[\text{KSD}_{\pi^*}(\bar{\mu}(\cdot; (\mathbf{x}_0^{(l)})_{l \in (KT)}, S = t) | \pi^*)^2 | (\mathbf{x}_0^{(l)})_{l \in (KT)} \right] \right] \\ &= \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[\text{KSD}_{\pi^*}(\mu_t(\cdot; \mathbf{x}_0^{(0)}, \cdot, \mathbf{x}_0^{(Kt-1)}) | \pi^*)^2 \right] \\ &= \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[\text{KSD}_{\pi^*}(\mu_t | \mathcal{F}_t | \pi^*)^2 \right] \\ &\leq \frac{2\text{KL}(\mu_0 | \mathcal{F}_0 | \pi^*)}{\gamma T} + \frac{\gamma(4+L)B\xi^2}{K} \end{aligned}$$

710 where we use the fact that $S \sim \text{Uniform}((T))$ is sampled independent of everything else. \square

711 C.5 High-Probability Guarantees

712 We establish the convergence guarantee for VP-SVGD which holds with high probability, when
713 conditioned on the virtual particles $\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(KT-1)}$

714 **Theorem 3 (VP-SVGD: High-Probability Rates).** *Let the assumptions and parameter settings of*
715 *Theorem 1 apply and let $\delta \in (0, 1)$. Then, the following holds with probability at least $1 - \delta$:*

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} \text{KSD}_{\pi^*}(\mu_t | \mathcal{F}_t | \pi^*)^2 &\leq \frac{4\text{KL}(\mu_0 | \mathcal{F}_0 | \pi^*)}{\gamma T} + \frac{2\gamma(4+L)B\xi^2}{K} \\ &\quad + \frac{32\xi^2 \log(2/\delta)}{KT} + 12\gamma(4+L)B\xi^2 \sqrt{\frac{\log(2/\delta)}{T}} \end{aligned}$$

716 Let $\bar{\mu}(\cdot; \mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(KT-1)}, S)$ be the probability kernel defined in the statement of Theorem 1. Then,
717 conditioned on $\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(KT-1)}, S$, the n particles output by VP-SVGD are i.i.d samples from

718 $\bar{\mu}(\cdot; \mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(KT-1)}, S)$. Furthermore, with probability at least $1 - \delta$

$$\begin{aligned} \mathbb{E}_S[\text{KSD}_{\pi^*}(\bar{\mu}(\cdot; \mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(KT-1)}, S) \|\pi^*)^2] &\leq \frac{4\text{KL}(\mu_0|\mathcal{F}_0\|\pi^*)}{\gamma T} + \frac{2\gamma(4+L)B\xi^2}{K} \\ &\quad + \frac{32\xi^2 \log(2/\delta)}{KT} + 12\gamma(4+L)B\xi^2 \sqrt{\frac{\log(2/\delta)}{T}} \end{aligned}$$

719 where \mathbb{E}_S denotes that the expectation is being taken only with respect to $S \sim \text{Uniform}((T))$

720 *Proof.* Following the same steps as Theorem 1, we note that the following holds almost surely.

$$\begin{aligned} \text{KL}(\mu_{t+1}|\mathcal{F}_{t+1}\|\pi^*) &\leq \text{KL}(\mu_t|\mathcal{F}_t\|\pi^*) - \gamma \langle h_{\mu_t|\mathcal{F}_t}, g_t \rangle_{\mathcal{H}} + \frac{\gamma^2(4+L)B}{2} \|g_t\|_{\mathcal{H}}^2 \\ &\leq \text{KL}(\mu_t|\mathcal{F}_t\|\pi^*) - \frac{\gamma}{2} \|h_{\mu_t|\mathcal{F}_t}\|_{\mathcal{H}}^2 + \gamma \langle h_{\mu_t|\mathcal{F}_t}, h_{\mu_t|\mathcal{F}_t} - g_t \rangle_{\mathcal{H}} \\ &\quad + \frac{\gamma^2(4+L)B\xi^2}{2K} + \frac{\gamma^2(4+L)B}{2} \left[\|g_t\|_{\mathcal{H}}^2 - \|h_{\mu_t|\mathcal{F}_t}\|_{\mathcal{H}}^2 - \frac{\xi^2}{K} \right] \end{aligned} \quad (21)$$

721 where the last inequality uses the fact that $\gamma \leq 1/(4+L)B$. We now define $\Delta_t^{(l)}$, Δ_t and r_t for
722 $l \in (K)$, $t \in (T)$ as follows:

$$\begin{aligned} \Delta_t^{(l)} &= \left\langle h_{\mu_t|\mathcal{F}_t}, h_{\mu_t|\mathcal{F}_t} - h(\cdot, \mathbf{x}_t^{(Kt+l)}) \right\rangle_{\mathcal{H}} \\ \Delta_t &= \frac{1}{K} \sum_{l=0}^{K-1} \Delta_t^{(l)} = \langle h_{\mu_t|\mathcal{F}_t}, h_{\mu_t|\mathcal{F}_t} - g_t \rangle_{\mathcal{H}} \\ r_t &= \|g_t\|_{\mathcal{H}}^2 - \|h_{\mu_t|\mathcal{F}_t}\|_{\mathcal{H}}^2 - \frac{\xi^2}{K} \end{aligned}$$

723 Substituting the above into (21), we obtain the following:

$$\text{KL}(\mu_{t+1}|\mathcal{F}_{t+1}\|\pi^*) \leq \text{KL}(\mu_t|\mathcal{F}_t\|\pi^*) - \frac{\gamma}{2} \|h_{\mu_t|\mathcal{F}_t}\|_{\mathcal{H}}^2 + \gamma \Delta_t + \frac{\gamma^2(4+L)B\xi^2}{2K} + \frac{\gamma^2(4+L)Br_t}{2}$$

724 Telescoping and averaging both sides, and using $\|h_{\mu_t|\mathcal{F}_t}\|_{\mathcal{H}}^2 = \text{KSD}_{\pi^*}(\mu_t|\mathcal{F}_t\|\pi^*)^2$, we obtain the
725 following:

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} \text{KSD}_{\pi^*}(\mu_t|\mathcal{F}_t\|\pi^*)^2 &\leq \frac{4\text{KL}(\mu_0|\mathcal{F}_0\|\pi^*)}{\gamma T} + \frac{2\gamma(4+L)B\xi^2}{K} \\ &\quad + \frac{4}{T} \sum_{t=0}^{T-1} \left(\Delta_t - \frac{\|h_{\mu_t|\mathcal{F}_t}\|_{\mathcal{H}}^2}{4} \right) + \frac{2\gamma(4+L)B}{T} \sum_{t=0}^{T-1} r_t \end{aligned} \quad (22)$$

726 We note that the first two terms are the same as that of the in-expectation guarantee for VP-SVGd in
727 Theorem 1. The third and fourth term are random quantities that vanish in expectation. The remainder
728 of our analysis upper bounds them with high probability.

729 We begin by deriving a high probability upper bound for the fourth term in (22). To this end, we
730 note that, since $\gamma \leq 1/2A_1L$, $\|h(\cdot, \mathbf{x}_t^{(Kt+l)})\|_{\mathcal{H}} \leq \xi$ for any $t \in (T)$, $l \in (K)$ as per Lemma 9.
731 Furthermore, since $\mathbb{E}[h(\cdot, \mathbf{x}_t^{(Kt+l)})|\mathcal{F}_t] = h_{\mu_t|\mathcal{F}_t}$ (both pointwise and in the sense of the Gelfand-
732 Pettis integral, see proof of Lemma 3 in Appendix C.3), it follows by Jensen's inequality that
733 $\|h_{\mu_t|\mathcal{F}_t}\|_{\mathcal{H}} \leq \xi$. This further implies that $|r_t| \leq 3\xi^2$. Moreover, r_t is \mathcal{F}_{t+1} measurable (as g_t is an
734 \mathcal{F}_{t+1} measurable random function) with $\mathbb{E}[r_t|\mathcal{F}_t] \leq 0$ (as per Lemma 3)

735 Thus, $S_t = \sum_{s=0}^{t-1} r_s$ is an \mathcal{F} -adapted supermartingale difference sequence with bounded increments.
736 Thus, by the Hoeffding-Azuma inequality, we conclude that the following holds with probability at
737 least $1 - \delta/2$

$$\frac{1}{T} \sum_{t=0}^{T-1} r_t \leq 6\xi^2 \sqrt{\frac{\log(2/\delta)}{T}} \quad (23)$$

738 We now proceed to control the third term in (22). Recall from the proof of Theorem 1 in Appendix
 739 C.4, that, for any fixed $t \in (T)$, $(\mathbf{x}_t^{(l)})_{l \in (KT)}$ are i.i.d when conditioned on \mathcal{F}_t . As discussed
 740 above, $\mathbb{E}[h(\cdot, \mathbf{x}_t^{(Kt+l)})] = h_{\mu_t|\mathcal{F}_t}$ in the sense of the Gelfand-Pettis integral, implying $\mathbb{E}[\Delta_t^{(l)}] = 0$.
 741 Moreover, $|\Delta_t^{(l)}|_t \leq 2\xi \|h_{\mu_t|\mathcal{F}_t}\|$. Thus, when conditioned on \mathcal{F}_t , $\Delta_t^{(l)}$ are independent zero-mean
 742 bounded random variables. Hence, we conclude the following by Hoeffding's Lemma

$$\mathbb{E}[e^{\theta \Delta_t} | \mathcal{F}_t] \leq \prod_{l=0}^{K-1} \mathbb{E}[e^{\frac{\theta \Delta_t^{(l)}}{K}} | \mathcal{F}_t] \leq e^{\frac{2\theta^2 \xi^2}{K} \|h_{\mu_t|\mathcal{F}_t}\|_{\mathcal{H}}^2}, \quad \forall \theta \in \mathbb{R} \quad (24)$$

743 We now define the sequence M_t as follows, where $\lambda = K/8\xi^2$

$$M_t = \exp\left(\sum_{s=0}^{t-1} \lambda \Delta_s - \frac{\lambda}{4} \|h_{\mu_s|\mathcal{F}_s}\|_{\mathcal{H}}^2\right)$$

744 Since g_t is \mathcal{F}_{t+1} measurable, so is Δ_t , which implies M_t is \mathcal{F}_{t+1} measurable. Furthermore,

$$\begin{aligned} \mathbb{E}[M_t | \mathcal{F}_t] &= M_{t-1} e^{-\frac{\lambda}{4} \|h_{\mu_t|\mathcal{F}_t}\|_{\mathcal{H}}^2} \mathbb{E}[e^{\lambda \Delta_t} | \mathcal{F}_t] \\ &\leq M_{t-1} e^{(-\frac{\lambda}{4} + \frac{2\lambda^2 \xi^2}{K}) \|h_{\mu_t|\mathcal{F}_t}\|_{\mathcal{H}}^2} \leq M_{t-1} \end{aligned}$$

745 Thus, M_t is an \mathcal{F} -adapted supermartingale sequence. Following the same steps, we conclude
 746 $E[M_1] \leq 1$, which implies $\mathbb{E}[M_T] \leq \mathbb{E}[M_1] \leq 1$. Thus, from Markov's Inequality

$$\mathbb{P}\left[\sum_{t=0}^{T-1} \Delta_t - \frac{1}{4} \|h_{\mu_t|\mathcal{F}_t}\|_{\mathcal{H}}^2 > x\right] \leq e^{-\lambda x} \mathbb{E}[M_T] \leq e^{-\lambda x}$$

747 Hence, the following holds with probability at least $1 - \delta/2$.

$$\sum_{t=0}^{T-1} \Delta_t - \frac{1}{4} \|h_{\mu_t|\mathcal{F}_t}\|_{\mathcal{H}}^2 \leq \frac{8\xi^2}{K} \log(2/\delta) \quad (25)$$

748 Substituting (24) and (25) into (20) and taking a union bound, we conclude that the following holds
 749 with probability at least $1 - \delta$:

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} \text{KSD}_{\pi^*}(\mu_t | \mathcal{F}_t | \pi^*)^2 &\leq \frac{4\text{KL}(\mu_0 | \mathcal{F}_0 | \pi^*)}{\gamma T} + \frac{2\gamma(4+L)B\xi^2}{K} \\ &\quad + \frac{32\xi^2 \log(2/\delta)}{KT} + 12\gamma(4+L)B\xi^2 \sqrt{\frac{\log(2/\delta)}{T}} \end{aligned} \quad (26)$$

750 Recall from the proof of Theorem 1 in Appendix C.4 that the outputs $(\mathbf{y}^{(l)})_{l \in (n)}$ of VP-
 751 SVGD are i.i.d samples from the random measure $\bar{\mu}(\cdot; \mathbf{x}_0^{(0)}, \dots, \mathbf{x}_{KT-1}^{(0)}, S)$ when conditioned on
 752 $\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_{KT-1}^{(0)}, S$. Furthermore, when conditioned on $S = t$, $\bar{\mu}(\cdot; \mathbf{x}_0^{(0)}, \dots, \mathbf{x}_{KT-1}^{(0)}, S = t) =$
 753 $\mu_t | \mathcal{F}_t$. Thus, from (26), we conclude that, upon taking an expectation over $S \sim \text{Uniform}((T))$ while
 754 conditioning on the virtual particles $\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(KT-1)}$, the following holds with probability at least
 755 $1 - \delta$:

$$\begin{aligned} \mathbb{E}_S[\text{KSD}_{\pi^*}(\bar{\mu}(\cdot; \mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(KT-1)}, S) | \pi^*)^2] &\leq \frac{1}{T} \sum_{t=0}^{T-1} \text{KSD}_{\pi^*}(\mu_t | \mathcal{F}_t | \pi^*)^2 \\ &\leq \frac{4\text{KL}(\mu_0 | \mathcal{F}_0 | \pi^*)}{\gamma T} + \frac{2\gamma(4+L)B\xi^2}{K} \\ &\quad + \frac{32\xi^2 \log(2/\delta)}{KT} + 12\gamma(4+L)B\xi^2 \sqrt{\frac{\log(2/\delta)}{T}} \end{aligned}$$

756

□

757 D Analysis of GB-SVGD

758 In this section, we present our analysis of GB-SVGD. For any $t \in (T)$, we use \tilde{g}_t to denote the
 759 random function $\tilde{g}_t(\mathbf{x}) = \frac{1}{K} \sum_{r \in \mathcal{K}_t} h(\mathbf{x}, \mathbf{x}_t^{(r)})$ where \mathcal{K}_t is the random batch of size K sampled at
 760 time-step t of GB-SVGD.

761 In order to prove Theorem 2, we first establish an almost-sure iterate bound for GB-SVGD which is
 762 similar to that of Lemma 9 for VP-SVGD.

763 **Lemma 10** (Almost-Sure Iterate Bounds). *Let Assumptions 1, 2, 3 and 4 be satisfied. Then, the*
 764 *following holds almost surely for any $s \in \mathbb{N} \cup \{0\}$ and $t \in (T + 1)$ whenever $\gamma \leq 1/2A_1L$*

$$\begin{aligned} \|\mathbf{x}_t^{(s)}\| &\leq \zeta_0 + \zeta_1(\gamma T)^{1/\alpha} + \zeta_2(\gamma^2 T)^{1/\alpha} + \zeta_3 R^{2/\alpha} \\ \|h(\cdot, \mathbf{x}_t^{(s)})\|_{\mathcal{H}} &\leq \zeta_0 + \zeta_1(\gamma T)^{1/\alpha} + \zeta_2(\gamma^2 T)^{1/\alpha} + \zeta_3 R^{2/\alpha} \\ \|\tilde{g}_t\|_{\mathcal{H}} &\leq \zeta_0 + \zeta_1(\gamma T)^{1/\alpha} + \zeta_2(\gamma^2 T)^{1/\alpha} + \zeta_3 R^{2/\alpha} \end{aligned}$$

765 where ζ_0, \dots, ζ_3 are problem-dependent constants that depend polynomially on
 766 $A_1, A_2, A_3, B, d_1, d_2, L$ for any fixed α .

767 *Proof.* The proof of this Lemma is identical to that of Lemma 9. To this end, let $c_t^{(s)} =$
 768 $\frac{1}{K} \sum_{r \in \mathcal{K}_t} k(\mathbf{x}_t^{(s)}, \mathbf{x}_t^{(r)})$. Note that by Assumption 3, $c_t^{(s)} \geq 0$ Since $\mathbf{x}_{t+1}^{(s)} = \mathbf{x}_t^{(s)} - \gamma \tilde{g}_t(\mathbf{x}_t^{(s)})$,
 769 it follows from the smoothness of F that,

$$F(\mathbf{x}_{t+1}^{(s)}) - F(\mathbf{x}_t^{(s)}) \leq -\gamma \left\langle \nabla F(\mathbf{x}_t^{(s)}), \tilde{g}_t(\mathbf{x}_t^{(s)}) \right\rangle + \frac{\gamma^2 L}{2} \|\tilde{g}_t(\mathbf{x}_t^{(s)})\|^2 \quad (27)$$

770 By Lemma 6, we note that,

$$\begin{aligned} -\gamma \left\langle \nabla F(\mathbf{x}_t^{(s)}), \tilde{g}_t(\mathbf{x}_t^{(s)}) \right\rangle &= -\frac{\gamma}{K} \sum_{r \in \mathcal{K}_t} \left\langle \nabla F(\mathbf{x}_t^{(s)}), h(\mathbf{x}_t^{(s)}, \mathbf{x}_t^{(r)}) \right\rangle \\ &\leq \frac{\gamma}{K} \sum_{r \in \mathcal{K}_t} \left[-\frac{1}{2} k(\mathbf{x}_t^{(s)}, \mathbf{x}_t^{(r)}) \|\nabla F(\mathbf{x}_t^{(s)})\|^2 + L^2 A_1 + A_3 \right] \\ &\leq -\frac{\gamma c_t^{(s)}}{2} \|\nabla F(\mathbf{x}_t^{(s)})\|^2 + \gamma L^2 A_1 + \gamma A_3 \end{aligned} \quad (28)$$

771 Moreover, by Jensen's Inequality and Lemma 6

$$\begin{aligned} \|\tilde{g}_t(\mathbf{x}_t^{(s)})\|^2 &\leq \frac{1}{K} \sum_{r \in \mathcal{K}_t} \|h(\mathbf{x}_t^{(s)}, \mathbf{x}_t^{(r)})\|^2 \\ &\leq \frac{1}{K} \sum_{r \in \mathcal{K}_t} 2(A_1 L/2 + A_2)^2 + 2k(\mathbf{x}_t^{(s)}, \mathbf{x}_t^{(r)})^2 \|F(\mathbf{x}_t^{(s)})\|^2 \\ &\leq \frac{1}{K} \sum_{r \in \mathcal{K}_t} 2(A_1 L/2 + A_2)^2 + 2A_1 k(\mathbf{x}_t^{(s)}, \mathbf{x}_t^{(r)}) \|F(\mathbf{x}_t^{(s)})\|^2 \\ &\leq 2(A_1 L/2 + A_2)^2 + 2A_1 c_t^{(s)} \|F(\mathbf{x}_t^{(s)})\|^2 \end{aligned} \quad (29)$$

772 Substituting (28) and (29) into (27), we obtain,

$$\begin{aligned} F(\mathbf{x}_{t+1}^{(s)}) - F(\mathbf{x}_t^{(s)}) &\leq -\frac{\gamma c_t^{(s)}}{2} \|\nabla F(\mathbf{x}_t^{(s)})\|^2 + \gamma L^2 A_1 + \gamma A_3 \\ &\quad + \gamma^2 L(A_1 L/2 + A_2)^2 + \gamma^2 L A_1 c_t^{(s)} \|F(\mathbf{x}_t^{(s)})\|^2 \\ &\leq -\frac{\gamma c_t^{(s)}}{2} (1 - 2A_1 L \gamma) \|\nabla F(\mathbf{x}_t^{(s)})\|^2 + \gamma A_3 + \gamma L^2 A_1 + \gamma^2 L(A_1 L/2 + A_2)^2 \\ &\leq \gamma A_3 + \gamma L^2 A_1 + \gamma^2 L(A_1 L/2 + A_2)^2 \end{aligned}$$

773 where the last inequality uses the fact that $c_t^{(s)} \geq 0$ and $\gamma \leq 1/2A_1L$. Now, iterating through the above
 774 inequality, we obtain the following for any $t \in [T]$, $s \in \mathbb{N} \cup \{0\}$

$$F(\mathbf{x}_t^{(s)}) \leq F(\mathbf{x}_0^{(s)}) + \gamma T L^2 A_1 + \gamma T A_3 + \gamma^2 T L(A_1 L/2 + A_2)^2 \quad (30)$$

775 Furthermore, by Assumption 1

$$\begin{aligned} F(\mathbf{x}_0^{(s)}) &\leq F(0) + \|\nabla F(0)\| \|\mathbf{x}_0^{(s)}\| + \frac{L}{2} \|\mathbf{x}_0^{(s)}\|^2 \\ &\leq F(0) + 1/2 + L \|\mathbf{x}_0^{(s)}\|^2 \end{aligned}$$

776 Substituting the above inequality into (30), and using Assumption 2, we obtain the following for any
777 $t \in [T]$, $s \in \mathbb{N} \cup \{0\}$

$$\begin{aligned} d_1 \|\mathbf{x}_t^{(s)}\|^\alpha - d_2 \leq F(\mathbf{x}_t^s) &\leq F(0) + 1/2 + L \|\mathbf{x}_0^{(s)}\|^2 + \gamma T L^2 A_1 + \gamma T A_3 \\ &\quad + \gamma^2 T L (A_1 L/2 + A_2)^2 \end{aligned}$$

778 Rearranging and applying Assumption 4, we obtain

$$\begin{aligned} \|\mathbf{x}_t^{(s)}\| &\leq d_1^{-1/\alpha} [F(0) + 1/2 + L R^2 + \gamma T L^2 A_1 + \gamma T A_3 + \gamma^2 T L (A_1 L + A_2)^2]^{1/\alpha} \\ &\leq \tilde{\zeta}_0 + \tilde{\zeta}_1 (\gamma T)^{1/\alpha} + \tilde{\zeta}_2 (\gamma^2 T)^{1/\alpha} + \tilde{\zeta}_3 R^{2/\alpha} \end{aligned}$$

779 where $\tilde{\zeta}_0, \dots, \tilde{\zeta}_3$ are constants that depend polynomially on L, A_1, A_2, A_3, R . We note that, since
780 $0 < \alpha \leq 2$, the above inequality also holds for $t = 0$.

781 Using the above inequality Lemma 6 and Assumption 1, we conclude that the following holds almost
782 surely for any $t \in (T + 1)$, $s \in \mathbb{N} \cup \{0\}$

$$\begin{aligned} \|h(\cdot, \mathbf{x}_t^{(s)})\|_{\mathcal{H}} &\leq B L \|\mathbf{x}_t^{(s)}\| + B \sqrt{L} + B \\ &\leq \tilde{\eta}_0 + \tilde{\eta}_1 (\gamma T)^{1/\alpha} + \tilde{\eta}_2 (\gamma^2 T)^{1/\alpha} + \tilde{\eta}_3 R^{2/\alpha} \end{aligned}$$

783 where $\tilde{\eta}_0, \dots, \tilde{\eta}_3$ are constants that depend polynomially on L, B, A_1, A_2, A_3, R . Using the above
784 inequality, we conclude that the following also holds for any $t \in (T + 1)$.

$$\|\tilde{g}_t\|_{\mathcal{H}} \leq \tilde{\eta}_0 + \tilde{\eta}_1 (\gamma T)^{1/\alpha} + \tilde{\eta}_2 (\gamma^2 T)^{1/\alpha} + \tilde{\eta}_3 R^{2/\alpha}$$

785 Taking $\zeta_i = \max\{\tilde{\zeta}_i, \tilde{\eta}_i\}$, the proof is complete. \square

786 D.1 Proof of Theorem 2

787 *Proof.* Let $\xi = \zeta_0 + \zeta_1 (\gamma T)^{1/\alpha} + \zeta_2 (\gamma^2 T)^{1/\alpha} + \zeta_3 R^{2/\alpha}$ where ζ_0, \dots, ζ_3 are constants as described
788 in Lemma 9 and Lemma 10. Since the assumptions and parameter settings of Theorem 1 holds,
789 $\gamma \leq 1/2 A_1 L$ and thus, by Lemma 9 and Lemma 10, the particles output by VP-SVGD and GB-
790 SVGD are bounded as $\|\mathbf{y}^{(l)}\| \leq \xi$ and $\|\bar{\mathbf{y}}^{(l)}\| \leq \xi$.

791 Let $\mathbf{Y} = (\mathbf{y}^{(0)}, \dots, \mathbf{y}^{(n-1)})$ and $\bar{\mathbf{Y}} = (\bar{\mathbf{y}}^{(0)}, \dots, \bar{\mathbf{y}}^{(n-1)})$ denote the outputs of VP-SVGD and
792 GB-SVGD. Let $\hat{\mu}^{(n)} = \frac{1}{n} \sum_{i=0}^{n-1} \delta_{\mathbf{y}^{(i)}}$ and $\hat{\nu}^{(n)} = \frac{1}{n} \sum_{i=0}^{n-1} \delta_{\bar{\mathbf{y}}^{(i)}}$ be their respective empirical
793 distributions. We shall now explicitly construct a coupling between the inputs of VP-SVGD and
794 GB-SVGD such that the first $n - KT$ particles of their respective outputs are equal. This in turn will
795 allow us to control the expected squared KSD between $\hat{\mu}^{(n)}$ and $\hat{\nu}^{(n)}$.

796 To this end, let \mathcal{E} denote the event that each random batch \mathcal{K}_t of GB-SVGD is disjoint and contains
797 unique elements for every $t \in (T)$. Subsequently, let \mathcal{K} denote the set of all indices that were chosen
798 to be part of some random batch \mathcal{K}_t . Let Λ be a uniformly random permutation over $\{0, \dots, n - 1\}$.
799 We note that, conditioned on \mathcal{E} , the distribution of the random set \mathcal{K} is the same as the distribution of
800 $\{\Lambda(0), \dots, \Lambda(KT - 1)\}$. We can couple a uniformly random permutation Λ and \mathcal{K}_t for $0 \leq t \leq T$
801 such that under the event \mathcal{E} , $\mathcal{K} = \{\Lambda(0), \dots, \Lambda(KT - 1)\}$ and $\{\Lambda(tK), \dots, \Lambda((t+1)K - 1)\}$ is the
802 random batch \mathcal{K}_t . Thus, under the event \mathcal{E} , one can couple a uniformly random permutation Λ and
803 \mathcal{K}_t for $t \in (T)$ such that $\mathcal{K} = \{\Lambda(0), \dots, \Lambda(KT - 1)\}$ and $\mathcal{K}_t = \{\Lambda(tK), \dots, \Lambda((t+1)K - 1)\}$

804 With this insight, we couple VP-SVGD and GB-SVGD as follows. We note that, the random batch \mathcal{K}_t
805 in GB-SVGD is sampled independently of the initial particles. To this end, let $\bar{\mathbf{x}}_0^{(0)}, \dots, \bar{\mathbf{x}}_0^{(n-1)}$ *i.i.d.*
806 μ_0 , and let the random batches \mathcal{K}_t and permutation Λ be jointly distributed as described above,
807 independently of $\bar{\mathbf{x}}_0^{(0)}, \dots, \bar{\mathbf{x}}_0^{(n-1)}$, i.e.

$$\Lambda \sim \text{Uniform}(\mathbb{S}_{(n)}), \quad \mathcal{K}_t = \{\Lambda(tK), \dots, \Lambda((t+1)K - 1)\}, \quad t \in (T)$$

808 We now define $\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(KT+n-1)}$ as:

$$\mathbf{x}_0^{(l)} := \begin{cases} = \bar{\mathbf{x}}_0^{(\Lambda(l))} & \text{for } 0 \leq l \leq n-1 \\ \sim \mu_0 \text{ independent of everything else} & \text{for } n \leq l \leq KT+n-1 \end{cases} \quad (31)$$

809 Let $\bar{\mathbf{x}}_0^{(0)}, \dots, \bar{\mathbf{x}}_0^{(n-1)}$ and \mathcal{K}_t as the initialization and random batches for GB-SVGD, and let
 810 $\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(KT+n-1)}$ be the initialization for VP-SVGD. We first show that this construction is
 811 indeed a valid coupling between VP-SVGD and GB-SVGD.

812 **Claim 1.** *Conditioned on \mathcal{E} , the inputs to VP-SVGD and GB-SVGD, as constructed above is a valid*
 813 *coupling, i.e., the marginal distribution of $\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(KT+n-1)}$ is equal to the distribution of initial*
 814 *particles in VP-SVGD, and the marginal distribution of $\bar{\mathbf{x}}_0^{(0)}, \dots, \bar{\mathbf{x}}_0^{(n-1)}, (\mathcal{K}_t)_{t \in (T)}$ is the same as*
 815 *the distribution of initial particles and random batches in \mathcal{K}_t*

816 *Proof.* By construction $\bar{\mathbf{x}}_0^{(0)}, \dots, \bar{\mathbf{x}}_0^{(n-1)} \stackrel{i.i.d.}{\sim} \mu_0$. Moreover, conditioned on \mathcal{E} , the distribution of
 817 $\mathcal{K}_t = \{\Lambda(tK), \dots, \Lambda((t+1)K-1)\}$, has the distribution of a uniform random batch of size K
 818 since $\Lambda \sim \text{Uniform}(\mathbb{S}_n)$. Furthermore, since Λ is sampled independently of $\bar{\mathbf{x}}_0^{(0)}, \dots, \bar{\mathbf{x}}_0^{(n-1)}$, \mathcal{K}_t
 819 is independent of $\bar{\mathbf{x}}_0^{(0)}, \dots, \bar{\mathbf{x}}_0^{(n-1)}$ for any $t \in (T)$. Thus, the coupling constructed above has the
 820 correct marginal with respect to GB-SVGD.

821 To establish the same for VP-SVGD, we note that by (31), $\mathbf{x}_0^{(n)}, \dots, \mathbf{x}_0^{(KT+n-1)} \stackrel{i.i.d.}{\sim} \mu_0$, *sam-*
 822 *pled independently of everything else.* Moreover, since $\bar{\mathbf{x}}_0^{(0)}, \dots, \bar{\mathbf{x}}_0^{(n-1)} \stackrel{i.i.d.}{\sim} \mu_0$, we infer that
 823 $\bar{\mathbf{x}}_0^{(\Lambda(0))}, \dots, \bar{\mathbf{x}}_0^{(\Lambda(n-1))} \stackrel{i.i.d.}{\sim} \mu_0$ for any arbitrary permutation $\Lambda \in \mathbb{S}_n$. From this, and (31), we
 824 conclude that $\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(KT+n-1)} \stackrel{i.i.d.}{\sim} \mu_0$. Hence, the coupling constructed above has the correct
 825 marginal with respect to VP-SVGD. \square

826 We now show that, under the constructed coupling, the time-evolution of the particles of VP-
 827 SVGD and GB-SVGD satisfy $\bar{\mathbf{x}}_t^{(\Lambda(l))} = \mathbf{x}_t^{(l)}$, $KT \leq l \leq n-1, t \in (T+1)$, when conditioned on
 828 the event \mathcal{E} .

829 **Claim 2.** *Let the inputs to VP-SVGD and GB-SVGD be coupled as per the construction above. Then,*
 830 *conditioned on the event \mathcal{E} , the particles $\mathbf{x}_t^{(s)}$ and $\bar{\mathbf{x}}_t^{(s)}$ of VP-SVGD and GB-SVGD respectively,*
 831 *satisfy $\bar{\mathbf{x}}_t^{(\Lambda(l))} = \mathbf{x}_t^{(l)}$ for every $KT \leq l \leq n-1$ and $0 \leq t \leq T$*

832 *Proof.* We prove this by an inductive argument. Clearly, the claim holds for $t = 0$ by the construction
 833 of our coupling. Assume it holds for some arbitrary $t \in (T)$. Now, writing the update equation for
 834 GB-SVGD for $KT \leq l \leq n-1$,

$$\begin{aligned} \bar{\mathbf{x}}_{t+1}^{(\Lambda(l))} &= \bar{\mathbf{x}}_t^{(\Lambda(l))} - \frac{\gamma}{K} \sum_{r \in \mathcal{K}_t} h(\bar{\mathbf{x}}_t^{(\Lambda(l))}, \bar{\mathbf{x}}_t^{(r)}) \\ &= \bar{\mathbf{x}}_t^{(\Lambda(l))} - \frac{\gamma}{K} \sum_{l=0}^{K-1} h(\bar{\mathbf{x}}_t^{(\Lambda(l))}, \bar{\mathbf{x}}_t^{(\Lambda(Kt+l))}) \\ &= \mathbf{x}_t^{(l)} - \frac{\gamma}{K} \sum_{l=0}^{K-1} h(\mathbf{x}_t^{(l)}, \mathbf{x}_t^{(Kt+l)}) = \mathbf{x}_t^{(l+1)} \end{aligned}$$

835 where the second equality uses the fact that $\mathcal{K}_t = \{\Lambda(tK), \dots, \Lambda((t+1)K-1)\}$ when conditioned
 836 on \mathcal{E} and the third equality uses the induction hypothesis $\bar{\mathbf{x}}_t^{(\Lambda(l))} = \mathbf{x}_t^{(l)}$ for $KT \leq l \leq n-1$. Hence,
 837 the claim is proven true by induction. \square

838 Equipped with the above coupling between the inputs of VP-SVGD and GB-SVGD, one can now
 839 couple their outputs by sampling an $S \sim \text{Uniform}((n))$ and using this sampled S as the random
 840 timestep chosen by both VP-SVGD (Step 6 in Algorithm 1) and GB-SVGD (Step 7 in Algorithm 2)
 841 that are run with the coupled input constructed above. It is easy to see that this results in a coupling

842 of the outputs \mathbf{Y} and $\bar{\mathbf{Y}}$ of VP-SVGD and GB-SVGD respectively. Furthermore, by Claim 2, we
 843 note that, conditioned on the event \mathcal{E} , $\mathbf{y}^{(l-TK)} = \bar{\mathbf{y}}^{(\Lambda(l))}$ for every $KT \leq l \leq n-1$. We now define
 844 the permutation $\tau \in \mathbb{S}_{(n)}$ as follows:

$$\tau(\Lambda(l)) = \begin{cases} l+n-KT & \text{for } 0 \leq l \leq KT-1 \\ l-KT & \text{for } KT \leq l \leq n-1 \end{cases} \quad (32)$$

845 It follows that $\bar{\mathbf{y}}^{\tau(l)} = \mathbf{y}^{(l)}$ for $KT \leq l \leq n-1$. Thus, by definition of Kernel Stein Discrepancy
 846 (Definition 1), we can infer that the following holds when conditioned on the event \mathcal{E}

$$\begin{aligned} \mathbb{E}[\text{KSD}_{\pi^*}^2(\hat{\nu}^{(n)} || \hat{\mu}^{(n)}) | \mathcal{E}] &= \mathbb{E}[\|h_{\hat{\nu}^{(n)}} - h_{\hat{\mu}^{(n)}}\|_{\mathcal{H}}^2 | \mathcal{E}] \\ &= \mathbb{E}\left[\left\|\frac{1}{n} \sum_{l=0}^{n-1} h(\cdot, \bar{\mathbf{y}}^{(l)}) - \frac{1}{n} \sum_{l=0}^{n-1} h(\cdot, \mathbf{y}^{(l)})\right\|_{\mathcal{H}}^2 | \mathcal{E}\right] \\ &= \frac{1}{n^2} \mathbb{E}\left[\left\|\sum_{l=0}^{n-1} h(\cdot, \bar{\mathbf{y}}^{(\tau(l))}) - h(\cdot, \mathbf{y}^{(l)})\right\|_{\mathcal{H}}^2 | \mathcal{E}\right] \\ &= \frac{1}{n^2} \mathbb{E}\left[\left\|\sum_{l=0}^{KT-1} h(\cdot, \bar{\mathbf{y}}^{(\tau(l))}) - h(\cdot, \mathbf{y}^{(l)})\right\|_{\mathcal{H}}^2 | \mathcal{E}\right] \\ &\leq \frac{KT}{n^2} \sum_{l=0}^{KT-1} \mathbb{E}\left[\|h(\cdot, \bar{\mathbf{y}}^{(\tau(l))}) - h(\cdot, \mathbf{y}^{(l)})\|_{\mathcal{H}}^2 | \mathcal{E}\right] \\ &\leq \frac{2K^2T^2\xi^2}{n^2} \end{aligned} \quad (33)$$

847 where the second step uses the permutation invariance of summation, the third step uses the fact that
 848 $\bar{\mathbf{y}}^{\tau(l)} = \bar{\mathbf{y}}^{(l)}$ for $KT \leq l \leq n-1$, the fourth step uses the convexity of $\|\cdot\|_{\mathcal{H}}^2$ and the last step uses
 849 the almost-sure iterate bounds of Lemma 9 and 10

850 Under the event \mathcal{E}^c , we directly apply the almost-sure iterate bounds of Lemma 9 and 10 to obtain
 851 the following:

$$\begin{aligned} \mathbb{E}[\text{KSD}_{\pi^*}^2(\hat{\nu}^{(n)} || \hat{\mu}^{(n)}) | \mathcal{E}^c] &= \mathbb{E}[\|h_{\hat{\nu}^{(n)}} - h_{\hat{\mu}^{(n)}}\|_{\mathcal{H}}^2 | \mathcal{E}^c] \\ &= \frac{1}{n^2} \mathbb{E}\left[\left\|\sum_{l=0}^{n-1} h(\cdot, \bar{\mathbf{y}}^{(l)}) - h(\cdot, \mathbf{y}^{(l)})\right\|_{\mathcal{H}}^2 | \mathcal{E}^c\right] \\ &\leq 2\xi^2 \end{aligned} \quad (34)$$

852 From Equations (33) and (34), it follows that:

$$\begin{aligned} \mathbb{E}[\text{KSD}_{\pi^*}^2(\hat{\nu}^{(n)} || \hat{\mu}^{(n)})] &= \mathbb{E}[\text{KSD}_{\pi^*}^2(\hat{\nu}^{(n)} || \hat{\mu}^{(n)}) | \mathcal{E}] \mathbb{P}(\mathcal{E}) + \mathbb{E}[\text{KSD}_{\pi^*}^2(\hat{\nu}^{(n)} || \hat{\mu}^{(n)}) | \mathcal{E}^c] \mathbb{P}(\mathcal{E}^c) \\ &\leq \frac{2K^2T^2\xi^2}{n^2} \mathbb{P}(\mathcal{E}) + 2\xi^2 \mathbb{P}(\mathcal{E}^c) \end{aligned}$$

853 Recall that $P(\mathcal{E}) = 1$ under sampling without replacement and $P(\mathcal{E}) = 1 - \frac{K^2T^2}{n}$ under sampling
 854 with replacement. Thus, we conclude that the following holds under the constructed coupling of \mathbf{Y}
 855 and $\bar{\mathbf{Y}}$

$$\mathbb{E}[\text{KSD}_{\pi^*}^2(\hat{\nu}^{(n)} || \hat{\mu}^{(n)})] \leq \begin{cases} \frac{2K^2T^2\xi^2}{n^2} & \text{(without replacement sampling)} \\ \frac{2K^2T^2\xi^2}{n^2} \left(1 - \frac{K^2T^2}{n}\right) + \frac{2K^2T^2\xi^2}{n} & \text{(with replacement sampling)} \end{cases}$$

856 □

857 E Finite-Particle Convergence Guarantees for VP-SVGD and GB-SVGD

858 In this section, we show that the empirical measure of the particles output by VP-SVGD and GB-
 859 SVGD rapidly converge to the target distribution π^* in KSD. To this end, we prove the finite-particle

860 convergence rates for VP-SVGD in Appendix E.1 and that of GB-SVGD in Appendix E.2. Finally, we
 861 compare the oracle complexity (i.e., the number of evaluations of ∇F) of VP-SVGD and GB-SVGD
 862 to that of SVGD in Appendix E.3

863 E.1 VP-SVGD

864 **Corollary 2 (VP-SVGD : Fast Finite-Particle Convergence).** *Let the assumptions and parameter*
 865 *settings of Theorem 1 be satisfied. Let $\hat{\mu}^{(n)}$ denote the empirical measure of the n particles output by*
 866 *VP-SVGD.*

$$\mathbb{E}[\text{KSD}_{\pi^*}^2(\hat{\mu}^{(n)}|\pi^*)] \leq \frac{\xi^2}{n} + \frac{2\text{KL}(\mu_0|\mathcal{F}_0|\pi^*)}{\gamma T} + \frac{\gamma B(4+L)\xi^2}{K}$$

867 where ξ is as defined in Theorem 1. Setting $R = \sqrt{d/L}, \gamma = O(\frac{(Kd)^\eta}{T^{1-\eta}})$ with $\eta = \frac{\alpha}{2(1+\alpha)}$ and

868 $KT = d^{2+\alpha} n^{\frac{2(1+\alpha)}{2+\alpha}}$ suffices to ensure,

$$\mathbb{E}[\text{KSD}_{\pi^*}^2(\hat{\mu}^{(n)}|\pi^*)] \leq O\left(\frac{d^{\frac{2}{2+\alpha}}}{n^{\frac{2+\alpha}{2+\alpha}}} + \frac{d^{2/\alpha}}{n}\right)$$

869 *Proof.* Recall from Algorithm 1 that the outputs of VP-SVGD are $\mathbf{x}_S^{(KT)}, \dots, \mathbf{x}_S^{(KT+n-1)}$ where $S \sim$
 870 $\text{Uniform}(\{0, \dots, T-1\})$. Hence, their empirical measure $\hat{\mu}^{(n)}$ is given by $\hat{\mu}^{(n)} = \frac{1}{n} \sum_{l=0}^{n-1} \delta_{\mathbf{x}_S^{(KT+l)}}$.
 871 From the definition of the Kernel Stein Discrepancy (Definition 1), it follows that,

$$\text{KSD}_{\pi^*}^2(\hat{\mu}^{(n)}|\pi^*) = \|h_{\hat{\mu}^{(n)}}\|_{\mathcal{H}}^2 = \left\| \frac{1}{n} \sum_{l=1}^N h(\cdot, \mathbf{x}_S^{(KT+l)}) \right\|_{\mathcal{H}}^2 \quad (35)$$

872 For the sake of clarity, only in this proof, we use \mathcal{C} to denote the conditioning on the virtual particles
 873 $\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(KT-1)}$. Now, consider any arbitrary $t \in \{0, \dots, T-1\}$. Taking conditional expectations
 874 on both sides of Equation (35) by conditioning on \mathcal{C} and the event $\{S = t\}$, we obtain the following:

$$\begin{aligned} \mathbb{E} \left[\text{KSD}_{\pi^*}^2(\hat{\mu}^{(n)}|\pi^*) \mid \mathcal{C}, S = t \right] &= \mathbb{E} \left[\left\| \frac{1}{n} \sum_{l=0}^{n-1} h(\cdot, \mathbf{x}_S^{(KT+l)}) \right\|_{\mathcal{H}}^2 \mid \mathcal{C}, S = t \right] \\ &= \mathbb{E} \left[\left\| \frac{1}{n} \sum_{l=0}^{n-1} h(\cdot, \mathbf{x}_t^{(KT+l)}) \right\|_{\mathcal{H}}^2 \mid (\mathbf{x}_0^{(s)})_{0 \leq s \leq KT-1} \right] \end{aligned} \quad (36)$$

875 Recall from Equation (14) in Appendix C.1 that for any $l \in \{0, \dots, n-1\}$ $\mathbf{x}_t^{(KT+l)}$ depends only
 876 on $\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(Kt-1)}$ and $\mathbf{x}_0^{(KT+l)}$. Furthermore, from Appendix C.1, we recall that the filtration \mathcal{F}_t
 877 is defined as $\mathcal{F}_t = \sigma(\{\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(Kt-1)}\})$. It follows that,

$$\begin{aligned} \mathbb{E} \left[\left\| \frac{1}{n} \sum_{l=0}^{n-1} h(\cdot, \mathbf{x}_t^{(KT+l)}) \right\|_{\mathcal{H}}^2 \mid (\mathbf{x}_0^{(s)})_{0 \leq s \leq KT-1} \right] &= \mathbb{E} \left[\left\| \frac{1}{n} \sum_{l=1}^N h(\cdot, \mathbf{x}_t^{(KT+l)}) \right\|_{\mathcal{H}}^2 \mid (\mathbf{x}_0^{(s)})_{0 \leq s \leq Kt-1} \right] \\ &= \mathbb{E} \left[\left\| \frac{1}{n} \sum_{l=0}^{n-1} h(\cdot, \mathbf{x}_t^{(KT+l)}) \right\|_{\mathcal{H}}^2 \mid \mathcal{F}_t \right] \end{aligned} \quad (37)$$

878 To control $\mathbb{E} \left[\left\| \frac{1}{n} \sum_{l=0}^{n-1} h(\cdot, \mathbf{x}_t^{(KT+l)}) \right\|_{\mathcal{H}}^2 \mid \mathcal{F}_t \right]$, we apply the arguments used in the proof of Lemma
 879 3. To this end, note that when conditioned on the virtual particles $\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(Kt-1)}$, the particles
 880 $\mathbf{x}_t^{(KT)}, \dots, \mathbf{x}_t^{(KT+n-1)}$ *i.i.d.* $\mu_t|\mathcal{F}_t$. Furthermore, since $\gamma \leq 1/2A_1L$ (as per the parameter settings
 881 of Theorem 1), $\|h(\cdot, \mathbf{x}_t^{(KT+l)})\|_{\mathcal{H}} \leq \xi \forall l \in (n)$ by Lemma 6. Finally, $\mathbb{E}[h(\mathbf{x}, \mathbf{x}_t^{(KT+l)})|\mathcal{F}_t] =$
 882 $h_{\mu_t|\mathcal{F}_t}(\mathbf{x}) \forall l \in (n), \mathbf{x} \in \mathbb{R}^d$. Hence, from Lemma 8, we conclude that $h_{\mu_t|\mathcal{F}_t}$ is the Gelfand-Pettis
 883 integral of the map $\mathbf{x} \rightarrow h(\mathbf{x}, \mathbf{x}_t^{(KT+l)})$ with respect to the measure $\mu_t|\mathcal{F}_t$, i.e.,

$$\mathbb{E}[\langle h(\cdot, \mathbf{x}_t^{(KT+l)}), f \rangle_{\mathcal{H}} \mid \mathcal{F}_t] = \langle h_{\mu_t|\mathcal{F}_t}, f \rangle \quad \forall f \in \mathcal{H} \quad (38)$$

884 To control $\mathbb{E} \left[\left\| \frac{1}{n} \sum_{l=0}^{n-1} h(\cdot, \mathbf{x}_t^{(KT+l)}) \right\|_{\mathcal{H}}^2 \mid \mathcal{F}_t \right]$, we proceed as follows:

$$\begin{aligned} \left\| \frac{1}{n} \sum_{l=0}^{n-1} h(\cdot, \mathbf{x}_t^{(KT+l)}) \right\|_{\mathcal{H}}^2 &= \frac{1}{n^2} \sum_{l_1, l_2=0}^{n-1} \left\langle h(\cdot, \mathbf{x}_t^{(KT+l_1)}), h(\cdot, \mathbf{x}_t^{(KT+l_2)}) \right\rangle_{\mathcal{H}} \\ &= \frac{1}{n^2} \sum_{l=0}^{n-1} \|h(\cdot, \mathbf{x}_t^{(KT+l)})\|_{\mathcal{H}}^2 + \frac{1}{n^2} \sum_{0 \leq l_1 \neq l_2 \leq n-1} \left\langle h(\cdot, \mathbf{x}_t^{(KT+l_1)}), h(\cdot, \mathbf{x}_t^{(KT+l_2)}) \right\rangle_{\mathcal{H}} \\ &\leq \frac{\xi^2}{n} + \frac{1}{n^2} \sum_{0 \leq l_1 \neq l_2 \leq n-1} \left\langle h(\cdot, \mathbf{x}_t^{(KT+l_1)}), h(\cdot, \mathbf{x}_t^{(KT+l_2)}) \right\rangle_{\mathcal{H}} \end{aligned}$$

885 where the last inequality uses the fact that $\|h(\cdot, \mathbf{x}_t^{(KT+l)})\|_{\mathcal{H}} \leq \xi$ almost surely as per Lemma 9.

886 To control the conditional expectation of the off-diagonal terms, let $i = KT + l_1$ and $j = KT + l_2$
887 for any arbitrary l_1, l_2 with $0 \leq l_1 \neq l_2 \leq n - 1$. Conditioned on \mathcal{F}_t , $\mathbf{x}_t^{(i)}$ and $\mathbf{x}_t^{(j)}$ are i.i.d samples
888 from $\mu_t \mid \mathcal{F}_t$. Thus, by Equation (38) and Fubini's Theorem,

$$\begin{aligned} \mathbb{E}_{\mathbf{x}_t^{(i)}, \mathbf{x}_t^{(j)}} \left[\left\langle h(\cdot, \mathbf{x}_t^{(i)}), h(\cdot, \mathbf{x}_t^{(j)}) \right\rangle_{\mathcal{H}} \mid \mathcal{F}_t \right] &= \mathbb{E}_{\mathbf{x}_t^{(i)}} \left[\mathbb{E}_{\mathbf{x}_t^{(j)}} \left[\left\langle h(\cdot, \mathbf{x}_t^{(i)}), h(\cdot, \mathbf{x}_t^{(j)}) \right\rangle_{\mathcal{H}} \mid \mathcal{F}_t \right] \right] \\ &= \mathbb{E}_{\mathbf{x}_t^{(i)}} \left[\left\langle h_{\mu_t \mid \mathcal{F}_t}, h(\cdot, \mathbf{x}_t^{(i)}) \right\rangle_{\mathcal{H}} \mid \mathcal{F}_t \right] \\ &= \|h_{\mu_t \mid \mathcal{F}_t}\|_{\mathcal{H}}^2 \end{aligned}$$

889 It follows that,

$$\mathbb{E} \left[\left\| \frac{1}{n} \sum_{l=0}^{n-1} h(\cdot, \mathbf{x}_t^{(KT+l)}) \right\|_{\mathcal{H}}^2 \mid \mathcal{F}_t \right] \leq \|h_{\mu_t \mid \mathcal{F}_t}\|_{\mathcal{H}}^2 + \frac{\xi^2}{n}$$

890 Substituting the above into equation 36 and equation 37, we obtain the following:

$$\mathbb{E} \left[\text{KSD}_{\pi^*}^2(\hat{\mu}^{(n)} \parallel \pi^*) \mid \mathcal{C}, S = t \right] \leq \frac{\xi^2}{n} + \|h_{\mu_t \mid \mathcal{F}_t}\|_{\mathcal{H}}^2 = \frac{\xi^2}{n} + \text{KSD}_{\pi^*}^2(\mu_t \mid \mathcal{F}_t \parallel \pi^*)$$

891 where the second step applies Definition 1. Finally, taking expectations with respect to \mathcal{C} and
892 $S \sim \text{Uniform}(\{0, \dots, T-1\})$ on both sides of the above inequality, we get:

$$\mathbb{E} \left[\text{KSD}_{\pi^*}^2(\hat{\mu}^{(n)} \parallel \pi^*) \right] \leq \frac{\xi^2}{n} + \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\text{KSD}_{\pi^*}^2(\mu_t \mid \mathcal{F}_t \parallel \pi^*)]$$

893 Substituting the bound from Theorem 1 into the above inequality, we conclude that:

$$\mathbb{E}[\text{KSD}_{\pi^*}^2(\hat{\mu}^{(n)} \parallel \pi^*)] \leq \frac{\xi^2}{n} + \frac{2\text{KL}(\mu_0 \mid \mathcal{F}_0 \parallel \pi^*)}{\gamma T} + \frac{\gamma B(4+L)\xi^2}{K}$$

894 We note that for $\gamma = O\left(\frac{(Kd)^\eta}{T^{1-\eta}}\right)$ and $R = \sqrt{d/L}$, $\text{KL}(\mu_0 \mid \mathcal{F}_0 \parallel \pi^*) = O(d)$ by Lemma 4 and

$$\xi^2 \leq 4\zeta_0 + 4\zeta_1(\gamma T)^{2/\alpha} + 4\zeta_2(\gamma^2 T)^{2/\alpha} + 4\zeta_3 R^{4/\alpha} \leq O\left((KdT)^{\frac{1}{1+\alpha}} + d^{2/\alpha}\right)$$

895 Furthermore,

$$\begin{aligned} \frac{2\text{KL}(\mu_0 \mid \mathcal{F}_0 \parallel \pi^*)}{\gamma T} + \frac{\gamma B(4+L)\xi^2}{K} &\leq O\left(\frac{d}{\gamma T} + \frac{\gamma B(4+L)\xi^2}{2K}\right) \leq O\left(\frac{d^{1-\eta}}{(KT)^\eta}\right) \\ &\leq O\left(\frac{d^{\frac{2+\alpha}{2(1+\alpha)}}}{(KT)^{\frac{\alpha}{2(1+\alpha)}}}\right) \end{aligned}$$

896 It follows that,

$$\mathbb{E}[\text{KSD}_{\pi^*}^2(\hat{\mu}^{(n)} \parallel \pi^*)] \leq O\left(\frac{d^{2/\alpha}}{n} + \frac{(KTd)^{\frac{1}{1+\alpha}}}{n} + \frac{d^{\frac{2+\alpha}{2(1+\alpha)}}}{(KT)^{\frac{\alpha}{2(1+\alpha)}}}\right)$$

897 $KT = d^{\frac{\alpha}{2+\alpha}} n^{\frac{2(1+\alpha)}{2+\alpha}}$, we conclude:

$$\mathbb{E}[\text{KSD}_{\pi^*}^2(\hat{\mu}^{(n)}|\pi^*)] \leq O\left(\frac{d^{\frac{2}{2+\alpha}}}{n^{\frac{\alpha}{2+\alpha}}} + \frac{d^{2/\alpha}}{n}\right)$$

898 □

899 E.2 GB-SVGD

900 **Corollary 3 (GB-SVGD : Fast Finite-Particle Convergence).** *Let the assumptions and parameter*
 901 *settings of Theorem 1 be satisfied. Let $\hat{\nu}^{(n)}$ denote the empirical measure of the n particles output by*
 902 *GB-SVGD. Then, under without-replacement sampling of the minibatches, the following holds:*

$$\mathbb{E}[\text{KSD}_{\pi^*}^2(\hat{\nu}^{(n)}|\pi^*)] \leq \frac{4K^2T^2\xi^2}{n^2} + \frac{2\xi^2}{n} + \frac{4\text{KL}(\mu_0|\mathcal{F}_0|\pi^*)}{\gamma T} + \frac{2\gamma B(4+L)\xi^2}{K}$$

903 *and the following holds under with-replacement sampling of the minibatches*

$$\mathbb{E}[\text{KSD}_{\pi^*}^2(\hat{\nu}^{(n)}|\pi^*)] \leq \frac{4K^2T^2\xi^2}{n^2}\left(1 - \frac{K^2T^2}{n}\right) + \frac{4K^2T^2\xi^2}{n} + \frac{2\xi^2}{n} + \frac{4\text{KL}(\mu_0|\mathcal{F}_0|\pi^*)}{\gamma T} + \frac{2\gamma B(4+L)\xi^2}{K}$$

904 *where ξ is as defined in Theorem 1. In particular, for GB-SVGD under without-replacement sampling*
 905 *of the minibatches, setting $R = \sqrt{d/L}$, $\gamma = O\left(\frac{(Kd)^\eta}{T^{1-\eta}}\right)$ with $\eta = \frac{\alpha}{2(1+\alpha)}$ and $KT = \sqrt{n}$ suffices to*
 906 *ensure the following*

$$\mathbb{E}[\text{KSD}_{\pi^*}^2(\hat{\nu}^{(n)}|\pi^*)] \leq O\left(\frac{d^{2/\alpha}}{n} + \frac{d^{\frac{1}{1+\alpha}}}{n^{\frac{1+2\alpha}{2(1+\alpha)}}} + \frac{d^{\frac{2+\alpha}{2(1+\alpha)}}}{n^{\frac{\alpha}{4(1+\alpha)}}}\right)$$

907 *Proof.* Let $\bar{\mathbf{Y}} = (\bar{\mathbf{y}}^{(0)}, \dots, \bar{\mathbf{y}}^{(n-1)})$ denote the n particles output by GB-SVGD and let $\hat{\nu}^{(n)} =$
 908 $\frac{1}{n} \sum_{l=0}^{n-1} \delta_{\bar{\mathbf{y}}^{(l)}}$ denote their empirical measure. Let \mathcal{E} denote the event that each random batch \mathcal{K}_t
 909 of GB-SVGD is disjoint and contains unique elements for every $t \in (T)$. Moreover, let $\mathbf{Y} =$
 910 $(\mathbf{y}^{(0)}, \dots, \mathbf{y}^{(n-1)})$ denote the n particles output by VP-SVGD, run with the parameter settings stated
 911 above, and coupled with $\bar{\mathbf{Y}}$ as per the coupling constructed in the proof of Theorem 2 in Appendix D.1.
 912 Let $\hat{\mu}^{(n)} = \frac{1}{n} \sum_{l=0}^{n-1} \delta_{\mathbf{y}^{(l)}}$ denote their empirical measure. By definition of Kernel Stein Discrepancy
 913 (Definition 1) and the convexity $\|\cdot\|_{\mathcal{H}}^2$, it follows that:

$$\begin{aligned} \mathbb{E}[\text{KSD}_{\pi^*}^2(\hat{\nu}^{(n)}|\pi^*)] &= \mathbb{E}[\|h_{\hat{\nu}^{(n)}}\|_{\mathcal{H}}^2] \\ &= \mathbb{E}[\|h_{\hat{\nu}^{(n)}} - h_{\hat{\mu}^{(n)}} + h_{\hat{\mu}^{(n)}}\|_{\mathcal{H}}^2] \\ &\leq 2\mathbb{E}[\|h_{\hat{\nu}^{(n)}} - h_{\hat{\mu}^{(n)}}\|_{\mathcal{H}}^2] + 2\mathbb{E}[\|h_{\hat{\mu}^{(n)}}\|_{\mathcal{H}}^2] \\ &= 2\mathbb{E}[\text{KSD}_{\pi^*}^2(\hat{\nu}^{(n)}|\hat{\mu}^{(n)})] + 2\mathbb{E}[\text{KSD}_{\pi^*}^2(\hat{\mu}^{(n)}|\pi^*)] \end{aligned}$$

914 Substituting the bounds of Theorem 2 and Corollary 2 into the above inequality, we conclude the
 915 following:

$$\mathbb{E}[\text{KSD}_{\pi^*}^2(\hat{\nu}^{(n)}|\pi^*)] \leq \frac{4K^2T^2\xi^2}{n^2}\mathbb{P}(\mathcal{E}) + 4\xi^2\mathbb{P}(\mathcal{E}^c) + \frac{2\xi^2}{n} + \frac{4\text{KL}(\mu_0|\mathcal{F}_0|\pi^*)}{\gamma T} + \frac{2\gamma B(4+L)\xi^2}{K}$$

916 We recall that, $\mathbb{P}(\mathcal{E}) = 1$ under without-replacement sampling of the random batches \mathcal{K}_t and
 917 $\mathbb{P}(\mathcal{E}) = 1 - K^2T^2/n$ under with-replacement sampling. Thus, under without-replacement sampling,
 918 the following holds:

$$\mathbb{E}[\text{KSD}_{\pi^*}^2(\hat{\nu}^{(n)}|\pi^*)] \leq \frac{4K^2T^2\xi^2}{n^2} + \frac{2\xi^2}{n} + \frac{4\text{KL}(\mu_0|\mathcal{F}_0|\pi^*)}{\gamma T} + \frac{2\gamma B(4+L)\xi^2}{K}$$

919 Moreover, the following holds under with-replacement sampling

$$\mathbb{E}[\text{KSD}_{\pi^*}^2(\hat{\nu}^{(n)}|\pi^*)] \leq \frac{4K^2T^2\xi^2}{n^2}\left(1 - \frac{K^2T^2}{n}\right) + \frac{4K^2T^2\xi^2}{n} + \frac{2\xi^2}{n} + \frac{4\text{KL}(\mu_0|\mathcal{F}_0|\pi^*)}{\gamma T} + \frac{2\gamma B(4+L)\xi^2}{K}$$

920 Now, let us consider GB-SVGD without replacement with $R = \sqrt{d/L}$, $\gamma = O(\frac{(Kd)^\eta}{T^{1-\eta}})$ and $KT =$
 921 $n^{1/2}$. It follows that $\text{KL}(\mu_0|\mathcal{F}_0||\pi^*) = O(d)$ by Lemma 4 and

$$\begin{aligned}\xi^2 &\leq 4\zeta_0 + 4\zeta_1(\gamma T)^{2/\alpha} + 4\zeta_2(\gamma^2 T)^{2/\alpha} + 4\zeta_3 R^{4/\alpha} \\ &\leq O\left((KdT)^{\frac{1}{1+\alpha}} + d^{2/\alpha}\right) \\ &\leq O\left(d^{2/\alpha} + d^{\frac{1}{1+\alpha}} n^{\frac{1}{2(1+\alpha)}}\right)\end{aligned}$$

922 Furthermore,

$$\begin{aligned}\frac{4\text{KL}(\mu_0|\mathcal{F}_0||\pi^*)}{\gamma T} + \frac{2\gamma B(4+L)\xi^2}{K} &\leq O\left(\frac{d}{\gamma T} + \frac{\gamma B(4+L)\xi^2}{2K}\right) \leq O\left(\frac{d^{1-\eta}}{(KT)^\eta}\right) \\ &\leq O\left(\frac{d^{\frac{2+\alpha}{2(1+\alpha)}}}{n^{\frac{\alpha}{4(1+\alpha)}}}\right)\end{aligned}$$

923 Hence, we conclude that,

$$\begin{aligned}\mathbb{E}[\text{KSD}_{\pi^*}^2(\hat{\nu}^{(n)}||\pi^*)] &\leq \frac{4K^2T^2\xi^2}{n^2} + \frac{2\xi^2}{n} + \frac{4\text{KL}(\mu_0|\mathcal{F}_0||\pi^*)}{\gamma T} + \frac{2\gamma B(4+L)\xi^2}{K} \\ &\leq \frac{6\xi^2}{n} + \frac{4\text{KL}(\mu_0|\mathcal{F}_0||\pi^*)}{\gamma T} + \frac{2\gamma B(4+L)\xi^2}{K} \\ &\leq O\left(\frac{d^{2/\alpha}}{n} + \frac{d^{\frac{1}{1+\alpha}}}{n^{\frac{1+2\alpha}{2(1+\alpha)}}} + \frac{d^{\frac{2+\alpha}{2(1+\alpha)}}}{n^{\frac{\alpha}{4(1+\alpha)}}}\right)\end{aligned}$$

924

□

925 E.3 Oracle Complexity of SVGD, VP-SVGD and GB-SVGD

926 We now compare the gradient oracle complexity, (i.e., the number of evaluations of ∇F) of VP-SVGD
 927 (as implied by Corollary 2) and GB-SVGD (as implied by Corollary 3) with that of SVGD as implied
 928 by the state-of-the-art finite particle guarantee of Shi and Mackey [37].

929 E.3.1 SVGD

930 From Equation (1), We note that T steps of SVGD run with n particles requires n^2T evaluations of
 931 ∇F .

932 **Subgaussian** π^* For subgaussian π^* , the finite-particle convergence rate obtained by Shi and
 933 Mackey [37] is $\text{KSD}_{\pi^*}(\hat{\mu}_{\text{SVGD}}^{(n)}||\pi^*) = \tilde{O}\left(\frac{\text{poly}(d)}{\sqrt{\log \log n^{\Theta(1/d)}}}\right)$, where $\hat{\mu}_{\text{SVGD}}^{(n)}$ denotes the empirical mea-
 934 sure of the n particles output by SVGD. By carefully following the analysis of Shi and Mackey
 935 [37], we infer that, to achieve $\text{KSD}_{\pi^*}(\hat{\mu}_{\text{SVGD}}^{(n)}||\pi^*) \leq \epsilon$, SVGD requires $T = \tilde{O}\left(\frac{\text{poly}(d)}{\epsilon^2}\right)$ and
 936 $n = \tilde{O}\left(\exp\left(\Theta\left(de^{\frac{\text{poly}(d)}{\epsilon^2}}\right)\right)\right)$. Thus the oracle complexity of SVGD (as implied by Shi and Mackey
 937 [37]) for achieving $\text{KSD}_{\pi^*}(\hat{\mu}_{\text{SVGD}}^{(n)}||\pi^*)$ is $\tilde{O}\left(\frac{\text{poly}(d)}{\epsilon^2} \cdot \exp\left(\Theta\left(de^{\frac{\text{poly}(d)}{\epsilon^2}}\right)\right)\right)$

938 E.3.2 VP-SVGD

939 From Algorithm 1, we note that T steps of VP-SVGD run with n particles and a batch-size of K
 940 requires $K^2T^2 + KTn$ evaluations of ∇F .

941 **Subgaussian** π^* For subgaussian π^* , Corollary 2 implies a finite-particle convergence rate of
 942 $\mathbb{E}[\text{KSD}_{\pi^*}^2(\hat{\mu}^{(n)}||\pi^*)] = O\left(\frac{d^{1/2}}{n^{1/2}} + \frac{d}{n}\right)$ (where $\hat{\mu}^{(n)}$ denotes the empirical measure of the n particles
 943 output by VP-SVGD) assuming $KT = d^{1/2}n^{3/2}$. Hence, to achieve $\mathbb{E}[\text{KSD}_{\pi^*}(\hat{\mu}^{(n)}||\pi^*)] \leq \epsilon$, VP-
 944 SVGD requires $n = O\left(\frac{d}{\epsilon^4}\right)$ and $KT = d^{1/2}n^{3/2} = \frac{d^2}{\epsilon^6}$. The resulting oracle complexity for achieving

945 $\mathbb{E}[\text{KSD}_{\pi^*}(\hat{\mu}^{(n)}|\pi^*)] \leq \epsilon$ is $O(\frac{d^4}{\epsilon^{12}})$. Compared to the oracle complexity of SVGD obtained above,
 946 this is a *double exponential improvement in both d and $1/\epsilon$* . Notably, the obtained oracle complexity
 947 guarantee *completely eliminates the curse of dimensionality*.

948 **Subexponential π^*** For subexponential π^* , Corollary 2 implies a finite-particle convergence rate of
 949 $\mathbb{E}[\text{KSD}_{\pi^*}^2(\hat{\mu}^{(n)}|\pi^*)] = O(\frac{d^{2/3}}{n^{1/3}} + \frac{d^2}{n})$ (where $\hat{\mu}^{(n)}$ denotes the empirical measure of the n particles
 950 output by VP-SVGD) assuming $KT = d^{1/3}n^{4/3}$. Hence, to achieve $\mathbb{E}[\text{KSD}_{\pi^*}(\hat{\mu}^{(n)}|\pi^*)] \leq \epsilon$, VP-
 951 SVGD requires $n = O(\frac{d^2}{\epsilon^6})$ and $KT = d^{1/3}n^{4/3} = \frac{d^3}{\epsilon^8}$. The resulting oracle complexity for achieving
 952 $\mathbb{E}[\text{KSD}_{\pi^*}(\hat{\mu}^{(n)}|\pi^*)] \leq \epsilon$ is $O(\frac{d^6}{\epsilon^{16}})$.

953 E.3.3 GB-SVGD

954 From Algorithm 2, we note that T steps of GB-SVGD run with n particles and a batch-size of K
 955 requires KTn evaluations of ∇F .

956 **Subgaussian π^*** For subgaussian π^* , Corollary 3 implies a finite-particle convergence rate of
 957 $\mathbb{E}[\text{KSD}_{\pi^*}^2(\hat{\nu}^{(n)}|\pi^*)] = O(\frac{d^{2/3}}{n^{1/6}} + \frac{d}{n})$ (where $\hat{\nu}^{(n)}$ denotes the empirical measure of the n particles
 958 output by GB-SVGD) assuming $KT = n^{1/2}$. Hence, to achieve $\mathbb{E}[\text{KSD}_{\pi^*}(\hat{\nu}^{(n)}|\pi^*)] \leq \epsilon$, GB-SVGD
 959 requires $n = \frac{d^4}{\epsilon^{12}}$ and $KT = \sqrt{n} = \frac{d^2}{\epsilon^6}$. Under this setting, the oracle complexity of GB-SVGD as
 960 implied by Corollary 3 is $O(\frac{d^6}{\epsilon^{18}})$. Compared to the oracle complexity of SVGD obtained above,
 961 this is a *double exponential improvement in both d and $1/\epsilon$* . Notably, the obtained oracle complexity
 962 guarantee *completely eliminates the curse of dimensionality*.

963 **Subexponential π^*** For subexponential π^* , Corollary 3 implies a finite-particle convergence rate of
 964 $\mathbb{E}[\text{KSD}_{\pi^*}^2(\hat{\nu}^{(n)}|\pi^*)] = O(\frac{d^{3/4}}{n^{1/8}} + \frac{d^2}{n})$ (where $\hat{\nu}^{(n)}$ denotes the empirical measure of the n particles
 965 output by GB-SVGD) assuming $KT = n^{1/2}$. Hence, to achieve $\mathbb{E}[\text{KSD}_{\pi^*}(\hat{\nu}^{(n)}|\pi^*)] \leq \epsilon$, GB-SVGD
 966 requires $n = \frac{d^6}{\epsilon^{16}}$ and $KT = \sqrt{n} = \frac{d^3}{\epsilon^8}$. Under this setting, the oracle complexity of GB-SVGD as
 967 implied by Corollary 3 is $O(\frac{d^9}{\epsilon^{24}})$.

968 F Literature Review

969 Initial works on the analysis of SVGD such as Liu [26], Lu et al. [30], Duncan et al. [14], Chewi
 970 et al. [7], Nüsken and Renger [33] consider the continuous-time population limit, i.e., the limit of
 971 infinite particles and vanishing step-sizes. In this regime, Liu [26], Lu et al. [30], Nüsken and Renger
 972 [33] show that the behavior of SVGD is characterized by a Partial Differential Equation (PDE), and
 973 established asymptotic convergence of this PDE to the target distribution. The work of Duncan et al.
 974 [14] proposes the Stein Logarithmic Sobolev Inequality which ensures exponential convergence of
 975 this PDE to the target distribution. However, characterizing the conditions under which this inequality
 976 holds is an open problem. The work of Chewi et al. [7] show that the PDE governing SVGD in
 977 the continuous-time population limit can be interpreted as an approximate Wasserstein gradient
 978 flow of the Chi-squared divergence. To this end, Chewi et al. [7] shows that the (exact) Wasserstein
 979 gradient flow of the Chi-squared divergence exhibits exponential convergence to the target distribution
 980 when π^* satisfies a Poincare Inequality. To the best of our knowledge, the first discrete-time non-
 981 asymptotic convergence result for population-limit SVGD was established in Korba et al. [23], where
 982 the authors interpreted population-limit SVGD as projected Wasserstein gradient descent. Their result
 983 relied on the assumption that the Kernel Stein Discrepancy to the target is uniformly bounded along
 984 the trajectory of SVGD, a condition which is hard to verify apriori. This result was significantly
 985 improved in Salim et al. [36], which established convergence of population-limit SVGD assuming the
 986 potential F is smooth the target $\pi^* \propto e^{-F}$ satisfies Talagrand's inequality \mathbb{T}_1 , an assumption which
 987 is equivalent to subgaussianity of π^* . This result was extended in Sun et al. [39] to accommodate for
 988 potentials F that satisfy a more general smoothness condition.

989 In comparison to prior works on population-limit SVGD, the literature on finite-particle SVGD
 990 is relatively sparse. The works of Liu [26] and Gorham et al. [18] establish that the dynamics of
 991 finite-particle SVGD asymptotically converge to that of population-limit SVGD in bounded Lipschitz

992 distance and Wasserstein-1 distance respectively, as the number of particles approaches infinity.
993 Under the stringent condition of bounded F (which is violated in various scenarios, e.g. log-strongly
994 concave π^*), Korba et al. [23] derived a non-asymptotic bound between the expected Wasserstein-2
995 distance between finite-particle SVGD and population-limit SVGD. To the best of our knowledge,
996 Shi and Mackey [37] is the only prior work that explicitly establishes a non-asymptotic convergence
997 guarantee of finite-particle SVGD to the target, which shows that the empirical measure of SVGD
998 run with n particles converges to the target density in KSD at a rate of $O\left(\sqrt{\frac{\text{poly}(d)}{\log \log n^{\Theta(1/d)}}}\right)$