

## MEASURING SAMPLE QUALITY WITH DIFFUSIONS<sup>1</sup>

BY JACKSON GORHAM<sup>\*</sup>, ANDREW B. DUNCAN<sup>†</sup>, SEBASTIAN J. VOLLMER<sup>‡</sup>  
AND LESTER MACKEY<sup>§</sup>

*Stanford University<sup>\*</sup>, Imperial College London<sup>†</sup>, University of Warwick<sup>‡</sup> and  
Microsoft Research New England<sup>§</sup>*

Stein’s method for measuring convergence to a continuous target distribution relies on an operator characterizing the target and *Stein factor* bounds on the solutions of an associated differential equation. While such operators and bounds are readily available for a diversity of univariate targets, few multivariate targets have been analyzed. We introduce a new class of characterizing operators based on Itô diffusions and develop explicit multivariate Stein factor bounds for any target with a fast-coupling Itô diffusion. As example applications, we develop computable and convergence-determining *diffusion Stein discrepancies* for log-concave, heavy-tailed and multimodal targets and use these quality measures to select the hyperparameters of biased Markov chain Monte Carlo (MCMC) samplers, compare random and deterministic quadrature rules and quantify bias-variance tradeoffs in approximate MCMC. Our results establish a near-linear relationship between diffusion Stein discrepancies and Wasserstein distances, improving upon past work even for strongly log-concave targets. The exposed relationship between Stein factors and Markov process coupling may be of independent interest.

**1. Introduction.** Consider a target probability distribution  $P$  with finite mean, continuously differentiable density  $p$  and support on all of  $\mathbb{R}^d$ . We will name the set of all such distributions  $\mathcal{P}_1$ . We assume that  $p$  can be evaluated up to its normalizing constant but that exact expectations under  $P$  are unattainable for most functions of interest. We will therefore use a *weighted sample*, represented as a discrete distribution  $Q_n = \sum_{i=1}^n q(x_i)\delta_{x_i}$ , to approximate intractable expectations  $\mathbb{E}_P[h(Z)]$  with tractable sample estimates  $\mathbb{E}_{Q_n}[h(X)] = \sum_{i=1}^n q(x_i)h(x_i)$ . Here, the support of  $Q_n$  is a collection of distinct sample points  $x_1, \dots, x_n \in \mathbb{R}^d$ , and the weight  $q(x_i)$  associated with each point is governed by a probability mass function  $q$ . We assume nothing about the process generating the sample points, so they may be the product of any random or deterministic mechanism.

---

Received February 2018; revised November 2018.

<sup>1</sup>This material is based upon work supported by EPSRC Grant EP/N000188/1, National Science Foundation DMS RTG Grant 1501767, National Science Foundation Graduate Research Fellowship Grant DGE-114747, the Frederick E. Terman Fellowship and the Lloyd’s Register Foundation programme on Data Centric engineering at the Alan Turing Institute.

*MSC2010 subject classifications.* Primary 60J60, 62-04, 62E17, 60E15, 65C60; secondary 62-07, 65C05, 68T05.

*Key words and phrases.* Multivariate Stein factors, Itô diffusion, Stein’s method, Stein discrepancy, sample quality, Wasserstein decay, Markov chain Monte Carlo.

Our ultimate goal is to develop a computable quality measure suitable for comparing any two samples approximating the same target distribution. More precisely, we seek to quantify how well  $\mathbb{E}_{Q_n}$  approximates  $\mathbb{E}_P$  in a manner that, at the very least, (i) indicates when a sample sequence is converging to  $P$ , (ii) identifies when a sample sequence is not converging to  $P$  and (iii) is computationally tractable. A natural starting point is to consider the maximum error incurred by the sample approximation over a class of scalar test functions  $\mathcal{H}$ ,

$$(1) \quad d_{\mathcal{H}}(Q_n, P) \triangleq \sup_{h \in \mathcal{H}} |\mathbb{E}_P[h(Z)] - \mathbb{E}_{Q_n}[h(X)]|.$$

When  $\mathcal{H}$  is convergence determining, the measure (1) is an *integral probability metric* (IPM) [69], and  $d_{\mathcal{H}}(Q_n, P)$  converges to zero only if the sample sequence  $(Q_n)_{n \geq 1}$  converges in distribution to  $P$ .

While a variety of standard probability metrics are representable as IPMs [69], the intractability of integration under  $P$  precludes us from computing most of these candidate quality measures. Recently, Gorham and Mackey [36] sidestepped this issue by constructing a class of test functions  $h$  known a priori to have zero mean under  $P$ . Their resulting quality measure—the *Langevin graph Stein discrepancy*—satisfied our computability and convergence detection requirements (Desiderata (i) and (iii)) and detected sample sequence nonconvergence (Desideratum (ii)) for strongly log concave targets with bounded third and fourth derivatives [65].

Our first contribution is to show that the Langevin Stein discrepancy in fact determines convergence for all smooth, *distantly dissipative* target distributions by explicitly lower and upper bounding the Langevin Stein discrepancy by standard Wasserstein distances. Distant dissipativity is a substantial relaxation of log concavity that covers a variety of common non-log concave targets like Gaussian mixtures and robust Student’s  $t$  regression posteriors. This contribution greatly extends the range of applicability of the Langevin Stein discrepancy.

Because heavy-tailed distributions are never distantly dissipative, as a second contribution, we extend the computable Stein discrepancy framework of [36] to accommodate heavy-tailed target distributions by introducing a new class of multivariate Stein operators based on general Itô diffusions. These operators can be used as drop-in replacements for the commonly used Langevin operator in applications.

As a third contribution, we establish a near linear relationship between the introduced *diffusion Stein discrepancies*  $\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|})$  and standard  $L^s$  Wasserstein distances  $W_{s, \|\cdot\|}(Q_n, P) \triangleq \inf_{X \sim Q_n, Z \sim P} \mathbb{E}[\|X - Z\|^s]^{1/s}$ . Namely,

$$W_{1, \|\cdot\|}(Q_n, P) \leq C_1 \mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|}) \max(1, \log(1/\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|}))) \quad \text{and}$$

$$\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|}) \leq C_2 W_{2, \|\cdot\|}(Q_n, P)$$

for constants  $C_1, C_2 > 0$  determined by Theorem 7 and Proposition 8. This improves upon prior analyses even in the case of strongly log concave targets.

Our primary contribution underlies these three advances. By relating Stein’s method to Markov process coupling rates in Section 2, we prove that every sufficiently fast coupling Itô diffusion gives rise to explicit, uniform multivariate *Stein factor* bounds on the derivatives of Stein equation solutions. Stein factor bounds are central to Stein’s method of measuring distributional convergence, and while a wealth of bounds are available for univariate targets (see, e.g., [10, 11, 89] for explicit bounds or [57] for a recent review), Stein factors for continuous multivariate distributions have largely been relegated to Gaussian [4, 9, 30, 38, 68, 70, 80], Dirichlet [29] and strongly log-concave [65] target distributions. Our approach, which exposes a general relationship between Stein factors and Markov process coupling times, extends the reach of Stein’s method to the stationary distributions of all fast coupling Itô diffusions.

In Section 3, we provide examples of practically checkable sufficient conditions for fast coupling and illustrate the process of verifying these conditions for canonical log-concave, heavy-tailed and multimodal targets. Section 4 describes a practical algorithm for computing diffusion Stein discrepancies using a geometric spanner and linear programming. In Section 5, we complement the principal theoretical contributions of this work with several simple numerical examples illustrating how diffusion Stein discrepancies can be deployed in practice. In particular, we use our discrepancies to select the hyperparameters of biased samplers, compare random and deterministic quadrature rules and quantify bias-variance tradeoffs in approximate Markov chain Monte Carlo. A discussion of related and future work follows in Section 6, and all proofs are deferred to the [Appendices](#).

NOTATION. For  $r \in [1, \infty]$ , let  $\|\cdot\|_r$  denote the  $\ell^r$  norm on  $\mathbb{R}^d$ . We will use  $\|\cdot\|$  as a generic norm on  $\mathbb{R}^d$  satisfying  $\|\cdot\| \geq \|\cdot\|_2$  and define the associated dual norms,  $\|v\|^* \triangleq \sup_{u \in \mathbb{R}^d: \|u\|=1} \langle u, v \rangle$  for vectors  $v \in \mathbb{R}^d$  and  $\|W\|^* \triangleq \sup_{u \in \mathbb{R}^d: \|u\|=1} \|Wu\|^*$  for matrices  $W \in \mathbb{R}^{d \times d}$ . Let  $e_j$  be the  $j$ th standard basis vector,  $\nabla_j$  be the partial derivative  $\frac{\partial}{\partial x_j}$ , and  $\lambda_{\min}(\cdot)$  and  $\lambda_{\max}(\cdot)$  be the smallest and largest eigenvalues of a symmetric matrix. For any real vector  $v$  and tensor  $T$ , let  $\|v\|_{\text{op}} \triangleq \|v\|_2$  and  $\|T\|_{\text{op}} \triangleq \sup_{\|u\|_2=1} \|T[u]\|_{\text{op}}$ . For each sufficiently differentiable vector- or matrix-valued function  $g$ , we define the bound  $M_0(g) \triangleq \sup_{x \in \mathbb{R}^d} \|g(x)\|_{\text{op}}$  and the  $k$ th order Hölder coefficients

$$M_k(g) \triangleq \sup_{x, y \in \mathbb{R}^d, x \neq y} \frac{\|\nabla^{\lceil k \rceil - 1} g(x) - \nabla^{\lceil k \rceil - 1} g(y)\|_{\text{op}}}{\|x - y\|_2^{\{k\}}}$$

where  $\{k\} \triangleq k - \lceil k - 1 \rceil$  and  $\nabla^0$  is the identity operator. For each differentiable matrix-valued function  $a$ , we let  $\langle \nabla, a(x) \rangle = \sum_j e_j \sum_k \nabla_k a_{jk}(x)$  represent the divergence operator applied to each row of  $a$  and define the Lipschitz coefficients  $F_k(a) \triangleq \sup_{x \in \mathbb{R}^d, \|v_1\|_2=1, \dots, \|v_k\|_2=1} \|\nabla^k a(x)[v_1, \dots, v_k]\|_F$  for  $\|\cdot\|_F$  the Frobenius norm. Finally, when the domain and range of a function  $f$  can be inferred from context, we write  $f \in C^k$  to indicate that  $f$  has  $k$  continuous derivatives.

**2. Stein’s method.** In the early 1970s, Charles Stein [88] introduced a powerful three-step approach to upper-bounding a reference IPM  $d_{\mathcal{H}}$ :

1. First, identify an operator  $\mathcal{T}$  that maps input functions<sup>2</sup>  $g : \mathbb{R}^d \rightarrow \mathbb{R}^d$  in a domain  $\mathcal{G}$  into mean-zero functions under  $P$ , that is,

$$\mathbb{E}_P[(\mathcal{T}g)(Z)] = 0 \quad \text{for all } g \in \mathcal{G}.$$

The operator  $\mathcal{T}$  and its domain  $\mathcal{G}$  define the *Stein discrepancy* [36],

$$\begin{aligned} \mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}) &\triangleq \sup_{g \in \mathcal{G}} |\mathbb{E}_{Q_n}[(\mathcal{T}g)(X)]| \\ (2) \qquad \qquad \qquad &= \sup_{g \in \mathcal{G}} |\mathbb{E}_{Q_n}[(\mathcal{T}g)(X)] - \mathbb{E}_P[(\mathcal{T}g)(Z)]| = d_{\mathcal{T}\mathcal{G}}(Q_n, P), \end{aligned}$$

a quality measure which takes the form of an integral probability metric while avoiding explicit integration under  $P$ .

2. Next, prove that, for each test function  $h$  in the reference class  $\mathcal{H}$ , the *Stein equation*

$$(3) \qquad \qquad \qquad h(x) - \mathbb{E}_P[h(Z)] = (\mathcal{T}g_h)(x)$$

admits a solution  $g_h \in \mathcal{G}$ . This step ensures that the reference metric  $d_{\mathcal{H}}$  lower bounds the Stein discrepancy (Desideratum (ii)) and, in practice, can be carried out simultaneously for large classes of target distributions.

3. Finally, use whatever means necessary to upper bound the Stein discrepancy and thereby establish convergence to zero under appropriate conditions (Desideratum (i)). Our general result, Proposition 8, suffices for this purpose.

While Stein’s method is traditionally used as analytical tool to establish rates of distributional convergence, we aim, following [36], to develop the method into a practical computational tool for measuring the quality of a sample. We begin by assessing the convergence properties of a broad class of Stein operators derived from Itô diffusions. Our efforts will culminate in Section 4, where we show how to explicitly compute the Stein discrepancy (2) given any sample measure  $Q_n$  and appropriate choices of  $\mathcal{T}$  and  $\mathcal{G}$ .

*2.1. Identifying a Stein operator.* To identify an operator  $\mathcal{T}$  that generates mean-zero functions under  $P$ , we will appeal to the elegant and widely applicable *generator method* construction of Barbour [3, 4] and Götze [38]. These authors note that if  $(Z_t)_{t \geq 0}$  is a Feller process with invariant measure  $P$ , then the *infinitesimal generator*  $\mathcal{A}$  of the process, defined pointwise by

$$(4) \qquad \qquad \qquad (\mathcal{A}u)(x) = \lim_{t \rightarrow 0} (\mathbb{E}[u(Z_t) \mid Z_0 = x] - u(x))/t$$

---

<sup>2</sup>Real-valued  $g$  are also common, but  $\mathbb{R}^d$ -valued  $g$  are more convenient for our purposes.

satisfies  $\mathbb{E}_P[(\mathcal{A}u)(Z)] = 0$  under very mild restrictions on  $u$  and  $\mathcal{A}$ . Gorham and Mackey [36] developed a *Langevin Stein operator* based on the generator of a specific Markov process—the Langevin diffusion described in (D1). Here, we will consider a broader class of continuous Markov processes known as *Itô diffusions*.

DEFINITION 1 (Itô diffusion [72], Definition 7.1.1). A (time-homogeneous) *Itô diffusion* with starting point  $x \in \mathbb{R}^d$ , Lipschitz drift coefficient  $b : \mathbb{R}^d \rightarrow \mathbb{R}^d$  and Lipschitz diffusion coefficient  $\sigma : \mathbb{R}^d \rightarrow \mathbb{R}^{d \times m}$  is a stochastic process  $(Z_{t,x})_{t \geq 0}$  solving the Itô stochastic differential equation

$$(5) \quad dZ_{t,x} = b(Z_{t,x}) dt + \sigma(Z_{t,x}) dW_t \quad \text{with } Z_{0,x} = x \in \mathbb{R}^d,$$

where  $(W_t)_{t \geq 0}$  is an  $m$ -dimensional Wiener process.

As the next theorem (distilled from [64], Theorem 2, and [75], Section 4.6) shows, it is straightforward to construct Itô diffusions with a given invariant measure  $P$  (see also [49, 52]).

THEOREM 2 ([64], Theorem 2, and [75], Section 4.6). Fix an Itô diffusion with  $C^1$  drift and diffusion coefficients  $b$  and  $\sigma$ , and define its covariance coefficient  $a(x) \triangleq \sigma(x)\sigma(x)^\top$ .  $P \in \mathcal{P}_1$  is an invariant measure of this diffusion if and only if  $b(x) = \frac{1}{2} \frac{1}{p(x)} \langle \nabla, p(x)a(x) \rangle + f(x)$  for a nonreversible component  $f \in C^1$  satisfying  $\langle \nabla, p(x)f(x) \rangle = 0$  for all  $x \in \mathbb{R}^d$ . If  $f$  is  $P$ -integrable, then

$$(6) \quad b(x) = \frac{1}{2} \frac{1}{p(x)} \langle \nabla, p(x)(a(x) + c(x)) \rangle$$

for  $c$  a differentiable  $P$ -integrable skew-symmetric  $d \times d$  matrix-valued function termed the stream coefficient [16, 56]. In this case, for all  $u \in C^2 \cap \text{dom}(\mathcal{A})$ , the infinitesimal generator (4) of the diffusion takes the form

$$(7) \quad (\mathcal{A}u)(x) = \frac{1}{2} \frac{1}{p(x)} \langle \nabla, p(x)(a(x) + c(x)) \nabla u(x) \rangle.^3$$

REMARKS. Theorem 2 does not require Lipschitz assumptions on  $b$  or  $\sigma$ . An example of a nonreversible component which is not  $P$ -integrable is  $f(x) = v/p(x)$  for any constant vector  $v \in \mathbb{R}^d$ . Prominent examples of  $P$ -targeted diffusions include:

---

<sup>3</sup>We have chosen an atypical form for the infinitesimal generator in (7), as it will give rise to a first-order differential operator (8) with more desirable properties. One can check, for instance, that the first-order operator  $(\mathcal{T}g)(x) = 2\langle b(x), g(x) \rangle + \langle a(x), \nabla g(x) \rangle$  derived from the standard form of the generator,  $(\mathcal{A}u)(x) = \langle b(x), \nabla u(x) \rangle + \frac{1}{2} \langle a(x), \nabla^2 u(x) \rangle$ , fails to satisfy Proposition 3 whenever the nonreversible component  $f(x) \neq 0$ .

(D1) the *(overdamped) Langevin diffusion* (also known as the *Brownian* or *Smoluchowski dynamics*) [75], Sections 6.5 and 4.5, where  $a \equiv I$  and  $c \equiv 0$ ;

(D2) the *preconditioned Langevin diffusion* [90], where  $c \equiv 0$  and  $a \equiv \sigma\sigma^\top$  for a constant diffusion coefficient  $\sigma \in \mathbb{R}^{d \times m}$ ;

(D3) the *Riemannian Langevin diffusion* [33, 52, 83], where  $c \equiv 0$  and  $a$  is not constant;

(D4) the *nonreversible preconditioned Langevin diffusion* (see, e.g., [20, 64, 81]), where  $a \equiv \sigma\sigma^\top$  for  $\sigma \in \mathbb{R}^{d \times m}$  constant and  $c$  not identically 0;

(D5) and the *second-order or underdamped Langevin diffusion* [45], where we target the joint distribution  $P \otimes \mathcal{N}(0, I)$  on  $\mathbb{R}^{2d}$  with

$$a \equiv 2 \begin{pmatrix} 0 & 0 \\ 0 & I \end{pmatrix} \quad \text{and} \quad c \equiv 2 \begin{pmatrix} 0 & -I \\ I & 0 \end{pmatrix}.$$

We will present detailed examples making use of these diffusion classes in Sections 3 and 5.

Theorem 2 forms the basis for our Stein operator of choice, the *diffusion Stein operator*  $\mathcal{T}$ , defined by substituting  $g$  for  $\frac{1}{2}\nabla u$  in the generator (7):

$$(8) \quad (\mathcal{T}g)(x) = \frac{1}{p(x)} \langle \nabla, p(x)(a(x) + c(x))g(x) \rangle.$$

$\mathcal{T}$  is an appropriate choice for our setting as it depends on  $P$  only through  $\nabla \log p$  and is therefore computable even when the normalizing constant of  $p$  is unavailable. One suitable domain for  $\mathcal{T}$  is the *classical Stein set* [36] of 1-bounded functions with 1-bounded, 1-Lipschitz derivatives:

$$\mathcal{G}_{\|\cdot\|} \triangleq \left\{ g : \mathbb{R}^d \rightarrow \mathbb{R} \mid \sup_{x \neq y \in \mathbb{R}^d} \max \left( \|g(x)\|^*, \|\nabla g(x)\|^*, \frac{\|\nabla g(x) - \nabla g(y)\|^*}{\|x - y\|} \right) \leq 1 \right\}.$$

Indeed, our next proposition, proved in Appendix A, shows that, on this domain, the diffusion Stein operator generates mean-zero functions under  $P$ .

**PROPOSITION 3.** *If  $\mathcal{T}$  is the diffusion Stein operator (8) for  $P \in \mathcal{P}_1$  with  $a, c \in C^1$  and  $a, c, b$  (6)  $P$ -integrable, then  $\mathbb{E}_P[(\mathcal{T}g)(Z)] = 0$  for all  $g \in \mathcal{G}_{\|\cdot\|}$ .*

Together,  $\mathcal{T}$  and  $\mathcal{G}_{\|\cdot\|}$  give rise to the *classical diffusion Stein discrepancy*  $\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|})$ , our primary object of study in Sections 2.2 and 2.3.

2.2. *Lower bounding the diffusion Stein discrepancy.* To establish that the classical diffusion Stein discrepancy detects nonconvergence (Desideratum (ii)), we will lower bound the discrepancy in terms of the  $L^1$  Wasserstein distance,  $d_{\mathcal{W}_{\|\cdot\|_2}} = W_{1, \|\cdot\|_2}$ , a standard reference IPM generated by

$$\mathcal{H} = \mathcal{W}_{\|\cdot\|_2} \triangleq \left\{ h : \mathbb{R}^d \rightarrow \mathbb{R} \mid \sup_{x \neq y \in \mathbb{R}^d} |h(x) - h(y)| \leq \|x - y\|_2 \right\}.$$

The first step is to show that, for each  $h \in \mathcal{W}_{\|\cdot\|_2}$ , the solution  $g_h$  to the Stein equation (3) with diffusion Stein operator (8) has low-order derivatives uniformly bounded by target-specific constants called *Stein factors*.

Explicit Langevin diffusion (D1) Stein factor bounds are readily available for a wide variety of univariate targets<sup>4</sup> (see, e.g., [10, 11, 89] for explicit bounds or [57] for a recent review). In contrast, in the multivariate setting, efforts to establish Stein factors have focused on Gaussian [4, 9, 30, 38, 68, 70, 80], Dirichlet [29] and strongly log-concave [65] targets with preconditioned Langevin (D2) operators. To extend the reach of the literature, we will derive multivariate Stein factors for targets with fast-coupling Itô diffusions. Our measure of coupling speed is the *Wasserstein decay rate*.

DEFINITION 4 (Wasserstein decay rate). Let  $(P_t)_{t \geq 0}$  be the *transition semigroup* of an Itô diffusion  $(Z_{t,x})_{t \geq 0}$  defined via

$$(P_t f)(x) \triangleq \mathbb{E}[f(Z_{t,x})] \quad \text{for all measurable } f, x \in \mathbb{R}^d, \text{ and } t \geq 0.$$

For any nonincreasing integrable function  $r : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}$ , we say that  $(P_t)_{t \geq 0}$  has *Wasserstein decay rate*  $r$  if

$$(9) \quad d_{\mathcal{W}_{\|\cdot\|_2}}(\delta_x P_t, \delta_y P_t) \leq r(t) d_{\mathcal{W}_{\|\cdot\|_2}}(\delta_x, \delta_y) \quad \text{for all } x, y \in \mathbb{R}^d \text{ and } t \geq 0,$$

where  $\delta_x P_t$  denotes the distribution of  $Z_{t,x}$ .

Our next result, proved in Appendix B, shows that the smoothness of a solution  $g_h$  to a Stein equation is controlled by the rate of Wasserstein decay, and hence by how quickly two diffusions with distinct starting points couple. The Stein factor bounds on the derivatives of  $u_h$  and  $g_h$  may be of independent interest for establishing rates of distributional convergence.

THEOREM 5 (Stein factors from Wasserstein decay). *Fix any Lipschitz  $h$ . If an Itô diffusion has invariant measure  $P \in \mathcal{P}_1$ , transition semigroup  $(P_t)_{t \geq 0}$ , Wasserstein decay rate  $r$  and infinitesimal generator  $\mathcal{A}$  (4), then*

$$(10) \quad u_h \triangleq \int_0^\infty \mathbb{E}_P[h(Z)] - P_t h \, dt$$

---

<sup>4</sup>The Langevin operator recovers Stein’s density method operator [89] when  $d = 1$ .



is twice continuously differentiable and satisfies

$$M_1(u_h) \leq M_1(h) \int_0^\infty r(t) dt \quad \text{and} \quad h - \mathbb{E}_P[h(Z)] = \mathcal{A}u_h.$$

Hence,  $g_h \triangleq \frac{1}{2} \nabla u_h$  solves the Stein equation (3) with diffusion Stein operator (8) whenever  $\mathcal{A}$  has the form (7). If the drift and diffusion coefficients  $b$  and  $\sigma$  have locally Lipschitz second derivatives and a right inverse  $\sigma^{-1}(x)$  for each  $x \in \mathbb{R}^d$  and  $h \in C^2$  with bounded second derivatives, then

$$(11) \quad M_2(u_h) \leq M_1(h)(\beta_1 + \beta_2),$$

where

$$\begin{aligned} \beta_1 &= r(0)(2M_0(\sigma^{-1}) + r(0)M_1(\sigma)M_0(\sigma^{-1}) + r(0)\sqrt{\alpha}), \quad \text{and} \\ \beta_2 &= r(0) \left( e^{\gamma^2} M_0(\sigma^{-1}) + e^{\gamma^2} M_1(\sigma)M_0(\sigma^{-1}) + \frac{2}{3} e^{\gamma^4} \sqrt{\alpha} \right) \int_0^\infty r(t) dt \end{aligned}$$

for  $\gamma_\rho \triangleq \rho M_1(b) + \frac{\rho^2 - 2\rho}{2} M_1(\sigma)^2 + \frac{\rho}{2} F_1(\sigma)^2$ ,  $\alpha \triangleq \frac{M_2(b)^2}{2M_1(b) + 4M_1(\sigma)^2} + 2F_2(\sigma)^2$ . If, additionally,  $\nabla^3 b$  and  $\nabla^3 \sigma$  are locally Lipschitz and  $h \in C^3$  with bounded third derivatives, then, for all  $\iota \in (0, 1)$ ,

$$(12) \quad M_{3-\iota}(u_h) \leq M_1(h) \frac{1}{K} \left( \frac{1}{\iota} + \int_0^\infty r(t) dt \right)$$

for  $K > 0$  a constant depending only on  $M_{1:3}(\sigma)$ ,  $M_{1:3}(b)$ ,  $M_0(\sigma^{-1})$  and  $r$ .

REMARK. Theorems 1 and 2 of Pardoux and Veretennikov [73] also bound the solutions of the Stein equation (3). However, for generic Lipschitz  $h$ , [73], Theorems 1 and 2, provide inexplicit constants; only guarantee the polynomial growth of  $g_h$  and its derivatives, not uniform boundedness; and require bounded  $\sigma$ , a strong assumption which rules out the heavy-tailed examples of Section 3.

A first consequence of Theorem 5, proved in Appendix D, concerns Stein operators (8) with constant covariance and stream matrices  $a$  and  $c$ . In this setting, fast Wasserstein decay implies that the diffusion Stein discrepancy converges to zero only if the Wasserstein distance does (Desideratum (ii)).

THEOREM 6 (Stein discrepancy lower bound: constant  $a$  and  $c$ ). Consider an Itô diffusion with diffusion Stein operator  $\mathcal{T}$  (8) for  $P \in \mathcal{P}_1$ , Wasserstein decay rate  $r$ , constant covariance and stream matrices  $a$  and  $c$  and Lipschitz drift  $b(x) = \frac{1}{2}(a + c)\nabla \log p(x)$ . If  $s_r \triangleq \int_0^\infty r(t) dt$ , then

$$(13) \quad d_{\mathcal{W}_{\|\cdot\|_2}}(Q_n, P) \leq 3s_r \max \left( S(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|}) \right. \\ \left. \sqrt[3]{S(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|}) \sqrt{2\mathbb{E}[\|G\|_2]^2} \left( 2M_1(b) + \frac{1}{s_r} \right)^2} \right),$$

where  $G \in \mathbb{R}^d$  is a standard normal vector and  $M_1(b) \leq \frac{1}{2}\|a + c\|_{\text{op}} M_2(\log p)$ .



Theorem 6 in fact provides an explicit upper bound on the Wasserstein distance in terms of the Stein discrepancy and the Wasserstein decay rate. Under additional smoothness assumptions on the coefficients, the explicit relationship between Stein discrepancy and Wasserstein distance can be improved and extended to diffusions with nonconstant diffusion coefficient, as our next result, proved in Appendix E, shows.

**THEOREM 7** (Stein discrepancy lower bound: nonconstant  $a$  and  $c$ ). *Consider an Itô diffusion for  $P \in \mathcal{P}_1$  with diffusion Stein operator  $\mathcal{T}$  (8), Wasserstein decay rate  $r$  and Lipschitz drift and diffusion coefficients  $b$  (6) and  $\sigma$  with locally Lipschitz second derivatives. If  $s_r \triangleq \int_0^\infty r(t) dt$ , then*

$$\begin{aligned} d_{\mathcal{W}_{\|\cdot\|_2}}(Q_n, P) &\leq 2 \max(\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|}) \max(s_r, \beta_1 + \beta_2), \\ &\quad \sqrt{\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|})} \sqrt{2/\pi} (\beta_1 + \beta_2) \zeta), \end{aligned}$$

for  $\beta_1, \beta_2$  defined in Theorem 5 and

$$\zeta \triangleq \mathbb{E}[\|G\|_2] (1 + 2M_1(b)s_r + M_1^*(m)(\beta_1 + \beta_2)),$$

where  $G \in \mathbb{R}^d$  is a standard normal vector,  $m \triangleq a + c$ , and  $M_1^*(m) \triangleq \sup_{x \neq y} \|m(x) - m(y)\|_{\text{op}}^* / \|x - y\|_2$ .

If, additionally,  $\nabla^3 b$  and  $\nabla^3 \sigma$  are locally Lipschitz, then

$$\begin{aligned} (14) \quad d_{\mathcal{W}_{\|\cdot\|_2}}(Q_n, P) &\leq 2\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|}) \max\left(\max(s_r, \beta_1 + \beta_2), \right. \\ &\quad \left. e \max\left(\frac{d^{1/4} \sqrt{\zeta}}{\sqrt{K}}, \frac{\sqrt{d}}{K}\right) (s_r + \max(\log(1/\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|})), 1))\right), \end{aligned}$$

for a constant  $K > 0$  depending only on  $M_{1:3}(\sigma)$ ,  $M_{1:3}(b)$ ,  $M_0(\sigma^{-1})$  and  $r$ .

**REMARK.** The  $\log(1/\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|}))$  term in (14) reflects the potential non-smoothness of the Stein equation solution  $g_h$  studied in Theorem 5. Indeed, for  $d \geq 2$  and standard multivariate Gaussian  $P$ , there exist Lipschitz  $h$  with infinite  $M_2(g_h)$  [78], Remark 2.

In Section 3, we will present practically checkable conditions implying fast Wasserstein decay and discuss both broad families and specific diffusion-target pairings covered by this theory.

2.3. *Upper bounding the diffusion Stein discrepancy.* In upper bounding the Stein discrepancy, one classically aims to establish rates of convergence to  $P$  for specific sequences  $(Q_n)_{n=1}^\infty$ . Since our interest is in explicitly computing Stein discrepancies for arbitrary sample sequences, our general upper bound in Proposition 8 serves principally to provide sufficient conditions under which the classical diffusion Stein discrepancy converges to zero.

PROPOSITION 8 (Stein discrepancy upper bound). *Let  $\mathcal{T}$  be the diffusion Stein operator (8) for  $P \in \mathcal{P}_1$ . If  $m \triangleq a + c$  and  $b$  (6) are  $P$ -integrable,*

$$\begin{aligned} S(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|}) &\leq \inf_{X \sim Q_n, Z \sim P} (\mathbb{E}[2\|b(X) - b(Z)\| + \|m(X) - m(Z)\|] \\ &\quad + \mathbb{E}[(2\|b(Z)\| + \|m(Z)\|) \min(\|X - Z\|, 2)]) \\ &\leq W_{s, \|\cdot\|}(Q_n, P)(2M_1^{\|\cdot\|}(b) + M_1^{\|\cdot\|}(m)) \\ &\quad + W_{s, \|\cdot\|}(Q_n, P)^t 2^{1-t} \mathbb{E}[(2\|b(Z)\| + \|m(Z)\|)^{s/(s-t)}]^{(s-t)/s} \end{aligned}$$

for any  $s \geq 1$  and  $t \in (0, 1]$ . Moreover, for  $\mu_0 \triangleq \mathbb{E}[e^{2\|b(Z)\| + \|m(Z)\|}]$ ,

$$\begin{aligned} S(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|}) &\leq W_{1, \|\cdot\|}(Q_n, P)(2M_1^{\|\cdot\|}(b) + M_1^{\|\cdot\|}(m)) \\ &\quad + \min(W_{1, \|\cdot\|}(Q_n, P), 2) \log((e\mu_0)/\min(W_{1, \|\cdot\|}(Q_n, P), 2)). \end{aligned}$$

This result, proved in Appendix F, complements the Wasserstein distance lower bounds of Section 2.2 and implies that, for Lipschitz and sufficiently integrable  $m$  and  $b$ , the diffusion Stein discrepancy converges to zero whenever  $Q_n$  converges to  $P$  in Wasserstein distance.

2.4. *Extension to nonuniform Stein sets.* For any  $c_1, c_2, c_3 > 0$ , our analyses and algorithms readily accommodate the nonuniform Stein set

$$\begin{aligned} \mathcal{G}_{\|\cdot\|}^{c_1, c_2, c_3} &\triangleq \left\{ g : \mathbb{R}^d \right. \\ &\quad \left. \rightarrow \mathbb{R}^d \mid \sup_{x \neq y \in \mathbb{R}^d} \max \left( \frac{\|g(x)\|^*}{c_1}, \frac{\|\nabla g(x)\|^*}{c_2}, \frac{\|\nabla g(x) - \nabla g(y)\|^*}{c_3 \|x - y\|} \right) \leq 1 \right\}. \end{aligned}$$

This added flexibility can be valuable when tight upper bounds on a reference IPM, like the Wasserstein distance, are available for a particular choice of Stein factors  $(c_1, c_2, c_3)$ . When such Stein factors are unknown or difficult to compute, we recommend the parameter-free classical Stein set and graph Stein set of the sequel as practical defaults, since the classical Stein discrepancy is strongly equivalent to any nonuniform Stein discrepancy.

PROPOSITION 9 (Equivalence of nonuniform Stein discrepancies). *For any  $c_1, c_2, c_3 > 0$ ,*

$$\begin{aligned} \min(c_1, c_2, c_3)\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|}) &\leq \mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|}^{c_1:c_3}) \\ &\leq \max(c_1, c_2, c_3)\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|}). \end{aligned}$$

REMARK. The short proof follows exactly as in [36], Proposition 4.

**3. Sufficient conditions for Wasserstein decay.** Since the Stein discrepancy lower bounds of Section 2 depend on the Wasserstein decay (9) of the chosen diffusion, we next provide examples of practically checkable sufficient conditions for Wasserstein decay and illustrate the process of verifying these conditions for a selection of diffusion-target pairings. These pedagogical examples serve to succinctly illustrate the process of verifying our assumptions and do not represent the full scope of applicability.

3.1. *Uniform dissipativity.* It is well known (see, e.g., [7], equation (7)) that the Langevin diffusion (D1) enjoys exponential Wasserstein decay whenever  $\log p$  is  $k$ -strongly log concave, that is, when the drift  $b = \frac{1}{2}\nabla \log p$  satisfies  $\langle b(x) - b(y), x - y \rangle \leq -\frac{k}{2}\|x - y\|_2^2$  for  $k > 0$ . An analogous *uniform dissipativity* condition gives explicit exponential decay for a generic Itô diffusion.

THEOREM 10 (Wasserstein decay: uniform dissipativity). *Fix  $k > 0$  and  $G \succ 0$ , and let  $\|w\|_G^2 \triangleq \langle w, Gw \rangle$ , for any vector or matrix  $w \in \mathbb{R}^{d \times d'}$ ,  $d' \geq 1$ . An Itô diffusion with drift and diffusion coefficients  $b$  and  $\sigma$  satisfying*

$$\begin{aligned} 2\langle b(x) - b(y), G(x - y) \rangle + \|\sigma(x) - \sigma(y)\|_G^2 \\ \leq -k\|x - y\|_G^2 \quad \text{for all } x, y \in \mathbb{R}^d \end{aligned}$$

*has Wasserstein decay rate (9)  $r(t) = e^{-kt/2} \sqrt{\lambda_{\max}(G)/\lambda_{\min}(G)}$ .*

REMARK. The proof of Theorem 10 in Appendix G holds even when the drift  $b$  is not Lipschitz, yields the same decay rate for  $\mathcal{W}_{2, \|\cdot\|_2}$ , and relies on a synchronous coupling of Itô diffusions, mimicking [7], Section 1.

Hence, if the drift  $b$  of an Itô diffusion is  $-k/2$ -one-sided Lipschitz, that is,

$$(15) \quad 2\langle b(x) - b(y), G(x - y) \rangle \leq -k\|x - y\|_G^2 \quad \text{for all } x, y \in \mathbb{R}^d$$

and some  $G \succ 0$ , and the diffusion coefficient  $\sigma$  is  $\sqrt{k'}$ -Lipschitz, that is,

$$\|\sigma(x) - \sigma(y)\|_G^2 \leq k'\|x - y\|_G^2 \quad \text{for all } x, y \in \mathbb{R}^d,$$

then, whenever  $k' < k$ , the diffusion exhibits exponential Wasserstein decay. with rate  $e^{-(k-k')t/2} \sqrt{\lambda_{\max}(G)/\lambda_{\min}(G)}$ .

EXAMPLE 1 (Bayesian logistic regression with Gaussian prior). A one-sided Lipschitz drift arises naturally in the setting of Bayesian logistic regression [32], a canonical model of binary outcomes  $y \in \{-1, 1\}$  given measured covariates  $v \in \mathbb{R}^d$ . Consider the log density of a Bayesian logistic regression posterior based on a dataset of  $L$  observations  $(v_l, y_l)$  and a  $\mathcal{N}(\mu, \Sigma)$  prior:

$$\log p(\beta) = \underbrace{-\frac{1}{2} \|\Sigma^{-1/2}(\beta - \mu)\|_2^2}_{\text{multivariate Gaussian prior}} - \underbrace{\sum_{l=1}^L \log(1 + \exp(-y_l \langle v_l, \beta \rangle))}_{\text{logistic regression likelihood}} + \text{const.}$$

Here, our inferential target is the unobserved parameter vector  $\beta \in \mathbb{R}^d$ . Since

$$\begin{aligned} -\Sigma^{-1} &\succcurlyeq \nabla^2 \log p(\beta) \\ &= -\Sigma^{-1} - \sum_{l=1}^L \frac{e^{y_l \langle v_l, \beta \rangle}}{(1 + e^{y_l \langle v_l, \beta \rangle})^2} v_l v_l^\top \succcurlyeq -\Sigma^{-1} - \frac{1}{4} \sum_{l=1}^L v_l v_l^\top, \end{aligned}$$

the  $P$ -targeted preconditioned Langevin diffusion (D2) drift  $b(\beta) = \frac{1}{2} \Sigma \nabla \log p(\beta)$  satisfies (15) with  $k = 1$  and  $G = \Sigma^{-1}$  and  $M_1(b) \leq \frac{1}{2} \|I + \frac{1}{4} \Sigma \sum_{l=1}^L v_l v_l^\top\|_{\text{op}}$ . Hence, the diffusion enjoys geometric Wasserstein decay (Theorem 10) and a Wasserstein lower bound on the Stein discrepancy (Theorem 6).

EXAMPLE 2 (Bayesian Huber regression with Gaussian prior). Huber’s least favorable distribution provides a robust error model for the regression of a continuous response  $y \in \mathbb{R}$  onto a vector of measured covariates  $v \in \mathbb{R}^d$  [46]. Given  $L$  observations  $(v_l, y_l)$  and a  $\mathcal{N}(\mu, \Sigma)$  prior on an unknown parameter vector  $\beta \in \mathbb{R}^d$ , the Bayesian Huber regression log posterior takes the form

$$\log p(\beta) = \underbrace{-\frac{1}{2} \|\Sigma^{-1/2}(\beta - \mu)\|_2^2}_{\text{multivariate Gaussian prior}} - \underbrace{\sum_{l=1}^L \rho_c(y_l - \langle v_l, \beta \rangle)}_{\text{Huber's least favorable likelihood}} + \text{const.}$$

where  $\rho_c(r) \triangleq \frac{1}{2} r^2 \mathbb{I}[|r| \leq c] + c(|r| - \frac{1}{2}c) \mathbb{I}[|r| > c]$  for fixed  $c > 0$ . Since  $\rho'_c(r) = \min(\max(r, -c), c)$  is 1-Lipschitz,  $\rho_c$  is convex, and the Hessian of the log prior is  $-\Sigma^{-1}$ , the  $P$ -targeted preconditioned Langevin diffusion (D2) drift  $b(\beta) = \frac{1}{2} \Sigma \nabla \log p(\beta)$  satisfies (15) with  $k = 1$  and  $G = \Sigma^{-1}$  and  $M_1(b) \leq \frac{1}{2} \|I + \Sigma \sum_{l=1}^L v_l v_l^\top\|_{\text{op}}$ . This is again sufficient for exponential Wasserstein decay and a Wasserstein lower bound on the Stein discrepancy.

3.2. *Distant dissipativity, constant  $\sigma$ .* When the diffusion coefficient  $\sigma$  is constant with  $a \triangleq \frac{1}{2} \sigma \sigma^\top$  invertible, Eberle [22] showed that a *distant dissipativity* condition is sufficient for exponential Wasserstein decay. Specifically, if we define a one-sided Lipschitz constant conditioned on a distance  $r > 0$ ,

$$-\kappa(r) = \sup\{2(b(x) - b(y))^\top a^{-1}(x - y)/r^2 : (x - y)^\top a^{-1}(x - y) = r^2\},$$

then [22], Corollary 2, establishes exponential Wasserstein decay whenever  $\kappa$  is continuous with  $\liminf_{r \rightarrow \infty} \kappa(r) > 0$  and  $\int_0^1 r\kappa(r)^- dr < \infty$ . For a Lipschitz drift, this holds whenever  $b$  is dissipative at large distances, that is, whenever, for some  $k > 0$ , we have  $\kappa(r) \geq k$  for all  $r$  sufficiently large [22], Lemma 1.

**EXAMPLE 3 (Gaussian mixture with common covariance).** Consider an  $m$ -component mixture density  $p(x) = \sum_{j=1}^m w_j \phi_j(x)$ , where the component weights  $w_j \geq 0$  sum to one and  $\phi_j$  is the density of a  $\mathcal{N}(\mu_j, \Sigma)$  distribution on  $\mathbb{R}^d$ . Fix any  $x, y \in \mathbb{R}^d$ . If  $\|\Sigma^{-1/2}(x - y)\|_2 = r$ , the  $P$ -targeted preconditioned Langevin diffusion (D2) with drift  $b(z) = \frac{1}{2}a \nabla \log p(z)$  and  $a = \Sigma$  satisfies

$$\begin{aligned} & 2(b(x) - b(y))^\top a^{-1}(x - y) \\ &= (\nabla \log p(x) - \nabla \log p(y))^\top (x - y) \\ &= -r^2 + \langle \Sigma^{-1/2}(\mu(x) - \mu(y)), \Sigma^{-1/2}(x - y) \rangle \leq -r^2 + r\Delta, \end{aligned}$$

by Cauchy–Schwarz and Jensen’s inequality, for  $\Delta \triangleq \sup_{j,k} \|\Sigma^{-1/2}(\mu_j - \mu_k)\|_2$ ,  $\mu(x) \triangleq \sum_{j=1}^m \pi_j(x)\mu_j$ , and  $\pi_j(x) \triangleq \frac{w_j \phi_j(x)}{p(x)}$ . Moreover, by the mean value theorem, Cauchy–Schwarz and Jensen’s inequality, we have, for each  $v \in \mathbb{R}^d$ ,

$$\begin{aligned} & 2\langle \Sigma^{-1/2}(b(x) - b(y)), v \rangle \\ &= \langle \Sigma^{-1/2}(\nabla \mu(z) - I)(x - y), v \rangle \\ &= \langle (\Sigma^{-1/2}S(z)\Sigma^{-1/2} - I)\Sigma^{-1/2}(x - y), v \rangle \leq \|v\|_2 \|\Sigma^{-1/2}(x - y)\|_2 L, \end{aligned}$$

for some  $z \in \mathbb{R}^d$ ,  $S(x) \triangleq \frac{1}{2} \sum_{j=1}^m \sum_{k=1}^m \pi_j(x)\pi_k(x)(\mu_j - \mu_k)(\mu_j - \mu_k)^\top$ , and  $L \triangleq \sup_{j,k} |1 - \|\Sigma^{-1/2}(\mu_j - \mu_k)\|_2^2/2|$ . Hence,  $b$  is Lipschitz, and  $\kappa(r) \geq \frac{1}{2}$  when  $r > 2\Delta$ , so our diffusion enjoys exponential Wasserstein decay [22], Lemma 1, and a Stein discrepancy upper bound on the Wasserstein distance.

**3.3. Distant dissipativity, nonconstant  $\sigma$ .** Using a combination of synchronous and reflection couplings, Wang [94], Theorem 2.6, showed that diffusions satisfying a distant dissipativity condition exhibit exponential Wasserstein decay, even when the diffusion coefficient  $\sigma$  is nonconstant. In Appendix H, we combine the coupling strategy of [94], Theorem 2.6, with the approach of [22] for diffusions with constant  $\sigma$  to obtain the following explicit Wasserstein decay rate for distantly dissipative diffusions with bounded  $\sigma^{-1}$ .

**THEOREM 11 (Wasserstein decay: distant dissipativity).** *Let  $(P_t)_{t \geq 0}$  be the transition semigroup of an Itô diffusion with drift and diffusion coefficients  $b$  and  $\sigma$ . Define the truncated diffusion coefficient*

$$\sigma_0(x) = (\sigma(x)\sigma(x)^\top - \lambda_0^2 I)^{1/2} \quad \text{for some } \lambda_0 \in [0, 1/M_0(\sigma^{-1})]$$

and the distance-conditional dissipativity function

$$(16) \quad \kappa(r) = \inf \left\{ -2\alpha \left( \langle b(x) - b(y), x - y \rangle + \frac{1}{2} \|\sigma_0(x) - \sigma_0(y)\|_F^2 - \frac{1}{2} \|(\sigma_0(x) - \sigma_0(y))^\top (x - y)\|_2^2 / r^2 \right) / r^2 : \|x - y\|_2 = r \right\}$$

for any  $m_0 \leq \inf_{x \neq y} \frac{\|(\sigma_0(x) - \sigma_0(y))^\top (x - y)\|_2}{\|x - y\|_2}$  and  $\alpha \triangleq 1/(\lambda_0^2 + m_0^2/4)$ .

If the constants  $R_0 = \inf\{R \geq 0 : \kappa(r) \geq 0, \forall r \geq R\}$  and  $R_1 = \inf\{R \geq R_0 : \kappa(r)R(R - R_0) \geq 8, \forall r \geq R\}$  satisfy  $R_0 \leq R_1 < \infty$ , then

$$(17) \quad d_{\mathcal{W}_{\|\cdot\|_2}}(\delta_x P_t, \delta_y P_t) \leq 2\varphi(R_0)^{-1} e^{-ct} d_{\mathcal{W}_{\|\cdot\|_2}}(\delta_x, \delta_y)$$

for all  $x, y \in \mathbb{R}^d$  and  $t \geq 0$ , where  $\frac{1}{c} = \alpha \int_0^{R_1} \int_0^s \exp(\frac{1}{4} \int_t^s u \kappa^-(u) du) dt ds$ ,  $\varphi(r) = e^{-\frac{1}{4} \int_0^r s \kappa^-(s) ds}$  and  $\kappa^-(s) = \max(-\kappa(s), 0)$ .

REMARK 1. Theorem 11 holds even when the drift  $b$  is not Lipschitz.

The Wasserstein decay rate (17) in Theorem 11 has a simple form when the diffusion is dissipative at large distances and  $\kappa$  is bounded below. These rates follow exactly as in [22], Lemma 1.

COROLLARY 12. Under the conditions of Theorem 11, suppose that, for  $R, L \geq 0$  and  $K > 0$ ,  $\kappa(r) \geq -L$  for  $r \leq R$  and  $\kappa(r) \geq K$  for  $r > R$ . Then

$$\alpha^{-1} c^{-1} \leq \begin{cases} \frac{e-1}{2} R^2 + e\sqrt{8K^{-1}R} + 4K^{-1} & \text{if } LR_0^2 \leq 8, \\ \frac{8\sqrt{2\pi}}{RL^{1/2}} (L^{-1} + K^{-1}) \exp\left(\frac{LR^2}{8}\right) + 32R^{-2}K^{-2} & \text{if } LR_0^2 > 8. \end{cases}$$

EXAMPLE 4 (Multivariate Student’s t regression with pseudo-Huber prior). The multivariate Student’s t distribution is also commonly employed as a robust error model for the linear regression of continuous responses  $y \in \mathbb{R}^L$  onto measured covariates  $V \in \mathbb{R}^{L \times d}$  [58, 95]. Under a pseudo-Huber prior [44], a Bayesian multivariate Student’s t regression posterior takes the form

$$p(\beta) \propto \underbrace{\exp\left(\delta^2 \left(1 - \sqrt{1 + \|\beta/\delta\|_2^2}\right)\right)}_{\text{pseudo-Huber prior}} \underbrace{\left(1 + \frac{1}{\nu} (y - V\beta)^\top \Sigma^{-1} (y - V\beta)\right)^{-(\nu+L)/2}}_{\text{multivariate Student's t likelihood}}$$

for fixed  $\delta, \nu > 0$  and  $\Sigma \succ 0$ . Introduce the shorthand  $\psi_\lambda(r) \triangleq 2\sqrt{1 + r^2/\delta^2} - \lambda^2$  for each  $\lambda \in [0, \sqrt{2})$  and  $\xi(\beta) \triangleq 1 + \frac{1}{\nu} (y - V\beta)^\top \Sigma^{-1} (y - V\beta)$ . Since

$$\nabla \log p(\beta) = -2\beta/\psi_0(\|\beta\|_2) + \left(1 + \frac{\nu}{L}\right) V^\top \Sigma^{-1} (y - V\beta) / \xi(\beta)$$

is bounded, no  $P$ -targeted preconditioned Langevin diffusion (D2) will satisfy the distant dissipativity conditions of Section 3.2. However, we will show that whenever  $V^\top V \succ 0$ , the Riemannian Langevin diffusion (D3) with  $\sigma(\beta) = \sqrt{\psi_0(\|\beta\|_2)}I \in \mathbb{R}^{d \times d}$ ,  $a(\beta) = \frac{1}{2}\psi_0(\|\beta\|_2)I$ , and  $b(\beta) = a(\beta)\nabla \log p(\beta) + \langle \nabla, a(\beta) \rangle$  satisfies the Wasserstein decay preconditions of Corollary 12.

Indeed, fix any  $\lambda_0 \in (0, 1/M_0(\sigma^{-1})) = (0, \sqrt{2})$ . Since  $M_1(\sqrt{\psi_\lambda}) \leq \frac{1}{\delta\sqrt{2-\lambda^2}}$ ,  $M_1(\psi_\lambda) \leq \frac{2}{\delta}$ , and  $M_2(\psi_\lambda) \leq \frac{2}{\delta^2}$ ,  $\sigma_0, \sigma, a$  and  $\nabla a$  are all Lipschitz. The drift  $b$  is also Lipschitz, since  $\nabla \log p$  and the product of  $a(\beta)$  and

$$\begin{aligned} &\nabla^2 \log p(\beta) \\ &= -2I/\psi_0(\|\beta\|_2) + 8\beta\beta^\top/(\delta^2\psi_0^3(\|\beta\|_2)) \\ &\quad + \left(1 + \frac{\nu}{L}\right)(2V^\top \Sigma^{-1}(y - V\beta)(y - V\beta)^\top \Sigma^{-1}V/\xi^2(\beta) \\ &\quad - V^\top \Sigma^{-1}V/\xi(\beta)) \end{aligned}$$

are bounded. Hence,  $\kappa$  (16) is bounded below. Moreover, the Hölder continuity of  $x \mapsto \sqrt{x}$ , Cauchy–Schwarz and the triangle inequality imply

$$\begin{aligned} \kappa(r) &\geq \inf_{\|\beta-\beta'\|_2=r} \frac{2\alpha}{r^2} \left( \langle b(\beta') - b(\beta), \beta - \beta' \rangle \right. \\ &\quad \left. - \frac{d-1}{2} \left| \sqrt{\psi_{\lambda_0}(\|\beta\|_2)} - \sqrt{\psi_{\lambda_0}(\|\beta'\|_2)} \right|^2 \right) \\ &\geq 2\alpha - \frac{2\alpha}{r} \left( \frac{d-1}{\delta} + M_1(\psi_0) \right. \\ &\quad \left. + \sup_{\beta} \left( 1 + \frac{\nu}{L} \right) \psi_0(\|\beta\|_2) \|V^\top \Sigma^{-1}(y - V\beta)\|_2 / \xi(\beta) \right) \\ &\geq 2\alpha - \frac{2\alpha}{r} \left( \frac{d+1}{\delta} \right. \\ &\quad \left. + \sup_s \left( 1 + \frac{\nu}{L} \right) \frac{2(1+s/\delta)(\|V^\top \Sigma^{-1}y\|_2 + s\|V^\top \Sigma^{-1}V\|_{\text{op}})}{1 + \frac{1}{\nu} \max(0, s/\|(V^\top \Sigma^{-1}V)^{-1}\|_{\text{op}} - \|\Sigma^{-1}y\|_2)^2} \right). \end{aligned}$$

Letting  $\zeta$  represent the supremum in the final inequality, we see that  $\kappa(r) \geq \alpha = 1/\lambda_0^2$  whenever  $r \geq 2(\frac{d+1}{\delta} + \zeta)$ . Hence, Corollary 12 delivers exponential Wasserstein decay. A Wasserstein lower bound on the Stein discrepancy now follows from Theorem 7, since  $M_2(\sqrt{\psi_0}) \leq \frac{1}{\sqrt{2}\delta^2}$ ,  $M_3(\psi_0) \leq \frac{96}{25\sqrt{5}\delta^3}$ , and  $a(\beta)\nabla^2 \log p(\beta)$  is Lipschitz, and hence  $M_2(\sigma)$  and  $M_2(b)$  are bounded.

**4. Computing Stein discrepancies.** In this section, we introduce a computationally tractable Stein discrepancy that inherits the favorable convergence properties established in Sections 2 and 3. We will directly port the spanner discrepancy



methodology developed and detailed in [36] and use our new diffusion operators as drop-in replacements for the overdamped Langevin operators advocated in [36]. While we only explicitly discuss target distributions supported on all of  $\mathbb{R}^d$ , constrained domains of the form  $(\alpha_1, \beta_1) \times \dots \times (\alpha_d, \beta_d)$  where  $-\infty \leq \alpha_i < \beta_i \leq \infty$  for all  $1 \leq i \leq d$  can be handled by introducing boundary constraints as in [36], Section 4.4.

4.1. *Spanner Stein discrepancies.* For any sample  $Q_n$ , Stein operator  $\mathcal{T}$ , and Stein set  $\mathcal{G}$ , the Stein discrepancy  $\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G})$  is recovered by solving an optimization problem over functions  $g \in \mathcal{G}$ . For example, if we write  $m \triangleq a + c$  and  $b(x) \triangleq \frac{1}{2} \frac{1}{p(x)} \langle \nabla, p(x)m(x) \rangle$ , the classical diffusion Stein discrepancy is the value

$$\begin{aligned} &\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|}) \\ &= \sup_g \sum_{i=1}^n q(x_i) (2\langle b(x_i), g(x_i) \rangle + \langle m(x_i), \nabla g(x_i) \rangle) \\ &\text{s.t. } \max \left( \|g(x)\|^*, \|\nabla g(x)\|^*, \frac{\|\nabla g(x) - \nabla g(y)\|^*}{\|x - y\|} \right) \leq 1, \quad \forall x, y \in \mathbb{R}^d. \end{aligned}$$

For all Stein sets, the diffusion Stein discrepancy objective is linear in  $g$  and only queries  $g$  and  $\nabla g$  at the  $n$  sample points underlying  $Q_n$ . However, the classical Stein set  $\mathcal{G}_{\|\cdot\|}$  constrains  $g$  at all points in its domain, resulting in an infinite-dimensional optimization problem.<sup>5</sup>

To obtain a finite-dimensional problem that is convergence-determining and straightforward to optimize, we will make use of the *graph Stein sets* of [36]. For a given graph  $G = (V, E)$  with  $V = \text{supp}(Q_n)$ , the graph Stein set,

$$\begin{aligned} \mathcal{G}_{\|\cdot\|, Q_n, G} = \left\{ g : \max \left( \|g(v)\|^*, \|\nabla g(v)\|^*, \frac{\|g(x) - g(y)\|^*}{\|x - y\|}, \right. \right. \\ \left. \frac{\|\nabla g(x) - \nabla g(y)\|^*}{\|x - y\|} \right) \leq 1, \\ \frac{\|g(x) - g(y) - \nabla g(x)(x - y)\|^*}{\frac{1}{2}\|x - y\|^2} \leq 1, \\ \left. \frac{\|g(x) - g(y) - \nabla g(y)(x - y)\|^*}{\frac{1}{2}\|x - y\|^2} \leq 1, \forall (x, y) \in E, v \in V \right\}, \end{aligned}$$

imposes boundedness constraints only at sample points and smoothness constraints only at pairs of sample points enumerated in the edge set  $E$ . The graph is termed a *t-spanner* [14, 76] if each edge  $(x, y) \in E$  is assigned the weight  $\|x - y\|$ , and,

---

<sup>5</sup>When  $d = 1$ , the problem reduces to a finite-dimensional convex quadratically constrained quadratic program with linear objective as in [36], Theorem 9.

for all  $x' \neq y' \in V$ , there exists a path between  $x'$  and  $y'$  in the graph with total path weight no greater than  $t\|x' - y'\|$ . Remarkably, for any linear Stein operator  $\mathcal{T}$ , a *spanner Stein discrepancy*  $\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|, Q_n, G_t})$  based on a  $t$ -spanner  $G_t$  is equivalent to the classical Stein discrepancy in the following strong sense, implying Desiderata (i) and (ii).

**PROPOSITION 13** (Equivalence of classical and spanner Stein discrepancies). *If  $G_t = (\text{supp}(Q_n), E)$  is a  $t$ -spanner for  $t \geq 1$ , then*

$$\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|}) \leq \mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|, Q_n, G_t}) \leq \kappa_d t^2 \mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|})$$

where  $\kappa_d$  is independent of  $(Q_n, P, \mathcal{T}, G_t)$  and depends only on  $d$  and  $\|\cdot\|$ .

**REMARK.** The proof relies on the Whitney–Glaeser extension theorem [86], Theorem 1.4, of Glaeser [34] and follows exactly as in [36], Propositions 5 and 6.

When  $d = 1$ , a  $t$ -spanner with exactly  $n - 1$  edges is obtained in  $O(n \log n)$  time for all  $t \geq 1$  by introducing edges just between sample points that are adjacent in sorted order. More generally, if  $\|\cdot\|$  is an  $\ell^p$  norm, one can construct a 2-spanner with  $O(\kappa'_d n)$  edges in  $O(\kappa'_d n \log(n))$  expected time where  $\kappa'_d$  is a constant that depends only on the norm  $\|\cdot\|$  and the dimension  $d$  [43]. Hence, a spanner Stein discrepancy can be computed by solving a finite-dimensional convex optimization problem with a linear objective,  $O(n)$  variables and  $O(\kappa'_d n)$  convex constraints, making it an appealing choice for a computable quality measure (Desideratum (iii)).

**4.2. Decoupled linear programs.** Moreover, if we choose the norm  $\|\cdot\| = \|\cdot\|_1$ , the graph Stein discrepancy optimization problem decouples into  $d$  independent linear programs (LPs) that can be solved in parallel using off-the-shelf solvers. Indeed, for any  $G = (\text{supp}(Q_n), E)$ ,  $\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|_1, Q_n, G})$  equals

$$\begin{aligned} & \sum_{j=1}^d \sup_{\psi_j \in \mathbb{R}^n, \Psi_j \in \mathbb{R}^{d \times n}} \sum_{i=1}^n q(x_i) \left( 2b_j(x_i) \psi_{ji} + \sum_{k=1}^d m_{jk}(x_i) \Psi_{jki} \right) \\ & \text{s.t. } \|\psi_j\|_\infty \leq 1, \|\Psi_j\|_\infty \leq 1, \quad \text{and for all } i \neq l, (x_i, x_l) \in E \\ (18) \quad & \max \left( \frac{|\psi_{ji} - \psi_{jl}|}{\|x_i - x_l\|_1}, \frac{\|\Psi_j(e_i - e_k)\|_\infty}{\|x_i - x_l\|_1}, \frac{|\psi_{ji} - \psi_{jl} - \langle \Psi_j e_i, x_i - x_l \rangle|}{\frac{1}{2} \|x_i - x_l\|_1^2}, \right. \\ & \left. \frac{|\psi_{ji} - \psi_{jl} - \langle \Psi_j e_i, x_l - x_i \rangle|}{\frac{1}{2} \|x_i - x_l\|_1^2} \right) \leq 1, \end{aligned}$$

where  $\psi_{ji}$  and  $\Psi_{jki}$  represent the values  $g_j(x_i)$  and  $\nabla_k g_j(x_i)$  respectively. Therefore, our recommended quality measure is the 2-spanner diffusion Stein discrepancy with  $\|\cdot\| = \|\cdot\|_1$ . Its computation is summarized in Algorithm 1. An efficient

**Algorithm 1** Spanner diffusion Stein discrepancy,  $\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|_1, Q_n, G_2})$ 


---

**input:** sample  $Q_n$ , target score  $\nabla \log p$ , covariance coefficient  $a$ , stream coefficient  $c$   
 $G_2 \leftarrow$  2-spanner of  $V = \text{supp}(Q_n)$   
**for**  $j = 1$  **to**  $d$  **do (in parallel)**  
     $\tau_j \leftarrow$  Optimal value of  $j$ th coordinate linear program (18) with graph  $G_2$   
**return**  $\sum_{j=1}^d \tau_j$

---

implementation of Algorithm 1, integrated with 11 linear program solver options, is publicly available via our Julia package.<sup>6</sup>

**5. Numerical illustrations.** In this section, we complement the principal theoretical contributions of this work with several simple numerical illustrations demonstrating how diffusion Stein discrepancies can be deployed in practice. We will use our proposed quality measures to select hyperparameters for biased samplers, to quantify a bias-variance trade-off for approximate MCMC, and to compare deterministic and random quadrature rules. In each case, we choose experimental settings in which a notion of surrogate ground truth is available for external validation. We solve all linear programs using Julia for Mathematical Programming [62] with the Gurobi 6.0.4 solver [40] and use the C++ greedy spanner implementation of Bouts et al. [5] to compute our 2-spanners. Our timings were obtained on a single core of an Intel Xeon CPU E5-2650 v2 @ 2.60 GHz. Code reconstructing all experiments is available on the Julia package site.<sup>5</sup>

5.1. *A simple example.* We first present a simple example to illustrate several Stein discrepancy properties. For a Gaussian mixture target  $P$  (Example 3) with  $p(x) \propto e^{-\frac{1}{2}(x-\frac{\Delta}{2})^2} + e^{-\frac{1}{2}(x+\frac{\Delta}{2})^2}$  and  $\Delta > 0$ , we simulate one i.i.d. sequence of sample points from  $P$  and a second i.i.d. sequence from  $\mathcal{N}(-\frac{\Delta}{2}, 1)$ , which represents only one component of  $P$ . For various mode separations  $\Delta$ , Figure 1 shows that the Langevin spanner Stein discrepancy (D1) applied to the first  $n$  Gaussian mixture sample points decreases to zero at a  $n^{-1/2}$  rate, while the discrepancy applied to the single mode sequence stays bounded away from zero. However, Figure 1 also indicates that larger sample sizes are needed to distinguish between the mixture and single mode sample sequences when  $\Delta$  is large. This accords with our theory (see Example 3, Corollary 12 and Theorem 6), which implies that both the Langevin diffusion Wasserstein decay rate and the bound relating Stein to Wasserstein degrade as the mixture mode separation  $\Delta$  increases.

---

<sup>6</sup><https://jgorham.github.io/SteinDiscrepancy.jl/>

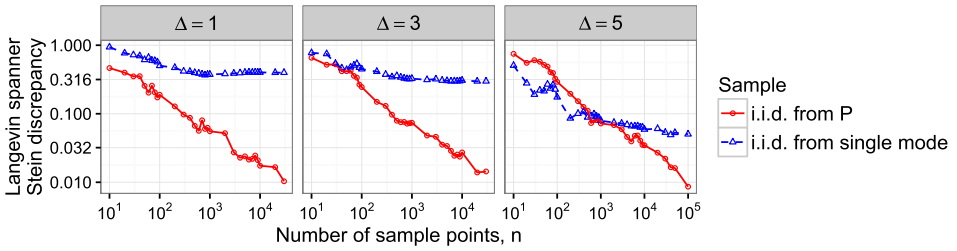


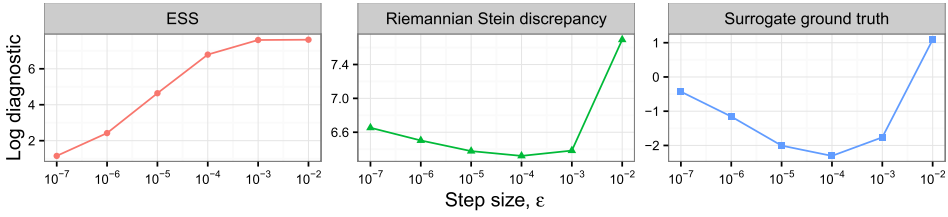
FIG. 1. Stein discrepancy for normal mixture target  $P$  with  $\Delta$  mode separation (Section 5.1).

5.2. *Selecting sampler hyperparameters.* Stochastic Gradient Riemannian Langevin Dynamics (SGRLD) [74] with a constant step size  $\epsilon$  is an approximate MCMC procedure designed to accelerate posterior inference. Unlike asymptotically correct MCMC algorithms, SGRLD has a stationary distribution that deviates increasingly from its target  $P$  as its step size  $\epsilon$  grows. On the other hand, if  $\epsilon$  is too small, SGRLD fails to explore the sample space sufficiently quickly. Hence, an appropriate setting of  $\epsilon$  is paramount for accurate inference.

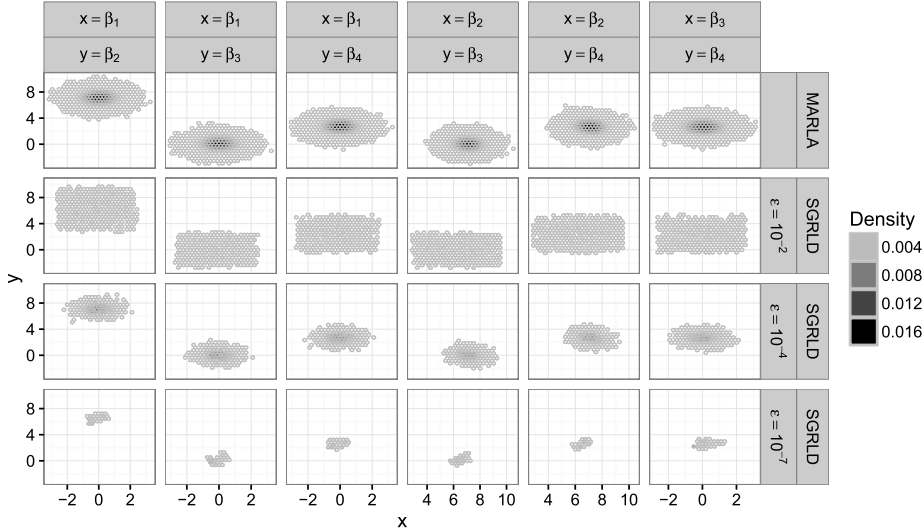
To demonstrate the value of diffusion Stein discrepancies for hyperparameter selection, we analyzed a biometric data set of  $L = 202$  athletes from the Australian Institute of Sport that was previously the focus of a heavy-tailed regression analysis [87]. In the notation of Example 4, we used SGRLD to conduct a Bayesian multivariate Student's  $t$  regression ( $\nu = 10$ ,  $\Sigma = I$ ) of athlete lean body mass onto red blood count, white blood count, plasma ferritin concentration and a constant regressor of value  $1/\sqrt{L}$  with a pseudo-Huber prior ( $\delta = 0.1$ ) on the unknown parameter vector  $\beta \in \mathbb{R}^4$ .

After standardizing the output variable and nonconstant regressors and initializing each chain with an approximate posterior mode found by L-BFGS started at the origin, we ran SGRLD with minibatch size 30, metric  $G(\beta) = 1/(2\sqrt{1 + \|\beta/\delta\|_2^2})I$ , and a variety of step sizes  $\epsilon$  to produce sample sequences of length 200,000 thinned to length 2000. We then selected the step size that delivered the highest quality sample—either the maximum effective sample size (ESS, a popular MCMC mixing diagnostic based on asymptotic variance [6]) or the minimum Riemannian Langevin spanner Stein discrepancy with  $a(\beta) = G^{-1}(\beta)$ . The longest discrepancy computation consumed 6s for spanner construction and 65s to solve a coordinate optimization problem. As a surrogate measure of ground truth, we also generated a sample  $Q^*$  of size  $2 \times 10^8$  from the Metropolis-adjusted Riemannian Langevin Algorithm (MARLA) [33] with metric  $G$  and compute the median bivariate marginal Wasserstein distance  $d_{\mathcal{W}_{\|\cdot\|_1}}$  between each SGRLD sample and  $Q^*$  thinned to 5000 points [41].

Figure 2(a) shows that ESS, which does not account for stationary distribution bias, selects the largest step size available,  $\epsilon = 10^{-2}$ . As seen in Figure 2(b), this choice results in samples that are greatly overdispersed when compared with the



(a) Step size selection criteria and surrogate ground truth (median marginal Wasserstein). ESS maximized at  $\epsilon = 10^{-2}$ . Stein discrepancy and ground truth minimized at  $\epsilon = 10^{-4}$ .



(b) Bivariate hexbin plots. **Top row:** surrogate ground truth sample ( $2 \times 10^8$  MARLA points). **Bottom 3 rows:** 2,000 SGRLD sample points for various step sizes  $\epsilon$ .

FIG. 2. Step size selection, stochastic gradient Riemannian Langevin dynamics (Section 5.2).

ground truth MARLA sample  $Q^*$ . At the other extreme, the selection  $\epsilon = 10^{-7}$  produces greatly underdispersed samples due to slow mixing. The Stein discrepancy chooses an intermediate value,  $\epsilon = 10^{-4}$ . The same value minimizes the surrogate ground truth Wasserstein measure and produces samples that most closely resemble the  $Q^*$  in Figure 2(b).

5.3. *Quantifying a bias-variance trade-off.* Approximate random walk Metropolis–Hastings (ARWMH) [54] with tolerance parameter  $\epsilon$  is a biased MCMC procedure that accelerates posterior inference by approximating the standard MH correction. Qualitatively, a smaller setting of  $\epsilon$  produces a more faithful approximation of the MH correction and less bias between the chain’s stationary distribution and the target distribution of interest. A larger setting of  $\epsilon$  leads to faster sampling and a more rapid reduction of Monte Carlo variance, as fewer datapoint likelihoods are computed per sampling step. We will quantify this bias-

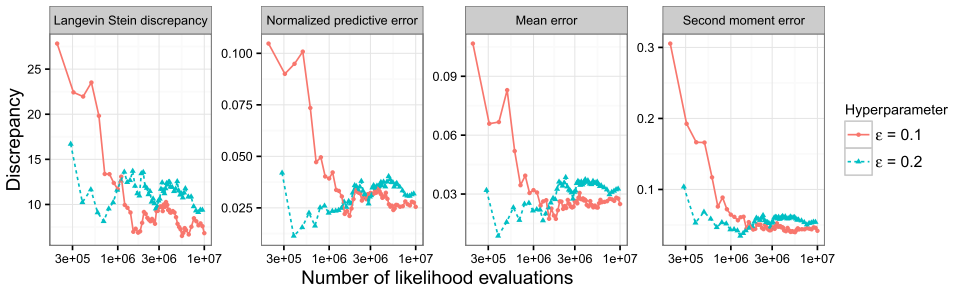


FIG. 3. Bias-variance trade-off curves for approximate random walk MH (Section 5.3).

variance trade-off as a function of sampling time using the Langevin spanner Stein discrepancy.

In the notation of Example 2, we conduct a Bayesian Huber regression analysis ( $c = 1$ ) of the log radon levels in 1190 Minnesota households [31] as a function of the log amount of uranium in the county, an indicator of whether the radon reading was performed in a basement and an intercept term. A  $\mathcal{N}(0, I)$  prior is placed on the coefficient vector  $\beta$ . We run ARWMH with minibatch size 5 and two settings of the tolerance threshold  $\epsilon$  (0.1 and 0.2) for  $10^7$  likelihood evaluations, discard the sample points from the first  $10^5$  evaluations, and thin the remaining points to sequences of length 1000. Figure 3 displays the Langevin spanner Stein discrepancy applied to the first  $n$  points in each sequence as a function of the likelihood evaluation count, which serves as a proxy for sampling time. As expected, the higher tolerance sample ( $\epsilon = 0.2$ ) is of higher Stein quality for a small computational budget but is eventually overtaken by the  $\epsilon = 0.1$  sample with smaller asymptotic bias. The longest discrepancy computation consumed 0.8s for the spanner and 20.1s for a coordinate LP.

To provide external support for the Stein discrepancy quantification, we generate a Metropolis-adjusted Langevin chain [84] of length  $10^8$  as a surrogate  $Q^*$  for the target  $P$  and display several measures of expectation error between  $X \sim Q_n$  and  $Z \sim Q^*$  in Figure 3: the normalized predictive error  $\max_l |\mathbb{E}[\langle X - Z, v_l / \|v_l\|_\infty \rangle]|$  for  $v_l$  the  $l$ th datapoint covariate vector, the mean error  $\frac{\max_j |\mathbb{E}[X_j - Z_j]|}{\max_j |\mathbb{E}_{Q^*}[Z_j]|}$  and the second moment error  $\frac{\max_{j,k} |\mathbb{E}[X_j X_k - Z_j Z_k]|}{\max_{j,k} |\mathbb{E}_{Q^*}[Z_j Z_k]|}$ . We see that the Stein discrepancy provides comparable results without the need for an additional surrogate chain.

5.4. Comparing quadrature rules. Stein discrepancies can also measure the quality of deterministic sample sequences designed to improve upon Monte Carlo sampling. For the Gaussian mixture target of Section 5.1, Figure 4 compares the median quality of 50 sample sequences generated from four quadrature rules recently studied in [55], Section 4.1: i.i.d. sampling from  $P$ , Quasi-Monte Carlo

(QMC) sampling using a deterministic quasirandom number generator, Frank–Wolfe (FW) kernel herding [2, 13] and fully-corrective Frank–Wolfe (FCFW) kernel herding [55]. The quality judgments of the Langevin spanner Stein discrepancy (D1) closely mimic those of the  $L^1$  Wasserstein distance  $d_{\mathcal{W}_{||\cdot||}}$ , which is computable for simple univariate targets [92]. Each Stein discrepancy was computed in under 0.03s.

Under both diagnostics and as previously observed in other metrics [55], the i.i.d. samples are typically of lower median quality than their deterministic counterparts. More surprisingly and in contrast to past work focused on very smooth function classes [55], FCFW underperforms FW and QMC in our diagnostics for larger sample sizes. Apparently FCFW, which is heavily optimized for smooth function integration, has sacrificed approximation quality for less smooth test functions. For example, Figure 4 shows that QMC offers a better quadrature estimate than FCFW for  $h_1(x) = \max\{0, 1 - \min_{j \in \{1,2\}} |x - \mu_j|\}$ , a 1-Lipschitz approximation to the indicator of being within one standard deviation of a mode.

In addition to providing a sample quality score, the Stein discrepancy optimization problem produces an optimal Stein function  $g^*$  and an associated test function  $h^* = \mathcal{T}g^*$  that is mean zero under  $P$  and best distinguishes the sample  $Q_n$  from the target  $P$ . Figure 4 gives examples of these maximally discriminative functions  $h^*$  for a target mode separation of  $\Delta = 5$  and length 200 sequences from each

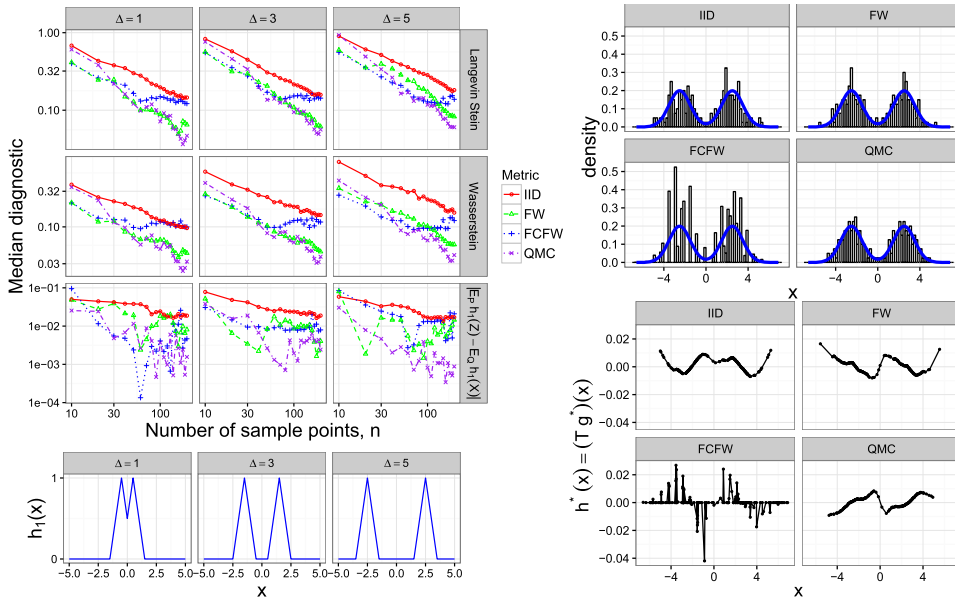


FIG. 4. Left: Quadrature rule quality comparison for Gaussian mixture targets  $P$  with mode separation  $\Delta$  (Section 5.4). Right: (Top) Sample histograms with  $p$  overlaid ( $\Delta = 5, n = 200$ ). (Bottom) Optimal discriminating test functions  $h^* = \mathcal{T}g^*$  from Stein program.



quadrature rule. We also display the associated sample histograms with overlaid target density. The optimal FCFW function reflects the jagged nature of the FCFW histogram.

**6. Connections and conclusions.** We developed quality measures suitable for comparing the fidelity of arbitrary “off-target” sample sequences by generating infinite collections of known target expectations.

*Alternative quality measures.* The score statistic of Fan et al. [25] and the Gibbs sampler convergence criteria of Zellner and Min [96] account for some sample biases but sacrifice differentiating power by exploiting only a finite number of known target expectations. For example, when  $P = \mathcal{N}(0, 1)$ , the score statistic [25] cannot differentiate two samples with the same means and variances. Maximum mean discrepancies (MMDs) over characteristic reproducing kernel Hilbert spaces [39] do detect arbitrary distributional biases but are only computable when the chosen kernel functions can be integrated under the target. In practice, one often approximates MMD using a sample from the target, but this requires a separate trustworthy sample from  $P$ .

While we have focused on the graph and classical Stein sets of [36], our diffusion Stein operators can also be paired with the reproducing kernel Hilbert space unit balls advocated in [15, 37, 60, 71] to form tractable *kernel diffusion Stein discrepancies* or with the random feature functions advocated in [48] to form *random feature diffusion Stein discrepancies*. We have also restricted our attention to Stein operators arising from diffusion generators. These take the form  $(\mathcal{T}g)(x) = \frac{1}{p(x)} \langle \nabla, p(x)m(x)g(x) \rangle$  with  $m = a + c$  for  $a(x)$  positive semidefinite and  $c(x)$  skew-symmetric. More generally, if the matrix  $m$  possesses eigenvalues having a negative real part, then the resulting operator need not correspond to a diffusion process. Such operators fall into the class of *pseudo-Fokker-Planck* operators which have been studied in the context of quantum optics [82]. As noted in [18, 19], it is possible to obtain corresponding stochastic dynamics in an extended state space by introducing complex-valued noise terms; these operators may merit further study in future work.

*Alternative inferential tasks.* While our chief motivation is sample quality measurement, our work is also directly applicable to a variety of inferential tasks that currently rely on the Langevin operator introduced by [36, 71], including control variate design [71], goodness-of-fit testing [15, 60], variational inference [12, 61, 79] and importance sampling [59]. The Stein factor bounds of Theorem 5 can also be used, in the manner of [42, 50, 67], to characterize the error of numerical discretizations of diffusions. These works convert bounds on the solutions of Poisson equations—Stein factors—into central limit theorems for  $\mathbb{E}_{Q_n}[h(X)] - \mathbb{E}_P[h(Z)]$ , confidence intervals for  $\mathbb{E}_P[h(Z)]$ , and mean-squared

error bounds for the estimate  $\mathbb{E}_{Q_n}[h(X)]$ . Teh et al. [91] and Vollmer et al. [93] extended these approaches to obtain error estimates for approximate discretizations of the Langevin diffusion on  $\mathbb{R}^d$ , while, independently of our work, Huggins and Zou [47] established error estimates for Itô diffusion approximations with biased drifts and constant diffusion coefficients. By Theorem 5, their results also hold for Itô diffusions with nonconstant diffusion coefficients. Following the release of the present paper and with the aim of analyzing discretization error for the overdamped Langevin diffusion, Fang et al. [26], Theorem 3.1, derived multivariate Stein factor bounds for a class of strongly log-concave distributions. Our Theorem 5 with the choice  $\iota = 1/\log(1/\epsilon)$  provides Stein factors of the same form but applies also to nonlog-concave targets and more general diffusions.

*Alternative targets.* Our exposition has focused on the Wasserstein distance  $d_{\mathcal{W}_{\|\cdot\|}}$ , which is only defined for distributions with finite means. A parallel development could be made for the Dudley metric [69] to target distributions with undefined mean. The work of Cerrai [8] also suggests that the Lipschitz condition on our drift and diffusion coefficients can be relaxed.

APPENDIX A: PROOF OF PROPOSITION 3

Fix any  $g \in \mathcal{G}_{\|\cdot\|}$ . Since  $g$  and  $\nabla g$  are bounded and  $b, a$ , and  $c$  are  $P$ -integrable,  $\mathbb{E}_P[(\mathcal{T}g)(Z)]$  is finite. Define the ball  $\mathcal{B}_r = \{x \in \mathbb{R}^d : \|x\|_2 \leq r\}$  with  $n_r(z)$  the outward facing unit normal vector for each  $z$  on the boundary  $\partial\mathcal{B}_r$ . Since  $z \mapsto p(z)(a(z) + c(z))g(z)$  is in  $C^1$ , we may apply the dominated convergence theorem and then the divergence theorem to obtain

$$\begin{aligned} \mathbb{E}_P[(\mathcal{T}g)(Z)] &= \lim_{r \rightarrow \infty} \int_{\mathcal{B}_r} \langle \nabla, p(z)(a(z) + c(z))g(z) \rangle dz \\ &= \lim_{r \rightarrow \infty} \int_{\partial\mathcal{B}_r} \langle n_r(z), (a(z) + c(z))g(z)p(z) \rangle dz. \end{aligned}$$

Let  $f(r) = M_0(g) \int_{\partial\mathcal{B}_r} \|a(z) + c(z)\|_{\text{op}} p(z) dz$ . Since  $g$  and  $n_r$  are bounded,

$$\int_{\partial\mathcal{B}_r} \langle n_r(z), (a(z) + c(z))g(z)p(z) \rangle dz \leq f(r).$$

The coarea formula [1] and the integrability of  $a$  and  $c$  further imply that

$$\int_0^\infty f(r) dr = \int_{\mathbb{R}^d} M_0(g) \|a(z) + c(z)\|_{\text{op}} p(z) dz < \infty.$$

Hence,  $\liminf_{r \rightarrow \infty} f(r) = 0$ , and therefore  $\mathbb{E}_P[(\mathcal{T}g)(Z)] = 0$ .

APPENDIX B: PROOF OF THEOREM 5

Fix any  $x \in \mathbb{R}^d$  and  $h \in \mathcal{W}_{\|\cdot\|_2}$  with  $\mathbb{E}_P[h(Z)] = 0$ . Since the drift and diffusion coefficients are Lipschitz, [53], Theorem 3.4, guarantees that the diffusion

$(Z_{t,x})_{t \geq 0}$  is well defined. Using the shorthand  $s_r \triangleq \int_0^\infty r(t) dt$ , we will show that the posited function  $u_h$  (10) exists and solves the *Poisson equation*

$$(19) \quad h = \mathcal{A}u_h$$

with infinitesimal generator  $\mathcal{A}$ , that  $u_h$  is Lipschitz, that  $u_h$  has a continuous Hessian, that  $u_h$  has a bounded and Hölder continuous Hessian under additional smoothness assumptions.

**Existence of  $u_h$  and solving the Poisson equation (19).** Consider the set  $L \triangleq (1 + \|x\|_2^2)C_0(\mathbb{R}^d) = \{(1 + \|x\|_2^2)f : f \in C_0(\mathbb{R}^d)\}$ , where  $C_0(\mathbb{R}^d)$  is the set of continuous functions vanishing at infinity. Equipped with the norm  $\|f\|_L = \sup_{x \in \mathbb{R}^d} |f(x)|/(1 + \|x\|_2^2)$ , the set  $L$  is a Banach space [85]. As noted in [17], the space  $L$  can also be characterized as the closure of the set of bounded continuous functions,  $C_b(\mathbb{R}^d)$ , in the set  $\{f : \mathbb{R}^d \rightarrow \mathbb{R} : \|f\|_L < \infty\}$ . To discuss the well-posedness of the Poisson equation (19), we first show that the transition semigroup of an Itô diffusion is strongly continuous on  $L$ .

**PROPOSITION 14.** *The transition semigroup  $(P_t)_{t \geq 0}$  of an Itô diffusion with Lipschitz drift and diffusion coefficients is strongly continuous on  $L$ .*

**PROOF.** Fix any  $f \in L$  and  $x \in \mathbb{R}^d$ . We first show that  $(P_t f)(x)$  converges pointwise to  $f(x)$  as  $t \rightarrow 0^+$ . Since the associated Itô process  $(Z_{t,x})_{t \geq 0}$  is almost surely pathwise continuous [53], Theorem 3.4, and  $f$  is continuous in a neighborhood of  $x$ , it follows that  $f(Z_{t,x}) \rightarrow f(x)$  as  $t \rightarrow 0^+$ , almost surely. Moreover, [28], Section 5, Corollary 1.2, implies that

$$\mathbb{E} \left[ \sup_{0 \leq t \leq 1} |f(Z_{t,x})| \right] \leq \|f\|_L \left( 1 + \mathbb{E} \left[ \sup_{0 \leq t \leq 1} \|Z_{t,x}\|_2^2 \right] \right) \leq C \|f\|_L (1 + \|x\|_2^2),$$

for some  $C > 0$  depending only on  $M_1(b)$  and  $M_1(\sigma)$ . The dominated convergence theorem now yields the desired pointwise convergence.

To prove the strong continuity of  $(P_t)_{t \geq 0}$ , it suffices, by [23], Theorem I.5.8, p. 40, to verify that  $(P_t)_{t \geq 0}$  is weakly continuous, that is, that  $l(P_t f) \rightarrow l(f)$ , as  $t \rightarrow 0^+$ , for all elements  $l$  of the dual space  $L^*$ . To this end, fix any  $l \in L^*$ . By the Riesz–Markov theorem for  $L$  [17], Theorem 2.4, there exists a finite signed Radon measure  $\mu$  such that

$$(20) \quad l(f) = \int_{\mathbb{R}^d} f(x) \mu(dx) \quad \text{and} \quad \int_{\mathbb{R}^d} (1 + \|x\|_2^2) |\mu|(dx) = \|l\|_{L^*},$$

for  $\|\cdot\|_{L^*}$  the dual norm. By Jensen’s inequality and [28], Section 5, Corollary 1.2,

$$\begin{aligned} \forall t, \quad \|(P_t f)(x)\|_2 &\leq \mathbb{E}[|f(Z_{t,x})|] \\ &\leq \|f\|_L \mathbb{E}[1 + \|Z_{t,x}\|_2^2] \leq C \|f\|_L (1 + \|x\|_2^2). \end{aligned}$$

Since  $1 + \|x\|_2^2$  is  $|\mu|$ -integrable by (20), dominated convergence gives

$$\lim_{t \rightarrow 0^+} l(P_t f) = \lim_{t \rightarrow 0^+} \int_{\mathbb{R}^d} (P_t f)(x) \mu(dx) = \int_{\mathbb{R}^d} f(x) \mu(dx) = l(f),$$

yielding the result.  $\square$

Consider the infinitesimal generator  $\mathcal{A}$  of the semigroup  $(P_t)_{t \geq 0}$  on  $L$  with

$$\text{dom}(\mathcal{A}) = \left\{ f \in L : \lim_{t \rightarrow 0^+} \frac{P_t f - f}{t} \text{ exists in the } \|\cdot\|_L \text{ norm} \right\}.$$

Since  $P_t$  is strongly continuous on  $L$  and  $h \in L$  with  $M_1(h) \leq 1$  and  $\mathbb{E}_P[h(Z)] = 0$ , [24], Proposition 1.5, implies that

$$h - P_t h = -\mathcal{A} \int_0^t P_s h ds = \mathcal{A} u_{h,t} \quad \text{for } u_{h,t} \triangleq - \int_0^t P_s h ds.$$

The stationarity of  $P$  and the definitions of  $d_{\mathcal{W}_{\|\cdot\|_2}}$  and  $r$  imply that

$$\begin{aligned} \|P_t h\|_L &= \|P_t h - E_P[h]\|_L \\ &= \sup_{x \in \mathbb{R}^d} |\mathbb{E}_P[P_t h(x) - P_t h(Z)]| / (1 + \|x\|_2^2) \\ &\leq \sup_{x \in \mathbb{R}^d} \frac{\mathbb{E}_P[d_{\mathcal{W}_{\|\cdot\|_2}}(\delta_x P_t, \delta_Z P_t)]}{1 + \|x\|_2^2} \\ &\leq r(t) \sup_{x \in \mathbb{R}^d} \frac{\mathbb{E}_P[\|x - Z\|_2]}{1 + \|x\|_2^2}, \end{aligned}$$

and hence  $\|P_t h\|_L \rightarrow 0$  as  $t \rightarrow \infty$ , since  $P$  has a finite mean, and  $r(t) \rightarrow 0$  as  $t \rightarrow \infty$  as  $r$  is integrable and monotonic. Arguing similarly,

$$\begin{aligned} \|u_{h,t} - u_{h,t'}\|_L &\leq \left\| \int_t^{t'} \mathbb{E}_P[d_{\mathcal{W}_{\|\cdot\|_2}}(\delta_x P_s, \delta_Z P_s)] ds \right\|_L \\ &\leq \sup_{x \in \mathbb{R}^d} \frac{\mathbb{E}_P[\|x - Z\|_2]}{1 + \|x\|_2^2} \int_t^{t'} r(s) ds. \end{aligned}$$

Thus, it follows that  $(u_{h,t})_{t > 0}$  is a Cauchy sequence in  $L$  with limit  $u_h = \int_0^\infty P_s h ds \in L$ . Thus,  $(h - P_t h, u_{h,t}) \rightarrow (h, u_h)$  in the graph norm on  $L \times L$ , and since  $\mathcal{A}$  is closed [24], Corollary 1.6,  $u_h \in \text{dom}(\mathcal{A})$  and  $h = \mathcal{A}u_h$ .

REMARK. The choice of the Banach space is crucial for the argument above. As noted in [66] and contrary to the claim in [4], the semigroup  $(P_t)_{t \geq 0}$  fails to be strongly continuous over the Banach space  $\tilde{L} \triangleq (1 + \|x\|_2^2)C_b(\mathbb{R}^d)$  when  $(Z_{t,x})_{t \geq 0}$  is an Ornstein–Uhlenbeck process, that is, a Langevin diffusion (D1) with a multivariate Gaussian invariant measure.

**Lipschitz continuity of  $u_h$ .** To demonstrate that  $u_h$  is Lipschitz, we choose an arbitrary  $v \in \mathbb{R}^d$ , and apply the definition of the Wasserstein distance, the assumed decay rate, and the integrability of  $r$  to obtain

$$\begin{aligned} \|u_h(x + v) - u_h(x)\|_2 &\leq \int_0^\infty \|\mathbb{E}[h(Z_{t,x}) - h(Z_{t,x+v})]\|_2 dt \\ &\leq \int_0^\infty d_{\mathcal{W}_{\|\cdot\|}}(\delta_x P_t, \delta_{x+v} P_t) dt \\ &\leq d_{\mathcal{W}_{\|\cdot\|}}(\delta_x, \delta_{x+v}) s_r = \|v\|_2 s_r < \infty. \end{aligned}$$

**Continuity of  $\nabla^2 u_h$ .** Since  $u_h \in \text{dom}(\mathcal{A})$  is a continuous solution of the Poisson equation (19), and since the infinitesimal generator agrees with the characteristic operator of a diffusion when both are defined [72], p. 129, Theorem 5.9 of [21] implies that  $u_h \in C^2$ .

**Boundedness of  $\nabla^2 u_h$ .** Instantiate the additional preconditions of (11), and assume that  $M_0(\sigma^{-1}), F_2(\sigma), M_2(b) < \infty$ , or else (11) is vacuous. Lemma 15, established in Appendix C, shows that the semigroup  $P_t h$  admits a bounded continuous Hessian, which is integrable in  $t$ .

LEMMA 15 (Semigroup Hessian estimate). *Suppose that the drift and diffusion coefficients  $b$  and  $\sigma$  of an Itô diffusion are Lipschitz with Lipschitz gradients and locally Lipschitz second derivatives. If the transition semigroup  $(P_t)_{t \geq 0}$  has Wasserstein decay rate  $r$ , and  $\sigma(x)$  has a right inverse  $\sigma^{-1}(x)$  for each  $x \in \mathbb{R}^d$ , then, for all  $t > 0$  and any  $f \in C^2$  with bounded first and second derivatives,  $P_t f$  is twice continuously differentiable with*

$$(21) \quad M_1(P_t f) \leq M_1(f)r(t) \quad \text{and}$$

$$\begin{aligned} (22) \quad M_2(P_t f) &\leq \inf_{t_0 \in (0,t]} M_1(f)r(t - t_0) \sqrt{\frac{1}{t_0}} e^{t_0 \gamma_2} M_0(\sigma^{-1}) \\ &\quad + M_1(f)r(t - t_0)r(0)e^{t_0 \gamma_2} M_1(\sigma)M_0(\sigma^{-1}) \\ &\quad + M_1(f)r(t - t_0)\sqrt{t_0}r(0)e^{t_0 \gamma_4} \frac{2}{3}\sqrt{\alpha} \end{aligned}$$

for  $\gamma_\rho \triangleq \rho M_1(b) + \frac{\rho^2 - 2\rho}{2} M_1(\sigma)^2 + \frac{\rho}{2} F_1(\sigma)^2, \alpha \triangleq \frac{M_2(b)^2}{2M_1(b) + 4M_1(\sigma)^2} + 2F_2(\sigma)^2$ .

The dominated convergence theorem now implies that the Hessian of  $u_h$  is obtained by differentiating twice under the integral sign. The advertised bound (11) on  $\nabla^2 u_h$  follows by replacing the infimum on the right-hand side of the semigroup bound (22) with the selection  $t_0 = \min(t, 1)$ , applying the bound  $e^{\min(t,1)\gamma_\rho} \leq e^{\gamma_\rho}$  for each  $\gamma_\rho$  and  $t$ , and integrating the result over  $t$ .

**Hölder continuity of  $\nabla^2 u_h$ .** Finally, instantiate the additional preconditions of (12), and fix any  $\iota \in (0, 1)$ . The integral representation (10) of  $u_h$ , the dominated convergence theorem, and Jensen’s inequality imply

$$M_{1-\iota}(\nabla^2 u_h) = M_{1-\iota}\left(-\int_0^\infty \nabla^2 P_t h dt\right) \leq \int_0^\infty M_{1-\iota}(\nabla^2 P_t h) dt.$$

When  $t \leq 1$ , a seminorm interpolation lemma (Lemma 19 in the Supplementary Material [35], which is based on results on seminorm interpolation, see e.g. [63]), a semigroup third derivative estimate (Lemma 20 in the Supplementary Material [35]) with  $t_0 = \min(t, 1)$  and the semigroup second derivative estimate of Lemma 15 with  $t_0 = \min(t, 1)$  imply

$$M_{1-\iota}(\nabla^2 P_t h) \leq M_1(h)2^\iota M_0(\nabla^2 P_t h)^\iota M_1(\nabla^2 P_t h)^{1-\iota} \leq M_1(h)t^{\iota/2-1}/K_1$$

for some constant  $K_1 > 0$  depending only on  $M_{1:3}(b)$ ,  $M_{1:3}(\sigma)$ ,  $M_0(\sigma^{-1})$ , and  $r$ . Thus  $\int_0^1 M_{1-\iota}(\nabla^2 P_t h) dt \leq \frac{2M_1(h)}{K_1^\iota}$ . For  $t > 1$ , Lemmas 19, 20 and 15 and the integrability of  $r$  yield

$$\int_1^\infty M_{1-\iota}(\nabla^2 P_t h) dt \leq M_1(h)\frac{2}{K_2} \int_1^\infty r(t-1) dt = M_1(h)\frac{2}{K_2}s_r$$

for a constant  $K_2 > 0$  again depending only on  $M_{1:3}(b)$ ,  $M_{1:3}(\sigma)$ ,  $M_0(\sigma^{-1})$ , and  $r$ . Combining these bounds and choosing  $K = \min(K_1, K_2)/2$  completes the proof. An explicit constant  $K$  can be obtained by tracing constants through the proof of Lemma 20.

APPENDIX C: PROOF OF LEMMA 15

Fix any  $x \in \mathbb{R}^d$  and  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  in  $C^2$  with bounded first and second derivatives, and let  $(Z_{t,x})_{t \geq 0}$  be an Itô diffusion solving the stochastic differential equation (5) with starting point  $Z_{0,x} = x$ , underlying Wiener process  $(W_t)_{t \geq 0}$ , and transition semigroup  $(P_t)_{t \geq 0}$ . Our proof is divided into five pieces establishing, for each  $t > 0$ , the Lipschitz continuity of  $P_t f$ , the Lipschitz continuity of  $\nabla P_t f$ , the continuity of  $\nabla^2 P_t f$ , an initial bound on  $\nabla^2 P_t f$ , and the infimal bound (22) on  $\nabla^2 P_t f$ .

**Lipschitz continuity of  $P_t f$ .** The semigroup gradient bound (21) follows from the Lipschitz continuity of  $f$  and the definitions of the Wasserstein decay rate and the Wasserstein distance, as, for any  $y \in \mathbb{R}^d$  and  $t \geq 0$ ,

$$\begin{aligned} (P_t f)(x) - (P_t f)(y) &= \mathbb{E}[f(Z_{t,x}) - f(Z_{t,y})] \\ &\leq M_1(f)d_{\mathcal{W}_{\|\cdot\|_2}}(\delta_x P_t, \delta_y P_t) \\ &\leq M_1(f)r(t)d_{\mathcal{W}_{\|\cdot\|_2}}(\delta_x, \delta_y) \\ &= M_1(f)r(t)\|x - y\|_2. \end{aligned}$$

**Lipschitz continuity of  $\nabla P_t f$ .** Fix any  $v, v' \in \mathbb{R}^d$ . Under our smoothness assumptions on  $b$  and  $\sigma$ , [77], Theorem V.40, implies that  $(Z_{t,x})_{t \geq 0}$  is twice continuously differentiable in  $x$ . The first directional derivative flow  $(V_{t,v})_{t \geq 0}$  solves the first variation equation,

$$(23) \quad dV_{t,v} = \nabla b(Z_{t,x})V_{t,v} dt + \nabla \sigma(Z_{t,x})V_{t,v} dW_t \quad \text{with } V_{0,v} = v,$$

obtained by formally differentiating the equation (5) defining  $(Z_{t,x})_{t \geq 0}$  with respect to  $x$  in the direction  $v$ . The second directional derivative flow  $(U_{t,v,v'})_{t \geq 0}$  solves the second variation equation,

$$(24) \quad \begin{aligned} dU_{t,v,v'} &= (\nabla b(Z_{t,x})U_{t,v,v'} + \nabla^2 b(Z_{t,x})[V_{t,v'}]V_{t,v}) dt \\ &+ (\nabla \sigma(Z_{t,x})U_{t,v,v'} \\ &+ \nabla^2 \sigma(Z_{t,x})[V_{t,v'}]V_{t,v}) dW_t \quad \text{with } U_{0,v,v'} = 0, \end{aligned}$$

obtained by differentiating (23) with respect to  $x$  in the direction  $v'$ .

Since  $f$  has bounded first and second derivatives, the dominated convergence theorem implies that, for each  $t \geq 0$ ,  $P_t f$  is twice differentiable with

$$(25) \quad \begin{aligned} \langle \nabla(P_t f)(x), v \rangle &= \mathbb{E}[\langle \nabla f(Z_{t,x}), V_{t,v} \rangle] \quad \text{and} \\ v'^T \nabla^2(P_t f)(x)v &= \mathbb{E}[V_{t,v'}^T \nabla^2 f(Z_{t,x})V_{t,v} + \langle \nabla f(Z_{t,x}), U_{t,v,v'} \rangle] \end{aligned}$$

obtained by differentiating under the integral sign. Lemma 16, proved in Section C.1, justifies the exchanges of derivative and expectation by ensuring that the derivative flows have moments bounded uniformly in  $x$ .

**LEMMA 16 (Derivative flow bounds).** *Suppose that  $(Z_{t,x})_{t \geq 0}$  is an Itô diffusion with starting point  $Z_{0,x} = x \in \mathbb{R}^d$ , driving Wiener process  $(W_t)_{t \geq 0}$ , and Lipschitz drift and diffusion coefficients  $b$  and  $\sigma$  with Lipschitz gradients and locally Lipschitz second derivatives. If  $(V_{t,v})_{t \geq 0}$  and  $(U_{t,v,v'})_{t \geq 0}$ , respectively, solve the stochastic differential equations (23) and (24) for  $v, v' \in \mathbb{R}^d$ , then, for any  $\rho \geq 2$ ,*

$$(26) \quad \mathbb{E}[\|V_{t,v}\|_2^\rho] \leq \|v\|_2^\rho e^{t\gamma_\rho} \quad \text{and}$$

$$(27) \quad \mathbb{E}[\|U_{t,v,v'}\|_2^2] \leq \alpha \|v\|_2^2 \|v'\|_2^2 e^{t\gamma_4}$$

for  $\gamma_\rho \triangleq \rho M_1(b) + \frac{\rho^2 - 2\rho}{2} M_1(\sigma)^2 + \frac{\rho}{2} F_1(\sigma)^2$  and  $\alpha \triangleq \frac{M_2(b)^2}{2M_1(b) + 4M_1(\sigma)^2} + 2F_2(\sigma)^2$ .

Since  $\nabla f$  and  $\nabla^2 f$  are bounded, and  $(V_{t,v})_{t \geq 0}$ ,  $(V_{t,v'})_{t \geq 0}$ , and  $(U_{t,v,v'})_{t \geq 0}$  have second moments bounded uniformly in  $x$  by Lemma 16, the Hessian formula (25) implies that  $\nabla^2 P_t f$  is bounded, and hence that  $\nabla P_t f$  is Lipschitz continuous for each  $t \geq 0$ .



**Continuity of  $\nabla^2 P_t f$ .** Hereafter, we assume that  $M_0(\sigma^{-1}) < \infty$ , as the semi-group Hessian bound (22) is otherwise vacuous.

The Lipschitz continuity of  $f$  and the Itô diffusion moment bound of [53], Theorem 3.4, part 4, together imply that

$$\mathbb{E}[f(Z_{t,x})^2] \leq \mathbb{E}[(|f(x)| + \|Z_{t,x} - x\|_2 M_1(f))^2] < \infty$$

for all  $t \geq 0$ . Since  $\sigma^{-1}$  is bounded, and  $\nabla b$  and  $\nabla \sigma$  are bounded and Lipschitz, [27], Proposition 3.2, gives the following Bismut–Elworthy–Li-type formula for the directional derivative of  $P_t f$  for each  $t > 0$ :

$$\langle \nabla(P_t f)(x), v \rangle = \frac{1}{t} \mathbb{E} \left[ f(Z_{t,x}) \int_0^t \langle \sigma^{-1}(Z_{s,x}) V_{s,v}, dW_s \rangle \right].$$

By interchanging derivative and integral, the dominated convergence theorem now delivers the Hessian expression

$$\begin{aligned} v'^\top \nabla^2(P_t f)(x)v &= \mathbb{E}[J_{1,x} + J_{2,x} + J_{3,x}] \quad \text{for} \\ J_{1,x} &\triangleq \frac{1}{t} \langle \nabla f(Z_{t,x}), V_{t,v'} \rangle \int_0^t \langle \sigma^{-1}(Z_{s,x}) V_{s,v}, dW_s \rangle, \\ J_{2,x} &\triangleq \frac{1}{t} f(Z_{t,x}) \int_0^t \langle \nabla \sigma^{-1}(Z_{s,x}) [V_{s,v'}] V_{s,v}, dW_s \rangle, \quad \text{and} \\ J_{3,x} &\triangleq \frac{1}{t} f(Z_{t,x}) \int_0^t \langle \sigma^{-1}(Z_{s,x}) U_{s,v,v'}, dW_s \rangle, \end{aligned} \tag{28}$$

for each  $t > 0$ , provided that  $J_{1,x}$ ,  $J_{2,x}$ , and  $J_{3,x}$  are continuous in  $x$ . The requisite continuity follows from the Lipschitz continuity of  $\nabla f$  and  $f$ , the boundedness of  $\sigma^{-1}$ ,  $\nabla \sigma$ , and  $\nabla^2 \sigma$ , and the controlled moment growth and Hölder continuity of  $(Z_{t,x})_{t \geq 0}$ ,  $(V_{t,v})_{t \geq 0}$ ,  $(V_{t,v'})_{t \geq 0}$  and  $(U_{t,v,v'})_{t \geq 0}$  as functions of  $x$  [77], Theorem V.40. The dominated convergence theorem further implies that  $\nabla^2 P_t f$  is continuous for each  $t > 0$ .

**Initial bound on  $\nabla^2 P_t f$ .** Now, we fix any  $t > 0$  and turn to bounding  $\nabla^2 P_t f$  in terms of  $M_1(f)$ , by bounding the expectations of  $J_{1,x}$ ,  $J_{2,x}$  and  $J_{3,x}$  of (28) in turn.

To control  $\mathbb{E}[J_{1,x}]$ , we apply Cauchy–Schwarz, the Itô isometry [28], equations (7.1) and (7.2), the derivative flow bound (26) and the fact  $e^{s\gamma/2} \leq e^{t\gamma/2}$  for all  $s \leq t$  to obtain

$$\begin{aligned} \mathbb{E}[J_{1,x}] &\leq \frac{1}{t} \sqrt{\mathbb{E}[\langle \nabla f(Z_{t,x}), V_{t,v'} \rangle^2]} \mathbb{E} \left[ \left( \int_0^t \langle \sigma^{-1}(Z_{s,x}) V_{s,v}, dW_s \rangle \right)^2 \right] \\ &\leq \frac{1}{t} M_1(f) \sqrt{\mathbb{E}[\|V_{t,v'}\|_2^2]} \int_0^t \mathbb{E}[\|\sigma^{-1}(Z_{s,x}) V_{s,v}\|_2^2] ds \end{aligned}$$

$$\begin{aligned} &\leq \frac{1}{t} M_1(f) M_0(\sigma^{-1}) \sqrt{\mathbb{E}[\|V_{t,v'}\|_2^2] \int_0^t \mathbb{E}[\|V_{s,v}\|_2^2] ds} \\ &\leq \frac{1}{t} M_1(f) M_0(\sigma^{-1}) \|v'\|_2 \|v\|_2 \sqrt{e^{t\gamma_2} \int_0^t e^{s\gamma_2} ds} \\ &\leq \sqrt{\frac{1}{t}} e^{t\gamma_2} M_1(f) M_0(\sigma^{-1}) \|v'\|_2 \|v\|_2, \end{aligned}$$

where we have adopted the definition of  $\gamma_\rho$  given in Lemma 16.

To control  $\mathbb{E}[J_{2,x}]$ , we will first rewrite the unbounded quantity  $f(Z_{t,x})$  in terms of more manageable semigroup gradients. To this end, we note that, since  $P_{t-s}f \in C^2$  for all  $s \in [0, t]$ , we may apply Itô’s formula [28], Theorem 7.1, to  $(s, x) \mapsto P_{t-s}f(x)$  to obtain the identity

$$(29) \quad f(Z_{t,x}) = (P_t f)(x) + \int_0^t \langle \nabla(P_{t-s}f)(Z_{s,x}), \sigma(Z_{s,x}) dW_s \rangle.$$

Now we may rewrite  $\mathbb{E}[J_{2,x}]$  as

$$\begin{aligned} \mathbb{E}[J_{2,x}] &= \frac{1}{t} \mathbb{E} \left[ (P_t f)(x) \int_0^t \langle \nabla \sigma^{-1}(Z_{s,x})[V_{s,v'}] V_{s,v}, dW_s \rangle \right. \\ &\quad \left. + \int_0^t \langle \nabla(P_{t-s}f)(Z_{s,x}), \sigma(Z_{s,x}) dW_s \rangle \right. \\ &\quad \left. \times \int_0^t \langle \nabla \sigma^{-1}(Z_{s,x})[V_{s,v'}] V_{s,v}, dW_s \rangle \right] \\ &= \frac{1}{t} \mathbb{E} \left[ \int_0^t \langle \nabla(P_{t-s}f)(Z_{s,x}), \sigma(Z_{s,x}) \nabla \sigma^{-1}(Z_{s,x})[V_{s,v'}] V_{s,v} \rangle ds \right] \\ &= -\frac{1}{t} \mathbb{E} \left[ \int_0^t \langle \nabla(P_{t-s}f)(Z_{s,x}), \nabla \sigma(Z_{s,x})[V_{s,v'}] \sigma^{-1}(Z_{s,x}) V_{s,v} \rangle ds \right], \end{aligned}$$

where we have used Dynkin’s formula [28], equation (7.11), the Itô isometry and the chain rule,

$$(30) \quad \nabla \sigma^{-1}(x)[v] = -\sigma^{-1}(x) \nabla \sigma(x)[v] \sigma^{-1}(x).$$

Finally, we bound  $\mathbb{E}[J_{2,x}]$  using Cauchy–Schwarz, the semigroup gradient bound (21), the derivative flow bound (26) and the fact that  $s \mapsto r(t-s)e^{s\gamma_2}$  is increasing:

$$\begin{aligned} \mathbb{E}[J_{2,x}] &\leq \frac{1}{t} M_1(\sigma) M_0(\sigma^{-1}) \int_0^t M_1(P_{t-s}f) \mathbb{E}[\|V_{s,v'}\|_2 \|V_{s,v}\|_2] ds \\ &\leq \frac{1}{t} M_1(\sigma) M_0(\sigma^{-1}) \int_0^t M_1(P_{t-s}f) \sqrt{\mathbb{E}[\|V_{s,v'}\|_2^2] \mathbb{E}[\|V_{s,v}\|_2^2]} ds \end{aligned}$$

$$\begin{aligned} &\leq \frac{1}{t} M_1(\sigma) M_0(\sigma^{-1}) M_1(f) \|v'\|_2 \|v\|_2 \int_0^t r(t-s) e^{s\gamma_2} ds \\ &\leq r(0) e^{t\gamma_2} M_1(\sigma) M_0(\sigma^{-1}) M_1(f) \|v'\|_2 \|v\|_2. \end{aligned}$$

To control  $\mathbb{E}[J_{3,x}]$ , we again appeal to Dynkin’s formula and the Itô isometry to obtain

$$\begin{aligned} \mathbb{E}[J_{3,x}] &= \frac{1}{t} \mathbb{E} \left[ (P_t f)(x) \int_0^t \langle \sigma^{-1}(Z_{s,x}) U_{s,v,v'}, dW_s \rangle \right. \\ &\quad \left. + \int_0^t \langle \nabla(P_{t-s} f)(Z_{s,x}), \sigma(Z_{s,x}) dW_s \rangle \int_0^t \langle \sigma^{-1}(Z_{s,x}) U_{s,v,v'}, dW_s \rangle \right] \\ &= \mathbb{E} \left[ \int_0^t \langle \nabla(P_{t-s} f)(Z_{s,x}), U_{s,v,v'} \rangle ds \right], \end{aligned}$$

and we bound this expression using Cauchy–Schwarz, Jensen’s inequality, the semigroup gradient bound (21), the second derivative flow bound (27) and the fact that  $s \mapsto r(t-s)e^{s\gamma_4}$  is increasing:

$$\begin{aligned} \mathbb{E}[J_{3,x}] &\leq \frac{1}{t} \int_0^t M_1(P_{t-s} f) \mathbb{E}[\|U_{s,v,v'}\|_2] ds \\ &\leq \frac{1}{t} \int_0^t M_1(P_{t-s} f) \sqrt{\mathbb{E}[\|U_{s,v,v'}\|_2^2]} ds \\ &\leq \frac{1}{t} M_1(f) \sqrt{\alpha} \|v'\|_2 \|v\|_2 \int_0^t r(t-s) \sqrt{s} e^{s\gamma_4} ds \\ &\leq \frac{2}{3} \sqrt{tr(0)} e^{t\gamma_4} M_1(f) \sqrt{\alpha} \|v'\|_2 \|v\|_2, \end{aligned}$$

where  $\alpha$  is defined in Lemma 16. The advertised result (22) for  $t_0 = t$  follows by summing the bounds developed for  $\mathbb{E}[J_{1,x}]$ ,  $\mathbb{E}[J_{2,x}]$  and  $\mathbb{E}[J_{3,x}]$ .

**Infimal bound on  $\nabla^2 P_t f$ .** To obtain the infimum over  $t_0 \in (0, t]$  in (22), we adapt an argument of [8], Proposition 1.5.1. Specifically, fix any  $t_0 \in (0, t]$ . Our work thus far shows that  $v'^\top \nabla^2(P_{t_0} \tilde{f})(x)v \leq M_1(\tilde{f})\zeta(t_0)$  for a real-valued function  $\zeta$  and  $\tilde{f} \in C^2$  with bounded first and second derivatives. Since we now know that  $P_{t-t_0} f \in C^2$  with bounded first and second derivatives, the Markov property of the diffusion and the first derivative bound (21) yield

$$\begin{aligned} v'^\top \nabla^2(P_t f)(x)v &= v'^\top \nabla^2(P_{t_0} P_{t-t_0} f)(x)v \\ &\leq M_1(P_{t-t_0} f)\zeta(t_0) \\ &\leq M_1(f)r(t-t_0)\zeta(t_0). \end{aligned}$$

**C.1. Proof of Lemma 16: Derivative flow bounds.** Fix any  $\rho \geq 2$  and  $v \in \mathbb{R}^d$ . Since Dynkin’s formula and Cauchy–Schwarz give

$$\begin{aligned} \mathbb{E}[\|V_{t,v}\|_2^\rho] &= \|v\|_2^\rho + \mathbb{E}\left[\int_0^t \rho \langle V_{s,v}, \|V_{s,v}\|_2^{\rho-2}, \nabla b(Z_{s,x}) V_{s,v} \rangle \right. \\ &\quad \left. + \frac{\rho}{2} \|V_{s,v}\|_2^{\rho-4} ((\rho - 2) \|V_{s,v}^\top \nabla \sigma(Z_{s,x}) [V_{s,v}]\|_2^2 \right. \\ &\quad \left. + \|V_{s,v}\|_2^2 \|\nabla \sigma(Z_{s,x}) [V_{s,v}]\|_F^2) ds \right] \\ &\leq \|v\|_2^\rho + \int_0^t \left( \rho M_1(b) + \frac{\rho^2 - 2\rho}{2} M_1(\sigma)^2 \right. \\ &\quad \left. + \frac{\rho}{2} F_1(\sigma)^2 \right) \mathbb{E}[\|V_{s,v}\|_2^\rho] ds, \end{aligned}$$

the advertised result (26) follows from Grönwall’s inequality.

Now fix any  $v, v' \in \mathbb{R}^d$ , and define  $U_t \triangleq U_{t,v,v'}$ . Dynkin’s formula and multiple applications of Cauchy–Schwarz and Young’s inequality give

$$\begin{aligned} \mathbb{E}[\|U_t\|_2^2] &= \mathbb{E}\left[\int_0^t 2 \langle U_s, \nabla b(Z_{s,x}) U_s + \nabla^2 b(Z_{s,x}) [V_{s,v'}] V_{s,v} \rangle \right. \\ &\quad \left. + \|\nabla \sigma(Z_{s,x}) [U_s] + \nabla^2 \sigma(Z_{s,x}) [V_{s,v'}] V_{s,v}\|_F^2 ds \right] \\ &\leq \mathbb{E}\left[\int_0^t 2 \|U_s\|_2^2 M_1(b) + 2 \|U_s\|_2 \|V_{s,v}\|_2 \|V_{s,v'}\|_2 M_2(b) \right. \\ &\quad \left. + 2 \|\nabla \sigma(Z_{s,x}) [U_s]\|_F^2 + 2 \|\nabla^2 \sigma(Z_{s,x}) [V_{s,v'}] V_{s,v}\|_F^2 ds \right] \\ &\leq \int_0^t (2M_1(b) + 2F_1(\sigma)^2 + \epsilon) \mathbb{E}[\|U_s\|_2^2] \\ &\quad + (M_2(b)^2/\epsilon + 2F_2(\sigma)^2) \mathbb{E}[\|V_{s,v}\|_2^2 \|V_{s,v'}\|_2^2] ds \end{aligned}$$

for any  $\epsilon > 0$ . Letting  $\gamma_\rho = \rho M_1(b) + \frac{\rho^2 - 2\rho}{2} M_1(\sigma)^2 + \frac{\rho}{2} F_1(\sigma)^2$ , we see that, by Cauchy–Schwarz and our derivative flow bound (26),

$$\begin{aligned} \int_0^t \mathbb{E}[\|V_{s,v}\|_2^2 \|V_{s,v'}\|_2^2] ds &\leq \int_0^t \sqrt{\mathbb{E}[\|V_{s,v}\|_2^4] \mathbb{E}[\|V_{s,v'}\|_2^4]} ds \\ &\leq \int_0^t \|v\|_2^2 \|v'\|_2^2 e^{s\gamma_4} ds \\ &= \|v\|_2^2 \|v'\|_2^2 \frac{e^{t\gamma_4} - 1}{\gamma_4}. \end{aligned}$$

Hence, if we choose  $\epsilon = \gamma_4 - (2M_1(b) + 2F_1(\sigma)^2)$  and define  $\alpha = M_2(b)^2/\epsilon + 2F_2(\sigma)^2$  we may write

$$\mathbb{E}[\|U_t\|_2^2] \leq \alpha \|v\|_2^2 \|v'\|_2^2 \frac{e^{t\gamma_4} - 1}{\gamma_4} + \int_0^t \gamma_4 \mathbb{E}[\|U_s\|_2^2] ds.$$

Gronwall's inequality now yields the result (27) via

$$\begin{aligned} \mathbb{E}[\|U_t\|_2^2] &\leq \alpha \|v\|_2^2 \|v'\|_2^2 \left( \frac{e^{t\gamma_4} - 1}{\gamma_4} + \int_0^t \frac{e^{s\gamma_4} - 1}{\gamma_4} \gamma_4 e^{(t-s)\gamma_4} ds \right) \\ &= \alpha \|v\|_2^2 \|v'\|_2^2 t e^{t\gamma_4}. \end{aligned}$$

APPENDIX D: PROOF OF THEOREM 6

We first derive the result for  $\|\cdot\| = \|\cdot\|_2$ . Without loss of generality, assume  $h \in \mathcal{W}_{\|\cdot\|_2}$  with  $\mathbb{E}_p[h(Z)] = 0$ . Our high-level strategy is to relate the Wasserstein distance to the Stein discrepancy via the Stein equation (3) with diffusion Stein operator  $\mathcal{T}$  (8). Since the infinitesimal generator  $\mathcal{A}$  (4) has the form (7) by Theorem 2, Theorem 5 implies that there exists a continuously differentiable solution  $g_h$  to the the Stein equation  $h(x) = (\mathcal{T}g_h)(x)$  satisfying  $M_0(g_h) \leq s_r M_1(h) \leq s_r$ . Since boundedness alone is insufficient to declare that  $g_h$  falls into a scaled copy of the classical Stein set  $\mathcal{G}_{\|\cdot\|}$ , we will develop a smoothed version of the Stein solution with greater regularity.

Since  $a$  and  $c$  are constant,  $b(x) = \frac{1}{2}(a + c)\nabla \log p(x)$ . Fix any  $s > 0$  and consider the convolution  $g_{h,s}(x) \triangleq \mathbb{E}[g_h(x + sG)]$ . If the smoothing level  $s$  is small, the Lipschitz continuity of  $h$  implies that that  $(\mathcal{T}g_{h,s})(x)$  provides a close approximation to  $h(x)$  for each  $x \in \mathbb{R}^d$ :

$$\begin{aligned} (31) \quad h(x) &\leq \mathbb{E}[h(x + sG)] + M_1(h)s\mathbb{E}[\|G\|_2] \\ &\leq \mathbb{E}\left[\frac{1}{p(x + sG)} \langle \nabla, p(x + sG)(a + c)g_h(x + sG) \rangle\right] \\ &\quad + s\mathbb{E}[\|G\|_2] \\ &\leq 2\mathbb{E}[\langle b(x + sG), g_h(x + sG) \rangle] \\ &\quad + \mathbb{E}[(a + c, \nabla g_h(x + sG))] + s\mathbb{E}[\|G\|_2] \\ &\leq (\mathcal{T}g_{h,s})(x) + s\mathbb{E}[\|G\|_2](1 + 2M_1(b)M_0(g_h)). \end{aligned}$$

Moreover, by our next lemma, proved in Section D.1, the smoothed Stein solution admits a bounded Lipschitz gradient  $\nabla g_{h,s}(x) = \mathbb{E}[\nabla g_h(x + sG)]$ .

LEMMA 17 (Smoothing by Gaussian convolution). *Let  $G \in \mathbb{R}^d$  be a standard normal random vector, and fix  $s > 0$ . If  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is bounded and measurable*

and  $f_s(x) \triangleq \mathbb{E}[f(x + sG)]$ , then

$$M_0(f_s) \leq M_0(f), \quad M_1(f_s) \leq \sqrt{\frac{2}{\pi}} \frac{M_0(f)}{s}, \quad \text{and} \quad M_2(f_s) \leq \sqrt{2} \frac{M_0(f)}{s^2}.$$

If, additionally,  $f \in C^1$ , then  $\nabla f_s(x) = \mathbb{E}[\nabla f(x + sG)]$ .

Indeed, for each nonzero  $w \in \mathbb{R}^d$ , we may apply Lemma 17 to the function  $f_w(x) \triangleq \langle w, g_h(x) \rangle / \|w\|_2$  with convolution  $f_{w,s}(x) = \langle w, g_{h,s}(x) \rangle / \|w\|_2$  to obtain the bounds

$$M_0(g_{h,s}) = \sup_{w \neq 0} M_0(f_{w,s}) \leq \sup_{w \neq 0} M_0(f_w) = M_0(g_h) \leq s_r,$$

$$M_1(g_{h,s}) = \sup_{w \neq 0} M_1(f_{w,s}) \leq \sup_{w \neq 0} \sqrt{\frac{2}{\pi}} \frac{M_1(f_w)}{s} = \sqrt{\frac{2}{\pi}} \frac{M_1(f_w)}{s} \leq \sqrt{\frac{2}{\pi}} \frac{s_r}{s},$$

and

$$M_2(g_{h,s}) = \sup_{w \neq 0} M_2(f_{w,s}) \leq \sup_{w \neq 0} \frac{\sqrt{2} M_2(f_w)}{s^2} = \frac{\sqrt{2} M_2(f_w)}{s^2} \leq \frac{\sqrt{2} s_r}{s^2}.$$

Hence, since our choice of  $h$  was arbitrary, and

$$\begin{aligned} \kappa_s &\triangleq \max\left(1, \frac{1}{s} \sqrt{\frac{2}{\pi}}, \frac{\sqrt{2}}{s^2}\right) \\ &= \max\left(1, \frac{\sqrt{2}}{s^2}\right) \geq \frac{\max(M_0(g_{h,s}), M_1(g_{h,s}), M_2(g_{h,s}))}{s_r}, \end{aligned}$$

we may take expectation under  $Q_n$  and supremum over  $h$  in (31) to reach

$$\begin{aligned} d_{\mathcal{W}_{\|\cdot\|_2}}(\mu, \nu) &\leq \inf_{s>0} \mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|_2}) s_r \kappa_s + s \mathbb{E}[\|G\|_2] (1 + 2M_1(b) s_r) \\ &\leq \max(\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|_2}) s_r, \eta) + 2\eta \\ &\leq 3 \max(\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|_2}) s_r, \eta), \end{aligned}$$

where we define  $\eta = \sqrt[3]{\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|_2}) \sqrt{2} s_r \mathbb{E}[\|G\|_2]^2 (1 + 2M_1(b) s_r)^2}$  and select  $s = \sqrt[3]{\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|_2}) 2\sqrt{2} s_r / (\mathbb{E}[\|G\|_2] (1 + 2M_1(b) s_r))}$  to produce the second inequality. The generic norm result now follows from the assumed norm domination property  $\|\cdot\| \geq \|\cdot\|_2$ , which implies  $\mathcal{G}_{\|\cdot\|_2} \subseteq \mathcal{G}_{\|\cdot\|}$ .

**D.1. Proof of Lemma 17: Smoothing by Gaussian convolution.** The conclusion  $M_0(f_s) \leq M_0(f)$  follows from Hölder’s inequality. Now, fix any  $x$  and

nonzero  $v_1, v_2 \in \mathbb{R}^d$ . Since  $f_s = f \star \phi_s$ , where  $\phi_s \in C^\infty$  is the density of  $sG$  and  $\star$  is the convolution operator, Leibniz’s rule implies that

$$\begin{aligned} \langle v_1, \nabla f_s(x) \rangle &= \langle v_1, (f \star \nabla \phi_s)(x) \rangle \\ &= \frac{1}{s^2} \int f(x - y) \langle v_1, y \rangle \phi_s(y) dy \\ &\leq \frac{M_0(f)}{s^2} \int |\langle v_1, y \rangle| \phi_s(y) dy = \sqrt{\frac{2}{\pi}} \frac{M_0(f)}{s} \|v_1\|_2, \end{aligned}$$

as  $\langle v_1, G \rangle / \|v_1\|_2$  has a standard normal distribution. Leibniz’s rule also gives

$$\begin{aligned} \nabla^2 f_s(x)[v_1, v_2] &= (f \star \nabla^2 \phi_s)(x)[v_1, v_2] \\ &\leq \frac{M_0(f)}{s^2} \int_{\mathbb{R}^d} |\langle v_1, zz^\top v_2 \rangle / s^2 - \langle v_1, v_2 \rangle| \phi_s(z) dz \\ &\leq \frac{M_0(f)}{s^2} \sqrt{\int_{\mathbb{R}^d} |\langle v_1, zz^\top v_2 \rangle / s^2 - \langle v_1, v_2 \rangle|^2 \phi_s(z) dz} \\ &= \frac{M_0(f)}{s^2} \sqrt{\langle v_1, v_2 \rangle^2 + \|v_1\|_2^2 \|v_2\|_2^2} \\ &\leq \frac{\sqrt{2} M_0(f)}{s^2} \|v_1\|_2 \|v_2\|_2, \end{aligned}$$

where the last equality follows by Isserlis’ theorem. Finally, when  $f \in C^1$ , Leibniz’s rule gives  $\nabla f_s = \nabla f \star \phi_s$ .

APPENDIX E: PROOF OF THEOREM 7

We will derive each inequality for  $\|\cdot\| = \|\cdot\|_2$ ; the generic norm results will then follow from the property  $\|\cdot\| \geq \|\cdot\|_2$ , which implies  $\mathcal{G}_{\|\cdot\|_2} \subseteq \mathcal{G}_{\|\cdot\|}$ .

Fix any  $h \in \mathcal{H} = \{h : \mathbb{R}^d \rightarrow \mathbb{R} \mid h \in C^3, M_1(h) \leq 1, M_2(h) < \infty, M_3(h) < \infty\}$  with  $\mathbb{E}_P[h(Z)] = 0$ . We assume that  $M_1(b), M_2(b), M_1(\sigma), F_2(\sigma), M_1^*(m)$  and  $M_0(\sigma^{-1})$  are all finite, or else the results are vacuous. Our high-level strategy is to relate the Wasserstein distance to the Stein discrepancy via the Stein equation (3) with diffusion Stein operator  $\mathcal{T}$  (8). By Theorem 5, we know that there exists a Lipschitz solution  $g_h$  to the Stein equation  $h(x) = (\mathcal{T} g_h)(x)$  satisfying  $M_0(g_h) \leq s_r M_1(h) \leq s_r$  and  $M_1(g_h) \leq \beta M_1(h) \leq \beta$ , for  $\beta \triangleq \beta_1 + \beta_2$ , where  $\beta_1$  and  $\beta_2$  are defined in Theorem 5. Since a Lipschitz gradient is also needed to declare that  $g_h$  falls into a scaled copy of the classical Stein set  $\mathcal{G}_{\|\cdot\|}$ , we will develop a smoothed version of the Stein solution with greater regularity.

For this purpose, fix any  $s > 0$  and consider the convolution  $g_{h,s}(x) \triangleq \mathbb{E}[g_h(x + sG)]$ . If the smoothing level  $s$  is small, the Lipschitz continuity of  $m$



and  $h$  implies that  $(\mathcal{T}g_{h,s})(x)$  closely approximates  $h(x)$  for each  $x \in \mathbb{R}^d$ :

$$\begin{aligned}
 h(x) &\leq \mathbb{E}[h(x + sG)] + M_1(h)s\mathbb{E}[\|G\|_2] \\
 &\leq 2\mathbb{E}[\langle b(x + sG), g_h(x + sG) \rangle + \langle m(x + sG), \nabla g_h(x + sG) \rangle] \\
 (32) \quad &\quad + s\mathbb{E}[\|G\|_2] \\
 &\leq (\mathcal{T}g_{h,s})(x) + s\zeta.
 \end{aligned}$$

**E.1. Proof of the first inequality.** Moreover, by an argument mirroring that of Theorem 6, Lemma 17 shows that  $g_{h,s}$  admits a Lipschitz gradient  $\nabla g_{h,s}(x) = \mathbb{E}[\nabla g_h(x + sG)]$  and satisfies the derivative bounds

$$\begin{aligned}
 (33) \quad M_0(g_{h,s}) &\leq M_0(g_h) \leq s_r, \\
 M_1(g_{h,s}) &= M_0(\nabla g_{h,s}) \leq M_0(\nabla g_h) \leq \beta, \quad \text{and} \\
 M_2(g_{h,s}) &= M_1(\nabla g_{h,s}) \leq \sqrt{\frac{2}{\pi}} \frac{M_0(\nabla g_h)}{s} \leq \sqrt{\frac{2}{\pi}} \frac{\beta}{s}.
 \end{aligned}$$

Let  $\eta \triangleq s^* \zeta$  for  $s^* = \sqrt{\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|_2}) \sqrt{2/\pi} \beta / \zeta}$ . Since  $\mathcal{H}$  is dense in  $\mathcal{W}_{\|\cdot\|_2}$ , we may take expectation under  $Q_n$  and supremum over  $h$  in (32) to reach

$$\begin{aligned}
 d_{\mathcal{W}_{\|\cdot\|_2}}(\mu, \nu) &\leq \inf_{s>0} \mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|_2}) \max\left(s_r, \beta, \sqrt{\frac{2}{\pi}} \frac{\beta}{s}\right) + s\zeta \\
 &\leq \max(\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|_2}) \max(s_r, \beta), \eta) + \eta \\
 &\leq 2 \max(\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|_2}) \max(s_r, \beta), \eta).
 \end{aligned}$$

**E.2. Proof of the second inequality.** Assume now that  $\nabla^3 b$  and  $\nabla^3 \sigma$  are bounded and locally Lipschitz. Fix any  $\iota \in (0, 1)$ . Lemma 17 and an auxiliary smoothing lemma (Lemma 18 in the Supplementary Material [35]) imply that  $M_2(g_{h,s}) = M_1(\nabla g_{h,s}) \leq \sqrt{d} \frac{M_{1-\iota}(\nabla g_h)}{s^\iota}$ . This improved dependence on  $s$  will allow us to establish a near-linear relationship between the Stein discrepancy and the Wasserstein distance. By Theorem 5,  $M_{1-\iota}(\nabla g_h) \leq \frac{1}{K}(\frac{1}{\iota} + s_r)$  for  $K$  depending only on  $M_{1:3}(\sigma)$ ,  $M_{1:3}(b)$ ,  $M_0(\sigma^{-1})$  and  $r$ . Hence,  $M_2(g_{h,s}) \leq C_\iota/s^\iota$  for  $C_\iota \triangleq \frac{\sqrt{d}}{K}(\frac{1}{\iota} + s_r)$ . Following the derivation in Section E.1 and choosing  $s^* = (\frac{\iota C_\iota \mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|_2})}{\zeta})^{\frac{1}{\iota+1}}$  and  $\eta \triangleq \frac{\zeta}{\iota} s^*$ , we obtain

$$\begin{aligned}
 (34) \quad d_{\mathcal{W}_{\|\cdot\|_2}}(P, Q_n) &\leq \inf_{s>0} \mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|_2}) \max(s_r, \beta, C_\iota s^{-\iota}) + s\zeta \\
 &\leq \max(\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|_2}) \max(s_r, \beta), \eta) + \eta \iota \\
 &\leq 2 \max(\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|_2}) \max(s_r, \beta), \eta).
 \end{aligned}$$

Now consider the case in which  $\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|}) < e^{-1}$  and the choice  $\iota = 1/\log(1/\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|})) \in (0, 1)$ . Since  $x^{1/(\log x - 1)} \leq e$  for all  $x \in (0, e^{-1})$ ,

$$\begin{aligned} \frac{1}{\iota} \mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|})^{\frac{1}{1+\iota}} &= \log(1/\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|})) \mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|})^{1+1/(\log \mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|})-1)} \\ &\leq e \log(1/\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|})) \mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|}). \end{aligned}$$

Introduce the shorthand  $c_0 = \frac{\sqrt{d}}{K\xi}$ . Since  $1/1+\iota \in (1/2, 1)$ , we have  $c_0^{\frac{1}{1+\iota}} \leq \max(\sqrt{c_0}, c_0)$ . Similarly,  $1 + s_r \iota > 1$ , so  $(1 + \iota s_r)^{\frac{1}{1+\iota}} \leq 1 + \iota s_r$ . Therefore,

$$\begin{aligned} &\frac{\xi}{\iota} \mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|})^{\frac{1}{1+\iota}} \left( \frac{1 + \iota s_r}{K\xi/\sqrt{d}} \right)^{\frac{1}{1+\iota}} \\ &\leq e\xi \mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|}) \log(1/\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|})) \\ &\quad \times \max\left( \frac{d^{1/4}}{\sqrt{K\xi}}, \frac{\sqrt{d}}{K\xi} \right) \left( 1 + \frac{s_r}{\log(1/\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|}))} \right) \\ &= e\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|}) \max\left( \frac{d^{1/4}\sqrt{\xi}}{\sqrt{K}}, \frac{\sqrt{d}}{K} \right) (s_r + \log(1/\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|}))). \end{aligned}$$

Next, fix any  $\iota \in (0, 1)$  and consider the case in which  $\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|}) \geq e^{-1}$  so that  $\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|})^{\frac{1}{1+\iota}} \leq \mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|}) e^{\frac{\iota}{1+\iota}}$ . Because  $\frac{1}{\iota} e^{\frac{\iota}{1+\iota}} \leq \frac{1}{2} e^{1/2} < e$  and  $(1 + \iota s_r)^{1/1+\iota} \leq 1 + s_r$ , we conclude that

$$\begin{aligned} &\frac{\xi}{\iota} \mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|})^{\frac{1}{1+\iota}} \left( \frac{1 + \iota s_r}{K\xi/\sqrt{d}} \right)^{\frac{1}{1+\iota}} \\ &\leq e\mathcal{S}(Q_n, \mathcal{T}, \mathcal{G}_{\|\cdot\|}) \max\left( \frac{d^{1/4}\sqrt{\xi}}{\sqrt{K}}, \frac{\sqrt{d}}{K} \right) (s_r + 1). \end{aligned}$$

The result follows from estimates of these two cases and the bound (34).

APPENDIX F: PROOF OF PROPOSITION 8

Fix any  $g \in \mathcal{G}_{\|\cdot\|}$ . Since  $\mathbb{E}_P[(\mathcal{T}g)(Z)] = 0$  by Proposition 3, we may write

$$\begin{aligned} &|\mathbb{E}_{Q_n}[(\mathcal{T}g)(X)]| \\ (35) \quad &= |\mathbb{E}_{Q_n}[(\mathcal{T}g)(X)] - \mathbb{E}_P[(\mathcal{T}g)(Z)]| \\ &= |2\mathbb{E}[\langle b(X) - b(Z), g(X) \rangle + \langle b(Z), g(X) - g(Z) \rangle] \\ &\quad + \mathbb{E}[\langle m(X) - m(Z), \nabla g(X) \rangle + \langle m(Z), \nabla g(X) - \nabla g(Z) \rangle]| \end{aligned}$$

for any coupling of  $X$  and  $Z$ . We obtain the first advertised inequality by repeatedly applying the Fenchel–Young inequality for dual norms, invoking the boundedness

and Lipschitz constraints on  $g$  and  $\nabla g$ , and taking a supremum over  $g \in \mathcal{G}_{\|\cdot\|}$ . The second inequality follows from the first by invoking Jensen’s inequality, the fact  $\min(x, y) \leq x^t y^{1-t}$  for all  $x, y \geq 0$ , Hölder’s inequality, and finally the definition of  $W_{s, \|\cdot\|}$ .

We prove the final claim by bounding the first advertised inequality in a second manner. Let  $(X, Z)$  be coupled so that  $c \triangleq \min(W_{1, \|\cdot\|}(Q_n, P), 2) = \min(\mathbb{E}[\|X - Z\|], 2)$ ,  $A = 2\|b(Z)\| + \|m(Z)\|$ , and  $B = \min(\|X - Z\|, 2)$ . The Fenchel–Young inequality ( $xy \leq e^x - y + y \log y$  for  $y \geq 0, x \in \mathbb{R}$ ), the concavity of  $x \mapsto \min(x, 2)$  and Jensen’s inequality now yield the result as

$$\begin{aligned} \mathbb{E}[AB] &= \mathbb{E}[(A - \log(\mu_0/c))B] + \mathbb{E}[B] \log(\mu_0/c) \\ &\leq \mathbb{E}[e^{A - \log(\mu_0/c)} - B + B \log(B)] + \mathbb{E}[B] \log(\mu_0/c) \\ &= c - \mathbb{E}[B \log(e/B)] + \mathbb{E}[B] \log(\mu_0/c) \\ &\leq c + c \log(\mu_0/c) = c \log(e\mu_0/c). \end{aligned}$$

APPENDIX G: PROOF OF THEOREM 10

Fix any  $x, y \in \mathbb{R}^d$ , and define two Itô diffusions solving  $dZ_{t,x} = b(Z_{t,x}) dt + \sigma(Z_{t,x}) dW_t$  with  $Z_{0,x} = x$  and  $dZ_{t,y} = b(Z_{t,y}) dt + \sigma(Z_{t,y}) dW_t$  with  $Z_{0,y} = y$ , for  $(W_t)_{t \geq 0}$  a shared Wiener process. Applying Dynkin’s formula to the function  $f(t, x) = e^{kt} \|x\|_G^2$  for the difference process  $Z_{t,x} - Z_{t,y}$  yields

$$\begin{aligned} \mathbb{E}[f(t, Z_{t,x} - Z_{t,y})] &= \|x - y\|_G^2 + \mathbb{E}\left[\int_0^t k e^{ks} \|Z_{s,x} - Z_{s,y}\|_G^2 ds\right] \\ &\quad + \mathbb{E}\left[\int_0^t e^{ks} (\|\sigma(Z_{s,x}) - \sigma(Z_{s,y})\|_G^2 \right. \\ &\quad \left. + 2\langle b(Z_{s,x}) - b(Z_{s,y}), G(Z_{s,x} - Z_{s,y}) \rangle) ds\right]. \end{aligned}$$

By the uniform dissipativity assumption, the right-hand side is at most  $\|x - y\|_G^2 = d_{\mathcal{W}_{\|\cdot\|_G}}(\delta_x, \delta_y)^2$ . For the transition semigroup  $(P_t)_{t \geq 0}$ ,

$$\mathbb{E}[f(t, Z_{t,x} - Z_{t,y})] = e^{kt} \mathbb{E}[\|Z_{t,x} - Z_{t,y}\|_G^2] \geq e^{kt} d_{\mathcal{W}_{\|\cdot\|_G}}(\delta_x P_t, \delta_y P_t)^2,$$

by Cauchy–Schwarz. The result now follows from the fact that  $\lambda_{\min}(G_1) \leq \|z\|_G^2 / \|z\|_2^2 \leq \lambda_{\max}(G_1)$  for all  $z \neq 0$ .

APPENDIX H: PROOF OF THEOREM 11

As in the proof of [94], Theorem 2.6, we fix two arbitrary starting points  $x, y \in \mathbb{R}^d$  and define a pair of coupled Itô diffusions  $(Z_{t,x})_{t \geq 0}$  and  $(Z_{t,y})_{t \geq 0}$ , each with

associated marginal semigroup  $(P_t)_{t \geq 0}$ . Specifically, we set  $Z_{0,x} = x$  and  $Z_{0,y} = y$  and let  $(Z_{t,x})_{t \geq 0}$  and  $(Z_{t,y})_{t \geq 0}$  solve the equations

$$\begin{aligned} dZ_{t,x} &= b(Z_{t,x}) dt + \sigma_0(Z_{t,x}) dW'_t + \lambda_0 dW''_t \\ dZ_{t,y} &= b(Z_{t,y}) dt + \sigma_0(Z_{t,y}) dW'_t \\ &\quad + \lambda_0 \left( I - 2 \frac{Z_{t,x} - Z_{t,y}}{\|Z_{t,x} - Z_{t,y}\|_2} \frac{Z_{t,x} - Z_{t,y}^\top}{\|Z_{t,x} - Z_{t,y}\|_2} \right) dW''_t, \end{aligned}$$

where  $(W'_t)_{t \geq 0}$  is an  $m$ -dimensional Wiener process and  $(W''_t)_{t \geq 0}$  is an independent  $d$ -dimensional Wiener process.

Following the argument of Eberle [22], Section 4, we define the difference process  $Y_t = Z_{t,x} - Z_{t,y}$ , its norm  $r_t = \|Y_t\|_2$ , and the one-dimensional Wiener process  $W_t = \int_0^t \langle Y_s / r_s, dW''_s \rangle$ , and apply the generalized Itô formula [51], Theorem 22.5, to obtain the stochastic differential equations

$$\begin{aligned} d\|Y_t\|_2^2 &= (2\langle Y_t, b(Z_{t,x}) - b(Z_{t,y}) \rangle + \|\sigma_0(Z_{t,x}) - \sigma_0(Z_{t,y})\|_F^2 + 4\lambda_0^2) dt \\ &\quad + 2\langle Y_t, (\sigma_0(Z_{t,x}) - \sigma_0(Z_{t,y})) dW'_t \rangle + 4\lambda_0 \|Y_t\|_2 dW_t \quad \text{and} \\ df(r_t) &= f'(r_t) / (r_t) \langle Y_t, (\sigma_0(Z_{t,x}) - \sigma_0(Z_{t,y})) dW'_t \rangle + 2\lambda_0 f'(r_t) dW_t \\ &\quad + \left( f''(r_t) \left( 2\lambda_0^2 + \frac{1}{2} \|(\sigma_0(Z_{t,x}) - \sigma_0(Z_{t,y}))^\top Y_t\|_2^2 / r_t^2 \right) \right. \\ &\quad \left. - \frac{1}{2\alpha} f'(r_t) \kappa(r_t) r_t \right) dt \end{aligned}$$

for any concave increasing  $f : [0, \infty) \mapsto [0, \infty)$  with absolutely continuous derivative,  $f(0) = 0$  and  $f'(0) = 1$ . Since the drift term in the latter equation is bounded above by  $\beta_t \triangleq (2/\alpha)(f''(r_t) - (1/4)f'(r_t)\kappa(r_t)r_t)$ , the argument of [22], p. 15, shows that the results of [22], Theorem 1 and Corollary 2, hold for our choice of  $\alpha$  and  $\kappa$ .

**Acknowledgments.** We thank Simon Lacoste-Julien for sharing his quadrature code, Martin Hairer for discussing interpolation inequalities, Andreas Eberle for reading an earlier version of this manuscript and Murat Erdogdu for identifying an important typographical error in an earlier version of this manuscript.

SUPPLEMENTARY MATERIAL

**Supplementary information for “Measuring sample quality with diffusions”** (DOI: 10.1214/19-AAP1467SUPP; .pdf).

REFERENCES

[1] AMBROSIO, L., FUSCO, N. and PALLARA, D. (2000). *Functions of Bounded Variation and Free Discontinuity Problems*. Oxford University Press. The Clarendon Press, New York. MR1857292

- [2] BACH, F., LACOSTE-JULIEN, S. and OBOZINSKI, G. (2012). On the equivalence between herding and conditional gradient algorithms. In *Proc. 29th ICML, ICML'12*.
- [3] BARBOUR, A. D. (1988). Stein's method and Poisson process convergence. *J. Appl. Probab. Special Vol. 25A* 175–184. A celebration of applied probability. [MR0974580](#)
- [4] BARBOUR, A. D. (1990). Stein's method for diffusion approximations. *Probab. Theory Related Fields* **84** 297–322. [MR1035659](#)
- [5] BOUTS, Q. W., TEN BRINK, A. P. and BUCHIN, K. (2014). A framework for computing the greedy spanner. In *Computational Geometry (SoCG'14)* 11–19. ACM, New York. [MR3382271](#)
- [6] BROOKS, S., GELMAN, A., JONES, G. L. and MENG, X.-L., eds. (2011). *Handbook of Markov Chain Monte Carlo*. CRC Press, Boca Raton, FL.
- [7] CATTIAUX, P. and GUILLIN, A. (2014). Semi log-concave Markov diffusions. In *Séminaire de Probabilités XLVI. Lecture Notes in Math.* **2123** 231–292. Springer, Cham. [MR3330820](#)
- [8] CERRAI, S. (2001). *Second Order PDE's in Finite and Infinite Dimension: A Probabilistic Approach. Lecture Notes in Math.* **1762**. Springer, Berlin. [MR1840644](#)
- [9] CHATTERJEE, S. and MECKES, E. (2008). Multivariate normal approximation using exchangeable pairs. *ALEA Lat. Am. J. Probab. Math. Stat.* **4** 257–283. [MR2453473](#)
- [10] CHATTERJEE, S. and SHAO, Q.-M. (2011). Nonnormal approximation by Stein's method of exchangeable pairs with application to the Curie–Weiss model. *Ann. Appl. Probab.* **21** 464–483. [MR2807964](#)
- [11] CHEN, L. H. Y., GOLDSTEIN, L. and SHAO, Q.-M. (2011). *Normal Approximation by Stein's Method. Probability and Its Applications (New York)*. Springer, Heidelberg. [MR2732624](#)
- [12] CHEN, W. Y., MACKEY, L., GORHAM, J., BRIOL, F.-X. and OATES, C. (2018). Stein points. In *Proc. 35th ICML, ICML'18*.
- [13] CHEN, Y., WELLING, M. and SMOLA, A. (2010). Super-samples from kernel herding. In *UAI*.
- [14] CHEW, P. (1986). There is a Planar Graph Almost As Good As the Complete Graph. In *Proc. 2nd SOCG* 169–177. ACM, New York.
- [15] CHWIAKOWSKI, K., STRATHMANN, H. and GRETTON, A. (2016). A kernel test of goodness of fit. In *Proc. 33rd ICML, ICML*.
- [16] CONCA, C. and VANNINATHAN, M. (2007). Periodic homogenization problems in incompressible fluid equations. *Handbook of Mathematical Fluid Dynamics* **4** 649–698.
- [17] DOERSEK, P. and TEICHMANN, J. (2010). A semigroup point of view on splitting schemes for stochastic (partial) differential equations. Available at <https://arxiv.org/abs/1011.2651>.
- [18] DRUMMOND, P. D. and GARDINER, C. W. (1980). Generalised  $P$ -representations in quantum optics. *J. Phys. A* **13** 2353–2368. [MR0578413](#)
- [19] DRUMMOND, P. D. and WALLS, D. F. (1980). Quantum theory of optical bistability. I. Non-linear polarisability model. *J. Phys. A* **13** 725.
- [20] DUNCAN, A. B., LELIÈVRE, T. and PAVLIOTIS, G. A. (2016). Variance reduction using non-reversible Langevin samplers. *J. Stat. Phys.* **163** 457–491. [MR3483241](#)
- [21] DYNKIN, E. B. (1965). *Markov Processes. Vols. I*. Springer, Berlin–Göttingen–Heidelberg.
- [22] EBERLE, A. (2016). Reflection couplings and contraction rates for diffusions. *Probab. Theory Related Fields* **166** 851–886. [MR3568041](#)
- [23] ENGEL, K.-J. and NAGEL, R. (2000). *One-Parameter Semigroups for Linear Evolution Equations. Graduate Texts in Mathematics* **194**. Springer, New York. [MR1721989](#)
- [24] ETHIER, S. N. and KURTZ, T. G. (1986). *Markov Processes. Wiley Series in Probability and Mathematical Statistics: Probability and Mathematical Statistics*. Wiley, New York. [MR0838085](#)
- [25] FAN, Y., BROOKS, S. P. and GELMAN, A. (2006). Output assessment for Monte Carlo simulations via the score statistic. *J. Comput. Graph. Statist.* **15** 178–206. [MR2269368](#)
- [26] FANG, X., SHAO, Q.-M. and XU, L. (2018). Multivariate approximations in Wasserstein distance by Stein's method and Bismut's formula. *Probab. Theory Related Fields* 1–35.

- [27] FOURNIÉ, E., LASRY, J.-M., LEBUCHOUX, J., LIONS, P.-L. and TOUZI, N. (1999). Applications of Malliavin calculus to Monte Carlo methods in finance. *Finance Stoch.* **3** 391–412. [MR1842285](#)
- [28] FRIEDMAN, A. (1975). *Stochastic Differential Equations and Applications*. Vol. 1. Academic Press, New York. [MR0494490](#)
- [29] GAN, H. L., RÖLLIN, A. and ROSS, N. (2017). Dirichlet approximation of equilibrium distributions in Cannings models with mutation. *Adv. in Appl. Probab.* **49** 927–959. [MR3694323](#)
- [30] GAUNT, R. E. (2016). Rates of convergence in normal approximation under moment conditions via new bounds on solutions of the Stein equation. *J. Theoret. Probab.* **29** 231–247. [MR3463084](#)
- [31] GELMAN, A. (2006). Multilevel (hierarchical) modeling: What it can and cannot do. *Technometrics* **48** 432–435. [MR2252307](#)
- [32] GELMAN, A., CARLIN, J. B., STERN, H. S., DUNSON, D. B., VEHTARI, A. and RUBIN, D. B. (2014). *Bayesian Data Analysis*, 3rd ed. *Texts in Statistical Science Series*. CRC Press, Boca Raton, FL. [MR3235677](#)
- [33] GIROLAMI, M. and CALDERHEAD, B. (2011). Riemann manifold Langevin and Hamiltonian Monte Carlo methods. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **73** 123–214. [MR2814492](#)
- [34] GLAESER, G. (1958). Étude de quelques algèbres tayloriennes. *J. Anal. Math.* **6** 1–124; erratum, insert to 6 (1958), no. 2. [MR0101294](#)
- [35] GORHAM, J., DUNCAN, A. B., VOLLMER, S. J and MACKEY, L. (2019). Supplement to “Measuring sample quality with diffusions.” DOI:[10.1214/19-AAP1467SUPP](https://doi.org/10.1214/19-AAP1467SUPP).
- [36] GORHAM, J. and MACKEY, L. (2015). Measuring sample quality with Stein’s method. *Adv. NIPS* **28** 226–234.
- [37] GORHAM, J. and MACKEY, L. (2017). Measuring sample quality with kernels. In *Proc. of 34th ICML, ICML’17*.
- [38] GÖTZE, F. (1991). On the rate of convergence in the multivariate CLT. *Ann. Probab.* **19** 724–739. [MR1106283](#)
- [39] GRETTON, A., BORGWARDT, K., RASCH, M., SCHÖLKOPF, B. and SMOLA, A. (2006). A kernel method for the two-sample-problem. *Adv. NIPS* **19** 513–520.
- [40] GU, Z., ROTHBERG, E. and BIXBY, R. (2015). Gurobi optimizer reference manual. Available at <http://www.gurobi.com>.
- [41] GUDMUNDSSON, J., KLEIN, O., KNAUER, C. and SMID, M. (2007). Small manhattan networks and algorithmic applications for the Earth movers distance. In *Proc. 23rd EuroCG* 174–177.
- [42] HAIRER, M., STUART, A. M. and VOLLMER, S. J. (2014). Spectral gaps for a Metropolis–Hastings algorithm in infinite dimensions. *Ann. Appl. Probab.* **24** 2455–2490. [MR3262508](#)
- [43] HAR-PELED, S. and MENDEL, M. (2006). Fast construction of nets in low-dimensional metrics and their applications. *SIAM J. Comput.* **35** 1148–1184. [MR2217141](#)
- [44] HARTLEY, R. and ZISSERMAN, A. (2003). *Multiple View Geometry in Computer Vision*, 2nd ed. Cambridge Univ. Press, Cambridge. [MR2059248](#)
- [45] HOROWITZ, A. M. (1987). The second order Langevin equation and numerical simulations. *Nuclear Phys. B* **280** 510–522. [MR0881123](#)
- [46] HUBER, P. J. and RONCHETTI, E. M. (2009). *Robust Statistics*, 2nd ed. *Wiley Series in Probability and Statistics*. Wiley, Hoboken, NJ. [MR2488795](#)
- [47] HUGGINS, J. and ZOU, J. (2017). Quantifying the accuracy of approximate diffusions and Markov chains. In *Proc. 20th AISTATS* 382–391.
- [48] HUGGINS, J. H. and MACKEY, L. (2018). Random feature Stein discrepancies. In *Adv. NIPS* **31**.

- [49] HWANG, C.-R., HWANG-MA, S.-Y. and SHEU, S. J. (1993). Accelerating Gaussian diffusions. *Ann. Appl. Probab.* **3** 897–913. [MR1233633](#)
- [50] JOULIN, A. and OLLIVIER, Y. (2010). Curvature, concentration and error estimates for Markov chain Monte Carlo. *Ann. Probab.* **38** 2418–2442. [MR2683634](#)
- [51] KALLENBERG, O. (2002). *Foundations of Modern Probability*, 2nd ed. *Probability and Its Applications (New York)*. Springer, New York. [MR1876169](#)
- [52] KENT, J. (1978). Time-reversible diffusions. *Adv. in Appl. Probab.* **10** 819–835. [MR0509218](#)
- [53] KHAMINSKII, R. (2012). *Stochastic Stability of Differential Equations*, 2nd ed. *Stochastic Modelling and Applied Probability* **66**. Springer, Heidelberg. [MR2894052](#)
- [54] KORATTIKARA, A., CHEN, Y. and WELLING, M. (2014). Austerity in MCMC land: Cutting the Metropolis–Hastings budget. In *Proc. of 31st ICML, ICML'14*.
- [55] LACOSTE-JULIEN, S., LINDSTEN, F. and BACH, F. (2015). Sequential kernel herding: Frank–Wolfe optimization for particle filtering. In *AISTATS*.
- [56] LANDIM, C., OLLA, S. and YAU, H. T. (1998). Convection–diffusion equation with space-time ergodic random flow. *Probab. Theory Related Fields* **112** 203–220. [MR1653837](#)
- [57] LEY, C., REINERT, G. and SWAN, Y. (2017). Stein’s method for comparison of univariate distributions. *Probab. Surv.* **14** 1–52. [MR3595350](#)
- [58] LIU, C. (1996). Bayesian robust multivariate linear regression with incomplete data. *J. Amer. Statist. Assoc.* **91** 1219–1227. [MR1424619](#)
- [59] LIU, Q. and LEE, J. (2017). Black-box importance sampling. In *Proc. 20th AISTATS* 952–961.
- [60] LIU, Q., LEE, J. and JORDAN, M. (2016). A kernelized Stein discrepancy for goodness-of-fit tests. In *Proc. of 33rd ICML, ICML* **48** 276–284.
- [61] LIU, Q. and WANG, D. (2016). Stein variational gradient descent: A general purpose Bayesian inference algorithm. *Adv. NIPS* **29** 2378–2386.
- [62] LUBIN, M. and DUNNING, I. (2015). Computing in operations research using Julia. *INFORMS J. Comput.* **27** 238–248. [MR3347876](#)
- [63] LUNARDI, A. (2018). *Interpolation Theory*. Springer, Berlin.
- [64] MA, Y., CHEN, T. and FOX, E. (2015). A complete recipe for stochastic gradient MCMC. *Adv. NIPS* **28** 2899–2907.
- [65] MACKEY, L. and GORHAM, J. (2016). Multivariate Stein factors for a class of strongly log-concave distributions. *Electron. Commun. Probab.* **21** 56. [MR3548768](#)
- [66] MANCA, L. (2008). *Kolmogorov Operators in Spaces of Continuous Functions and Equations for Measures. Tesi. Scuola Normale Superiore di Pisa (Nuova Series) [Theses of Scuola Normale Superiore di Pisa (New Series)]* **10**. Edizioni della Normale, Pisa. [MR2487956](#)
- [67] MATTINGLY, J. C., STUART, A. M. and TRETYAKOV, M. V. (2010). Convergence of numerical time-averaging and stationary measures via Poisson equations. *SIAM J. Numer. Anal.* **48** 552–577. [MR2669996](#)
- [68] MECKES, E. (2009). On Stein’s method for multivariate normal approximation. In *High Dimensional Probability V: The Luminy Volume. Inst. Math. Stat. (IMS) Collect.* **5** 153–178. IMS, Beachwood, OH. [MR2797946](#)
- [69] MÜLLER, A. (1997). Integral probability metrics and their generating classes of functions. *Adv. in Appl. Probab.* **29** 429–443. [MR1450938](#)
- [70] NOURDIN, I., PECCATI, G. and RÉVEILLAC, A. (2010). Multivariate normal approximation using Stein’s method and Malliavin calculus. *Ann. Inst. Henri Poincaré Probab. Stat.* **46** 45–58. [MR2641769](#)
- [71] OATES, C. J., GIROLAMI, M. and CHOPIN, N. (2017). Control functionals for Monte Carlo integration. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **79** 695–718. [MR3641403](#)
- [72] ØKSENDAL, B. (2003). *Stochastic Differential Equations: An Introduction with Applications*, 6th ed. *Universitext*. Springer, Berlin. [MR2001996](#)
- [73] PARDOUX, E. and VERETENNIKOV, A. Y. (2001). On the Poisson equation and diffusion approximation. I. *Ann. Probab.* **29** 1061–1085. [MR1872736](#)



- [74] PATTERSON, S. and TEH, Y. (2013). Stochastic gradient Riemannian Langevin dynamics on the probability simplex. *Adv. NIPS* **26** 3102–3110.
- [75] PAVLIOTIS, G. A. (2014). *Stochastic Processes and Applications. Texts in Applied Mathematics* **60**. Springer, New York. Diffusion processes, the Fokker–Planck and Langevin equations. [MR3288096](#)
- [76] PELEG, D. and SCHÄFFER, A. A. (1989). Graph spanners. *J. Graph Theory* **13** 99–116. [MR0982872](#)
- [77] PROTTER, P. E. (2005). *Stochastic Integration and Differential Equations. Stochastic Modelling and Applied Probability* **21**. Springer, Berlin. Second edition. Version 2.1, Corrected third printing. [MR2273672](#)
- [78] RAIČ, M. (2004). A multivariate CLT for decomposable random vectors with finite second moments. *J. Theoret. Probab.* **17** 573–603. [MR2091552](#)
- [79] RANGANATH, R., TRAN, D., ALTOSAAR, J. and BLEI, D. (2016). Operator variational inference. In *Advances in Neural Information Processing Systems* 496–504.
- [80] REINERT, G. and RÖLLIN, A. (2009). Multivariate normal approximation with Stein’s method of exchangeable pairs under a general linearity condition. *Ann. Probab.* **37** 2150–2173. [MR2573554](#)
- [81] REY-BELLET, L. and SPILIOPOULOS, K. (2015). Irreversible Langevin samplers and variance reduction: A large deviations approach. *Nonlinearity* **28** 2081–2103. [MR3366637](#)
- [82] RISKEN, H. (1989). *The Fokker–Planck Equation: Methods of Solution and Applications*, 2nd ed. *Springer Series in Synergetics* **18**. Springer, Berlin. [MR0987631](#)
- [83] ROBERTS, G. O. and STRAMER, O. (2002). Langevin diffusions and Metropolis–Hastings algorithms. *Methodol. Comput. Appl. Probab.* **4** 337–357. [MR2002247](#)
- [84] ROBERTS, G. O. and TWEEDIE, R. L. (1996). Exponential convergence of Langevin distributions and their discrete approximations. *Bernoulli* **2** 341–363. [MR1440273](#)
- [85] RÖCKNER, M., SOBOL, Z., et al. (2006). Kolmogorov equations in infinite dimensions: Well-posedness and regularity of solutions, with applications to stochastic generalized Burgers equations. *Ann. Probab.* **34** 663–727. [MR2223955](#)
- [86] SHVARTSMAN, P. (2008). The Whitney extension problem and Lipschitz selections of set-valued mappings in jet-spaces. *Trans. Amer. Math. Soc.* **360** 5529–5550. [MR2415084](#)
- [87] SOFFRITTI, G. and GALIMBERTI, G. (2011). Multivariate linear regression with non-normal errors: A solution based on mixture models. *Stat. Comput.* **21** 523–536. [MR2826690](#)
- [88] STEIN, C. (1972). A bound for the error in the normal approximation to the distribution of a sum of dependent random variables. In *Proc. Sixth Berkeley Symposium on Mathematical Statistics and Probability, Vol. II: Probability Theory* 583–602. Univ. California Press, Berkeley, CA. [MR0402873](#)
- [89] STEIN, C., DIACONIS, P., HOLMES, S. and REINERT, G. (2004). Use of exchangeable pairs in the analysis of simulations. In *Stein’s Method: Expository Lectures and Applications. Institute of Mathematical Statistics Lecture Notes—Monograph Series* **46** 1–26. IMS, Beachwood, OH. [MR2118600](#)
- [90] STUART, A. M., VOSS, J., WIBERG, P., et al. (2004). Conditional path sampling of SDEs and the Langevin MCMC method. *Commun. Math. Sci.* **2** 685–697.
- [91] TEH, Y. W., THIERY, A. H. and VOLLMER, S. J. (2016). Consistency and fluctuations for stochastic gradient Langevin dynamics. *J. Mach. Learn. Res.* **17** 7. [MR3482927](#)
- [92] VALLANDER, S. S. (1973). Calculations of the Vasserstein distance between probability distributions on the line. *Theory Probab. Appl.* **18** 824–827. [MR0328982](#)
- [93] VOLLMER, S. J., ZYGALAKIS, K. C. and TEH, Y. W. (2016). Exploration of the (non-) asymptotic bias and variance of stochastic gradient Langevin dynamics. *J. Mach. Learn. Res.* **17** 159. [MR3555050](#)
- [94] WANG, F. (2016). Exponential contraction in Wasserstein distances for diffusion semigroups with negative curvature. Available at <https://arxiv.org/abs/1603.05749>.



- [95] ZELLNER, A. (1976). Bayesian and non-Bayesian analysis of the regression model with multivariate Student- $t$  error terms. *J. Amer. Statist. Assoc.* **71** 400–405. [MR0405699](#)
- [96] ZELLNER, A. and MIN, C. (1995). Gibbs sampler convergence criteria. *J. Amer. Statist. Assoc.* **90** 921–927.

J. GORHAM  
STANFORD UNIVERSITY  
SEQUOIA HALL  
STANFORD, CALIFORNIA 94305  
USA  
E-MAIL: [jacksongorham@gmail.com](mailto:jacksongorham@gmail.com)

S. J. VOLLMER  
MATHEMATICS INSTITUTE  
UNIVERSITY OF WARWICK  
COVENTRY CV4 7AL  
UNITED KINGDOM  
E-MAIL: [svollmer@turing.ac.uk](mailto:svollmer@turing.ac.uk)

A. B. DUNCAN  
DEPARTMENT OF MATHEMATICS  
IMPERIAL COLLEGE LONDON  
LONDON SW7 2AZ  
UNITED KINGDOM  
E-MAIL: [a.duncan@imperial.ac.uk](mailto:a.duncan@imperial.ac.uk)

L. MACKEY  
MICROSOFT RESEARCH NEW ENGLAND  
CAMBRIDGE, MASSACHUSETTS 02474  
USA  
E-MAIL: [lmackey@microsoft.com](mailto:lmackey@microsoft.com)