IMPLICIT REGULARISATION IN DIFFUSION MODELS: AN ALGORITHM-DEPENDENT GENERALISATION ANALYSIS

Anonymous authors

000

001

002

004

006

008 009 010

011

013

014

015

016

017

018

019

021

023

024

025

026

027

028

029

031

033 034

035

037

038

040

041 042

043

044

045

046 047

048

051

052

Paper under double-blind review

ABSTRACT

The success of denoising diffusion models raises important questions regarding their generalisation behaviour, particularly in high-dimensional settings. Notably, it has been shown that when training and sampling are performed perfectly, these models memorise training data—implying that some form of regularisation is essential for generalisation. Existing theoretical analyses primarily rely on algorithmindependent techniques such as uniform convergence, heavily utilising model structure to obtain generalisation bounds. In this work, we instead leverage the algorithmic aspects that promote generalisation in diffusion models, developing a general theory of algorithm-dependent generalisation for this setting. Borrowing from the framework of algorithmic stability, we introduce the notion of score stability, which quantifies the sensitivity of score-matching algorithms to dataset perturbations. We derive generalisation bounds in terms of score stability, and apply our framework to several fundamental learning settings, identifying sources of regularisation. In particular, we consider denoising score matching with early stopping (denoising regularisation), sampler-wide coarse discretisation (sampler regularisation) and optimising with SGD (optimisation regularisation). By grounding our analysis in algorithmic properties rather than model structure, we identify multiple sources of implicit regularisation unique to diffusion models that have so far been overlooked in the literature.

1 Introduction

Diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song et al., 2021) are a class of generative models that have achieved state-of-the-art performance across image, audio, video, and protein synthesis tasks (Rombach et al., 2022; Saharia et al., 2022; Ramesh et al., 2022; Watson et al., 2023; Esser et al., 2024). Their ability to generate high-quality samples from complex, high-dimensional distributions with limited data motivates the need for a theoretical understanding of the mechanisms underpinning their strong generalisation capabilities.

The goal of diffusion models is to generate new synthetic samples from a data distribution $\nu_{\rm data}$ using a finite set of N data points $\{x_i\}_{i=1}^N$. Central to the methodology is a unique approach to generating data, formulating it as the iterative transformation of noise into data, or equivalently, the reversal of a diffusion process (Song et al., 2021). This diffusion process, called the *forward process*, is defined by the stochastic differential equation (SDE),

$$dX_t = -\alpha X_t dt + \sqrt{2} dW_t, \qquad X_0 \sim \nu_{\text{data}}, \qquad t \in [0, T], \tag{1}$$

for some $\alpha \geq 0$, where W_t denotes the Brownian motion in \mathbb{R}^d and T > 0 is the terminal time. It can then be shown that the time-reversal of this process, $Y_t := X_{T-t}$ admits a weak formulation as a solution to the SDE,

$$dY_t = \alpha Y_t dt + 2\nabla \log p_{T-t}(Y_t) dt + \sqrt{2} dW_t, \qquad Y_0 \sim p_T, \qquad t \in [0, T), \tag{2}$$

where p_t denotes the marginal density of X_t (Haussmann & Pardoux, 1986). Therefore, simulating samples from $\nu_{\rm data} = p_0$ can be achieved by solving the diffusion process in (2), which requires

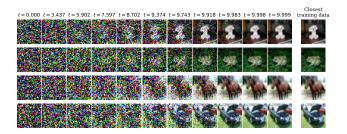


Figure 1: Samples generated using the empirical score function on CIFAR-10 compared to the closest image in the dataset, illustrating memorisation of the training data.

an approximation of the score function, $\nabla \log p_t$. This is achieved by fitting a time-dependent deep neural network to minimise a weighted L^2 -distance called the *(population) score matching loss*:

$$\ell_{\text{sm}}(s;\tau) := \int \mathbb{E}_{X_t}[\|s(X_t, t) - \nabla \log p_t(X_t)\|^2] \, \tau(dt), \tag{3}$$

where τ is a probability measure over (0,T] that determines the weighting of the timepoints. Since $\nabla \log p_t$ is unknown, typically the (population) denoising score matching loss $\ell_{\rm dsm}$, which differs from $\ell_{\rm sm}(s)$ only by a constant, is used instead and is then approximated using the dataset, forming the empirical denoising score matching loss $\ell_{\rm dsm}$ (see equations (6) and (8)). The score network s(x,t) is trained on this objective using standard stochastic optimisation methods relying on minibatching. Once an approximation is obtained, samples are generated by numerically solving the reverse-time SDE, (2). Both score matching and sampling introduce distinct challenges and design choices that impact the quality of model output (Karras et al., 2022).

Score matching presents a key difference from standard supervised learning. In the space of all L^2 score functions, the empirical objective $\hat{\ell}_{\rm dsm}$ possesses a unique minimiser—the empirical score function—as a result of the integration over $X_t|X_0$ (see Lemma 1). This contrasts with traditional supervised learning, where the empirical risk minimisation problem can have infinitely many solutions (e.g., in overparameterised regression) and often requires regularisation to be well-posed. As shown in Figure 1, sampling with this empirical score leads to exact recovery of the training data (Pidstrigach, 2022). This behaviour is distinct from 'benign overfitting', a phenomenon from the deep learning literature where interpolating the data does not necessarily prevent generalisation (Bartlett et al., 2021; Zhang et al., 2021). This divergence suggests that existing theory may be insufficient to explain the success of diffusion models, highlighting the need for new frameworks tailored to this setting.

Recently, there has been a drive towards developing theory for better understanding the unique structure of diffusion models. The most developed subset of this work focuses on connecting sample quality to score matching by deriving upper bounds on distribution error (e.g. KL divergence, total variation, or Wasserstein distance) between model samples and the data distribution, controlling it by the population score matching loss (De Bortoli et al., 2021; De Bortoli, 2022; Lee et al., 2022; Chen et al., 2023; Benton et al., 2024; Potaptchik et al., 2024). These results, often referred to as *convergence bounds*, typically take the form,

Distribution error
$$\lesssim \ell_{\rm sm}(s) + \Delta$$
,

where Δ is the discretisation error of the sampling scheme, which can be made small with sufficiently fine discretisation. However, since $\ell_{\rm sm}$ is not computable, these bounds say little about performance under empirical guarantees—that is, their *generalisation* properties. One line of work, initiated by Oko et al. (2023) and extended in (Azangulov et al., 2024; Tang & Yang, 2024), applies classical uniform convergence theory to bound the generalisation gap from the decomposition,

$$\ell_{\rm sm}(s) = \hat{\ell}_{\rm sm}(s) + \underbrace{\ell_{\rm sm}(s) - \hat{\ell}_{\rm sm}(s)}_{\text{generalisation gap}},\tag{4}$$

where $\hat{\ell}_{sm}$ denotes the empirical counterpart to ℓ_{sm} . These results rely on covering number bounds for specific classes of neural networks and, while informative, they are limited to carefully chosen model classes and do not account for algorithmic properties. An alternative approach by De Bortoli

(2022) uses a decomposition of the Wasserstein distance that leverages convergence properties of the empirical measure. Though more model-agnostic, this method overlooks how diffusion models generate *novel* data. Both lines of work are fundamentally algorithm-independent, in that they lack any utilisation of the algorithmic aspects that uniquely define diffusion models. Recent efforts aim to incorporate algorithmic effects by restricting the problem. For instance, Shah et al. (2023); Chen et al. (2024) consider Gaussian mixture targets, while Li et al. (2023); Yang (2022) study random feature models. These settings allow for finer analysis of the role of the score matching algorithm, but remain limited in scope, leaving open the challenge of developing a more general algorithm-dependent theory of generalisation in diffusion models.

As noted earlier, if the empirical score matching loss was completely minimised and sampling was performed perfectly, the diffusion model would simply return training data, failing to generalise. Therefore, the observed success of diffusion models in producing novel data implies that, in practice, they either avoid completely minimising $\hat{\ell}_{\rm sm}$ or must avoid perfectly sampling. This suggests that (implicit) regularisation in the score matching or sampling algorithm is crucial for generalisation, making algorithmic considerations essential for understanding diffusion models.

1.1 OUR CONTRIBUTIONS

We introduce score stability, a general, algorithm-dependent framework for analysing diffusion model generalisation based on the classical approach of algorithmic stability. This framework quantifies an algorithm's dependence on individual training examples, from which we derive expected generalisation gap bounds for score matching losses. Using the score stability framework, we then analyse several examples of score matching algorithms, identifying three distinct sources of implicit regularisation in diffusion model training and sampling: noising, sampler, and optimisation-induced regularisation.

Denoising regularisation. To begin with, we consider the empirical risk minimisation algorithm (ERM) that minimises $\hat{\ell}_{\rm dsm}$ over a hypothesis class \mathcal{H} . Through a score stability analysis, we reveal a regularisation source within this objective when early stopping of the forward process is used—a standard practice in the diffusion model literature. Utilising properties of the noising forward process, we obtain generalisation gap bounds with near-linear rate, $\epsilon^{-d^*/4}(\epsilon^{-d^*/2}N^{-2} + \min_{\mathcal{H}}\hat{\ell}_{\rm sm})^{c/2}$ for any c < 1, where $\epsilon > 0$ is the early stopping time and d^* is the dimension of the data support.

Sampler regularisation. We then apply this analysis to discrete-time sampling algorithms, deriving statistical guarantees for the expected KL divergence between the true data distribution and samples generated by the diffusion model. The bound we derive is formed of two stages: we obtain generic rates $\epsilon^{-1/2}(\epsilon^{-d^*/2}N^{-2} + \min_{\mathcal{H}}\hat{\ell}_{\rm sm})^{c/d^*}$ but when N^{-2} and $\min_{\mathcal{H}}\hat{\ell}_{\rm sm}$ are sufficiently small relative to ϵ , we obtain bounds with rates $\epsilon^{-d^*/4}(\epsilon^{-d^*/2}N^{-2} + \min_{\mathcal{H}}\hat{\ell}_{\rm sm})^{c/2}$ that are faster in N and $\min_{\mathcal{H}}\hat{\ell}_{\rm sm}$. To derive this bound, we utilise regularisation brought about by the coarseness of the discretisation. We find that by increasing discretisation coarseness, we can improve the generalisation gap bound at the expense of worsening the discretisation error term.

Optimisation regularisation. Finally, we consider the role of the optimisation scheme, analysing stochastic gradient descent (SGD) with gradient clipping and weight decay. On the model class, we assume only structural assumptions typical in the optimisation literature, including non-global Lipschitz and smoothness assumptions. While this initially yields bounds that grow with the number of iterations, we more closely inspect the impact of the high-variance gradient estimator used in diffusion training. We show this gradient noise induces a contractive behaviour in the training dynamics, which we harness to obtain stability bounds that do not grow with the number of iterations (Proposition 14), showing that the noisy dynamics enable tighter generalisation guarantees.

2 BACKGROUND

Suppose that the data distribution $\nu_{\rm data}$ is on \mathbb{R}^d and we are provided a finite dataset of samples $S=\{x_1,...,x_N\}$ which we assume are sampled independently and identically (i.i.d.) from $\nu_{\rm data}$. As discussed in the introduction, diffusion models are formed of two distinct stages. The first stage, score matching, consists of learning an approximation to the score function $\nabla \log p_t$ using the dataset S. In this work, we take a score function to be any function belonging to the set $L^0(\mathbb{R}^d \times [0,T];\mathbb{R}^d)$, the

set of Borel measurable functions of the form $\mathbb{R}^d \times [0,T] \to \mathbb{R}^d$. Then, a *score matching algorithm* is taken to be any mapping of the form $A_{\mathrm{sm}}: (\cup_{N=1}^\infty (\mathbb{R}^d)^{\otimes N}) \times \Omega \to \mathcal{H}$ where \mathcal{H} is a measurable subset of $L^0(\mathbb{R}^d \times [0,T];\mathbb{R}^d)$. Here, Ω is the event space belonging to a probability space $(\Omega,\mathcal{F},\mathbb{P})$.

The second stage of diffusion models, sampling, consists of generating samples with the learned score function. We take a sampling algorithm to be a mapping of the form, $A_{\mathrm{samp}}: \mathcal{H} \to \mathcal{P}(\mathbb{R}^d)$ where $\mathcal{P}(\mathbb{R}^d)$ denotes the set of Borel measures on \mathbb{R}^d . Typically, sampling is performed using an approximation to the reverse process given in (2), replacing $\nabla \log p_t$ with the learned score function $s(\cdot,t)$ and replacing its initial distribution, p_T with a data-independent prior, $p_{\mathrm{prior}} = \mathcal{N}(\mathbf{0}, \sigma_{\mathrm{prior}}^2 I)$. In the case of $\alpha > 0$, we choose $\sigma_{\mathrm{prior}}^2 = \alpha^{-1}$ so that the prior coincides with the stationary distribution of the forward process, and when $\alpha = 0$ we simply set $\sigma_{\mathrm{prior}}^2 = 2T$. With this, we arrive at the SDE,

$$d\hat{Y}_t = \alpha \hat{Y}_t dt + 2s(\hat{Y}_t, T - t) dt + \sqrt{2} dW_t, \quad \hat{Y}_0 \sim p_{\text{prior}}.$$
 (5)

Thus, a sample is generated by sampling from \hat{Y}_T , or more commonly, the process is terminated early, sampling from $\hat{Y}_{T-\epsilon}$ for some small $\epsilon > 0$. Therefore, diffusion models are density estimation algorithms formed from the composition $A_{\mathrm{samp}} \circ A_{\mathrm{sm}}$.

Denoising score matching and overfitting. As stated in the introduction, computing $\ell_{\rm sm}$ requires access to the population score function, $\nabla \log p_t$. So instead, the *(population) denoising score matching loss* is used in its place:

$$\ell_{\text{dsm}}(s;\tau) := \mathbb{E}_{X_0 \sim \nu} \left[\int \mathbb{E}_{X_t \mid X_0} [\|s(X_t, t) - \nabla \log p_{t|0}(X_t \mid X_0)\|^2 |X_0|] \tau(dt) \right], \tag{6}$$

which differs from $\ell_{\rm sm}(s)$ only by a constant $C_{\rm sm}$, (see Lemma 16) whilst being easier to approximate without access to $\nabla \log p_t$ (Hyvärinen, 2005). Since $p_{t|0}$ is a Gaussian kernel, its score is given by,

$$\nabla_y \log p_{t|0}(y|x) = \frac{\mu_t x - y}{\sigma_t^2}, \qquad \mu_t = e^{-\alpha t}, \qquad \sigma_t^2 = \alpha^{-1} (1 - \mu_t^2). \tag{7}$$

In practice, the objective in (6) is further approximated via Monte Carlo estimation using the dataset which leads to the *empirical denoising score matching loss*,

$$\hat{\ell}_{\text{dsm}}(s; S, \tau) := \frac{1}{N} \sum_{i=1}^{N} \int \mathbb{E}_{X_t \mid X_0} [\|s(X_t, t) - \nabla \log p_{t \mid 0}(X_t \mid x_i)\|^2 |X_0 = x_i] \, \tau(dt). \tag{8}$$

In the following lemma, we highlight the important property that this can equivalently be defined as the denoising score matching objective for the process \hat{X}_t which evolves as in (1) but with the initial distribution given by the empirical distribution, $\hat{X}_0 \sim \frac{1}{N} \sum_{i=1}^N \delta_{x_i}(dx)$.

Lemma 1. The objective $\hat{\ell}_{dsm}(s; S, \tau)$ is identical, up to a constant, to the objective

$$\hat{\ell}_{\text{sm}}(s; S, \tau) := \int \mathbb{E}[\|s(\hat{X}_t, t) - \nabla \log \hat{p}_t(\hat{X}_t)\|^2 |S| \tau(dt), \tag{9}$$

where \hat{p}_t is the marginal density of \hat{X}_t . Therefore, any minimiser of $\hat{\ell}_{dsm}(\cdot; S, \tau)$ on $L^0(\mathbb{R}^d \times [0, T]; \mathbb{R}^d)$ is identical to $\nabla \log \hat{p}_t$ a.e. for any $t \in \text{supp}(\tau)$.

See Appendix A.2 for the proof. This lemma shows that, unlike in traditional supervised learning problems, the empirical objective here admits a single unique minimiser, the *empirical score function*, $\nabla \log \hat{p}_t$. The nature of this score function and the samples it generates has been the focus of several recent studies, notably (Pidstrigach, 2022) which shows that with perfect sampling, any score function sufficiently close to $\nabla \log \hat{p}_t$ recovers the training data.

Other notation. When the score matching algorithm $A_{\rm sm}$ is random, we use $A_{\rm sm}(S)$ as shorthand for the random score function $(x,t,\omega)\mapsto A_{\rm sm}(S,\omega)(x,t)$. Given two random score functions s,s', we let $\Gamma(s,s')$ denote the set of all couplings of these random functions (Appendix A.1 for details).

3 Score stability and generalisation

Algorithmic stability is a classical technique in learning theory used to understand the generalisation properties of a variety of important learning algorithms (Kearns & Ron, 1999; Devroye & Wagner,

1979; Bousquet & Elisseeff, 2002; Hardt et al., 2016). While there are various formulations, they all share the common aim of connecting properties of a learning algorithm to its robustness under changes in the dataset. Its use has primarily been focused around regression and classification problems—in this section, we propose a notion of stability that applies specifically to diffusion models.

We introduce the notion of *score stability* which quantifies how sensitive a score matching algorithm $A_{\rm sm}$ is to individual changes in the dataset. We do this by defining the adjacent dataset $S^i := \{x_1,...,x_{i-1},\tilde{x},x_{i+1},...,x_N\}$ where $\tilde{x} \sim \nu_{\rm data}$, independent from S, and then measuring the similarity between the score functions $\hat{s} = A_{\rm sm}(S)$ and $\hat{s}^i = A_{\rm sm}(S^i)$.

Definition 2. A score matching algorithm A_{sm} is score stable with constant $\varepsilon_{stab} > 0$ if for any $i \in [N]$ it holds that,

$$\mathbb{E}_{S,\tilde{x}} \left[\inf_{(\hat{s},\hat{s}^i) \in \Gamma_i} \int \mathbb{E}[\|\hat{s}(X_t,t) - \hat{s}^i(X_t,t)\|^2 | X_0 = \tilde{x}, S, \tilde{x}] \, \tau(dt) \right] \leq \varepsilon_{stab}^2,$$

where $\Gamma_i = \Gamma(A_{\rm sm}(S), A_{\rm sm}(S^i)).$

Since $A_{\rm sm}$ may be random, we define score stability in terms of the best-case coupling of the random score functions \hat{s}, \hat{s}^i . We recall that $\Gamma(\cdot, \cdot)$ denotes the set of couplings between two random score functions, and when it is not random, it is given by the singleton $\Gamma_i = \{(A_{\rm sm}(S), A_{\rm sm}(S^i))\}$. In the following theorem, we connect score stability to generalisation by controlling the expected generalisation gap by the score stability constant.

Theorem 3. Suppose that the score matching algorithm A_{sm} is score stable with constant ε_{stab} . Then, with $\hat{s} = A_{sm}(S)$, it holds that

$$\left| \mathbb{E} \left[\ell_{\mathrm{dsm}}(\hat{s}; \tau) \right]^{1/2} - \mathbb{E} \left[\hat{\ell}_{\mathrm{dsm}}(\hat{s}; S, \tau) \right]^{1/2} \right| \le \varepsilon_{\mathit{stab}}. \tag{10}$$

Furthermore, it holds that

$$\mathbb{E}[\ell_{\rm sm}(\hat{s};\tau)] - \mathbb{E}[\hat{\ell}_{\rm sm}(\hat{s};S,\tau)] \le 2\,\varepsilon_{\it stab}\,\mathbb{E}[\hat{\ell}_{\rm dsm}(\hat{s};S,\tau)]^{1/2} + \varepsilon_{\it stab}^2. \tag{11}$$

With Theorem 3, we obtain that the generalisation gap for both the denoising score matching loss and the score matching loss decays at the same rate as score stability. We can further simplify the bound for the score matching loss using the fact that $\hat{\ell}_{dsm}$ and $\hat{\ell}_{sm}$ are identical up to a constant, to obtain,

$$\mathbb{E}\big[\ell_{\mathrm{sm}}(\hat{s};\tau)\big] \lesssim \mathbb{E}\big[\hat{\ell}_{\mathrm{sm}}(\hat{s};\tau)\big] + \varepsilon_{\mathrm{stab}}\,C_{\mathrm{sm}}^{1/2} + \varepsilon_{\mathrm{stab}}^2.$$

One should expect that if the score matching algorithm is effective, both \hat{s} and \hat{s}^i converge to the ground truth as N grows, and thus $\varepsilon_{\text{stab}}$ should decrease to 0. Ascertaining the rate at which N decreases requires an analysis of the algorithm at hand, hence the categorisation of algorithmic stability as an algorithm-dependent approach. This contrasts with uniform learning, which utilises control over the hypothesis class, providing a worst-case bound that is independent from the algorithm.

In the following sections, we apply the framework of score stability to some common learning settings for diffusion models. We derive estimates of the score stability constant for these algorithms and identify features that promote generalisation.

4 EMPIRICAL SCORE MATCHING AND IMPLICIT REGULARISATION

We begin our examples by considering the score matching algorithm that minimises the empirical denoising score matching loss. Given a hypothesis class $\mathcal{H} \subseteq L^0(\mathbb{R}^d \times [0,T];\mathbb{R}^d)$, we define this algorithm by,

$$A_{\operatorname{erm}}(S) = \operatorname{argmin}_{s \in \mathcal{H}} \hat{\ell}_{\operatorname{dsm}}(s; S, \tau).$$

While this algorithm is not often used in practice, it is the natural analogue to empirical risk minimisation from traditional supervised learning and thus serves as a canonical example. We consider the setting of the manifold hypothesis where the data distribution is supported on a submanifold of \mathbb{R}^d .

Assumption 4. Suppose that ν_{data} is supported on a smooth submanifold of \mathbb{R}^d that has dimension d^* and reach $\tau_{\text{reach}} > 0$. Furthermore, its density on the submanifold, p_{ν} , satisfies $c_{\nu} := \inf p_{\nu} > 0$.

The reach describes the maximum distance where the projection to the manifold is uniquely defined and therefore, it quantifies the maximum curvature of the manifold. We refer to Appendix A.3 for the full definition. Several recent works have considered the assumption that $\nu_{\rm data}$ lies on a submanifold of \mathbb{R}^d . These works argue that d^* can often be far smaller than d and so dependence with respect to d^* over d is favourable (De Bortoli, 2022; Pidstrigach, 2022; Loaiza-Ganem et al., 2024; Potaptchik et al., 2024; Huang et al., 2024). The assumption that the density is bounded from below has also appeared in several of these works (Potaptchik et al., 2024; Huang et al., 2024). We also make the following assumption about the class of score networks.

Assumption 5. Suppose there exists $D_{\mathcal{H}} \geq 0$ such that for any $s, s' \in \mathcal{H}$, it holds that

$$||s(\cdot,t)-s'(\cdot,t)||_{L^{\infty}} \le D_{\mathcal{H}}/\sigma_t^2, \quad \text{for all } t \in \operatorname{supp}(\tau).$$

Under these assumptions, we obtain the following estimate for the stability constant.

Proposition 6. Suppose that assumptions 4 and 5 hold and that $\epsilon := \inf \operatorname{supp}(\tau) \in (0, \tau_{reach}^2)$, then for any $c \in (0, 1)$ and sufficiently large N, the score matching algorithm A_{erm} is score stable with,

$$\varepsilon_{\textit{stab}}^2 \lesssim C \big(C C_{\text{sm}} N^{-2} + \mathbb{E}[\hat{\ell}_{\text{sm}}(\hat{s})] \big)^c, \qquad C = \frac{D_{\mathcal{H}}^2}{\sigma_{\epsilon}^4} \vee \frac{1}{c_{\nu} \sigma_{\epsilon}^{d^*}}.$$

An interesting feature of Proposition 6 is that generalisation bounds under only basic assumptions about the structure of the hypothesis class *without any additional regularisation*. This contrasts with algorithmic stability in the setting of traditional supervised learning, where empirical risk minimisation is stable only when restricting the hypothesis class or with the use of explicit regularisation (Zhang et al., 2021; Bousquet & Elisseeff, 2002). Here, we show that the denoising score matching loss possesses the unique property that it is stable without the need for additional regularisation, suggesting that the denoising score matching loss possesses a form of *implicit regularisation*.

When $d^*>4$, for ϵ sufficiently small, we have that $C=\mathcal{O}(c_{\nu}^{-1}\epsilon^{-d^*/2}), C_{\mathrm{sm}}=\mathcal{O}(d^*\epsilon^{-1})$. Since the bound only depends on d^* and not d, this suggests that diffusion models are automatically manifold-adaptive. The bound also heavily depends on ϵ , with it being smaller for larger ϵ and growing exponentially fast as ϵ approaches zero, indicating that the natural regularisation present in the score matching objective is more prevalent at larger noise scales. The requirement to have $\epsilon>0$ is closely related to the technique of early stopping which is frequently used in the diffusion model literature (Song & Kingma, 2021; Karras et al., 2022). This is where the backwards process \hat{Y}_t is terminated early by some small amount of time to avoid irregularity issues of the score function when close to convergence. Other theoretical works have also identified the importance of early stopping in the generalisation properties of diffusion models (Oko et al., 2023; Azangulov et al., 2024).

Proof summary We now provide a brief summary of the proof of Proposition 6. The first step of the proof technique utilises a fundamental property of the empirical denoising score matching objective, $\hat{\ell}_{\rm dsm}(s;S,\tau)$: that it is strongly convex in s in a data-dependent weighted L^2 -space. Strong convexity is often used in algorithmic stability analyses, especially in deriving stability bounds for linear models—here we borrow a similar approach, but we analyse the stability of the algorithm in function space. With this, we arrive at the following inequality (see Lemma 19):

$$\int \mathbb{E}[\|\hat{s}(\hat{X}_t, t) - \hat{s}^i(\hat{X}_t, t)\|^2] \tau(dt) \lesssim \mathbb{E}[\hat{\ell}_{sm}(\hat{s})] + \frac{\varepsilon_{stab}}{N} (C_{sm}^{1/2} + \varepsilon_{stab}), \tag{12}$$

where $\varepsilon_{\text{stab}}$ is the (yet-to-be bounded) score stability constant of A_{erm} .

The second step of the proof technique utilises a characteristic property of the heat kernel—that it smooths out functions. In particular, we utilise the celebrated Harnack inequality of Wang (1997) that captures this property by showing that for any positive measurable $\phi: \mathbb{R}^d \to \mathbb{R}_+$, $x,y \in \mathbb{R}^d$, it holds that

$$\mathbb{E}[\phi(X_t)|X_0 = x] \le \mathbb{E}[\phi(X_t)^p | X_0 = y]^{1/p} \exp\left(\frac{\mu_t^2 ||x - y||^2}{2(p - 1)\sigma_t^2}\right),$$

for any t > 0, p > 1. Utilising this bound, we convert the upper bound in (12) to a bound on the stability constant. The full proof can be found in Appendix C.

5 STOCHASTIC SAMPLING AND SCORE STABILITY

In practice, the backwards process in (5) cannot be sampled exactly, so we instead rely on approximations based on numerical integration schemes. In this section, we investigate how algorithmic

stability interacts with these sampling schemes. We consider the Euler-Maruyama-type sampling scheme proposed in (Benton et al., 2024; Potaptchik et al., 2024) which discretises at the timesteps $(t_k)_{k=0}^K$, where $t_k = T - (1+\kappa)^{\frac{T-1}{\kappa}-k}$ for large k. The quantity $\kappa > 0$ is chosen freely and we choose $K \sim \log(\epsilon^{-1})/\log(1+\kappa)$ so that $t_K \approx \epsilon$ (see Appendix D for details). By sampling its terminating iterate \hat{y}_K , we obtain a *sampling algorithm*, $A_{\rm em}$, that maps a score function s to the distribution $\log(\hat{y}_K)$, which approximates the distribution $\log(\hat{Y}_{\epsilon})$.

In the previous section, we identified that early stopping of the backwards process benefits generalisation. In the present section, we will consider how coarseness of the discretisation scheme produces similar benefits. It is often the case that the score function is trained only at those time steps considered by the sampler, i.e. using the time-weighting, $\hat{\tau}_{\kappa}(dt) = \frac{1}{K} \sum_{k=0}^{K-1} \delta_{T-t_k}(dt)$ (Ho et al., 2020). As a result, the effective stopping time of the algorithm can be much larger than the early stopping time, ϵ . In the following proposition, we demonstrate how this benefits generalisation.

Proposition 7. Consider the setting of Proposition 6 with $\alpha = 1$ and set $\tau = \hat{\tau}_{\kappa}$, then for sufficiently large N, $\kappa \leq \epsilon^{-1}/4$ and any $c \in (0,1)$, we have that for $q_K = A_{\rm em} \circ A_{\rm erm}(S)$,

$$\begin{split} \mathbb{E}[D(p_{\epsilon}\|q_K)] \lesssim \mathbb{E}[\hat{\ell}_{\mathrm{sm},\kappa}^{\star}] + B_{\kappa}^{\frac{1}{2}}(1+\kappa)^{-d^*} + \frac{B_{\kappa}}{C_{\mathrm{sm}}}(1+\kappa)^{-2d^*} + \kappa(1+\kappa)d^*\log(\epsilon^{-1})^2 + de^{-2T}, \\ where \ B_{\kappa} &= \frac{C_{\mathrm{sm}}}{c_{\nu}} \big(\frac{C_{\mathrm{sm}}}{c_{\nu}}N^{-2} + \mathbb{E}[\hat{\ell}_{\mathrm{sm},\kappa}^{\star}]\big)^c \epsilon^{-d^*}, \ \hat{\ell}_{\mathrm{sm},\kappa}^{\star} &:= \inf_{\mathcal{H}} \hat{\ell}_{\mathrm{sm}}(h; S, \hat{\tau}_{\kappa}). \end{split}$$

The second and third terms of the bound in Proposition 7 are due to the score stability of the ERM algorithm and decay as κ increases. The fourth term of the bound captures the discretisation error and therefore increases with κ . What this result captures is that there is a trade-off between sampler accuracy and generalisation that is managed by the discretisation of the diffusion model. In the following corollary, this trade-off is optimised.

Corollary 8. Consider the setting of Proposition 7, then for any $c \in (0,1)$ and sufficiently small ϵ , there exists $\kappa > 0$ such that with $q_K = A_{\rm em} \circ A_{\rm erm}(S)$

$$\mathbb{E}[D(p_{\epsilon}\|q_K)] \lesssim \begin{cases} B_{\kappa}^{\frac{1}{2}} + C_{\text{sm}}^{-1} B_{\kappa}, & \text{if } B_{\kappa} \leq \log(\epsilon^{-1})^2, \\ \log(\epsilon^{-1}) B_{\kappa}^{\frac{1}{2(d^*+1)}} + (C_{\text{sm}}^{-1} + d^*) \log(\epsilon^{-1})^2 B_{\kappa}^{\frac{1}{d^*+1}} + de^{-2T}, & \text{otherwise.} \end{cases}$$

The primary strength of this result over (Oko et al., 2023; Azangulov et al., 2024) is that we assume little about the hypothesis class. Their results require carefully constrained network architectures and a specific early stopping time to control complexity. In contrast, our result holds for any sufficiently small early stopping time, relying instead on a carefully chosen discretisation scheme, which is usually tuned in practice (Karras et al., 2022; Williams et al., 2024). The main drawback is that our general approach does not exploit the model class to adapt to smoothness properties of the underlying measure, which we leave for future work.

6 STOCHASTIC OPTIMISATION AND IMPLICIT REGULARISATION

To learn the score function, it is common to choose it from a parametric hypothesis class $\{s_{\theta} : \theta \in \mathbb{R}^n\}$ (e.g. a deep neural network) by minimising $\hat{\ell}_{\rm dsm}$ via stochastic optimisation (Karras et al., 2024). In this section, we consider the score stability of this setting, focusing on stochastic gradient descent (SGD) with gradient clipping and weight decay. We consider the standard gradient estimator: given the mini-batch $(x_i')_{i=1}^{N_B}$ of size $N_B \ll N$ we define the random estimator,

$$G(\theta, (x_i')_{i=1}^{N_B}) = \frac{1}{N_B P} \sum_{i=1}^{N_B} \sum_{j=1}^{P} w_{t_{i,j}} \nabla_{\theta} \|s_{\theta}(X_{i,j}, t_{i,j}) - \nabla \log p_{t_{i,j}}|_0(X_{i,j}|x)\|^2,$$
(13)

where we define the random variables $X_{i,j} = \mu_{t_{i,j}} x_i' + \sigma_{t_{i,j}} \xi_{i,j}$, $t_{i,j} \sim w_t^{-1} \tau(dt)$, $\xi_{i,j} \sim N(0,I_d)$. The additional variance introduced by the random variables $\xi_{i,j}$ and $t_{i,j}$ leads to a gradient estimator with significantly higher variance than in standard supervised learning. This presents several challenges during training, and various strategies have been proposed to mitigate this issue (Karras et al., 2024; Song & Kingma, 2021). For example, the weighting function $w:[0,T]\to\mathbb{R}_+$ can be tuned to reduce variance (Karras et al., 2022) or the number of resamples $P\in\mathbb{N}$ can be increased. We consider the following iterative scheme, defined for a given weight decay constant $\lambda>0$ and clipping value C>0:

$$\theta_{k+1} = (1 - \eta_k \lambda)\theta_k - \eta_k \operatorname{Clip}_C(G_k(\theta_k, (x_i)_{i \in B_k})), \tag{14}$$

where $\eta_k > 0$ and $B_k \subset [N]$ are the learning rates and mini-batch indices for each iteration $k \in \mathbb{N}$ and we define the clipping operator $\mathrm{Clip}_C(v) = (1 \wedge (C\|v\|^{-1}))v$. Both gradient clipping and weight decay are widely used in diffusion model training and are typically motivated by their stabilising effect on optimisation, minimising the impact of the high variance of the gradient estimator (Song et al., 2021; Ho et al., 2020). Throughout this section, we take the mini-batch B_k to be i.i.d. and uniformly sampled from [N] without replacement. For the sake of simplicity, we suppose that the iterative scheme is terminated after $K \in \mathbb{N}$ iterations, where K is fixed and independent of the data.

6.1 STABILITY OF SGD WITH WEIGHT DECAY AND CLIPPING

In our analysis, we avoid restricting the score network to a specific parametric class and instead make structural assumptions based on its smoothness properties. We recall that a function is Lipschitz with constant $L \ge 0$ if it is differentiable and its directional derivatives are uniformly bounded by L.

Assumption 9 (Smoothness of the score network). There exists $L: \mathbb{R}^d \times (0,T] \to \mathbb{R}_+$ and $M: \mathbb{R}^d \times (0,T] \to \mathbb{R}_+$ such that for almost all $x \in \mathbb{R}^d$, $t \in (0,T]$, $s_{\theta}(x,t)$ is Lipschitz and smooth (gradient-Lipschitz) in $\theta \in \mathbb{R}^n$ with constants L(x,t) and M(x,t), respectively. Furthermore, there exists constants $\overline{L}, \overline{M} \geq 0$ such that for any $x \in \operatorname{supp}(\nu_{\text{data}})$,

$$\int \mathbb{E}[L(X_t, t)^2 | X_0 = x] \, \tau(dt) \le \overline{L}^2, \quad \int \mathbb{E}[M(X_t, t)^2 | X_0 = x] \, \tau(dt) \le \overline{M}^2.$$

The use of Lipschitz and smoothness assumptions is commonplace in the analysis of optimisation schemes (Nesterov, 2018; Hardt et al., 2016). However, the assumption differs slightly from the usual in that we only require these properties to hold *almost everywhere* with respect to the input distribution and we allow the Lipschitz and smoothness constants to vary with the input, provided their square averages remain bounded. This relaxation enables us to accommodate common models that would otherwise violate global smoothness assumptions, such as ReLU networks.

Assumption 10. Suppose there exists $B_{\ell} > 0$ such that for any $\theta \in \mathbb{R}^n$, it holds that

$$\hat{\ell}_{\text{dsm}}(s_{\theta}; \{x\}, \delta_t) \le B_{\ell}^2 / \sigma_t^4, \qquad \text{for each } x \in \text{supp}(\nu_{\text{data}}), t \in \text{supp}(\tau).$$
 (15)

This property requires that the supported score functions are made of denoising functions that are concentrated on a compact set. To highlight that this can be achieved quite easily, we note that with the naive estimate $s(x,t) = -x/\sigma_t^2$, (15) is satisfied with $B_\ell^2 = \mathbb{E}[\|X_0\|^2]$.

In the following proposition we demonstrate score stability bounds in the case that the step size is decaying with a rate of 1/k.

Proposition 11. Consider the score matching algorithm $A_{sm}: S \mapsto s_{\theta_K}$ for some fixed $K \in \mathbb{N}$ where $(\theta_k)_k$ is as given in (14). Suppose that assumptions 9 and 10 hold and $\eta_k \leq \bar{\eta}/k$ for all k < K, for some $\bar{\eta} \in (0, \lambda^{-1})$. Then, we obtain that A_{sm} is score stable with constant,

$$\varepsilon_{\mathit{stab}}^2 \lesssim \left(\frac{C}{\lambda} \vee R\right)^{1 + \frac{\bar{\eta} v}{\bar{\eta} v + 1}} \frac{\overline{L}^2}{(\bar{\eta} v) \vee 1} \left(\frac{C}{\bar{\eta}}\right)^{\frac{1}{\bar{\eta} v + 1}} \frac{N_B K^{\frac{\bar{\eta} v}{\bar{\eta} v + 1}}}{N},$$

where
$$R^2 = \mathbb{E}[\|\theta_0\|^2]$$
, $\upsilon = (\overline{M}B_\ell C_\tau^{1/2} + \overline{L}^2 - \lambda) \lor 0$ and $C_\tau = \int \sigma_t^{-4} \tau(dt)$.

Since the score matching algorithm is random, to control the stability constant we construct a coupling of the random score functions $A_{\rm sm}(S)$ and $A_{\rm sm}(S^i)$ through a coupling of the optimisation trajectories associated with training on S versus S^i .

6.2 Utilising noise in the gradient estimator

The primary drawback of Proposition 11 is that the bound grows with the number of iterations. This is particularly problematic since diffusion models often require numerous steps due to the high-variance gradient estimator. In this section, we improve this dependence by explicitly leveraging the noise in the gradient estimator. The idea that stochasticity in optimisation can act as a form of implicit regularisation has motivated the development of numerous learning algorithms and theoretical works in recent years (Srivastava et al., 2014; Bishop, 1995; Mou et al., 2018; Pensia et al., 2018). Here, we investigate how the noise intrinsic to the gradient estimator for $\hat{\ell}_{\rm dsm}$ can play a similar role in promoting generalisation in diffusion models.

To incorporate the effects of the gradient noise, we consider a simplified model in which the noise from the stochastic gradient estimator is approximated with a second-order Gaussian approximation:

$$\theta_{k+1} = (1 - \eta \lambda)\theta_k - \eta \mathbb{E}\left[\operatorname{Clip}_C(G_k) \middle| \theta_k, B_k, S\right] + \eta \operatorname{Cov}\left(\operatorname{Clip}_C(G_k) \middle| \theta_k, B_k, S\right)^{1/2} \xi_k, \quad (16)$$

where $\xi_k \in \mathbb{R}^d$ is a standard Gaussian and we use $G_k := G_k(\theta_k, B_k)$. This approximation can be justified by observing that the inner summation in (13) is over conditionally i.i.d. variables, once conditioned on θ , B and S. Therefore, the gradient estimator G becomes approximately Gaussian as P grows large. For this analysis, we assume the following lower bound on the gradient noise.

Assumption 12. There exists a positive semi-definite matrix $\overline{\Sigma} \in \mathbb{R}^{n \times n}$ such that for any $x \in \operatorname{supp}(\nu)$ and $\theta \in \mathbb{R}^n$,

$$\operatorname{Cov}_{t \sim \tau, X_t \mid X_0} \left(\operatorname{Clip}_C(\nabla_{\theta} || s_{\theta}(X_t, t) - x ||^2) | X_0 = x \right) \succeq \overline{\Sigma}.$$

Furthermore, the eigenvalues of $\overline{\Sigma}$, $(\lambda_i)_{i=1}^n$, possess the spectral gap $\lambda_{gap} := \min_{\lambda_i \neq 0} \lambda_i > 0$.

We use the matrix $\overline{\Sigma}$ to dictate the geometry on which we perform our analysis. In particular, we consider the weighted norm $\|v\|_{\overline{\Sigma}^+} := v^T \overline{\Sigma}^+ v$ where $\overline{\Sigma}^+$ is the pseudoinverse matrix.

Assumption 13. For almost all $x \in \mathbb{R}^d$, $t \in (0,T]$, $s_{\theta}(x,t)$ is Lipschitz and smooth (gradient-Lipschitz) in $\theta \in \mathbb{R}^n$ with respect to the seminorm $\|\cdot\|_{\overline{\Sigma}^+}$ and with constants L(x,t) and M(x,t), respectively. Furthermore, there exists constants $\overline{L}, \overline{M} \geq 0$ such that for any $x \in \text{supp}(\nu_{\text{data}})$,

$$\int \mathbb{E}[L(X_t, t)^4 | X_0 = x] \, \tau(dt) \le \overline{L}^4, \quad \int \mathbb{E}[M(X_t, t)^4 | X_0 = x] \, \tau(dt) \le \overline{M}^4.$$

By requiring that the Lipschitz and smoothness properties hold with respect to $\|\cdot\|_{\overline{\Sigma}^+}$, we effectively require that the gradient estimator adds noise in all directions aside from those that do not change the function (e.g. along symmetries in the parameter space). With this, we arrive at our time-convergent score stability bound for SGD.

Proposition 14. Consider the score matching algorithm $A_{sm}: S \mapsto s_{\theta_K}$ for some fixed $K \in \mathbb{N}$ where $(\theta_k)_k$ is as given in (16). Suppose that assumptions 10, 12 and 13 hold, then there exists some $\bar{\eta} > 0$ such that, if $\sup_p \eta_p \leq \bar{\eta}$, we obtain that A_{sm} is score stable with constant

$$\varepsilon_{\textit{stab}}^2 \lesssim \frac{\overline{L}^2 C^2 (P+n)}{\lambda_{\textit{gap}} N} \min \bigg\{ \frac{\eta_{\min} \lambda_{\textit{gap}} \lambda^2}{P N_B C} \sum_{k=0}^{K-1} \eta_k, \exp \bigg(\tilde{c} \frac{P N_B C}{\eta_{\min} \lambda_{\textit{gap}} \lambda^2} \bigg) \bigg\},$$

where
$$\tilde{c} \lesssim (\overline{M}_4 B_\ell C_\tau^{1/2} + \overline{L}_4^2) (P N_B \lambda_{gap})^{-1/2} \vee 1$$
, $\eta_{\min} = \min_k \eta_k$.

In this bound, we recover the $\frac{1}{\sqrt{N}}$ score stability bounds from Proposition 11 while also introducing the property that the bound does not grow endlessly with the number of iterations. This property is obtained using the noise in the gradient estimator and is not possible without additional noise. Through this analysis, we identify the generalisation benefit of a property unique to diffusion models and how they interact with SGD.

7 CONCLUSION AND FUTURE WORK

In this paper, we propose a general algorithm-dependent framework for analysing the generalisation capabilities of diffusion models. We introduce *score stability*, which quantifies an algorithm's sensitivity to the dataset, and use it to derive expected generalisation gap bounds. Applying this framework to several common algorithms, we derive closed-form bounds and identify several previously overlooked sources of implicit regularisation in diffusion models. First, our analysis of empirical risk minimisation finds that the denoising score matching objective provides inherent stability guarantees without further regularisation (denoising regularisation). We then analyse how score stability interacts with discrete-time samplers, identifying that coarse discretisation can improve generalisation guarantees (sampler regularisation). Finally, we consider stochastic optimisation schemes for score matching, obtaining stability guarantees (optimisation regularisation).

This work opens several avenues for future research. Key directions include developing high-probability bounds, developing bounds on privacy and memorisation, tightening our analysis by incorporating data or model properties, like smoothness, and extending the framework to compare different sampling algorithms, such as the probability flow ODE.

REFERENCES

- Eddie Aamari. *Convergence Rates for Geometric Inference*. PhD thesis, Université Paris-Saclay, September 2017.
- Eddie Aamari, Jisu Kim, Frédéric Chazal, Bertrand Michel, Alessandro Rinaldo, and Larry Wasserman. Estimating the reach of a manifold. *Electronic Journal of Statistics*, 13(1):1359–1399, 2019.
 - Amit Attia and Tomer Koren. Uniform stability for first-order empirical risk minimization. In Po-Ling Loh and Maxim Raginsky (eds.), *Proceedings of Thirty Fifth Conference on Learning Theory*, volume 178 of *Proceedings of Machine Learning Research*, pp. 3313–3332. PMLR, 2022.
 - Iskander Azangulov, George Deligiannidis, and Judith Rousseau. Convergence of diffusion models under the manifold hypothesis in high-dimensions. *arXiv preprint arXiv:2409.18804*, 2024.
 - Dominique Bakry, Ivan Gentil, and Michel Ledoux. *Analysis and Geometry of Markov Diffusion Operators*. Springer International Publishing, 2014.
 - Peter L Bartlett, Andrea Montanari, and Alexander Rakhlin. Deep learning: a statistical viewpoint. *Acta Numerica*, 30:87–201, 2021.
 - Joe Benton, Valentin De Bortoli, Arnaud Doucet, and George Deligiannidis. Nearly *d*-linear convergence bounds for diffusion models via stochastic localization. In *International Conference on Learning Representations*, 2024.
 - Chris M Bishop. Training with noise is equivalent to tikhonov regularization. *Neural Computation*, 7 (1):108–116, 1995.
 - Olivier Bousquet and Andre Elisseeff. Stability and generalization. *The Journal of Machine Learning Research*, 2:499–526, 2002.
 - Zachary Charles and Dimitris Papailiopoulos. Stability and generalization of learning algorithms that converge to global optima. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 745–754. PMLR, 2018.
 - Sitan Chen, Sinho Chewi, Jerry Li, Yuanzhi Li, Adil Salim, and Anru R Zhang. Sampling is as easy as learning the score: theory for diffusion models with minimal data assumptions. In *International Conference on Learning Representations*, 2023.
 - Sitan Chen, Vasilis Kontonis, and Kulin Shah. Learning general Gaussian mixtures with efficient score matching. *arXiv* preprint arXiv:2404.18893, 2024.
 - Valentin De Bortoli. Convergence of denoising diffusion models under the manifold hypothesis. *Transactions on Machine Learning Research*, 2022.
 - Valentin De Bortoli, James Thornton, Jeremy Heng, and Arnaud Doucet. Diffusion Schrödinger bridge with applications to score-based generative modeling. In *Advances in Neural Information Processing Systems*, 2021.
 - L Devroye and T Wagner. Distribution-free performance bounds for potential function rules. *IEEE Transactions on Information Theory*, 25(5):601–604, September 1979.
 - Andreas Eberle. Reflection couplings and contraction rates for diffusions. *Probability Theory and Related Fields*, 166(3):851–886, December 2016.
 - Andreas Eberle and Mateusz B Majka. Quantitative contraction rates for Markov chains on general state spaces. *Electronic Journal of Probability*, 24:1–36, January 2019.
 - Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first international conference on machine learning*, 2024.

- Moritz Hardt, Ben Recht, and Yoram Singer. Train faster, generalize better: Stability of stochastic gradient descent. In *International Conference on Machine Learning*, pp. 1225–1234, 2016.
- U G Haussmann and E Pardoux. Time reversal of diffusions. *The Annals of Probability*, 14(4): 1188–1205, 1986.
 - Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, 2020.
 - Zhihan Huang, Yuting Wei, and Yuxin Chen. Denoising diffusion probabilistic models are optimally adaptive to unknown low dimensionality. *arXiv* [cs.LG], October 2024.
 - Aapo Hyvärinen. Estimation of Non-Normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(24):695–709, 2005.
 - Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. In S. Koyejo and S. Mohamed and A. Agarwal and D. Belgrave and K. Cho and A. Oh (ed.), *Advances in Neural Information Processing Systems*, 2022.
 - Tero Karras, Miika Aittala, Jaakko Lehtinen, Janne Hellsten, Timo Aila, and Samuli Laine. Analyzing and improving the training dynamics of diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 24174–24184, June 2024.
 - M Kearns and D Ron. Algorithmic stability and sanity-check bounds for leave-one-out cross-validation. *Neural Computation*, 11(6):1427–1453, 1999.
 - Holden Lee, Jianfeng Lu, and Yixin Tan. Convergence for score-based generative modeling with polynomial complexity. *Advances in Neural Information Processing Systems*, 35:22870–22882, 2022.
 - Puheng Li, Zhong Li, Huishuai Zhang, and Jiang Bian. On the generalization properties of diffusion models. In *Advances in Neural Information Processing Systems*, 2023.
 - Gabriel Loaiza-Ganem, Brendan Leigh Ross, Rasa Hosseinzadeh, Anthony L Caterini, and Jesse C Cresswell. Deep generative models through the lens of the manifold hypothesis: A survey and new connections. *arXiv* [cs.LG], April 2024.
 - Mateusz B Majka, Aleksandar Mijatovic, and Lukasz Szpruch. Nonasymptotic bounds for sampling algorithms without log-concavity. *Ann. Appl. Probab.*, 30(4):1534–1581, August 2020.
 - Wenlong Mou, Liwei Wang, Xiyu Zhai, and Kai Zheng. Generalization bounds of sgld for non-convex learning: Two theoretical viewpoints. In Bubeck, Sébastien and Perchet, Vianney and Rigollet, Philippe (ed.), *Proceedings of the 31st Conference On Learning Theory*, volume 75 of *Proceedings of Machine Learning Research*, pp. 605–638. PMLR, 2018.
 - Yurii Nesterov. *Lectures on convex optimization*. Springer optimization and its applications. Springer International Publishing, Cham, Switzerland, 2 edition, December 2018.
 - Kazusato Oko, Shunta Akiyama, and Taiji Suzuki. Diffusion models are minimax optimal distribution estimators. In *International Conference on Machine Learning*, 2023.
 - Ankit Pensia, Varun Jog, and Po-Ling Loh. Generalization error bounds for noisy, iterative algorithms. In 2018 IEEE International Symposium on Information Theory (ISIT), pp. 546–550, June 2018.
 - Jakiw Pidstrigach. Score-based generative models detect manifolds. In *Advances in Neural Information Processing Systems*, 2022.
 - Peter Potaptchik, Iskander Azangulov, and George Deligiannidis. Linear convergence of diffusion models under the manifold hypothesis. *arXiv preprint arXiv:2410.09046*, 2024.
 - Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.

- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Bjorn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, June 2022.
- Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. Photorealistic Text-to-Image diffusion models with deep language understanding. In *Advances in Neural Information Processing Systems*, 2022.
- Kulin Shah, Sitan Chen, and Adam Klivans. Learning mixtures of Gaussians using the DDPM objective. In *Advances in Neural Information Processing Systems*, 2023.
- Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, 2015.
- Yang Song and Diederik P Kingma. How to train your energy-based models. *arXiv preprint* arXiv:2101. 03288, 2021.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-Based generative modeling through stochastic differential equations. In *9International Conference on Learning Representations*, 2021.
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(56):1929–1958, 2014.
- Rong Tang and Yun Yang. Adaptivity of diffusion models to manifold structures. In *International Conference on Artificial Intelligence and Statistics*, 2024.
- Simon Vary, David Martínez-Rubio, and Patrick Rebeschini. Black-box uniform stability for non-euclidean empirical risk minimization. *arXiv* [cs.LG], December 2024.
- Martin J Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2019.
- Feng-Yu Wang. Logarithmic sobolev inequalities on noncompact riemannian manifolds. *Probability Theory and Related Fields*, 109(3):417–424, November 1997.
- Joseph L Watson, David Juergens, Nathaniel R Bennett, Brian L Trippe, Jason Yim, Helen E Eisenach, Woody Ahern, Andrew J Borst, Robert J Ragotte, Lukas F Milles, et al. De novo design of protein structure and function with RFdiffusion. *Nature*, 620(7976):1089–1100, 2023.
- Christopher Williams, Andrew Campbell, Arnaud Doucet, and Saifuddin Syed. Score-optimal diffusion schedules. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- Hongkang Yang. A mathematical framework for learning probability distributions. *arXiv* preprint *arXiv*:2212.11481, 2022.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Google Brain, Benjamin Recht, and Oriol Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021.

A FURTHER BACKGROUND

We begin with some further details on notation and lemmas used throughout this work and provide proofs for the lemmas in Section 2.

A.1 RANDOM SCORE MATCHING ALGORITHMS

We begin with some additional details on how random score matching algorithms are defined in this work. Recalling the probability space $(\Omega, \mathcal{F}, \mathbb{P})$, we define the set of random score functions,

$$\mathcal{S} := \left\{ s : \mathbb{R}^d \times [0, T] \times \Omega : s(\cdot, \cdot, \omega) \in L^0(\mathbb{R}^d \times [0, T]; \mathbb{R}^d) \right\}.$$

For any random score matching algorithm $A_{\mathrm{sm}}: (\cup_{N=1}^{\infty}(\mathbb{R}^d)^{\otimes N}) \times \Omega \to L^0(\mathbb{R}^d \times [0,T];\mathbb{R}^d)$, we use $A_{\mathrm{sm}}(S)$ as shorthand for the random score function $(\omega,x,t) \mapsto A_{\mathrm{sm}}(S,\omega)(x,t)$ belonging to \mathcal{S} .

Given two random score functions s, s', let $\Gamma(s, s')$ denote the set of all couplings of these functions which we define as,

$$\Gamma(s,s') := \Big\{ (\tilde{s},\tilde{s}') \in \mathcal{S} \times \mathcal{S} : \tilde{s} \simeq s, \tilde{s}' \simeq s' \Big\},$$

where $\tilde{s} \simeq s$ denotes the fact that for any bounded measurable test function $\phi: L^0(\mathbb{R}^d \times [0,T];\mathbb{R}^d) \to \mathbb{R}$, it holds that,

$$\int \phi(s(\cdot,\cdot,\omega))d\mathbb{P} = \int \phi(\tilde{s}(\cdot,\cdot,\omega))d\mathbb{P}.$$

A.2 PRELIMINARY LEMMAS

For the score matching loss bound, we begin with the fact that the score matching loss is equivalent to the denoising score matching loss up to an added constant Song et al. (2021); Hyvärinen (2005).

Lemma 15. For any t > 0, $y \in \mathbb{R}^d$, we have

$$\nabla \log p_t(y) = \frac{\mu_t \mathbb{E}[X_0 | X_t = y] - y}{\sigma_t^2}, \qquad \nabla \log \hat{p}_t(y) = \frac{\mu_t \mathbb{E}[\hat{X}_0 | \hat{X}_t = y, S] - y}{\sigma_t^2}. \tag{17}$$

Proof. We begin by showing that the conditional score is an unbiased estimate of $\nabla \log p_t$. For any $x \in \mathbb{R}^d$, t > 0, we have

$$\mathbb{E}[\nabla \log p_{t|0}(X_t|X_0)|X_t = x] = \int \nabla_x \log p_{t|0}(x|y) \, p_{0|t}(y|x) dy$$

$$= \int \nabla \log p_{t|0}(x|y) \, \frac{p_{t|0}(x|y)p_0(y)}{p_t(x)} dy$$

$$= \int \nabla p_{t|0}(x|y) \, \frac{p_0(y)}{p_t(x)} dy.$$

Therefore, using the exchangeability of gradients and integrals (note that $p_{t|0}$ is C^{∞}), we arrive at

$$\mathbb{E}[\nabla \log p_{t|0}(X_t|X_0)|X_t = x] = \frac{\nabla p_t(x)}{p_t(x)}$$

$$= \nabla \log p_t(x). \tag{18}$$

Alternatively, using (7), we obtain that the left-hand side takes the form,

$$\mathbb{E}[\nabla \log p_{t|0}(X_t|X_0)|X_t = x] = \frac{\mu_t \mathbb{E}[X_0|X_t = x] - x}{\sigma_t^2},$$

completing the proof of the first equality in (17). For the second equality, concerning that empirical score function, the proof follows similarly once the empirical measure $\frac{1}{N} \sum_{i=1}^{N} \delta_{x_i}$ is considered in place of ν_{data} .

Lemma 16. For any integrable score function s, it holds that

$$\ell_{\rm dsm}(s;\tau) = \ell_{\rm sm}(s;\tau) + C_{\rm sm},$$

where, given $s^*(x,t) := \nabla \log p_t(x)$, we define

$$C_{\rm sm} := \int \frac{\mu_t^2}{\sigma_t^4} \mathbb{E}[\operatorname{Tr} \operatorname{Cov}(X_0|X_t)] \tau(dt) = \ell_{\rm dsm}(s^*; \tau). \tag{20}$$

Proof. Let s be any score function. Using the equality in (19), we obtain the following bias-variance decomposition of $\ell_{\rm dsm}(s;\tau)$:

 $\ell_{\rm dsm}(s;\tau)$

$$\begin{split} &= \int \mathbb{E} \Big[\|s(X_t, t) - \nabla \log p_{t|0}(X_t | X_0)\|^2 \Big] \tau(dt) \\ &= \int \mathbb{E} \Big[\|s(X_t, t) - \nabla \log p_t(X_t)\|^2 \Big] \tau(dt) + \int \mathbb{E} \Big[\|\nabla \log p_{t|0}(X_t | X_0) - \nabla \log p_t(X_t)\|^2 \Big] \tau(dt) \\ &= \ell_{\text{sm}}(s; \tau) + \int \mathbb{E} \Big[\text{Tr Cov} \left(\nabla \log p_{t|0}(X_t | X_0) \Big| X_t \right) \Big] \tau(dt). \end{split}$$

Once we note that,

$$\operatorname{Tr} \operatorname{Cov} \left(\nabla \log p_{t|0}(X_t|X_0) \middle| X_t \right) = \operatorname{Tr} \operatorname{Cov} \left(\frac{\mu_t X_0 - x}{\sigma_t^2} \middle| X_t \right)$$
$$= \frac{\mu_t^2}{\sigma_t^4} \operatorname{Tr} \operatorname{Cov}(X_0|X_t),$$

we obtain the bound $\ell_{\rm dsm}(s;\tau)=\ell_{\rm sm}(s;\tau)+C_{\rm sm}$ from the statement. To derive the equality $C_{\rm sm}=\ell_{\rm dsm}(s^\star;\tau)$, we use that $\ell_{\rm sm}(s^\star;\tau)=0$ and so we obtain $\ell_{\rm dsm}(s^\star;\tau)=0+C_{\rm sm}$.

Similarly, there is an equivalence between the empirical forms of the denoising score matching loss and the score matching loss,

$$\hat{\ell}_{\rm dsm}(s; S, \tau) = \hat{\ell}_{\rm sm}(s; S, \tau) + \hat{C}_{\rm sm},\tag{21}$$

where

$$\hat{C}_{\mathrm{sm}} := \int \frac{\mu_t^2}{\sigma_t^4} \mathbb{E}[\operatorname{Tr} \operatorname{Cov}(\hat{X}_0 | \hat{X}_t, S) | S] \tau(dt) = \hat{\ell}_{\mathrm{dsm}}(\hat{s}^{\star}; S, \tau), \tag{22}$$

and $\hat{s}^{\star}(x,t) = \nabla \hat{p}_t(x)$. This follows immediately from the above proof once the empirical measure $\frac{1}{N} \sum_{i=1}^{N} \delta_{x_i}$ is considered in place of ν_{data} . This effectively completes the proof of Lemma 1 in Section 2.

Lemma 1. The objective $\hat{\ell}_{dsm}(s; S, \tau)$ is identical, up to a constant, to the objective

$$\hat{\ell}_{sm}(s; S, \tau) := \int \mathbb{E}[\|s(\hat{X}_t, t) - \nabla \log \hat{p}_t(\hat{X}_t)\|^2 |S| \tau(dt), \tag{23}$$

where \hat{p}_t is the marginal density of \hat{X}_t . Therefore, any minimiser of $\hat{\ell}_{dsm}(\cdot; S, \tau)$ on $L^0(\mathbb{R}^d \times [0, T]; \mathbb{R}^d)$ is identical to $\nabla \log \hat{p}_t$ a.e. for any $t \in \text{supp}(\tau)$.

Proof. The proof follows nearly immediately from (21). Since $p_{t|0}$ is C^{∞} , $\nabla \log p_{t|0}$ is measurable and thus its empirical average $\nabla \log \hat{p}_t$ must be also. Therefore, the score function $s^{\star}(x,t) = \nabla \log \hat{p}_t(x)$ satisfies $\hat{s}^{\star} \in L^0(\mathbb{R}^d \times [0,T];\mathbb{R}^d)$ as well as,

$$\hat{\ell}_{\rm sm}(\hat{s}^{\star}; S, \tau) = 0.$$

Now let $s \in L^0(\mathbb{R}^d \times [0,T];\mathbb{R}^d)$ be any minimiser of $\hat{\ell}_{\mathrm{dsm}}(\cdot;S,\tau)$. Through the equivalence of $\hat{\ell}_{\mathrm{dsm}}$ and $\hat{\ell}_{\mathrm{sm}}$ up to a constant, it follows that s must also be a minimiser of $\hat{\ell}_{\mathrm{sm}}(\cdot;S,\tau)$ and, due to the existence of \hat{s}^\star , must satisfy $\hat{\ell}_{\mathrm{sm}}(s;S,\tau)=0$ also. Letting $t\in\mathrm{supp}(t)$, we note that since t>0, we must have that $p_{t|0}$ has full support and thus, $s(\cdot,t)=s^\star(\cdot,t)$ almost everywhere.

A.3 MANIFOLDS

We also introduce some basic properties of smooth manifolds, primarily referencing Aamari et al. (2019). We define the manifold reach and include a known property of this quantity.

Definition 17. The reach of a set $A \subset \mathbb{R}^d$, is defined by $\tau_A = \inf_{p \in A} d(p, Med(A))$, where we define the set,

$$Med(A) = \Big\{z \in \mathbb{R}^d: \exists p,q \in A \text{ s.t. } p \neq q, \|p-z\| = \|q-z\|\Big\}.$$

Lemma 18. Suppose that the measure μ is supported on a manifold M with reach $\tau_M > 0$ and dimension d^* . Then, for any $r \leq \tau_M$, we have

$$\mu(B_r(x)) \ge \left| \inf_{B_r(x)} p_\mu \right| r^{d^*},$$

where p_{μ} denotes the density of μ with respect to the volume measure on M.

For the proof of this lemma, we refer to the proof of Proposition 4.3 in Aamari et al. (2019) or Lemma III.23 in Aamari (2017).

B PROOFS FOR THE GENERALISATION GAP BOUNDS

We now provide provide the proof of theorem 3 that bound the generalisation gap under score stability guarantees. For the sake of brevity, throughout this section we suppress the notation for the time weighting, for example, using the shorthand $\hat{\ell}_{sm}(s;S)$ in place of $\hat{\ell}_{sm}(s;S,\tau)$.

Theorem 3. Suppose that the score matching algorithm A_{sm} is score stable with constant ε_{stab} . Then, with $\hat{s} = A_{sm}(S)$, it holds that

$$\left| \mathbb{E} \left[\ell_{\rm dsm}(\hat{s}; \tau) \right]^{1/2} - \mathbb{E} \left[\hat{\ell}_{\rm dsm}(\hat{s}; S, \tau) \right]^{1/2} \right| \le \varepsilon_{stab}. \tag{24}$$

Furthermore, it holds that

$$\mathbb{E}\left[\ell_{\mathrm{sm}}(\hat{s};\tau)\right] - \mathbb{E}\left[\hat{\ell}_{\mathrm{sm}}(\hat{s};S,\tau)\right] \le 2\,\varepsilon_{\mathit{stab}}\,\mathbb{E}\left[\hat{\ell}_{\mathrm{dsm}}(\hat{s};S,\tau)\right]^{1/2} + \varepsilon_{\mathit{stab}}^{2}.\tag{25}$$

Proof. Setting $\hat{s} = A_{\rm sm}(S)$ and $\hat{s}^i = A_{\rm sm}(S^i)$, we use the property that (\hat{s}, \tilde{x}) and (\hat{s}^i, x_i) are distributed identically to obtain that,

$$\begin{split} \mathbb{E}[\ell_{\mathrm{dsm}}(\hat{s};\tau)] &= \mathbb{E}[\hat{\ell}_{\mathrm{dsm}}(\hat{s};\{\tilde{x}\})] \\ &= \mathbb{E}\Big[\frac{1}{N}\sum_{i=1}^{N}\hat{\ell}_{\mathrm{dsm}}(\hat{s}^{i};\{x_{i}\})\Big] \\ &= \mathbb{E}\Big[\frac{1}{N}\sum_{i=1}^{N}\int\mathbb{E}_{X_{t}}[\|\hat{s}^{i}(X_{t},t,\omega) - \nabla\log p_{t|0}(X_{t}|x_{i})\|^{2}|X_{0} = x_{i},S]\,\tau(dt)\,\Big]. \end{split}$$

Therefore, it follows from the triangle inequality in L^2 -norm that

$$\left| \mathbb{E}[\ell_{\mathrm{dsm}}(\hat{s};\tau)]^{1/2} - \mathbb{E}[\hat{\ell}_{\mathrm{dsm}}(\hat{s};S)]^{1/2} \right| \leq \mathbb{E}\left[\frac{1}{N} \sum_{i=1}^{N} \int \mathbb{E}[\|\hat{s}(X_{t},t) - \hat{s}^{i}(X_{t},t)\|^{2} |X_{0} = x_{i},S] \, \tau(dt) \right]^{1/2}$$

Note that if the algorithm $A_{\rm sm}$ is stochastic, the right-hand side would hold regardless of how $\hat{s}|S, \tilde{x}$ and $\hat{s}^i|S, \tilde{x}$ were coupled. Therefore the most efficient coupling can be chosen, leading to the bound,

$$\left| \mathbb{E}[\ell_{\mathrm{dsm}}(\hat{s};\tau)]^{1/2} - \mathbb{E}[\hat{\ell}_{\mathrm{dsm}}(\hat{s};S)]^{1/2} \right|$$

$$\leq \mathbb{E}\left[\inf_{(\hat{s},\hat{s}^i)\in\Gamma_i} \frac{1}{N} \sum_{i=1}^N \int \mathbb{E}[\|\hat{s}(X_t,t) - \hat{s}^i(X_t,t)\|^2 | X_0 = x_i, S] \tau(dt) \right]^{1/2}$$

$$\leq \varepsilon_{\mathrm{stab}},$$

$$(27)$$

completing the proof of the bound in (24).

To obtain the bound in (25), we use Lemma 16 to derive

$$\mathbb{E}[\ell_{\rm sm}(\hat{s};\tau)] = \mathbb{E}[\hat{\ell}_{\rm sm}(\hat{s};S)] + \mathbb{E}[\ell_{\rm dsm}(\hat{s};\tau) - \hat{\ell}_{\rm dsm}(\hat{s};S)] + \mathbb{E}[\hat{\ell}_{\rm dsm}(\nabla \log \hat{p}_t;S)] - \ell_{\rm dsm}(\nabla \log p_t;\tau). \tag{28}$$

Since $\hat{\ell}_{dsm}(\cdot; S)$ is a unbiased estimator of $\ell_{sm}(\cdot; \tau)$, we have that

$$\ell_{\rm dsm}(\nabla \log p_t; \tau) = \mathbb{E}[\hat{\ell}_{\rm dsm}(\nabla \log p_t; S)] \ge \mathbb{E}[\hat{\ell}_{\rm dsm}(\nabla \log \hat{p}_t; S)], \tag{29}$$

where the inequality follows from the fact that $\nabla \log \hat{p}_t$ minimises $\hat{\ell}_{\rm dsm}$. Furthermore, using (27), we deduce the bound,

$$\begin{aligned} |\mathbb{E}[\ell_{\mathrm{dsm}}(\hat{s};\tau) - \hat{\ell}_{\mathrm{dsm}}(\hat{s};S)]| \\ &= \left(\mathbb{E}[\ell_{\mathrm{dsm}}(\hat{s};\tau)]^{1/2} + \mathbb{E}[\hat{\ell}_{\mathrm{dsm}}(\hat{s};S)]^{1/2} \right) \left| \mathbb{E}[\ell_{\mathrm{dsm}}(\hat{s};S)]^{1/2} - \mathbb{E}[\hat{\ell}_{\mathrm{dsm}}(\hat{s};S)]^{1/2} \right| \\ &\leq \left(2\mathbb{E}[\hat{\ell}_{\mathrm{dsm}}(\hat{s};S)]^{1/2} + \varepsilon_{\mathrm{stab}} \right) \varepsilon_{\mathrm{stab}} \\ &= 2\varepsilon_{\mathrm{stab}} \mathbb{E}[\hat{\ell}_{\mathrm{dsm}}(\hat{s};S)]^{1/2} + \varepsilon_{\mathrm{stab}}^{2}. \end{aligned}$$
(30)

Thus, substituting (29) and (30) in to (28) recovers the bound in (25) in the statement. \Box

We obtain upper bounds relying on the fact that the constant separating the score matching loss from the denoising score matching loss is larger on average in the empirical case. One could obtain lower bounds through our techniques but this would require an analysis of the rate of convergence of this constant which is beyond the scope of this paper.

C Proofs for stability of empirical denoising score matching

In this section, we provide the proof for Theorem 3, where the algorithm that minimises $\hat{\ell}_{\rm dsm}(\cdot; S, \tau)$ over some class of score functions \mathcal{H} is shown to be score stable.

C.1 On-average stability of the ERM algorithm

We begin with an important lemma that shows that under minimal assumptions, $\hat{s} = A_{\rm erm}(S)$ and $\hat{s}^i = A_{\rm erm}(S)$ are close in L^2 space, averaged over the full dataset. The first half of this proof utilises the fact that $\hat{\ell}_{\rm dsm}$ is 1-strongly convex in a weighted L^2 space, exploiting a well-known relationship between strong-convexity and algorithmic stability (e.g. see (Bousquet & Elisseeff, 2002; Charles & Papailiopoulos, 2018; Vary et al., 2024; Attia & Koren, 2022)).

Lemma 19. Suppose that A_{erm} is score stable with constant $\varepsilon_{\text{stab}}$, then for any $i \in [N]$, we obtain,

$$\mathbb{E}\left[\int\int \|\hat{s}^{i}(y,t) - \hat{s}(y,t)\|^{2} \,\hat{p}_{t}(dy) \,\tau(dt)\right] \leq 8\mathbb{E}[\hat{\ell}_{\mathrm{sm}}(\hat{s})] + \frac{8}{N} \varepsilon_{\mathit{stab}}(C_{\mathrm{sm}}^{1/2} + \varepsilon_{\mathit{stab}}) \tag{31}$$

where $\hat{s} = A_{\text{erm}}(S), \hat{s}^i = A_{\text{erm}}(S)$.

Proof. Choose $i \in [N]$ and let $\hat{s} = A_{\mathrm{erm}}(S), \hat{s}^i = A_{\mathrm{erm}}(S^i)$ so that $\hat{s} \in \operatorname{argmin}_{\mathcal{H}} \hat{\ell}_{\mathrm{dsm}}(\cdot; S, \tau), \hat{s}^i \in \operatorname{argmin}_{\mathcal{H}} \hat{\ell}_{\mathrm{dsm}}(\cdot; S^i, \tau)$. The proof begins with the following simple expression, that holds for all $j \in [N]$:

$$2\int \left\langle \hat{s}^{i}(y,t) - \hat{s}(y,t), \hat{s}^{i} - \nabla \log p_{t|0}(y|x_{j}) \right\rangle p_{t|0}(dy|x_{j})$$

$$= \int \|\hat{s}^{i}(y,t) - \nabla \log p_{t|0}(y|x_{j})\|^{2} p_{t|0}(dy|x_{j}) - \int \|\hat{s}(y,t) - \nabla \log p_{t|0}(y|x_{j})\|^{2} p_{t|0}(dy|x_{j})$$

$$+ \int \|\hat{s}^{i}(y,t) - \hat{s}(y,t)\|^{2} p_{t|0}(dy|x_{j}).$$

By averaging over $j \in [N]$ and integrating with respect to $\tau(dt)$, we arrive at the upper bound,

$$\frac{2}{N} \sum_{j \in [N]} \int \int \left\langle \hat{s}^{i}(y, t) - \hat{s}(y, t), \hat{s}^{i} - \nabla \log p_{t|0}(y|x_{j}) \right\rangle p_{t|0}(dy|x_{j}) \tau(dt)
= \hat{\ell}_{dsm}(\hat{s}^{i}; S, \tau) - \hat{\ell}_{dsm}(\hat{s}; S, \tau) + \int \int \|\hat{s}^{i}(y, t) - \hat{s}(y, t)\|^{2} \hat{p}_{t}(dy) \tau(dt)
\geq \int \int \|\hat{s}^{i}(y, t) - \hat{s}(y, t)\|^{2} \hat{p}_{t}(dy) \tau(dt),$$
(32)

where the inequality follows from the fact that $\hat{\ell}_{\mathrm{dsm}}(\hat{s}; S, \tau) \leq \hat{\ell}_{\mathrm{dsm}}(s; S, \tau)$ for any score function $s \in \mathcal{H}$. Additionally, the left-hand side is upper bounded using the Cauchy-Schwarz inequality to obtain,

$$\frac{2}{N} \sum_{x \in S} \int \int \langle \hat{s}^{i}(y,t) - \hat{s}(y,t), \hat{s}^{i} - \nabla \log p_{t|0}(y|x) \rangle p_{t|0}(dy|x) \tau(dt)
= \frac{2}{N} \sum_{x \in S^{i}} \int \int \langle \hat{s}^{i}(y,t) - \hat{s}(y,t), \hat{s}^{i}(y,t) - \nabla \log p_{t|0}(y|x) \rangle p_{t|0}(dy|x) \tau(dt)
+ \frac{2}{N} \int \int \langle \hat{s}^{i}(y,t) - \hat{s}(y,t), \hat{s}^{i}(y,t) - \nabla \log p_{t|0}(y|x_{i}) \rangle p_{t|0}(dy|x_{i}) \tau(dt)
- \frac{2}{N} \int \int \langle \hat{s}^{i}(y,t) - \hat{s}(y,t), \hat{s}^{i}(y,t) - \nabla \log p_{t|0}(y|\tilde{x}) \rangle p_{t|0}(dy|\tilde{x}) \tau(dt)
\leq 2\hat{\ell}_{sm}(\hat{s}^{i}; S^{i}, \tau)^{1/2} \left(\int \int \|\hat{s}^{i}(y,t) - \hat{s}(y,t)\|^{2} \hat{p}_{t}^{i}(dy) \tau(dt) \right)^{1/2}
+ \frac{2}{N} \hat{\ell}_{dsm}(\hat{s}^{i}; \{x_{i}\}, \tau)^{1/2} \left(\int \int \|\hat{s}^{i}(y,t) - \hat{s}(y,t)\|^{2} p_{t|0}(dy|x_{i}) \tau(dt) \right)^{1/2}
+ \frac{2}{N} \hat{\ell}_{dsm}(\hat{s}^{i}; \{\tilde{x}\}, \tau)^{1/2} \left(\int \int \|\hat{s}^{i}(y,t) - \hat{s}(y,t)\|^{2} p_{t|0}(dy|x_{i}) \tau(dt) \right)^{1/2}, \quad (33)$$

where $\hat{p}_t^i(dy) = \frac{1}{N} \sum_{x \in S^i} p_{t|0}(dy|x)$. Combining the expressions in (32) and (33) and taking the expectation, we derive the bound,

$$\begin{split} \mathbb{E} \bigg[\int \int \| \hat{s}^i(y,t) - \hat{s}(y,t) \|^2 \, \hat{p}_t(dy) \, \tau(dt) \bigg] \\ & \leq 2 \mathbb{E} [\hat{\ell}_{\mathrm{sm}}(\hat{s}^i;S^i,\tau)]^{1/2} \mathbb{E} \bigg[\int \int \| \hat{s}^i(y,t) - \hat{s}(y,t) \|^2 \, \hat{p}_t^i(dy) \, \tau(dt) \bigg]^{1/2} \\ & \quad + \frac{2}{N} \mathbb{E} [\hat{\ell}_{\mathrm{dsm}}(\hat{s}^i;\{x_i\},\tau)]^{1/2} \mathbb{E} \bigg[\int \int \| \hat{s}^i(y,t) - \hat{s}(y,t) \|^2 \, p_{t|0}(dy|x_i) \, \tau(dt) \bigg]^{1/2} \\ & \quad + \frac{2}{N} \mathbb{E} [\hat{\ell}_{\mathrm{dsm}}(\hat{s}^i;\{\tilde{x}\},\tau)]^{1/2} \mathbb{E} \bigg[\int \int \| \hat{s}^i(y,t) - \hat{s}(y,t) \|^2 \, p_{t|0}(dy|\tilde{x}) \, \tau(dt) \bigg]^{1/2} \\ & \leq 2 \mathbb{E} [\hat{\ell}_{\mathrm{sm}}(\hat{s};S,\tau)]^{1/2} \mathbb{E} \bigg[\int \int \| \hat{s}^i(y,t) - \hat{s}(y,t) \|^2 \, \hat{p}_t(dy) \, \tau(dt) \bigg]^{1/2} \\ & \quad + \frac{2}{N} \varepsilon_{\mathrm{stab}} \bigg(\mathbb{E} [\hat{\ell}_{\mathrm{dsm}}(\hat{s};S,\tau)]^{1/2} + \mathbb{E} [\ell_{\mathrm{dsm}}(\hat{s},)]^{1/2} \bigg), \end{split}$$

where we recall that $\varepsilon_{\mathrm{stab}}$ is the stability constant for A_{erm} . Here, we have used the fact that (\hat{s}, S) has the same law as (\hat{s}^i, S^i) and also $\mathbb{E}[\hat{\ell}_{\mathrm{dsm}}(\hat{s}^i; \{\tilde{x}\})] = \mathbb{E}[\hat{\ell}_{\mathrm{dsm}}(\hat{s}; S)]$ and $\mathbb{E}[\hat{\ell}_{\mathrm{dsm}}(\hat{s}^i; \{x_i\})] = \mathbb{E}[\hat{\ell}_{\mathrm{dsm}}(\hat{s}; S)]$

 $\mathbb{E}[\ell_{\mathrm{dsm}}(\hat{s})]$. By solving the quadratic equation, we deduce that the above inequality implies that,

$$\begin{split} \mathbb{E}\bigg[\int\int \|\hat{s}^i(y,t) - \hat{s}(y,t)\|^2 \, \hat{p}_t(dy) \, \tau(dt)\bigg] \\ &\leq \left(\mathbb{E}[\hat{\ell}_{\mathrm{sm}}(\hat{s};S,\tau)]^{\frac{1}{2}} + \sqrt{\mathbb{E}[\hat{\ell}_{\mathrm{sm}}(\hat{s};S,\tau)] + \frac{2}{N}\varepsilon_{\mathrm{stab}}(\mathbb{E}[\ell_{\mathrm{dsm}}(\hat{s};\tau)]^{\frac{1}{2}} + \mathbb{E}[\hat{\ell}_{\mathrm{dsm}}(\hat{s};S,\tau)]^{1/2})}\right)^2 \\ &\leq 4\mathbb{E}[\hat{\ell}_{\mathrm{sm}}(\hat{s};S,\tau)] + \frac{4}{N}\varepsilon_{\mathrm{stab}}(\mathbb{E}[\ell_{\mathrm{dsm}}(\hat{s};\tau)]^{\frac{1}{2}} + \mathbb{E}[\hat{\ell}_{\mathrm{dsm}}(\hat{s};S,\tau)]^{\frac{1}{2}}). \end{split}$$

We simplify the above expression further using Theorem 3. Using the stability assumption, it follows from (24) that $\mathbb{E}[\ell_{\rm dsm}(\hat{s})]^{1/2} \leq \mathbb{E}[\hat{\ell}_{\rm dsm}(\hat{s})]^{1/2} + \varepsilon$. Furthermore, from Lemma 16, we have

$$\begin{split} \mathbb{E}[\hat{\ell}_{\mathrm{dsm}}(\hat{s})] &= \mathbb{E}[\hat{\ell}_{\mathrm{sm}}(\hat{s})] + \mathbb{E}[\hat{C}_{\mathrm{sm}}] \\ &\leq \mathbb{E}[\hat{\ell}_{\mathrm{sm}}(\hat{s})] + C_{\mathrm{sm}}, \end{split}$$

where we recall the definitions of $\hat{C}_{\rm sm}$ and $C_{\rm sm}$ from (22) and (20) and recall that $\mathbb{E}[\hat{C}_{\rm sm}] \leq C_{\rm sm}$ from (29). Thus, from Young's inequality, we obtain the bound

$$\begin{split} \mathbb{E}\bigg[\int\int \|\hat{s}^i(y,t) - \hat{s}(y,t)\|^2 \, \hat{p}_t(dy) \, \tau(dt)\bigg] \\ & \leq 4\mathbb{E}[\hat{\ell}_{\mathrm{sm}}(\hat{s})] + \frac{4}{N} \varepsilon_{\mathrm{stab}} (2\mathbb{E}[\hat{\ell}_{\mathrm{sm}}(\hat{s})]^{1/2} + 2C_{\mathrm{sm}}^{1/2} + \varepsilon_{\mathrm{stab}}) \\ & \leq 8\mathbb{E}[\hat{\ell}_{\mathrm{sm}}(\hat{s})] + \frac{4}{N} \varepsilon_{\mathrm{stab}} (\varepsilon_{\mathrm{stab}}/N + 2C_{\mathrm{sm}}^{1/2} + \varepsilon_{\mathrm{stab}}) \\ & \leq 8\mathbb{E}[\hat{\ell}_{\mathrm{sm}}(\hat{s})] + \frac{8}{N} \varepsilon_{\mathrm{stab}} (C_{\mathrm{sm}}^{1/2} + \varepsilon_{\mathrm{stab}}). \end{split}$$

C.2 PROOF OF PROPOSITION 6

To obtain the stability bound in Proposition 6, we convert the result in Lemma 19, which is a bound in $L^2(\hat{p}_t)$, to a bound in $L^2(p_{t|0}(\cdot|\tilde{x}))$ which is required of score stability. For this, we rely on two further lemmas, the first of which is a fundamental property of the Ornstein-Uhlenbeck process, captured by the Harnack inequality of Wang (1997) (see Theorem 5.6.1 Bakry et al. (2014)).

Lemma 20 (Wang's Harnack inequality). For each positive measurable function $\phi : \mathbb{R}^d \to \mathbb{R}$, every t > 0, p > 1 and every $x, y \in \mathbb{R}^d$, it holds that

$$\mathbb{E}[\phi(X_t)|X_0 = x] \le \mathbb{E}[\phi(X_t)^p | X_0 = y]^{1/p} \exp\left(\frac{\mu_t^2 ||x - y||^2}{2(p - 1)\sigma_t^2}\right).$$

This result describes the stability of the diffusion semigroup under changes in initial position and shows that as t grows, the distribution of X_t depends less on X_0 . The second lemma, for which we provide a proof, controls the empirical measure,

$$\hat{\nu}(dx) = \frac{1}{N} \sum_{i=1}^{N} \delta_{x_i}(dx),$$

on balls around training examples.

Lemma 21. Suppose that Assumption 4 is satisfied, then for any $i \in [N], r \in (0, \tau_{reach}]$ and any decreasing function $\phi : (0, \infty) \to \mathbb{R}_+$, we have the bound

$$\mathbb{E}\Big[\phi\Big(\hat{\nu}(B_r(x_i))\Big)\Big] \le \phi(N^{-1})\exp(-c_{\nu}N^2r^{d^*}) + \phi(c_{\nu}r^{d^*}/2),$$

whenever $N \geq 4c_{\nu}^{-1}r^{-d^*}$, where $c_{\nu} = \inf p_{\nu}$.

Proof. We rewrite the object $\hat{\nu}(B_r(x_i))$ as an empirical average of Bernoulli random variables

$$\hat{\nu}(B_r(x_i)) = \frac{1}{N} \sum_{j=1}^N \mathbb{1}_{x_j \in B_r(x_i)} = \frac{1}{N} + \frac{1}{N} \sum_{j \neq i} \mathbb{1}_{x_j \in B_r(x_i)}.$$

When conditioned on x_i , the random variables $(\mathbb{1}_{x_j \in B_r(x_i)})_{j \neq i}$ are independently and identically distributed Bernoulli random variable with probability $\mu = \nu_{\text{data}}(B_r(x_i))$. To utilise concentration of the empirical process, we first rewrite the probability

$$\mathbb{P}\Big(\hat{\nu}(B_r(x_i)) \le \mu/2 \Big| x_i\Big) \le \mathbb{P}\Big(S_{N-1} \le \frac{N\mu}{2} - 1 \Big| x_i\Big), \qquad S_{N-1} = \sum_{i \ne i} \mathbb{1}_{x_i \in B_r(x_i)}.$$

Therefore, by Chernoff's inequality we obtain

$$\mathbb{P}\Big(\hat{\nu}(B_r(x_i)) \le \mu/2 \Big| x_i\Big) \le \exp\Big(-\mu^{-1}(N\mu/2 - 1)^2\Big)$$

$$\le \exp(-N^2\mu/16),$$

where the last bound holds when $N \ge 4\mu^{-1}$. Therefore, using the above bound as well as the trivial bound $\hat{\nu}(B_r(x_i)) \ge N^{-1}$ we apply the law of total expectation to obtain,

$$\mathbb{E}\Big[\phi\Big(\hat{\nu}(B_r(x_i))\Big)\Big|x_i\Big] = \mathbb{E}\Big[\phi\Big(\hat{\nu}(B_r(x_i))\Big)\Big|\hat{\nu}(B_r(x_i)) > \mu/2\Big] + \mathbb{P}\Big(\hat{\nu}(B_r(x_i)) \le \mu/2\Big|x_i\Big)\phi(N^{-1})$$

$$\le \phi(\mu/2) + \exp(-N^2\mu/16)\phi(N^{-1}).$$

To control μ , we use Lemma 18 which asserts that $\mu \geq c_{\nu} r^{d^*}$.

This now brings us to the proof of the proposition, which we first restate.

Proposition 6. Suppose that assumptions 4 and 5 hold and that $\epsilon := \inf \operatorname{supp}(\tau) \in (0, \tau_{reach}^2)$, then for any $c \in (0, 1)$ and sufficiently large N, the score matching algorithm A_{erm} is score stable with,

$$\varepsilon_{stab}^2 \lesssim C \left(C C_{sm} N^{-2} + \mathbb{E}[\hat{\ell}_{sm}(\hat{s})] \right)^c, \qquad C = \frac{D_{\mathcal{H}}^2}{\sigma^4} \vee \frac{1}{c_{cr} \sigma^{d^*}}.$$

Proof. We use the shorthand $\hat{\ell}_{\rm sm}(s) = \hat{\ell}_{\rm sm}(s;S,\tau), \hat{\ell}_{\rm dsm}(s) = \hat{\ell}_{\rm dsm}(s;S,\tau), \ell_{\rm sm}(s) = \ell_{\rm sm}(s;\tau)$ for the sake of brevity. We start from Lemma 19 which provides a bound on the difference between \hat{s}^i and \hat{s} in $L^2(\hat{p}_t)$ and use it to develop a bound in $L^2(\hat{p}_{t|0}(\cdot|\tilde{x}))$, as required by score stability. In particular, we define the quantity

$$\varepsilon^2 = \mathbb{E} \left[\int \int \|\hat{s}^i(y,t) - \hat{s}(y,t)\|^2 p_{t|0}(dy|x_i) \tau(dt) \right],$$

so that, by the symmetric of the algorithm, $A_{\rm erm}$ is score stable with constant ε (we have that $\varepsilon < \infty$ from Assumption 5. Therefore, from Lemma 19, we have

$$\mathbb{E}\left[\int\int \|\hat{s}^i(y,t) - \hat{s}(y,t)\|^2 \,\hat{p}_t(dy) \,\tau(dt)\right] \leq 8\mathbb{E}[\hat{\ell}_{\rm sm}(\hat{s})] + \frac{8}{N}\varepsilon(C_{\rm sm}^{1/2} + \varepsilon).$$

We proceed using Lemma 20 with $\phi(y) = \|\hat{s}^i(y,t) - \hat{s}(y,t)\|^2$ to obtain that for any $j \in [N]$, p > 1,

$$\int \|\hat{s}^{i}(y,t) - \hat{s}(y,t)\|^{2} p_{t|0}(dy|x_{i})$$

$$\leq \left(\int \|\hat{s}^{i}(y,t) - \hat{s}(y,t)\|^{2p} p_{t|0}(dy|x_{j})\right)^{1/p} \exp\left(\frac{\mu_{t}^{2} \|x_{i} - x_{j}\|^{2}}{2(p-1)\sigma_{t}^{2}}\right)$$

Given any subset of the dataset $B \subset S$ with $x_i \in B$ we can average over the above bound to obtain,

$$\int \|\hat{s}^{i}(y,t) - \hat{s}(y,t)\|^{2} p_{t|0}(dy|x_{i})
\leq \frac{1}{|B|} \sum_{x \in B} \left(\int \|\hat{s}^{i}(y,t) - \hat{s}(y,t)\|^{2p} p_{t|0}(dy|x) \right)^{1/p} \exp\left(\frac{\mu_{t}^{2} \operatorname{diam}(B)^{2}}{2(p-1)\sigma_{t}^{2}} \right)
\leq \left(\frac{1}{|B|} \sum_{x \in B} \int \|\hat{s}^{i}(y,t) - \hat{s}(y,t)\|^{2p} p_{t|0}(dy|x) \right)^{1/p} \exp\left(\frac{\mu_{t}^{2} \operatorname{diam}(B)^{2}}{2(p-1)\sigma_{t}^{2}} \right)
\leq \hat{\nu}(B)^{-1/p} \left(\int \|\hat{s}^{i}(y,t) - \hat{s}(y,t)\|^{2p} \hat{p}_{t}(dy) \right)^{1/p} \exp\left(\frac{\mu_{t}^{2} \operatorname{diam}(B)^{2}}{2(p-1)\sigma_{t}^{2}} \right)
\leq (D_{\mathcal{H}}/\sigma_{t}^{2})^{2(1-1/p)} \hat{\nu}(B)^{-1/p} \left(\int \|\hat{s}^{i}(y,t) - \hat{s}(y,t)\|^{2} \hat{p}_{t}(dy) \right)^{\frac{1}{p}} \exp\left(\frac{\mu_{t}^{2} \operatorname{diam}(B)^{2}}{2(p-1)\sigma_{t}^{2}} \right),$$

where in the final inequality we use the L^{∞} bound in Assumption 5. Integrating with respect to τ and taking the expectation, we obtain,

$$\varepsilon^{2} \leq (D_{\mathcal{H}}/\sigma_{\epsilon}^{2})^{2/q} \mathbb{E}\Big[\hat{\nu}(B)^{-1/p}\Big] \mathbb{E}\Big[\int \int \|\hat{s}^{i}(y,t) - \hat{s}(y,t)\|^{2} \hat{p}_{t}(dy)\tau(dt)\Big] \exp\left(\frac{\mu_{\epsilon}^{2} \operatorname{diam}(B)^{2}}{2(p-1)\sigma_{\epsilon}^{2}}\right)$$

$$\leq (D_{\mathcal{H}}/\sigma_{\epsilon}^{2})^{2/q} \mathbb{E}\Big[\hat{\nu}(B)^{-q/p}\Big]^{1/q} \left(8\mathbb{E}[\hat{\ell}_{sm}(\hat{s})] + \frac{8}{N}\varepsilon(C_{sm}^{1/2} + \varepsilon)\right)^{1/p} \exp\left(\frac{\mu_{\epsilon}^{2} \operatorname{diam}(B)^{2}}{2(p-1)\sigma^{2}}\right),$$

where we define $q := (1 - 1/p)^{-1}$. Using Young's inequality, it follows that for any $\lambda > 0$,

$$\varepsilon^2 \leq \frac{D_{\mathcal{H}}^2}{\sigma_{\epsilon}^4 \lambda^q q} \mathbb{E}\Big[\hat{\nu}(B)^{-q/p}\Big] \exp\bigg(\frac{q\mu_{\epsilon}^2 \operatorname{diam}(B)^2}{2(p-1)\sigma_{\epsilon}^2}\bigg) + \frac{\lambda^p}{p} \bigg(8\mathbb{E}[\hat{\ell}_{\mathrm{sm}}(\hat{s})] + \frac{8}{N}\varepsilon(C_{\mathrm{sm}}^{1/2} + \varepsilon)\bigg).$$

Setting $\kappa := 8\lambda^p/pN$, we can rearrange this to obtain the quadratic inequality.

$$(1-\kappa)\varepsilon^2 - C_{\mathrm{sm}}^{1/2}\kappa\varepsilon \le \left(\frac{8}{Np\kappa}\right)^{q/p} \frac{D_{\mathcal{H}}^2}{\sigma_{\epsilon}^4 q} \mathbb{E}\Big[\hat{\nu}(B)^{-q/p}\Big] \exp\left(\frac{q\mu_{\epsilon}^2 \operatorname{diam}(B)^2}{2(p-1)\sigma_{\epsilon}^2}\right) + N\kappa \mathbb{E}[\hat{\ell}_{\mathrm{sm}}(\hat{s})].$$

Requiring that $\kappa \leq 1/2$, we solve the quadratic to obtain the inequality,

$$\frac{\varepsilon^2}{4} \le C_{\rm sm} \kappa^2 + \left(\frac{8}{Np\kappa}\right)^{q/p} \frac{D_{\mathcal{H}}^2}{\sigma_{\epsilon}^4 q} \mathbb{E}\left[\hat{\nu}(B)^{-q/p}\right] \exp\left(\frac{q\mu_{\epsilon}^2 \operatorname{diam}(B)^2}{2(p-1)\sigma_{\epsilon}^2}\right) + N\kappa \mathbb{E}\left[\hat{\ell}_{\rm sm}(\hat{s})\right]. \tag{34}$$

Next, we optimise B by setting $B=B_{\sigma_{\epsilon}}(x_i)\cap S$. We apply Lemma 21 with $\phi(r)=r^{-q/p}$ to obtain that whenever $\sigma_{\epsilon}\leq \tau_{\text{reach}}$ we obtain,

$$\mathbb{E}\Big[\hat{\nu}(B)^{-q/p}\Big] \le N^{q/p} \exp(-c_{\nu}N^{2}r^{d^{*}}) + \left(\frac{2}{c_{\nu}r^{d^{*}}}\right)^{q/p}$$

$$\le 2\left(\frac{2}{c_{\nu}\sigma_{\epsilon}^{d^{*}}}\right)^{q/p},$$

where the second inequality holds whenever $N \ge q/2p$. Returning to (34), it follows from the above that

$$\frac{\varepsilon^2}{4} \le C_{\rm sm} \kappa^2 + \left(\frac{16}{Npc_{\nu}\sigma_{\epsilon}^{d^*}\kappa}\right)^{q/p} \frac{2D_{\mathcal{H}}^2}{\sigma_{\epsilon}^4 q} \exp\left(\frac{2q}{p-1}\right) + N\kappa \mathbb{E}[\hat{\ell}_{\rm sm}(\hat{s})]. \tag{35}$$

We now choose κ by optimising the second two terms of this bound, by which we arrive at the choice

$$\kappa^{q/p+1} = \frac{2D_{\mathcal{H}}^2}{\sigma_{\epsilon}^4 p N \gamma} \exp{\left(\frac{2q}{p-1}\right)} \left(\frac{16}{N p c_{\nu} \sigma_{\epsilon}^{d^*}}\right)^{q/p},$$

for some $\gamma > 0$. Substituting this in to (35), we arrive at the bound

$$\begin{split} \frac{\varepsilon^2}{4} &\leq C_{\mathrm{sm}}(Np)^{-2} \bigg(\frac{2D_{\mathcal{H}}^2}{\sigma_{\epsilon}^4}\bigg)^{2/q} \exp\bigg(\frac{4}{p-1}\bigg) \bigg(\frac{16}{c_{\nu}\sigma_{\epsilon}^{d^*}}\bigg)^{2/p} \gamma^{-1/q} \\ &+ \bigg(\frac{2D_{\mathcal{H}}^2}{\sigma_{\epsilon}^4}\bigg)^{1/q} \exp\bigg(\frac{2}{p-1}\bigg) \bigg(\frac{16}{c_{\nu}\sigma_{\epsilon}^{d^*}}\bigg)^{1/p} \bigg(\frac{\gamma^{1/p}}{q} + \frac{1}{p\gamma^{1/q}} \mathbb{E}[\hat{\ell}_{\mathrm{sm}}(\hat{s})]\bigg). \end{split}$$

Optimising γ leads to the bound,

$$\frac{\varepsilon^{2}}{4} \leq \left(\frac{2D_{\mathcal{H}}^{2}}{\sigma_{\epsilon}^{4}}\right)^{\frac{1}{q}} \exp\left(\frac{2}{p-1}\right) \left(\frac{16}{c_{\nu}\sigma_{\epsilon}^{d^{*}}}\right)^{\frac{1}{p}} \left(C_{\text{sm}}N^{-2}\left(\frac{2D_{\mathcal{H}}^{2}}{\sigma_{\epsilon}^{4}}\right)^{\frac{1}{q}} \exp\left(\frac{2}{p-1}\right) \left(\frac{16}{c_{\nu}\sigma_{\epsilon}^{d^{*}}}\right)^{\frac{1}{p}} + \mathbb{E}[\hat{\ell}_{\text{sm}}(\hat{s})]\right)^{\frac{1}{p}} \\
\leq \left(\frac{2D_{\mathcal{H}}^{2}}{\sigma_{\epsilon}^{4}} \vee \frac{16}{c_{\nu}\sigma_{\epsilon}^{d^{*}}}\right) \exp\left(\frac{4}{p-1}\right) \left(\left(\frac{2D_{\mathcal{H}}^{2}}{\sigma_{\epsilon}^{4}} \vee \frac{16}{c_{\nu}\sigma_{\epsilon}^{d^{*}}}\right) C_{\text{sm}}N^{-2} + \mathbb{E}[\hat{\ell}_{\text{sm}}(\hat{s})]\right)^{1/p}$$

Optimising p, we obtain

$$\frac{\varepsilon^2}{4} \lesssim \left(\frac{2D_{\mathcal{H}}^2}{\sigma_{\epsilon}^4} \vee \frac{16}{c_{\nu}\sigma_{\epsilon}^{d^*}}\right) \exp\left(\frac{5}{2\sqrt{2}}\log(\alpha^{-1})^{1/2} - 2\right) \alpha, \quad \alpha = \left(\frac{2D_{\mathcal{H}}^2}{\sigma_{\epsilon}^4} \vee \frac{16}{c_{\nu}\sigma_{\epsilon}^{d^*}}\right) C_{\mathrm{sm}} N^{-2} + \mathbb{E}[\hat{\ell}_{\mathrm{sm}}(\hat{s})],$$

from which the bound in the statement follows. To obtain that $\kappa \leq 1/2$ and $N \geq q/2q$, it is sufficient to require that N is sufficiently large.

D PROOFS FOR SAMPLING AND SCORE STABILITY

In this section, we provide details for the discretisation scheme considered in Section 5 and give the proof for Proposition 7 and Corollary 8. In the work of Potaptchik et al. (2024), they consider the following discretisation scheme, based on the scheme of (Benton et al., 2024):

$$\hat{y}_{k+1} = \mu_{t_{k+1}-t_k}^{-1} \hat{y}_k + \frac{\sigma_{t_{k+1}-t_k}^2}{\mu_{t_{k+1}-t_k}} s(\hat{y}_k, T - t_k) + \sigma_{t_{k+1}-t_k} \frac{\sigma_{T-t_{k+1}}}{\sigma_{T-t_k}} \zeta_k, \qquad k \in \{0, ..., K-1\},$$

where $\zeta_k \sim N(0, I_d)$ and we recall that the timesteps $(t_k)_{k=0}^K$ are given by,

$$t_k = \begin{cases} \kappa k, & \text{if } k < \frac{T-1}{\kappa}, \\ T - (1+\kappa)^{\frac{T-1}{\kappa} - k}, & \text{if } \frac{T-1}{\kappa} \le k \le K, \end{cases}$$

where $L = \frac{T-1}{\kappa} > 0$, $K = \lfloor L + \log(\epsilon^{-1})/\log(1+\kappa) \rfloor$ and $\kappa > 0, T \ge 1$ is chosen freely. We recall the following result from Potaptchik et al. (2024).

Lemma 22. Suppose that $\alpha = 1$ and Assumption 4 holds with diam $supp(\nu_{data}) \leq 1$. Then, it holds that,

$$D(p_{\epsilon}||A_{\mathrm{em}}(s)) \lesssim \ell_{\mathrm{sm}}(s;\hat{\tau}) + D(p_T||p_{\infty}) + \Delta_{\kappa,K},$$

$$\Delta_{\kappa,K} = \kappa + d^*\kappa^2 (K - L)(\log(\epsilon^{-1}) + \sup|\log(p_{\nu})|),$$

where we define the measure,

$$\hat{\tau}(dt) = \frac{1}{K} \sum_{k=0}^{K-1} \delta_{T-t_k}(dt).$$

D.1 COARSE DISCRETISATION AND REGULARISATION

Fix $\epsilon > 0$ and suppose that κ is such that $\log(\epsilon^{-1})/\log(1+\kappa)$ is an integer. Set $K = L + \log(\epsilon^{-1})/\log(1+\kappa)$ so that, according to the discretisation scheme,

$$t_K = T - (1 + \kappa)^{-\log(\epsilon^{-1})/\log(1+\kappa)} = T - \epsilon.$$

Proof of Proposition 7. Let $\hat{s} = A_{\text{erm}}(S)$. We begin with Lemma 22, which provides the bound,

$$\mathbb{E}[D(p_{\epsilon}||A_{\mathrm{em}}(\hat{s}))] \lesssim \mathbb{E}[\ell_{\mathrm{sm}}(\hat{s};S,\hat{\tau})] + D(p_T||p_{\infty}) + \Delta_{\kappa,K}.$$

For ϵ sufficiently small we have the bound,

$$\Delta_{\kappa,K} = \kappa + d^* \kappa^2 \frac{\log(\epsilon^{-1})}{\log(1+\kappa)} (\log(\epsilon^{-1}) + \sup|\log(p_{\nu})|)$$

$$\leq \kappa (1+\kappa) d^* \log(\epsilon^{-1})^2.$$

Using Theorem 3, we obtain that if the algorithm is $\varepsilon_{\text{stab}}$ -score stable, we have

$$egin{aligned} \mathbb{E}[\ell_{\mathrm{sm}}(\hat{s};\hat{ au})] &\lesssim \mathbb{E}[\hat{\ell}_{\mathrm{sm}}(\hat{s};S,\hat{ au})] + arepsilon_{\mathrm{stab}} \mathbb{E}[\hat{\ell}_{\mathrm{dsm}}(\hat{s};S,\hat{ au})]^{1/2} + arepsilon_{\mathrm{stab}}^2 \ &\lesssim \mathbb{E}[\hat{\ell}_{\mathrm{sm}}(\hat{s};\hat{ au})] + arepsilon_{\mathrm{stab}} C_{\mathrm{sm}}^{1/2} + arepsilon_{\mathrm{stab}}^2 \end{aligned}$$

Using Proposition 6 we obtain that with $\tau = \hat{\tau}$, $A_{\rm erm}$ is score stable, with constant,

$$\begin{split} \varepsilon_{\mathrm{stab}}^2 &\lesssim C \Big(C C_{\mathrm{sm}} N^{-2} + \mathbb{E}[\hat{\ell}_{\mathrm{sm}}(\hat{s})] \Big)^c \\ &\lesssim c_{\nu}^{-1} \sigma_{T-t_{K-1}}^{-d^*} \Big(c_{\nu}^{-1} \sigma_{T-t_{K-1}}^{-d^*} C_{\mathrm{sm}} N^{-2} + \mathbb{E}[\hat{\ell}_{\mathrm{sm}}(\hat{s})] \Big)^c. \end{split}$$

Now by definition, we have that

$$T - t_{K-1} = (1 + \kappa)^{L-K+1} = \epsilon (1 + \kappa),$$

so if we take ϵ, κ sufficiently small so that $\epsilon(1+\kappa) \leq \frac{1}{2}$, we also have $\sigma_{\epsilon(1+\kappa)}^2 \geq \epsilon(1+\kappa)$ and thus we obtain,

$$\begin{split} \varepsilon_{\text{stab}}^2 &\lesssim c_{\nu}^{-1} \epsilon^{-d^*/2} (1+\kappa)^{-d^*/2} \Big(c_{\nu}^{-1} \epsilon^{-d^*/2} (1+\kappa)^{-d^*/2} C_{\text{sm}} N^{-2} + \mathbb{E}[\hat{\ell}_{\text{sm}}(\hat{s})] \Big)^c . \\ &\lesssim c_{\nu}^{-1} \epsilon^{-d^*} (1+\kappa)^{-d^*} \Big(c_{\nu}^{-1} C_{\text{sm}} N^{-2} + \mathbb{E}[\hat{\ell}_{\text{sm}}(\hat{s})] \Big)^c , \end{split}$$

where in the last inequality, we use that $\epsilon(1+\kappa) \leq 1/2$.

We now proceed by proving Corollary 8 in which the bound in Proposition 7 is optimised.

Proof of Corollary 8. Let $\tilde{\tau}_{\epsilon}$ denote the weak limit of the measure τ_{κ} as $\kappa \to 0^+$. Since $\operatorname{supp}(\tilde{\tau}_{\epsilon}) \subseteq [\epsilon, T]$ and $\epsilon > 0$, we know that $\inf_{\mathcal{H}} \hat{\ell}_{sm}(\cdot; S, \tilde{\tau}_{\epsilon}) < \infty$. From this, we deduce that $\lim_{\kappa \to 0^+} B_{\kappa} < \infty$.

With this there exists $\kappa^* \geq 1$ which is the smallest quantity satisfying,

$$(1 + \kappa^*)^{2d^* + 2} = \frac{B_{\kappa^*}}{\log(\epsilon^{-1})^2} \vee 1.$$

In the case that $B_{\kappa^*} > \log(\epsilon^{-1})$, we have that

$$\begin{split} B_{\kappa^*}^{1/2} (1+\kappa^*)^{-d^*} + \frac{B_{\kappa^*}}{C_{\mathrm{sm}}} (1+\kappa^*)^{-2d^*} + \kappa^* (1+\kappa^*) d^* \log(\epsilon^{-1})^2 \\ = B_{\kappa^*}^{\frac{1}{2(d^*+1)}} \log(\epsilon^{-1})^{\frac{d^*}{d^*+1}} + (C_{\mathrm{sm}}^{-1} + d^*) B_{\kappa^*}^{\frac{1}{d^*+1}} \\ \leq B_{\kappa^*}^{\frac{1}{2(d^*+1)}} \log(\epsilon^{-1}) + (C_{\mathrm{sm}}^{-1} + d^*) B_{\kappa^*}^{\frac{1}{d^*+1}} \log(\epsilon^{-1})^2. \end{split}$$

Plus, if $B_{\kappa^*} \leq \log(\epsilon^{-1})$ and therefore $\kappa^* = 1$, then there exists κ such that,

$$B_{\kappa}^{1/2}(1+\kappa)^{-d^*} + \frac{B_{\kappa}}{C_{\rm sm}}(1+\kappa)^{-2d^*} + \kappa(1+\kappa)d^*\log(\epsilon^{-1})^2 \lesssim B_{\kappa}^{1/2} + \frac{B_{\kappa}}{C_{\rm sm}} + de^{-T}.$$

Combining these leads to the bound in the statement.

E PROOFS FOR STABILITY OF SGD

In this section, we analyse the stochastic optimisation scheme in (14), deriving the score stability bounds given in Proposition 11. We begin with a basic lemma that follows from weight decay and gradient clipping.

Lemma 23. Suppose that $\eta_k < \lambda^{-1}$ for all $k \in \mathbb{N}$, then for any $K \in \mathbb{N}$, it holds that

$$\|\theta_K\| \le \frac{Ce}{\lambda} \vee \|\theta_0\|.$$

Proof. We begin with the bound,

$$\|\theta_{k+1}\| \le (1 - \eta_k \lambda) \|\theta_k\| + \eta_k \|\operatorname{Clip}_C(G_k(\theta_k, \{x_i\}_{i \in B_k}))\|$$

$$\le (1 - \eta_k \lambda) \|\theta_k\| + \eta_k C.$$

By comparison, this leads to the bound

$$\begin{aligned} \|\theta_{k}\| &\leq C \sum_{k=0}^{K-1} \eta_{k} \prod_{i=k+1}^{K-1} (1 - \eta_{k} \lambda) + \prod_{k=0}^{K-1} (1 - \eta_{k} \lambda) \|\theta_{0}\| \\ &\leq C \sum_{k=0}^{K-1} \eta_{k} \exp\left(\lambda \sum_{i=0}^{k} \eta_{k}\right) + \exp\left(-\lambda \sum_{i=0}^{K-1} \eta_{k}\right) \|\theta_{0}\| \\ &\leq C \exp\left(-\lambda \sum_{i=0}^{K-1} \eta_{k} + \lambda \max_{k} \eta_{k}\right) \sum_{k=0}^{K-1} \eta_{k} \exp\left(\lambda \sum_{i=0}^{k-1} \eta_{k}\right) + \exp\left(-\lambda \sum_{i=0}^{K-1} \eta_{k}\right) \|\theta_{0}\| \end{aligned}$$

Since the sum forms a left Riemann sum, approximating an integral of an increasing function, we can upper bound it by the integral over $\exp(\lambda t)$. Furthermore, we have that $\lambda \max_k \eta_k \le 1$, which leads to the bound.

$$\|\theta_k\| \le Ce \exp\left(-\lambda \sum_{i=0}^{K-1} \eta_k\right) \int_0^{\sum_{k=0}^{K-1} \eta_k} \exp(\lambda t) dt + \exp\left(-\lambda \sum_{i=0}^{K-1} \eta_k\right) \|\theta_0\|$$

$$\le \frac{Ce}{\lambda} \left(1 - \exp\left(-\lambda \sum_{k=0}^{K-1} \eta_k\right)\right) + \exp\left(-\lambda \sum_{k=0}^{K-1} \eta_k\right) \|\theta_0\|$$

$$\le \frac{Ce}{\lambda} \vee \|\theta_0\|.$$

We are now ready to prove Proposition 11.

Proposition 11. Consider the score matching algorithm $A_{sm}: S \mapsto s_{\theta_K}$ for some fixed $K \in \mathbb{N}$ where $(\theta_k)_k$ is as given in (14). Suppose that assumptions 9 and 10 hold and $\eta_k \leq \bar{\eta}/k$ for all k < K, for some $\bar{\eta} \in (0, \lambda^{-1})$. Then, we obtain that A_{sm} is score stable with constant,

$$\varepsilon_{stab}^2 \lesssim \left(\frac{C}{\lambda} \vee R\right)^{1 + \frac{\bar{\eta}v}{\bar{\eta}v + 1}} \frac{\overline{L}^2}{(\bar{\eta}v) \vee 1} \left(\frac{C}{\bar{\eta}}\right)^{\frac{1}{\bar{\eta}v + 1}} \frac{N_B K^{\frac{\bar{\eta}v}{\bar{\eta}v + 1}}}{N},$$

where
$$R^2 = \mathbb{E}[\|\theta_0\|^2]$$
, $v = (\overline{M}B_{\ell}C_{\tau}^{1/2} + \overline{L}^2 - \lambda) \vee 0$ and $C_{\tau} = \int \sigma_t^{-4} \tau(dt)$.

Proof. Since the stochastic mini-batch scheme, and therefore the resulting score matching algorithm, is symmetric to dataset permutations, we consider stability under changes in the N^{th} entry of the dataset, without loss of generality. Let θ_k be the process given in (14), using the dataset S and let $\tilde{\theta}_k$ be the same process using S^N instead of S:

$$\tilde{\theta}_{k+1} = (1 - \eta \lambda)\tilde{\theta}_k - \eta_k \operatorname{Clip}_C(G_k(\tilde{\theta}_p, \{\tilde{x}_i\}_{i \in B_k})), \qquad \tilde{\theta}_0 = \theta_0,$$

where $\tilde{x}_i = x_i$ for $i \neq N$, $\tilde{x}_N = \tilde{x}$. By having the processes share the same mini-batch indices B_k and gradient approximation G_k (i.e. sharing the same random time variables $t_{i,j}$ and noise $\xi_{i,j}$), we couple the processes θ_k and $\tilde{\theta}_k$.

 We proceed by first controlling the stability of the gradient estimator, computing the bound,

$$\begin{aligned} &\|G_k(\theta_k,(x_i)_{i\in B_k}) - G_k(\tilde{\theta}_k,(x_i)_{i\in B_k})\| \\ &1245 \\ &1246 \\ &1247 \\ &1248 \\ &1249 \\ &1249 \\ &1249 \\ &1250 \\ &1251 \\ &1252 \\ &1252 \\ &1253 \\ &1254 \\ &1255 \\ &1255 \\ &1255 \\ &1255 \\ &1256 \\ &1257 \\ &1258 \end{aligned} \qquad \begin{aligned} &\|G_k(\theta_k,(x_i)_{i\in B_k}) - G_k(\tilde{\theta}_k,(x_i)_{i\in B_k})\| \\ &\leq \frac{1}{N_PP} \sum_{i\in B_k} \sum_{j=1}^P w_{t_{i,j}} \Big(\|\nabla s_{\theta_k}(X_{t_{i,j}},t_{i,j}) - \nabla \log p_{t_{i,j}}|_0(X_{t_{i,j}}|x_i)) \| \\ &\leq \frac{1}{N_PP} \sum_{i\in B_k} \sum_{j=1}^P w_{t_{i,j}} \Big(\|\nabla s_{\theta_k}(X_{t_{i,j}},t_{i,j}) - \nabla s_{\tilde{\theta}_k}(X_{t_{i,j}},t_{i,j}) \| \|s_{\theta_k}(X_{t_{i,j}},t_{i,j}) \| \\ &\leq \frac{1}{N_PP} \sum_{i\in B_k} \sum_{j=1}^P w_{t_{i,j}} \Big(\|\nabla s_{\theta_k}(X_{t_{i,j}},t_{i,j}) - \nabla s_{\tilde{\theta}_k}(X_{t_{i,j}},t_{i,j}) \| \|s_{\theta_k}(X_{t_{i,j}},t_{i,j}) - S_{\tilde{\theta}_k}(X_{t_{i,j}},t_{i,j}) \| \Big) \\ &\leq \frac{1}{N_PP} \sum_{i\in B_k} \sum_{j=1}^P w_{t_{i,j}} \Big(M(X_{t_{i,j}},t_{i,j}) \| s_{\theta_k}(X_{t_{i,j}},t_{i,j}) - \nabla \log p_{t_{i,j}|0}(X_{t_{i,j}}|x_j) \| \\ &\leq \frac{1}{N_PP} \sum_{i\in B_k} \sum_{j=1}^P w_{t_{i,j}} \Big(M(X_{t_{i,j}},t_{i,j}) \| s_{\theta_k}(X_{t_{i,j}},t_{i,j}) - \nabla \log p_{t_{i,j}|0}(X_{t_{i,j}}|x_j) \| \\ &\leq \frac{1}{N_PP} \sum_{i\in B_k} \sum_{j=1}^P w_{t_{i,j}} \Big(M(X_{t_{i,j}},t_{i,j}) \| s_{\theta_k}(X_{t_{i,j}},t_{i,j}) - \nabla \log p_{t_{i,j}|0}(X_{t_{i,j}}|x_j) \| \\ &\leq \frac{1}{N_PP} \sum_{i\in B_k} \sum_{j=1}^P w_{t_{i,j}} \Big(M(X_{t_{i,j}},t_{i,j}) \| s_{\theta_k}(X_{t_{i,j}},t_{i,j}) - \nabla \log p_{t_{i,j}|0}(X_{t_{i,j}}|x_j) \| \\ &\leq \frac{1}{N_PP} \sum_{i\in B_k} \sum_{j=1}^P w_{t_{i,j}} \Big(M(X_{t_{i,j}},t_{i,j}) \| s_{\theta_k}(X_{t_{i,j}},t_{i,j}) - \nabla \log p_{t_{i,j}|0}(X_{t_{i,j}},t_{i,j}) \| \\ &\leq \frac{1}{N_PP} \sum_{i\in B_k} \sum_{j=1}^P w_{t_{i,j}} \Big(M(X_{t_{i,j}},t_{i,j}) \| s_{\theta_k}(X_{t_{i,j}},t_{i,j}) - \nabla \log p_{t_{i,j}|0}(X_{t_{i,j}},t_{i,j}) \| \\ &\leq \frac{1}{N_PP} \sum_{i\in B_k} \sum_{j=1}^P w_{t_{i,j}} \Big(M(X_{t_{i,j}},t_{i,j}) \| s_{\theta_k}(X_{t_{i,j}},t_{i,j}) - \nabla \log p_{t_{i,j}|0}(X_{t_{i,j}},t_{i,j}) \| \\ &\leq \frac{1}{N_PP} \sum_{i\in B_k} \sum_{j=1}^P w_{t_{i,j}} \Big(M(X_{t_{i,j}},t_{i,j}) \| s_{\theta_k}(X_{t_{i,j}},t_{i,j}) - \nabla \log p_{t_{i,j}|0}(X_{t_{i,j}},t_{i,j}) \| s_{\theta_k}(X_{t_{i,j}},t_{i,j}) - \nabla \log p_{t_{i,j}|0}(X_{t_{i,j}},t_{i,j}) \| s_{\theta_k}(X_{t_{i,j}},t_{i,j}) - \nabla \log p_{t_{i,j}$$

We control the expectation of this by first noting that,

$$\mathbb{E}\Big[w_{t_{i,j}}\Big(M(X_{t_{i,j}},t_{i,j})\|s_{\theta_k}(X_{t_{i,j}},t_{i,j}) - \nabla \log p_{t_{i,j}|0}(X_{t_{i,j}}|x_j)\| + L(X_{t_{i,j}},t_{i,j})^2\Big)\Big|\theta_k,\tilde{\theta}_k,S,\tilde{x}\Big] \\
\leq \left(\int \mathbb{E}[M(X_t,t)^2|X_0=x_i]\tau(dt)\right)^{1/2} \left(\int \hat{\ell}_{dsm}(s_{\theta_k};\{x_i\},\delta_t)\tau(dt)\right)^{1/2} \\
+ \int \mathbb{E}[L(X_t,t)^2|X_0=x_i]\tau(dt) \\
\leq \overline{M}B_{\ell}C_{\tau}^{1/2} + \overline{L}^2,$$

where we define the quantity $C_{\tau} := \int \sigma_t^{-4} \tau(dt)$. From this, it follows that

$$\mathbb{E}\Big[\|G_k(\theta_k,(x_i)_{i\in B_k}) - G_k(\tilde{\theta}_k,(x_i)_{i\in B_k})\|\Big|\theta_k,\tilde{\theta}_k,S,\tilde{x}\Big] \le \left(\overline{M}B_\ell C_\tau^{1/2} + \overline{L}^2\right)\|\theta_k - \tilde{\theta}_k\|.$$

Furthermore, we can control the difference between $G_k(\tilde{\theta}_k,(x_i)_{i\in B_k})$ and $G_k(\tilde{\theta}_k,(\tilde{x}_i)_{i\in B_k})$, using the fact that they are identical whenever $N \notin B_k$. Thus, obtaining,

$$\begin{split} \mathbb{E}\Big[\| \mathrm{Clip}_C(G(\theta_k, (x_i)_{i \in B_k})) - \mathrm{Clip}_C(G(\tilde{\theta}_k, (\tilde{x}_i)_{i \in B_k})) \| \Big| \theta_k, \tilde{\theta}_k, S, \tilde{x} \Big] \\ &\leq \mathbb{E}\Big[\| G(\theta_k, (x_i)_{i \in B_k}) - G(\tilde{\theta}_k, (x_i)_{i \in B_k}) \| \Big| \theta_k, \tilde{\theta}_k, S, \tilde{x} \Big] \\ &+ \mathbb{E}\Big[\| \mathrm{Clip}_C(G(\tilde{\theta}_k, (x_i)_{i \in B_k})) - \mathrm{Clip}_C(G(\tilde{\theta}_k, (\tilde{x}_i)_{i \in B_k})) \| \Big| \theta_k, \tilde{\theta}_k, S, \tilde{x} \Big] \\ &\leq \Big(\overline{M} B_\ell C_\tau^{1/2} + \overline{L}^2 \Big) \| \theta_k - \tilde{\theta}_k \| + 2C \frac{N_B}{N}, \end{split}$$

where we have used the fact that $\mathbb{P}(N \in B_k) = \frac{N_B}{N}$. Thus, using (14), we obtain that for any $k_0 \leq k$,

$$\begin{split} \mathbb{E}\Big[\|\theta_{k+1} - \tilde{\theta}_{k+1}\| \Big| \theta_{k_0}, \tilde{\theta}_{k_0}, S, \tilde{x} \Big] \\ & \leq \Big(1 + \eta_k \Big(\overline{M} B_\ell C_\tau^{1/2} + \overline{L}^2 - \lambda \Big) \Big) \, \mathbb{E}\Big[\|\theta_k - \tilde{\theta}_k\| \Big| \theta_{k_0}, \tilde{\theta}_{k_0}, S, \tilde{x} \Big] + 2\eta_k C \frac{N_B}{N}. \\ & \leq (1 + \eta_k \upsilon) \, \mathbb{E}\Big[\|\theta_k - \tilde{\theta}_k\| \Big| \theta_{k_0}, \tilde{\theta}_{k_0}, S, \tilde{x} \Big] + 2\eta_k C \frac{N_B}{N}, \end{split}$$

where $v = \overline{M}B_{\ell}C_{\tau}^{1/2} + \overline{L}^2 - \lambda$. Thus, by comparison, we obtain,

$$\mathbb{E}\Big[\|\theta_K - \tilde{\theta}_K\| \Big| \theta_{k_0}, \tilde{\theta}_{k_0}, S, \tilde{x}\Big] \leq \sum_{i=k_0}^{K-1} 2\eta_i C \frac{N_B}{N} \prod_{j=i+1}^{K-1} (1 + \eta_j v) + \|\theta_{k_0} - \tilde{\theta}_{k_0}\| \prod_{j=k_0}^{K-1} (1 + \eta_j v).$$

From this we obtain the following:

$$\mathbb{E}[\|\theta_K - \tilde{\theta}_K\||\theta_{k_0} = \tilde{\theta}_{k_0}, S, \tilde{x}] \leq 2C \frac{N_B}{N} \sum_{i=k_0}^{K-1} \eta_i \exp\left(\sum_{j=i+1}^{K-1} \eta_j \upsilon\right)$$

$$\leq \frac{2CN_B \bar{\eta}}{N} \sum_{i=k_0}^{K-1} \frac{1}{i} \left(\frac{K}{i}\right)^{\bar{\eta}\upsilon}$$

$$\lesssim \frac{CN_B}{N\upsilon} \left(\frac{K}{k_0}\right)^{\bar{\eta}\upsilon},$$

where we use the fact that $\sum_{j=i+1}^{K-1} \frac{1}{j} \leq \log(K) - \log(i)$. By the law of total probability, we have

$$\begin{split} \mathbb{E}[\|\theta_K - \tilde{\theta}_K\||\theta_0] \\ &= \mathbb{E}[\|\theta_K - \tilde{\theta}_K\||\theta_{k_0} = \tilde{\theta}_{k_0}]\mathbb{P}(\theta_{k_0} = \tilde{\theta}_{k_0}|\theta_0) + \mathbb{E}[\|\theta_K - \tilde{\theta}_K\||\theta_{k_0} \neq \tilde{\theta}_{k_0}, \theta_0]\mathbb{P}(\theta_{k_0} \neq \tilde{\theta}_{k_0}|\theta_0) \\ &\lesssim \frac{CN_B}{N\upsilon} \left(\frac{K}{k_0}\right)^{\bar{\eta}\upsilon} + \left(\frac{Ce}{\lambda} \vee \|\theta_0\|\right) \frac{k_0N_B}{N}, \end{split}$$

where in the second inequality, we use Lemma 23. Thus, optimising k_0 leads to the bound,

$$\mathbb{E}[\|\theta_K - \tilde{\theta}_K\||\theta_0] \lesssim \left(\frac{C}{c}\right)^{\frac{1}{v+1}} (1 + 1/cv) \left(\frac{Ce}{\lambda} \vee \|\theta_0\|\right)^{\frac{cv}{cv+1}} \frac{N_B}{N} K^{\frac{cv}{cv+1}}.$$

Finally, we obtain score stability using the fact that

$$\int \mathbb{E}[\|s_{\theta_K}(X_t, t) - s_{\tilde{\theta}_K}(X_t, t)\|^2 | X_0 = \tilde{x}, S] \tau(dt)
\leq \mathbb{E}\Big[\bar{L}^2 \|\theta_K - \tilde{\theta}_K\|^2\Big]
\leq 2\mathbb{E}\Big[\bar{L}^2 \Big(\frac{Ce}{\lambda} \vee \|\theta_0\|\Big) \|\theta_K - \tilde{\theta}_K\|\Big]
\lesssim \bar{L}^2 \Big(\frac{Ce}{\lambda} \vee R\Big)^{1 + \frac{cv}{cv+1}} \Big(\frac{C}{c}\Big)^{\frac{1}{cv+1}} (1 + 1/cv) \frac{N_B}{N} K^{\frac{cv}{cv+1}},$$

where $R^2 = \mathbb{E}\|\theta_0\|^2$.

F WASSERSTEIN CONTRACTIONS

In this section, we derive the Wasserstein contraction result used in the proof of Proposition 14. We begin with the more abstract problem of deriving Wasserstein contractions for a discrete time diffusion process with anisotropic non-constant volatility. We consider stochastic processes given by the discrete-time update,

$$x_{k+1} = (1 - \eta \lambda)x_k + \eta b(x_k) + \sqrt{\eta}\sigma(x_k)\xi_k,$$
 (36)

$$y_{k+1} = (1 - \eta \lambda)y_k + \eta \tilde{b}(y_k) + \sqrt{\eta} \tilde{\sigma}(y_k) \xi_k, \tag{37}$$

for some $b, \tilde{b}: \mathbb{R}^d \to \mathbb{R}^d, \sigma, \tilde{\sigma}: \mathbb{R}^d \to \mathbb{R}^{d \times d}$ where $\xi_k \sim N(0, I_d)$, and we show that the laws of x_k and y_k contract in Wasserstein distance. We borrow the strategy developed by Eberle (2016) and extended in (Eberle & Majka, 2019; Majka et al., 2020), constructing a coupling and a metric for which exponential contractions of the coupling can be obtained. However, these works are restricted to the setting of isotropic noise with constant volatility (i.e. $\sigma(x) = cI_d$) and so some careful modification to the strategy is required. In particular, we analyse this process with respect to the seminorm $\|\cdot\|_{G^+}$ given by $\|x\|_{G^+}^2 = x^T G^+ x$, where G^+ denotes the Moore-Penrose pseudoinverse of the matrix G. Furthermore, we allow for x_k and y_k to have different bias and volatility terms and so controlling for this will also require some modifications to the proof technique.

To define our coupling we first suppose that there exists a symmetric positive semi-definite matrix $G \in \mathbb{R}^{d \times d}$ such that $\sigma(x), \tilde{\sigma}(y) \succcurlyeq G^{1/2}$ for all $x \in \mathbb{R}^d$, and to couple an update from the above process starting at $x, y \in \mathbb{R}^d$, we first define the update,

$$\tilde{x} = (1 - \eta \lambda)x + \eta b(x), \qquad \tilde{y} = (1 - \eta \lambda)y + \eta \tilde{b}(y),$$

$$\hat{x} = \tilde{x} + \sqrt{\eta}(\sigma(x) - G^{1/2})Z', \qquad \hat{y} = \tilde{y} + \sqrt{\eta}(\tilde{\sigma}(y) - G^{1/2})Z',$$

where $Z' \sim N(0, I_d)$. We then define the synchronous coupled processes,

$$X' = \hat{x} + \sqrt{\eta} G^{1/2} Z$$

$$Y'_{s} = \hat{y} + \sqrt{\eta} G^{1/2} Z,$$

with $Z \sim N(0, I_d)$. We also consider the reflection coupling,

$$Y_r' = \hat{y} + \sqrt{\eta} G^{1/2} \Big(I - 2(G^{1/2})^+ ee^T (G^{1/2})^+ \Big) Z, \quad \text{with } e = (\hat{x} - \hat{y}) / \|\hat{x} - \hat{y}\|_{G^+}$$
 (38)

which has the noise act in the mirrored direction. We combine these couplings to arrive at the final coupling (X', Y'):

$$Y' = \begin{cases} X', & \text{if } \zeta \leq \phi_{\hat{y},\eta G}(X')/\phi_{\hat{x},\eta G}(X'), |\langle e,Z\rangle|^2 < m^2/\eta \text{ and } \hat{r} \leq r_1 \\ Y'_r, & \text{if } \zeta > \phi_{\hat{y},\eta G}(X')/\phi_{\hat{x},\eta G}(X'), |\langle e,Z\rangle|^2 < m^2/\eta \text{ and } \hat{r} \leq r_1 \\ Y'_s, & \text{otherwise,} \end{cases}$$
(39)

for some fixed m > 0.

 We assume the following regularity properties.

Assumption 24. Suppose that b is bounded, satisfying $B := \sup_{x \in \mathbb{R}^n} \|b(x)\|_{G^+} < \infty$ and we have the Lipschitz property, $\|b(x) - b(y)\|_{G^+} \le L_b \|x - y\|_{G^+}$ and $\|\sigma(x) - \sigma(y)\|_{op,G^+} \le L_\sigma \|x - y\|_{G^+}$ for all $x, y \in \mathbb{R}^n$ and for some $L_b, L_\sigma \ge 0$.

We also allow for $b \neq \tilde{b}$ and $\sigma \neq \tilde{\sigma}$, making the following assumption.

Assumption 25. Suppose that b, \tilde{b} satisfy $||b(x) - \tilde{b}(x)||_{G^+} \leq \tilde{B}_b, ||\sigma(x) - \tilde{\sigma}(x)||_{op,G^+} \leq \tilde{B}_{\sigma}$ for all $x \in \mathbb{R}^n$ and for some $\tilde{B}_b, \tilde{B}_{\sigma} \geq 0$.

We define the objects,

$$R = \|x - y\|_{G^+}, \qquad \tilde{r} = \|\tilde{x} - \tilde{y}\|_{G^+}, \qquad \hat{r} = \|\hat{x} - \hat{y}\|_{G^+}, \qquad R' = \|X' - Y'\|_{G^+}.$$

We wish to show that R' contracts in expectation, i.e. it is less than R on average. We modify the metric to guarantee this is possible. We define the function,

$$f(r) = \begin{cases} \frac{1}{a}(1 - e^{-ar}), & \text{if } r \le r_2, \\ \frac{1}{a}(1 - e^{-ar_2}) + \frac{1}{2r_2}e^{-ar_2}(r^2 - r_2^2), & \text{otherwise,} \end{cases}$$

where $a = 6L_b r_1/c_0$, $r_1 = 4(1 + \eta_0 L_b)B/\lambda$, $r_2 = r_1 + \sqrt{\eta_0}$ and c_0, η_0 are defined below. The coupling and the strategy for proving contractions is closely based on an analysis in Majka et al. (2020) and for the sake of comparison, we rely on similar notation. We will also heavily borrow properties of the function f that are proven in this work.

By allowing σ to be non-constant, we run in to additional complications that are controlled by making the following assumption about the scale of L_{σ} .

Assumption 26. *Suppose that the following three inequalities hold:*

$$n-1, (\lambda^2/16L_{\sigma}^2-1)^2 \ge 32\log\left(\frac{8L_{\sigma}(6\vee(4a))\kappa_0^{1/2}}{\sqrt{\eta}(1-e^{-ar_2})c}\right), \qquad L_{\sigma}^2 \le \lambda/8n,$$

for some universal constant κ_0 .

Under these assumptions, we obtain exponential contractions.

Proposition 27. Suppose that assumptions 24, 25 and 26 hold and $m = \sqrt{\eta_0}/2$, then for any $\eta \le \eta_0$ and $x, y \in \mathbb{R}^d$, it holds that

$$\mathbb{E}[f(R')] \le (1 - \eta c/4)f(r) + \frac{3}{2r_2}e^{-ar_2}(\eta^2 \tilde{B}^2 + \eta n \tilde{B}_{\sigma}^2),$$

where

$$c := \min \left\{ e^{-ar_2} \frac{\lambda}{16}, \frac{\frac{1}{2}e^{-ar_2}r_2}{\frac{1}{a}(1 - e^{-ar_2})} \frac{\lambda}{16}, \frac{9L^2r_1^2}{2c_0} e^{-6Lr_1^2/c_0}, \frac{3Lr_1}{16\sqrt{\eta_0}} \right\},$$

$$\eta_0 := \min \left\{ \frac{\lambda}{4L^2}, \frac{16}{\lambda}, \frac{1}{2L}, \frac{2c_0 \log(3/2)\lambda^2}{432L^2B^2}, \frac{4B^2}{\lambda^2}, \frac{c_0^2(\log(2))^2\lambda^2}{2304L^2B^2} \right\},$$

for some universal c_0 and $L = 2(L_b - \lambda)_+ + 4\eta^{-1/2}L_\sigma\sqrt{2(n-1)}$.

F.1 THE COUPLING

Before we provide the proof of Proposition 27, we provide an explanation of how the coupling is arrived at. We begin by discussing the one-dimensional coupling of the Gaussian distribution that the construction is ultimately based on. Consider the following coupling of $\mathcal{N}(t,\eta)$ and $\mathcal{N}(s,\eta)$ for $t,s\in\mathbb{R}$: with $z\sim\mathcal{N}(0,1)$,

$$t' = t + \sqrt{\eta}z,\tag{40}$$

$$s' = \begin{cases} t', & \text{if } \zeta \le \phi_{s,\eta}(t')/\phi_{t,\eta}(t'), |\sqrt{\eta}z| < \tilde{m}, \text{ and } |t-s| \le r_1, \\ s - \sqrt{\eta}z, & \text{if } \zeta > \phi_{s,\eta}(t')/\phi_{t,\eta}(t'), |\sqrt{\eta}z| < \tilde{m}, \text{ and } |t-s| \le r_1, \\ s + \sqrt{\eta}z, & \text{otherwise.} \end{cases}$$

$$(40)$$

This coupling has the following property given in lemmas 3.1 and 3.2 of Majka et al. (2020).

Lemma 28. For the coupling defined in (40) and (41), we have

$$\mathbb{E}[|t' - s'|] = |t - s|,$$

and if $\eta \leq 4\tilde{m}^2$, we have

$$\mathbb{E}\Big[(|t'-s'|-|t-s|)^2\mathbb{1}_{|t'-s'|\in I_{|t-s|}}\Big] \geq \frac{1}{2}c_0\min(\sqrt{\eta},|t-s|)\sqrt{\eta},$$
where
$$I_r = \begin{cases} (0,r+\sqrt{\eta}), & \text{if } r \leq \sqrt{\eta}, \\ (r-\sqrt{\eta},r), & \text{otherwise,} \end{cases}$$

for some universal constant $c_0 > 0$.

Thus, through the second bound, we have control of the probability that |t'-s'| contracts below |t-s|. The coupling proposed in (39) is a multivariate analogue of this that also accounts for the diffusion coefficient $G^{1/2}$. Let the vector $e \in \mathbb{R}^d$ be as defined in (38), then we obtain that,

$$\langle e, G^+ X' \rangle = \langle e, G^+ \hat{x} \rangle + \sqrt{h} \langle (G^{1/2})^+ e, Z \rangle,$$

$$\langle e, G^+ Y'_s \rangle = \langle e, G^+ \hat{y} \rangle + \sqrt{h} \langle (G^{1/2})^+ e, Z \rangle.$$

Therefore, $\langle e, G^+X' \rangle, \langle e, G^+Y'_s \rangle$ are a synchronous coupling of $\mathcal{N}(\langle e, G^+\hat{x} \rangle, h)$ and $\mathcal{N}(\langle e, G^+\hat{y} \rangle, h)$. Furthermore, we have

$$\begin{split} \langle e, G^{+}Y'_{r} \rangle &= \langle e, G^{+}\hat{y} \rangle + \sqrt{h} \langle (G^{1/2})^{+}e, (I - 2(G^{1/2})^{+}ee^{T}(G^{1/2})^{+})Z \rangle \\ &= \langle e, G^{+}\hat{y} \rangle + \sqrt{h} \langle (G^{1/2})^{+}e, Z \rangle - 2\sqrt{h} \langle e, G^{+}e \rangle \langle (G^{1/2})^{+}e, Z \rangle \\ &= \langle e, G^{+}\hat{y} \rangle - \sqrt{h} \langle (G^{1/2})^{+}e, Z \rangle, \end{split}$$

and so $\langle e, G^+X' \rangle$, $\langle e, G^+Y'_r \rangle$ is the one-dimensional reflection coupling. Finally we obtain,

$$\begin{split} \frac{\phi_{\hat{y},\eta G}(X')}{\phi_{\hat{x},\eta G}(X')} &= \frac{\phi_{(G^{1/2})^{+}(\hat{y}-\hat{x}),\eta(G^{1/2})^{+}G^{1/2}}(\sqrt{\eta}Z)}{\phi_{\mathbf{0},\eta(G^{1/2})^{+}G^{1/2}}(\sqrt{\eta}Z)} \\ &= \exp\left(-\frac{1}{2\eta}\|\sqrt{\eta}Z - (G^{1/2})^{+}(\hat{y}-\hat{x})\|_{(G^{1/2})^{+}G^{1/2}}^{2} + \frac{1}{2\eta}\|\sqrt{\eta}Z\|_{(G^{1/2})^{+}G^{1/2}}^{2}\right) \\ &= \exp\left(-\frac{1}{2\eta}\|\hat{y}-\hat{x}\|_{G^{+}}^{2} + \frac{1}{\eta}\sqrt{\eta}\langle(G^{1/2})^{+}(\hat{y}-\hat{x}),Z\rangle\right) \\ &= \exp\left(-\frac{1}{2\eta}|\langle e,G^{+}(\hat{y}-\hat{x})\rangle|^{2} + \frac{1}{\eta}\sqrt{\eta}\langle e,G^{+}(\hat{y}-\hat{x})\rangle\langle(G^{1/2})^{+}e,Z\rangle\right) \\ &= \exp\left(-\frac{1}{2\eta}(\sqrt{\eta}\langle e,G^{+}Z\rangle - \langle e,G^{+}(\hat{y}-\hat{x})\rangle)^{2} + \frac{|\sqrt{\eta}\langle(G^{1/2})^{+}e,Z\rangle|^{2}}{2\eta}\right) \\ &= \frac{\phi_{\langle e,G^{+}(\hat{y}-\hat{x})\rangle,\eta}(\sqrt{\eta}\langle(G^{1/2})^{+}e,Z\rangle)}{\phi_{0,\eta}(\sqrt{\eta}\langle(G^{1/2})^{+}e,Z\rangle)} \\ &= \frac{\phi_{\langle e,G^{+}(\hat{y}),\eta}(\langle e,G^{+}X'\rangle)}{\phi_{\langle e,G^{+}\hat{x}\rangle,\eta}(\langle e,G^{+}X'\rangle)}. \end{split}$$

From this, we deduce that $\langle e, G^+X' \rangle$, $\langle e, G^+Y' \rangle$ are coupled as in (40), (41). The equivalence follows by setting

$$t' = \langle e, G^+ X' \rangle, \qquad s' = \langle e, G^+ Y' \rangle$$
 (42)

$$t = \langle e, G^+ \hat{x} \rangle, \qquad s = \langle e, G^+ \hat{y} \rangle, \qquad z = \langle (G^{1/2})^+ e, Z \rangle.$$
 (43)

Through this equivalence, we can extend the previous lemma to obtain the following result about the high dimensional coupling.

Lemma 29. For the coupling defined in (39), we obtain that for $\eta \leq 4m^2$, we have the following:

$$\mathbb{E}[R'] = \hat{r}, \qquad \mathbb{E}\left[(R' - \hat{r})^2 \mathbb{1}_{R' \in I_{\hat{r}}} \right] \ge \frac{1}{2} c_0 \min(\sqrt{\eta}, \hat{r}) \sqrt{\eta},$$

where c_0 and I_r is as in Lemma 28.

Proof. Let $\{e_i\}_{i=1}^n$ be a basis of \mathbb{R}^n with respect to the inner product $\langle \cdot, \cdot \rangle_{G^+}$ with $e_1 = e$. Then, we have that

$$(R')^{2} = \sum_{i=1}^{n} |\langle e_{i}, G^{+}(X' - Y') \rangle|^{2}$$

$$= |t' - s'|^{2} + \sum_{i=2}^{n} |\langle e_{i}, G^{+}(X' - Y') \rangle|^{2},$$
(44)

where t', s' are as defined in (42). For any $i \neq 1$, we can use that $e_i \perp e$, to obtain that

$$\langle e_i, G^+ Y_r' \rangle = \langle e_i, G^+ \hat{y} \rangle + \sqrt{h} \langle e_i, e \rangle + 2\sqrt{h} \langle e_i, Z \rangle$$
$$= \langle e_i, G^+ \hat{y} \rangle + 2\sqrt{h}z.$$

From this, we obtain that,

$$\langle e_i, G^+(X'-Y'_r)\rangle = \langle e_i, G^+(\hat{x}-\hat{y})\rangle = 0.$$

This also holds for the synchronous coupling and hence we obtain $\langle e_i, G^+(X'-Y') \rangle = 0$. Combined with (44), we obtain that R' = |t' - s'|. Similarly it can be shown that $\hat{r} = |t - s|$ and thus, from Lemma 28, the statement of the lemma follows.

F.2 Proof of Proposition 27

We begin by considering the setting where Z' is truncated Gaussian noise and that $b=\tilde{b},\,\sigma=\tilde{\sigma}.$ We will then extend this to the more general setting in Section F.2.3. We begin by decomposing $\xi\sim N(0,I_d)$ in to directions parallel and perpendicular to the radial vector,

$$\xi_1 = vv^T \xi, \qquad \xi_2 = (I - vv^T)\xi, \qquad v = \frac{\tilde{x} - \tilde{y}}{\|\tilde{x} - \tilde{y}\|_{G^+}}.$$

We then clip each direction according to constants $\bar{z}_1, \bar{z}_2 > 0$ and add them together:

$$Z' = (1 \wedge \bar{z}_1 \|\xi_1\|_{G^+}^{-1})\xi_1 + (1 \wedge \bar{z}_2 \|\xi_2\|_{G^+}^{-1})\xi_2. \tag{45}$$

To prove that the process is contractive, we consider two cases based on the initial distance r.

F.2.1 The case of $r \ge r_1$

When r is large, we can rely on contractive properties following from the weight decay. For this, we obtain the following.

Lemma 30. Suppose that Assumption 24 holds and that $4\bar{z}_1 \leq \lambda L_{\sigma}^{-1} \sqrt{\eta}, 2\bar{z}_2 \leq \sqrt{\lambda} L_{\sigma}, \eta \leq \lambda^{-1}$. Then whenever $r \geq 4B/\lambda$, we have

$$\hat{r} \le \left(1 - \frac{\eta \lambda}{8}\right) r,\tag{46}$$

and when $r < 4B/\lambda$,

$$\hat{r} \le (1 + \eta L)r,\tag{47}$$

where $L = 2(L_b - \lambda)_+ + 4\eta^{-1/2}L_{\sigma}\bar{z}_1$.

Proof. From the triangle inequality and the Lipschitz property of b, we obtain

$$\tilde{r} \le (1 - \eta \lambda) \|x - y\|_{G^+} + \eta \|b(x) - b(y)\|_{G^+}$$

 $\le (1 + \eta(L_b - \lambda)_+)r.$

Alternatively, we can use the fact that $\|b\|_{G^+} \leq B$ to obtain, $\tilde{r} \leq (1 - \eta \lambda)r + 2\eta B$. In particular, if $r \geq 4B/\lambda$, we obtain $\tilde{r} \leq (1 - \eta \lambda/2)r$. Next we bound \hat{r} using the decomposition,

$$\hat{r}^{2} = \|\tilde{x} - \tilde{y} + \sqrt{\eta}(\sigma(x) - \sigma(y))Z'\|_{G^{+}}^{2}$$

$$= \|\tilde{x} - \tilde{y} + \sqrt{\eta}(\sigma(x) - \sigma(y))(1 \wedge \bar{z}_{1}\|\xi_{1}\|_{G^{+}}^{-1})\xi_{1}\|_{G^{+}}^{2} + \|\sqrt{\eta}(\sigma(x) - \sigma(y))(1 \wedge \bar{z}_{2}\|\xi_{2}\|_{G^{+}}^{-1})\xi_{2}\|_{G^{+}}^{2}$$

$$\leq \|\tilde{x} - \tilde{y} + \sqrt{\eta}(\sigma(x) - \sigma(y))(1 \wedge \bar{z}_{1}\|\xi_{1}\|_{G^{+}}^{-1})\xi_{1}\|_{G^{+}}^{2}$$
(48)

The second term is then bounded by,

$$\|\sqrt{\eta}(\sigma(x) - \sigma(y))(1 \wedge \bar{z}_2\|\xi_2\|_{G^+}^{-1})\xi_2\|_{G^+}^2 \le \eta \|\sigma(x) - \sigma(y)\|_{op,G^+} (1 \wedge \bar{z}_2\|\xi_2\|_{G^+}^{-1})^2 \|\xi_2\|_{G^+}^2 \le \eta L_\sigma^2 \bar{z}_2^2 r^2, \tag{49}$$

and the first term is bounded by,

$$\|\tilde{x} - \tilde{y} + \sqrt{\eta}(\sigma(x) - \sigma(y))(1 \wedge \bar{z}_1 \|\xi_1\|_{G^+}^{-1})\xi_1\|_{G^+}^2 \le \tilde{r}^2 + \eta L_{\sigma}^2 \bar{z}_1^2 r^2 + 2\sqrt{\eta} L_{\sigma} \langle v, G^+ \xi_1 \rangle \tilde{r}^2$$

$$\le (1 + 2\sqrt{\eta} L_{\sigma} \bar{z}_1) \tilde{r}^2 + \eta L_{\sigma}^2 \bar{z}_1^2 r^2.$$
 (50)

We substitute (49) and (50) in to (48) to obtain

$$\hat{r}^{2} \leq (1 + 2\sqrt{\eta}L_{\sigma}\bar{z}_{1})\hat{r}^{2} + \eta L_{\sigma}^{2}(\bar{z}_{1}^{2} + \bar{z}_{2}^{2})r^{2}
\leq (1 + \eta(L_{b} - \lambda)_{+})^{2}(1 + 2\sqrt{\eta}L_{\sigma}\bar{z}_{1})r^{2} + \eta L_{\sigma}^{2}(\bar{z}_{1}^{2} + \bar{z}_{2}^{2})r^{2}
\leq (1 + \eta(L_{b} - \lambda)_{+} + 2\eta^{3/2}(L_{b} - \lambda)_{+}L_{\sigma}\bar{z}_{1} + 2\sqrt{\eta}L_{\sigma}\bar{z}_{1} + \eta L_{\sigma}^{2}(\bar{z}_{1}^{2} + \bar{z}_{2}^{2}))r^{2}
\leq (1 + 2\eta(L_{b} - \lambda)_{+} + 4\sqrt{\eta}L_{\sigma}\bar{z}_{1})r^{2},$$

where we have used that $2\eta^{1/2}L_{\sigma}\bar{z}_1 \leq 1$, $\eta^{1/2}L_{\sigma}(z_1^2+z_2^2) \leq \bar{z}_1$, producing the bound in (47). In the case that $r \leq 4B/\lambda$, we can use the fact that $2\eta^{1/2}L_{\sigma}\bar{z}_1 \leq \eta\lambda/2$ and $L_{\sigma}^2(\bar{z}_1^2+\bar{z}_2^2) \leq \lambda/4$ to refine this bound:

$$\hat{r}^2 \le (1 - \eta \lambda/2)^2 (1 + 2\sqrt{\eta} L_{\sigma} \bar{z}_1) r^2 + \eta L_{\sigma}^2 (z_1^2 + z_2^2) r^2$$

$$\le (1 - \eta \lambda/2) (1 - \eta^2 \lambda^2/4)^2 r^2 + \eta L_{\sigma}^2 (z_1^2 + z_2^2) r^2$$

$$< (1 - \eta \lambda/4) r^2.$$

Using the fact that $(1 - \eta \lambda/4)^{1/2} \le 1 - \eta \lambda/8$, we obtain the bound in (46).

We will also need a property of f given in Majka et al. (2020).

Lemma 31. The function f satisfies the property that for all $r \geq r_2$,

$$f\left(\left(1 - \frac{\eta K}{2}\right)r\right) - f(r) \le -\eta c f(r).$$

Using the fact that f is increasing, it follows from lemmas 30 and 31 that,

$$f(\hat{r}) \le f\left(\left(1 - \frac{\eta K}{2}\right)r\right) \le (1 - \eta c)f(r).$$

Thus, to obtain contractions of $\mathbb{E}[f(R')]$ when $r \geq r_1$, it is sufficient to show that $\mathbb{E}[f(R')|Z'] \leq f(\hat{r})$. Note that when $\hat{r} \geq r_1$ or $\|\sqrt{\eta}Z\| \geq m$, the synchronous coupling is used and so $R' = \hat{r}$. Furthermore, if $\hat{r} < r_1$ and $\|\sqrt{\eta}Z\| < m$, we have that $R' \leq r_2$ and thus, using the concavity of f, we deduce that

$$\mathbb{E}[f(R')|Z'] - f(\hat{r}) \le f'(\hat{r})(\mathbb{E}[R'|Z'] - \hat{r}) = 0.$$

Thus, we have shown that whenever $r \geq r_1$, $\mathbb{E}[f(R')|Z'] \leq f(\hat{r})$.

F.2.2 The case of $r < r_1$

When r is small we no longer have contractions due to weight decay and must instead rely on properties of the coupling and function. From Taylor's theorem we have the following:

$$f(R') - f(\hat{r}) = f'(\hat{r})(R' - \hat{r}) + \frac{1}{2} \sup_{\theta} f''(\theta)(R' - \hat{r})^2.$$

where the supremum is between all $\theta \geq 0$ between R' and \hat{r} . We note that in the present setting, $\hat{r} \leq r_1$ also (this follows from Lemma 30) and furthermore $R' - \hat{r} \leq 2m \leq r_2$. Therefore, we can use that f is concave between R' and \hat{r} and so f'' is negative. Using this fact, as well as the fact that $\mathbb{E}[R'|Z'] = \hat{r}$, we obtain the bound,

$$\mathbb{E}[f(R')|Z'] - f(\hat{r}) \leq \frac{1}{2} \mathbb{E}\Big[\sup_{\theta} f''(\theta) (R' - \hat{r})^2 \mathbb{1}_{R' \in I_{\hat{r}}}\Big]$$

$$\leq \frac{1}{2} \sup_{\theta \in I_{\hat{r}}} f''(\theta) \mathbb{E}\Big[(R' - \hat{r})^2 \mathbb{1}_{R' \in I_{\hat{r}}}\Big|Z'\Big]$$

$$\leq \frac{1}{4} \sup_{\theta \in I_{\hat{r}}} f''(\theta) c_0 \min(\sqrt{\eta}, \hat{r}) \sqrt{\eta}.$$

Furthermore, we analyse the contractions between \hat{r} and r using the fact that the function is concave between these values, obtaining,

$$f(\hat{r}) - f(r) \le f'(r)(\hat{r} - r) \le f'(r)\eta Lr.$$

Since we have the derivative $f'(r) = e^{-ar} = f'(\hat{r})e^{-a(r-\hat{r})} \le f'(\hat{r})e^{a\eta Lr_1}$, it holds that

$$f(\hat{r}) - f(r) < f'(\hat{r})e^{a\eta L r_1} \eta L \hat{r},\tag{51}$$

where we have used that $f(\hat{r}) - f(r) \le 0$ holds trivially whenever $r \ge \hat{r}$. Putting these together, we obtain the bound,

$$\mathbb{E}[f(R')|Z'] - f(r) \le f'(\hat{r})e^{a\eta L r_1}\eta L \hat{r} + \frac{1}{4} \sup_{\theta \in I_{\hat{r}}} f''(\theta)c_0 \min(\sqrt{\eta}, \hat{r})\sqrt{\eta}.$$

To complete the analysis of this case, we borrow a result from Majka et al. (2020).

Lemma 32. The function f, satisfies the property that for all $\hat{r} \in [0, r_1]$,

$$f'(\hat{r})e^{a\eta Lr_1}\eta L\hat{r} + \frac{1}{4}c_0\min(\sqrt{\eta},\hat{r})\sqrt{\eta}\sup_{I_{\hat{r}}}f''(\hat{r}) \le -chf(\hat{r}).$$

Between this section and the previous, we have shown that for any $x, y \in \mathbb{R}^n$,

$$\mathbb{E}[f(R')] \le (1 - \eta c/2) f(r),$$

in the setting where Z' is the truncated Gaussian defined in (45) and $b = \tilde{b}, \sigma = \tilde{\sigma}$.

F.2.3 FULL NOISE AND INACCURATE DRIFT

We now consider the more general case of $b \neq b$, $\sigma \neq \tilde{\sigma}$ necessarily and also set $Z' = \xi$, so that it is Gaussian distributed. We do this by borrowing the contraction analysis above. We use the notation $R'' = \|X' - Y'\|_{G^+}$ to not confuse it with R' used above. We obtain,

$$R'' \leq R' + \eta \|b(y) - \tilde{b}(y)\|_{G^{+}} + \sqrt{\eta} \|(\sigma(y) - \tilde{\sigma}(y))\xi\|_{G^{+}} + \sqrt{\eta} \|(\sigma(x) - \sigma(y))(Z' - \xi_{1} - \xi_{2})\|_{G^{+}}$$

$$\leq R' + \eta \tilde{B} + \sqrt{\eta} \tilde{B}_{\sigma} \|\xi\|_{G^{+}}$$

$$+ \sqrt{\eta} \|\sigma(x) - \sigma(y)\|_{op,G^{+}} \|Z' - (1 \wedge \bar{z}_{1} \|\xi_{1}\|_{G^{+}}^{-1})\xi_{1} - (1 \wedge \bar{z}_{2} \|\xi_{2}\|^{-1})\xi_{2} \|_{G^{+}}$$

$$\leq R' + \eta \tilde{B} + \sqrt{\eta} \tilde{B}_{\sigma} \|\xi\|_{G^{+}} + \sqrt{\eta} L_{\sigma} r(\|\xi_{1}\|_{G^{+}} \mathbb{1}_{\|\xi_{1}\|_{G^{+}}} \geq \bar{z}_{1} + \|\xi_{2}\|_{G^{+}} \mathbb{1}_{\|\xi_{2}\|_{G^{+}}} \geq \bar{z}_{2}).$$

We use the following stability bound for the function f given in the proof of Theorem 2.5 in Majka et al. (2020).

Lemma 33. For any $t, s \ge 0$, we have

$$f(t) - f(s) \le (r_2^{-1}e^{-ar_2}(t \lor s) + 1)|t - s|.$$

Thus, the difference between f(R'') and f(R') is given by,

$$f(R'') - f(R')$$

$$\leq f(R' + \eta \tilde{B} + \sqrt{\eta} \tilde{B}_{\sigma} \|\xi\|_{G^{+}} + \sqrt{\eta} L_{\sigma} r(\|\xi_{1}\|_{G^{+}} \mathbb{1}_{\|\xi_{1}\|_{G^{+}} \geq \bar{z}_{1}} + \|\xi_{2}\|_{G^{+}} \mathbb{1}_{\|\xi_{2}\|_{G^{+}} \geq \bar{z}_{2}})) - f(R')$$

$$\leq (r_{2}^{-1} e^{-ar_{2}} (R' + \eta \tilde{B} + \sqrt{\eta} \tilde{B}_{\sigma} \|\xi\|_{G^{+}} + \sqrt{\eta} L_{\sigma} r(\|\xi_{1}\|_{G^{+}} \mathbb{1}_{\|\xi_{1}\|_{G^{+}} \geq \bar{z}_{1}} + \|\xi_{2}\|_{G^{+}} \mathbb{1}_{\|\xi_{2}\|_{G^{+}} \geq \bar{z}_{2}}))$$

$$+ 1)(\eta \tilde{B} + \sqrt{\eta} \tilde{B}_{\sigma} \|\xi\|_{G^{+}} + \sqrt{\eta} L_{\sigma} r(\|\xi_{1}\|_{G^{+}} \mathbb{1}_{\|\xi_{1}\|_{G^{+}} \geq \bar{z}_{1}} + \|\xi_{2}\|_{G^{+}} \mathbb{1}_{\|\xi_{2}\|_{G^{+}} \geq \bar{z}_{2}})). \tag{52}$$

We now control the expected value of this. Using concentration of the χ^2 distribution (see Example 2.11 of Wainwright (2019)), we obtain that for any $\bar{z}_1 = \sqrt{2\lambda_{\rm gap}(G)^{-1}(n-1)}$,

$$\mathbb{E}[\|\xi_2\|_{G^+}^2 \mathbb{1}_{\|\xi_2\|_{G^+}} \geq \bar{z}_2]$$

$$\leq \lambda_{\mathrm{gap}}(G)^{-1} \mathbb{E}[\|\xi_2\|^2 \mathbb{1}_{\|\xi_2\| \geq \lambda_{\mathrm{gap}}(G)^{1/2} \bar{z}_2}]$$

$$\leq \lambda_{\mathrm{gap}}(G)^{-1} \int_{\lambda_{\mathrm{gap}}(G) \bar{z}_2^2}^{\infty} \mathbb{P}(\|\xi_2\|^2 \geq r) \, dr$$

$$+ \lambda_{\mathrm{gap}}(G)^{-1} \int_{0}^{\lambda_{\mathrm{gap}}(G) \bar{z}_2^2}^{\infty} \mathbb{P}(\|\xi_2\| \geq \bar{z}_2) \, dr$$

$$\leq \lambda_{\mathrm{gap}}(G)^{-1} \int_{\lambda_{\mathrm{gap}}(G) \bar{z}_2^2}^{\infty} \mathbb{P}(\|\xi_2\| \geq \bar{z}_2) \, dr$$

$$\leq \lambda_{\mathrm{gap}}(G)^{-1} \int_{\lambda_{\mathrm{gap}}(G) \bar{z}_2^2}^{\infty} \exp\left(-\frac{(r-(n-1))^2}{8n}\right) dr$$

$$+ \lambda_{\mathrm{gap}}(G)^{-1} \exp\left(-\frac{(\lambda_{\mathrm{gap}}(G) \bar{z}_2^2 - (n-1))^2}{8n}\right) \bar{z}^2$$

$$\leq \lambda_{\mathrm{gap}}(G)^{-1} \exp\left(-\frac{(\lambda_{\mathrm{gap}}(G) \bar{z}_2^2 - (n-1))^2}{8n}\right) \bar{z}^2$$

$$\leq \lambda_{\mathrm{gap}}(G)^{-1} (\sqrt{8(n-1)\pi} + \lambda_{\mathrm{gap}}(G) \bar{z}_2^2) \exp\left(-\frac{(\lambda_{\mathrm{gap}}(G) \bar{z}_2^2 - (n-1))^2}{8(n-1)}\right)$$

$$\leq \lambda_{\mathrm{gap}}(G)^{-1} \left(\sqrt{8(n-1)\pi} \exp(-(n-1)/16)\right)$$

$$\leq \lambda_{\mathrm{gap}}(G)^{-2} \exp\left(-\frac{\bar{z}_2^4}{64}\right) \exp\left(-\frac{(\lambda_{\mathrm{gap}}(G) \bar{z}_2^2 - (n-1))^2}{16(n-1)}\right)$$

$$\leq \kappa_0 \lambda_{\mathrm{gap}}(G)^{-1} \exp\left(-\frac{(\lambda_{\mathrm{gap}}(G) \bar{z}_2^2 - (n-1))^2}{16(n-1)}\right)$$

for some universal constant $\kappa_0 \geq 1$ (independent of n and \bar{z}). Similarly, we have

$$\mathbb{E}[\|\xi_1\|_{G^+}^2 \mathbb{1}_{\|\xi_1\|_{G^+} \ge \bar{z}_1}] \le \kappa_0 \lambda_{\text{gap}}(G)^{-1} \exp\left(-\frac{(\lambda_{\text{gap}}(G)\bar{z}_1^2 - 1)^2}{16}\right),$$

for any $\bar{z}_1 \geq \lambda_{\text{gap}}(G)^{-1/2}$. Therefore, we choose $\bar{z}_1 = \frac{\lambda}{4}L_{\sigma}^{-1}\sqrt{\eta}$ We now return to (52) using these bounds as well as the fact that $\mathbb{E}[R'|Z'] = \hat{r}$. Defining the quantity,

$$A := \kappa_0^{1/2} \lambda_{\text{gap}}(G)^{-1/2} \exp(-(n-1)/32) + \kappa_0^{1/2} \lambda_{\text{gap}}(G)^{-1/2} \exp(-(\lambda_{\text{gap}}(G)\bar{z}_1^2 - 1)^2/32),$$

we obtain that for $\eta \leq \min\{\tilde{B}/2, d\tilde{B}_{\sigma}^2/4, 1/2L, 1/2L_{\sigma}^2A^2\}$

$$\mathbb{E}[f(R'') - f(R')]$$

$$= (r_2^{-1} e^{-ar_2}) (\mathbb{E}[\hat{r}^2]^{1/2} + \eta \tilde{B} + \sqrt{\eta d} \tilde{B}_{\sigma} + \sqrt{\eta} L_{\sigma} r A) + 1) (\eta \tilde{B} + \sqrt{\eta d} \tilde{B}_{\sigma} + \sqrt{\eta} L_{\sigma} r A)$$

$$\leq (r_2^{-1} e^{-ar_2}) (1 + \eta L + \sqrt{\eta} L_{\sigma} A) r + 1) \sqrt{\eta} L_{\sigma} r A + \frac{1}{r_2} e^{-ar_2} (\eta^2 \tilde{B}^2 + \eta d \tilde{B}_{\sigma}^2)$$

$$+ (r_2^{-1} e^{-ar_2}) (1 + \eta L + \sqrt{\eta} L_{\sigma} A) r + 1) (\eta \tilde{B} + \sqrt{\eta} d \tilde{B}_{\sigma})$$

$$+ r_2^{-1} e^{-ar_2} (\eta \tilde{B} + \sqrt{\eta} d \tilde{B}_{\sigma}) \sqrt{\eta} L_{\sigma} r^2 A$$

$$\leq (4r_2^{-1} e^{-ar_2} r + 1) \sqrt{\eta} L_{\sigma} r A + \frac{3}{2r_2} e^{-ar_2} (\eta^2 \tilde{B}^2 + \eta d \tilde{B}_{\sigma}^2).$$

When $r \leq r_2$, we have

$$(r_2^{-1}e^{-ar_2}(1+\eta L+3\sqrt{\eta}L_{\sigma}A)r+1)r$$

$$\leq (e^{-ar_2}(1+\eta L+\sqrt{\eta}L_{\sigma}A)+1)r$$

$$\leq (e^{-ar_2}(1+\eta L+\sqrt{\eta}L_{\sigma}A)+1)(a^{-1}(1-e^{-ar_2}))^{-1}a^{-1}(1-e^{-ar})$$

$$\leq 4(a^{-1}(1-e^{-ar_2}))^{-1}f(r),$$

where in the final line, we used $\eta \leq L^{-1}$ and $\sqrt{\eta} L_{\sigma} \kappa_0^{1/2} \leq 1$. When $r > r_2$, we have

$$(r_2^{-1}e^{-ar_2}(1+\eta L+\sqrt{\eta}L_{\sigma}A)r+1)r \leq (e^{-ar_2}(1+\eta L+\sqrt{\eta}L_{\sigma}A)+1)r_2^{-1}r^2 \leq 2(2+e^{ar_2})f(r).$$

Thus, we obtain,

$$\mathbb{E}[f(R'')] \leq \mathbb{E}[f(R')] + \sqrt{\eta} L_{\sigma} A_{\frac{6\vee(4a)}{1-e^{-ar_2}}} f(r) + \frac{1}{2r_2} e^{-ar_2} (\eta^2 \tilde{B}^2 + \eta d\tilde{B}_{\sigma}^2)$$

$$\leq (1 - \eta c/2 + \sqrt{\eta} L_{\sigma} A_{\frac{6\vee(4a)}{1-e^{-ar_2}}}) f(r) + \frac{3}{2r_2} e^{-ar_2} (\eta^2 \tilde{B}^2 + \eta d\tilde{B}_{\sigma}^2),$$

where we used $AL_{\sigma} \frac{6\vee(4a)}{1-e^{-ar_2}} \leq \sqrt{\eta}c/4$.

G Proofs for the stability of the noisy gradient estimator

Using the Wasserstein contraction obtained in the previous section, we will now prove Proposition 14.

Proposition 14. Consider the score matching algorithm $A_{\rm sm}: S \mapsto s_{\theta_K}$ for some fixed $K \in \mathbb{N}$ where $(\theta_k)_k$ is as given in (16). Suppose that assumptions 10, 12 and 13 hold, then there exists some $\bar{\eta} > 0$ such that, if $\sup_{\eta} \eta_p \leq \bar{\eta}$, we obtain that $A_{\rm sm}$ is score stable with constant

$$\varepsilon_{stab}^2 \lesssim \frac{\overline{L}^2 C^2 (P+n)}{\lambda_{gap} N} \min \left\{ \frac{\eta_{\min} \lambda_{gap} \lambda^2}{P N_B C} \sum_{k=0}^{K-1} \eta_k, \exp \left(\tilde{c} \frac{P N_B C}{\eta_{\min} \lambda_{gap} \lambda^2} \right) \right\},$$

where
$$\tilde{c} \lesssim (\overline{M}_4 B_\ell C_\tau^{1/2} + \overline{L}_4^2) (P N_B \lambda_{gap})^{-1/2} \vee 1$$
, $\eta_{\min} = \min_k \eta_k$.

The proof of Proposition follows from an application of Proposition 27 to the process in (16). Similar to the proof of Proposition 11, we obtain stability estimates by analysing the trajectories θ_k and $\tilde{\theta}_k$ trained on S and S^N with coupled minibatch indices. In particular, given a set of minibatch indices $B \subset [N]$ with $|B| = N_B$, if we set

$$b(\theta) := \mathbb{E}\Big[\mathrm{Clip}_C(G(\theta, (x_i)_{i \in B})) \Big| \theta, B, S\Big], \qquad \tilde{b}(\theta) := \mathbb{E}\Big[\mathrm{Clip}_C(G(\theta, (\tilde{x}_i)_{i \in B})) \Big| \theta, B, S^N\Big]$$
$$\sigma(\theta) := \sqrt{\eta} \, \Sigma_S(\theta, B)^{1/2}, \qquad \tilde{\sigma}(\theta) := \sqrt{\eta} \, \Sigma_{S^N}(\theta, B)^{1/2},$$

where we use $(\tilde{x}_i)_{i=1}^N$ to denote the dataset S^N (i.e. $\tilde{x}_i = x_i$ for all $i \neq N$ and $\tilde{x}_N = \tilde{x}$), then the trajectories θ_k and $\tilde{\theta}_k$ are updated as in (36), (37). Using the shorthand, $v_{i,j}(\theta) = w_{t_{(i,j)}} \nabla_{\theta} \|s_{\theta}(X_{(i,j)}, t_{(i,j)}) - \nabla \log p_{t_{(i,j)}}|_{\theta}(X_{(i,j)}|_{x_i})\|^2$, we obtain the bound,

$$\Sigma_{S}(\theta, B) \succcurlyeq \operatorname{Cov}\left(\frac{1}{PN_{B}} \sum_{i \in B} \sum_{j=1}^{P} \operatorname{Clip}_{C}(v_{i,j}(\theta)) \middle| \theta, B, S\right)$$

$$\succcurlyeq \frac{1}{P} \frac{1}{N_{B}^{2}} \sum_{i \in B} \operatorname{Cov}\left(\operatorname{Clip}_{C}(v_{i,j}(\theta)) \middle| \theta, B, S\right)$$

$$\succcurlyeq \frac{1}{PN_{B}} \bar{\Sigma}.$$

Therefore, we have $\sigma(\theta) \succcurlyeq \sqrt{\eta/PN_B} \, \bar{\Sigma}^{1/2} =: G^{1/2}$, and similarly, $\tilde{\sigma}(\theta) \succcurlyeq G^{1/2}$. The weighted norm $\|\cdot\|_{G^+}$ satisfies the property,

$$\|\theta\|_{G^+} \le \lambda_{\max}(G^+)^{1/2} \|\theta\| \le \sqrt{\frac{PN_B}{\eta \lambda_{\text{gap}}}} \|\theta\|$$

Therefore, due to the gradient clipping, we have $\|b(\theta)\|_{G^+} \le \sqrt{PN_B/\eta\lambda_{\rm gap}}C =: B$. Furthermore, by Assumption 13, we apply the same argument used in the proof of Proposition 11 to obtain

$$||b(\theta) - b(\theta')||_{G^+} \le (\overline{M}_4 B_\ell C_\tau^{1/2} + \overline{L}_4^2) ||\theta - \theta'||_{G^+}$$

so $L_b=\overline{M}_4B_\ell C_ au^{1/2}+\overline{L}_4^2$. To obtain the Lipschitz constant for the volatility matrix, we first obtain,

$$\sigma(\theta) - \sigma(\theta') \preceq \sqrt{\eta} \operatorname{Cov} \left(\frac{1}{PN_B} \sum_{i \in B} \sum_{j=1}^{P} \left((1 \vee (C \| v_{i,j}(\theta) \|^{-1})) v_{i,j}(\theta) - (1 \vee (C \| v_{i,j}(\theta') \|^{-1})) v_{i,j}(\theta') \right) \middle| \theta, B, S \right)^{1/2}.$$

From this, we deduce,

$$\|\sigma(\theta) - \sigma(\theta')\|_{op,G^{+}} \leq \sqrt{\eta} \sup_{\|v\|_{G^{+}} = 1} \operatorname{Var}\left(\left\langle G^{+}v, \frac{1}{PN_{B}} \sum_{i \in B} \sum_{j=1}^{P} \left((1 \vee (C\|v_{i,j}(\theta)\|^{-1}))v_{i,j}(\theta) - (1 \vee (C\|v_{i,j}(\theta')\|^{-1}))v_{i,j}(\theta') \right) \right) \Big| \theta, B, S \Big)^{1/2}$$

$$\leq \sqrt{\eta} \left(\frac{1}{PN_{B}^{2}} \sum_{i \in B} \operatorname{Var}\left(\|v_{i,j}(\theta) - v_{i,j}(\theta')\|_{G^{+}} \Big| \theta, B, S \right) \right)^{1/2}.$$

To control this further, we use the Lipschitz assumption on to show that v is Lipschitz also:

$$\begin{aligned} \|v_{i,j}(\theta) - v_{i,j}(\theta')\|_{G^{+}} \\ &\leq 2\|s_{\theta}(X_{(i,j)}, t_{(i,j)}) - s_{\theta'}(X_{(i,j)}, t_{(i,j)})\|_{G^{+}} \|\nabla_{\theta}s_{\theta}(X_{(i,j)}, t_{(i,j)})\|_{op,G^{+}} \\ &\quad + 2\|s_{\theta}(X_{(i,j)}, t_{(i,j)}) - \nabla \log p_{t_{(i,j)}|0}(X_{(i,j)}|x_{i})\|_{G^{+}} \|\nabla_{\theta}s_{\theta}(X_{(i,j)}, t_{(i,j)}) \\ &\quad - \nabla_{\theta}s_{\theta'}(X_{(i,j)}, t_{(i,j)})\|_{op,G^{+}} \\ &\leq 2L(X_{(i,j)}, t_{(i,j)})^{2} \|\theta - \theta'\|_{G^{+}} + \frac{2c^{1/2}\mu_{t}}{\sigma_{s}^{2}} M(X_{(i,j)}, t_{(i,j)}) \|\theta - \theta'\|_{G^{+}}. \end{aligned}$$

Computing the variance of this leads to the bound,

$$\|\sigma(\theta) - \sigma(\theta')\|_{op,G^{+}} \leq 2\sqrt{\frac{\eta}{PN_{B}\lambda_{gap}}} (\overline{M}_{4}B_{\ell}C_{\tau}^{1/2} + \overline{L}_{4}^{2}) \|\theta - \theta'\|_{G^{+}} =: L_{\sigma}\|\theta - \theta'\|_{G^{+}}.$$

Next, we use a similar argument to the proof of Proposition 11 to obtain

$$||b(\theta) - \tilde{b}(\theta)||_{G^+} \le \sqrt{\frac{PN_B}{\eta \lambda_{\text{gap}}}} ||b(\theta) - \tilde{b}(\theta)|| \le \sqrt{\frac{PN_B}{\eta \lambda_{\text{gap}}}} \frac{2C}{N_B} \mathbb{1}_{N \in B} =: \tilde{B}_b.$$

Therefore we have satisfied all assumptions of Proposition 27 aside from Assumption 26. To satisfy this assumption we use that $L_{\sigma} \sim \sqrt{\eta/P}$ and so if η is sufficiently small, or P is sufficiently large, this assumption is satisfied once n is sufficiently large also.

Using Proposition 27, we obtain the contraction,

$$\mathbb{E}[d(\theta_{k+1}, \tilde{\theta}_{k+1}) | \theta_k, \tilde{\theta}_k, B_k] \le (1 - \eta c/2) d(\theta_k, \tilde{\theta}_k) + \frac{3}{2r_2} e^{-ar_2} \Big(\eta \frac{4PC^2}{\lambda_{\text{gap}} N_B} \mathbb{1}_{N \in B} + \frac{\eta n}{N_B \lambda_{\text{gap}}} C^2 \mathbb{1}_{N \in B} \Big).$$

Using the fact that $\mathbb{P}(N \in B_k) = N_B/N$, we obtain,

$$\mathbb{E}[d(\theta_{k+1},\tilde{\theta}_{k+1})] \leq (1 - \eta c/2) \mathbb{E}[d(\theta_k,\tilde{\theta}_k)] + \eta \frac{3}{2r_2} e^{-ar_2} \bigg(\frac{4PC^2}{\lambda_{\text{gap}}} + \frac{n}{\lambda_{\text{gap}}} C^2 \bigg) \frac{1}{N}.$$

Thus, by comparison, we obtain the bound,

$$\mathbb{E}[d(\theta_K, \tilde{\theta}_K)] \leq \frac{3}{2r_2} e^{-ar_2} \left(\frac{4PC^2}{\lambda_{\text{gap}}} + \frac{n}{\lambda_{\text{gap}}} C^2 \right) \frac{1}{N} \eta \sum_{k=0}^{K-1} (1 - \eta c/2)^k$$

$$= \frac{3}{2r_2} e^{-ar_2} \left(\frac{4PC^2}{\lambda_{\text{gap}}} + \frac{n}{\lambda_{\text{gap}}} C^2 \right) \frac{1 - (1 - \eta c/2)^K}{Nc/2}$$

$$\leq \frac{3}{2r_2} e^{-ar_2} (4P + n) \frac{C^2}{\lambda_{\text{gap}} N} (\eta K \wedge 2/c).$$

By the definition of f(r), we have that it dominates r^2 up to a multiplicative constant:

$$\begin{split} f(r) &\geq \left(\left(\frac{1}{a} (1 - e^{-ar_2}) \right) \wedge \left(\frac{1}{2r_2} e^{-ar_2} \right) \right) r^2 \\ &\geq \frac{1}{2r_2} e^{-ar_2} \left(\left(\frac{2r_2}{a} (e^{ar_2} - 1) \right) \wedge 1 \right) r^2 \\ &\geq \frac{1}{2r_2} e^{-ar_2} ((2r_2^2) \wedge 1) r^2. \end{split}$$

Thus, using assumption 13, it follows that

$$\int \mathbb{E}[\|s_{\theta_{K}}(X_{t},t) - s_{\tilde{\theta}_{K}}(X_{t},t)\|^{2}|X_{0} = \tilde{x},S]\tau(dt)
\leq \bar{L}^{2}\mathbb{E}\Big[\|\theta_{K} - \tilde{\theta}_{K}\|_{G^{+}}^{2}\Big]
\leq \bar{L}^{2}\Big(\frac{1}{2r_{2}}e^{-ar_{2}}((2r_{2}^{2})\wedge 1)\Big)^{-1}\mathbb{E}\Big[d(\theta_{K},\tilde{\theta}_{K})\Big]
\leq 3\bar{L}^{2}((2r_{2}^{2})^{-1}\vee 1)(4P+n)\frac{C^{2}}{\lambda_{\text{gap}}N}(\eta K \wedge 2/c).$$

We then use the fact that when η is sufficiently small, we obtain the estimate $\eta_0 \gtrsim \lambda^{-1}$ and therefore,

$$r_1^2, r_2^2 \gtrsim \frac{PN_BC}{\eta\lambda_{\rm gap}\lambda^2}, \qquad L \lesssim (\overline{M}_4B_\ell C_\tau^{1/2} + \overline{L}_4^2)(PN_B\lambda_{\rm gap})^{-1/2} \vee 1$$

and since L and r_1 explode as $\eta \to 0^+$, we also have,

$$r_2^2 c \gtrsim L^2 r_1^4 \exp(-6Lr_1^2/c_0)$$

 $\gtrsim \exp(-6Lr_1^2/c_0).$