

Stochastic Modified Flows, Mean-Field Limits and Dynamics of Stochastic Gradient Descent

Benjamin Gess^{*†}, Sebastian Kassing^{*}, Vitalii Konarovskyi^{*‡}

February 14, 2023

Abstract

We propose new limiting dynamics for stochastic gradient descent in the small learning rate regime called stochastic modified flows. These SDEs are driven by a cylindrical Brownian motion and improve the so-called stochastic modified equations by having regular diffusion coefficients and by matching the multi-point statistics. As a second contribution, we introduce distribution dependent stochastic modified flows which we prove to describe the fluctuating limiting dynamics of stochastic gradient descent in the small learning rate - infinite width scaling regime.

Keywords. Stochastic gradient descent, machine learning, overparametrization, stochastic modified equation, fluctuation mean field limit.

1 Introduction

Stochastic gradient descent algorithms (SGD), going back to [26], are the most common way to train neural networks. Due to the non-convexity and non-smoothness of the corresponding loss landscapes, the analysis of the optimization dynamics is highly challenging. The analysis of the implicit, algorithmic bias of SGD in overparameterized networks is one of the key open problems in the understanding of the empirically observed good generalization properties of networks trained by SGD. Since the dynamics of SGD depend

^{*}Fakultät für Mathematik, Universität Bielefeld, 33615 Bielefeld, Germany.

[†]Max Planck Institute for Mathematics in the Sciences, 04103 Leipzig, Germany

[‡]Institute of Mathematics of NAS of Ukraine, 01024 Kyiv, Ukraine

E-mails: benjamin.gess@math.uni-bielefeld.de, skassing@math.uni-bielefeld.de, vitalii.konarovskyi@math.uni-bielefeld.de

Mathematics Subject Classification (2020): Primary 60J05, 60H15, 68T07; Secondary 60G46, 60G57, 46G05

on many choices, like the choice of the loss function, the architecture of the network and the training data, their systematic understanding relies on the identification of universal structures that are invariant to these many degrees of freedoms, while retaining the essential properties of SGD. In recent years, several of such scaling limits and corresponding limiting dynamics have been identified. Among these, solutions to SDEs have been obtained as universal continuum objects in the small learning rate regime [10, 19], while (stochastic) Wasserstein gradient flows have been found in infinite width overparameterized limits [4, 5, 14, 16, 22, 23, 24, 27, 30, 32, 33]. In the present work, we introduce a new form of stochastic limiting dynamics which solves simultaneously three challenges met in previous works: (1) The irregularity of diffusion coefficients, (2) matching multi-point statistics, and (3) incorporating overparameterized limits.

Before we comment on each of these aspects in a few more details, let us recall the principle setup of SGD in supervised learning. For a given training data set $\Theta \subseteq \mathbb{R}^{n_0}$ sampled from a probability distribution ϑ , one aims to minimize the empirical risk

$$R(z) := \mathbb{E}_{\vartheta} \tilde{R}(z, \theta), \quad z \in \mathbb{R}^d,$$

where $\tilde{R} : \mathbb{R}^d \times \Theta \rightarrow \mathbb{R}$ is a loss function. Let $\theta_n, n \in \mathbb{N}_0 (:= \mathbb{N} \cup \{0\})$, be i.i.d. samples of training data drawn from ϑ . Then, the SGD dynamics is given by

$$Z_{n+1}^\eta(x) = Z_n^\eta(x) - \eta \nabla \tilde{R}(Z_n^\eta(x), \theta_n), \quad n \in \mathbb{N}_0, \quad (1.1)$$

where $Z_0(x) = x, x \in \mathbb{R}^d$ and $\eta > 0$. In particular, $Z_n^\eta, n \in \mathbb{N}_0$, allows to analyze the training dynamics of different initializations x subject to the same choice of training data.

We next address the above mentioned challenges in a few more details.

(1) The irregularity of diffusion coefficients: In the regime of small learning rate, the foundational works of Li, Tai and E [19, 20] have suggested stochastic modified equations (SME) as universal continuum limits that capture both the average gradient descent performed by SGD and its fluctuations. More precisely, it is shown that the SGD dynamics $Z_n^\eta, n \in \mathbb{N}_0$, with learning rate η can be approximated to higher order in η by solutions to SMEs

$$dY_t^\eta(x) = -\nabla \left(R(Y_t^\eta(x)) + \frac{\eta}{4} |\nabla R(Y_t^\eta(x))|^2 \right) dt + \sqrt{\eta} \Sigma(Y_t^\eta(x))^{1/2} dW_t \quad (1.2)$$

where $Y_0^\eta(x) = x$ for $x \in \mathbb{R}^d, W_t, t \geq 0$, is a Brownian motion in \mathbb{R}^d and $\Sigma : \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$ is the matrix defined by

$$\Sigma(y) = \mathbb{E}_{\vartheta} \left[(\nabla_y \tilde{R}(y, \theta) - \nabla R(y)) \otimes (\nabla_y \tilde{R}(y, \theta) - \nabla R(y)) \right], \quad y \in \mathbb{R}^d. \quad (1.3)$$

Indeed, this incorporates a certain degree of universality of (1.2), since the noise in SGD is represented in (1.2) by Brownian noise. In machine learning, and, in particular, in overparameterized settings, the covariance matrix Σ is typically degenerate. As a result, the square root $\Sigma^{1/2}$ appearing in (1.2) has limited regularity properties³, which makes the

³The simple example $\Sigma(y) = y^2, \Sigma^{1/2}(y) = |y|$ shows that not more than Lipschitz continuity can be expected from $\Sigma^{1/2}$ in general.

analysis of (1.2) challenging, and leads to assumptions on $\Sigma^{1/2}$ that are in general not known to hold. The first contribution of this work is to resolve this issue by introducing a new model for the stochastic limiting dynamics, which we name stochastic modified flow (SMF),

$$\begin{aligned} dX_t^\eta(x) &= -\nabla \left(R(X_t^\eta(x)) + \frac{\eta}{4} |\nabla R(X_t^\eta(x))|^2 \right) dt + \sqrt{\eta} \int_{\Theta} G(X_t^\eta(x), \theta) W(d\theta, dt), \\ X_0^\eta(x) &= x, \quad x \in \mathbb{R}^d, \end{aligned} \tag{1.4}$$

where $G(x, \theta) = \nabla \tilde{R}(x, \theta) - \nabla R(x)$ and W is a cylindrical Wiener process on the space $L_2((\Theta, \vartheta); \mathbb{R})$. It is important to notice that (1.4) satisfies the same martingale problem as (1.2), while avoiding the appearance of $\Sigma^{\frac{1}{2}}$, thereby bypassing the resulting irregularity of the diffusion coefficients. In contrast, only regularity assumptions on the individual losses \tilde{R} are needed. More precisely, we get the following result.

Theorem 1.1 (see Theorem 3.3 and Corollary 3.5). *Let $\tilde{R}(\cdot, \theta)$ be regular enough for ϑ -a.e. $\theta \in \Theta$ and let $T > 0$. Then for every $f \in C_b^4(\mathbb{R}^d)$, one has*

$$\sup_{x \in \mathbb{R}^d} \sup_{n: n\eta \leq T} \left| \mathbb{E}f(X_{n\eta}^\eta(x)) - \mathbb{E}f(Z_n^\eta(x)) \right| \lesssim \eta^2.$$

(2) Matching multi-point statistics: In a variety of works, a dynamical systems approach to the dynamics of SGD has been introduced [31, 35]. This aims at using the concepts of attractors, Lyapunov exponents, stochastic synchronization etc. in the analysis of SGD dynamics, for example, in order to analyze asymptotic global stability, that is, if

$$|Z_n^\eta(x) - Z_n^\eta(y)| \rightarrow 0 \quad \text{for } n \rightarrow \infty \tag{1.5}$$

in probability. As before, the systematic analysis of such dynamical behavior of SGD relies on the identification of appropriate universal limiting models. It is thus tempting to analyze the dynamical features of SGD by means of those of (1.2). However, this is not correct, since (1.2) only captures the single-point motion of SGD, while dynamical features like stability (1.5) are properties of the multi-point motions. More precisely, (1.2) captures the limiting behavior of the law of single motions $\text{Law}(Z_n^\eta(x))$, but not the joint multi-point laws $\text{Law}(Z_n^\eta(x_1), \dots, Z_n^\eta(x_m))$ (see also Example 3.7). As a second main contribution, in this work we prove that (SMF), in contrast to (1.2), captures the correct multi-point distributions of SGD, and therefore opens the way for an analysis of the dynamical properties of its (stochastic) flow.

Theorem 1.2 (see Theorem 3.3 and Corollary 3.6). *Under the assumption of Theorem 1.1, for every $\Phi \in C_b^4(\mathcal{P}_2(\mathbb{R}^d))$ one has*

$$\sup_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} \sup_{n: n\eta \leq T} \left| \mathbb{E}\Phi(\mu \circ (X_{n\eta}^\eta)^{-1}) - \mathbb{E}\Phi(\mu \circ (Z_n^\eta)^{-1}) \right| \lesssim \eta^2,$$

where $\mu \circ f^{-1}$ denotes the push forward of the measure μ under a map f . Furthermore, for every $m \in \mathbb{N}$ and $f \in C_b^4(\mathbb{R}^{dm})$,

$$\sup_{x_1, \dots, x_m \in \mathbb{R}^d} \sup_{n: n\eta \leq T} \left| \mathbb{E}f(X_{n\eta}^\eta(x_1), \dots, X_{n\eta}^\eta(x_m)) - \mathbb{E}f(Z_n^\eta(x_1), \dots, Z_n^\eta(x_m)) \right| \lesssim \eta^2.$$

(3) Overparameterized limits: As a third main contribution, we extend the small learning rate limit to also incorporate the infinite width limit. We here consider networks with quadratic loss function. Let $\mathcal{D} \subseteq \mathbb{R}^{n_0} \times \mathbb{R}^{k_0}$ be a given training data set⁴ with inputs $\Theta = \{\theta : (\theta, f(\theta)) \in \mathcal{D}\}$ and labels $\{f(\theta) : (\theta, f(\theta)) \in \mathcal{D}\}$. For the approximation of f we choose a parameterized hypotheses space $\mathcal{M} := \{f^M(z, \cdot) : z \in \mathbb{R}^{Md}\}$, $M, d \in \mathbb{N}$, where

$$f^M(z, \theta) = \frac{1}{M} \sum_{i=1}^M \Psi(z^i, \theta), \quad \theta \in \Theta, \quad (1.6)$$

with $\Psi : \mathbb{R}^d \times \Theta \rightarrow \mathbb{R}^{k_0}$, $z = (z^i)_{i \in [M]}$ and $[M] := \{1, \dots, M\}$. For example, one can choose \mathcal{M} to be the space of response functions of fully connected feed-forward neural networks with one hidden layer containing M hidden neurons. In that case, we choose a function $\phi : \mathbb{R} \rightarrow \mathbb{R}$, the activation function, and we write $z = (z^i)_{i \in [M]}$ with $z^i = (c^i, U^i, b^i) \in \mathbb{R}^{k_0} \times \mathbb{R}^{n_0} \times \mathbb{R}$ and $\Psi(z^i, \theta) = c^i \phi(U^i \cdot \theta + b^i)$. Then,

$$f^M(z, \theta) = \frac{1}{M} \sum_{i=1}^M c^i \phi(U^i \cdot \theta + b^i), \quad \theta \in \Theta.$$

The aim of risk minimization (with respect to the square loss) is to select a suitable model $f^M(z, \cdot)$ minimizing the risk $R(z) = \mathbb{E}_\vartheta \tilde{R}(z, \theta)$, $z \in \mathbb{R}^{Md}$, for

$$\tilde{R}(z, \theta) = \frac{1}{2} |f(\theta) - f^M(z, \theta)|^2, \quad z \in \mathbb{R}^{Md}, \theta \in \Theta.$$

As before, this optimization task is executed by the stochastic gradient descent algorithm (1.1) with the starting value $Z_0^\eta = (Z_0^{i,\eta})_{i \in [M]}$ being a tuple of i.i.d. random variables with distribution $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ that are independent of θ_n , $n \in \mathbb{N}_0$.

A simple computation gives that

$$R(z) = C_f - \frac{1}{M} \sum_{i=1}^M F(z^i) + \frac{1}{2M^2} \sum_{i,j=1}^M K(z^i, z^j),$$

where $C_f = \frac{1}{2} \mathbb{E}_\vartheta |f(\theta)|^2$ and

$$F(z^i) = \mathbb{E}_\vartheta [f(\theta) \cdot \Psi(z^i, \theta)], \quad K(z^i, z^j) = \mathbb{E}_\vartheta [\Psi(z^i, \theta) \cdot \Psi(z^j, \theta)]. \quad (1.7)$$

⁴For simplicity we assume that the ground-truth is given by a function $f : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{k_0}$.

Taking

$$\begin{aligned}
V(\nu, z^i) &= \nabla F(z^i) - \int_{\mathbb{R}^d} \nabla_{z^i} K(z^i, y) \nu(dy), \\
G(\nu, z^i, \theta) &= \left(f(\theta) - \int_{\mathbb{R}^d} \Psi(y, \theta) \nu(dy) \right) \nabla_{z^i} \Psi(z^i, \theta) \\
&\quad - \mathbb{E}_{\vartheta} \left[\left(f(\theta) - \int_{\mathbb{R}^d} \Psi(y, \theta) \nu(dy) \right) \nabla_{z^i} \Psi(z^i, \theta) \right]
\end{aligned} \tag{1.8}$$

and replacing η in (1.1) by $M\eta$, we can rewrite the expression for the dynamics of $Z_n^\eta = (Z_n^{i,\eta})_{i \in [M]}$, $n \in \mathbb{N}_0$, as follows

$$\begin{aligned}
Z_{n+1}^{i,\eta} &= Z_n^{i,\eta} + \eta V(\Gamma_n^{M,\eta}, Z_n^{i,\eta}) + \eta G(\Gamma_n^{M,\eta}, Z_n^{i,\eta}, \theta_n), \\
\Gamma_n^{M,\eta} &= \frac{1}{M} \sum_{j=1}^M \delta_{Z_n^{j,\eta}}, \quad i \in [M], \quad n \in \mathbb{N}_0,
\end{aligned} \tag{1.9}$$

where δ_z denotes the δ -measure in z .

We obtain quantified estimates on the approximation of the dynamics of the empirical measure $\Gamma_n^{M,\eta}$, $n \in \mathbb{N}_0$, of SGD by the solution to a distribution dependent stochastic modified flow (DDSMF)

$$\begin{aligned}
dX_t^\eta(x) &= \left[V(\Lambda_t^\eta, X_t^\eta(x)) - \frac{\eta}{4} \nabla |V(\Lambda_t^\eta, X_t^\eta(x))|^2 - \frac{\eta}{4} \langle D|V(\Lambda_t^\eta, X_t^\eta(x))|^2, \Lambda_t^\eta \rangle \right] dt \\
&\quad + \sqrt{\eta} \int_{\Theta} G(\Lambda_t^\eta, X_t^\eta(x), \theta) W(d\theta, dt), \\
X_0^\eta(x) &= x, \quad \Lambda_t^\eta = \mu \circ (X_t^\eta)^{-1}, \quad x \in \mathbb{R}^d, \quad t \geq 0,
\end{aligned} \tag{1.10}$$

where D denotes the differentiation with respect to the measure dependent argument in the sense of Lions⁵, $\langle \varphi, \nu \rangle$ denote the integration of a function $\varphi : \mathbb{R}^d \rightarrow \mathbb{R}$ with respect to a measure ν and W is a cylindrical Wiener process on $L_2((\Theta, \vartheta); \mathbb{R})$. We remark that

$$\frac{1}{2} \langle D|V(\nu, z)|^2, \nu \rangle = V(\nu, z) \langle \nabla_x \nabla_z K(z, x), \nu(dx) \rangle,$$

according to the form of V in (1.8) and properties of Lions derivative.

Theorem 1.3 (see Theorem 3.3, Corollary 3.8 and Remark 3.9). *Let Ψ be regular enough and $T > 0$. Then for every $\Phi \in C_b^4(\mathcal{P}_2(\mathbb{R}^d))$ and $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ with a finite p moment with $p > 2$, one has*

$$\sup_{n:\eta \leq T} |\mathbb{E}\Phi(\Lambda_{n\eta}^\eta) - \mathbb{E}\Phi(\Gamma_n^{M,\eta})| \lesssim \eta^2$$

for every $\eta > 0$ and M large enough.

⁵For more details see Section 2.2

This extends the framework of SMEs and SMFs to (1.10) which can thus serve as the starting point to analyze the stochastic dynamics of SGD in large, shallow networks.

Overview of the literature. Stochastic modified equations as limiting objects of SGD in the regime of small learning rates have been introduced by Li, Tai and E in [19, 20]. Following these original papers several results were derived for diffusion approximations with SMEs, e.g. generator based proofs [11, 15] and uniform-in-time estimates for strongly convex objective functions [18]. For a discussion on the validity of the diffusion approximation for finite (non-infinitesimal) learning rate see [21].

In [4, 16, 22, 28, 33], the convergence of gradient descent dynamics for overparameterized neural networks to a Wasserstein gradient flow has been analyzed. The conservative SPDE describing the mean-field limit that incorporates the fluctuations of the stochastic gradient descent was suggested by Rotskoff and Vanden-Eijnden in [29, 30]. The rigorous study of the well-posedness of this conservative SPDE and proof of quantified central limit theorem has been done in [14], using the observation that its solutions can be described by the SDE with interaction (2.1) below, which was investigated e.g., in [1, 7, 9, 25, 34] (see also [3, 6, 17, 34] for its connection with McKean–Vlasov SDEs with common noise). It should be noted that the stochastic modified flows proposed in this work are of a particular form of the SDE with interaction (2.1). In [28, 32], a linear SPDE has been rigorously identified in the context of central limit fluctuations of stochastic gradient descent in the overparameterized regime.

The paper is organized as follows: In Section 2, we introduce a stochastic differential equation with interaction (see (2.1)) that covers both the SMF (1.4) and the DDSMF (1.10) and recall existence and uniqueness results assuming Lipschitz-continuity of its coefficients. Moreover, we state a result for the continuous dependence of solutions to the SDE with interaction with respect to its initial distribution, as well as an analog of Kolmogorov’s equation in the setting of SDEs with interaction. Section 3 is devoted to the main result of this article, Theorem 3.3, which compares the dynamics of a discrete time Markov chain with those of a solution to a corresponding SDE with interaction. Theorem 1.1, Theorem 1.2 and Theorem 1.3 then follow as consequences of Theorem 3.3, see Corollary 3.5, Corollary 3.6 and Corollary 3.8, respectively.

2 Measure-valued diffusion and stochastic modified flows

The goal of this section is to prove the well-posedness for stochastic modified flows and investigate some properties of the associated semigroup. We recall that $\mathcal{P}_2(\mathbb{R}^d)$ denotes the space of probability measures μ on \mathbb{R}^d such that

$$\int_{\mathbb{R}^d} |x|^2 \mu(dx) < \infty$$

with the Wasserstein distance defined by

$$\mathcal{W}_2(\mu, \nu) = \inf_{\chi \in \Pi(\mu, \nu)} \left(\int_{\mathbb{R}^d} \int_{\mathbb{R}^d} |x - y|^2 \chi(dx, dy) \right),$$

where $\Pi(\mu, \nu)$ is the set of all probability measures on $\mathbb{R}^d \times \mathbb{R}^d$ with marginals μ and ν . It is well known that $\mathcal{P}_2(\mathbb{R}^d)$ equipped with the Wasserstein distance \mathcal{W}_2 is a Polish space.

Let $L_2((E, \nu); \mathbb{R}^k)$ be the space of all 2-integrable functions from a measure space (E, \mathcal{E}, ν) to \mathbb{R}^k with the usual inner product $\langle \cdot, \cdot \rangle_\nu$ and the associated norm $\| \cdot \|_\nu$. We will further fix a measure space $(\Theta, \mathcal{G}, \vartheta)$ such that ϑ is a finite measure and the space $L_2((\Theta, \vartheta); \mathbb{R})$ is separable. We will also consider a cylindrical Wiener process $W_t, t \geq 0$, on $L_2((\Theta, \vartheta); \mathbb{R})$ defined on a filtered complete probability space $(\Omega, \mathcal{F}, (\mathcal{F}_t)_{t \geq 0}, \mathbb{P})$, that is,

- (i) for every $t \geq 0$, the map $W_t : L_2((\Theta, \vartheta); \mathbb{R}) \rightarrow L_2((\Omega, \mathbb{P}); \mathbb{R})$ is linear;
- (ii) for every $h \in L_2((\Theta, \vartheta); \mathbb{R})$, $W_t(h), t \geq 0$, is an $(\mathcal{F}_t)_{t \geq 0}$ -Brownian motion with $\text{Var } W_t(h) = \|h\|_\vartheta^2 t$.

We will assume that $(\mathcal{F}_t)_{t \geq 0}$ is the complete right-continuous filtration generated by $W_t, t \geq 0$. For an $(\mathcal{F}_t)_{t \geq 0}$ -progressively measurable $L_2((\Theta, \vartheta); \mathbb{R}^k)$ -valued process $g(t, \cdot) = \{g(t, \theta), \theta \in \Theta\}, t \geq 0$, with

$$\int_0^t \|g(s, \cdot)\|_\vartheta^2 ds < \infty$$

a.s. for every $t \geq 0$, we will write⁶

$$\int_0^t \int_\Theta g(s, \theta) W(d\theta, ds) := \int_0^t \Upsilon(s) dW_s$$

for $\Upsilon(s)h = \langle g(s, \cdot), h \rangle_\vartheta = (\langle g_i(s, \cdot), h \rangle_\vartheta)_{i \in [k]}, h \in L_2((\Theta, \vartheta); \mathbb{R})$.

2.1 Stochastic modified flows

For measurable functions $B : [0, \infty) \times \mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d \rightarrow \mathbb{R}^d$, $G : [0, \infty) \times \mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d \rightarrow L_2((\Theta, \vartheta); \mathbb{R}^d)$ and a probability measure $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, we consider the following stochastic differential equation

$$\begin{aligned} dX_t(x) &= B(t, \Lambda_t, X_t(x))dt + \int_\Theta G(t, \Lambda_t, X_t(x), \theta)W(d\theta, dt), \\ X_0(x) &= x, \quad \Lambda_t = \mu \circ X_t^{-1}, \quad x \in \mathbb{R}^d, \quad t \geq 0. \end{aligned} \tag{2.1}$$

It is clear that the equations (1.4) and (1.10) can be written in the form of (2.1). Therefore, in this section we will only focus on (2.1) which is called the stochastic differential equation with interaction and was studied, e.g. in [8, 25, 34]. Let $\mathcal{B}(E)$ denote the Borel σ -algebra on a topological space E . Following the definition from [7, Definition 2.1.1] or [14, Definition 2.5], we introduce the notion of a solution to (2.1).

⁶For the definition of the integral with respect to a cylindrical Wiener process see, e.g., [13, Section 2.2.4]

Definition 2.1. A family of continuous processes $\{X_t(x), t \geq 0\}$, $x \in \mathbb{R}^d$, is called a (*strong*) *solution* to the SDE with interaction (2.1) with initial mass distribution $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ if, for each $t \geq 0$ the restriction of X to the time interval $[0, t]$ is $\mathcal{B}([0, t]) \otimes \mathcal{B}(\mathbb{R}^d) \otimes \mathcal{F}_t$ -measurable, $\Lambda_t = \mu \circ X_t^{-1}$, $t \geq 0$, is a continuous process in $\mathcal{P}_2(\mathbb{R}^d)$ and for every $x \in \mathbb{R}^d$, a.s.,

$$X_t(x) = x + \int_0^t B(s, \Lambda_s, X_s(x)) ds + \int_0^t \int_{\Theta} G(s, \Lambda_s, X_s(x), \theta) W(d\theta, ds),$$

for all $t \geq 0$. For convenience, we will also call the measure-valued process Λ_t , $t \geq 0$, a solution to (2.1).

Let $\phi_p(x) = |x|^p$, $x \in \mathbb{R}^d$. The following theorem was proved in [14]. See Theorem 2.9 and Corollary 2.10 for the well-posedness and the estimates; the existence of a continuous modification of X was observed in the proof of Theorem 2.9 *ibid*.

Theorem 2.2. *Assume that the coefficients B, G of (2.1) are Lipschitz continuous with respect to $(\mu, x) \in \mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d$, that is, for every $T > 0$ there exists $L > 0$ such that for each $t \in [0, T]$, $\mu, \nu \in \mathcal{P}_2(\mathbb{R}^d)$ and $x, y \in \mathbb{R}^d$*

$$\begin{aligned} |B(t, \mu, x) - B(t, \nu, y)| + \|G(t, \mu, x, \cdot) - G(t, \nu, y, \cdot)\|_{\vartheta} \\ \leq L (\mathcal{W}_2(\mu, \nu) + |x - y|), \end{aligned} \quad (2.2)$$

and

$$|B(t, \delta_0, 0)| + \|G(t, \delta_0, 0, \cdot)\|_{\vartheta} \leq L, \quad (2.3)$$

where δ_0 denotes the δ -measure at 0 on \mathbb{R}^d . Then, for every $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, there exists a unique strong solution $X_t(x)$, $t \geq 0$, $x \in \mathbb{R}^d$, to the SDE with interaction (2.1). Moreover, there exists a version of $X_t(x)$, $x \in \mathbb{R}^d$, that is a continuous in (t, x) , and for every $T > 0$ and $p \geq 2$ there exists a constant $C > 0$ such that

$$\mathbb{E} \sup_{t \in [0, T]} |X_t(x)|^p \leq C(1 + \langle \phi_p, \mu \rangle + |x|^p),$$

for all $x \in \mathbb{R}^d$. In particular,

$$\mathbb{E} \sup_{t \in [0, T]} \langle \phi_p, \Lambda_t \rangle \leq C(1 + \langle \phi_p, \mu \rangle),$$

where $\Lambda_t = \mu \circ X_t^{-1}$.

From now on, we will only consider the version $X_t(x)$, $t \geq 0$, $x \in \mathbb{R}^d$, of a solution to the SDE with interaction (2.1) which is continuous in (x, t) . In order to reflect the dependency on the initial mass distribution we will write $X_t(\mu, x)$ and $\Lambda_t(\mu)$ instead of $X_t(x)$ and Λ_t . We next recall the result on the continuous dependence of $\Lambda_t(\mu)$, $t \geq 0$, with respect to the initial condition μ , that was obtained in [14, Theorem 2.14].

Proposition 2.3. *Under the assumption of Theorem 2.2, for every $T > 0$ there exists a constant $C > 0$ depending only on T and the Lipschitz constant L such that*

$$\mathbb{E} \sup_{t \in [0, T]} |X_t(\mu, x) - X_t(\nu, y)|^2 \leq C (\mathcal{W}_2^2(\mu, \nu) + |x - y|^2)$$

and

$$\mathbb{E} \sup_{t \in [0, T]} \mathcal{W}_2^2(\Lambda_t(\mu), \Lambda_t(\nu)) \leq C \mathcal{W}_2^2(\mu, \nu)$$

for all $\mu, \nu \in \mathcal{P}_2(\mathbb{R}^d)$ and $x, y \in \mathbb{R}^d$.

2.2 Measure-valued diffusion

The goal of this section is to obtain an analog of Kolmogorov's equation for the process $\Lambda_t(\mu)$, $t \geq 0$, given in (2.1). For this purpose, we need to recall the notion of Lions derivative according to [2]. We say that a function $f : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}^k$ is L -differentiable at μ , if there exists an element $Df(\mu)$ in $L_2((\mathbb{R}^d, \mu); \mathbb{R}^k \times \mathbb{R}^d)$ such that

$$\lim_{\|h\|_\mu \rightarrow 0} \frac{f(\mu \circ (\text{id} + h)^{-1}) - f(\mu) - \langle Df(\mu), h \rangle_\mu}{\|h\|_\mu} = 0,$$

where id denotes the identity map on \mathbb{R}^d and the limit is taken over $h \in L_2((\mathbb{R}^d, \mu); \mathbb{R}^d)$. In this case, $Df(\mu)$ is called the L -derivative of f at μ . We write $f \in C^1(\mathcal{P}_2(\mathbb{R}^d))$ if f is L -differentiable at every point $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ and, for every μ , the derivative has a μ -version $Df(\mu, x)$ such that $Df(\mu, x)$ is jointly continuous in $(\mu, x) \in \mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d$. The set of all functions $f \in C^1(\mathcal{P}_2(\mathbb{R}^d))$ such that $f(x)$ and $Df(\mu, x)$ are bounded in $(\mu, x) \in \mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d$ will be denoted by $C_b^1(\mathcal{P}_2(\mathbb{R}^d))$. Moreover, we define $C_b^{1,1}(\mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^m)$ as the set of all functions from $\mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^m$ to \mathbb{R}^k such that $f(\cdot, x) \in C^1(\mathcal{P}_2(\mathbb{R}^d))$, $f(\mu, \cdot) \in C^1(\mathbb{R}^m)$ for all $x \in \mathbb{R}^m$, $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, and $f(\mu, x)$, $\nabla f(\mu, x)$, $Df(\mu, x, y)$ are jointly continuous and bounded in $(\mu, x, y) \in \mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^m \times \mathbb{R}^d$. We iteratively define $C_b^{l,l}(\mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^m)$ as a set of all functions $f : \mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^m \rightarrow \mathbb{R}^k$ such that $f, \nabla f \in C_b^{l-1, l-1}(\mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^m)$ and $Df \in C_b^{l-1, l-1}(\mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^{m+d})$. In particular, $f \in C_b^{2,2}(\mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^m)$ provided the function f , and all its derivatives up to the second order, i.e. ∇f , Df , $\nabla^2 f$, ∇Df , $D\nabla f$, $D^2 f$, exist, are bounded and jointly continuous. If $f : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}^k$ belongs to $C_b^{l,l}(\mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^m)$, we will write $f \in C_b^l(\mathcal{P}_2(\mathbb{R}^d))$.

Similarly to $C_b^{l,l}(\mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^m)$, we define the class $\tilde{C}_b^{l,l}(\mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^m)$ as the set of continuous and bounded functions $f : \mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^m \rightarrow L_2((\Theta, \vartheta); \mathbb{R}^k)$ such that for ϑ -a.e. $\theta \in \Theta$ we have $f(\cdot, \cdot, \theta) \in C_b^{l,l}(\mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^m)$ and all its derivatives up to the l -th order are continuous and bounded as $L_2((\Theta, \vartheta); \mathbb{R}^k)$ -valued functions.

Example 2.4. If $\varphi_i \in C_b^1(\mathbb{R}^d)$, $i \in [n]$, and $h \in C^1(\mathbb{R}^n)$ then the function $f(\mu) = h(\langle \varphi_1, \mu \rangle, \dots, \langle \varphi_n, \mu \rangle)$, $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, belongs to $C_b^1(\mathcal{P}_2(\mathbb{R}^d))$ and

$$Df(\mu, x) = \sum_{i=1}^n \partial_i h(\langle \varphi_1, \mu \rangle, \dots, \langle \varphi_n, \mu \rangle) \nabla \varphi_i(x), \quad x \in \mathbb{R}^d, \quad \mu \in \mathcal{P}_2(\mathbb{R}^d).$$

For the coefficients B and G of the equation (2.1), we define the following second-order differential operator

$$\begin{aligned} \mathcal{L}_t f(\mu) &:= \frac{1}{2} \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} \tilde{A}(t, \mu, x, y) : D^2 f(\mu, x, y) \mu(dx) \mu(dy) \\ &+ \frac{1}{2} \int_{\mathbb{R}^d} A(t, \mu, x) : \nabla D f(\mu, x) \mu(dx) \\ &+ \int_{\mathbb{R}^d} B(t, \mu, x) \cdot D f(\mu, x) \mu(dx), \end{aligned} \quad (2.4)$$

for $f \in C_b^2(\mathcal{P}_2(\mathbb{R}^d))$, where

$$\begin{aligned} \tilde{A}(t, \mu, x, y) &= \mathbb{E}_\vartheta [G(t, \mu, x, \theta) \otimes G(t, \mu, y, \theta)] \\ &= (\langle G_i(t, \mu, x, \cdot), G_j(t, \mu, y, \cdot) \rangle_\vartheta)_{i,j \in [d]}, \\ A(t, \mu, x) &= \tilde{A}(t, \mu, x, x) \end{aligned}$$

and we use the notation $C : D = \sum_{i,j=1}^d c_{i,j} d_{i,j}$ for $C = (c_{i,j})_{i,j \in [d]}$, $D = (d_{i,j})_{i,j \in [d]}$ and $a \cdot b = \sum_{i=1}^d a_i b_i$ for $a = (a_i)_{i \in [d]}$, $b = (b_i)_{i \in [d]}$.

Let us introduce some additional notation. We will write that a function $f : [0, T] \times \mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^m \rightarrow \mathbb{R}^k$ belongs to $C_b^{0,1,1}([0, T] \times \mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^m)$ if, for every $t \in [0, T]$, $f(t, \cdot, \cdot) \in C_b^{1,1}(\mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^m)$ and the maps f , $\nabla f(t, \mu, x)$, $Df(t, \mu, x, y)$ are jointly continuous and bounded in (t, μ, x, y) . Iteratively, we write $f \in C_b^{0,l,l}([0, T] \times \mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^m)$ provided that $f, \nabla f \in C_b^{0,l-1,l-1}([0, T] \times \mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^m)$ and $Df \in C_b^{0,l-1,l-1}([0, T] \times \mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^{m+d})$. The set of all functions $f : [0, T] \times \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}^k$ such that $f \in C_b^{0,2,2}([0, T] \times \mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^m)$ will be denoted by $C_b^{0,2}([0, T] \times \mathcal{P}_2(\mathbb{R}^d))$.

Similarly to $C_b^{0,l,l}([0, T] \times \mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^m)$, we also define the class $\tilde{C}_b^{0,l,l}([0, T] \times \mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^m)$ as the set of continuous and bounded functions $f : [0, T] \times \mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^m \rightarrow L_2((\Theta, \vartheta); \mathbb{R}^k)$ such that for ϑ -a.e. $\theta \in \Theta$ we have $f(\cdot, \cdot, \cdot, \theta) \in C_b^{0,l,l}([0, T] \times \mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^m)$ and all corresponding derivatives are continuous and bounded as $L_2((\Theta, \vartheta); \mathbb{R}^k)$ -valued functions.

We next provide the well posedness of the Kolmogorov equation associated to (2.1). This result can be obtained as in the proof of Theorem 3.1 in [34] with slight changes, where a similar equation driven by a finite dimensional noise was considered.

Proposition 2.5 (Kolmogorov equation). *Let $T > 0$ and the coefficients B, G of (2.1) belong to $C_b^{0,2,2}([0, T] \times \mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d)$ and $\tilde{C}_b^{0,2,2}([0, T] \times \mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d)$, respectively. For $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, let $\Lambda_t(\mu)$, $t \in [0, T]$, be a solution to (2.1) with initial mass distribution $\Lambda_0(\mu) = \mu$. Then, for every $\Phi \in C_b^2(\mathcal{P}_2(\mathbb{R}^d))$, the function*

$$U(t, \mu) = \mathbb{E} \Phi(\Lambda_t(\mu)), \quad (t, \mu) \in [0, T] \times \mathcal{P}_2(\mathbb{R}^d),$$

is a unique solution to the equation

$$\begin{aligned} \partial_t U(t, \mu) &= \mathcal{L}_t U(t, \mu), \\ U(0, \mu) &= \Phi(\mu), \quad (t, \mu) \in [0, T] \times \mathcal{P}_2(\mathbb{R}^d), \end{aligned} \quad (2.5)$$

in the class $C_b^{0,2}([0, T] \times \mathcal{P}_2(\mathbb{R}^d))$ with $\partial_t U \in C([0, T] \times \mathcal{P}_2(\mathbb{R}^d))$.

If, additionally, $B \in C_b^{0,l,l}([0, T] \times \mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d)$, $G \in \tilde{C}_b^{0,l,l}([0, T] \times \mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d)$ and $\Phi \in C_b^l(\mathcal{P}_2(\mathbb{R}^d))$, for some $l > 2$, then $U \in C_b^{0,l}([0, T] \times \mathcal{P}_2(\mathbb{R}^d))$.

3 Diffusion approximation via stochastic modified flows

The goal of this section is to prove the theorems stated in the introduction. For this, we first show a general result comparing the dynamics of a Markov chain defined below with a corresponding SDE with interaction and cylindrical noise. Then, we show that the results given in the introduction immediately follow from the general comparison statement. We fix measurable functions $V : \mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ and $G : \mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d \rightarrow L_2((\Theta, \vartheta); \mathbb{R}^d)$ such that $\mathbb{E}_\vartheta G(\mu, x, \theta) = 0$ for all $(\mu, x) \in \mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d$. For $\eta > 0$ and $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, we consider a Markov chain defined by

$$\begin{aligned} Z_{n+1}^\eta(z) &= Z_n^\eta(z) + \eta V(\Gamma_n^\eta, Z_n^\eta(z)) + \eta G(\Gamma_n^\eta, Z_n^\eta(z), \theta_n), \\ Z_0^\eta(z) &= z, \quad \Gamma_n^\eta = \mu \circ (Z_n^\eta)^{-1}, \quad z \in \mathbb{R}^d, \quad n \in \mathbb{N}_0, \end{aligned} \quad (3.1)$$

where $\theta_n, n \in \mathbb{N}_0$, are i.i.d. sampled from the distribution ϑ . We remark that, e.g., the SGD dynamics in the overparameterized shallow neural network in (1.9) can be written in form of (3.1) by taking $\mu = \frac{1}{M} \sum_{i=1}^M \delta_{Z_0^i}$ and $Z_n^{i,\eta} = Z_n^\eta(Z_0^{i,\eta})$, $i \in [M]$. We will approximate $\Gamma_n^\eta, n \in \mathbb{N}_0$, by solutions to the DDSMF

$$\begin{aligned} dX_t^\eta(x) &= \left[V(\Lambda_t^\eta, X_t^\eta(x)) - \frac{\eta}{4} \nabla |V(\Lambda_t^\eta, X_t^\eta(x))|^2 - \frac{\eta}{4} \langle D|V(\Lambda_t^\eta, X_t^\eta(x))|^2, \Lambda_t^\eta \rangle \right] dt \\ &\quad + \sqrt{\eta} \int_{\Theta} G(\Lambda_t^\eta, X_t^\eta(x)) W(d\theta, dt), \\ X_0^\eta(x) &= x, \quad \Lambda_t^\eta = \mu \circ (X_t^\eta)^{-1}, \quad x \in \mathbb{R}^d, \quad t \geq 0, \end{aligned} \quad (3.2)$$

where W is a cylindrical Wiener process on $L_2((\Theta, \vartheta); \mathbb{R})$. We first prove some auxiliary statements that will imply the well-posedness of the DDSMF (3.2).

Lemma 3.1. *Let $\varepsilon > 0$ and $\xi_s, s \in [0, \varepsilon]$, be a family of square integrable random variables on \mathbb{R}^k defined on a probability space $(\Omega, \mathcal{F}, \mathbb{P})$. If*

$$\xi_0' := \lim_{s \rightarrow 0^+} \frac{\xi_s - \xi_0}{s}$$

exists in $L_2((\Omega, \mathbb{P}); \mathbb{R}^k)$, then for every $f \in C^1(\mathcal{P}_2(\mathbb{R}^k))$ one has

$$\lim_{s \rightarrow 0^+} \frac{f(\text{Law}(\xi_s)) - f(\text{Law}(\xi_0))}{s} = \mathbb{E} [Df(\text{Law}(\xi_0), \xi_0) \cdot \xi_0'] .$$

Proof. This statement was obtained in [34, Lemma 2.4]. □

Lemma 3.2. *Let the functions V and G belong to $C_b^{1,1}(\mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d)$ and $\tilde{C}_b^{1,1}(\mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d)$, respectively. Then, for every $x, y \in \mathbb{R}^d$ and $\mu, \nu \in \mathcal{P}_2(\mathbb{R}^d)$, we have*

$$\begin{aligned} |V(\mu, x) - V(\nu, y)| + \|G(\mu, x, \cdot) - G(\nu, y, \cdot)\|_{\vartheta} \\ \leq L(\mathcal{W}_2(\mu, \nu) + |x - y|), \end{aligned}$$

with

$$\begin{aligned} L = & \sup_{x, y \in \mathbb{R}^d, \mu \in \mathcal{P}_2(\mathbb{R}^d)} (|\nabla V(\mu, x)| + |DV(\mu, x, y)|) \\ & + \sup_{x, y \in \mathbb{R}^d, \mu \in \mathcal{P}_2(\mathbb{R}^d)} (\|\nabla G(\mu, x, \cdot)\|_{\vartheta} + \|DG(\mu, x, y, \cdot)\|_{\vartheta}). \end{aligned}$$

Proof. Let $x, y \in \mathbb{R}^d$ and $\mu, \nu \in \mathcal{P}_2(\mathbb{R}^d)$ be fixed. We take an arbitrary probability measure χ on $\mathbb{R}^d \times \mathbb{R}^d$ with marginals μ, ν and consider random variables ζ_0, ζ_1 on the probability space $(\mathbb{R}^d \times \mathbb{R}^d, \mathcal{B}(\mathbb{R}^d \times \mathbb{R}^d), \chi)$ defined by $\zeta_0(x, y) = y$ and $\zeta_1(x, y) = x$ for all $(x, y) \in \mathbb{R}^d \times \mathbb{R}^d$. Then $\text{Law}(\zeta_0) = \nu$ and $\text{Law}(\zeta_1) = \mu$. Set $\xi_s = (1 - s)\zeta_0 + s\zeta_1$, $s \in [0, 1]$, and note that $\xi_i = \zeta_i$, for $i \in \{0, 1\}$, and $\xi'_s = (\zeta_1 - \zeta_0)$, for all $s \in [0, 1]$. We have

$$|V(\mu, x) - V(\nu, y)| \leq |V(\mu, x) - V(\mu, y)| + |V(\mu, y) - V(\nu, y)|$$

and we can bound the terms on the right hand side of the inequality as follows. With the mean-value theorem, the first term can be bounded by $\sup_{z \in \mathbb{R}^d, \rho \in \mathcal{P}_2(\mathbb{R}^d)} |\nabla V(\rho, z)| |x - y|$. To bound the second term, we will use Lemma 3.1 and the mean-value theorem:

$$\begin{aligned} |V(\mu, y) - V(\nu, y)| &= |V(\text{Law}(\zeta_1), y) - V(\text{Law}(\zeta_0), y)| \\ &\leq \sup_{s \in [0, 1]} \left| \frac{d}{ds} V(\text{Law}(\xi_s), y) \right| \\ &\leq \sup_{s \in [0, 1]} \left| \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} DV(\text{Law}(\xi_s), y, \xi_s(z_1, z_2)) \cdot \xi'_s(z_1, z_2) \chi(dz_1, dz_2) \right| \\ &\leq \sup_{z_1, z_2 \in \mathbb{R}^d, \rho \in \mathcal{P}_2(\mathbb{R}^d)} |DV(\rho, z_1, z_2)| \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} |\zeta_1(z_1, z_2) - \zeta_0(z_1, z_2)| \chi(dz_1, dz_2) \\ &\leq \sup_{z_1, z_2 \in \mathbb{R}^d, \rho \in \mathcal{P}_2(\mathbb{R}^d)} |DV(\rho, z_1, z_2)| \left(\int_{\mathbb{R}^d} \int_{\mathbb{R}^d} |z_2 - z_1|^2 \chi(dz_1, dz_2) \right)^{\frac{1}{2}}. \end{aligned}$$

Taking the infimum over all probability measures χ on $\mathbb{R}^d \times \mathbb{R}^d$ with marginals μ, ν , we obtain

$$|V(\mu, y) - V(\nu, y)| \leq \sup_{z_1, z_2 \in \mathbb{R}^d, \rho \in \mathcal{P}_2(\mathbb{R}^d)} |DV(\rho, z_1, z_2)| \mathcal{W}_2(\mu, \nu).$$

The estimate for $\|G(x, \mu) - G(y, \nu)\|_{\vartheta}$ can be obtained similarly. \square

Now we are ready to proof the main result of this work.

Theorem 3.3. *Let $V \in C_b^{5,5}(\mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d)$, $G \in \tilde{C}_b^{4,4}(\mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d)$ and $\mathbb{E}_\vartheta G(\mu, x, \theta) = 0$ for all $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, $x \in \mathbb{R}^d$. For $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ and $\eta > 0$, let $\Gamma_n^\eta(\mu)$, $n \in \mathbb{N}_0$, and $\Lambda_t^\eta(\mu)$, $t \geq 0$, be defined by (3.1), and (3.2), respectively. Then, for every $\Phi \in C_b^4(\mathcal{P}_2(\mathbb{R}^d))$ and $T > 0$ there exists a constant C independent of η such that*

$$\sup_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} \sup_{n: n\eta \leq T} |\mathbb{E}\Phi(\Lambda_{n\eta}^\eta(\mu)) - \mathbb{E}\Phi(\Gamma_n^\eta(\mu))| \leq C\eta^2, \quad (3.3)$$

for all $\eta > 0$.

Proof. We first remark that the measure-valued process $\Lambda_t^\eta(\mu)$, $t \geq 0$, is uniquely defined due to Theorem 2.2 and Lemma 3.2. Without loss of generality, we consider $\eta \leq T$. The proof of this theorem relies on the comparison of the generators associated with the processes $\Gamma_n^\eta(\mu)$, $n \in \mathbb{N}_0$, and $\Lambda_t^\eta(\mu)$, $t \geq 0$, up to a certain order of η . We first demonstrate how such a bound on their difference can be used to conclude the proof.

We start from the definition of the transition semigroup for the process $\Gamma_n^\eta(\mu)$, $n \in \mathbb{N}_0$. For convenience of notation, we will drop the superscript η in Γ_n^η and Λ_t^η and simply write Γ_n and Λ_t , respectively. Note that $\Gamma_{n+1} = \Gamma_n \circ Y_n^{-1}(\Gamma_n, \cdot)$, where

$$Y_n(\mu, y) = y + \eta V(\mu, y) + \eta G(\mu, y, \theta_n), \quad \mu \in \mathcal{P}_2(\mathbb{R}^d), \quad y \in \mathbb{R}^d.$$

Indeed, by (3.1), $Z_{n+1}(z) = Y_n(\Gamma_n, Z_n(z))$, $z \in \mathbb{R}^d$, and, hence,

$$\begin{aligned} \Gamma_n \circ Y_n^{-1}(\Gamma_n, \cdot) &= \mu \circ Z_n^{-1} \circ Y_n^{-1}(\Gamma_n, \cdot) \\ &= \mu \circ Y_n(\Gamma_n, Z_n(\cdot))^{-1} = \mu \circ Z_{n+1}^{-1} = \Gamma_{n+1}, \end{aligned}$$

for all $n \in \mathbb{N}_0$. Therefore, defining the linear operator \mathcal{S} on the set of all bounded measurable functions $\Psi : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$ by

$$\mathcal{S}\Psi(\mu) = \mathbb{E}_\vartheta \Psi(\mu \circ Y_1^{-1}(\mu, \cdot)), \quad \mu \in \mathcal{P}_2(\mu),$$

we conclude that

$$\begin{aligned} \mathbb{E}_\vartheta \Psi(\Gamma_n(\mu)) &= \mathbb{E}_\vartheta \Psi(\Gamma_{n-1}(\mu) \circ Y_{n-1}^{-1}(\Gamma_{n-1}(\mu), \cdot)) \\ &= \mathbb{E}_\vartheta \left[\mathbb{E}_\vartheta \left[\Psi(\Gamma_{n-1}(\mu) \circ Y_{n-1}^{-1}(\Gamma_{n-1}(\mu), \cdot)) \middle| \Gamma_{n-1}(\mu) \right] \right] \\ &= \mathbb{E}_\vartheta \mathcal{S}\Psi(\Gamma_{n-1}(\mu)) = \dots = \mathcal{S}^n \Psi(\mu), \end{aligned} \quad (3.4)$$

for all $n \in \mathbb{N}$. Hence, defining $U(t, \mu) = \mathbb{E}\Phi(\Lambda_t(\mu))$, $t \geq 0$, $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, and using (3.4), we get for each $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ and $n \in \mathbb{N}$

$$\begin{aligned} \mathbb{E}\Phi(\Gamma_n(\mu)) - \mathbb{E}\Phi(\Lambda_{n\eta}(\mu)) &= \mathcal{S}^n \Phi(\mu) - U(t_n, \mu) \\ &= \sum_{i=0}^{n-1} \mathcal{S}^{n-i-1} (\mathcal{S}U(t_i, \mu) - U(t_{i+1}, \mu)), \end{aligned} \quad (3.5)$$

where $t_i := i\eta$.

Thus, by (3.5), and by the inequality

$$\sup_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} |\mathcal{S}\Psi(\mu)| \leq \sup_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} |\Psi(\mu)|,$$

we deduce that there exists a constant $C > 0$, such that for all $n \in \mathbb{N}$ with $n\eta \leq T$

$$\sup_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} |\mathbb{E}\Phi(\Lambda_{n\eta}(\mu)) - \mathbb{E}\Phi(\Gamma_n(\mu))| \leq \sup_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} \sum_{i=0}^{n-1} |\mathcal{S}U(t_i, \mu) - U(t_{i+1}, \mu)|. \quad (3.6)$$

In conclusion, to prove (3.3), it remains to compare $\mathcal{S}U(t_i, \mu)$ with $U(t_{i+1}, \mu)$. For this, we will expand the generators associated with the processes $\Gamma_n^\eta(\mu)$, $n \in \mathbb{N}_0$, and $\Lambda_t^\eta(\mu)$, $t \geq 0$, with respect to η up to the second order.

To obtain the expansion of $\mathcal{S}\Psi(\mu)$ for $\Psi \in C_b^3(\mathcal{P}_2(\mathbb{R}^d))$, we fix $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, $\theta \in \Theta$ and consider $Y(\mu, y) = y + \eta V(\mu, y) + \eta G(\mu, y, \theta)$, $y \in \mathbb{R}^d$, as a random variable on the probability space $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d), \mu)$. Define

$$\xi_s(y) = (1-s)y + sY(\mu, y), \quad y \in \mathbb{R}^d, \quad s \in [0, 1].$$

Then $\xi_0(y) = y$, $\xi_1(y) = Y(\mu, y)$, $\xi'_s(y) = \eta(V(\mu, y) + G(\mu, y, \theta))$ and $\text{Law}(\xi_s) := \mu \circ \xi_s^{-1}$ for all $y \in \mathbb{R}^d$, $s \in [0, 1]$. Using Taylor's formula, we obtain

$$\begin{aligned} \Psi(\mu \circ Y^{-1}(\mu, \cdot)) &= \Psi(\text{Law}(\xi_1)) = \Psi(\text{Law}(\xi_0)) + \frac{d}{ds} \Psi(\text{Law}(\xi_s)) \Big|_{s=0} \\ &\quad + \frac{1}{2} \frac{d^2}{ds^2} \Psi(\text{Law}(\xi_s)) \Big|_{s=0} + \frac{1}{2} \int_0^1 \frac{d^3}{ds^3} \Psi(\text{Law}(\xi_s)) (1-s)^3 ds. \end{aligned} \quad (3.7)$$

We next compute the derivatives appearing in the expression above. By Lemma 3.1, we get

$$\frac{d}{ds} \Psi(\text{Law}(\xi_s)) = \eta \int_{\mathbb{R}^d} \text{D}\Psi(\text{Law}(\xi_s), \xi_s(x)) \cdot (V(\mu, x) + G(\mu, x, \theta)) \mu(dx)$$

and

$$\begin{aligned} \frac{d^2}{ds^2} \Psi(\text{Law}(\xi_s)) &= \eta \frac{d}{ds} \int_{\mathbb{R}^d} \text{D}\Psi(\text{Law}(\xi_s), \xi_s(x)) \cdot (V(\mu, x) + G(\mu, x, \theta)) \mu(dx) = \\ &= \eta^2 \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} \text{D}^2 \Psi(\text{Law}(\xi_s), \xi_s(x), \xi_s(y)) \\ &\quad : (V(\mu, x) + G(\mu, x, \theta)) \otimes (V(\mu, y) + G(\mu, y, \theta)) \mu(dx) \mu(dy) \\ &\quad + \eta^2 \int_{\mathbb{R}^d} \nabla \text{D}\Psi(\text{Law}(\xi_s), \xi_s(x)) \\ &\quad : (V(\mu, x) + G(\mu, x, \theta)) \otimes (V(\mu, x) + G(\mu, x, \theta)) \mu(dx). \end{aligned}$$

The third derivative $\frac{d^3}{ds^3} \Psi(\text{Law}(\xi_s))$ can be computed analogously. Since its precise form is not needed, we omit its computation and note only that $\frac{d^3}{ds^3} \Psi(\text{Law}(\xi_s))$, $s \in [0, 1]$, is uniformly bounded by $C \|\Psi\|_{C_b^3}$ for some constant $C > 0$.

Taking the expectation of (3.7) with respect to ϑ and using the dominated convergence theorem, the equalities $\xi_0(x) = x$, $\text{Law}(\xi_0) = \mu$, $\mathbb{E}_\vartheta G(\mu, x, \theta) = 0$ and the fact that $\Psi \in C_b^3(\mathcal{P}_2(\mathbb{R}^d))$, we obtain

$$\begin{aligned}
\mathcal{S}\Psi(\mu) &= \mathbb{E}_\vartheta \Psi(\text{Law}(\xi_1)) = \Psi(\mu) + \eta \int_{\mathbb{R}^d} D\Psi(\mu, x) \cdot V(\mu, x) \mu(dx) \\
&+ \frac{\eta^2}{2} \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} D^2\Psi(\mu, x, y) : V(\mu, x) \otimes V(\mu, y) \mu(dx) \mu(dy) \\
&+ \frac{\eta^2}{2} \int_{\mathbb{R}^d} \nabla D\Psi(\mu, x) : V(\mu, x) \otimes V(\mu, x) \mu(dx) \\
&+ \frac{\eta^2}{2} \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} D^2\Psi(\mu, x, y) : \tilde{A}(\mu, x, y) \mu(dx) \mu(dy) \\
&+ \frac{\eta^2}{2} \int_{\mathbb{R}^d} \nabla D\Psi(\mu, x) : A(\mu, x) \mu(dx) + \eta^3 R_1(\Psi, \mu),
\end{aligned} \tag{3.8}$$

where $\sup_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} |R_1(\Psi, \mu)| \leq C \|\Psi\|_{C_b^3}$, for a constant $C > 0$ and

$$\tilde{A}(\mu, x, y) = \mathbb{E}_\vartheta [G(\mu, x, \theta) \otimes G(\mu, y, \theta)], \quad A(\mu, x) = \tilde{A}(\mu, x, x).$$

We next expand the generator of the process $\Lambda_t^\eta(\mu)$, $t \geq 0$. Recall that $U(t, \mu) = \mathbb{E}\Phi(\Lambda_t(\mu))$, $t \geq 0$, $\mu \in \mathcal{P}_2(\mathbb{R}^d)$. According to Proposition 2.5, we can conclude that for every $t \geq t_i$

$$U(t, \mu) = U(t_i, \mu) + \int_{t_i}^t \mathcal{L}U(r, \mu) dr, \tag{3.9}$$

where $\mathcal{L} = \mathcal{L}_1 + \eta \mathcal{L}_2$ and

$$\begin{aligned}
\mathcal{L}_1 U(r, \mu) &:= \int_{\mathbb{R}^d} V(\mu, x) \cdot DU(r, \mu, x) \mu(dx), \\
\mathcal{L}_2 U(r, \mu) &:= \frac{1}{2} \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} \tilde{A}(\mu, x, y) : D^2 U(r, \mu, x, y) \mu(dx) \mu(dy) \\
&+ \frac{1}{2} \int_{\mathbb{R}^d} A(\mu, x) : \nabla DU(r, \mu, x) \mu(dx) \\
&- \frac{1}{4} \int_{\mathbb{R}^d} \nabla |V(\mu, x)|^2 \cdot DU(r, \mu, x) \mu(dx) \\
&- \frac{1}{4} \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} D|V(\mu, x)|^2(y) \cdot DU(r, \mu, x) \mu(dx) \mu(dy).
\end{aligned}$$

Iterating the equality (3.9) as in the proof of Lemma 3 in [18], we obtain

$$U(t_{i+1}, \mu) = U(t_i, \mu) + \eta \mathcal{L}_1 U(t_i, \mu) + \eta^2 \left(\mathcal{L}_2 + \frac{1}{2} \mathcal{L}_1^2 \right) U(t_i, \mu) + \eta^3 R_2(\mu), \tag{3.10}$$

where $\sup_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} |R_2(\mu)| \leq C \|U\|_{C_b^{0,4}([0,T] \times \mathcal{P}_2(\mathbb{R}^d))}$ for a constant $C > 0$.

In order to compare $SU(t_i, \mu)$ and $U(t_{i+1}, \mu)$, we next express $\mathcal{L}_2 + \frac{1}{2}\mathcal{L}_1^2$ in terms of the coefficients of the equation (3.2). Note that, according to Example 2.4, we have

$$\begin{aligned} D\mathcal{L}_1U(r, \mu, x) &= \nabla [V(\mu, x) \cdot DU(r, \mu, x)] + \int_{\mathbb{R}^d} D[V(\mu, y) \cdot DU(r, \mu, y)](x)\mu(dy) \\ &= DU(r, \mu, x)\nabla V(\mu, x) + V(\mu, x)\nabla DU(r, \mu, x) \\ &\quad + \int_{\mathbb{R}^d} DU(r, \mu, y)DV(\mu, y, x)\mu(dy) + \int_{\mathbb{R}^d} V(\mu, y)D^2U(r, \mu, y, x)\mu(dy). \end{aligned}$$

Thus, using the equality $\frac{1}{2}\nabla|V(\mu, x)|^2 = V(\mu, x)\nabla V(\mu, x)$ and $\frac{1}{2}D|V(\mu, x)|^2(y) = V(\mu, x)DV(\mu, x, y)$, we get

$$\begin{aligned} \mathcal{L}_1^2U(r, \mu, x) &= \frac{1}{2} \int_{\mathbb{R}^d} \nabla|V(\mu, x)|^2 \cdot DU(r, \mu, x)\mu(dx) \\ &\quad + \int_{\mathbb{R}^d} \nabla DU(r, \mu, x) : V(\mu, x) \otimes V(\mu, x)\mu(dx) \\ &\quad + \frac{1}{2} \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} D|V(\mu, x)|^2(y) \cdot DU(r, \mu, x)\mu(dx)\mu(dy) \\ &\quad + \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} D^2U(r, \mu, x, y) : V(\mu, x) \otimes V(\mu, y)\mu(dx)\mu(dy). \end{aligned}$$

Consequently,

$$\begin{aligned} \left(\mathcal{L}_2 + \frac{1}{2}\mathcal{L}_1^2\right)U(r, \mu) &= \frac{1}{2} \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} \tilde{A}(\mu, x, y) : D^2U(r, \mu, x, y)\mu(dx)\mu(dy) \\ &\quad + \frac{1}{2} \int_{\mathbb{R}^d} A(\mu, x) : \nabla DU(r, \mu, x)\mu(dx) \\ &\quad - \frac{1}{4} \int_{\mathbb{R}^d} \nabla|V(\mu, x)|^2 \cdot DU(r, \mu, x)\mu(dx) \\ &\quad - \frac{1}{4} \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} D|V(\mu, x)|^2(y) \cdot DU(r, \mu, x)\mu(dx)\mu(dy) \\ &\quad + \frac{1}{4} \int_{\mathbb{R}^d} \nabla|V(\mu, x)|^2 \cdot DU(r, \mu, x)\mu(dx) \\ &\quad + \frac{1}{2} \int_{\mathbb{R}^d} \nabla DU(r, \mu, x) : V(\mu, x) \otimes V(\mu, x)\mu(dx) \\ &\quad + \frac{1}{4} \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} D|V(\mu, x)|^2(y) \cdot DU(r, \mu, x)\mu(dx)\mu(dy) \\ &\quad + \frac{1}{2} \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} D^2U(r, \mu, x, y) : V(\mu, x) \otimes V(\mu, y)\mu(dx)\mu(dy) \\ &= \frac{1}{2} \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} \tilde{A}(\mu, x, y) : D^2U(r, \mu, x, y)\mu(dx)\mu(dy) \\ &\quad + \frac{1}{2} \int_{\mathbb{R}^d} A(\mu, x) : \nabla DU(r, \mu, x)\mu(dx) \end{aligned}$$

$$\begin{aligned}
& + \frac{1}{2} \int_{\mathbb{R}^d} \nabla D U(r, \mu, x) : V(\mu, x) \otimes V(\mu, x) \mu(dx) \\
& + \frac{1}{2} \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} D^2 U(r, \mu, x, \mu) : V(\mu, x) \otimes V(\mu, y) \mu(dx) \mu(dy).
\end{aligned}$$

Comparing (3.10) with (3.8) for $\Psi = U(t_i, \cdot)$, we conclude that

$$\begin{aligned}
\mathcal{S}U(t_i, \mu) & = U(t_i, \mu) + \eta \mathcal{L}_1 U(t_i, \mu) + \eta^2 \left(\mathcal{L}_2 + \frac{1}{2} \mathcal{L}_1^2 \right) U(t_i, \mu) + \eta^3 R_1(U(t_i, \cdot), \mu) \\
& = U(t_{i+1}, \mu) + \eta^3 R_1(U(t_i, \cdot), \mu) - \eta^3 R_2(\mu).
\end{aligned}$$

Inserting into (3.5), and using the fact that $U \in C_b^{0,4}([0, T] \times \mathcal{P}_2(\mathbb{R}^d))$ (see Proposition (2.5)), yields, for all $n \in \mathbb{N}$ with $n\eta \leq T$

$$\begin{aligned}
\sup_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} |\mathbb{E}\Phi(\Lambda_{n\eta}(\mu)) - \mathbb{E}\Phi(\Gamma_n(\mu))| & \leq \sup_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} \sum_{i=0}^{n-1} \eta^3 |R_1(U(t_i, \cdot), \mu) - R_2(\mu)| \\
& \leq Cn\eta^3 \leq CT\eta^2.
\end{aligned}$$

This completes the proof of the theorem. \square

Remark 3.4. From the proof of Theorem 3.3 one can see that for every $\Phi \in C_b^4(\mathcal{P}_2(\mathbb{R}^d))$ and $T > 0$ there exists a constant $C > 0$ such that

$$\sup_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} \sup_{n: n\eta \leq T} |\mathbb{E}\Phi(\Lambda_{n\eta}(\mu)) - \mathbb{E}\Phi(\Gamma_n(\mu))| \leq C\eta,$$

for all $\eta > 0$, if $\Lambda_t = \Lambda_t(\mu)$, $t \geq 0$, is defined by the SDE with interaction

$$\begin{aligned}
dX_t(x) & = V(\Lambda_t, X_t(x))dt + \sqrt{\eta} \int_{\Theta} G(\Lambda_t, X_t(x)) W(d\theta, dt), \\
X_0(x) & = x, \quad \Lambda_t = \mu \circ X_t^{-1}, \quad x \in \mathbb{R}^d, \quad t \geq 0.
\end{aligned}$$

We now apply Theorem 3.3 to the comparison of the SGD dynamics and stochastic modified flows considered in the introduction. First, we recover a variant of the statement for stochastic modified equations.

Corollary 3.5. *Let $Z_n^\eta(x)$, $n \in \mathbb{N}_0$, be defined by (1.1) for a loss function \tilde{R} and $X_t^\eta(x)$, $t \geq 0$, be a solution to (1.4). Let also $\tilde{R}(\cdot, \theta) \in C_b^6(\mathbb{R}^d)$ for ϑ -a.e. $\theta \in \Theta$ and assume that*

$$\int_{\Theta} \|\tilde{R}(\cdot, \theta)\|_{C_b^6}^2 \vartheta(d\theta) < \infty.$$

Then, for every $f \in C_b^4(\mathbb{R}^d)$ and $T > 0$, there exists a constant $C > 0$ independent of η such that

$$\sup_{x \in \mathbb{R}^d} \sup_{n: n\eta \leq T} |\mathbb{E}f(X_{n\eta}^\eta(x)) - \mathbb{E}f(Z_n^\eta(x))| \leq C\eta^2$$

for all $\eta > 0$.

Proof. Using the dominated convergence theorem it is easily seen that the functions $V := -\nabla R$ and $G := \nabla \tilde{R} - \nabla R$ belong to $C_b^{5,5}(\mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d)$ and $\tilde{C}_b^{4,4}(\mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d)$, respectively, where $R = \mathbb{E}_\vartheta \tilde{R}$. Hence, applying Theorem 3.3 to the function $\Phi(\mu) = \langle f, \mu \rangle$, $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, that trivially belongs to $C_b^4(\mathcal{P}_2(\mathbb{R}^d))$, we obtain

$$\begin{aligned} & \sup_{x \in \mathbb{R}^d} \sup_{n: n\eta \leq T} \left| \mathbb{E}f(X_{n\eta}^\eta(x)) - \mathbb{E}f(Z_n^\eta(x)) \right| \\ &= \sup_{\mu = \delta_x, x \in \mathbb{R}^d} \sup_{n: n\eta \leq T} \left| \mathbb{E}\langle f(X_{n\eta}^\eta), \mu \rangle - \mathbb{E}\langle f(Z_n^\eta), \mu \rangle \right| \\ &\leq \sup_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} \sup_{n: n\eta \leq T} \left| \mathbb{E}\Phi(\Lambda_{n\eta}^\eta(\mu)) - \mathbb{E}\Phi(\Gamma_n^\eta(\mu)) \right| \leq C\eta^2, \end{aligned}$$

for all $\eta > 0$ and some constant $C > 0$ independent of η , where $\Lambda_t^\eta(\mu) = \mu \circ (X_t^\eta)^{-1}$ and $\Gamma_n^\eta(\mu) = \mu \circ (Z_n^\eta)^{-1}$. This completes the proof of the statement. \square

Corollary 3.6. *Under the assumptions of Corollary 3.5, for every $m \in \mathbb{N}$, $f \in C_b^4(\mathbb{R}^{dm})$, $\Phi \in C_b^4(\mathcal{P}_2(\mathbb{R}^d))$ and $T > 0$ there exists a constant $C > 0$ independent of η such that*

$$\sup_{x_1, \dots, x_m \in \mathbb{R}^d} \sup_{n: n\eta \leq T} \left| \mathbb{E}f(X_{n\eta}^\eta(x_1), \dots, X_{n\eta}^\eta(x_m)) - \mathbb{E}f(Z_n^\eta(x_1), \dots, Z_n^\eta(x_m)) \right| \leq C\eta^2 \quad (3.11)$$

and

$$\sup_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} \sup_{n: n\eta \leq T} \left| \mathbb{E}\Phi(\mu \circ (X_{n\eta}^\eta)^{-1}) - \mathbb{E}\Phi(\mu \circ (Z_n^\eta)^{-1}) \right| \leq C\eta^2 \quad (3.12)$$

for all $\eta > 0$.

Proof. The estimate (3.12) can be obtained by the same argument as in the proof of Corollary 3.5. To prove (3.11), we will apply Corollary 3.5 to the function

$$\tilde{R}^{\text{ext}}(z, \theta) = \tilde{R}(z_1, \theta) + \dots + \tilde{R}(z_m, \theta), \quad z = (z_i)_{i \in [m]} \in \mathbb{R}^{dm}, \quad \theta \in \Theta.$$

Note that

$$\nabla \tilde{R}^{\text{ext}}(z, \theta) = \left(\nabla_{z_i} \tilde{R}(z_i, \theta) \right)_{i \in [m]}$$

for all $z = (z_i)_{i \in [m]} \in \mathbb{R}^{dm}$ and $\theta \in \Theta$. Defining $Z_n^{\text{ext}, \eta}(x)$, $n \in \mathbb{N}_0$, by (1.1) with \tilde{R} and \mathbb{R}^d replaced by \tilde{R}_{ext} and \mathbb{R}^{dm} , respectively, it is easily seen that

$$Z_n^{\text{ext}, \eta}(x) = (Z_n^\eta(x_i))_{i \in [m]}, \quad n \in \mathbb{N}_0,$$

for all $x = (x_i)_{i \in [m]} \in \mathbb{R}^{dm}$.

We next set $R^{\text{ext}}(z) = \mathbb{E}_\vartheta \tilde{R}^{\text{ext}}(z, \theta)$, $z = (z_i)_{i \in [m]}$. Then

$$\nabla R^{\text{ext}}(z) = (\nabla_{z_i} R(z_i))_{i \in [m]}$$

and

$$\nabla |\nabla R^{\text{ext}}(z)|^2 = (\nabla_{z_i} |\nabla_{z_i} R(z_i)|^2)_{i \in [m]}$$

for all $z = (z_i)_{i \in [m]} \in \mathbb{R}^{dm}$. Moreover,

$$G^{\text{ext}}(z, \theta) := \nabla \tilde{R}^{\text{ext}}(z, \theta) - \nabla R^{\text{ext}}(z, \theta) = (G(z_i, \theta))_{i \in [m]},$$

where G is the coefficient of (1.4) that equals $\nabla \tilde{R} - \nabla R$. Under the assumptions of the corollary, equation (1.4) with R and G replaced by R^{ext} and G^{ext} , respectively, has a unique solution $X_t^{\text{ext}, \eta}(x)$, $x \in \mathbb{R}^{dm}$, $t \geq 0$. Moreover,

$$X_t^{\text{ext}, \eta}(x) = (X_t^\eta(x_i))_{i \in [m]}, \quad t \geq 0,$$

a.s. for all $x = (x_i)_{i \in [m]}$. Since \tilde{R}^{ext} satisfies the assumptions of Corollary 3.5, one gets for every $f \in C_b^4(\mathbb{R}^{dm})$

$$\begin{aligned} & \sup_{x \in \mathbb{R}^{dm}} |\mathbb{E}f(X_{n\eta}^{\text{ext}, \eta}(x)) - f(Z_n^{\text{ext}, \eta}(x))| \\ &= \sup_{x_1, \dots, x_m \in \mathbb{R}^d} \sup_{n: n\eta \leq T} |\mathbb{E}f(X_{n\eta}^\eta(x_1), \dots, X_{n\eta}^\eta(x_m)) - \mathbb{E}f(Z_n^\eta(x_1), \dots, Z_n^\eta(x_m))| \leq C\eta^2 \end{aligned}$$

for a constant $C > 0$ independent of η . This completes the proof of the statement. \square

In the next example, we show that Corollary 3.6 cannot hold for the solution to the classical stochastic modified equation (1.2), since the distribution of the two-point motion is different from the distribution of the two-point motion of (1.4).

Example 3.7. The covariation of the two-point motion $(X_t^\eta(x), X_t^\eta(\bar{x}))$, $t \geq 0$, from the SMF (1.4) equals

$$[X^\eta(x), X^\eta(\bar{x})]_t = \eta \int_0^t \tilde{A}(X_s^\eta(x), X_s^\eta(\bar{x})) ds, \quad t \geq 0, \quad (3.13)$$

where $\tilde{A}(x, y) = \langle G(x, \cdot) \otimes G(y, \cdot) \rangle_\vartheta$. However, the covariation of the two-point motion $(Y_t^\eta(x), Y_t^\eta(\bar{x}))$, $t \geq 0$, obtained from the SDE (1.2), is given by

$$[Y^\eta(x), Y^\eta(\bar{x})]_t = \eta \int_0^t \Sigma(Y_s^\eta(x))^{1/2} \Sigma(Y_s^\eta(\bar{x}))^{1/2} ds, \quad t \geq 0, \quad (3.14)$$

for $\Sigma(x) = \tilde{A}(x, x)$. This implies that the processes $(X^\eta(x), X^\eta(\bar{x}))$ and $(Y^\eta(x), Y^\eta(\bar{x}))$ have different distributions in general. We further notice that the covariance of the one step SGD dynamics defined by (1.1) satisfies

$$\text{cov}(Z_1^\eta(x), Z_1^\eta(y)) = \eta^2 \tilde{A}(x, y),$$

which is comparable with (3.13), but not with (3.14).

Next, we consider the SGD scheme $Z_n^\eta = (Z_n^{i, \eta})_{i \in [M]}$, $n \in \mathbb{N}_0$, incorporating the infinite width limit that is defined by (1.9), where $Z_0^{i, \eta}$, $i \in [M]$, are i.i.d. random variables sampled from a measure $\mu \in \mathcal{P}_2(\mathbb{R}^d)$. We prove the convergence of the empirical distribution process $\Gamma_n^{M, \eta} = \frac{1}{M} \sum_{i=1}^M \delta_{Z_n^{i, \eta}}$, $n \in \mathbb{N}_0$, to a mean-field solution $\Lambda_t^\eta = \mu \circ (X_t^\eta)^{-1}$, $t \geq 0$, of the DDSMF defined by (1.10).

Corollary 3.8. Let $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ and $\mu^M = \frac{1}{M} \sum_{j=1}^M \delta_{Z_0^{j,\eta}}$, where $Z_0^{j,\eta}$, $j \in [M]$, are i.i.d. random variables with distribution μ . Let $\Gamma_n^{M,\eta}$, $n \in \mathbb{N}_0$, and Λ_t^η , $t \geq 0$, be as in (1.9) and (1.10), respectively, with $\Gamma_0^{M,\eta} = \mu^M$ and $\Lambda_0^\eta = \mu$. Assume that the function Ψ in (1.6) satisfies: $\Psi(\cdot, \theta) \in C_b^6(\mathbb{R}^d)$ for ϑ -a.e. $\theta \in \Theta$ and

$$\int_{\Theta} \left(\|\Psi(\cdot, \theta)\|_{C_b^6}^2 + |f(\theta)|^2 \right) \|\Psi(\cdot, \theta)\|_{C_b^6}^2 \vartheta(d\theta) < \infty.$$

Then, for every $\Phi \in C_b^4(\mathcal{P}_2(\mathbb{R}^d))$ there exists a constant $C > 0$ independent of η and M such that

$$\sup_{n:n\eta \leq T} |\mathbb{E}\Phi(\Lambda_{n\eta}^\eta) - \mathbb{E}\Phi(\Gamma_n^{M,\eta})| \leq C\eta^2 + C\sqrt{\mathbb{E}\mathcal{W}_2^2(\mu, \mu^M)} \quad (3.15)$$

for all $\eta > 0$ and $M \in \mathbb{N}$. In particular, if μ has finite p th moment for some $p > 2$, with $p \neq 4$ for $d \leq 4$ and $p \neq \frac{d}{d-2}$ for $d \geq 5$, then for every $a > 0$ there exists a constant $C > 0$ independent of η and M such that

$$\sup_{n:n\eta \leq T} |\mathbb{E}\Phi(\Lambda_{n\eta}^\eta) - \mathbb{E}\Phi(\Gamma_n^{M,\eta})| \leq C\eta^2 \quad (3.16)$$

for all $\eta > 0$ and $M \in \mathbb{N}$ satisfying $\frac{K(M)}{\eta^4} \leq a$, where

$$K(M) = \begin{cases} M^{-\frac{1}{2}} + M^{-\frac{p-2}{2}} & \text{if } d \leq 3, \\ M^{-\frac{1}{2}} \ln(1+M) + M^{-\frac{p-2}{2}} & \text{if } d = 4, \\ M^{-\frac{2}{d}} + M^{-\frac{p-2}{2}} & \text{if } d \geq 5. \end{cases}$$

Proof. First, we show that $V \in C_b^{5,5}(\mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d)$ and $G \in \tilde{C}_b^{4,4}(\mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d)$, where V and G are given by (1.8). Analogously to the proof of Corollary 3.5, we get that $F \in C_b^6(\mathbb{R}^d)$, where $F(z) = \mathbb{E}_\vartheta[f(\theta) \cdot \Psi(z, \theta)]$, $z \in \mathbb{R}^d$, and, thus, $\nabla F \in C_b^5(\mathbb{R}^d)$. For $K(z^1, z^2) = \mathbb{E}_\vartheta[\Psi(z^1, \theta) \cdot \Psi(z^2, \theta)]$, $z_1, z_2 \in \mathbb{R}^d$, we use the dominated convergence theorem to get that $K \in C_b^6(\mathbb{R}^{2d})$. Using Example 2.4, we get, for $\tilde{K}(\mu, z) = \langle \nabla_z K(z, \cdot), \mu \rangle$, $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, $z \in \mathbb{R}^d$, that

$$D\tilde{K}(\mu, z^1, z^2) = \nabla_{z^2} \nabla_{z^1} K(z^1, z^2),$$

with analogous expressions for higher derivatives. Thus, $\tilde{K} \in C_b^{5,5}(\mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d)$ and, therefore, $V \in C_b^{5,5}(\mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d)$. To see that $G \in \tilde{C}_b^{4,4}(\mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d)$ note that

$$\tilde{G}(\mu, z, \theta) = (f(\theta) - \langle \Psi(\cdot, \theta), \mu \rangle) \nabla_z \Psi(z, \theta), \quad \mu \in \mathcal{P}_2(\mathbb{R}^d), z \in \mathbb{R}^d, \theta \in \Theta,$$

satisfies $\tilde{G} \in \tilde{C}_b^{4,4}(\mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d)$ and $\mathbb{E}_\vartheta[\tilde{G}(\cdot, \cdot, \theta)] \in C_b^{4,4}(\mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d)$.

Note that one needs to check the estimate (3.15) only for $\eta \in (0, T]$. Let $\Lambda_t^\eta(\mu)$, $t \geq 0$, be defined by (1.10) for every $\mu \in \mathcal{P}_2(\mathbb{R}^d)$. We next fix $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ and consider the empirical distribution $\mu^M = \frac{1}{M} \sum_{i=1}^M \delta_{Z_0^{i,\eta}}$ associated with the family of i.i.d. random

variables $Z_0^{i,\eta}$, $i \in [M]$, sampled from the distribution μ . By Theorem 3.3, there exists a constant $C > 0$ independent of η such that

$$\sup_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} \sup_{n: n\eta \leq T} |\mathbb{E}\Phi(\Lambda_{n\eta}^\eta(\mu)) - \mathbb{E}\Phi(\Gamma_n^\eta(\mu))| \leq C\eta^2$$

for all $\eta \in (0, T]$, where $\Gamma_n^\eta(\mu)$, $n \in \mathbb{N}_0$, is determined by (3.1) with V and G given by (1.8). Therefore, using the equality $\Gamma_n^{M,\eta} = \Gamma_n^\eta(\mu^M)$ for all $n \in \mathbb{N}_0$, one has

$$\begin{aligned} \sup_{n: n\eta \leq T} |\mathbb{E}\Phi(\Lambda_{n\eta}^\eta(\mu^M)) - \mathbb{E}\Phi(\Gamma_n^{M,\eta})| &= \sup_{n: n\eta \leq T} |\mathbb{E}\Phi(\Lambda_{n\eta}^\eta(\mu^M)) - \mathbb{E}\Phi(\Gamma_n^\eta(\mu^M))| \\ &= \sup_{n: n\eta \leq T} \left| \mathbb{E} \left[\mathbb{E} \left[\Phi(\Lambda_{n\eta}^\eta(\mu^M)) \middle| \mathcal{A} \right] - \mathbb{E} \left[\Phi(\Gamma_n^\eta(\mu^M)) \middle| \mathcal{A} \right] \right] \right| \\ &\leq \mathbb{E} \left[\sup_{n: n\eta \leq T} \left| \mathbb{E} \left[\Phi(\Lambda_{n\eta}^\eta(\mu^M)) \middle| \mathcal{A} \right] - \mathbb{E} \left[\Phi(\Gamma_n^\eta(\mu^M)) \middle| \mathcal{A} \right] \right| \right] \\ &\leq \sup_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} \sup_{n: n\eta \leq T} |\mathbb{E}\Phi(\Lambda_{n\eta}^\eta(\mu)) - \mathbb{E}\Phi(\Gamma_n^\eta(\mu))| \leq C\eta^2 \end{aligned}$$

for all $\eta \in (0, T]$ and $M \in \mathbb{N}$, where $\mathcal{A} = \sigma(Z_0^{i,\eta}, i \in [M])$.

We next compare $\mathbb{E}\Phi(\Lambda_{n\eta}^\eta(\mu))$ with $\mathbb{E}\Phi(\Lambda_{n\eta}^\eta(\mu^M))$. Applying Lemma 3.2 to $V = \Phi$ and $G = 0$, we can estimate

$$|\mathbb{E}\Phi(\Lambda_{n\eta}^\eta(\mu)) - \mathbb{E}\Phi(\Lambda_{n\eta}^\eta(\mu^M))|^2 \leq \|\Phi\|_{C_1}^2 \mathbb{E}\mathcal{W}_2^2(\Lambda_{n\eta}^\eta(\mu), \Lambda_{n\eta}^\eta(\mu^M)).$$

Since the coefficients of the SDE (1.10) are Lipschitz continuous, where the Lipschitz constant can be chosen independently of $\eta \in (0, T]$ due to the assumptions of the corollary and Lemma 3.2, we can apply Proposition 2.3 to bound $\mathbb{E}\mathcal{W}_2^2(\Lambda_{n\eta}^\eta(\mu), \Lambda_{n\eta}^\eta(\mu^M))$. Thus, there exists a constant $C > 0$ independent of η , M and n such that

$$\mathbb{E}\mathcal{W}_2^2(\Lambda_{n\eta}^\eta(\mu), \Lambda_{n\eta}^\eta(\mu^M)) \leq C\mathbb{E}\mathcal{W}_2^2(\mu, \mu^M)$$

for all $\eta \in (0, T]$, $M \in \mathbb{N}$ and $n \in \mathbb{N}_0$ with $n\eta \leq T$. This completes the proof of the first part of the corollary.

If μ has finite p th moment for $p > 2$ such that $p \neq 4$ for $d \leq 4$ and $p \neq \frac{d}{d-2}$ for $d \geq 5$, then, by Theorem 1 in [12],

$$\mathbb{E}\mathcal{W}_2^2(\mu, \mu^M) \leq C_1 \langle \phi_p, \mu \rangle^{\frac{2}{p}} K(M),$$

where $\phi_p(x) = |x|^p$, $x \in \mathbb{R}^d$, and $C_1 > 0$ depends only on p and d . Assuming that $\frac{K(M)}{\eta^4} \leq a$ for some $a > 0$, we get

$$\sup_{n: n\eta \leq T} |\mathbb{E}\Phi(\Lambda_{n\eta}^\eta) - \mathbb{E}\Phi(\Gamma_n^{M,\eta})| \leq C\eta^2 + C\sqrt{\mathbb{E}\mathcal{W}_2^2(\mu, \mu^M)} \leq C\eta^2 + C\sqrt{aC_1} \langle \phi_p, \mu \rangle^{\frac{1}{p}} \eta^2.$$

This completes the proof of the second part of the statement. \square

Remark 3.9. Assume that the measure $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ has all finite moments in Corollary 3.8. Then we can choose p so large that the first term in every case of the definition of the constant $K(M)$ dominates. Therefore, the estimate (3.16) holds for all $\eta > 0$ and $M \geq \frac{a}{\eta^q}$, where $q = 8$ for $d \leq 3$, $q = 2d$ for $d \geq 5$ and any $q > 8$ for $d = 4$, since $\frac{K(M)}{\eta^4} \leq a$ is satisfied for some $a > 0$ and large enough p .

Acknowledgements

The authors were supported by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) – SFB 1283/2 2021 – 317210226. BG acknowledges support by the Max Planck Society through the Research Group ”Stochastic Analysis in the Sciences (SAiS)”. The third author thanks the Max Planck Institute for Mathematics in the Sciences for its warm hospitality, where a part of this research was carried out.

References

- [1] M. A. Belozerova, *Asymptotic behavior of solutions to stochastic differential equations with interaction*, Theory Stoch. Process. **25** (2020), no. 2, 1–8. MR 4354470
- [2] Pierre Cardaliaguet, *Notes on mean field games*, P.-L. Lions lectures at College de France. (2013), Online at <https://www.ceremade.dauphine.fr/~cardaliaguet/MFG20130420.pdf>.
- [3] René Carmona, François Delarue, and Daniel Lacker, *Mean field games with common noise*, Ann. Probab. **44** (2016), no. 6, 3740–3803. MR 3572323
- [4] Lénaïc Chizat and Francis Bach, *On the global convergence of gradient descent for over-parameterized models using optimal transport*, Advances in Neural Information Processing Systems, vol. 31, Curran Associates, Inc., 2018.
- [5] ———, *Implicit bias of gradient descent for wide two-layer neural networks trained with the logistic loss*, arXiv:2002.04486 (2020).
- [6] A. A. Dorogovtsev, *Measure-valued Markov processes and stochastic flows on abstract spaces*, Stoch. Stoch. Rep. **76** (2004), no. 5, 395–407. MR 2096728
- [7] ———, *Meroznachnye protsessy i stokhasticheskie potoki [Measurevalued processes and stochastic flows]*, vol. 66, Proceedings of Institute of Mathematics of NAS of Ukraine. Mathematics and its Applications, Institut Matematiki, Kiev, 2007 (Russian). MR 2375817
- [8] A. A. Dorogovtsev and P. Kotelenez, *Smooth stationary solutions of quasilinear stochastic differential equations. Finite mass*, Preprint No. 97. Department of Mathematics Case Western Reserv University Cleveland, Ohio, 1997.
- [9] A. A. Dorogovtsev and O. V. Ostapenko, *Large deviations for flows of interacting Brownian motions*, Stoch. Dyn. **10** (2010), no. 3, 315–339. MR 2671379
- [10] Weinan E, Chao Ma, and Lei Wu, *Machine learning from a continuous viewpoint, I*, Sci. China Math. **63** (2020), no. 11, 2233–2266. MR 4170870

- [11] Yuanyuan Feng, Lei Li, and Jian-Guo Liu, *Semigroups of stochastic gradient descent and online principal component analysis: properties and diffusion approximations*, Communications in Mathematical Sciences **16** (2018), no. 3.
- [12] Nicolas Fournier and Arnaud Guillin, *On the rate of convergence in Wasserstein distance of the empirical measure*, Probab. Theory Related Fields **162** (2015), no. 3–4, 707–738. MR 3383341
- [13] Leszek Gawarecki and Vidyadhar Mandrekar, *Stochastic differential equations in infinite dimensions with applications to stochastic partial differential equations*, Probability and its Applications (New York), Springer, Heidelberg, 2011. MR 2560625
- [14] Benjamin Gess, Rishabh S. Gvalani, and Vitalii Konarovskiy, *Conservative SPDEs as fluctuating mean field limits of stochastic gradient descent*, arXiv:2207.05705 (2022).
- [15] Wenqing Hu, Chris Junchi Li, Lei Li, and Jian-Guo Liu, *On the diffusion approximation of nonconvex stochastic gradient descent*, Annals of Mathematical Sciences and Applications **4** (2019), no. 1.
- [16] Adel Javanmard, Marco Mondelli, and Andrea Montanari, *Analysis of a two-layer neural network via displacement convexity*, Ann. Statist. **48** (2020), no. 6, 3619–3642. MR 4185822
- [17] Thomas G. Kurtz and Jie Xiong, *Particle representations for a class of nonlinear SPDEs*, Stochastic Process. Appl. **83** (1999), no. 1, 103–126. MR 1705602
- [18] Lei Li and Yuliang Wang, *On uniform-in-time diffusion approximation for stochastic gradient descent*, arXiv:2207.04922 (2022).
- [19] Qianxiao Li, Cheng Tai, and Weinan E, *Stochastic modified equations and dynamics of stochastic gradient algorithms I: mathematical foundations*, J. Mach. Learn. Res. **20** (2019), Paper No. 40, 47. MR 3948080
- [20] Qianxiao Li, Cheng Tai, and E Weinan, *Stochastic modified equations and adaptive stochastic gradient algorithms*, International Conference on Machine Learning, PMLR, 2017, pp. 2101–2110.
- [21] Zhiyuan Li, Sadhika Malladi, and Sanjeev Arora, *On the validity of modeling SGD with stochastic differential equations (SDEs)*, Advances in Neural Information Processing Systems **34** (2021), 12712–12725.
- [22] Song Mei, Andrea Montanari, and Phan-Minh Nguyen, *A mean field view of the landscape of two-layer neural networks*, Proc. Natl. Acad. Sci. USA **115** (2018), no. 33, E7665–E7671. MR 3845070

- [23] Phan-Minh Nguyen, *Mean field limit of the learning dynamics of multilayer neural networks*, arXiv:1902.02880 (2019).
- [24] Atsushi Nitanda and Taiji Suzuki, *Stochastic particle gradient descent for infinite ensembles*, arXiv:1712.05438 (2017).
- [25] A. Yu. Pilipenko, *Support theorem on stochastic flows with interaction*, Theory Stoch. Process. **12** (2006), no. 1-2, 127–141. MR 2316293
- [26] Herbert Robbins and Sutton Monro, *A stochastic approximation method*, The Annals of Mathematical Statistics **22** (1951), no. 3, 400–407. MR MR42668
- [27] Grant Rotskoff, Samy Jelassi, Joan Bruna, and Eric Vanden-Eijnden, *Neuron birth-death dynamics accelerates gradient descent and converges asymptotically*, Proceedings of the 36th International Conference on Machine Learning (Kamalika Chaudhuri and Ruslan Salakhutdinov, eds.), Proceedings of Machine Learning Research, vol. 97, PMLR, 09–15 Jun 2019, pp. 5508–5517.
- [28] Grant Rotskoff and Eric Vanden-Eijnden, *Neural networks as interacting particle systems: Asymptotic convexity of the loss landscape and universal scaling of the approximation error*, CoRR **abs/1805.00915** (2018).
- [29] ———, *Parameters as interacting particles: Long time convergence and asymptotic error scaling of neural networks*, Advances in Neural Information Processing Systems, vol. 31, Curran Associates, Inc., 2018.
- [30] ———, *Trainability and accuracy of artificial neural networks: an interacting particle system approach*, Comm. Pure Appl. Math. **75** (2022), no. 9, 1889–1935. MR 4465905
- [31] Yuzuru Sato, Daiji Tsutsui, and Akio Fujiwara, *Noise-induced degeneration in online learning*, Physica D: Nonlinear Phenomena **430** (2022), 133095.
- [32] Justin Sirignano and Konstantinos Spiliopoulos, *Mean field analysis of neural networks: a central limit theorem*, Stochastic Process. Appl. **130** (2020), no. 3, 1820–1852. MR 4058290
- [33] ———, *Mean field analysis of neural networks: a law of large numbers*, SIAM J. Appl. Math. **80** (2020), no. 2, 725–752. MR 4074020
- [34] Feng-Yu Wang, *Image-dependent conditional McKean-Vlasov SDEs for measure-valued diffusion processes*, J. Evol. Equ. **21** (2021), no. 2, 2009–2045. MR 4278420
- [35] Lei Wu, Chao Ma, and Weinan E, *How SGD selects the global minima in over-parameterized learning: a dynamical stability perspective*, Advances in Neural Information Processing Systems, vol. 31, Curran Associates, Inc., 2018.