

# A trajectorial approach to relative entropy dissipation of McKean–Vlasov diffusions: Gradient flows and HWBI inequalities

BERTRAM TSCHIDERER<sup>1,a</sup> and LANE CHUN YEUNG<sup>2,b</sup>

<sup>1</sup>*Faculty of Mathematics, University of Vienna, Austria, [bertram.tschiderer@univie.ac.at](mailto:bertram.tschiderer@univie.ac.at)*

<sup>2</sup>*Department of Industrial Engineering and Operations Research, Columbia University, New York, USA,*

<sup>b</sup>*[l.yeung@columbia.edu](mailto:l.yeung@columbia.edu)*

We formulate a trajectorial version of the relative entropy dissipation identity for McKean–Vlasov diffusions, extending recent results which apply to non-interacting diffusions. Our stochastic analysis approach is based on time-reversal of diffusions and Lions’ differential calculus over Wasserstein space. It allows us to compute explicitly the rate of relative entropy dissipation along every trajectory of the underlying diffusion via the semimartingale decomposition of the corresponding relative entropy process. As a first application, we obtain a new interpretation of the gradient flow structure for the granular media equation, generalizing a formulation developed recently for the linear Fokker–Planck equation. Secondly, we show how the trajectorial approach leads to a new derivation of the HWBI inequality, which relates relative entropy (H), Wasserstein distance (W), barycenter (B) and Fisher information (I).

**Keywords:** Relative entropy dissipation; gradient flow; McKean–Vlasov diffusion; granular media equation; HWBI inequality

## 1. Introduction

We are interested in the relative entropy dissipation of McKean–Vlasov stochastic differential equations of the form

$$dX_t = -(\nabla V(X_t) + \nabla(W * P_t)(X_t)) dt + \sqrt{2} dB_t, \quad 0 \leq t \leq T, \quad (1.1)$$

where  $X_0$  has some given initial distribution  $P_0$  on  $\mathbb{R}^n$ . Here, the functions  $V, W: \mathbb{R}^n \rightarrow [0, \infty)$  play the roles of confinement and interaction potentials and are assumed to be suitably regular,  $P_t := \text{Law}(X_t)$  denotes the distribution of the random vector  $X_t$ , the symbol  $*$  stands for the standard convolution operator, and  $(B_t)_{0 \leq t \leq T}$  is a standard  $n$ -dimensional Brownian motion. In particular, this SDE is non-local (or non-linear) in the sense that the drift term depends on the distribution of the state variable. Non-local equations of this form arise in the modeling of weakly interacting diffusion equations, after the seminal work of McKean [47].

Since the work of Carrillo–McCann–Villani [13, 14], relative entropy dissipation has been known to be an effective method for studying convergence rates to equilibrium and propagation of chaos of McKean–Vlasov equations. Some notable examples include the works [8, 11, 12, 16, 45, 54]. In a broader context, [34, 48] recently applied entropy methods to the mean-field theory of neural networks.

We denote by  $\mathcal{P}_{\text{ac}}(\mathbb{R}^n)$  the set of absolutely continuous probability measures on  $\mathbb{R}^n$ , which we will often identify with their corresponding probability density functions with respect to Lebesgue measure. The *free energy functional*

$$\mathcal{P}_{\text{ac}}(\mathbb{R}^n) \ni p \longmapsto \mathcal{F}(p) := \mathcal{U}(p) + \mathcal{V}(p) + \mathcal{W}(p) \quad (1.2)$$

is defined as the sum of the energy functionals

$$\mathcal{U}(p) := \int_{\mathbb{R}^n} p(x) \log p(x) dx, \quad \mathcal{V}(p) := \int_{\mathbb{R}^n} V(x) p(x) dx, \quad \mathcal{W}(p) := \frac{1}{2} \int_{\mathbb{R}^n} (W * p)(x) p(x) dx \quad (1.3)$$

corresponding to internal ( $\mathcal{U}$ ), potential ( $\mathcal{V}$ ) and interaction ( $\mathcal{W}$ ) energy, respectively. Defining the *relative entropy dissipation functional*

$$\mathcal{D}_{\text{ac}}(\mathbb{R}^n) \ni p \mapsto \mathcal{D}(p) := \int_{\mathbb{R}^n} |\nabla \log p(x) + \nabla V(x) + \nabla (W * p)(x)|^2 p(x) dx, \quad (1.4)$$

the well-known relative entropy dissipation identity takes the form

$$\mathcal{F}(p_t) - \mathcal{F}(p_{t_0}) = - \int_{t_0}^t \mathcal{D}(p_u) du. \quad (1.5)$$

This identity is of a deterministic nature: it only depends on the curve of probability density functions  $(p_t)_{0 \leq t \leq T}$ , but not on the trajectories of the underlying process  $(X_t)_{0 \leq t \leq T}$  itself. It is then natural to ask whether there is a process-level analogue of the relative entropy dissipation identity (1.5), depending directly on the trajectories of the McKean–Vlasov process  $(X_t)_{0 \leq t \leq T}$ . The main contribution of this paper is to give an affirmative answer to this question, by formulating a trajectorial version of the relative entropy dissipation identity via a stochastic analysis approach.

Before going into details, let us briefly describe the main ideas. We draw inspiration from prior literature [29,40] based on a simpler (linear) setting without interaction, i.e.,  $W \equiv 0$ . In this case, the McKean–Vlasov SDE (1.1) reduces to a Langevin–Smoluchowski diffusion equation of the form

$$dX_t = -\nabla V(X_t) dt + \sqrt{2} dB_t, \quad 0 \leq t \leq T. \quad (1.6)$$

In particular, the drift term does not depend on the distribution of  $X_t$ . Moreover, there is an explicit stationary distribution (also known as the Gibbs distribution [30,35,52]) with density proportional to  $\mathbb{R}^n \ni x \mapsto q(x) := e^{-V(x)}$ . Defining the likelihood ratio function (or Radon–Nikodym derivative)  $\ell_t(x) := p_t(x)/q(x)$ , the free energy at time  $t$  can be expressed as  $\mathcal{F}(p_t) = \mathbb{E}[\log \ell_t(X_t)]$ , and the resulting stochastic process

$$\log \ell_t(X_t) = \log p_t(X_t) + V(X_t), \quad 0 \leq t \leq T \quad (1.7)$$

is called *free energy* or *relative entropy process*. As shown in [29,40], the time-reversal

$$(\log \ell_{T-s}(X_{T-s}))_{0 \leq s \leq T}$$

of this process is a submartingale, and Itô calculus can be used to obtain its Doob–Meyer decomposition

$$\log \ell_{T-s}(X_{T-s}) - \log \ell_T(X_T) = M_{T-s} + F_{T-s}. \quad (1.8)$$

Here,  $(M_{T-s})_{0 \leq s \leq T}$  is a martingale and  $(F_{T-s})_{0 \leq s \leq T}$  is an increasing process of finite first variation, both with explicit expressions. This decomposition describes exactly the rate of relative entropy dissipation along every trajectory of the Langevin–Smoluchowski diffusion. Therefore, it can be viewed as a trajectorial analogue of the (deterministic) relative entropy dissipation identity (1.5).

Let us now return to our McKean–Vlasov setting. In order to take into account the interaction potential  $W$ , it is natural to consider a generalized relative entropy process of the form

$$\log p_t(X_t) + V(X_t) + \frac{1}{2}(W * p_t)(X_t), \quad 0 \leq t \leq T. \quad (1.9)$$

The task is now to compute the semimartingale decomposition of this process. We will provide a detailed analysis of this extension, which is subtler than might appear at first sight. The main difficulty is that, even when it exists, the stationary distribution of the McKean–Vlasov diffusion does not have a closed-form expression and is not even unique in general; see the works [6,20,21,33,54–57]. This prevents us from defining the likelihood ratio function in a straightforward manner as in the setting of Langevin–Smoluchowski diffusions, where one can rely on the invariant Gibbs distribution. An appropriate definition of the generalized likelihood ratio function turns out to be that (1.9) should be viewed as a function of the form  $\log \ell_t(X_t, P_t)$ , depending explicitly on the distribution  $P_t$  of  $X_t$  itself, in addition to the state  $X_t$ . This form of generalized likelihood ratio function allows us to take the  $L$ -derivative with respect to the probability distribution  $P_t$ . The notion of  $L$ -differentiation for functions of probability measures was introduced by Lions [43]. We refer to the monograph [10, Chapter 5] for a detailed discussion of differential calculus and stochastic analysis over spaces of probability measures. In particular, we will use a generalized form of Itô’s formula for functions of curves of measures, to derive the dynamics of the time-reversal of the relative entropy process (1.9), in terms of the semimartingale decomposition

$$\log \ell_{T-s}(X_{T-s}, P_{T-s}) - \log \ell_T(X_T, P_T) = M_{T-s} + F_{T-s}, \quad 0 \leq s \leq T, \quad (1.10)$$

where  $(M_{T-s})_{0 \leq s \leq T}$  is a martingale and  $(F_{T-s})_{0 \leq s \leq T}$  is a process of finite first variation, both of which will be explicitly computed. Similar to the case of Langevin–Smoluchowski dynamics, this decomposition can be viewed as the trajectorial rate of relative entropy dissipation. The classical (deterministic) identity (1.5) can then be recovered by taking expectations.

## 1.1. Gradient flow structure of the granular media equation

As a first application of our trajectorial approach we obtain a new interpretation of the gradient flow structure of the *granular media equation*

$$\partial_t p_t(x) = \operatorname{div} \left( \nabla p_t(x) + p_t(x) \nabla V(x) + p_t(x) \nabla (W * p_t)(x) \right), \quad (t, x) \in (0, T) \times \mathbb{R}^n, \quad (1.11)$$

which describes the evolution of the curve of probability density functions  $(p_t)_{0 \leq t \leq T}$  corresponding to the McKean–Vlasov diffusion  $(X_t)_{0 \leq t \leq T}$  of (1.1). When  $n = 1$ , this PDE appears in the modeling of the time evolution of granular media [7,17,59]; in that context, the granular medium is modeled as system of particles performing inelastic collisions, and  $p_t(x)$  is regarded as the velocity of a representative particle in the system at time  $t$  and position  $x$ , while  $V$  and  $W$  represent the friction and the inelastic collision forces, respectively. Note that in the interaction-free case  $W \equiv 0$ , the equation (1.11) reduces to a linear *Fokker–Planck equation*. As is well known from [13,14], this curve of probability densities can be characterized as a *gradient flow* in  $\mathcal{P}_{ac,2}(\mathbb{R}^n)$ , the space of absolutely continuous probability measures with finite second moments. Roughly speaking, this is an optimality property stating that the curve  $(p_t)_{0 \leq t \leq T}$  evolves in the direction of steepest possible descent for the free energy functional (1.2) with respect to the quadratic Wasserstein distance

$$W_2(\mu, \nu) = \left( \inf_{Y \sim \mu, Z \sim \nu} \mathbb{E}|Y - Z|^2 \right)^{1/2}, \quad \mu, \nu \in \mathcal{P}_2(\mathbb{R}^n). \quad (1.12)$$

The Wasserstein gradient flow structure of the linear Fokker–Planck equation was first discovered by Jordan, Kinderlehrer and Otto in the seminal work [36]. In the paper [50], Otto and Villani developed a formal Riemannian structure on the space of probability measures with finite second moments, leading to heuristic proofs of gradient flow properties as in [49], where the porous medium equation was

studied. This pioneering approach is often referred to as “Otto calculus”. Later, a rigorous framework based on minimizing movement schemes and curves of maximal slope was introduced in [3]. Recently, a trajectorial approach to the gradient flow properties of Langevin–Smoluchowski diffusions [40] and Markov chains [38] was established. We will follow this approach and adapt it to our McKean–Vlasov setting. For gradient flows of McKean–Vlasov equations on discrete spaces we refer to [24].

Returning to the setting of this paper, our main result leads to a new formulation of the gradient flow property of the granular media equation. To show this steepest descent property, the main idea is to consider a *perturbed* McKean–Vlasov diffusion of the form

$$dX_t = -(\nabla V(X_t) + \nabla \beta(X_t) + \nabla(W * P_t^\beta)(X_t)) dt + \sqrt{2} dB_t^\beta, \quad t_0 \leq t \leq T \quad (1.13)$$

which is constructed by adding a perturbation  $\beta: \mathbb{R}^n \rightarrow \mathbb{R}$  to the confinement potential<sup>1</sup> of the original McKean–Vlasov SDE (1.1). In other words, from time  $t_0$  onward, the perturbed diffusion drifts in a direction different from that of the original diffusion, hence the perturbed curve of time-marginal distributions  $(P_t^\beta)_{t_0 \leq t \leq T}$  also evolves differently from the unperturbed curve  $(P_t)_{t_0 \leq t \leq T}$ . In parallel with the unperturbed case, we may compute the dynamics of the perturbed relative entropy process associated with (1.13). As a consequence, we derive the rate of relative entropy dissipation for the perturbed McKean–Vlasov diffusion. On the other hand, the rate of change of the Wasserstein distance along the perturbed curve  $(P_t^\beta)_{t_0 \leq t \leq T}$  can be computed based on the general theory of metric derivative of absolutely continuous curves, see [3]. Finally, comparing these two rates in both the perturbed and unperturbed settings, allows us to establish the gradient flow property.

## 1.2. The HWBI inequality

The second application of our trajectorial approach deals with the HWBI inequality [2, Theorem 4.2], which is an extension of the HWI inequality [50]. It relates not only relative entropy (H), Wasserstein distance (W), and relative Fisher information (I), but also barycenter (B). These quantities are defined as follows: for two probability measures  $\nu, \mu \in \mathcal{P}(\mathbb{R}^n)$ , the *relative entropy* of  $\nu$  with respect to  $\mu$  is defined by

$$H(\nu | \mu) := \begin{cases} \int_{\mathbb{R}^n} \frac{d\nu}{d\mu} \log \left( \frac{d\nu}{d\mu} \right) d\mu, & \text{if } \nu \ll \mu \\ +\infty, & \text{otherwise,} \end{cases} \quad (1.14)$$

the *relative Fisher information* of  $\nu$  with respect to  $\mu$  is given by

$$I(\nu | \mu) := \begin{cases} \int_{\mathbb{R}^n} \left| \nabla \log \left( \frac{d\nu}{d\mu} \right) \right|^2 d\mu, & \text{if } \nu \ll \mu \\ +\infty, & \text{otherwise,} \end{cases} \quad (1.15)$$

and the *barycenter* of a probability measure  $\nu \in \mathcal{P}_2(\mathbb{R}^n)$  is defined as  $b(\nu) := \int_{\mathbb{R}^n} x d\nu(x) \in \mathbb{R}^n$ , where the integral is understood as a Bochner integral. Informally, the HWBI inequality then states that any two probability measures  $\nu_0, \nu_1 \in \mathcal{P}_2(\mathbb{R}^n)$  satisfy

$$H(\nu_0 | \mu_0) - H(\nu_1 | \mu_1) \leq \sqrt{I(\nu_0 | \mu_0^\uparrow)} W_2(\nu_0, \nu_1) - \frac{\kappa_V + \kappa_W}{2} W_2^2(\nu_0, \nu_1) + \frac{\kappa_W}{2} |b(\nu_0) - b(\nu_1)|^2, \quad (1.16)$$

<sup>1</sup>As we will see, the steepest descent property is already visible by perturbing the confinement potential from  $V$  to  $V + \beta$ , thus we avoid complicating the setup further by adding another perturbation to the interaction potential  $W$ .

where  $\mu_0, \mu_1, \mu_0^\uparrow$  are some appropriate  $\sigma$ -finite reference measures depending on the potentials  $V, W$  (see [Subsection 3.3](#) for the details), and  $\kappa_V, \kappa_W \in \mathbb{R}$  are the moduli of uniform convexity for  $V, W$ . This inequality describes the evolution of the relative entropy along the *displacement interpolation*  $(\nu_t)_{0 \leq t \leq 1}$  between  $\nu_0$  and  $\nu_1$ . Compared with the HWI inequality, there are two additional terms on the right-hand side of (1.16) contributed by the interaction energy functional  $\mathcal{W}$  of (1.3). Intuitively, the  $\kappa_W$ -uniform convexity of  $W$  leads to the first additional term  $-\frac{\kappa_W}{2} W_2^2(\nu_0, \nu_1)$ , which alone would correspond to the  $\kappa_W$ -uniform displacement convexity of  $\mathcal{W}$  along  $(\nu_t)_{0 \leq t \leq 1}$ . But since  $\mathcal{W}(p)$  is invariant under any translation of  $p$ , the functional  $\mathcal{W}$  might fail to be uniformly displacement convex when the barycenter shifts. This suggests that the barycentric shift along  $(\nu_t)_{0 \leq t \leq 1}$  should be factored out of the consideration of the displacement convexity of  $\mathcal{W}$ , which is intuitively why the second additional term  $\frac{\kappa_W}{2} |b(\nu_0) - b(\nu_1)|^2$  in (1.16) appears.

Coming back to our second application, we illustrate how our approach yields a trajectorial proof of the inequality (1.16), in the slightly strengthened form of [18, Theorem 4.1] and [25, Theorem D.50]. Much of this consists of arguments similar in spirit to our main result (1.10), but with one key difference: instead of the time-marginals of the McKean–Vlasov diffusion, we apply the trajectorial approach to the displacement interpolation  $(\nu_t)_{0 \leq t \leq 1}$ . In this regard, our derivation can be seen as a generalization of the trajectorial proof of the HWI inequality in [40, Section 4.2]; see also [38, Section 9.4], where the same idea was used to derive a discrete version of the HWI inequality in a Riemannian-geometric framework. Let us also point out that for the proof of the HWBI inequality we shall impose convexity assumptions (see [Assumptions 3.18](#)) on the potentials  $V, W$ . We do not require these assumptions in the rest of the paper.

In the literature, similar trajectorial approaches have also been applied in the context of martingale inequalities [1, 5], functional inequalities [4, 15, 31, 41], and their stability estimates [22, 23]. In particular, we refer to [4, Corollary 1.4] for a related HWI inequality derived from the entropic interpolation of the mean-field Schrödinger problem.

### 1.3. Organization of the paper

We set up the probabilistic framework and discuss some regularity assumptions in [Section 2](#). In [Section 3](#) we state our main trajectorial results, [Theorem 3.1](#) and [Theorem 3.9](#), and develop two explicit examples for illustration. As immediate consequences, we derive the classical relative entropy dissipation identities in [Corollary 3.4](#) and [Corollary 3.10](#). Building on these results, we formulate the gradient flow property of the granular media equation in [Theorem 3.15](#). The HWBI inequality is then stated in [Theorem 3.19](#). The proofs of the trajectorial results and of the HWBI inequality are developed in [Section 4](#). Some proofs of auxiliary results postponed in previous sections are contained in [Section 5](#).

## 2. The probabilistic framework

### 2.1. The setting

We fix a terminal time  $T \in (0, \infty)$  and let  $\Omega := C([0, T]; \mathbb{R}^n)$  be the path space of  $\mathbb{R}^n$ -valued continuous functions defined on  $[0, T]$ . We denote by  $(X_t)_{0 \leq t \leq T}$  the canonical process defined by  $X_t(\omega) := \omega(t)$  for  $\omega \in \Omega$ , and fix a probability distribution  $P_0 \in \mathcal{P}_{ac, 2}(\mathbb{R}^n)$ .

As will be shown in [Lemma 2.2](#), under the [Assumptions 2.1](#) below, the SDE (1.1) with initial distribution  $P_0$  has a unique strong solution, when it is posed on an arbitrary filtered probability space. This implies that there exists a probability measure  $\mathbb{P}$  on  $\Omega$  and a  $\mathbb{P}$ -Brownian motion  $(B_t)_{0 \leq t \leq T}$  such that

the SDE (1.1) holds. We write  $\mathbb{F} = (\mathcal{F}_t)_{0 \leq t \leq T}$  for the right-continuous augmentation of the canonical filtration.

For each time  $t \in [0, T]$ , we denote by  $P_t := \mathbb{P} \circ X_t^{-1}$  the distribution of  $X_t$  under  $\mathbb{P}$ , and by  $p_t$  the corresponding probability density function on  $\mathbb{R}^n$ . The density functions  $(p_t)_{0 \leq t \leq T}$  then solve the granular media equation (1.11).

## 2.2. Regularity assumptions

**Assumptions 2.1.** The following regularity assumptions will be used frequently.

- (i) The functions  $V, W: \mathbb{R}^n \rightarrow [0, \infty)$  are smooth and have Lipschitz continuous gradients with Lipschitz constants  $\|\nabla V\|_{\text{Lip}}, \|\nabla W\|_{\text{Lip}}$ . All derivatives of  $V$  and  $W$  grow at most exponentially as  $|x|$  tends to infinity, and the first derivatives are of linear growth. The latter condition means that there exists a constant  $C > 0$  such that

$$\forall x \in \mathbb{R}^n: \quad |\nabla V(x)| \leq C(1 + |x|), \quad |\nabla W(x)| \leq C(1 + |x|). \quad (2.1)$$

Furthermore, the function  $W$  is even (in other words, symmetric), i.e.,  $W(x) = W(-x)$  for all  $x \in \mathbb{R}^n$ .

- (ii) The probability distribution  $P_0$  is an element of the space  $\mathcal{P}_{\text{ac},2}(\mathbb{R}^n)$  and the corresponding probability density function  $\mathbb{R}^n \ni x \mapsto p_0(x)$  is strictly positive. Moreover, the initial free energy  $\mathcal{F}(p_0)$  is finite.

These assumptions ensure that the equation (1.1) belongs to a broad class of strongly solvable McKean–Vlasov SDEs. We relegate the proof of the following result to [Subsection 5.1](#).

**Lemma 2.2.** *Suppose [Assumptions 2.1](#) hold. Then on an arbitrary filtered probability space, the McKean–Vlasov SDE (1.1) has a pathwise unique, strong solution  $(X_t)_{0 \leq t \leq T}$  satisfying*

$$\mathbb{E} \left[ \sup_{0 \leq t \leq T} |X_t|^2 \right] < \infty. \quad (2.2)$$

*Moreover, its marginal distributions  $(P_t)_{0 \leq t \leq T}$  belong to  $\mathcal{P}_{\text{ac},2}(\mathbb{R}^n)$ , and the corresponding curve of probability density functions  $(p_t)_{0 \leq t \leq T}$  is a classical solution of the granular media equation (1.11).*

## 2.3. Probabilistic representations of gradient flow functionals

To set up our framework, the first step is to express the free energy as well as the relative entropy dissipation functional in probabilistic terms. To this end, we introduce the generalized potential  $\Psi: \mathbb{R}^n \times \mathcal{P}_2(\mathbb{R}^n) \rightarrow [0, \infty)$  and its close relative  $\Psi^\uparrow$  given by

$$\Psi(x, \mu) := V(x) + \frac{1}{2}(W * \mu)(x), \quad \Psi^\uparrow(x, \mu) := V(x) + (W * \mu)(x) \quad (2.3)$$

for  $(x, \mu) \in \mathbb{R}^n \times \mathcal{P}_2(\mathbb{R}^n)$ . Furthermore, we define the density functions

$$q(x, \mu) := e^{-\Psi(x, \mu)}, \quad q^\uparrow(x, \mu) := e^{-\Psi^\uparrow(x, \mu)}, \quad q^\downarrow(x) := e^{-V(x)} \quad (2.4)$$

and the corresponding generalized likelihood ratio functions

$$\ell_t(x, \mu) := \frac{p_t(x)}{q(x, \mu)}, \quad \ell_t^\uparrow(x, \mu) := \frac{p_t(x)}{q^\uparrow(x, \mu)}, \quad \ell_t^\downarrow(x) := \frac{p_t(x)}{q^\downarrow(x)} \quad (2.5)$$

for  $t \in [0, T]$ . Note that if  $W \equiv 0$ , these three likelihood ratio functions coincide.

For each time  $t \in [0, T]$ , we introduce  $\sigma$ -finite measures on the Borel sets of  $\mathbb{R}^n$ , given by

$$Q_t(A) := \int_A q(x, P_t) dx, \quad Q_t^\uparrow(A) := \int_A q^\uparrow(x, P_t) dx, \quad A \in \mathcal{B}(\mathbb{R}^n). \quad (2.6)$$

Intuitively, these measures are (unnormalized) time-dependent Gibbs distributions. If  $W \equiv 0$ , they coincide with the true Gibbs distribution of the Langevin–Smoluchowski equation (1.6), which is also its stationary distribution (when normalized to a probability measure).

With these definitions, we can now write the gradient flow functionals  $\mathcal{F}$  and  $\mathcal{D}$ , introduced in (1.2) and (1.4), in probabilistic terms: the relative entropy (defined in (1.14)) of  $P_t$  with respect to  $Q_t$  and the relative Fisher information (defined in (1.15)) of  $P_t$  with respect to  $Q_t^\uparrow$  can be expressed respectively as

$$H(P_t | Q_t) = \mathbb{E}_{\mathbb{P}} [\log \ell_t(X_t, P_t)], \quad I(P_t | Q_t^\uparrow) = \mathbb{E}_{\mathbb{P}} [|\nabla \log \ell_t^\uparrow(X_t, P_t)|^2]; \quad (2.7)$$

and we have the relations  $H(P_t | Q_t) = \mathcal{F}(P_t)$  as well as  $I(P_t | Q_t^\uparrow) = \mathcal{D}(P_t)$ . In particular, the relative entropy  $H(P_t | Q_t)$  can be written as the  $\mathbb{P}$ -expectation of the *relative entropy process*

$$\log \ell_t(X_t, P_t) = \log p_t(X_t) + V(X_t) + \frac{1}{2}(W * P_t)(X_t), \quad 0 \leq t \leq T. \quad (2.8)$$

The dynamics of this stochastic process, together with its perturbed counterpart to be introduced in Subsection 3.2 below, will be our main objects of interest.

**Remark 2.3.** If the reference measure  $Q_t$  in (2.6) is a *probability* measure, then the expression (2.7) matches the classical definition of relative entropy given in (1.14). In the general case when  $Q_t$  is a  $\sigma$ -finite measure, the definition (1.14) is also valid under the condition that  $P_t$  has finite second moment, with the only difference that the range of the function  $t \mapsto H(P_t | Q_t)$  is extended from  $[0, \infty]$  to  $(-\infty, \infty]$ ; we refer to [39, Appendix C] or [19, Section 3] for the details.

## 3. Main results

### 3.1. Trajectorial dissipation of relative entropy for McKean–Vlasov diffusions

Our first main result is the semimartingale decomposition of the relative entropy process (2.8). It describes the dissipation of relative entropy along every trajectory of a particle undergoing the McKean–Vlasov dynamics (1.1). In the same spirit as the trajectorial approaches of [29] and [40], we shall study the dynamics of the relative entropy process in the backward direction of time. Concretely, we consider for arbitrary, fixed  $T \in (0, \infty)$  the time-reversed canonical process

$$\bar{X}_s := X_{T-s}, \quad 0 \leq s \leq T \quad (3.1)$$

on the filtered probability space  $(\Omega, \mathbb{G}, \mathbb{P})$ , where  $\mathbb{G} = (\mathcal{G}_s)_{0 \leq s \leq T}$  is the  $\mathbb{P}$ -augmented filtration generated by  $(\bar{X}_s)_{0 \leq s \leq T}$ .

In order to formulate [Theorem 3.1](#) below, we introduce the time-reversed *Fisher information process*

$$\begin{aligned} \bar{I}_s := & \left( |\nabla \log \bar{\ell}_s^\downarrow|^2 + \frac{1}{2} |\nabla(W * \bar{P}_s)|^2 + \left\langle \frac{1}{2} \nabla(W * \bar{P}_s), 2\nabla \log \bar{\ell}_s^\downarrow + \nabla V \right\rangle \right) (\bar{X}_s) \\ & - \mathbb{E}_{\bar{\mathbb{P}}} \left[ \left\langle \frac{1}{2} \nabla W(\bar{X}_s - \bar{Y}_s), \left( 2\nabla \log \bar{\ell}_s^\downarrow - \nabla V + \nabla(W * \bar{P}_s) \right) (\bar{Y}_s) \right\rangle \right] \end{aligned} \quad (3.2)$$

for  $0 \leq s \leq T$ . Here, the process  $(\bar{Y}_s)_{0 \leq s \leq T}$  is defined on another probability space  $(\tilde{\Omega}, \tilde{\mathbb{G}}, \tilde{\mathbb{P}})$  such that the tuple  $(\tilde{\Omega}, \tilde{\mathbb{G}}, \tilde{\mathbb{P}}, (\bar{Y}_s)_{0 \leq s \leq T})$  is an exact copy of  $(\Omega, \mathbb{G}, \mathbb{P}, (\bar{X}_s)_{0 \leq s \leq T})$ . A bar over a letter means that time is reversed as in (3.1).

We also define the time-reversed *cumulative Fisher information process* as the time integral

$$\bar{F}_s := \int_0^s \bar{I}_u \, du, \quad 0 \leq s \leq T. \quad (3.3)$$

This process will act as the compensator in the semimartingale decomposition of the relative entropy process (2.8). Its relation with the relative Fisher information (2.7) will be given in (3.7) below.

**Theorem 3.1.** Suppose [Assumptions 2.1](#) hold. On  $(\Omega, \mathbb{G}, \mathbb{P})$ , the time-reversed relative entropy process

$$\log \bar{\ell}_s(\bar{X}_s, \bar{P}_s) = \log \ell_{T-s}(X_{T-s}, P_{T-s}), \quad 0 \leq s \leq T \quad (3.4)$$

admits the semimartingale decomposition

$$\log \bar{\ell}_s(\bar{X}_s, \bar{P}_s) - \log \bar{\ell}_0(\bar{X}_0, \bar{P}_0) = \bar{M}_s + \bar{F}_s. \quad (3.5)$$

Here  $(\bar{M}_s)_{0 \leq s \leq T}$  is the  $L^2(\mathbb{P})$ -bounded martingale

$$\bar{M}_s := \int_0^s \left\langle \nabla \log \bar{\ell}_u(\bar{X}_u, \bar{P}_u), \sqrt{2} \, d\bar{B}_u \right\rangle, \quad (3.6)$$

with  $(\bar{B}_s)_{0 \leq s \leq T}$  a  $\mathbb{P}$ -Brownian motion of the backward filtration  $\mathbb{G}$ , and the compensator  $(\bar{F}_s)_{0 \leq s \leq T}$  satisfies

$$\mathbb{E}_{\mathbb{P}}[\bar{F}_s] = \int_0^s I(\bar{P}_u | \bar{Q}_u^\uparrow) \, du = \mathbb{E}_{\mathbb{P}} \left[ \int_0^s |\nabla \log \bar{\ell}_u^\uparrow(\bar{X}_u, \bar{P}_u)|^2 \, du \right] < \infty. \quad (3.7)$$

### 3.1.1. Examples

We give two concrete examples to illustrate [Theorem 3.1](#).

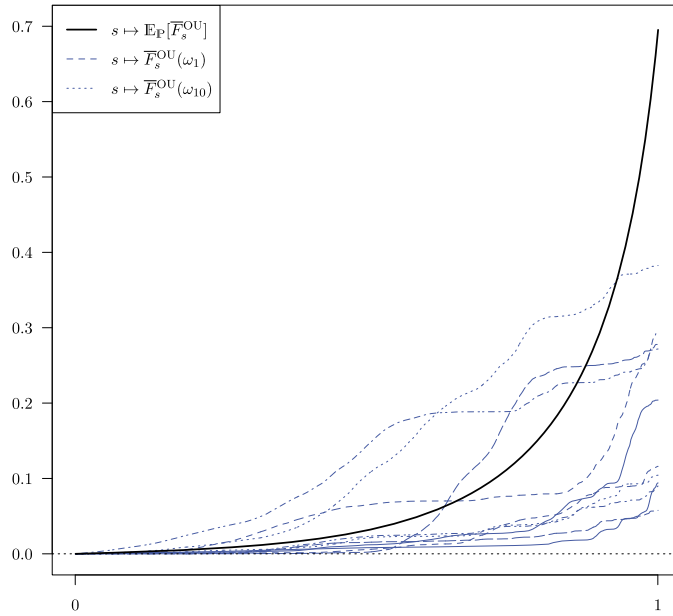
**Example 3.2.** We set  $n = 1$  and specialize [Theorem 3.1](#) to the case of quadratic confinement potential  $V(x) = \frac{x^2}{2}$  and no interaction potential  $W \equiv 0$ . The initial position  $X_0$  is chosen to be independent of  $(B_t)_{0 \leq t \leq T}$  and to be normally distributed with mean 0 and variance  $\sigma_0^2 > 0$ . In this case, the SDE of (1.1) becomes

$$dX_t = -X_t \, dt + \sqrt{2} \, dB_t, \quad 0 \leq t \leq T \quad (3.8)$$

and its solution is given by the Ornstein–Uhlenbeck process

$$X_t = e^{-t} X_0 + \sqrt{2} \int_0^t e^{u-t} \, dB_u, \quad 0 \leq t \leq T \quad (3.9)$$





**Figure 1.** Simulations of the cumulative Fisher information process (3.11) for the Ornstein–Uhlenbeck diffusion (3.8).

with probability density function

$$p_t(x) = \frac{1}{\sqrt{2\pi\sigma_t^2}} \exp\left(-\frac{x^2}{2\sigma_t^2}\right), \quad \sigma_t^2 := 1 + e^{-2t}(\sigma_0^2 - 1). \quad (3.10)$$

Recalling (2.3) – (2.5) and using (3.10), we see that in this setting the cumulative Fisher information process of (3.3) is explicitly given by

$$\bar{F}_s^{\text{OU}} = \int_0^s \left( \nabla \log \bar{\ell}_u^{\downarrow}(\bar{X}_u) \right)^2 du = \int_0^s \left( \nabla \log \bar{p}_u(\bar{X}_u) + \bar{X}_u \right)^2 du = \int_0^s \left( 1 - \frac{1}{\sigma_u^2} \right)^2 \bar{X}_u^2 du \quad (3.11)$$

for  $0 \leq s \leq T$ . In particular, the non-negativity of the integrand in (3.11) implies that the relative entropy decreases along almost every trajectory.

Now we set  $T = 1$  and  $\sigma_0^2 = 0.1$ . The blue lines in Figure 1 represent ten simulated trajectories  $s \mapsto \bar{F}_s^{\text{OU}}(\omega_i)$ , for  $i = 1, \dots, 10$ . The thick black line plots the expected path  $s \mapsto \mathbb{E}_P[\bar{F}_s^{\text{OU}}]$  of all possible trajectories.

**Example 3.3.** We set again  $n = 1$  and now consider the case of no confinement potential  $V \equiv 0$ , quadratic interaction potential  $W(x) = \frac{x^2}{2}$ , and a centered Gaussian initial position  $X_0$  with variance  $\sigma_0^2 > 0$ , which is independent of the Brownian motion  $(B_t)_{t \geq 0}$ . In this case, for any  $t \geq 0$ , the drift term of the SDE of (1.1) is

$$-\nabla(W * P_t)(X_t) = - \int_{\mathbb{R}^n} \nabla W(X_t - y) p_t(y) dy = -(X_t - \mathbb{E}[X_t]). \quad (3.12)$$

In particular, the drift term depends on the distribution  $P_t$  only through its mean. Substituting it into (1.1), this SDE reduces to

$$dX_t = -(X_t - \mathbb{E}[X_t]) dt + \sqrt{2} dB_t, \quad 0 \leq t \leq T. \quad (3.13)$$

This type of nonlinear, self-interacting SDE has been studied since [6], where it was shown that its solution is also given by the Ornstein–Uhlenbeck process of (3.9). Therefore, similar computations as in Example 3.2 show that in this setting the cumulative Fisher information process is given by

$$\overline{F}_s^{\text{NL}} = \int_0^s \left( \left( \frac{1}{\overline{\sigma}_u^4} + \frac{1}{2} - \frac{1}{\overline{\sigma}_u^2} \right) \overline{X}_u^2 + \mathbb{E}_{\tilde{\mathbb{P}}} \left[ \frac{1}{2} (\overline{X}_u - \overline{Y}_u) \left( \frac{2}{\overline{\sigma}_u^2} \overline{Y}_u - \overline{Y}_u \right) \right] \right) du \quad (3.14)$$

for  $0 \leq s \leq T$ . Using the fact that  $(\overline{X}_u)_\#(\tilde{\mathbb{P}}) = (\overline{Y}_u)_\#(\tilde{\mathbb{P}}) = \mathcal{N}(0, \overline{\sigma}_u^2)$ , which we have from (3.10), we can compute the expectation appearing in (3.14) and obtain

$$\overline{F}_s^{\text{NL}} = \int_0^s \left( \left( \frac{1}{\overline{\sigma}_u^4} + \frac{1}{2} - \frac{1}{\overline{\sigma}_u^2} \right) \overline{X}_u^2 + \frac{\overline{\sigma}_u^2}{2} - 1 \right) du. \quad (3.15)$$

Clearly,  $\overline{F}_s^{\text{OU}} \neq \overline{F}_s^{\text{NL}}$ . In particular, the integrand in (3.15) is non-negative if and only if

$$\overline{X}_u^2 \geq \left( \frac{1}{\overline{\sigma}_u^4} + \frac{1}{2} - \frac{1}{\overline{\sigma}_u^2} \right)^{-1} \left( 1 - \frac{\overline{\sigma}_u^2}{2} \right).$$

In other words, as opposed to Example 3.2, relative entropy only decreases along a trajectory if  $\overline{X}_u$  is far from its mean. However, after taking expectations in (3.11) and (3.15), we see that the expected rate of relative entropy dissipation in both cases is equal to

$$\mathbb{E}_{\mathbb{P}}[\overline{F}_s^{\text{OU}}] = \mathbb{E}_{\mathbb{P}}[\overline{F}_s^{\text{NL}}] = \int_0^s \left( \overline{\sigma}_u - \frac{1}{\overline{\sigma}_u} \right)^2 du, \quad 0 \leq s \leq T. \quad (3.16)$$

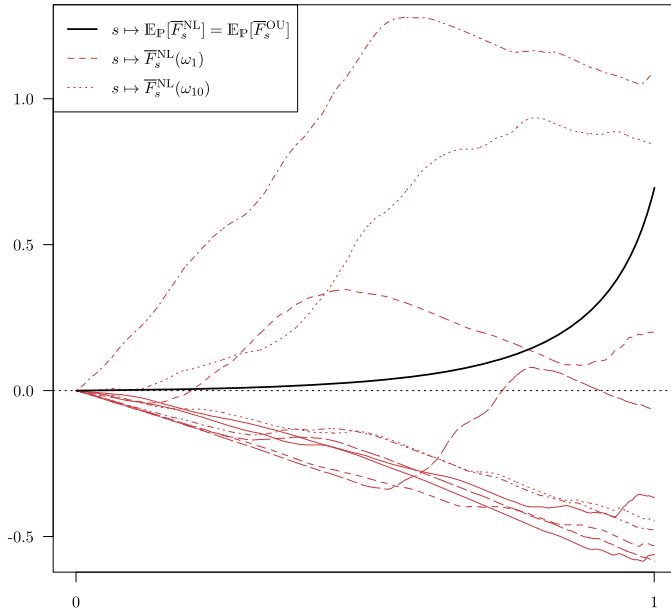
Now we set again  $T = 1$  and  $\sigma_0^2 = 0.1$ . In the same vein as in Figure 1, we plot in Figure 2 the paths of ten simulated trajectories  $s \mapsto \overline{F}_s^{\text{NL}}(\omega_i)$ , for  $i = 1, \dots, 10$ . We observe that some of the red lines describing the paths of these trajectories indeed take negative values. In other words, the cumulative Fisher information process of (3.15), and hence its integrand, can both be negative. Finally, the thick black line in Figure 2 follows the expected path  $s \mapsto \mathbb{E}_{\mathbb{P}}[\overline{F}_s^{\text{NL}}]$  of all possible trajectories. According to (3.16), this is the same black line as in Figure 1.

### 3.1.2. Consequences of Theorem 3.1

We now return to the general statement of Theorem 3.1 and deduce some direct consequences. By averaging the trajectorial result of Theorem 3.1 according to the path measure  $\mathbb{P}$ , we derive the well-known *relative entropy identity* (3.17) and the *dissipation of relative entropy* (3.18) below. A sketch of proof for the latter result was first given in [13, Proposition 2.1].

**Corollary 3.4.** *Suppose Assumptions 2.1 hold. For all  $t, t_0 \in [0, T]$ , we have the relative entropy identity*

$$H(P_t | Q_t) - H(P_{t_0} | Q_{t_0}) = - \int_{t_0}^t I(P_u | Q_u^\uparrow) du. \quad (3.17)$$



**Figure 2.** Simulations of the cumulative Fisher information process (3.15) for the nonlinear, self-interacting diffusion (3.13).

In particular, the relative entropy function  $t \mapsto H(\mathbb{P}_t | \mathbb{Q}_t)$  is monotonically decreasing. Furthermore, for Lebesgue-a.e.  $t \in [0, T]$ , the relative Fisher information  $I(\mathbb{P}_t | \mathbb{Q}_t^\uparrow)$  is finite, and the rate of relative entropy dissipation is given by

$$\frac{d}{dt} H(\mathbb{P}_t | \mathbb{Q}_t) = -I(\mathbb{P}_t | \mathbb{Q}_t^\uparrow). \quad (3.18)$$

**Proof.** The identity (3.17) follows by taking expectations with respect to the probability measure  $\mathbb{P}$  in (3.5), recalling the definitions of (2.7), using (3.7), and invoking the fact that the  $\mathbb{P}$ -expectation of the martingale (3.6) vanishes. Finally, applying the Lebesgue differentiation theorem to the monotone function  $t \mapsto H(\mathbb{P}_t | \mathbb{Q}_t)$  gives (3.18).  $\square$

**Remark 3.5.** The relation (3.18) describes the temporal dissipation of relative entropy at the *ensemble* level. It asserts that the rate of decay of the relative entropy  $t \mapsto H(\mathbb{P}_t | \mathbb{Q}_t)$  is given by the relative Fisher information  $I(\mathbb{P}_t | \mathbb{Q}_t^\uparrow)$ .

Finally, let us place ourselves again on the filtered probability space  $(\Omega, \mathbb{G}, \mathbb{P})$  as in Theorem 3.1 and emphasize that this trajectorial result is valid along almost every trajectory  $s \mapsto \bar{X}_s(\omega)$  of the underlying McKean–Vlasov process. As a consequence, by taking conditional expectations, we can generalize (3.18) and deduce the following trajectorial rate of relative entropy dissipation.

**Corollary 3.6.** Suppose Assumptions 2.1 hold and  $\int_0^T \mathbb{E}_{\mathbb{P}}[|\bar{I}_u|] du < \infty$ . For  $\mathbb{P}$ -a.e.  $\omega \in \Omega$  there exists a Lebesgue null set  $N_\omega \subseteq [0, T]$  such that for any  $s_0 \in [0, T] \setminus N_\omega$  we have

$$\lim_{s \downarrow s_0} \frac{\mathbb{E}_{\mathbb{P}}[\log \bar{\ell}_s(\bar{X}_s, \bar{\mathbb{P}}_s) | \mathcal{G}_{s_0}](\omega) - \log \bar{\ell}_{s_0}(\bar{X}_{s_0}(\omega), \bar{\mathbb{P}}_{s_0})}{s - s_0} = \bar{I}_{s_0}(\omega). \quad (3.19)$$

**Remark 3.7.** Recalling (3.7), we observe that  $\mathbb{E}_{\mathbb{P}}[\bar{I}_{s_0}] = I(\bar{\mathbb{P}}_{s_0} | \bar{Q}_{s_0}^\uparrow)$ . Therefore the limiting assertion (3.19) can indeed be viewed as a trajectorial version of the deterministic relative entropy dissipation identity (3.18).

**Proof.** We let  $0 \leq s_0 \leq s \leq T$ . By (3.5), (3.3) and Fubini's theorem, we have for  $\mathbb{P}$ -a.e.  $\omega \in \Omega$

$$\mathbb{E}_{\mathbb{P}}[\log \bar{\ell}_s(\bar{X}_s, \bar{\mathbb{P}}_s) | \mathcal{G}_{s_0}](\omega) - \log \bar{\ell}_{s_0}(\bar{X}_{s_0}(\omega), \bar{\mathbb{P}}_{s_0}) = \int_{s_0}^s \mathbb{E}_{\mathbb{P}}[\bar{I}_u | \mathcal{G}_{s_0}](\omega) du. \quad (3.20)$$

Furthermore, Jensen's inequality gives

$$\mathbb{E}_{\mathbb{P}}\left[\int_{s_0}^s |\mathbb{E}_{\mathbb{P}}[\bar{I}_u | \mathcal{G}_{s_0}](\omega)| du\right] \leq \int_{s_0}^s \mathbb{E}_{\mathbb{P}}[|\bar{I}_u|] du < \infty, \quad (3.21)$$

which implies

$$\int_{s_0}^s |\mathbb{E}_{\mathbb{P}}[\bar{I}_u | \mathcal{G}_{s_0}](\omega)| du < \infty \quad \text{for } \mathbb{P}\text{-a.e. } \omega \in \Omega. \quad (3.22)$$

By the Lebesgue differentiation theorem, for every such  $\omega$  there exists a Lebesgue null set  $N_\omega \subseteq [0, T]$  so that the limiting assertion

$$\lim_{s \downarrow s_0} \frac{\int_{s_0}^s \mathbb{E}_{\mathbb{P}}[\bar{I}_u | \mathcal{G}_{s_0}](\omega) du}{s - s_0} = \mathbb{E}_{\mathbb{P}}[\bar{I}_{s_0} | \mathcal{G}_{s_0}](\omega) = \bar{I}_{s_0}(\omega) \quad (3.23)$$

holds for every  $s_0 \in [0, T] \setminus N_\omega$ . Finally, combining (3.20) and (3.23) proves (3.19).  $\square$

### 3.2. Gradient flow structure of the granular media equation

In this subsection we apply the trajectorial approach of Subsection 3.1 in order to formulate the gradient flow property of the granular media equation (1.11). To this end, we consider a function  $\beta: \mathbb{R}^n \rightarrow \mathbb{R}$ , which will be treated as a perturbation potential. We denote by  $V^\beta := V + \beta$  the perturbed confinement potential and invoke the following regularity assumptions.

**Assumptions 3.8.** The function  $\beta: \mathbb{R}^n \rightarrow \mathbb{R}$  is smooth and compactly supported, and we require that Assumptions 2.1 are still satisfied if we replace  $V$  by  $V^\beta$ .

Note that Assumptions 3.8 allow us to apply Lemma 2.2 to the “perturbed” McKean–Vlasov SDE

$$dX_t = -(\nabla V^\beta(X_t) + \nabla(W * \mathbb{P}_t^\beta)(X_t)) dt + \sqrt{2} d\mathcal{B}_t^\beta, \quad t_0 \leq t \leq T \quad (3.24)$$

starting at time  $t_0 \geq 0$ , with  $X_{t_0}$  having initial distribution  $\mathbb{P}_{t_0}^\beta = \mathbb{P}_{t_0}$ . Therefore, by analogy with Subsection 2.1, we can construct a probability measure  $\mathbb{P}^\beta$  on  $\Omega := C([t_0, T]; \mathbb{R}^n)$ , under which the canonical process  $(X_t)_{t_0 \leq t \leq T}$  satisfies the SDE (3.24), with  $(\mathcal{B}_t^\beta)_{t_0 \leq t \leq T}$  being a  $\mathbb{P}^\beta$ -Brownian motion.

For each time  $t \in [t_0, T]$ , we denote by  $\mathbb{P}_t^\beta := \mathbb{P}^\beta \circ X_t^{-1}$  the probability distribution and by  $p_t^\beta$  the probability density function of  $X_t$  under  $\mathbb{P}^\beta$ . The “perturbed” curve of density functions  $(p_t^\beta)_{t_0 \leq t \leq T}$  then satisfies the *perturbed granular media equation*

$$\begin{cases} \partial_t p_t^\beta(x) = \operatorname{div} \left( \nabla p_t^\beta(x) + p_t^\beta(x) \nabla V^\beta(x) + p_t^\beta(x) \nabla(W * p_t^\beta)(x) \right), & (t, x) \in (t_0, T) \times \mathbb{R}^n, \\ p_{t_0}^\beta(x) = p_{t_0}(x), & x \in \mathbb{R}^n. \end{cases} \quad (3.25)$$

By analogy with (2.3), we define the perturbed potentials

$$\Psi^\beta(x, \mu) := V^\beta(x) + \frac{1}{2}(W * \mu)(x), \quad \Psi^{\beta\uparrow}(x, \mu) := V^\beta(x) + (W * \mu)(x), \quad \Psi^{\beta\downarrow} := V^\beta \quad (3.26)$$

for  $(x, \mu) \in \mathbb{R}^n \times \mathcal{P}_2(\mathbb{R}^n)$ . In parallel to (2.5), we introduce the *perturbed likelihood ratio functions*

$$\ell_t^\beta(x, \mu) := \frac{p_t^\beta(x)}{q(x, \mu)}, \quad \ell_t^{\beta\uparrow}(x, \mu) := \frac{p_t^\beta(x)}{q^\uparrow(x, \mu)}, \quad \ell_t^{\beta\downarrow}(x) := \frac{p_t^\beta(x)}{q^\downarrow(x)} \quad (3.27)$$

for  $t \in [t_0, T]$ . Finally, we define the  $\sigma$ -finite measures

$$Q_t^\beta(A) := \int_A q(x, P_t^\beta) dx, \quad Q_t^{\beta\uparrow}(A) := \int_A q^\uparrow(x, P_t^\beta) dx, \quad A \in \mathcal{B}(\mathbb{R}^n). \quad (3.28)$$

They are the perturbed versions of the measures  $Q_t$  and  $Q_t^\uparrow$  defined in (2.6). The relative entropy of  $P_t^\beta$  with respect to  $Q_t^\beta$  is then given by

$$H(P_t^\beta | Q_t^\beta) = \mathbb{E}_{\mathbb{P}^\beta} [\log \ell_t^\beta(X_t, P_t^\beta)] = \mathcal{F}(P_t^\beta) \quad (3.29)$$

and the relative Fisher information of  $P_t^\beta$  with respect to  $Q_t^{\beta\uparrow}$  equals

$$I(P_t^\beta | Q_t^{\beta\uparrow}) = \mathbb{E}_{\mathbb{P}^\beta} [|\nabla \log \ell_t^{\beta\uparrow}(X_t, P_t^\beta)|^2] = \mathcal{D}(P_t^\beta). \quad (3.30)$$

The following trajectorial result, [Theorem 3.9](#) below, provides the semimartingale decomposition of the *perturbed relative entropy process*

$$\log \ell_t^\beta(X_t, P_t^\beta) = \log p_t^\beta(X_t) + V(X_t) + \frac{1}{2}(W * P_t^\beta)(X_t), \quad t_0 \leq t \leq T. \quad (3.31)$$

In line with its unperturbed counterpart, [Theorem 3.1](#), we shall formulate this result in the reverse direction of time. We first introduce the perturbed analogues of (3.2) and (3.3): the *time-reversed perturbed Fisher information process* is defined as

$$\bar{I}_s^\beta := \left( |\nabla \log \bar{\ell}_s^{\beta\downarrow}|^2 + \frac{1}{2} |\nabla (W * \bar{P}_s^\beta)|^2 + \left\langle \frac{1}{2} \nabla (W * \bar{P}_s^\beta), 2 \nabla \log \bar{\ell}_s^{\beta\downarrow} + \nabla V^\beta \right\rangle \right) (\bar{X}_s) \quad (3.32)$$

$$- \mathbb{E}_{\mathbb{P}^\beta} \left[ \left\langle \frac{1}{2} \nabla W (\bar{X}_s - \bar{Y}_s), \left( 2 \nabla \log \bar{\ell}_s^{\beta\downarrow} - \nabla V + \nabla (W * \bar{P}_s^\beta) + \nabla \beta \right) (\bar{Y}_s) \right\rangle \right] \quad (3.33)$$

$$+ \left( \langle \nabla V, \nabla \beta \rangle - \Delta \beta \right) (\bar{X}_s) \quad (3.34)$$

for all  $0 \leq s \leq T - t_0$ , where  $(\bar{Y}_s)_{0 \leq s \leq T - t_0}$  is a copy of the process  $(\bar{X}_s)_{0 \leq s \leq T - t_0}$  on a copy  $(\tilde{\Omega}, \tilde{\mathbb{G}}, \tilde{\mathbb{P}}^\beta)$  of the original probability space  $(\Omega, \mathbb{G}, \mathbb{P}^\beta)$ ; the *perturbed cumulative Fisher information process* is defined as

$$\bar{F}_s^\beta := \int_0^s \bar{I}_u^\beta du, \quad 0 \leq s \leq T - t_0. \quad (3.35)$$

**Theorem 3.9.** Suppose [Assumptions 3.8](#) hold. On  $(\Omega, \mathbb{G}, \mathbb{P}^\beta)$ , the time-reversed perturbed relative entropy process

$$\log \bar{\ell}_s^\beta(\bar{X}_s, \bar{P}_s^\beta) = \log \ell_{T-s}^\beta(X_{T-s}, P_{T-s}^\beta), \quad 0 \leq s \leq T - t_0 \quad (3.36)$$

admits the semimartingale decomposition

$$\log \bar{\ell}_s^\beta(\bar{X}_s, \bar{\mathbb{P}}_s^\beta) - \log \bar{\ell}_0^\beta(\bar{X}_0, \bar{\mathbb{P}}_0^\beta) = \bar{M}_s^\beta + \bar{F}_s^\beta. \quad (3.37)$$

Here  $(\bar{M}_s^\beta)_{0 \leq s \leq T-t_0}$  is the  $L^2(\mathbb{P}^\beta)$ -bounded martingale

$$\bar{M}_s^\beta := \int_0^s \left\langle \nabla \log \bar{\ell}_u^\beta(\bar{X}_u, \bar{\mathbb{P}}_u^\beta), \sqrt{2} d\bar{B}_u^\beta \right\rangle, \quad (3.38)$$

with  $(\bar{B}_s^\beta)_{0 \leq s \leq T-t_0}$  a  $\mathbb{P}^\beta$ -Brownian motion of the backward filtration  $\mathbb{G}$ , and the compensator (3.35) satisfies

$$\mathbb{E}_{\mathbb{P}^\beta} [\bar{F}_s^\beta] = \int_0^s \left( I(\bar{\mathbb{P}}_u^\beta | \bar{\mathbb{Q}}_u^{\beta\uparrow}) + \mathbb{E}_{\mathbb{P}^\beta} \left[ \left( \left\langle \nabla V + \nabla(W * \bar{\mathbb{P}}_u^\beta), \nabla \beta \right\rangle - \Delta \beta \right) (\bar{X}_u) \right] \right) du < \infty. \quad (3.39)$$

With the dynamics of the time-reversed perturbed relative entropy process at hand, we repeat the same procedure which was carried out for the unperturbed case. Taking expectations with respect to the probability measure  $\mathbb{P}^\beta$ , we arrive at the perturbed relative entropy identity (3.40), and applying the Lebesgue differentiation theorem gives the perturbed relative entropy production identity (3.41).

**Corollary 3.10.** *Suppose Assumptions 3.8 hold. For all  $0 \leq t_0 \leq t \leq T$ , we have the perturbed relative entropy identity*

$$\begin{aligned} H(\mathbb{P}_t^\beta | \mathbb{Q}_t^\beta) - H(\mathbb{P}_{t_0}^\beta | \mathbb{Q}_{t_0}^\beta) \\ = - \int_{t_0}^t \left( I(\mathbb{P}_u^\beta | \mathbb{Q}_u^{\beta\uparrow}) + \mathbb{E}_{\mathbb{P}^\beta} \left[ \left( \left\langle \nabla V + \nabla(W * \mathbb{P}_u^\beta), \nabla \beta \right\rangle - \Delta \beta \right) (X_u) \right] \right) du. \end{aligned} \quad (3.40)$$

For Lebesgue-a.e.  $t_0 \in [0, T]$ , the perturbed rate of relative entropy dissipation is given by

$$\left. \frac{d}{dt} \right|_{t=t_0}^+ H(\mathbb{P}_t^\beta | \mathbb{Q}_t^\beta) = - \left( I(\mathbb{P}_{t_0}^\beta | \mathbb{Q}_{t_0}^\beta) + \mathbb{E}_{\mathbb{P}^\beta} \left[ \left( \left\langle \nabla V + \nabla(W * \mathbb{P}_{t_0}^\beta), \nabla \beta \right\rangle - \Delta \beta \right) (X_{t_0}) \right] \right). \quad (3.41)$$

Similarly, we have the following trajectorial rate of relative entropy dissipation for the perturbed diffusion.

**Corollary 3.11.** *Suppose Assumptions 3.8 hold and  $\int_0^{T-t_0} \mathbb{E}_{\mathbb{P}^\beta} [|\bar{I}_u^\beta|] du < \infty$ . For  $\mathbb{P}^\beta$ -a.e.  $\omega \in \Omega$  there exists a Lebesgue null set  $N_\omega^\beta \subseteq [0, T - t_0]$  such that for any  $s_0 \in [0, T - t_0] \setminus N_\omega^\beta$  we have*

$$\lim_{s \downarrow s_0} \frac{\mathbb{E}_{\mathbb{P}^\beta} [\log \bar{\ell}_s^\beta(\bar{X}_s, \bar{\mathbb{P}}_s^\beta) | \mathcal{G}_{s_0}] (\omega) - \log \bar{\ell}_{s_0}^\beta(\bar{X}_{s_0}(\omega), \bar{\mathbb{P}}_{s_0}^\beta)}{s - s_0} = \bar{I}_{s_0}^\beta(\omega). \quad (3.42)$$

**Proof.** The proof proceeds almost verbatim as the proof of Corollary 3.6. The only difference is that we now use the semimartingale decomposition (3.37) and the  $\mathbb{P}^\beta$ -martingale property of the process (3.38) in Theorem 3.9.  $\square$

We now turn to the computation of the rate of change of the Wasserstein distance along the curve of probability distributions  $(\mathbb{P}_t^\beta)_{t_0 \leq t \leq T}$ . To this end, we set

$$v_t^\beta(x) := -(\nabla \log p_t^\beta + \nabla V^\beta + \nabla(W * p_t^\beta))(x), \quad (t, x) \in [t_0, T] \times \mathbb{R}^n, \quad (3.43)$$

so that the perturbed granular media equation (3.25) can be viewed as a *continuity equation*

$$\partial_t p_t^\beta(x) + \operatorname{div}(v_t^\beta(x) p_t^\beta(x)) = 0, \quad (t, x) \in (t_0, T) \times \mathbb{R}^n, \quad (3.44)$$

with  $v_t^\beta(\cdot)$  as the corresponding *velocity field*. We recall the definition of the tangent space (see Definition 8.4.1 in [3])

$$\operatorname{Tan}_\mu \mathcal{P}_2(\mathbb{R}^n) := \overline{\{\nabla \varphi : \varphi \in C_c^\infty(\mathbb{R}^n; \mathbb{R})\}}^{L^2(\mu)} \quad (3.45)$$

of  $\mathcal{P}_2(\mathbb{R}^n)$  at the point  $\mu \in \mathcal{P}_2(\mathbb{R}^n)$ , and impose the following additional assumptions.

**Assumptions 3.12.** In addition to Assumptions 3.8, we suppose that

$$v_t(\cdot) \in \operatorname{Tan}_{P_t} \mathcal{P}_2(\mathbb{R}^n) \quad \text{for Lebesgue-a.e. } t \in [0, T], \quad (3.46)$$

where  $v_t(\cdot)$  is obtained by taking  $\beta \equiv 0$  and  $t_0 = 0$  in (3.43).

**Remark 3.13.** For example, we know from [3, Theorem 10.4.13] that the condition (3.46) is satisfied if, in addition to Assumptions 3.8,  $V$  is uniformly convex, i.e.,  $\operatorname{Hess}(V) \geq \kappa_V I_n$  for some real constant  $\kappa_V$ , and  $W$  is a convex function satisfying the *doubling condition*

$$\exists C_W > 0 \text{ such that } \forall x, y \in \mathbb{R}^n: \quad W(x + y) \leq C_W(1 + W(x) + W(y)). \quad (3.47)$$

The proof of the following result is based on the general theory of Wasserstein metric derivatives of absolutely continuous curves in  $\mathcal{P}_{\operatorname{ac}, 2}(\mathbb{R}^n)$ ; for a thorough discussion, we refer to Chapter 8 in [3].

**Lemma 3.14.** Suppose Assumptions 3.12 hold. For Lebesgue-a.e.  $t_0 \in [0, T]$ , the Wasserstein metric derivative of the perturbed curve  $(P_t^\beta)_{t_0 \leq t \leq T}$  is equal to

$$\lim_{t \downarrow t_0} \frac{W_2(P_t^\beta, P_{t_0}^\beta)}{t - t_0} = \|v_{t_0}^\beta(X_{t_0})\|_{L^2(\mathbb{P})} = \|\nabla \log \ell_{t_0}^\uparrow(X_{t_0}, P_{t_0}) + \nabla \beta(X_{t_0})\|_{L^2(\mathbb{P})}. \quad (3.48)$$

**Proof.** Without loss of generality we can set  $\beta \equiv 0$ . Note that from (3.7) we have

$$\int_0^T \int_{\mathbb{R}^n} |v_t(x)|^2 dp_t(x) dt = \mathbb{E}_{\mathbb{P}} \left[ \int_0^T |\nabla \log \ell_t^\uparrow(X_t, P_t)|^2 dt \right] < \infty, \quad (3.49)$$

which implies that  $v_t(\cdot) \in L^2(P_t)$  for Lebesgue-a.e.  $t \in [0, T]$ . Therefore we can apply Theorem 8.3.1 and Proposition 8.4.5 of [3] to the absolutely continuous curve  $(P_t)_{0 \leq t \leq T}$ , which yields (3.48).  $\square$

We now have all the ingredients to formulate the gradient flow property of the granular media equation. The Wasserstein metric slope of the free energy functional  $\mathcal{F}$  along the McKean–Vlasov curve  $(P_t)_{t_0 \leq t \leq T}$  is defined as

$$|\partial \mathcal{F}|_{W_2}(P_{t_0}) := \lim_{t \downarrow t_0} \frac{H(P_t | Q_t) - H(P_{t_0} | Q_{t_0})}{W_2(P_t, P_{t_0})}. \quad (3.50)$$

In order to show that this is the slope of steepest descent, we will compare it with the slope

$$|\partial \mathcal{F}^\beta|_{W_2}(P_{t_0}^\beta) := \lim_{t \downarrow t_0} \frac{H(P_t^\beta | Q_t^\beta) - H(P_{t_0}^\beta | Q_{t_0}^\beta)}{W_2(P_t^\beta, P_{t_0}^\beta)} \quad (3.51)$$

along the perturbed curve  $(P_t^\beta)_{t_0 \leq t \leq T}$ .

**Theorem 3.15.** *Suppose Assumptions 3.12 hold. Then the following assertions hold for Lebesgue-a.e.  $t_0 \in [0, T]$ : the random variables*

$$\mathcal{L}_{t_0}^\uparrow := \nabla \log \ell_{t_0}^\uparrow(X_{t_0}, P_{t_0}) \quad \text{and} \quad \mathcal{B}_{t_0} := \nabla \beta(X_{t_0}) \quad (3.52)$$

are elements of  $L^2(\mathbb{P})$ , and the Wasserstein metric slope of the free energy functional  $\mathcal{F}$  along the McKean–Vlasov curve  $(P_t)_{t_0 \leq t \leq T}$  is given by

$$|\partial \mathcal{F}|_{W_2}(P_{t_0}) = -\|\mathcal{L}_{t_0}^\uparrow\|_{L^2(\mathbb{P})}. \quad (3.53)$$

If  $\mathcal{L}_{t_0}^\uparrow + \mathcal{B}_{t_0} \neq 0$ , the metric slope along the perturbed curve  $(P_t^\beta)_{t_0 \leq t \leq T}$  is equal to

$$|\partial \mathcal{F}^\beta|_{W_2}(P_{t_0}^\beta) = -\left\langle \mathcal{L}_{t_0}^\uparrow, \frac{\mathcal{L}_{t_0}^\uparrow + \mathcal{B}_{t_0}}{\|\mathcal{L}_{t_0}^\uparrow + \mathcal{B}_{t_0}\|_{L^2(\mathbb{P})}} \right\rangle_{L^2(\mathbb{P})}. \quad (3.54)$$

In particular,

$$|\partial \mathcal{F}|_{W_2}(P_{t_0}) \leq |\partial \mathcal{F}^\beta|_{W_2}(P_{t_0}^\beta) \quad (3.55)$$

with equality if and only if  $\mathcal{L}_{t_0}^\uparrow + \mathcal{B}_{t_0}$  is a positive multiple of  $\mathcal{L}_{t_0}^\uparrow$ .

**Proof.** The equality (3.53) follows from (3.18) and by taking  $\beta \equiv 0$  in (3.48). For the proof of (3.54), we first observe that from (3.41) and (3.48) we obtain the equality

$$|\partial \mathcal{F}^\beta|_{W_2}(P_{t_0}^\beta) = -\frac{\|\mathcal{L}_{t_0}^\uparrow\|_{L^2(\mathbb{P})}^2 + \mathbb{E}_{\mathbb{P}} \left[ \left( \langle \nabla V + \nabla(W * P_{t_0}), \nabla \beta \rangle - \Delta \beta \right)(X_{t_0}) \right]}{\|\mathcal{L}_{t_0}^\uparrow + \mathcal{B}_{t_0}\|_{L^2(\mathbb{P})}} \quad (3.56)$$

for Lebesgue-a.e.  $t_0 \in [0, T]$ . Integrating by parts and recalling the notations in (2.3) – (2.5), we find that the expectation in the numerator of (3.56) is equal to

$$\int_{\mathbb{R}^n} \left\langle \log \nabla \ell_{t_0}^\uparrow(x), \nabla \beta(x) \right\rangle p_{t_0}(x) dx = \left\langle \mathcal{L}_{t_0}^\uparrow, \mathcal{B}_{t_0} \right\rangle_{L^2(\mathbb{P})}. \quad (3.57)$$

Now (3.55) follows by the Cauchy–Schwarz inequality.  $\square$

### 3.3. A trajectorial proof of the HWBI inequality

In this subsection, we show how our trajectorial approach can be adapted to give a simple proof of the HWBI inequality. While the techniques that will be used are similar, the setting of this section is independent from the rest of the paper. In particular, we shall impose convexity assumptions on the potentials  $V, W$ .

We fix two probability measures  $\nu_0$  and  $\nu_1$  in  $\mathcal{P}_{ac,2}(\mathbb{R}^n)$ . By Brenier’s theorem [9], there exists a convex function  $\varphi: \mathbb{R}^n \rightarrow \mathbb{R}$  such that

$$W_2^2(\nu_0, \nu_1) = \int_{\mathbb{R}^n} |x - \nabla \varphi(x)|^2 d\nu_0(x). \quad (3.58)$$



The displacement interpolation of McCann [46] between  $\nu_0$  and  $\nu_1$  is given by

$$\nu_t := (T_t)_\# \nu_0, \quad T_t(x) := (1-t)x + t\nabla\varphi(x), \quad 0 \leq t \leq 1. \quad (3.59)$$

In particular, since the endpoints  $\nu_0$  and  $\nu_1$  belong to  $\mathcal{P}_{ac,2}(\mathbb{R}^n)$ , each  $\nu_t$  has a probability density function  $\rho_t$ ; see, e.g., [58, Remarks 5.13 (i)].

As before, we consider a confinement potential  $V$  and an interaction potential  $W$ . For each  $t \in [0, 1]$ , we then define by analogy with (2.6), the  $\sigma$ -finite measures

$$\mu_t(A) := \int_A q(x, \nu_t) dx, \quad \mu_t^\uparrow(A) := \int_A q^\uparrow(x, \nu_t) dx, \quad A \in \mathcal{B}(\mathbb{R}^n), \quad (3.60)$$

where we recall the definitions of the density functions  $q$  and  $q^\uparrow$  in (2.4). In parallel to the likelihood ratio functions in (2.5), we define

$$r_t(x, \nu) := \frac{\rho_t(x)}{q(x, \nu)}, \quad r_t^\uparrow(x, \nu) := \frac{\rho_t(x)}{q^\uparrow(x, \nu)}, \quad (t, x, \nu) \in [0, 1] \times \mathbb{R}^n \times \mathcal{P}_2(\mathbb{R}^n). \quad (3.61)$$

Then the relative entropy of  $\nu_t$  with respect to  $\mu_t$  is given by

$$H(\nu_t | \mu_t) = \int_{\mathbb{R}^n} \rho_t(x) \log r_t(x, \nu_t) dx \quad (3.62)$$

and the relative Fisher information of  $\nu_t$  with respect to  $\mu_t^\uparrow$  is equal to

$$I(\nu_t | \mu_t^\uparrow) = \int_{\mathbb{R}^n} |\nabla \log r_t^\uparrow(x, \nu_t)|^2 \rho_t(x) dx. \quad (3.63)$$

We impose the following regularity conditions for Proposition 3.17, noting that the strong assumptions placed on  $\rho_0$  and  $\rho_1$  are only temporary and will be removed in Assumptions 3.18 of Theorem 3.19.

**Assumptions 3.16.** The functions  $V, W: \mathbb{R}^n \rightarrow [0, \infty)$  are smooth and  $W$  is symmetric. The probability density functions  $\rho_0$  and  $\rho_1$  are smooth, compactly supported and strictly positive in the interior of their respective supports.

**Proposition 3.17.** Suppose Assumptions 3.16 hold. Along the displacement interpolation  $(\nu_t)_{0 \leq t \leq 1}$ , the rate of relative entropy dissipation at time  $t = 0$ , with respect to the “reference curve of probability measures”  $(\mu_t)_{0 \leq t \leq 1}$ , is given by

$$\left. \frac{d}{dt} \right|_{t=0}^+ H(\nu_t | \mu_t) = \int_{\mathbb{R}^n} \left\langle \nabla \log r_0^\uparrow(x, \nu_0), \nabla \varphi(x) - x \right\rangle \rho_0(x) dx. \quad (3.64)$$

Combining Proposition 3.17 with the displacement convexity results of McCann [46], we obtain the following generalization of the HWBI inequality. Equivalent versions of this inequality can be found in [18, Theorem 4.1] and [25, Theorem D.50].

**Assumptions 3.18.** The functions  $V, W: \mathbb{R}^n \rightarrow [0, \infty)$  are smooth and  $W$  is symmetric. Furthermore,  $V$  and  $W$  are uniformly convex, i.e., there exist real constants  $\kappa_V$  and  $\kappa_W$  such that

$$\text{Hess}(V) \geq \kappa_V I_n, \quad \text{Hess}(W) \geq \kappa_W I_n. \quad (3.65)$$

**Theorem 3.19.** Suppose [Assumptions 3.18](#) hold and the relative entropy  $H(\nu_1 | \mu_1)$  is finite. Then

$$H(\nu_0 | \mu_0) - H(\nu_1 | \mu_1) \leq - \int_{\mathbb{R}^n} \left\langle \nabla \log r_0^\uparrow(x, \nu_0), \nabla \varphi(x) - x \right\rangle \rho_0(x) dx \quad (3.66)$$

$$- \frac{\kappa_V + \kappa_W}{2} W_2^2(\nu_0, \nu_1) + \frac{\kappa_W}{2} |b(\nu_0) - b(\nu_1)|^2. \quad (3.67)$$

**Remark 3.20.** By the Cauchy–Schwarz inequality, the right-hand side of (3.66) can be bounded from above by

$$\sqrt{\int_{\mathbb{R}^n} |\nabla \log r_0^\uparrow(x, \nu_0)|^2 \rho_0(x) dx} \sqrt{\int_{\mathbb{R}^n} |\nabla \varphi(x) - x|^2 \rho_0(x) dx} = \sqrt{I(\nu_0 | \mu_0^\uparrow)} W_2(\nu_0, \nu_1), \quad (3.68)$$

and we obtain the usual form of the HWBI inequality (1.16); see also [2, Theorem 4.2].

## 4. Proofs of the main results

This section is devoted to the proofs of the results stated in [Section 3](#). We shall first prove the main trajectorial results: [Theorem 3.1](#) and its “perturbed” counterpart, [Theorem 3.9](#).

### 4.1. The proofs of [Theorem 3.1](#) and [Theorem 3.9](#)

Since [Theorem 3.1](#) follows immediately from [Theorem 3.9](#) by setting the perturbation  $\beta: \mathbb{R}^n \rightarrow \mathbb{R}$  to be the zero function, we start with the general setting of [Theorem 3.9](#). We first recall a classical result concerning the time reversal of diffusions.

**Lemma 4.1** ([[32](#), Theorem 2.1], [[39](#), Theorems G.2, G.5]). Suppose [Assumptions 3.8](#) hold. On the probability space  $(\Omega, \mathbb{G}, \mathbb{P}^\beta)$ , the process

$$\bar{B}_s^\beta := B_{T-s}^\beta - B_T^\beta - \sqrt{2} \int_0^s \nabla \log \bar{p}_u^\beta(\bar{X}_u^\beta) du, \quad 0 \leq s \leq T - t_0 \quad (4.1)$$

is a Brownian motion. Moreover, the time-reversed canonical process  $(\bar{X}_s)_{0 \leq s \leq T-t_0}$  satisfies

$$d\bar{X}_s = \left( 2\nabla \log \bar{\ell}_s^{\beta\downarrow} - \nabla V + \nabla(W * \bar{\mathbb{P}}_s^\beta) + \nabla\beta \right)(\bar{X}_s) ds + \sqrt{2} d\bar{B}_s^\beta. \quad (4.2)$$

By means of [Lemma 4.1](#), the first step in the proof of [Theorem 3.9](#) is to compute the dynamics of the time-reversed perturbed relative entropy process (3.36). For the reader’s convenience, we recall the following characterization of the L-derivative in [[10](#), pp. 383].

**Definition 4.2.** Let  $f: \mathcal{P}_2(\mathbb{R}^n) \rightarrow \mathbb{R}$  and  $\mu_0 \in \mathcal{P}_2(\mathbb{R}^n)$ . On a probability space  $(\Omega, \mathbb{F}, \mathbb{P})$ , let  $X_0$  be a random variable with distribution  $\mu_0$ . We define  $\partial_\mu f(\mu_0): \mathbb{R}^n \rightarrow \mathbb{R}^n$  as the L-derivative of  $f$  at  $\mu_0$ , if for any  $\mu \in \mathcal{P}_2(\mathbb{R}^n)$  and any random variable  $X$  with distribution  $\mu$ ,

$$f(\mu) = f(\mu_0) + \mathbb{E}_{\mathbb{P}} \left[ \left\langle \partial_\mu f(\mu_0)(X_0), X - X_0 \right\rangle \right] + o(\|X - X_0\|_{L^2(\mathbb{P})}).$$

**Remark 4.3.** The above characterization of the L-derivative depends neither on the choice of the probability space  $(\Omega, \mathbb{F}, \mathbb{P})$ , nor of the random variables  $X$  and  $X_0$  used to represent  $\mu$  and  $\mu_0$ , respectively. Moreover, if the L-derivative exists, it is uniquely defined up to  $\mu_0$ -equivalence. We refer to Proposition 5.25 and Remark 5.26 in [10] for the details.

**Proposition 4.4.** Suppose *Assumptions 3.8* hold. On  $(\Omega, \mathbb{G}, \mathbb{P}^\beta)$ , the time-reversed perturbed relative entropy process (3.36) satisfies

$$d \log \bar{\ell}_s^\beta(\bar{X}_s, \bar{\mathbb{P}}_s^\beta) = \left\langle \nabla \log \bar{\ell}_s^\beta(\bar{X}_s, \bar{\mathbb{P}}_s^\beta), \sqrt{2} d\bar{B}_s^\beta \right\rangle + \left( |\nabla \log \bar{\ell}_s^{\beta\downarrow}|^2 + \frac{1}{2} |\nabla(W * \bar{\mathbb{P}}_s^\beta)|^2 \right) (\bar{X}_s) ds \quad (4.3)$$

$$+ \left( \left\langle \frac{1}{2} \nabla(W * \bar{\mathbb{P}}_s^\beta), 2\nabla \log \bar{\ell}_s^{\beta\downarrow} + \nabla V^\beta \right\rangle + \langle \nabla V, \nabla \beta \rangle - \Delta \beta \right) (\bar{X}_s) ds \quad (4.4)$$

$$- \mathbb{E}_{\bar{\mathbb{P}}^\beta} \left[ \left\langle \frac{1}{2} \nabla W(\bar{X}_s - \bar{Y}_s), \left( 2\nabla \log \bar{\ell}_s^{\beta\downarrow} - \nabla V + \nabla(W * \bar{\mathbb{P}}_s^\beta) + \nabla \beta \right) (\bar{Y}_s) \right\rangle \right] ds, \quad (4.5)$$

where  $(\bar{Y}_s)_{0 \leq s \leq T-t_0}$  is a copy of the process  $(\bar{X}_s)_{0 \leq s \leq T-t_0}$  on a copy  $(\tilde{\Omega}, \tilde{\mathbb{G}}, \tilde{\mathbb{P}}^\beta)$  of the original probability space  $(\Omega, \mathbb{G}, \mathbb{P}^\beta)$ .

**Proof.** Applying a generalized version of Itô's formula for McKean–Vlasov diffusions [10, Proposition 5.102] and using the backward dynamics in (4.2), we obtain

$$d \log \bar{\ell}_s^\beta(\bar{X}_s, \bar{\mathbb{P}}_s^\beta) = \left\langle \nabla \log \bar{\ell}_s^\beta(\bar{X}_s, \bar{\mathbb{P}}_s^\beta), \sqrt{2} d\bar{B}_s^\beta \right\rangle + \left( \partial_s \log \bar{\ell}_s^\beta + \Delta \log \bar{\ell}_s^\beta \right) (\bar{X}_s, \bar{\mathbb{P}}_s^\beta) ds \quad (4.6)$$

$$+ \left\langle \nabla \log \bar{\ell}_s^\beta(\bar{X}_s, \bar{\mathbb{P}}_s^\beta), \left( 2\nabla \log \bar{\ell}_s^{\beta\downarrow} - \nabla V + \nabla(W * \bar{\mathbb{P}}_s^\beta) + \nabla \beta \right) (\bar{X}_s) \right\rangle ds \quad (4.7)$$

$$+ \mathbb{E}_{\bar{\mathbb{P}}^\beta} \left[ \left\langle \left( \partial_\mu \log \bar{\ell}_s^\beta(\bar{X}_s, \bar{\mathbb{P}}_s^\beta) \right), \left( 2\nabla \log \bar{\ell}_s^{\beta\downarrow} - \nabla V + \nabla(W * \bar{\mathbb{P}}_s^\beta) + \nabla \beta \right) (\bar{Y}_s) \right\rangle \right] ds \quad (4.8)$$

$$+ \mathbb{E}_{\bar{\mathbb{P}}^\beta} \left[ \text{trace} \left( \partial_y \partial_\mu \log \bar{\ell}_s^\beta(\bar{X}_s, \bar{\mathbb{P}}_s^\beta) (\bar{Y}_s) \right) \right] ds, \quad (4.9)$$

where  $(\bar{Y}_s)_{0 \leq s \leq T-t_0}$  is a copy of the process  $(\bar{X}_s)_{0 \leq s \leq T-t_0}$  on a copy  $(\tilde{\Omega}, \tilde{\mathbb{G}}, \tilde{\mathbb{P}}^\beta)$  of the original probability space  $(\Omega, \mathbb{G}, \mathbb{P}^\beta)$ . The L-derivative appearing in (4.8) and (4.9) is calculated to be

$$\left( \partial_\mu \log \bar{\ell}_s^\beta(x, \mu) \right) (y) = \frac{1}{2} (\partial_\mu (W * \mu)(x)) (y) = -\frac{1}{2} \nabla W(x - y) \quad (4.10)$$

for  $(x, \mu, y) \in \mathbb{R}^n \times \mathcal{P}_2(\mathbb{R}^n) \times \mathbb{R}^n$ , see [10, Section 5.2.2, Example 1] for the computation of the L-derivative of a function which is linear in the distribution variable. Consequently, we have

$$\text{trace} \left( \partial_y \partial_\mu \log \bar{\ell}_s^\beta(x, \mu)(y) \right) = -\frac{1}{2} \text{trace}(\partial_y \nabla W(x - y)) = \frac{1}{2} \Delta W(x - y). \quad (4.11)$$

Putting (4.10) and (4.11) into (4.8) and (4.9), respectively, as well as using the identities

$$\partial_s \log \bar{\ell}_s^\beta(x, \mu) = \partial_s \log \bar{\ell}_s^{\beta\downarrow}(x), \quad (4.12)$$

$$\nabla \log \bar{\ell}_s^\beta(x, \mu) = \nabla \log \bar{\ell}_s^{\beta\downarrow}(x) + \frac{1}{2} \nabla(W * \mu)(x), \quad (4.13)$$

$$\Delta \log \bar{\ell}_s^\beta(x, \mu) = \Delta \log \bar{\ell}_s^{\beta\downarrow}(x) + \frac{1}{2} \Delta(W * \mu)(x), \quad (4.14)$$

we obtain

$$d \log \bar{\ell}_s^\beta(\bar{X}_s, \bar{P}_s^\beta) = \left\langle \nabla \log \bar{\ell}_s^\beta(\bar{X}_s, \bar{P}_s^\beta), \sqrt{2} d\bar{B}_s^\beta \right\rangle + \left( \partial_s \log \bar{\ell}_s^{\beta\downarrow} + \Delta \log \bar{\ell}_s^{\beta\downarrow} \right) (\bar{X}_s) ds \quad (4.15)$$

$$+ \left\langle \left( \nabla \log \bar{\ell}_s^{\beta\downarrow} + \frac{1}{2} \nabla(W * \bar{P}_s^\beta) \right), \left( 2\nabla \log \bar{\ell}_s^{\beta\downarrow} - \nabla V + \nabla(W * \bar{P}_s^\beta) + \nabla \beta \right) \right\rangle (\bar{X}_s) ds \quad (4.16)$$

$$- \mathbb{E}_{\bar{\mathbb{P}}^\beta} \left[ \left\langle \frac{1}{2} \nabla W(\bar{X}_s - \bar{Y}_s), \left( 2\nabla \log \bar{\ell}_s^{\beta\downarrow} - \nabla V + \nabla(W * \bar{P}_s^\beta) + \nabla \beta \right) (\bar{Y}_s) \right\rangle \right] ds \quad (4.17)$$

$$+ \frac{1}{2} \left( \Delta(W * \bar{P}_s^\beta)(\bar{X}_s) + \mathbb{E}_{\bar{\mathbb{P}}^\beta} \left[ \Delta W(\bar{X}_s - \bar{Y}_s) \right] \right) ds. \quad (4.18)$$

Regarding the expression of (4.18), we observe that  $\Delta(W * \bar{P}_s^\beta)(\bar{X}_s) = \mathbb{E}_{\bar{\mathbb{P}}^\beta} [\Delta W(\bar{X}_s - \bar{Y}_s)]$ . Finally, elementary computations based on (2.4), (3.25) and (3.27) show that the perturbed log-likelihood ratio function  $(s, x) \mapsto \log \bar{\ell}_s^{\beta\downarrow}(x)$  of (3.27) satisfies

$$\begin{aligned} \partial_s \log \bar{\ell}_s^{\beta\downarrow} &= \left\langle \nabla \log \bar{\ell}_s^{\beta\downarrow}, \nabla V - \nabla(W * \bar{P}_s^\beta) - \nabla \beta \right\rangle - |\nabla \log \bar{\ell}_s^{\beta\downarrow}|^2 - \Delta \log \bar{\ell}_s^{\beta\downarrow} \\ &\quad + \left\langle \nabla V, \nabla(W * \bar{P}_s^\beta) + \nabla \beta \right\rangle - \Delta(W * \bar{P}_s^\beta) - \Delta \beta \end{aligned} \quad (4.19)$$

on  $(0, T - t_0) \times \mathbb{R}^n$ , with terminal condition  $\log \bar{\ell}_{T-t_0}^{\beta\downarrow} = \log \bar{\ell}_{T-t_0}^\downarrow$ . Inserting (4.19) into (4.15), we obtain (4.3) – (4.5).  $\square$

Setting the perturbation  $\beta$  to be the zero function, we obtain the following result.

**Corollary 4.5.** *Suppose Assumptions 2.1 hold. On  $(\Omega, \mathbb{G}, \mathbb{P})$ , the time-reversed relative entropy process (3.4) satisfies*

$$d \log \bar{\ell}_s(\bar{X}_s, \bar{P}_s) = \left\langle \nabla \log \bar{\ell}_s(\bar{X}_s, \bar{P}_s), \sqrt{2} d\bar{B}_s \right\rangle \quad (4.20)$$

$$+ \left( |\nabla \log \bar{\ell}_s^\downarrow|^2 + \frac{1}{2} |\nabla(W * \bar{P}_s)|^2 + \left\langle \frac{1}{2} \nabla(W * \bar{P}_s), 2\nabla \log \bar{\ell}_s^\downarrow + \nabla V \right\rangle \right) (\bar{X}_s) ds \quad (4.21)$$

$$- \mathbb{E}_{\bar{\mathbb{P}}} \left[ \left\langle \frac{1}{2} \nabla W(\bar{X}_s - \bar{Y}_s), \left( 2\nabla \log \bar{\ell}_s^\downarrow - \nabla V + \nabla(W * \bar{P}_s) \right) (\bar{Y}_s) \right\rangle \right] ds. \quad (4.22)$$

Here, the process

$$\bar{B}_s := B_{T-s} - B_T - \sqrt{2} \int_0^s \nabla \log \bar{p}_u(\bar{X}_u) du, \quad 0 \leq s \leq T \quad (4.23)$$

is a  $\mathbb{P}$ -Brownian motion with respect to the backward filtration  $\mathbb{G}$ , and  $(\bar{Y}_s)_{0 \leq s \leq T}$  is a copy of the process  $(\bar{X}_s)_{0 \leq s \leq T}$  on a copy  $(\bar{\Omega}, \bar{\mathbb{G}}, \bar{\mathbb{P}})$  of the original probability space  $(\Omega, \mathbb{G}, \mathbb{P})$ .

Before turning to the final part of the proof of Theorem 3.1, we state a classical result based on the general theory of the Cameron–Martin–Maruyama–Girsanov transformation [44]. The connection between relative entropy (the left-hand side of (4.25) below) and energy (the right-hand side of (4.25)) is the foundation of Föllmer’s entropy approach to the time reversal of diffusion processes on Wiener space [26–28]. We denote by  $\mathbb{W}_x$  the Wiener measure on  $\Omega = C([0, T]; \mathbb{R}^n)$  with starting point  $x \in \mathbb{R}^n$ , and define by

$$\mathbb{W}_{x,2}(A) := \mathbb{W}_x \left( \omega \in \Omega : (\sqrt{2}X)(\omega) \in A \right), \quad A \in \mathcal{B}(\Omega) \quad (4.24)$$

the Wiener measure with starting point  $x$  and variance 2.

**Lemma 4.6.** *The relative entropy of  $\mathbb{P}$  with respect to  $\mathbb{W}_{P_0,2} := \int_{\mathbb{R}^n} \mathbb{W}_{x,2} dP_0(x)$  is given by*

$$H(\mathbb{P} | \mathbb{W}_{P_0,2}) = \mathbb{E}_{\mathbb{P}} \left[ \int_0^T |\nabla V(X_t) + \nabla(W * P_t)(X_t)|^2 dt \right] < \infty. \quad (4.25)$$

**Proof.** Recalling (2.3), the drift of the McKean–Vlasov dynamics (1.1) can be expressed as

$$-\nabla \Psi^\uparrow(x, P_t) = -(\nabla V(x) + \nabla(W * P_t)(x)), \quad (t, x) \in [0, T] \times \mathbb{R}^n. \quad (4.26)$$

For any  $t \in [0, T]$ , using the elementary inequality  $(a + b)^2 \leq 2a^2 + 2b^2$ , we have

$$\mathbb{E}_{\mathbb{P}} \left[ |\nabla \Psi^\uparrow(X_t, P_t)|^2 \right] \leq 2 \mathbb{E}_{\mathbb{P}} \left[ |\nabla V(X_t)|^2 \right] + 2 \mathbb{E}_{\mathbb{P}} \left[ |\nabla(W * P_t)(X_t)|^2 \right]. \quad (4.27)$$

Using the linear growth condition (2.1) from Assumptions 2.1 (i), we find

$$\mathbb{E}_{\mathbb{P}} \left[ |\nabla V(X_t)|^2 \right] \leq 2C^2 \left( 1 + \mathbb{E}_{\mathbb{P}} \left[ |X_t|^2 \right] \right) \leq 2C^2 \left( 1 + \mathbb{E}_{\mathbb{P}} \left[ \sup_{0 \leq t \leq T} |X_t|^2 \right] \right). \quad (4.28)$$

Similarly, by Jensen's inequality and (2.1), we obtain

$$\mathbb{E}_{\mathbb{P}} \left[ |\nabla(W * P_t)(X_t)|^2 \right] \leq \int_{\mathbb{R}^n \times \mathbb{R}^n} |\nabla W(x - y)|^2 p_t(y) p_t(x) dy dx \quad (4.29)$$

$$\leq 2C^2 \left( 1 + \int_{\mathbb{R}^n \times \mathbb{R}^n} |x - y|^2 p_t(y) p_t(x) dy dx \right) \quad (4.30)$$

$$\leq 2C^2 \left( 1 + 2 \int_{\mathbb{R}^n \times \mathbb{R}^n} (|x|^2 + |y|^2) p_t(y) p_t(x) dy dx \right) \quad (4.31)$$

$$= 2C^2 \left( 1 + 4 \mathbb{E}_{\mathbb{P}} \left[ |X_t|^2 \right] \right) \leq 8C^2 \left( 1 + \mathbb{E}_{\mathbb{P}} \left[ \sup_{0 \leq t \leq T} |X_t|^2 \right] \right). \quad (4.32)$$

Altogether, we get

$$\mathbb{E}_{\mathbb{P}} \left[ \int_0^T |\nabla \Psi^\uparrow(X_t, P_t)|^2 dt \right] \leq 20T \left( 1 + \mathbb{E}_{\mathbb{P}} \left[ \sup_{0 \leq t \leq T} |X_t|^2 \right] \right) < \infty, \quad (4.33)$$

where the finiteness of this expression follows from the uniform second moment property (2.2) of Lemma 2.2. From [44, Section 7.6.4] we now conclude that  $\mathbb{P}$  is absolutely continuous with respect to  $\mathbb{W}_{P_0,2}$ , and the Radon–Nikodym derivatives are given by

$$\frac{d\mathbb{P}}{d\mathbb{W}_{P_0,2}} \Big|_{\mathcal{F}_t} = \exp \left( -2 \int_0^t \langle \nabla \Psi^\uparrow(X_u, P_u), \sqrt{2} dB_u \rangle + \int_0^t |\nabla \Psi^\uparrow(X_u, P_u)|^2 du \right), \quad 0 \leq t \leq T. \quad (4.34)$$

The integrability property (4.33) implies that the  $\mathbb{P}$ -expectation of the stochastic integral in (4.34) vanishes, and we obtain

$$H(\mathbb{P} | \mathbb{W}_{P_0,2}) = \mathbb{E}_{\mathbb{P}} \left[ \log \left( \frac{d\mathbb{P}}{d\mathbb{W}_{P_0,2}} \right) \right] = \mathbb{E}_{\mathbb{P}} \left[ \int_0^T |\nabla \Psi^\uparrow(X_t, P_t)|^2 dt \right] < \infty, \quad (4.35)$$

which shows (4.25).  $\square$

Denoting  $n$ -dimensional Lebesgue measure by  $\lambda$ , we consider on  $\Omega = C([0, T]; \mathbb{R}^n)$  the  $\sigma$ -finite measure  $\mathbb{W}_{\lambda,2} := \int_{\mathbb{R}^n} \mathbb{W}_{x,2} d\lambda(x)$ , which is known as the law of the *reversible Brownian motion* on  $\mathbb{R}^n$  with variance 2; see [19, 42]. The fundamental property of reversible Brownian motion is that it is invariant under time reversal. This property can be formalized as follows. Let  $R: \Omega \rightarrow \Omega$  be the pathwise time reversal operator on  $\Omega$ , given by  $X_s \circ R = X_{T-s}$  for  $s \in [0, T]$ . For any measure  $\mu$  on  $\Omega$ , we denote its time reversal by  $\bar{\mu} := R_{\#}\mu$ . Then we have the invariance property  $\bar{\mathbb{W}}_{\lambda,2} = \mathbb{W}_{\lambda,2}$ . Let us also consider the probability measure  $\mathbb{W}_{P_T,2} := \int_{\mathbb{R}^n} \mathbb{W}_{x,2} dP_T(x)$  and its time reversal given by  $\bar{\mathbb{W}}_{P_T,2} = \int_{\mathbb{R}^n} \bar{\mathbb{W}}_{x,2} dP_T(x)$ . Then, as already noted in [26, Remarks 3.7], we have the following result.

**Lemma 4.7.** *We have the relative entropy relations*

$$H(\mathbb{P} | \mathbb{W}_{\lambda,2}) = H(P_0 | \lambda) + H(\mathbb{P} | \mathbb{W}_{P_0,2}) \quad (4.36)$$

and

$$H(\mathbb{P} | \bar{\mathbb{W}}_{\lambda,2}) = H(P_T | \lambda) + H(\mathbb{P} | \bar{\mathbb{W}}_{P_T,2}). \quad (4.37)$$

Furthermore, all these relative entropies are finite.

**Proof.** For any  $x \in \mathbb{R}^n$ , we let  $\mathbb{P}_x(\cdot) := \mathbb{P}(\cdot | X_0 = x)$  denote (a version of) the conditional probability measure  $\mathbb{P}$  given  $X_0 = x$ . By the chain rule for relative entropy [19, Theorem 2] we have

$$H(\mathbb{P} | \mathbb{W}_{\lambda,2}) = H(P_0 | \lambda) + \int_{\mathbb{R}^n} H(\mathbb{P}_x | \mathbb{W}_{x,2}) dP_0(x) \quad (4.38)$$

and at the same time

$$H(\mathbb{P} | \mathbb{W}_{P_0,2}) = H(P_0 | P_0) + \int_{\mathbb{R}^n} H(\mathbb{P}_x | \mathbb{W}_{x,2}) dP_0(x) = \int_{\mathbb{R}^n} H(\mathbb{P}_x | \mathbb{W}_{x,2}) dP_0(x), \quad (4.39)$$

implying the first identity (4.36). Regarding the finite entropy assertions, we recall (1.2), (1.3) and observe that

$$H(P_0 | \lambda) = \int_{\mathbb{R}^n} p_0(x) \log p_0(x) dx = \mathcal{U}(P_0) \leq \mathcal{F}(P_0) < \infty, \quad (4.40)$$

where the finiteness follows from Assumptions 2.1 (ii). Furthermore, from Lemma 4.6 we know that  $H(\mathbb{P} | \mathbb{W}_{P_0,2}) < \infty$ .

By the same arguments as above, (4.37) follows again by the chain rule for relative entropy. Using the invariance property  $\bar{\mathbb{W}}_{\lambda,2} = \mathbb{W}_{\lambda,2}$ , and as we already know that  $H(\mathbb{P} | \mathbb{W}_{\lambda,2}) < \infty$ , it follows that  $H(\mathbb{P} | \bar{\mathbb{W}}_{\lambda,2}) < \infty$ . Let us recall now that  $P_T \in \mathcal{P}_{ac,2}(\mathbb{R}^n)$  by Lemma 2.2. On the one hand, since  $P_T$  has finite second moment,  $H(P_T | \lambda)$  cannot take the value  $-\infty$  as noted in Remark 2.3. On the other hand, the absolute continuity of  $P_T$  implies that  $H(P_T | \lambda)$  cannot take the value  $+\infty$ . Therefore, we conclude that  $H(\mathbb{P} | \bar{\mathbb{W}}_{P_T,2}) < \infty$ .  $\square$

We have assembled now all the ingredients needed for the proof of Theorem 3.1.

**Proof of Theorem 3.1.** Recalling the definition of the stochastic integral process  $(\bar{M}_s)_{0 \leq s \leq T}$  in (3.6) and of the cumulative Fisher information process  $(\bar{F}_s)_{0 \leq s \leq T}$  in (3.3), we see that the stochastic differential of (4.20) – (4.22) can be expressed as claimed in (3.5).

Since  $(\bar{M}_s)_{0 \leq s \leq T}$  is a stochastic integral process, it is a continuous local martingale. In order to show that it is an  $L^2(\mathbb{P})$ -bounded martingale, it suffices to show the integrability condition

$$\mathbb{E}_{\mathbb{P}} \left[ \langle \bar{M}, \bar{M} \rangle_T \right] = \mathbb{E}_{\mathbb{P}} \left[ 2 \int_0^T |\nabla \log \bar{\ell}_u(\bar{X}_u, \bar{P}_u)|^2 du \right] < \infty; \quad (4.41)$$

see, e.g. [51, Corollary IV.1.25]. On  $(\Omega, \mathbb{G}, \mathbb{P})$ , the time-reversed canonical process  $(\bar{X}_s)_{0 \leq s \leq T}$  has backward dynamics

$$d\bar{X}_s = \bar{\vartheta}_s(\bar{X}_s) ds + \sqrt{2} d\bar{B}_s, \quad 0 \leq s \leq T \quad (4.42)$$

with initial distribution  $\mathbb{P}_T$ , where the drift term is given by

$$\bar{\vartheta}_s(x) := \left( 2\nabla \log \bar{\ell}_s^\downarrow - \nabla V + \nabla(W * \bar{P}_s) \right)(x) = 2\nabla \log \bar{\ell}_s(x, \bar{P}_s) - \nabla V(x) \in \mathbb{R}^n \quad (4.43)$$

for  $(s, x) \in [0, T] \times \mathbb{R}^n$ . Therefore, in order to prove (4.41), it suffices to show the two integrability conditions

$$\mathbb{E}_{\mathbb{P}} \left[ \int_0^T |\nabla V(\bar{X}_u)|^2 du \right] < \infty \quad \text{and} \quad \mathbb{E}_{\mathbb{P}} \left[ \int_0^T |\bar{\vartheta}_u(\bar{X}_u)|^2 du \right] < \infty. \quad (4.44)$$

The first condition is a direct consequence of (4.28). From [26, Lemma 2.6] we conclude that the expectation of the second condition is bounded by the relative entropy  $H(\mathbb{P} | \bar{\mathbb{W}}_{\mathbb{P}_T, 2})$ , which is finite on account of Lemma 4.7.

In order to complete the proof of Theorem 3.1, it remains to show (3.7). To begin with, we take expectation with respect to  $\mathbb{P}$  in (3.3) and invoke Fubini's theorem to interchange the  $\mathbb{P}$ -expectation and the time integral. Applying once more Fubini's theorem, we swap the  $\mathbb{P}$ -expectation with the  $\tilde{\mathbb{P}}$ -expectation appearing in (3.2). Next, we recall Assumptions 2.1 (i) and use the symmetry of the interaction potential, which implies that  $\nabla W(-x) = -\nabla W(x)$  for all  $x \in \mathbb{R}^n$ . Furthermore, as the distribution of  $\bar{Y}_u$  under  $\tilde{\mathbb{P}}$  is the same as the distribution of  $\bar{X}_u$  under  $\mathbb{P}$ , we deduce that

$$\mathbb{E}_{\mathbb{P}} [\bar{F}_s] = \int_0^s \mathbb{E}_{\mathbb{P}} \left[ \left( |\nabla \log \bar{\ell}_u^\downarrow|^2 + \frac{1}{2} |\nabla(W * \bar{P}_u)|^2 + \left\langle \frac{1}{2} \nabla(W * \bar{P}_u), 2\nabla \log \bar{\ell}_u^\downarrow + \nabla V \right\rangle \right) (\bar{X}_u) \right] du \quad (4.45)$$

$$+ \int_0^s \mathbb{E}_{\mathbb{P}} \left[ \left\langle \frac{1}{2} \nabla(W * \bar{P}_u), \left( 2\nabla \log \bar{\ell}_u^\downarrow - \nabla V + \nabla(W * \bar{P}_u) \right) (\bar{X}_u) \right\rangle \right] du \quad (4.46)$$

for  $0 \leq s \leq T$ . Recalling the definitions in (2.3) – (2.5), we obtain

$$\mathbb{E}_{\mathbb{P}} [\bar{F}_s] = \int_0^s \mathbb{E}_{\mathbb{P}} \left[ |\nabla \log \bar{\ell}_u^\uparrow(\bar{X}_u, \bar{P}_u)|^2 \right] du = \int_0^s I(\bar{P}_u | \bar{Q}_u^\uparrow) du < \infty, \quad 0 \leq s \leq T, \quad (4.47)$$

where the second equality is immediate from (2.7), and the finiteness of the expression in (4.47) is justified as follows. Again, from (2.3) – (2.5) we find

$$|\nabla \log \bar{\ell}_s^\uparrow(\bar{X}_s, \bar{P}_s)|^2 = |\nabla \log \bar{\ell}_s(\bar{X}_s, \bar{P}_s) + \frac{1}{2} \nabla(W * \bar{P}_s)(\bar{X}_s)|^2 \quad (4.48)$$

$$\leq 2|\nabla \log \bar{\ell}_s(\bar{X}_s, \bar{P}_s)|^2 + \frac{1}{2} |\nabla(W * \bar{P}_s)(\bar{X}_s)|^2. \quad (4.49)$$

In light of (4.41) and (4.29) – (4.32), we see that the expression in (4.47) is finite, which in turn justifies a posteriori the former applications of Fubini's theorem.  $\square$

**Remark 4.8.** In the above proofs of [Lemmas 4.6, 4.7](#) and [Theorem 3.1](#), the linear growth condition (2.1) from [Assumptions 2.1](#) (i) was crucial. For a different proof in the (linear) setting without interaction (i.e.,  $W \equiv 0$ ), we refer to [\[40\]](#), where the confinement potential  $V$  (which is called  $\Psi$  in [\[40\]](#)) is not necessarily of linear growth, but instead satisfies a weaker coercivity condition. The key in that setting is to apply [Lemma 2.48](#) in [\[37\]](#).

The proof of [Theorem 3.9](#) is now an easy consequence.

**Proof of Theorem 3.9.** Recalling the definition of the process  $(\overline{M}_s^\beta)_{0 \leq s \leq T}$  in (3.38) and of the perturbed cumulative Fisher information process  $(\overline{F}_s^\beta)_{0 \leq s \leq T}$  in (3.35), we see that the stochastic differential of (4.3) – (4.5) can be expressed as claimed in (3.37).

As in the proof of [Theorem 3.1](#) we will now argue that

$$\mathbb{E}_{\mathbb{P}^\beta} \left[ \langle \overline{M}^\beta, \overline{M}^\beta \rangle_{T-t_0} \right] = \mathbb{E}_{\mathbb{P}^\beta} \left[ 2 \int_0^{T-t_0} |\nabla \log \overline{\ell}_u^\beta(\overline{X}_u, \overline{\mathbf{P}}_u^\beta)|^2 du \right] < \infty, \quad (4.50)$$

which will then imply that the stochastic integral process  $(\overline{M}_s^\beta)_{0 \leq s \leq T-t_0}$  is an  $L^2(\mathbb{P}^\beta)$ -bounded martingale. To this end, we define the density  $q^\beta(x, \mu) := e^{-\Psi^\beta(x, \mu)}$  for  $(x, \mu) \in \mathbb{R}^n \times \mathcal{P}_2(\mathbb{R}^n)$ , and consider the “doubly perturbed” likelihood ratio function

$$\ell_t^{\beta, \beta}(x, \mu) := \frac{p_t^\beta(x)}{q^\beta(x, \mu)}, \quad (t, x) \in [t_0, T] \times \mathbb{R}^n. \quad (4.51)$$

As the [Assumptions 2.1](#) are invariant under the passage from the potential  $V$  to  $V^\beta = V + \beta$ , we can apply [Theorem 3.1](#) to the potential  $V^\beta$  and obtain

$$\mathbb{E}_{\mathbb{P}^\beta} \left[ 2 \int_0^{T-t_0} |\nabla \log \overline{\ell}_u^{\beta, \beta}(\overline{X}_u, \overline{\mathbf{P}}_u^\beta)|^2 du \right] < \infty. \quad (4.52)$$

Now, since  $\ell_t^\beta(x, \mu) / \ell_t^{\beta, \beta}(x, \mu) = e^{\beta(x)}$ , we observe that the difference

$$\nabla \log \ell_t^\beta(x, \mu) - \nabla \log \ell_t^{\beta, \beta}(x, \mu) = \nabla \beta(x) \quad (4.53)$$

is a bounded function. Together with (4.52), this implies (4.50).

It remains to check (3.39). A similar calculation as in the proof of [Theorem 3.1](#) leads to the identity

$$\mathbb{E}_{\mathbb{P}^\beta} [\overline{F}_s^\beta] = \int_0^s \mathbb{E}_{\mathbb{P}^\beta} \left[ |\nabla \log \overline{\ell}_u^{\beta \uparrow}(\overline{X}_u, \overline{\mathbf{P}}_u^\beta)|^2 + \left( \langle \nabla V + \nabla(W * \overline{\mathbf{P}}_u^\beta), \nabla \beta \rangle - \Delta \beta \right)(\overline{X}_u) \right] du \quad (4.54)$$

for  $0 \leq s \leq T - t_0$ . Repeating the reasoning of the previous paragraph for the function  $\ell_t^{\beta \uparrow}$  instead of  $\ell_t^\beta$ , we find that

$$\mathbb{E}_{\mathbb{P}^\beta} \left[ \int_0^{T-t_0} |\nabla \log \overline{\ell}_u^{\beta \uparrow}(\overline{X}_u, \overline{\mathbf{P}}_u^\beta)|^2 du \right] < \infty. \quad (4.55)$$

Since the function

$$[0, T - t_0] \times \mathbb{R}^n \ni (t, x) \mapsto \left\langle \nabla V + \nabla(W * \overline{\mathbf{P}}_t^\beta), \nabla \beta \right\rangle(x) - \Delta \beta(x) \quad (4.56)$$

is bounded, we conclude that the quantity of (4.54) is finite. Finally, recalling the definition (3.30), we arrive at (3.39).  $\square$



## 4.2. The proofs of Proposition 3.17 and Theorem 3.19

**Proof of Proposition 3.17.** The first step is to view the probability density functions  $(\rho_t)_{0 \leq t \leq 1}$ , corresponding to the displacement interpolation  $(v_t)_{0 \leq t \leq 1}$  of (3.59), as a solution to a continuity equation. Recalling the convex function  $\varphi: \mathbb{R}^n \rightarrow \mathbb{R}$  of (3.58), we define a function  $u_0: \mathbb{R}^n \rightarrow \mathbb{R}$  by  $u_0(x) := \varphi(x) - |x|^2/2$ ; and for each  $t \in (0, 1]$ , we let the function  $u_t: \mathbb{R}^n \rightarrow \mathbb{R}$  be defined by the Hopf–Lax formula

$$u_t(x) := \inf_{y \in \mathbb{R}^n} \left( u_0(y) + \frac{|x - y|^2}{2t} \right). \quad (4.57)$$

For all  $t \in [0, 1)$ , we denote the gradient of  $u_t$  by  $v_t := \nabla u_t$ . For  $t = 0$ , it is clear that  $v_0 = \nabla \varphi - \text{Id}$  is well-defined. For  $t \in (0, 1)$ , the gradient  $v_t$  is defined Lebesgue-a.e. by [58, Theorem 5.51 (i)], and

$$v_t(x) = \nabla u_0 \circ (T_t)^{-1}(x), \quad \text{for all } x \in T_t(\mathbb{R}^n), \quad (4.58)$$

where  $T_t$  is defined in (3.59). Note that the inverse of  $T_t$  is well-defined because  $T_t$  is injective; see [58, Section 5.4.8]. From (4.58) we see that  $(v_t)_{0 \leq t < 1}$  is the *velocity field* associated with the trajectories  $(T_t)_{0 \leq t < 1}$ , i.e.,

$$T_t(x) = x + \int_0^t v_s(T_s(x)) \, ds, \quad 0 \leq t < 1. \quad (4.59)$$

By [58, Theorem 5.51 (ii)], the curve of probability density functions  $(\rho_t)_{0 < t < 1}$  satisfies the continuity equation

$$\partial_t \rho_t(x) + \text{div}(\rho_t(x) v_t(x)) = 0, \quad (t, x) \in (0, 1) \times \mathbb{R}^n. \quad (4.60)$$

On a sufficiently rich probability space  $(\mathcal{S}, \mathcal{S}, \mathbb{P})$ , we let  $Z_0: \mathcal{S} \rightarrow \mathbb{R}^n$  be a random variable with probability distribution  $\nu_0$ . For each  $0 < t \leq 1$ , we let  $Z_t := T_t(Z_0)$ . From (3.59) we see that the random variable  $Z_t$  has distribution  $\nu_t$ , and (4.59) yields the representation

$$Z_t = Z_0 + \int_0^t v_s(Z_s) \, ds, \quad 0 \leq t < 1. \quad (4.61)$$

In conjunction with (4.60), we deduce

$$d\rho_t(Z_t) = \partial_t \rho_t(Z_t) + \langle \nabla \rho_t(Z_t), dZ_t \rangle = -\rho_t(Z_t) \text{div}(v_t(Z_t)) \, dt, \quad (4.62)$$

and thus

$$d \log \rho_t(Z_t) = -\text{div}(v_t(Z_t)) \, dt. \quad (4.63)$$

Recalling the definition of the density function  $q$  in (2.4), a similar argument as in (4.10) shows that

$$(\partial_\nu \log q(x, \nu))(y) = \frac{1}{2} (\partial_\nu (W * \nu))(x)(y) = -\frac{1}{2} \nabla W(x - y) \quad (4.64)$$

for  $(x, \nu, y) \in \mathbb{R}^n \times \mathcal{P}_2(\mathbb{R}^n) \times \mathbb{R}^n$ . Applying a generalized version of Itô's formula for McKean–Vlasov diffusions [10, Proposition 5.102], and using the dynamics (4.61) as well as the L-derivative (4.64), we obtain

$$d \log q(Z_t, \nu_t) = -\left\langle \nabla W + \frac{1}{2} \nabla (W * \nu_t), v_t \right\rangle (Z_t) \, dt + \frac{1}{2} \mathbb{E}_{\tilde{\mathbb{P}}} \left[ \left\langle \nabla W(Z_t - \tilde{Z}_t), v_t(\tilde{Z}_t) \right\rangle \right] \, dt \quad (4.65)$$

for  $0 < t < 1$ . Here, the process  $(\tilde{Z}_t)_{0 < t < 1}$  is defined on another probability space  $(\tilde{S}, \tilde{\mathcal{S}}, \tilde{\mathbb{P}})$  such that the tuple  $(\tilde{S}, \tilde{\mathcal{S}}, \tilde{\mathbb{P}}, (\tilde{Z}_t)_{0 < t < 1})$  is an exact copy of  $(\tilde{S}, \tilde{\mathcal{S}}, \tilde{\mathbb{P}}, (\tilde{Z}_t)_{0 < t < 1})$ . Now taking the difference between (4.63) and (4.65) gives the dynamics

$$\log r_t(Z_t, \nu_t) - \log r_0(Z_0, \nu_0) = \int_0^t \left( \left\langle \nabla V + \frac{1}{2} \nabla(W * \nu_s), \nu_s \right\rangle(Z_s) - \operatorname{div}(\nu_s(Z_s)) \right) ds \quad (4.66)$$

$$- \frac{1}{2} \int_0^t \mathbb{E}_{\tilde{\mathbb{P}}} \left[ \left\langle \nabla W(Z_s - \tilde{Z}_s), \nu_s(\tilde{Z}_s) \right\rangle \right] ds \quad (4.67)$$

of the relative entropy process  $(\log r_t(Z_t, \nu_t))_{0 < t < 1}$ . Next, let us make two observations. Firstly, integration by parts yields

$$\mathbb{E}_{\mathbb{P}} \left[ \operatorname{div}(\nu_t(Z_t)) \right] = -\mathbb{E}_{\mathbb{P}} \left[ \left\langle \nabla \log \rho_t(Z_t), \nu_t(Z_t) \right\rangle \right]. \quad (4.68)$$

Secondly, by applying Fubini's theorem, and using that  $W$  is an even function as well as  $(\tilde{Z}_t)_{\#} \tilde{\mathbb{P}} = \nu_t$ , we obtain the identity

$$\mathbb{E}_{\mathbb{P}} \left[ \mathbb{E}_{\tilde{\mathbb{P}}} \left[ \left\langle \nabla W(Z_t - \tilde{Z}_t), \nu_t(\tilde{Z}_t) \right\rangle \right] \right] = -\mathbb{E}_{\mathbb{P}} \left[ \left\langle \nabla(W * \nu_t)(Z_t), \nu_t(Z_t) \right\rangle \right]. \quad (4.69)$$

Returning to (4.66), (4.67), we take  $\mathbb{P}$ -expectations and use (4.68), (4.69) to obtain

$$H(\nu_t | \mu_t) - H(\nu_0 | \mu_0) = \int_0^t \mathbb{E}_{\mathbb{P}} \left[ \left\langle \nabla \log \rho_s + \nabla V + \nabla(W * \nu_s), \nu_s \right\rangle(Z_s) \right] ds \quad (4.70)$$

$$= \int_0^t \mathbb{E}_{\mathbb{P}} \left[ \left\langle \nabla \log r_s^\uparrow(Z_s, \nu_s), \nu_s(Z_s) \right\rangle \right] ds, \quad (4.71)$$

where for the second equality we recall the notations in (3.61) and (2.4). Finally, letting  $t \downarrow 0$ , we get

$$\frac{d}{dt} \Big|_{t=0}^+ H(\nu_t | \mu_t) = \int_{\mathbb{R}^n} \left\langle \nabla \log r_0^\uparrow(x, \nu_0), \nu_0(x) \right\rangle \rho_0(x) dx; \quad (4.72)$$

and since  $\nu_0 = \nabla \varphi - \operatorname{Id}$ , we arrive at (3.64).  $\square$

**Proof of Theorem 3.19.** Without loss of generality, we assume that the probability density functions  $\rho_0$  and  $\rho_1$  satisfy the strong regularity Assumptions 3.16. The general case then follows by a density argument. We will not provide the details here, but refer to [58, Chapter 9.4], where this regularization is carried out in the simpler setting of the HWI inequality.

Let us recall the energy functionals  $\mathcal{U}$ ,  $\mathcal{V}$ ,  $\mathcal{W}$  defined in (1.3), and introduce the functions

$$f(t) := \mathcal{U}(\rho_t), \quad g(t) := \mathcal{V}(\rho_t), \quad h(t) := \mathcal{W}(\rho_t), \quad 0 \leq t \leq 1, \quad (4.73)$$

where  $(\rho_t)_{0 \leq t \leq 1}$  is the curve of probability density functions corresponding to the displacement interpolation  $(\nu_t)_{0 \leq t \leq 1}$  of (3.59). Then the sum  $F := f + g + h$  of these functions satisfies the relation  $F(t) = H(\nu_t | \mu_t)$ . In light of [58, Theorem 5.15 (i)], the internal energy functional  $\mathcal{U}$  is displacement convex, i.e.,

$$f''(t) \geq 0, \quad 0 \leq t \leq 1. \quad (4.74)$$

By [Assumptions 3.18](#), the confinement potential  $V: \mathbb{R}^n \rightarrow [0, \infty)$  is  $\kappa_V$ -uniformly convex. Therefore, [\[58, Theorem 5.15 \(ii\)\]](#) implies that the potential energy functional  $\mathcal{V}$  is  $\kappa_V$ -uniformly displacement convex. In other words,

$$g''(t) \geq \kappa_V W_2^2(v_0, v_1), \quad 0 \leq t \leq 1. \quad (4.75)$$

Again from [Assumptions 3.18](#), the interaction potential  $W: \mathbb{R}^n \rightarrow [0, \infty)$  is assumed to be symmetric and  $\kappa_W$ -uniformly convex. Therefore, a similar argument as in the proof of [\[58, Theorem 5.15 \(iii\)\]](#) leads to the  $\kappa_W(W_2^2(v_0, v_1) - |b(v_0) - b(v_1)|^2)$ -uniform convexity of  $h$ , so

$$h''(t) \geq \kappa_W \left( W_2^2(v_0, v_1) - |b(v_0) - b(v_1)|^2 \right), \quad 0 \leq t \leq 1. \quad (4.76)$$

The details of the proof of [\(4.76\)](#) are postponed to [Subsection 5.2](#). By combining the estimates [\(4.74\) – \(4.76\)](#), we deduce that the relative entropy function  $[0, 1] \ni t \mapsto F(t) = H(v_t | \mu_t)$  satisfies

$$F''(t) \geq (\kappa_V + \kappa_W) W_2^2(v_0, v_1) - \kappa_W |b(v_0) - b(v_1)|^2. \quad (4.77)$$

Furthermore, from [Proposition 3.17](#) we have

$$F'(0^+) = \int_{\mathbb{R}^n} \left\langle \nabla \log r_0^\uparrow(x, v_0), \nabla \varphi(x) - x \right\rangle \rho_0(x) dx. \quad (4.78)$$

In conjunction with [\(4.77\)](#) and [\(4.78\)](#), the Taylor formula  $F(1) = F(0) + F'(0^+) + \int_0^1 (1-t)F''(t) dt$  now yields the inequality [\(3.66\) – \(3.67\)](#).  $\square$

## 5. Proofs of auxiliary results

### 5.1. Proof of Lemma 2.2

The generalized potential  $\Psi^\uparrow$  of [\(2.3\)](#) allows us to cast the McKean–Vlasov dynamics of [\(1.1\)](#) in the more compact form

$$dX_t = -\nabla \Psi^\uparrow(X_t, P_t) dt + \sqrt{2} dB_t, \quad 0 \leq t \leq T. \quad (5.1)$$

Then, for any two pairs  $(x, \mu), (x', \mu') \in \mathbb{R}^n \times \mathcal{P}_2(\mathbb{R}^n)$ , using the Lipschitz continuity of  $\nabla V$  in [Assumptions 2.1 \(i\)](#) yields

$$|\nabla \Psi^\uparrow(x, \mu) - \nabla \Psi^\uparrow(x', \mu')| \leq \|\nabla V\|_{\text{Lip}} |x - x'| + |\nabla(W * \mu)(x) - \nabla(W * \mu')(x')|. \quad (5.2)$$

For the convolution term, using Jensen's inequality and the Lipschitz continuity of  $\nabla W$  in [Assumptions 2.1 \(i\)](#) leads to

$$|\nabla(W * \mu)(x) - \nabla(W * \mu')(x')| \leq \|\nabla W\|_{\text{Lip}} |x - x'| + |\nabla(W * \mu)(x') - \nabla(W * \mu')(x')|. \quad (5.3)$$

For the last term above, by the Kantorovich–Rubinstein theorem [\[58, Theorem 1.14\]](#), we have

$$\begin{aligned} \left| \int_{\mathbb{R}^n} \nabla W(x' - \cdot) d(\mu - \mu') \right| &\leq \|\nabla W\|_{\text{Lip}} \sup \left\{ \int_{\mathbb{R}^n} \varphi d(\mu - \mu') : \varphi \in L^1(|\mu - \mu'|), \|\varphi\|_{\text{Lip}} \leq 1 \right\} \\ &= \|\nabla W\|_{\text{Lip}} W_1(\mu, \mu') \leq \|\nabla W\|_{\text{Lip}} W_2(\mu, \mu'), \end{aligned} \quad (5.4)$$

where

$$W_1(\mu, \mu') = \inf_{Y \sim \mu, Z \sim \mu'} \mathbb{E}|Y - Z|, \quad \mu, \mu' \in \mathcal{P}_1(\mathbb{R}^n) \quad (5.5)$$

denotes the 1-Wasserstein-distance, and the inequality in (5.4) follows from Jensen's inequality. Altogether, we obtain

$$|\nabla \Psi^\uparrow(x, \mu) - \nabla \Psi^\uparrow(x', \mu')| \leq (\|\nabla V\|_{\text{Lip}} + \|\nabla W\|_{\text{Lip}}) |x - x'| + \|\nabla W\|_{\text{Lip}} W_2(\mu, \mu'). \quad (5.6)$$

In particular, this shows that the function  $-\nabla \Psi^\uparrow$  is Lipschitz continuous on the product metric space  $(\mathbb{R}^n, |\cdot|) \times (\mathcal{P}_2(\mathbb{R}^n), W_2)$ . In conjunction with [Assumptions 2.1 \(ii\)](#), [[10](#), Theorem 4.21] implies that the McKean–Vlasov SDE (5.1) has a pathwise unique, strong solution satisfying the uniform second moment condition (2.2). Now we can linearize (5.1) by fixing the time-marginals  $(P_t)_{0 \leq t \leq T}$ , so that the drift term can be viewed as a function  $(t, x) \mapsto \nabla \Psi^\uparrow(x, P_t)$ , and (5.1) becomes an ordinary SDE with a time-inhomogeneous drift coefficient.

The absolute continuity of the time-marginals  $(P_t)_{0 \leq t \leq T}$  is immediate from [Lemma 4.6](#). A standard argument using the classical Itô's formula shows that the curve of probability density functions  $(p_t)_{0 \leq t \leq T}$  is a weak solution of the granular media equation (1.11). Finally, we turn to the regularity of this solution. From (5.6), we see that the drift  $x \mapsto \nabla \Psi^\uparrow(x, P_t)$  is Lipschitz continuous for every  $t \in [0, T]$ , and [Assumptions 2.1 \(i\)](#) implies that the drift is also of linear growth. The desired smoothness of  $(p_t)_{0 \leq t \leq T}$  now follows from a straightforward adaptation of the theorem in [[53](#)], see also Remarks (i) – (ii) therein.  $\square$

## 5.2. Proof of (4.76)

We first rewrite the interaction energy functional  $\mathcal{W}$  along the displacement interpolation  $(\nu_t)_{0 \leq t \leq 1}$ . Using (3.59), for any  $t \in [0, 1]$ , we have

$$h(t) = \frac{1}{2} \int_{\mathbb{R}^n \times \mathbb{R}^n} W(x - y) \nu_t(dx) \nu_t(dy) \quad (5.7)$$

$$= \frac{1}{2} \int_{\mathbb{R}^n \times \mathbb{R}^n} W(T_t(x) - T_t(y)) \nu_0(dx) \nu_0(dy) \quad (5.8)$$

$$= \frac{1}{2} \int_{\mathbb{R}^n \times \mathbb{R}^n} W\left(x - y - t(\theta(x) - \theta(y))\right) \nu_0(dx) \nu_0(dy), \quad (5.9)$$

where  $\theta: \mathbb{R}^n \rightarrow \mathbb{R}^n$  is defined as  $\theta(x) := x - \nabla \varphi(x)$ . Now, for any  $t_1, t_2, \sigma \in [0, 1]$ , by the  $\kappa_W$ -uniform convexity of  $W$  in [Assumptions 3.18](#), we obtain

$$\sigma h(t_1) + (1 - \sigma)h(t_2) - h(\sigma t_1 + (1 - \sigma)t_2) \quad (5.10)$$

$$= \frac{1}{2} \int_{\mathbb{R}^n \times \mathbb{R}^n} \left( \sigma W\left(x - y - t_1(\theta(x) - \theta(y))\right) + (1 - \sigma)W\left(x - y - t_2(\theta(x) - \theta(y))\right) \right. \quad (5.11)$$

$$\left. - W\left(x - y - (\sigma t_1 + (1 - \sigma)t_2)(\theta(x) - \theta(y))\right) \right) \nu_0(dx) \nu_0(dy) \quad (5.12)$$

$$\geq \frac{1}{4} \kappa_W \sigma(1 - \sigma)(t_1 - t_2)^2 \int_{\mathbb{R}^n \times \mathbb{R}^n} |\theta(x) - \theta(y)|^2 \nu_0(dx) \nu_0(dy). \quad (5.13)$$

Next, we express the integral in (5.13) as

$$\begin{aligned}
 & \frac{1}{2} \int_{\mathbb{R}^n \times \mathbb{R}^n} |\theta(x) - \theta(y)|^2 \nu_0(dx) \nu_0(dy) \\
 &= \int_{\mathbb{R}^n} |\theta(x)|^2 \nu_0(dx) - \left| \int_{\mathbb{R}^n} \theta(x) \nu_0(dx) \right|^2 \\
 &= \int_{\mathbb{R}^n} |x - \nabla \varphi(x)|^2 \nu_0(dx) - \left| \int_{\mathbb{R}^n} x \nu_0(dx) - \int_{\mathbb{R}^n} x \nu_1(dx) \right|^2 \\
 &= W_2^2(\nu_0, \nu_1) - |b(\nu_0) - b(\nu_1)|^2.
 \end{aligned}$$

Putting this back into (5.13), we deduce that  $h$  is uniformly convex, with constant

$$\kappa_W \left( W_2^2(\nu_0, \nu_1) - |b(\nu_0) - b(\nu_1)|^2 \right). \quad (5.14)$$

□

## Acknowledgments

We are grateful to Ioannis Karatzas and Walter Schachermayer for suggesting this problem and giving us generous advice. We thank Robert Fernholz, Miguel Garrido, Tomoyuki Ichiba, Donghan Kim, Kasper Larsen, and Mete Soner for helpful comments during the INTECH research meetings. Thanks also go to an associate editor and an anonymous referee for their valuable comments and suggestions.

## Funding

B. Tschiderer acknowledges support by the Austrian Science Fund (FWF) under grant P28661, by the Vienna Science and Technology Fund (WWTF) through project MA16-021, and additionally appreciates travel support through the National Science Foundation (NSF) under grant NSF-DMS-14-05210. L.C. Yeung acknowledges support under grant NSF-DMS-20-04997.

## References

- [1] Acciaio, B., Beiglböck, M., Penkner, F., Schachermayer, W. and Temme, J. (2013). A trajectorial interpretation of Doob’s martingale inequalities. *Ann. Appl. Probab.* **23** 1494–1505. [MR3098440](#) <https://doi.org/10.1214/12-aap878>
- [2] Agueh, M., Ghoussoub, N. and Kang, X. (2004). Geometric inequalities via a general comparison principle for interacting gases. *Geom. Funct. Anal.* **14** 215–244. [MR2053603](#) <https://doi.org/10.1007/s00039-004-0455-x>
- [3] Ambrosio, L., Gigli, N. and Savaré, G. (2008). *Gradient Flows in Metric Spaces and in the Space of Probability Measures*, 2nd ed. *Lectures in Mathematics ETH Zürich*. Basel: Birkhäuser. [MR2401600](#)
- [4] Backhoff, J., Conforti, G., Gentil, I. and Léonard, C. (2020). The mean field Schrödinger problem: Ergodic behavior, entropy estimates and functional inequalities. *Probab. Theory Related Fields* **178** 475–530. [MR4146543](#) <https://doi.org/10.1007/s00440-020-00977-8>
- [5] Beiglböck, M. and Siorpaes, P. (2015). Pathwise versions of the Burkholder–Davis–Gundy inequality. *Bernoulli* **21** 360–373. [MR3322322](#) <https://doi.org/10.3150/13-BEJ570>

- [6] Benachour, S., Roynette, B., Talay, D. and Vallois, P. (1998). Nonlinear self-stabilizing processes. I. Existence, invariant probability, propagation of chaos. *Stochastic Process. Appl.* **75** 173–201. [MR1632193](#) [https://doi.org/10.1016/S0304-4149\(98\)00018-0](https://doi.org/10.1016/S0304-4149(98)00018-0)
- [7] Benedetto, D., Caglioti, E., Carrillo, J.A. and Pulvirenti, M. (1998). A non-Maxwellian steady distribution for one-dimensional granular media. *J. Stat. Phys.* **91** 979–990. [MR1637274](#) <https://doi.org/10.1023/A:1023032000560>
- [8] Bolley, F., Gentil, I. and Guillin, A. (2013). Uniform convergence to equilibrium for granular media. *Arch. Ration. Mech. Anal.* **208** 429–445. [MR3035983](#) <https://doi.org/10.1007/s00205-012-0599-z>
- [9] Brenier, Y. (1991). Polar factorization and monotone rearrangement of vector-valued functions. *Comm. Pure Appl. Math.* **44** 375–417. [MR1100809](#) <https://doi.org/10.1002/cpa.3160440402>
- [10] Carmona, R. and Delarue, F. (2018). *Probabilistic Theory of Mean Field Games with Applications. I. Mean Field FBSDEs, Control, and Games. Probability Theory and Stochastic Modelling* **83**. Cham: Springer. [MR3752669](#)
- [11] Carrillo, J.A., Gvalani, R.S., Pavliotis, G.A. and Schlichting, A. (2020). Long-time behaviour and phase transitions for the McKean-Vlasov equation on the torus. *Arch. Ration. Mech. Anal.* **235** 635–690. [MR4062483](#) <https://doi.org/10.1007/s00205-019-01430-4>
- [12] Carrillo, J.A., Jüngel, A., Markowich, P.A., Toscani, G. and Unterreiter, A. (2001). Entropy dissipation methods for degenerate parabolic problems and generalized Sobolev inequalities. *Monatsh. Math.* **133** 1–82. [MR1853037](#) <https://doi.org/10.1007/s006050170032>
- [13] Carrillo, J.A., McCann, R.J. and Villani, C. (2003). Kinetic equilibration rates for granular media and related equations: Entropy dissipation and mass transportation estimates. *Rev. Mat. Iberoam.* **19** 971–1018. [MR2053570](#) <https://doi.org/10.4171/RMI/376>
- [14] Carrillo, J.A., McCann, R.J. and Villani, C. (2006). Contractions in the 2-Wasserstein length space and thermalization of granular media. *Arch. Ration. Mech. Anal.* **179** 217–263. [MR2209130](#) <https://doi.org/10.1007/s00205-005-0386-1>
- [15] Cattiaux, P. (2004). A pathwise approach of some classical inequalities. *Potential Anal.* **20** 361–394. [MR2032116](#) <https://doi.org/10.1023/B:POTA.0000009847.84908.6f>
- [16] Cattiaux, P. and Guillin, A. (2014). Semi log-concave Markov diffusions. In *Séminaire de Probabilités XLVI. Lecture Notes in Math.* **2123** 231–292. Cham: Springer. [MR3330820](#) [https://doi.org/10.1007/978-3-319-11970-0\\_9](https://doi.org/10.1007/978-3-319-11970-0_9)
- [17] Cattiaux, P., Guillin, A. and Malrieu, F. (2008). Probabilistic approach for granular media equations in the non-uniformly convex case. *Probab. Theory Related Fields* **140** 19–40. [MR2357669](#) <https://doi.org/10.1007/s00440-007-0056-3>
- [18] Cordero-Erausquin, D., Gangbo, W. and Houdré, C. (2004). Inequalities for generalized entropy and optimal transportation. In *Recent Advances in the Theory and Applications of Mass Transport. Contemp. Math.* **353** 73–94. Providence, RI: Amer. Math. Soc. [MR2079071](#) <https://doi.org/10.1090/conm/353/06433>
- [19] Léonard, C. (2014). Some Properties of Path Measures. In *Sémin. Probab. XLVI*, (C. Donati-Martin, A. Lejay and A. Rouault, eds.). *Lecture Notes in Math.* **2123** 207–230. Switzerland: Springer International Publishing.
- [20] Duong, M.H. and Tugaut, J. (2016). Stationary solutions of the Vlasov-Fokker-Planck equation: Existence, characterization and phase-transition. *Appl. Math. Lett.* **52** 38–45. [MR3416384](#) <https://doi.org/10.1016/j.aml.2015.08.003>
- [21] Duong, M.H. and Tugaut, J. (2018). The Vlasov-Fokker-Planck equation in non-convex landscapes: Convergence to equilibrium. *Electron. Commun. Probab.* **23** Paper No. 19, 10. [MR3779816](#) <https://doi.org/10.1214/18-ECP116>
- [22] Eldan, R., Lehec, J. and Shenfeld, Y. (2020). Stability of the logarithmic Sobolev inequality via the Föllmer process. *Ann. Inst. Henri Poincaré Probab. Stat.* **56** 2253–2269. [MR4116725](#) <https://doi.org/10.1214/19-AIHP1038>
- [23] Eldan, R. and Mikulincer, D. (2020). Stability of the Shannon-Stam inequality via the Föllmer process. *Probab. Theory Related Fields* **177** 891–922. [MR4126934](#) <https://doi.org/10.1007/s00440-020-00967-w>
- [24] Erbar, M., Fathi, M., Laschos, V. and Schlichting, A. (2016). Gradient flow structure for McKean-Vlasov equations on discrete spaces. *Discrete Contin. Dyn. Syst.* **36** 6799–6833. [MR3567821](#) <https://doi.org/10.3934/dcds.2016096>

- [25] Feng, J. and Kurtz, T.G. (2006). *Large Deviations for Stochastic Processes*. *Mathematical Surveys and Monographs* **131**. Providence, RI: Amer. Math. Soc. [MR2260560](#) <https://doi.org/10.1090/surv/131>
- [26] Föllmer, H. (1985). An entropy approach to the time reversal of diffusion processes. In *Stochastic Differential Systems (Marseille-Luminy, 1984)*. *Lect. Notes Control Inf. Sci.* **69** 156–163. Berlin: Springer. [MR0798318](#) <https://doi.org/10.1007/BFb0005070>
- [27] Föllmer, H. (1986). Time reversal on Wiener space. In *Stochastic Processes—Mathematics and Physics (Bielefeld, 1984)*. *Lecture Notes in Math.* **1158** 119–129. Berlin: Springer. [MR0838561](#) <https://doi.org/10.1007/BFb0080212>
- [28] Föllmer, H. (1988). Random fields and diffusion processes. In *École D’Été de Probabilités de Saint-Flour XV–XVII, 1985–87*. *Lecture Notes in Math.* **1362** 101–203. Berlin: Springer. [MR0983373](#) <https://doi.org/10.1007/BFb0086180>
- [29] Fontbona, J. and Jourdain, B. (2016). A trajectorial interpretation of the dissipations of entropy and Fisher information for stochastic differential equations. *Ann. Probab.* **44** 131–170. [MR3456334](#) <https://doi.org/10.1214/14-AOP969>
- [30] Gardiner, C. (2009). *Stochastic Methods: A Handbook for the Natural and Social Sciences*, 4th ed. *Springer Series in Synergetics*. Berlin: Springer. [MR2676235](#)
- [31] Gentil, I., Léonard, C., Ripani, L. and Tamanini, L. (2020). An entropic interpolation proof of the HWI inequality. *Stochastic Process. Appl.* **130** 907–923. [MR4046525](#) <https://doi.org/10.1016/j.spa.2019.04.002>
- [32] Haussmann, U.G. and Pardoux, É. (1986). Time reversal of diffusions. *Ann. Probab.* **14** 1188–1205. [MR0866342](#)
- [33] Herrmann, S. and Tugaut, J. (2010). Non-uniqueness of stationary measures for self-stabilizing processes. *Stochastic Process. Appl.* **120** 1215–1246. [MR2639745](#) <https://doi.org/10.1016/j.spa.2010.03.009>
- [34] Hu, K., Ren, Z., Šiška, D. and Szpruch, Ł. (2021). Mean-field Langevin dynamics and energy landscape of neural networks. *Ann. Inst. Henri Poincaré Probab. Stat.* **57** 2043–2065. [MR4328560](#) <https://doi.org/10.1214/20-aihp1140>
- [35] Jordan, R. and Kinderlehrer, D. (1996). An extended variational principle. In *Partial Differential Equations and Applications. Lecture Notes in Pure and Applied Mathematics* **177** 187–200. New York: Dekker. [MR1371591](#) <https://doi.org/10.5006/1.3292113>
- [36] Jordan, R., Kinderlehrer, D. and Otto, F. (1998). The variational formulation of the Fokker-Planck equation. *SIAM J. Math. Anal.* **29** 1–17. [MR1617171](#) <https://doi.org/10.1137/S0036141096303359>
- [37] Karatzas, I. and Kardaras, C. (2021). *Portfolio Theory and Arbitrage: A Course in Mathematical Finance. Grad. Stud. Math.* **214**. Providence, Rhode Island: Amer. Math. Soc.
- [38] Karatzas, I., Maas, J. and Schachermayer, W. (2021). Trajectorial dissipation and gradient flow for the relative entropy in Markov chains. *Commun. Inf. Syst.* **21** 481–536. [MR4273512](#) <https://doi.org/10.4310/CIS.2021.v21.n4.a1>
- [39] Karatzas, I., Schachermayer, W. and Tschiderer, B. (2020). Trajectorial Otto calculus. [arXiv:1811.08686](#).
- [40] Karatzas, I., Schachermayer, W. and Tschiderer, B. (2021). A trajectorial approach to the gradient flow properties of Langevin-Smoluchowski diffusions. *Teor. Veroyatn. Primen.* **66** 839–888. [MR4331222](#) <https://doi.org/10.4213/tvp55>
- [41] Lehec, J. (2013). Representation formula for the entropy and functional inequalities. *Ann. Inst. Henri Poincaré Probab. Stat.* **49** 885–899. [MR3112438](#) <https://doi.org/10.1214/11-aihp464>
- [42] Léonard, C. (2014). A survey of the Schrödinger problem and some of its connections with optimal transport. *Discrete Contin. Dyn. Syst.* **34** 1533–1574. [MR3121631](#) <https://doi.org/10.3934/dcds.2014.34.1533>
- [43] Lions, P.L. (2007–2008). Théorie des jeux de champ moyen et applications. Lectures at the Collège de France. <https://www.college-de-france.fr/site/pierre-louis-lions/course-2007-2008.htm>.
- [44] Liptser, R.S. and Shiryaev, A.N. (2001). *Statistics of Random Processes. I: General Theory*, expanded ed. *Applications of Mathematics (New York)* **5**. Berlin: Springer. [MR1800857](#)
- [45] Malrieu, F. (2003). Convergence to equilibrium for granular media equations and their Euler schemes. *Ann. Appl. Probab.* **13** 540–560. [MR1970276](#) <https://doi.org/10.1214/aoap/1050689593>
- [46] McCann, R.J. (1997). A convexity principle for interacting gases. *Adv. Math.* **128** 153–179. [MR1451422](#) <https://doi.org/10.1006/aima.1997.1634>
- [47] McKean, H.P. Jr. (1966). A class of Markov processes associated with nonlinear parabolic equations. *Proc. Natl. Acad. Sci. USA* **56** 1907–1911. [MR0221595](#) <https://doi.org/10.1073/pnas.56.6.1907>



- [48] Mei, S., Montanari, A. and Nguyen, P.-M. (2018). A mean field view of the landscape of two-layer neural networks. *Proc. Natl. Acad. Sci. USA* **115** E7665–E7671. [MR3845070](#) <https://doi.org/10.1073/pnas.1806579115>
- [49] Otto, F. (2001). The geometry of dissipative evolution equations: The porous medium equation. *Comm. Partial Differential Equations* **26** 101–174. [MR1842429](#) <https://doi.org/10.1081/PDE-100002243>
- [50] Otto, F. and Villani, C. (2000). Generalization of an inequality by Talagrand and links with the logarithmic Sobolev inequality. *J. Funct. Anal.* **173** 361–400. [MR1760620](#) <https://doi.org/10.1006/jfan.1999.3557>
- [51] Revuz, D. and Yor, M. (1999). *Continuous Martingales and Brownian Motion*, 3rd ed. *Grundlehren der Mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences]* **293**. Berlin: Springer. [MR1725357](#) <https://doi.org/10.1007/978-3-662-06400-9>
- [52] Risken, H. (1984). *The Fokker-Planck Equation: Methods of Solution and Applications*. *Springer Series in Synergetics* **18**. Berlin: Springer. [MR0749386](#) <https://doi.org/10.1007/978-3-642-96807-5>
- [53] Rogers, L.C.G. (1985). Smooth transition densities for one-dimensional diffusions. *Bull. Lond. Math. Soc.* **17** 157–161. [MR0806242](#) <https://doi.org/10.1112/blms/17.2.157>
- [54] Tugaut, J. (2013). Convergence to the equilibria for self-stabilizing processes in double-well landscape. *Ann. Probab.* **41** 1427–1460. [MR3098681](#) <https://doi.org/10.1214/12-AOP749>
- [55] Tugaut, J. (2014). Phase transitions of McKean-Vlasov processes in double-wells landscape. *Stochastics* **86** 257–284. [MR3180036](#) <https://doi.org/10.1080/17442508.2013.775287>
- [56] Tugaut, J. (2014). Self-stabilizing processes in multi-wells landscape in  $\mathbb{R}^d$ -invariant probabilities. *J. Theoret. Probab.* **27** 57–79. [MR3174216](#) <https://doi.org/10.1007/s10959-012-0435-2>
- [57] Veretennikov, A.Yu. (2006). On ergodic measures for McKean-Vlasov stochastic equations. In *Monte Carlo and Quasi-Monte Carlo Methods 2004* 471–486. Berlin: Springer. [MR2208726](#) [https://doi.org/10.1007/3-540-31186-6\\_29](https://doi.org/10.1007/3-540-31186-6_29)
- [58] Villani, C. (2003). *Topics in Optimal Transportation*. *Graduate Studies in Mathematics* **58**. Providence, RI: Amer. Math. Soc. [MR1964483](#) <https://doi.org/10.1090/gsm/058>
- [59] Villani, C. (2006). Mathematics of granular materials. *J. Stat. Phys.* **124** 781–822. [MR2264625](#) <https://doi.org/10.1007/s10955-006-9038-6>

Received July 2021 and revised February 2022