# GRADIENT FLOW STRUCTURE FOR MCKEAN-VLASOV EQUATIONS ON DISCRETE SPACES

MATTHIAS ERBAR, MAX FATHI, VAIOS LASCHOS, AND ANDRÉ SCHLICHTING

ABSTRACT. In this work, we show that a family of non-linear mean-field equations on discrete spaces can be viewed as a gradient flow of a natural free energy functional with respect to a certain metric structure we make explicit. We also prove that this gradient flow structure arises as the limit of the gradient flow structures of a natural sequence of N-particle dynamics, as N goes to infinity.

#### 1. Introduction

In this work, we are interested in the gradient flow structure of McKean-Vlasov equations on finite discrete spaces. They take the form

$$\dot{c}(t) = c(t)Q(c(t)) \tag{1.1}$$

where c(t) is a flow of probability measures on a fixed finite set  $\mathcal{X} = \{1, \ldots, d\}$ , and  $Q_{xy}(\mu)$  is collection of parametrized transition rates, that is for each  $\mu \in \mathcal{P}(\mathcal{X})$ ,  $Q(\mu)$  is a Markov transition kernel.

Such non-linear equations arise naturally as the scaling limit for the evolution of the empirical measure of a system of N particles undergoing a linear Markov dynamics with mean field interaction. Here the interaction is of mean field type if the transition rate  $Q_{i;x,y}^N$  for the i-th particle to jump from site x to y only depends on the empirical measure of the particles.

Mean-field systems are commonly used in physics and biology to model the evolution of a system where the influence of all particles on a tagged particle is the average of the force exerted by each particle on the tagged particle. In the recent work [7], it was shown that whenever Q satisfies suitable conditions a free energy of the form

$$\mathcal{F}(\mu) = \sum_{x \in \mathcal{X}} \mu_x \log \mu_x + \sum_{x \in \mathcal{X}} \mu_x K_x(\mu)$$
 (1.2)

for some appropriate potential  $K : \mathcal{P}(\mathcal{X}) \times \mathcal{X} \to \mathbf{R}$  (see Definition 2.3) is a Lyapunov function for the evolution equation (1.1), i.e. it decreases along the flow.

In this work, we show that this monotonicity is actually a consequence of a more fundamental underlying structure. Namely, we exibit a novel geometric structure on the space of probability measure  $\mathcal{P}(\mathcal{X})$  that allows to view the evolution equation (1.1) as the gradient flow of the free energy  $\mathcal{F}$ .

Date: October 26, 2016.

<sup>2010</sup> Mathematics Subject Classification. Primary: 60J27; Secondary: 34A34, 49J40, 49J45.

Key words and phrases. Gradient flow structure, weakly interacting particles systems, nonlinear Markov chains, McKean-Vlasov dynamics, mean-field limit, evolutionary Gamma convergence, transportation metric.

This gradient flow structure is a natural non-linear extension of the discrete gradient flow structures that were discovered in [31] and [34] in the context of linear equations describing Markov chains or more generally in [21] in the context of Lévy processes.

Moreover, we shall show that our new gradient flow structure for the non-linear equation arises as the limit of the gradient flow structures associated to a sequence of mean-field N particle Markov chains. As an application, we use the stability of gradient flows to show convergence of these mean-field dynamics to solutions of the non-linear equation (1.1), see Theorem 3.10.

1.1. **Gradient flows in spaces of probability measures.** Classically, a gradient flow is an ordinary differential equation of the form

$$\dot{x}(t) = -\nabla \mathcal{F}(x(t)).$$

By now there exists an extensive theory, initiated by De Giorgi and his collaborators [15], giving meaning to the notion of gradient flow when the curve x takes values in a metric space.

Examples of these generalized gradient flows are the famous results by Otto [28, 37] stating that many diffusive partial differential equations can be interpreted as gradient flows of appropriate energy functionals on the space of probability measures on  $\mathbb{R}^d$  equipped with the  $L^2$  Wasserstein distance. These include the Fokker-Planck and the pourous medium equations. An extensive treatment of these examples in the framework of De Giorgi was accomplished in [2].

Gradient flow structures allow to better understand the role of certain Lyapunov functionals as thermodynamic free energies. Recently, also unexpected connections of gradient flows to large deviations have been unveiled [1] [16], [19], [24], [25].

Since the heat equation is the PDE that governs the evolution of the Brownian motion, a natural question was whether a similar structure can be uncovered for reversible Markov chains on discrete spaces. This question was answered positively in works of Maas [31] and Mielke [34], which state that the evolution equations associated to reversible Markov chains on finite spaces can be reformulated as gradient flows of the entropy (with respect to the invariant measure of the chain) for a certain metric structure on the space of probability measures over the finite space. In [23], a gradient flow structure for discrete porous medium equations was also uncovered, based on similar ideas.

In Section 2, we shall highlight a gradient flow structure for (1.1), which is a natural non-linear generalization of the structure discovered in [31] and [34] for such non-linear Markov processes. This structure explains why the non-linear entropies of [7] are Lyapunov functions for the non-linear ODE. Moreover, we shall show in Section 3 that this structure is compatible with those of [31] and [34], in the sense that it arises as a limit of gradient flow structures for N-particle systems as N goes to infinity.

1.2. Convergence of gradient flows. Gradient flows have proven to be particularly useful for the study of convergence of sequences of evolution equations to some limit since they provide a very rigid structure. Informally, the philosophy can be summarized as follows: consider a sequence of gradient flows, each associated to some energy functional  $\mathcal{F}^N$  and some metric structure. If the sequence  $\mathcal{F}^N$  converges in some sense to a limit  $\mathcal{F}^\infty$  and if the metric structures converge to some limiting metric, then one would expect

the sequence of gradient flows to converge to a limit that can be described as a gradient flow of  $\mathcal{F}^{\infty}$  for the asymptotic metric structure.

There are several ways of rigorously implementing this philosophy to actually prove convergence in concrete situations. The one we shall be using in this work is due to Sandier and Serfaty in [38], and was later generalized in [40]. Other methods, based on discretization schemes, have been developed in [3] and [13]. See also the recent survey [35] for an extension of the theory to generalized gradient systems. In the context of diffusion equations, arguments similar to those of [40] have been used in [25] to study large deviations.

In the discrete setting, we can combine the framework of [31] and [34] with the method of [40] to study scaling limits of Markov chains on discrete spaces. In this work, we shall use this method to study scaling limits of N-particle mean-field dynamics on finite spaces. While the convergence statement could be obtained through more classical techniques, such as those of [36, 41], our focus here is on justifying that the gradient flow structure we present is the natural one, since it arises as the limit of the gradient flow structures for the N-particle systems.

While we were writing this article, we have been told by Maas and Mielke that they have also successfully used this technique to study the evolution of concentrations in chemical reactions. We also mention the work [27], which showed that the metric associated to the gradient flow structure for the simple random walk on the discrete torus  $\mathbf{Z}/N\mathbf{Z}$  converges to the Wasserstein structure on  $\mathcal{P}(\mathbf{T})$ , establishing compatibility of the discrete and continuous settings in a typical example. The technique can also be used to prove convergence of interacting particle systems on lattices, such as the simple exclusion process (see [26]). The technique is not restricted to the evolution of probability measures by Wasserstein-type gradient flows, but can be also applied for instance to coagulation-fragmentation processes like the Becker-Döring equations, where one can prove macroscopic limits (see [39]).

1.3. Continuous mean-field particle dynamics. Let us briefly compare the scaling limit for discrete mean-field dynamics considered in this paper with the more classical analogous scaling limit for particles in the continuum decribed by McKean-Vlasov equations.

N-particle mean-field dynamics describe the behavior of N particles given by some spatial positions  $x_1(t), \ldots, x_N(t)$ , where each particle is allowed to interact through the empirical measure of all other particles.

In nice situations, when the number of particles goes to infinity, the *empirical measure* of the system  $\frac{1}{N} \sum \delta_{x_i(t)}$  converges to some probability measure  $\mu(t)$ , whose evolution is described by a McKean-Vlasov equation. In the continuous setting, with positions in  $\mathbf{R}^d$ , this can be for example a PDE of the form

$$\partial_t \mu(t) = \Delta \mu(t) + \operatorname{div}(\mu(t)(\nabla W * \mu(t)))$$

where  $\nabla W * \mu$  is the convolution of  $\mu$  with an interaction that derives from a potential W. The according free energy in this case is given by

$$\mathcal{F}(\mu) = \begin{cases} \int \frac{d\mu}{dx}(x) \log \frac{d\mu}{dx}(x) \ dx + \frac{1}{2} \int \int W(x-y)\mu(dx)\mu(dy), & \mu \ll \mathcal{L}, \\ \infty, & \text{otherwise,} \end{cases}$$

i.e. formally  $K_x(\mu) = \frac{1}{2}(W * \mu)(x)$  in (1.2). More general PDEs, involving diffusion coefficients and confinement potentials, are also possible. We refer to [17, 41] for more information on convergence of N-particle dynamics to McKean-Vlasov equations. We also refer to [14, 12] for the large deviations behavior. An important consequence of this convergence is that, for initial conditions for which the particles are exchangeable, there is propagation of chaos: the laws of two different tagged particles become independent as the number of particles goes to infinity [41, Proposition 2.2].

It has been first noted in [9] that McKean-Vlasov equations on  $\mathbb{R}^d$  can be viewed as gradient flows of the free energy in the space of probability measures endowed with the Wasserstein metric. This fact has been useful in the study of the long-time behavior of these equations (cf. [9, 10, 11, 32] among others). The study of long-time behavior of particle systems on finite spaces has attracted recent interest (see for example [30] for the mean-field Ising model), and we can hope that curvature estimates for such systems may be useful to tackle this problem, as they have been in the continuous setting. Since lower bounds on curvature are stable, the study of curvature bounds for the mean field limit (which is defined as convexity of the free energy in the metric structure for which the dynamics is a gradient flow, see for example [22]) can shed light on this problem. We leave this issue for future work. We must also mention that Wasserstein distances have also been used to quantify the convergence of mean-field systems to their scaling limit, see for example [5].

1.4. **Outline.** In Section 2, we introduce the gradient flow structure of the mean-field system (1.1) on discrete spaces. In Section 3, we will obtain this gradient flow structure as the limit of the linear gradient flow structure associated with an N-particle system with mean-field interaction. The nonlinear gradient flow structure comes with a metric, whose properties will be studied in Section 4. We close the paper with two Appendices A and B, in which auxiliary results for the passage to the limit are provided.

## 2. Gradient flow structure of mean-field systems on discrete spaces

In this section, we derive the gradient flow formulation for the mean-field system (1.1). First, we introduce the metric concept of gradient flows in Section 2.1. Then in Section 2.2 we turn to the discrete setting. We give the precise assumptions on the Markov transition kernel  $Q(\mu)$  in the system (1.1) and state several necessary definitions for later use. In Section 2.3, we define curves of probability measures via a continuity equation and associate to them an action measuring their velocity. Based on this, we can introduce in Section 2.4 a transportation distance on the space of probability measure on  $\mathcal{X}$ . The gradient flow formulation is proven in Section 2.5 as curves of maximal slope. Finally, the gradient structure is lifted to the space of randomized probability measures in Section 2.6, which is a preparation for the passage to the limit.

2.1. **Gradient flows in a metric setting.** Let briefly recall the basic notions concerning gradient flows in metric spaces. For an extensive treatment we refer to [2].

Let (M,d) be a complete metric space. A curve  $(a,b) \ni t \mapsto u(t) \in M$  is said to be locally *p-absolutely continuous* if there exists  $m \in L^p_{loc}((a,b))$  such that

$$\forall a \le s < t \le b: \qquad d(u(s), u(t)) \le \int_s^t m(r) \, dr. \tag{2.1}$$

We write for short  $u \in AC^p_{loc}((a, b), (M, d))$ . For any such curve the metric derivative is defined by

$$|u'(t)| = \lim_{s \to t} \frac{d(u(s), u(t))}{|s - t|}.$$

The limit exists for a.e.  $t \in (a, b)$  and is the smallest m in (2.1), see [2, Thm. 1.1.2].

Now, let  $\Phi: M \to \mathbf{R}$  be lower semicontinuous function. The metric analogue of the modulus of the gradient of  $\Phi$  is given by the following definition.

**Definition 2.1** (Strong upper gradient). A function  $G: M \to [0, \infty]$ , is a strong upper gradient for  $\Phi$  if for every absolutely continuous curve  $u: (a, b) \to M$ , the function G(u) is Borel and

$$|\Phi(u(s)) - \Phi(u(t))| \le \int_s^t G(u(r))|u'|(r)dr, \quad \forall a < s \le t < b.$$

By Young's inequality, we see that the last inequality implies that

$$\Phi(u(s)) - \Phi(u(t)) \le \frac{1}{2} \int_{s}^{t} |u'|^{2}(r) dr + \frac{1}{2} \int_{s}^{t} G^{2}(u(r)) dr ,$$

for any absolutely continuous curve u provided the G is a strong upper gradient.

The following definition formalizes what it means for a curve to be a gradient flow of the function  $\Phi$  in the metric space (M, d). Shortly, it is a curve that saturates the previous inequality.

**Definition 2.2** (Curve of maximal slope). A locally absolutely continuous curve  $u:(a,b)\to M$  is called a curve of maximal slope for  $\Phi$  with respect to its strong upper gradient G if for all  $a\leq s\leq t\leq b$  we have the energy identity

$$\Phi(u(s)) - \Phi(u(t)) = \frac{1}{2} \int_{s}^{t} |u'|^{2}(r)dr + \frac{1}{2} \int_{s}^{t} G^{2}(u(r))dr.$$
 (2.2)

When  $\Phi$  is bounded below one has a convenient estimate on the modulus of continuity of a curve of maximal slope u. By Hölder's inequality and (2.2) we infer that for all s < t we have

$$d(u(s), u(t)) \le \int_{s}^{t} |u'|(r)dr \le \sqrt{t - s} \left( \int_{s}^{t} |u'|^{2}(r)dr \right)^{\frac{1}{2}}$$
  
 
$$\le \sqrt{t - s} \sqrt{2(\Phi(u(0)) - \Phi_{\min})}.$$

2.2. **Discrete setting.** Let us now introduce the setting for the discrete McKean–Vlasov equations that we consider.

In the sequel, we will denote with  $\mathcal{P}(\mathcal{X})$ , the space of probability measures on  $\mathcal{X}$ , and  $\mathcal{P}^*(\mathcal{X})$  the set of all measures that are strictly positive, i.e.

$$\mu \in \mathcal{P}^*(\mathcal{X})$$
 iff  $\forall x \in \mathcal{X} : \mu_x > 0$ ,

and finally with  $\mathcal{P}^a(\mathcal{X})$ , the set of all measures that have everywhere mass bigger than a, i.e.

$$\mu \in \mathcal{P}^a(\mathcal{X})$$
 iff  $\forall x \in \mathcal{X} : \mu_x \geq a$ .

As in [6, 7], we shall consider equations of the form (1.1) where Q is Gibbs with some potential function K. Here is the definition of such transition rates, taken from [7]:

**Definition 2.3.** Let  $K : \mathcal{P}(\mathcal{X}) \times \mathcal{X} \to \mathbf{R}$  be such that for each  $x \in \mathcal{X}, K_x : \mathcal{P}(\mathcal{X}) \to \mathbf{R}$  is a twice continuously differentiable function on  $\mathcal{P}(\mathcal{X})$ . A family of matrices  $\{Q(\mu) \in \mathbf{R}^{\mathcal{X} \times \mathcal{X}}\}_{\mu \in \mathcal{P}(\mathcal{X})}$  is Gibbs with potential function K, if for each  $\mu \in \mathcal{P}(\mathcal{X}), Q(\mu)$  is the rate matrix of an irreducible, reversible ergodic Markov chain with respect to the probability measure

$$\pi_x(\mu) = \frac{1}{Z(\mu)} \exp(-H_x(\mu)) ,$$
(2.3)

with

$$H_x(\mu) = \frac{\partial}{\partial \mu_x} U(\mu) , \quad U(\mu) = \sum_{x \in \mathcal{X}} \mu_x K_x(\mu) .$$

In particular  $Q(\mu)$  satisfies the detailed balance condition w.r.t.  $\pi(\mu)$ , that is for all  $x, y \in \mathcal{X}$ 

$$\pi_x(\mu)Q_{xy}(\mu) = \pi_y(\mu)Q_{yx}(\mu) \tag{2.4}$$

holds. Moreover, we assume that for each  $x, y \in \mathcal{X}$  the map  $\mu \mapsto Q_{xy}(\mu)$  is Lipschitz continuous over  $\mathcal{P}(\mathcal{X})$ .

In the above definition and in the following we use the convention that a function  $F(\cdot): \mathcal{P}(\mathcal{X}) \to \mathbf{R}$  is regular if and only if it can be extended in an open neighborhood of  $\mathcal{P}(\mathcal{X}) \subset \mathbf{R}^d$  in which it is regular in the usual sense. Hereby, regular could be continuous, Lipschitz or differentiable. In particular, we use this for the (twice) continuously differentiable function  $K_x$  and the Lipschitz continuous function  $Q_{xy}$  from above.

Remark 2.4. There are many ways of building a Markov kernel that is reversible with respect to a given probability measure. The most widely used method is the Metropolis algorithm, first introduced in [33]:

$$Q_{xy}^{\mathrm{MH}}(\mu) := \min\left(\frac{\pi_y(\mu)}{\pi_x(\mu)}, 1\right) = e^{-(H_y(\mu) - H_x(\mu))_+}, \quad \text{with} \quad (a)_+ := \max\left\{0, a\right\}.$$

By this choice of the rates it is only necessary to calculate  $H(\mu)$  and not the partition sum  $Z(\mu)$  in (2.3), which often is a costly computational problem.

A general scheme for obtaining rates satisfying the detailed balance condition with respect to  $\pi$  in (2.3) is to consider

$$Q_{xy}(\mu) = \frac{\sqrt{\pi_y(\mu)}}{\sqrt{\pi_x(\mu)}} A_{xy}(\mu),$$

where  $\{A(\mu)\}_{\mu\in\mathcal{P}(\mathcal{X})}$  is a family of irreducible symmetric matrices. If we choose  $A_{xy}(\mu)=\alpha_{x,y}\min\left(\frac{\sqrt{\pi_y(\mu)}}{\sqrt{\pi_x(\mu)}},\frac{\sqrt{\pi_x(\mu)}}{\sqrt{\pi_y(\mu)}}\right)$  with  $\alpha\in\{0,1\}^{\mathcal{X}\times\mathcal{X}}$  an irreducible symmetric adjacency matrix, we recover the Metropolis algorithm on the corresponding graph.

We will be interested in the non-linear evolution equation

$$\dot{c}_x(t) = \sum_{y \in \mathcal{X}} c_y(t) Q_{yx}(c(t)) , \qquad (2.5)$$

with the convention  $Q_{xx}(\mu) = -\sum_{y\neq x} Q_{xy}(\mu)$ . By the Lipschitz assumption on Q this equation has a unique solution.

One goal will be to express this evolution as the gradient flow of the associated free energy functional  $\mathcal{F}: \mathcal{P}(\mathcal{X}) \to \mathbf{R}$  defined by

$$\mathcal{F}(\mu) := \sum_{x \in \mathcal{X}} \mu_x \log \mu_x + U(\mu), \quad \text{with} \quad U(\mu) := \sum_{x \in \mathcal{X}} \mu_x K_x(\mu) . \quad (2.6)$$

To this end, it will be convenient to introduce the so-called *Onsager operator*  $\mathcal{K}(\mu)$ :  $\mathbf{R}^{\mathcal{X}} \to \mathbf{R}^{\mathcal{X}}$ . It is defined as follows:

Let  $\Lambda: \mathbf{R}_+ \times \mathbf{R}_+ \to \mathbf{R}_+$ , denote the logarithmic mean given by

$$\Lambda(s,t) := \int_0^1 s^{\alpha} t^{1-\alpha} d\alpha = \frac{s-t}{\log s - \log t} .$$

 $\Lambda$  is continuous, increasing in both variables, jointly concave and 1-homogeneous. See for example [31, 22] for more about properties of this logarithmic mean. In the sequel we are going to use the following notation

$$w_{xy}(\mu) := \Lambda(\mu_x Q_{xy}(\mu), \mu_y Q_{yx}(\mu)) \tag{2.7}$$

since this term will appear very often. By the definition of  $\Lambda$  and the properties of Q we get that  $w_{xy}$  is uniformly bounded on  $\mathcal{P}(\mathcal{X})$ , by a constant  $C_w$ .

Now, we can define

$$\mathcal{K}(\mu) := \frac{1}{2} \sum_{x,y} w_{xy}(\mu) \ (e_x - e_y) \otimes (e_x - e_y),$$
 (2.8)

where  $\{e_x\}_{x\in\mathcal{X}}$  is identified with the standard basis of  $\mathbf{R}^{\mathcal{X}}$ . More explicitly, we have for  $\psi \in \mathbf{R}^{\mathcal{X}}$ :

$$(\mathcal{K}(\mu)\psi)_x = \sum_{y} w_{xy}(\mu)(\psi_x - \psi_y).$$

With this in mind, we can formally rewrite the evolution (2.5) in gradient flow form:

$$\dot{c}(t) = -\mathcal{K}(c(t))D\mathcal{F}(c(t)), \tag{2.9}$$

where  $D\mathcal{F}(\mu) \in \mathbf{R}^{\mathcal{X}}$  is the differential of  $\mathcal{F}$  given by  $D\mathcal{F}(\mu)_x = \partial_{\mu_x} \mathcal{F}(\mu)$ .

Finally, let us introduce the Fisher information  $\mathcal{I}: \mathcal{P}(\mathcal{X}) \to [0, \infty]$  defined for  $\mu \in \mathcal{P}^*(\mathcal{X})$  by

$$\mathcal{I}(\mu) := \frac{1}{2} \sum_{(x,y) \in E_{\mu}} w_{xy}(\mu) \left( \log(\mu_x Q_{xy}(\mu)) - \log(\mu_y Q_{yx}(\mu)) \right)^2$$
 (2.10)

where, for  $\mu \in \mathcal{P}(\mathcal{X})$ , we define the edges of possible jumps by

$$E_{\mu} := \{ (x, y) \in \mathcal{X} \times \mathcal{X} : Q_{xy}(\mu) > 0 \}.$$
 (2.11)

For  $\mu \in \mathcal{P}(\mathcal{X}) \setminus \mathcal{P}^*(\mathcal{X})$  we set  $\mathcal{I}(\mu) = +\infty$ .

 $\mathcal{I}$  gives the dissipation of  $\mathcal{F}$  along the evolution, namely, if c is a solution to (2.5) then

$$\frac{d}{dt}\mathcal{F}(c(t)) = -\mathcal{I}(c(t)) .$$

2.3. Continuity equation and action. In the sequel we shall use the notation for the discrete gradient. Given a function  $\psi \in \mathbf{R}^{\mathcal{X}}$  we define  $\nabla \psi \in \mathbf{R}^{\mathcal{X} \times \mathcal{X}}$  via  $\nabla_{xy} \psi := \psi_y - \psi_x$ for  $x, y \in \mathcal{X}$ . We shall also use a notion of discrete divergence, given  $v \in \mathbf{R}^{\mathcal{X} \times \mathcal{X}}$ , we define  $\delta v \in \mathbf{R}^{\mathcal{X}}$  via  $(\delta v)_x = \frac{1}{2} \sum_{y} (v_{xy} - v_{yx})$ .

**Definition 2.5** (Continuity equation). Let T>0 and  $\mu,\nu\in\mathcal{P}(\mathcal{X})$ . A pair (c,v) is called a solution to the continuity equation, for short  $(c,v) \in \overrightarrow{CE}_T(\mu,\nu)$ , if

- (i)  $c \in C^0([0,T], \mathcal{P}(\mathcal{X}))$ , i.e.  $\forall x \in \mathcal{X} : t \mapsto c_x(t) \in C^0([0,T],[0,1])$ ;
- (ii)  $c(0) = \mu$ ;  $c(T) = \nu$ ; (iii)  $v : [0, T] \to \mathbf{R}^{\mathcal{X} \times \mathcal{X}}$  is measurable and integrable;
- (iv) The pair (c, v) satisfies the continuity equation for  $t \in (0, T)$  in the weak form, i.e. for all  $\varphi \in C_c^1((0,T), \mathbf{R})$  and all  $x \in \mathcal{X}$  holds

$$\int_0^T \left[ \dot{\varphi}(t) \ c_x(t) - \varphi(t) \ (\delta v)_x(t) \right] dt = 0.$$
 (2.12)

In a similar way, we shall write  $(c, \psi) \in CE_T(\mu, \nu)$  if  $(c, w(c)\nabla\psi) \in \overrightarrow{CE}_T(\mu, \nu)$  for  $\psi:[0,T]\to\mathbf{R}^{\mathcal{X}}$  and  $(w(c)\nabla\psi)_{xy}(t):=w_{xy}(c(t))\nabla_{xy}\psi(t)$  defined pointwise. In the case T=1 we will often neglect the time index in the notation setting  $\vec{CE}(\mu,\nu) := \vec{CE}_1(\mu,\nu)$ . Also, the endpoints  $(\mu, \nu)$  will often be suppressed in the notation.

To define the action of a curve it will be convenient to introduce the function  $\alpha$ :  $\mathbf{R} \times \mathbf{R}_+ \to \mathbf{R}_+$  defined by

$$\alpha(v,w) := \begin{cases} \frac{v^2}{w} & , w > 0\\ 0 & , v = 0 = w\\ +\infty & , \text{else} \end{cases}$$
 (2.13)

Note that  $\alpha$  is convex and lower semicontinuous.

**Definition 2.6** (Curves of finite action). Given  $\mu \in \mathcal{P}(\mathcal{X})$ ,  $v \in \mathbf{R}^{\mathcal{X} \times \mathcal{X}}$  and  $\psi \in \mathbf{R}^{\mathcal{X}}$ , we define the action of  $(\mu, \nu)$  and  $(\mu, \psi)$  via

$$\vec{\mathcal{A}}(\mu, v) := \frac{1}{2} \sum_{x,y} \alpha(v_{xy}, w_{xy}(\mu)) , \qquad (2.14)$$

$$\mathcal{A}(\mu, \psi) := \vec{\mathcal{A}}(\mu, w(\mu)\nabla\psi) = \frac{1}{2} \sum_{x,y} (\psi_y - \psi_x)^2 w_{xy}(\mu) .$$

Moreover, a solution to the continuity equation  $(c, v) \in CE_T$  is called a curve of finite action if

$$\int_0^T \vec{\mathcal{A}}(c(t),v(t)) \ dt < \infty \ .$$

It will be convenient to note that for a given solution (c, v) to the continuity equation we can find a vector field  $\tilde{v} = w(c)\nabla\psi$  of gradient form such that  $(c, \tilde{v})$  still solves the continuity equation and has lower action.

**Proposition 2.7** (Gradient fields). Let  $(c, v) \in \vec{\mathrm{CE}}_T(\mu, \nu)$  be a curve of finite action, then there exists  $\psi : [0, T] \to \mathbb{R}^{\mathcal{X}}$  measurable such that  $(c, \psi) \in \mathbf{CE}_T(\mu, \nu)$  and

$$\int_0^T \mathcal{A}(c(t), \psi(t)) dt \le \int_0^T \vec{\mathcal{A}}(c(t), v(t)) dt.$$
(2.15)

*Proof.* Given  $c \in \mathcal{P}(\mathcal{X})$  we will endow  $\mathbf{R}^{\mathcal{X} \times \mathcal{X}}$  with the weighted inner product

$$\langle \Psi, \Phi \rangle_{\mu} := \frac{1}{2} \sum_{x,y} \Psi_{xy} \Phi_{xy} w_{xy}(\mu) ,$$

such that  $\vec{\mathcal{A}}(\mu, v) = |\Psi|_{\mu}^2$  if  $v_{xy} = w_{xy}\Psi_{xy}$ . Denote by  $\operatorname{Ran}(\nabla) := \{\nabla \psi : \psi \in \mathbf{R}^{\mathcal{X}}\} \subset \mathbf{R}^{\mathcal{X} \times \mathcal{X}}$  the space of gradient fields. Moreover, denote by

$$\operatorname{Ker}(\nabla_{\mu}^{*}) := \left\{ \Psi \in \mathbf{R}^{\mathcal{X} \times \mathcal{X}} : \sum_{x,y} (\Psi_{yx} - \Psi_{xy}) w_{xy}(\mu) = 0 \right\}$$

the space of divergence free vector fields. Note that we have the orthogonal decomposition

$$\mathbf{R}^{\mathcal{X} \times \mathcal{X}} = \operatorname{Ker}(\nabla_{\mu}^*) \oplus^{\perp} \operatorname{Ran}(\nabla) \ .$$

Now, given  $(c, v) \in \vec{\mathrm{CE}}_T(\mu, \nu)$ , we have  $\vec{\mathcal{A}}(c(t), v(t)) < \infty$  for a.e.  $t \in [0, T]$ . Thus, from (2.13) we see that for a.e. t and all x, y we have that  $v_{xy}(t) = 0$  whenever  $w_{xy}(c(t)) = 0$ . Hence, we can define

$$\Psi_{xy}(t) := \frac{v_{xy}(t)}{w_{xy}(c(t))} \quad \text{for a.e. } t \in [0, T].$$

Then  $\psi:[0,T]\to\mathbf{R}^{\mathcal{X}}$  can be given by setting  $\nabla\psi(t)$  to be the orthogonal projection of  $\Psi(t)$  onto  $\mathrm{Ran}(\nabla)$  w.r.t.  $\langle\cdot,\cdot\rangle_{c(t)}$ . The orthogonal decomposition above then implies immediately that  $(c,w\nabla\psi)\in \vec{\mathrm{CE}}_T(\mu,\nu)$  and that  $|\nabla\psi(t)|^2_{c(t)}\leq |\Psi(t)|^2_{c(t)}=\vec{\mathcal{A}}(c(t),v(t))$  for a.e.  $t\in[0,T]$ . This yields (2.15).

2.4. **Metric.** We shall now introduce a new transportation distance on the space  $\mathcal{P}(\mathcal{X})$ , which will provide the underlying geometry for the gradient flow interpretation of the mean field evolution equation (1.1).

**Definition 2.8** (Transportation distance). Given  $\mu, \nu \in \mathcal{P}(\mathcal{X})$ , we define

$$W^{2}(\mu,\nu) := \inf \left\{ \int_{0}^{1} \mathcal{A}(c(t),\psi(t)) \ dt : (c,\psi) \in CE_{1}(\mu,\nu) \right\}.$$
 (2.16)

**Remark 2.9.** As a consequence of Proposition 2.7 and the fact that for any  $\mu \in \mathcal{P}(\mathcal{X})$  and  $\psi \in \mathbf{R}^{\mathcal{X}}$  give rise to  $v \in \mathbf{R}^{\mathcal{X} \times \mathcal{X}}$  via  $v_{xy} = w_{xy}(\mu) \nabla_{xy} \psi$  such that  $\mathcal{A}(\mu, \psi) = \vec{\mathcal{A}}(\mu, v)$  we obtain an equivalent reformulation of the function  $\mathcal{W}$ :

$$\mathcal{W}^2(\mu,\nu) = \inf \left\{ \int_0^1 \vec{\mathcal{A}}(c(t),v(t)) \ dt : (c,v) \in \vec{\mathrm{CE}}_1(\mu,\nu) \right\}.$$

It turns out that W is indeed a distance.

**Proposition 2.10.** The function W defined in Definition 2.8 is a metric and the metric space  $(\mathcal{P}(\mathcal{X}), \mathcal{W})$  is separable and complete. Moreover, any two points  $\mu, \nu \in \mathcal{P}(\mathcal{X})$  can be joined by a constant speed W-geodesic, i.e. there exists a curve  $(\gamma_t)_{t \in [0,1]}$  with  $\gamma_0 = \mu$  and  $\gamma_1 = \nu$  satisfying  $W(\gamma_s, \gamma_t) = |t - s|W(\mu, \nu)$  for all  $s, t \in [0,1]$ .

We defer the proof of this statement until Section 4. Let us give a characterization of absolutely continuous curves w.r.t. W.

**Proposition 2.11.** A curve  $c:[0,T] \to \mathcal{P}(\mathcal{X})$  is absolutely continuous w.r.t.  $\mathcal{W}$  if and only if there exists  $\psi:[0,T] \times \mathcal{X} \to \mathbf{R}$  such that  $(c,\psi) \in \mathrm{CE}_T$ , and  $\int_0^T \sqrt{\mathcal{A}(c(t),\psi(t))} \, dt < \infty$ . Moreover, we can choose  $\psi$  such that the metric derivative of c is given as  $|c'(t)| = \sqrt{\mathcal{A}(c(t),\psi(t))}$  for a.e. t.

*Proof.* The proof is identical to the one of [18, Thm. 5.17].  $\Box$ 

2.5. **Gradient flows.** In this section, we shall present the interpretation of the discrete McKean-Vlasov equation as a gradient flow with respect to the distance  $\mathcal{W}$ . We will use the abstract framework introduced in Section 2.1 above, where  $(M, d) = (\mathcal{P}(\mathcal{X}), \mathcal{W})$  and  $\Phi = \mathcal{F}$ .

**Lemma 2.12.** Let  $\mathcal{I}: \mathcal{P}(\mathcal{X}) \to [0, \infty]$  defined in (2.10) denote the Fisher information and let  $\mathcal{F}: \mathcal{P}(\mathcal{X}) \to \mathbf{R}$  defined in (2.6) denote the free energy. Then,  $\sqrt{\mathcal{I}}$  is a strong upper gradient for  $\mathcal{F}$  on  $(\mathcal{P}(\mathcal{X}), \mathcal{W})$ , i.e. for  $(c, \psi) \in \text{CE}_T$  and  $0 \le t_1 < t_2 \le T$  holds

$$|\mathcal{F}(c(t_2)) - \mathcal{F}(c(t_1))| \le \int_{t_1}^{t_2} \sqrt{\mathcal{I}(c(t))} \sqrt{\mathcal{A}(c(t), \psi(t))} dt.$$
 (2.17)

*Proof.* Let  $c:(a,b) \to (\mathcal{P}(\mathcal{X}),\mathcal{W})$  be a  $\mathcal{W}$ -absolutely continuous curve with  $\psi$  the associated gradient potential such that  $(c,\psi) \in \mathrm{CE}(c(a),c(b))$  and  $|c'|(t) = \sqrt{\mathcal{A}(c(t),\psi(t))}$  for a.e.  $t \in (a,b)$ . We can assume w.l.o.g. that the r.h.s. of (2.17) is finite. For the proof, we are going to define the auxiliary functions

$$\mathcal{F}_{\delta}(\mu) = \sum_{x \in \mathcal{X}} (\mu_x + \delta) \log(\mu_x + \delta) + U(\mu).$$

The function  $\mathcal{F}_{\delta}(\mu)$  are Lipschitz continuous and converge uniformly to  $\mathcal{F}$ , as  $\delta \to 0$ . By Lemma 4.2, c is also absolutely continuous with respect to Euclidean distance. Therewith, since  $\mathcal{F}_{\delta}$  are Lipschitz continuous, we have that  $\mathcal{F}_{\delta}(c):(a,b)\to\mathbb{R}$  is absolutely continuous and hence

$$\mathcal{F}_{\delta}(c(t_2)) - \mathcal{F}_{\delta}(c(t_1)) = \int_{t_1}^{t_2} \frac{d}{dt} \mathcal{F}_{\delta}(c)(t) dt = \int_{t_1}^{t_2} D\mathcal{F}_{\delta}(c(t)) \dot{c}(t) dt,$$

where  $D\mathcal{F}_{\delta}(c(t))$  is well-defined for a.e. t and given in terms of

$$\partial_{e_x} \mathcal{F}_{\delta}(c(t)) = H_x(c(t)) + \log(c_x(t) + \delta) + 1.$$

Now, we have by using the Cauchy-Schwarz inequality

$$\int_{t_1}^{t_2} |D\mathcal{F}_{\delta}(c(t))\dot{c}(t)| dt \leq \int_{t_1}^{t_2} \left| \frac{1}{2} \sum_{x,y \in \mathcal{X}} (\psi_x - \psi_y) (\partial_{e_x} \mathcal{F}_{\delta}(c) - \partial_{e_y} \mathcal{F}_{\delta}(c)) w_{xy}(c) \right| dt$$

$$\leq \int_{t_1}^{t_2} \sqrt{\frac{1}{2} \sum_{x,y \in \mathcal{X}} (\nabla_{xy} \psi)^2 w_{xy}(c)} \sqrt{\frac{1}{2} \sum_{x,y \in \mathcal{X}} (\partial_{e_x} \mathcal{F}_{\delta}(c) - \partial_{e_y} \mathcal{F}_{\delta}(c))^2 w_{xy}(c)} dt$$

$$= \int_{t_1}^{t_2} \sqrt{\mathcal{A}(c, \psi)} \times \sqrt{\frac{1}{2} \sum_{x,y \in \mathcal{X}} (\log(c_x + \delta) + H_x(c) - \log(c_y + \delta) - H_y(c))^2 w_{xy}(c)} dt.$$

$$\leq \int_{t_1}^{t_2} \sqrt{\mathcal{A}(c, \psi)} \sqrt{2\mathcal{I}(c) + C_H^2} \sum_{x,y \in \mathcal{X}} w_{xy}(c) dt,$$

where we dropped the t-dependence on c and  $\psi$ . For the last inequality, we observe that since  $H_x$  for  $x \in \mathcal{X}$  are uniformly bounded and for  $a < b, \delta > 0$  it holds  $\frac{b}{a} \ge \frac{b+\delta}{a+\delta}$ , it is easy to see that the quantity  $|\log(c_x(t) + \delta) + H_x(c(t)) - \log(c_y(t) + \delta) - H_y(c(t))|$  is bounded by  $|\log(c_x(t)) + H_x(c(t)) - \log(c_y(t)) - H_y(c(t))| + C_H$  with  $C_H$  only depending on H. Moreover, we observe that by definitions of  $H_x$  and  $\pi$  from (2.3), it holds

$$\begin{aligned} &|\log(\mu_x) + H_x(\mu) - \log(\mu_y) - H_y(\mu)| = \left|\log\frac{\mu_x}{\pi_x(\mu)} - \log\frac{\mu_y}{\pi_y(\mu)}\right| \\ &= \left|\log(\mu_x Q_{xy}(\mu)) - \log(\pi_x(\mu) Q_{xy}(\mu)) - \log(\pi_y(\mu) Q_{yx}(\mu)) - \log(\mu_y Q_{yx}(\mu))\right|. \end{aligned}$$

Then, by the detailed balance condition (2.4) the two middle terms cancel and we arrive at  $\mathcal{I}(\mu)$ . Since, we assumed  $\int_{t_1}^{t_2} \sqrt{\mathcal{A}(c(t), \psi(t))} \sqrt{\mathcal{I}(c(t))} dt$  to be finite, we can apply the dominated convergence theorem and get the conclusion.

**Proposition 2.13.** For any absolutely continuous curve  $(c(t))_{t\in[0,T]}$  in  $\mathcal{P}(\mathcal{X})$  holds

$$\mathcal{J}(c) := \mathcal{F}(c(T)) - \mathcal{F}(c(0)) + \frac{1}{2} \int_0^T \mathcal{I}(c(t)) \ dt + \frac{1}{2} \int_0^T \mathcal{A}(c(t), \psi(t)) \ dt \ge 0$$
 (2.18)

Moreover, equality is attained if and only if  $(c(t))_{t\in[0,T]}$  is a solution to (1.1). In this case  $c(t) \in \mathcal{P}^*(\mathcal{X})$  for all t > 0.

In other words, solution to (1.1) are the only gradient flow curves (i.e. curves of maximal slope) of  $\mathcal{F}$ .

*Proof.* The first statement follows as above by Young's inequality from the fact that  $\mathcal{I}$  is strong upper gradient for  $\mathcal{F}$ .

Now let us assume that for a curve c,  $\mathcal{J}(c) \leq 0$  holds. Then since (2.17) holds for every curve we can deduce that we actually have

$$\mathcal{F}(c(t_2)) - \mathcal{F}(c(t_1)) + \frac{1}{2} \int_{t_1}^{t_2} \mathcal{I}(c(t)) dt + \frac{1}{2} \int_{t_1}^{t_2} \mathcal{A}(c(t), \psi(t)) dt = 0, \quad 0 \le t_1 \le t_2 \le T.$$

Since  $\int_0^T \mathcal{I}(c(t))dt < \infty$ , we can find a sequence  $\epsilon_n$ , converging to zero, such that  $\mathcal{I}(c(\epsilon_n)) < \infty$ . By continuity of c, we can find  $a, \epsilon > 0$ , such that  $c(t) \in \mathcal{P}^a(\mathcal{X})$ , for

 $t \in [\epsilon_n, \epsilon_n + \epsilon]$ . Now, since  $\mathcal{I}$  is Lipschitz in  $\mathcal{P}^a(\mathcal{X})$ , we can apply the chain rule for  $\epsilon_n \leq t_1 \leq t_2 \leq \epsilon_n + \epsilon$  and get

$$\mathcal{F}(c(t_1)) - \mathcal{F}(c(t_2)) = \int_{t_1}^{t_2} \langle D\mathcal{F}(c(t)), \mathcal{K}(c(t)) \nabla \psi(t) \rangle$$
$$= \frac{1}{2} \int_{t_1}^{t_2} \mathcal{A}(c(t), \psi(t)) dt + \frac{1}{2} \int_{t_1}^{t_2} \mathcal{I}(c(t)) dt,$$

by comparison we get

$$\langle D\mathcal{F}(c(t)), \mathcal{K}(c(t))\nabla\psi(t)\rangle = \sqrt{\mathcal{A}(c(t), \psi(t)) \mathcal{I}(c(t))} = \mathcal{A}(c(t), \psi(t)) = \mathcal{I}(c(t)),$$

for  $t \in [\epsilon_n, \epsilon_n + \epsilon]$ . From which, by an application of the inverse of the Cauchy-Schwarz inequality, we get that  $\psi_x(t) - \psi_y(t) = \partial_{e_x} \mathcal{F}(c(t)) - \partial_{e_y} \mathcal{F}(c(t))$ . Now we have

$$\dot{c}(t) = -\mathcal{K}(c(t))D\mathcal{F}(c(t)) 
= -\frac{1}{2} \sum_{x,y} w_{xy}(c(t)) \left( e_x - e_y \right) \left( \partial_{e_x} \mathcal{F}(c(t)) - \partial_{e_y} \mathcal{F}(c(t)) \right) 
= -\frac{1}{2} \sum_{x,y} \left( Q_{xy}(c(t)) c_x(t) - Q_{yx}(c(t)) c_y(t) \right) \left( e_x - e_y \right) 
= -\sum_{x} \left( \sum_{y} \left( Q_{xy}(c(t)) c_x(t) - Q_{yx}(c(t)) c_y(t) \right) \right) e_x = c(t) Q(c(t))$$
(2.19)

on the interval  $[\epsilon_n, \epsilon_n + \epsilon]$ . We actually have that c(t) is a solution to  $\dot{c}(t) = c(t)Q(c(t))$  on  $[\epsilon_n, T]$ . Indeed, let  $T_n = \sup\{t' \leq T : \dot{c}(t) = c(t)Q(c(t)), \forall t \in [\epsilon_n, t']\}$ . We have  $c(T_n) \in \mathcal{P}^b(\mathcal{X})$ , for some b > 0, because c is a solution to  $\dot{c}(t) = c(t)Q(c(t))$ , on  $[\epsilon_n, T_n)$  and the dynamics are irreducible. Now if we apply the same argument for  $T_n$ , that we used for  $\epsilon_n$ , we can extent the solution beyond  $T_n$ . If  $T_n < T$ , then we will get a contradiction, Therefore  $T_n = T$ . Now by sending  $\epsilon_n$  to zero we get that c is a solution to  $\dot{c}(t) = c(t)Q(c(t))$ , on [0,T].

Now on the other hand if c is a solution to  $\dot{c}(t) = c(t)Q(c(t))$  on [0,T], we can get that for every  $\epsilon > 0$ , there exists a > 0, such that  $c(t) \in \mathcal{P}^a(\mathcal{X})$  on  $[\epsilon,T]$ . The choice  $\psi(t) = DF(t)$ , satisfies the continuity equation (see (2.19)), and by applying the chain rule, we get that

$$\mathcal{F}(c(T)) - \mathcal{F}(c(\epsilon)) + \frac{1}{2} \int_{\epsilon}^{T} \mathcal{I}(c(t))dt + \frac{1}{2} \int_{\epsilon}^{T} \mathcal{A}(c(t), \psi(t))dt = 0.$$

Sending  $\epsilon$  to zero concludes the proof.

**Remark 2.14.** Note that the formulation above contains the usual entropy entropy-production relation for gradient flows. If c is a solution to (1.1), then  $\psi(t) = -D\mathcal{F}(c(t))$  and especially it holds that  $\mathcal{A}(c(t), -D\mathcal{F}(c(t))) = \mathcal{I}(c(t))$ . Therewith, (2.18) becomes

$$\mathcal{F}(c(T)) + \int_0^T \mathcal{I}(c(t)) dt = \mathcal{F}(c(0)).$$

2.6. Lifted dynamics on the space of random measures. It is possible to lift the evolution  $\dot{c}(t) = c(t)Q(c(t))$  in  $\mathcal{P}(\mathcal{X})$  to an evolution for measures  $\mathbb{C}$  on  $\mathcal{P}(\mathcal{X})$ . This is convenient, if one does not want to start from a deterministic point but consider random initial data. The evolution is then formally given by

$$\partial_t \mathbb{C}(t,\nu) + \operatorname{div}_{\mathcal{P}(\mathcal{X})} \left( \mathbb{C}(t,\nu) \left( \nu Q(\nu) \right) \right) = 0, \quad \text{with} \quad \operatorname{div}_{\mathcal{P}(\mathcal{X})} = \sum_{x \in \mathcal{X}} \partial_{e_x}.$$
 (2.20)

**Notation.** In the following, all quantities connected to the space  $\mathcal{P}(\mathcal{P}(\mathcal{X}))$  will be denoted by blackboard-bold letters, like for instance random probability measures  $\mathbb{M} \in \mathcal{P}(\mathcal{P}(\mathcal{X}))$  or functionals  $\mathbb{F} : \mathcal{P}(\mathcal{P}(\mathcal{X})) \to \mathbf{R}$ .

The evolution (2.20) also has a natural gradient flow structure that is obtained by lifting the gradient flow structure of the underlying dynamics. In fact, (2.20) will turn out to be a gradient flow w.r.t. to the classical  $L^2$ -Wasserstein distance on  $\mathcal{P}(\mathcal{P}(\mathcal{X}))$ , which is build from the distance  $\mathcal{W}$  on the base space  $\mathcal{P}(\mathcal{X})$ . To establish this gradient structure, we need to introduce lifted analogues of the continuity equation and the action of a curve as well as a probabilistic representation result for the continuity equation.

**Definition 2.15** (Lifted continuity equation). A pair  $(\mathbb{C}, \mathbb{V})$  is called a solution to the lifted continuity equation, for short  $(\mathbb{C}, \mathbb{V}) \in \vec{\mathbb{CE}}_T(\mathbb{M}, \mathbb{N})$ , if

- (i)  $[0,T] \ni t \mapsto \mathbb{C}(t) \in \mathcal{P}(\mathcal{P}(\mathcal{X}))$  is weakly\* continuous,
- (ii)  $\mathbb{C}(0) = \mathbb{M}; \mathbb{C}(T) = \mathbb{N};$
- (iii)  $\mathbb{V}: [0,T] \times \mathcal{P}(\mathcal{X}) \to \mathbf{R}^{\mathcal{X} \times \mathcal{X}}$  is measurable and integrable w.r.t.  $\mathbb{C}(t,d\mu)dt$ ,
- (iv) The pair  $(\mathbb{C}, \mathbb{V})$  satisfies the continuity equation for  $t \in (0, T)$  in the weak form, i.e. for all  $\varphi \in C_c^1((0, T) \times \mathcal{P}(\mathcal{X}))$  holds

$$\int_{0}^{T} \int_{\mathcal{P}(\mathcal{X})} (\dot{\varphi}(t,\nu) - \langle \nabla \varphi(t,\nu), \delta \mathbb{V}(t,\nu) \rangle) \, \mathbb{C}(t,d\nu) \, dt = 0, \tag{2.21}$$

where  $\delta \mathbb{V} : \mathcal{P}(\mathcal{X}) \to \mathbf{R}^{\mathcal{X} \times \mathcal{X}}$ , is given by  $\delta \mathbb{V}(\nu)_x := \frac{1}{2} \sum_{u} (\mathbb{V}_{xy}(\nu) - \mathbb{V}_{yx}(\nu))$ .

Here we consider  $\mathcal{P}(\mathcal{X})$  as a subset of Euclidean space  $\mathbf{R}^{\mathcal{X}}$  with  $\langle \cdot, \cdot \rangle$  the usual inner product. In particular,  $\nabla \varphi(t, \mu) = (\partial_{\mu_x} \varphi(t, \mu))_{x \in \mathcal{X}}$  denotes the usual gradient on  $\mathbf{R}^{\mathcal{X}}$  and we have explicitly

$$\begin{split} \langle \nabla \varphi(t,\nu), \delta \mathbb{V}(t,\nu) \rangle &= \sum_{x \in \mathcal{X}} \partial_{\mu_x} \varphi(t,\mu) (\delta \mathbb{V}(t,\mu))_x \\ &= \frac{1}{2} \sum_{x,y \in \mathcal{X}} \partial_{\mu_x} \varphi(t,\mu) \Big( \mathbb{V}_{xy}(t,\mu) - \mathbb{V}_{yx}(t,\mu) \Big). \end{split}$$

Thus, (2.21) is simply the weak formulation of the classical continuity equation in  $\mathbf{R}^{\mathcal{X}}$ . In a similar way, we shall write  $(\mathbb{C}, \Psi) \in \mathbb{CE}_T(\mathbb{M}, \mathbb{N})$  if  $\Psi : [0, T] \times \mathcal{P}(\mathcal{X}) \to \mathbf{R}^{\mathcal{X}}$  is a function such that  $(\mathbb{C}, \tilde{\mathbb{V}}) \in \widetilde{\mathbb{CE}}_T(\mathbb{M}, \mathbb{N})$  with  $\tilde{\mathbb{V}}_{xy}(t, \mu) = w_{xy}(\mu)\nabla_{xy}\Psi(t, \mu)$ . In this case we have that  $\delta \mathbb{V}(\mu) = \mathcal{K}(\mu)\Psi(\mu)$ , where  $\mathcal{K}(\mu)$  is the Onsager operator defined in (2.8). Solutions to (2.20) are understood as weak solutions like in Definition (2.15). That is  $\mathbb{C}$  is a weak solution to (2.20) if  $(\mathbb{C}, \Psi^*) \in \mathbb{CE}_T$  with  $\Psi^*(\nu) := -D\mathcal{F}(\nu)$ . This leads, via the formal calculation

$$\delta \mathbb{V}^*(\nu) := \mathcal{K}(\nu) \mathbb{\Psi}^*(\nu) = -\mathcal{K}(\nu) D \mathcal{F}(\nu) = \nu Q(\nu),$$

to the formulation: For all  $\varphi \in C_c^1([0,T] \times \mathcal{P}(\mathcal{X}))$  we have

$$\int_0^T \int_{\mathcal{P}(\mathcal{X})} (\dot{\varphi}(t,\nu) - \langle \nabla \varphi(t,\nu), \nu Q(\nu) \rangle) \, \mathbb{C}(t,d\nu) \, dt = 0 \,. \tag{2.22}$$

By the Lipschitz assumption on Q the vector field  $\nu Q(\nu)$  given by the components  $(\nu Q(\nu))_x = \sum_y \nu_y Q_{xy}(\nu)$  is also Lipschitz. Then standard theory implies that equation (2.22) has a unique solution (cf. [2, Chapter 8]).

**Definition 2.16** (Lifted action). Given  $\mathbb{M} \in \mathcal{P}(\mathcal{P}(\mathcal{X}))$ ,  $\mathbb{V} : \mathcal{P}(\mathcal{X}) \to \mathbf{R}^{\mathcal{X} \times \mathcal{X}}$  and  $\mathbb{V} : \mathcal{P}(\mathcal{X}) \to \mathbf{R}^{\mathcal{X}}$ , we define the action of  $(\mathbb{M}, \mathbb{V})$  and  $(\mathbb{M}, \mathbb{V})$  by

$$\begin{split} \vec{\mathbb{A}}(\mathbb{M},\mathbb{V}) &:= \int_{\mathcal{P}(\mathcal{X})} \vec{\mathcal{A}}(\nu,\mathbb{V}(\nu)) \ \mathbb{M}(d\nu) \ , \\ \mathbb{A}(\mathbb{M},\mathbb{\Psi}) &:= \int_{\mathcal{P}(\mathcal{X})} \mathcal{A}(\nu,\mathbb{\Psi}(\nu)) \ \mathbb{M}(d\nu) \ . \end{split}$$

The next result tell us that is is sufficient to consider only gradient vector fields. It is the analog of Proposition 2.7.

**Proposition 2.17** (Gradient fields for Liouville equation). If  $(\mathbb{C}, \mathbb{V}) \in \mathbb{CE}_T$  is a curve of finite action, then there exists  $\Psi : [0, T] \times \mathcal{P}(\mathcal{X}) \to \mathbf{R}^{\mathcal{X}}$  measurable such that  $(\mathbb{C}, \Psi) \in \mathbb{CE}_T$  and

$$\int_{0}^{T} \mathbb{A}(\mathbb{C}(t), \Psi(t)) dt \le \int_{0}^{T} \vec{\mathbb{A}}(\mathbb{C}(t), \mathbb{V}(t)) dt. \tag{2.23}$$

*Proof.* Given a solution  $(\mathbb{C}, \mathbb{V}) \in \mathbb{CE}_T$ , for each t and  $\nu \in \mathcal{P}(\mathcal{X})$  we apply the contruction in the proof of Proposition 2.7 to  $\mathbb{V}(\nu)$  to obtain  $\mathbb{V}(t,\nu)$  with  $\mathcal{A}(\nu,\mathbb{V}(t,\nu)) \leq \vec{\mathcal{A}}(\nu,\mathbb{V}(t,\nu))$ . It is readily checked that  $(\mathbb{C},\mathbb{V}) \in \mathbb{CE}_T$ . Integration against  $\mathbb{C}$  and dt yields (2.23).  $\square$ 

**Definition 2.18** (Lifted distance). Given  $\mathbb{M}, \mathbb{N} \in \mathcal{P}(\mathcal{P}(\mathcal{X}))$  we define

$$\mathbb{W}^2(\mathbb{M}, \mathbb{N}) := \inf \left\{ \int_0^1 \mathbb{A}(\mathbb{C}(t), \Psi(t)) \ dt : (\mathbb{C}, \Psi) \in \mathbb{CE}_1(\mathbb{M}, \mathbb{N}) \right\} \ . \tag{2.24}$$

Analogously to Remark 2.9 we obtain an equivalent formulation of  $\mathbb{W}$ :

$$\mathbb{W}^{2}(\mathbb{M}, \mathbb{N}) = \inf \left\{ \int_{0}^{1} \vec{\mathbb{A}}(\mathbb{C}(t), \mathbb{V}(t)) \ dt : (\mathbb{C}, \mathbb{V}) \in \vec{\mathbb{CE}}_{1}(\mathbb{M}, \mathbb{N}) \right\} . \tag{2.25}$$

The following result is a probabilistic representation via characteristics for the continuity equation. It is a variant of [2, Prop. 8.2.1] adapted to our setting.

**Proposition 2.19.** For a given  $\mathbb{M}, \mathbb{N} \in \mathcal{P}(\mathcal{P}(\mathcal{X}))$  let  $(\mathbb{C}, \Psi) \in \mathbb{CE}_T(\mathbb{M}, \mathbb{N})$  be a solution of the continuity equation with finite action.

Then there exists a probability measure  $\Theta$  on  $\mathcal{P}(\mathcal{X}) \times \mathrm{AC}([0,T],\mathcal{P}(\mathcal{X}))$  such that:

(1) Any  $(\mu, c) \in \text{supp }\Theta$  is a solution of the ODE

$$\dot{c}(t) = \mathcal{K}(c(t))\Psi(t, c(t)) \quad \text{for a.e. } t \in [0, T] ,$$
  
$$c(0) = \mu .$$

(2) For any  $\varphi \in C_b^0(\mathcal{P}(\mathcal{X}))$  and any  $t \in [0,T]$  holds

$$\int_{\mathcal{P}(\mathcal{X})} \varphi(\nu) \ \mathbb{C}(t, d\nu) = \int_{\mathcal{P}(\mathcal{X}) \times AC([0, T], \mathcal{P}(\mathcal{X}))} \varphi(c(t)) \ \Theta(d\mu_0, dc). \tag{2.26}$$

Conversely any  $\Theta$  satisfying (1) and

$$\int_{\mathcal{P}(\mathcal{X})\times\mathrm{AC}([0,T],\mathcal{P}(\mathcal{X}))}\int_{0}^{T}\mathcal{A}\left(c(t),\Psi(t,c(t))\right)dt\ \Theta(d\mu,dc)<\infty$$

induces a family of measures  $\mathbb{C}(t)$  via (2.26) such that  $(\mathbb{C}, \Psi) \in \mathbb{CE}_T(\mathbb{M}, \mathbb{N})$ .

We will also use the measure  $\bar{\Theta}$  on  $AC([0,T],\mathcal{P}(\mathcal{X}))$  given by

$$\bar{\Theta}(dc) = \int_{\mathcal{P}(\mathcal{X})} \Theta(d\mu, dc) . \qquad (2.27)$$

Therewith, note that (2.26) can be rewritten as the pushforward  $\mathbb{C}(t) = (e_t)_{\#}\bar{\Theta}$  under the evaluation map  $e_t : AC([0,T],\mathcal{P}(\mathcal{X})) \ni c \mapsto c(t) \in \mathcal{P}(\mathcal{X})$  defined for any  $\varphi \in C_b^0(\mathcal{P}(\mathcal{X}))$  by

$$\int_{\mathcal{P}(\mathcal{X})} \varphi(\nu) \ \mathbb{C}(t, d\nu) = \int_{\mathrm{AC}([0, T], \mathcal{P}(\mathcal{X}))} \varphi(c(t)) \ \bar{\Theta}(dc).$$

Proof. Let  $(\mathbb{C}, \Psi) \in \mathbb{C}\mathbb{E}_T(\mathbb{M}, \mathbb{N})$  be a solution of the continuity equation with finite action. Define  $\mathbb{V}: [0,T] \times \mathcal{P}(\mathcal{X}) \to \mathbf{R}^{\mathcal{X} \times \mathcal{X}}$  via  $\mathbb{V}_{xy}(t,\nu) = w_{xy}(\nu) \nabla_{xy} \Psi(t,\nu)$  and note that  $\delta \mathbb{V}(t,\nu) = \mathcal{K}(\nu) \Psi(t,\nu)$ . We view  $\mathcal{P}(\mathcal{X})$  as a subset of  $\mathbf{R}^{\mathcal{X}}$  and  $\delta \mathbb{V}$  as a time-dependent vector field on  $\mathbf{R}^{\mathcal{X}}$  and note that  $(\mathbb{C}, \delta \mathbb{V})$  is a solution to the classical continuity equation in weak form

$$\int_0^T \int_{\mathbf{R}^{\mathcal{X}}} \left( \dot{\varphi}(t, \nu) - \nabla \varphi(t, \nu) \delta \mathbb{V}(t, \nu) \right) \mathbb{C}(t, d\nu) \ dt = 0$$

for all  $\varphi \in C_c^1((0,T) \times \mathbf{R}^{\mathcal{X}})$ . Moreover, note that for any  $\Psi \in \mathbf{R}^{\mathcal{X}}$  we have by Jensens inequality

$$\begin{aligned} |\mathcal{K}(\nu)\Psi|^2 &= \sum_{x \in \mathcal{X}} \left| \sum_{y \in \mathcal{X}} w_{xy}(\nu) (\Psi_x - \Psi_y) \right|^2 \\ &\leq \sum_{x,y \in \mathcal{X}} C_w w_{xy}(\nu) \left| (\Psi_x - \Psi_y) \right|^2 = C_w \mathcal{A}(\nu, \Psi(\nu)) , \end{aligned}$$

with

$$C_w := \max_{x,y \in \mathcal{X}} \sup_{\nu \in \mathcal{P}(\mathcal{X})} w_{xy}(\nu) = \max_{x,y \in \mathcal{X}} \sup_{\nu \in \mathcal{P}(\mathcal{X})} \Lambda(\nu_x Q_{xy}(\nu), \nu_y Q_{yx}(\nu)).$$

Since  $Q: \mathcal{P}(\mathcal{X}) \to \mathbf{R}_+$  is continuous,  $C_w$  is finite. This yields the integrability estimate

$$\int_0^T \int_{\mathbb{R}^d} |\delta \mathbb{V}(t,\nu)|^2 \ d\mathbb{C}(t,\nu) \ dt \le C \int_0^T \mathbb{A}(\mathbb{C}(t), \Psi(t)) \ dt < \infty \ .$$

Now, by the representation result [2, Proposition 8.2.1] for the classical continuity equation there exists a probability measure  $\Theta$  on  $\mathbf{R}^{\mathcal{X}} \times \mathrm{AC}([0,T],\mathbf{R}^{\mathcal{X}})$  such that any  $(\mu,c) \in \mathrm{supp}\,\Theta$  satisfies  $c(0) = \mu$  and  $\dot{c}(t) = \delta \mathbb{V}(t,c(t))$  in the sense weak sense (2.12) and moreover, (2.26) holds with  $\mathcal{P}(\mathcal{X})$  replaced by  $\mathbf{R}^{\mathcal{X}}$ . Since,  $\mathbb{C}(t)$  is supported on  $\mathcal{P}(\mathcal{X})$  we find that  $\Theta$  is actually a measure on  $\mathcal{P}(\mathcal{X}) \times \mathrm{AC}([0,T],\mathcal{P}(\mathcal{X}))$ , where absolute continuity

is understood still w.r.t. the Euclidean distance. To see that  $\Theta$  is the desired measure it remains to check that for  $\Theta$ -a.e.  $(\mu, c)$  we have that c is a curve of finite action. But this follows by observing that (2.26) implies

$$\int \int_0^T \mathcal{A}(c(t), \Psi(t, c(t))) \ dt \ \Theta(d\mu, dc) = \int_0^T \mathbb{A}(\mathbb{C}(t), \Psi(t)) \ dt < \infty \ .$$

This finishes the proof of the first statement.

The converse, statement follows in the same way as in [2, Proposition 8.2.1].

**Proposition 2.20** (Identification with Wasserstein distance). The distance  $\mathbb{W}$  defined in (2.24) coincides with the  $L^2$ -Wasserstein distance on  $\mathcal{P}(\mathcal{P}(\mathcal{X}))$  w.r.t. the distance  $\mathbb{W}$  on  $\mathcal{P}(\mathcal{X})$ . More precisely, for  $\mathbb{M}, \mathbb{N} \in \mathcal{P}(\mathcal{P}(\mathcal{X}))$  there holds

$$\mathbb{W}^2(\mathbb{M},\mathbb{N}) = W^2_{\mathcal{W}}(\mathbb{M},\mathbb{N}) := \inf_{\mathbb{G} \in \Pi(\mathbb{M},\mathbb{N})} \left\{ \int_{\mathcal{P}(\mathcal{X}) \times \mathcal{P}(\mathcal{X})} \mathcal{W}^2(\mu,\nu) \ d\mathbb{G}(\mu,\nu) \right\},$$

where  $\Pi(\mathbb{M}, \mathbb{N})$  is the set of all probability measures on  $\mathcal{P}(\mathcal{X}) \times \mathcal{P}(\mathcal{X})$  with marginals  $\mathbb{M}$  and  $\mathbb{N}$ .

*Proof.* We first show the inequality " $\geq$ ". For  $\varepsilon > 0$  let  $(\mathbb{C}, \Psi)$  be a solution to the continuity equation such that  $\int_0^1 \mathbb{A}(\mathbb{C}(t), \Psi(t)) dt \leq \mathbb{W}^2(\mathbb{M}, \mathbb{N}) + \varepsilon$  and let  $\bar{\Theta}$  be the measure on  $\mathrm{AC}([0,T],\mathcal{P}(\mathcal{X}))$  given by the previous Proposition. Then we obtain a coupling  $\mathbb{G} \in \Pi(\mathbb{M}, \mathbb{N})$  by setting  $\mathbb{G} = (e_0, e_1)_{\#}\bar{\Theta}$ . This yields

$$\begin{split} W^2_{\mathcal{W}}(\mathbb{M},\mathbb{N}) &\leq \int \mathcal{W}^2(\mu,\nu) \ d\mathbb{G}(\mu,\nu) = \int \mathcal{W}^2(c(0),c(1))d\bar{\Theta} \\ &\leq \int \int_0^1 \mathcal{A}(c(t),\Psi(t,c(t)))dtd\bar{\Theta}(c) = \int_0^1 \mathbb{A}(\mathbb{C}(t),\Psi(t))dt \leq \mathbb{W}^2(\mathbb{M},\mathbb{N}) + \varepsilon \ . \end{split}$$

Since  $\varepsilon$  was arbitrary this yields the inequality " $\geq$ ".

To prove the converse inequality " $\leq$ ", fix an optimal coupling  $\mathbb{G}$ , fix  $\varepsilon > 0$  and choose for  $\mathbb{G}$ -a.e.  $(\mu, \nu)$  a couple  $(c^{\mu,\nu}, v^{\mu,\nu}) \in \overrightarrow{\mathrm{CE}}_1(\mu, \nu)$  such that

$$\int_0^1 \vec{\mathcal{A}}(c^{\mu,\nu}(t), v^{\mu,\nu}(t)) dt \le \mathcal{W}(\mu,\nu) + \varepsilon.$$

Now, define a family of measures  $\mathbb{C}:[0,1]\to\mathcal{P}(\mathcal{P}(\mathcal{X}))$  and a family of vector-valued measures  $V:[0,1]\to\mathcal{P}(\mathcal{P}(\mathcal{X});\mathbf{R}^{\mathcal{X}\times\mathcal{X}})$  via

$$d\mathbb{C}(t,\tilde{\nu}) = \int_{\mathcal{P}(\mathcal{X})\times\mathcal{P}(\mathcal{X})} d\delta_{c^{\mu,\nu}(t)}(\tilde{\nu}) d\mathbb{G}(\mu,\nu) ,$$

$$V(t,\tilde{\nu}) = \int_{\mathcal{P}(\mathcal{X})\times\mathcal{P}(\mathcal{X})} v^{\mu,\nu}(t) d\delta_{c^{\mu,\nu}(t)}(\tilde{\nu}) d\mathbb{G}(\mu,\nu) .$$

Note that  $V(t) \ll \mathbb{C}(t)$  and define  $\mathbb{V} : [0,1] \times \mathcal{P}(\mathcal{X}) \to \mathbf{R}^{\mathcal{X} \times \mathcal{X}}$  as the density of V w.r.t.  $\mathbb{C}$ . By linearity of the continuity equations have that  $(\mathbb{C}, \mathbb{V}) \in \mathbb{CE}_1(\mathbb{M}, \mathbb{N})$ . Moreover, we find

$$\begin{split} \int_0^1 \vec{\mathbb{A}}(\mathbb{C}(t), \mathbb{V}(t)) dt &= \int_0^1 \int \vec{\mathcal{A}}(c^{\mu,\nu}(t), v^{\mu,\nu}(t)) d\mathbb{G}(\mu, \nu) dt \\ &\leq \int \mathcal{W}^2(\mu, \nu) d\mathbb{G}(\mu, \nu) + \varepsilon = W_{\mathcal{W}}^2(\mathbb{M}, \mathbb{N}) + \varepsilon \;. \end{split}$$

Since  $\varepsilon$  was arbitrary, in view of (2.25) this finishes the proof.

Finally, we can obtain a gradient flow structure for the Liouville equation (2.20) in a straightforward manner by averaging the gradient flow structure of the underlying dynamical system.

To this end, given  $\mathbb{M} \in \mathcal{P}(\mathcal{P}(\mathcal{X}))$  define the free energy by

$$\mathbb{F}(\mathbb{M}) := \int_{\mathcal{P}(\mathcal{X})} \mathcal{F}(\nu) \, \mathbb{M}(d\nu) \;,$$

and define the Fisher information by

$$\mathbb{I}(\mathbb{M}) := \mathbb{A}(\mathbb{M}, -D\mathbb{F}) = \int_{\mathcal{P}(\mathcal{X})} \mathcal{I}(\nu) \ \mathbb{M}(d\nu).$$

**Proposition 2.21** (Gradient flow structure for Liouville equation). The Liouville equation (2.20) is the gradient flow of  $\mathbb{F}$  w.r.t.  $\mathbb{W}$ . Moreover precisely,  $\sqrt{\mathbb{I}}$  is a strong upper gradient for  $\mathbb{F}$  on the metric space  $(\mathcal{P}(\mathcal{P}(\mathcal{X})), \mathbb{W})$  and the curves of maximal slope are precisely the solutions to (2.20). In other words, for any absolutely continuous curve  $\mathbb{C}$  in  $\mathcal{P}(\mathcal{P}(\mathcal{X}))$  holds

$$\mathbb{J}(\mathbb{C}) := \mathbb{F}(\mathbb{C}(T)) - \mathbb{F}(\mathbb{C}(0)) + \frac{1}{2} \int_0^T \mathbb{I}(\mathbb{C}(t)) \ dt + \frac{1}{2} \int_0^T \mathbb{A}(\mathbb{C}(t), \Psi(t)) \ dt \ge 0 \qquad (2.28)$$

with  $(\mathbb{C}(t), \Psi(t)) \in \mathbb{C}\mathbb{E}_T$ . Moreover,  $\mathbb{J}(\mathbb{C}) = 0$  if and only if  $\mathbb{C}$  solves (2.22).

*Proof.* Let  $\bar{\Theta}$  be the disintegration of  $\mathbb{C}$  from Proposition 2.19 defined in (2.27). The fact that  $\sqrt{\mathbb{I}}$  is a strong upper gradient of  $\mathbb{F}$  can be seen by integrating its defining inequality on the underlying level (2.17) w.r.t.  $\bar{\Theta}$ 

$$|\mathbb{F}(\mathbb{C}(t_2)) - \mathbb{F}(\mathbb{C}(t_1))| \leq \int_{AC([0,T];\mathcal{P}(\mathcal{X}))} |\mathcal{F}(c(t_2)) - \mathcal{F}(c(t_1))| \ \bar{\Theta}(dc)$$
  
$$\leq \int_{AC([0,T];\mathcal{P}(\mathcal{X}))} \int_{t_1}^{t_2} \sqrt{\mathcal{I}(c(t))} \sqrt{\mathcal{A}(c(t),\psi(t))} \ dt \ \bar{\Theta}(dc).$$

Then, using Jensen's inequality on the concave function  $(a,b) \mapsto \sqrt{ab}$ , we get the strong upper gradient property for  $\sqrt{\mathbb{I}}$ .

The "if" part of the last claim is easily verified from the definition. Now, assume that  $\mathbb{J}(\mathbb{C})=0$ . Since  $\mathbb{C}$  is absolutely continuous, we can apply Proposition 2.19 and obtain the probabilistic representation  $\bar{\Theta}\in\mathcal{P}\left(\mathrm{AC}([0,T]\times\mathcal{P}(\mathcal{X}))\right)$  (2.27) such that  $\mathbb{C}(t)=(e_t)_\#\bar{\Theta}$ . Then, (2.28) can be obtained by just integrating  $\mathcal{J}$  from (2.18) along  $\bar{\Theta}$  and it holds  $\mathcal{J}(c)=0$  for  $\bar{\Theta}$ -a.e.  $c\in\mathrm{AC}([0,T],\mathcal{P}(\mathcal{X}))$ . These c are by Proposition 2.13 solutions to (2.9) and satisfy  $c(t)\in\mathcal{P}^*(\mathcal{X})$  for all t>0. Then, we can conclude by the converse statement of Proposition 2.19 that  $\mathcal{K}(c(t))\Psi(t,c(t))=-\mathcal{K}(c(t))D\mathcal{F}(c(t))$ , which implies since  $c(t)\in\mathcal{P}^*(\mathcal{X})$  for t>0 up to a constant that  $\Psi(t,c(t))=-D\mathcal{F}(c(t))$  and hence  $\mathbb{C}$  solves (2.22).

### 3. From weakly interacting particle systems to mean field systems

In this section, we will show how the gradient flow structure we described in the previous sections arises as the limit of gradient flow structures for N-particle systems with mean field interactions, in the limit  $N \to \infty$ . Moreover, we show that the empirical distribution of the N-particle dynamics converges to a solution of the non-linear equation (1.1).

**Notation.** For N an integer bold face letters are elements connected to the space  $\mathcal{X}^N$  and hence implicitly depending on N. Examples are vectors  $\boldsymbol{x}, \boldsymbol{y} \in \mathcal{X}^N$ , matrices  $\boldsymbol{Q} \in \mathbf{R}^{\mathcal{X}^N \times \mathcal{X}^N}$  or measures  $\boldsymbol{\mu} \in \mathcal{P}(\mathcal{X}^N)$ . For  $i \in \{1, \dots, N\}$  let  $\boldsymbol{e}^i$  be the placeholder for i-th particle, such that  $\boldsymbol{x} \cdot \boldsymbol{e}^i = x_i \in \mathcal{X}$  is the position of the i-th particle. For  $\boldsymbol{x} \in \mathcal{X}^N$  and  $\boldsymbol{y} \in \mathcal{X}$  we denote by  $\boldsymbol{x}^{i;y}$  the particle system obtained from  $\boldsymbol{x}$  where the i-th particle jumped to site  $\boldsymbol{y}$ 

$$\mathbf{x}^{i;y} := \mathbf{x} - (x_i - y)\mathbf{e}^i = (x_1, \dots, x_{i-1}, y, x_{i+1}, \dots, x_N).$$

 $L^{N}(\boldsymbol{x})$  will denote the empirical distribution for  $\boldsymbol{x} \in \mathcal{X}^{N}$ , defined by

$$L^{N}(\boldsymbol{x}) := \frac{1}{N} \sum_{i=1}^{N} \delta_{x_{i}}$$

$$(3.1)$$

We introduce the discretized simplex  $\mathcal{P}_N(\mathcal{X}) \subset \mathcal{P}(\mathcal{X})$ , given by

$$\mathcal{P}_N(\mathcal{X}) := \left\{ L^N(\boldsymbol{x}) : \boldsymbol{x} \in \mathcal{X}^N \right\}.$$

Let us introduce a natural class of mean-field dynamics for the N-particle system. We follow the standard procedure outlined in Remark 2.4.

In analog to Definition 2.3, we fix  $K : \mathcal{P}(\mathcal{X}) \times \mathcal{X} \to \mathbf{R}$  such that for each  $x \in \mathcal{X}, K_x$  is a twice continuously differentiable function on  $\mathcal{P}(\mathcal{X})$  and set  $U(\mu) := \sum_{x \in \mathcal{X}} \mu_x K_x(\mu)$ . For every natural number N define the probability measure  $\boldsymbol{\pi}^N$  for  $\boldsymbol{x} \in \mathcal{X}^N$  by

$$\pi_{\boldsymbol{x}}^{N} := \frac{1}{Z^{N}} \exp\left(-NU\left(L^{N}\boldsymbol{x}\right)\right),$$
 (3.2)

and  $\mathbf{Z}^N := \sum_{\mathbf{x} \in \mathcal{X}^N} \exp\left(-NU\left(L^N\mathbf{x}\right)\right)$  is the partition sum. This shall be the invariant measure of the particle system and is already of mean-field form.

To introduce the dynamics, we use a family  $\left\{A^N(\mu) \in \mathbf{R}^{\mathcal{X} \times \mathcal{X}}\right\}_{\mu \in \mathcal{P}_N(\mathcal{X})}$  of irreducible

To introduce the dynamics, we use a family  $\{A^N(\mu) \in \mathbf{R}^{\mathcal{X} \times \mathcal{X}}\}_{\mu \in \mathcal{P}_N(\mathcal{X})}$  of irreducible symmetric matrices and define the rate matrices  $\{Q^N(\mu) \in \mathbf{R}^{\mathcal{X} \times \mathcal{X}}\}_{\mu \in \mathcal{P}_N(\mathcal{X})}$  for any  $\mathbf{x} \in \mathcal{X}^N, y \in \mathcal{X}, i \in \{1, ... N\}$  by

$$Q_{x_{i},y}^{N}(L^{N}\boldsymbol{x}) := \sqrt{\frac{\boldsymbol{\pi}_{\boldsymbol{x}^{i;y}}^{N}}{\boldsymbol{\pi}_{\boldsymbol{x}}^{N}}} A_{x,y}^{N}(L^{N}\boldsymbol{x})$$

$$= \exp\left(-\frac{N}{2}\left(U(L^{N}\boldsymbol{x}^{i;y}) - U(L^{N}\boldsymbol{x})\right)\right) A_{x,y}^{N}(L^{N}\boldsymbol{x}).$$
(3.3)

Finally, the actual rates of the N-particle system are given in terms of the rate matrix

$$\boldsymbol{Q}^{N} \in \mathbf{R}^{\mathcal{X}^{N} \times \mathcal{X}^{N}} : \qquad \boldsymbol{Q}_{\boldsymbol{x}, \boldsymbol{x}^{i;y}}^{N} := Q_{x_{i}, y}^{N}(L^{N}(\boldsymbol{x})). \tag{3.4}$$

By construction  $Q^N$  is irreducible and reversible w.r.t. the unique invariant measure  $\pi^N$ .

Remark 3.1. The irreducible family of matrices  $\{A^N(\mu)\}_{\mu\in\mathcal{P}(\mathcal{X})}$  encodes the underlying graph structure of admissible jumps and also the rates of the jumps. For instance,  $A^N_{x,y}(\nu) = \alpha_{x,y}$  for any symmetric adjacency matrix  $\alpha \in \{0,1\}^{\mathcal{X} \times \mathcal{X}}$  corresponds to Glauber dynamics on the corresponding graph. Another choice is  $A^N_{x,y}(L^N x) := \exp\left(-\frac{N}{2}\left|U(L^N x) - U(L^N x^{i;y})\right|\right)$ , which corresponds to Metropolis dynamics on the complete graph. In particular, all of these examples satisfy Assumption 1.

**Assumption 1** (Lipschitz assumptions on rates). There exists a family of irreducible symmetric matrices  $\{A(\mu)\}_{\mu\in\mathcal{P}(\mathcal{X})}$  such that  $\mu\mapsto A(\mu)$  is Lipschitz continuous on  $\mathcal{P}(\mathcal{X})$  and the family  $\{A^N(\mu)\}_{\mu\in\mathcal{P}(\mathcal{X}),N\in\mathbf{N}}$  of irreducible symmetric matrices satisfies

$$\forall x, y \in \mathcal{X}: A_{x,y}^N \to A_{x,y} \quad on \, \mathcal{P}(\mathcal{X}) \, as \, N \to \infty.$$

**Lemma 3.2.** Assume  $\{Q^N(\mu) \in \mathbf{R}^{\mathcal{X} \times \mathcal{X}}\}_{\mu \in \mathcal{P}_N(\mathcal{X})}$  is given by (3.3) with  $A^N$  satisfying Assumption 1, then for all  $x, y \in \mathcal{X}$ 

$$Q_{x,y}^N \to Q_{x,y} \quad on \mathcal{P}(\mathcal{X})$$
 (3.5)

with  $Q_{x,y}(\mu) = \sqrt{\frac{\pi_y(\mu)}{\pi_x(\mu)}} A_{x,y}(\mu)$  with  $\pi$  given in (2.3). In particular,  $\mu \mapsto Q_{x,y}(\mu)$  is Lipschitz continuous on  $\mathcal{P}(\mathcal{X})$  for all  $x, y \in \mathcal{X}$ .

*Proof.* By [6, Lemma 4.1] holds for  $x \in \mathcal{X}^N$  with  $\mu = L^N x$ ,  $y \in \mathcal{X}$  and  $i \in \{1, ..., N\}$ 

$$\frac{\boldsymbol{\pi}_{\boldsymbol{x}^{i:y}}^{N}}{\boldsymbol{\pi}_{\boldsymbol{x}}^{N}} = N\left(U(L^{N}\boldsymbol{x}) - U(L^{N}\boldsymbol{x}^{i;y})\right) = \partial_{\mu_{\boldsymbol{x}}}U(\mu) - \partial_{\mu_{\boldsymbol{y}}}U(\mu) + O(N^{-1})$$
$$= \frac{\pi_{\boldsymbol{y}}(\mu)}{\pi_{\boldsymbol{x}}(\mu)} + O(N^{-1}),$$

which shows by Assumption 1 the convergence statement. The Lipschitz continuity follows, since A is assumed Lipschitz and the function  $\mu \mapsto \partial_{\mu_x} U(\mu) = \mu_x + \sum_y \mu_y \partial_{\mu_x} K_y(\mu)$  is continuously differentiable, since K is assumed twice continuously differentiable.  $\square$ 

Remark 3.3. The mean-field behavior is manifested in the convergence statement (3.5). The typical example we have in mind, as presented in Section 4 of [6], is as follows: the mean-field model is described by

$$K_x(\mu) := V(x) + \sum_y W(x, y) \mu_y$$

where V is a potential energy and W an interaction energy between particles on sites x and y. For the N particle system, we can use a Metropolis dynamics, where possible jumps are those between configurations that differ by the position of a single particle, and reversible with respect to the measure

$$\pi_{\boldsymbol{x}}^N = A^{-1} \exp(-U^N(\boldsymbol{x})); \quad U^N(\boldsymbol{x}) := \sum_{i=1}^N V(x_i) + \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^N W(x_i, x_j).$$

Note, the by the definition of the  $L^N$ , we have the identity

$$U^{N}(\boldsymbol{x}) = NU(L^{N}\boldsymbol{x})$$
 with  $U(\mu) = \sum_{x} \mu_{x} K_{x}(\mu)$ ,

which makes it consistent with Definition 3.2. This is a typical class of mean-field spin systems from statistical mechanics.

In particular, the Curie-Weiss mean-field spin model for ferromagnetism is obtained by choosing  $\mathcal{X} = \{-, +\}$ , V(-) = V(+) = W(-, -) = W(+, +) = 0 and  $W(-, +) = W(+, -) = \beta > 0$ . This is among the simplest models of statistical mechanics showing a phase transition in the free energy

$$\mathcal{F}_{\beta}(\mu) := \sum_{\sigma \in \{-,+\}} (\log \mu_{\sigma} + K_{\sigma}(\mu)) \, \mu_{\sigma} = \mu_{-} \log \mu_{-} + \mu_{+} \log \mu_{+} + 2\beta \mu_{-} \mu_{+}$$

at  $\beta = 1$ . For  $\beta \leq 1$  the free energy is convex whereas for  $\beta > 1$  it is non-convex on  $\mathcal{P}(\mathcal{X})$ . We will investigate this phase transition on the level of curvature for the mean-field system as well as for the finite particle system in future work.

3.1. Gradient flow structure of interacting N-particle systems. The N-particle dynamics on  $\mathcal{X}^N$  is now defined by the rate matrix  $\mathbf{Q}^N$  given as in (3.4) with the generator

$$\mathcal{L}^{N} f := \sum_{i=1}^{N} \sum_{y \in \mathcal{X}} (f(\boldsymbol{x}^{i;y}) - f(\boldsymbol{x})) \boldsymbol{Q}_{\boldsymbol{x}, \boldsymbol{x}^{i;y}}^{N}.$$
(3.6)

Likewise the evolution of an initial density  $\mu_0 \in \mathcal{P}(\mathcal{X}^N)$  satisfies

$$\dot{\boldsymbol{c}}(t) = \boldsymbol{c}(t)\boldsymbol{Q}^{N}. \tag{3.7}$$

Since by construction the rate matrix  $Q^N$  defined in (3.4) satisfy the detailed balance condition w.r.t.  $\pi^N$  (3.2), this is the generator of a reversible Markov process w.r.t.  $\pi^N$  on the finite space  $\mathcal{X}^N$ . Hence, we can use the framework developed in [31] and [34] to view this dynamics as a gradient flow of the relative entropy with respect to its invariant measure. Let us introduce the relevant quantities.

We define the relative entropy  $\mathcal{H}(\boldsymbol{\mu} \mid \boldsymbol{\pi}^N)$  for  $\boldsymbol{\mu}, \boldsymbol{\pi}^N \in \mathcal{P}(\mathcal{X}^N)$  by setting

$$\boldsymbol{\mathcal{F}}^N(\boldsymbol{\mu}) := \mathcal{H}(\boldsymbol{\mu} \mid \boldsymbol{\pi}^N) = \sum_{\boldsymbol{x} \in \mathcal{X}^N} \boldsymbol{\mu_x} \log \frac{\boldsymbol{\mu_x}}{\boldsymbol{\pi}_{\boldsymbol{x}}^N} \;.$$

Furthermore we define the action of a pair  $\mu \in \mathcal{P}(\mathcal{X}^N)$ ,  $\psi \in \mathbf{R}^{\mathcal{X}^N}$  by

$$\mathcal{A}^N(\boldsymbol{\mu},\boldsymbol{\psi}) = \frac{1}{2} \sum_{\boldsymbol{x},\boldsymbol{y}} (\boldsymbol{\psi}_{\boldsymbol{y}} - \boldsymbol{\psi}_{\boldsymbol{x}})^2 \boldsymbol{w}_{\boldsymbol{x},\boldsymbol{y}}^N(\boldsymbol{\mu}),$$

where the weights  $\boldsymbol{w}_{\boldsymbol{x},\boldsymbol{y}}^{N}(\boldsymbol{\mu})$  are defined like in (2.7) as follows

$$\boldsymbol{w}_{\boldsymbol{x},\boldsymbol{y}}^{N}(\boldsymbol{\mu}) := \Lambda \left( \boldsymbol{\mu}_{\boldsymbol{x}} \boldsymbol{Q}^{N}(\boldsymbol{x},\boldsymbol{y}), \boldsymbol{\mu}_{\boldsymbol{y}} \boldsymbol{Q}^{N}(\boldsymbol{y},\boldsymbol{x}) \right). \tag{3.8}$$

Then, a distance  $\mathcal{W}^N$  on  $\mathcal{P}(\mathcal{X}^N)$  is given by

$$\mathcal{W}^{N}(\boldsymbol{\mu}, \boldsymbol{\nu})^{2} := \inf_{(\boldsymbol{c}(t), \boldsymbol{\psi}(t))} \int_{0}^{1} \mathcal{A}^{N}(\boldsymbol{c}(t), \boldsymbol{\psi}(t)) dt$$
 (3.9)

where the infimum runs over all pairs such that c is a path from  $\mu$  to  $\nu$  in  $\mathcal{P}(\mathcal{X}^N)$ , and such that the continuity equation

$$\dot{c}_{x}(t) + \sum_{y} (\psi_{y}(t) - \psi_{x}(t)) w_{x,y}^{N}(c(t)) = 0$$
 (3.10)

holds. For details of the construction and the proof that this defines indeed a distance we refer to [31]. In particular, we note that for any absolutely continuous curve  $c:[0,T]\to (\mathcal{P}(\mathcal{X}^N),\mathcal{W}^N)$  there exist a function  $\psi:[0,T]\to \mathbf{R}^{\mathcal{X}^N}$  such that the continuity equation (3.10) holds.

Finally, we define the N-particle Fisher information by

$$\mathcal{I}^{N}(\boldsymbol{\mu}) := \begin{cases} \frac{1}{2} \sum_{(\boldsymbol{x}, \boldsymbol{y}) \in E_{\boldsymbol{\mu}}} \boldsymbol{w}_{\boldsymbol{x}\boldsymbol{y}}^{N}(\boldsymbol{\mu}) (\log(\boldsymbol{\mu}_{\boldsymbol{x}} \boldsymbol{Q}_{\boldsymbol{x}\boldsymbol{y}}^{N}(\boldsymbol{\mu})) - \log(\boldsymbol{\mu}_{\boldsymbol{y}} \boldsymbol{Q}_{\boldsymbol{y}\boldsymbol{x}}^{N}(\boldsymbol{\mu})))^{2} & \boldsymbol{\mu} \in \mathcal{P}^{*}(\mathcal{X}^{N}) \\ \infty. & \text{otherwise} \end{cases}$$

We formulate the statement that (3.7) is the gradient flow of  $\mathcal{F}^N$  w.r.t.  $\mathcal{W}^N$  again in terms of curves of maximal slope.

**Proposition 3.4.** For any absolutely continuous curve  $c : [0,T] \to (\mathcal{P}(\mathcal{X}^N), \mathcal{W}^N)$  the function  $\mathcal{J}^N$  given by

$$\mathcal{J}^{N}(\boldsymbol{c}) := \mathcal{F}^{N}(\boldsymbol{c}(T)) - \mathcal{F}^{N}(\boldsymbol{c}(0)) + \frac{1}{2} \int_{0}^{T} \mathcal{I}^{N}(\boldsymbol{c}(t)) + \mathcal{A}^{N}(\boldsymbol{c}(t), \boldsymbol{\psi}(t)) dt$$
(3.11)

is non-negative, where  $\psi_t$  is such that the continuity equation (3.10) holds. Moreover, a curve  $\mathbf{c}$  is a solution to  $\dot{\mathbf{c}}(t) = \mathbf{c}(t)\mathbf{Q}^N$  if and only if  $\mathbf{\mathcal{J}}^N(\mathbf{c}) = 0$ .

*Proof.* The proof is exactly the same as for Proposition 2.13, so we omit it.  $\Box$ 

3.2. Convergence of gradient flows. In this section we prove convergence of the empirical distribution of the N-particle system (3.7) to a solution of the non-linear equation (1.1). This will be done by using the gradient flow structure exibited in the previous sections together with the techniques developed in [40] on convergence of gradient flows.

Heuristically, consider a sequence of gradient flows associated to a senquence of metric spaces and engergy functionals. Then to prove convergence of the flows it is sufficient to establish convergence of the metrics and the energy functionals in the sense that functionals of the type (3.11) satisfy a suitable notion of  $\Gamma$  –  $\lim$  inf estimate.

In the following Theorem 3.6 we adapt the result in [40] to our setting.

We consider the sequence of metric spaces  $S^N := (\mathcal{P}(\mathcal{X}^N), \mathcal{W}^N)$  with  $\mathcal{W}^N$  defined in (3.9) and the limiting metric space  $\mathbb{S} := (\mathcal{P}(\mathcal{P}(\mathcal{X})), \mathbb{W})$  with  $\mathbb{W}$  defined in (2.24). The following notion of convergence will provide the correct topology in our setting.

**Definition 3.5** (Convergence of random measures). A sequence  $\boldsymbol{\mu}^N \in \mathcal{P}(\mathcal{X}^N)$  converges in  $\tau$  topology to a point  $\mathbb{M} \in \mathcal{P}(\mathcal{P}(\mathcal{X}))$  if and only if  $L_{\#}^N(\boldsymbol{\mu}^N) \in \mathcal{P}(\mathcal{P}(\mathcal{X}))$  converges in distribution to  $\mathbb{M}$ , where  $L^N: \mathcal{X}^N \to \mathcal{P}(\mathcal{X})$  is defined in (3.1). Likewise, for  $(\boldsymbol{c}^N(t))_{t \in [0,T]}$  with  $\boldsymbol{c}_t^N \in \mathcal{P}(\mathcal{X}^N)$ :  $\boldsymbol{c}^N \stackrel{\tau}{\to} \mathbb{C}$  if for all  $t \in [0,T]$ ,  $L_{\#}^N \boldsymbol{c}^N(t) \to \mathbb{C}(t)$ .

**Theorem 3.6** (Convergence of gradient flows à la [40]). Assume there exists a topology  $\tau$  such that whenever a sequence  $\mathbf{c}^N \in \mathrm{AC}([0,T],\mathbf{S}^N)$  converges pointwise w.r.t.  $\tau$  to a limit  $\mathbb{C} \in \mathrm{AC}([0,T],\mathbb{S})$ , then this convergence is compatible with the energy functionals, that is

$$c^N \stackrel{\tau}{\to} \mathbb{C} \qquad \Rightarrow \qquad \liminf_{N \to \infty} \frac{1}{N} \mathcal{F}^N(c^N(T)) \ge \mathbb{F}(\mathbb{C}(T)) - \mathcal{F}_0,$$
 (3.12)

for some finite constant  $\mathcal{F}_0 \in \mathbf{R}$ . In addition, assume the following holds

(1) lim inf-estimate of metric derivatives:

$$\liminf_{N \to \infty} \frac{1}{N} \int_0^T \mathcal{A}^N(\boldsymbol{c}^N(t), \boldsymbol{\psi}^N(t)) dt \ge \int_0^T \mathbb{A}(\mathbb{C}(t), \mathbb{\Psi}(t)) dt, \tag{3.13}$$

where  $(\mathbf{c}^N, \psi^N)$  and  $(\mathbb{C}(t), \Psi(t))$  are related via the respective continuity equations in  $\mathbf{S}^N$  and  $\mathbb{S}$ .

(2)  $\liminf$ -estimate of the slopes pointwise in  $t \in [0, T]$ :

$$\lim_{N \to \infty} \inf_{N} \frac{1}{N} \mathcal{I}^{N}(\boldsymbol{c}^{N}(t)) \ge \mathbb{I}(\mathbb{C}(t)). \tag{3.14}$$

Let  $\mathbf{c}^N$  be a curve of maximal slope on (0,T) for  $\mathbf{\mathcal{J}}^N$  (3.11) such that  $\mathbf{c}^N(0) \stackrel{\tau}{\to} \mathbb{C}(0)$  which is well-prepared in the sense that  $\lim_{N\to\infty} \mathbf{\mathcal{F}}^N(\mathbf{c}^N(0)) = \mathbb{F}(\mathbb{C}(0))$ . Then  $\mathbb{C}$  is a curve of maximal slope for  $\mathbb{J}$  (2.28) and

$$\forall t \in [0,T), \quad \lim_{N \to \infty} \frac{1}{N} \mathcal{F}^N(\boldsymbol{c}^N(t)) = \mathbb{F}(\mathbb{C}(t))$$

$$\frac{1}{N} \mathcal{A}^N(\boldsymbol{c}^N, \boldsymbol{\psi}^N) \to \mathbb{A}(\mathbb{C}, \mathbb{\Psi}) \quad in \quad L^2[0,T]$$

$$\frac{1}{N} \mathcal{I}^N(\boldsymbol{c}^N) \to \mathbb{I}(\mathbb{C}) \quad in \quad L^2[0,T]$$

*Proof.* Let us sketch the proof. The assumptions (3.12), (3.13), (3.14) and the well-preparedness of the initial data allow to pass in the limit in the individual terms of  $\frac{1}{N}\mathcal{J}^N$  (3.11) to obtain

$$\liminf_{N\to\infty}rac{1}{N}\mathcal{J}^N(oldsymbol{c}^N)\geq \mathbb{J}(\mathbb{C}).$$

Hence, if each  $c^N$  is a curve of maximal slope w.r.t.  $\frac{1}{N}\mathcal{J}^N(c^N)$  then so is  $\mathbb{C}$  w.r.t.  $\mathbb{J}$ . The other statements also can be directly adapted from [40].

3.3. **Application.** To apply Theorem 3.6, we first have to show the convergence of the energy (3.12).

**Proposition 3.7** (lim inf inequality for the relative entropy). Let a sequence  $\mu^N \in \mathcal{P}(\mathcal{X}^N)$  be given such that  $\mu^N \stackrel{\tau}{\to} \mathbb{M}$  as  $N \to \infty$ , then

$$\liminf_{N\to\infty} \frac{1}{N} \mathcal{H}(\boldsymbol{\mu}^N \mid \boldsymbol{\pi}^N) \ge \int_{\mathcal{P}(\mathcal{X})} (\mathcal{F}(\nu) - \mathcal{F}_0) \, \, \mathbb{M}(d\nu) = \mathbb{F}(\mathbb{M}) - \mathcal{F}_0,$$

where

$$\mathcal{F}_0 := \inf_{\mu \in \mathcal{P}(\mathcal{X})} \mathcal{F}(\mu).$$

*Proof.* First we note that the relative entropy can be decomposed as follows:

$$\mathcal{H}(\boldsymbol{\mu}^{N} \mid \boldsymbol{\pi}^{N}) = \int \mathcal{H}(\boldsymbol{\mu}^{N}(\cdot \mid L^{N} = \nu) \mid \boldsymbol{\pi}^{N}(\cdot \mid L^{N} = \nu)) dL_{\#}^{N} \boldsymbol{\mu}^{N}(\nu)$$
$$+ \mathcal{H}(L_{\#}^{N} \boldsymbol{\mu}^{N} \mid L_{\#}^{N} \boldsymbol{\pi}^{N})$$

Let us denote by  $M^N$  the uniform probability measure on  $\mathcal{P}_N(\mathcal{X})$ . Using the fact that relative entropy w.r.t. a probability measure is non-negative we arrive a the estimate

$$\mathcal{H}(\boldsymbol{\mu}^{N} \mid \boldsymbol{\pi}^{N}) \geq \mathcal{H}(L_{\#}^{N} \boldsymbol{\mu}^{N} \mid L_{\#}^{N} \boldsymbol{\pi}^{N})$$

$$= \mathcal{H}(L_{\#}^{N} \boldsymbol{\mu}^{N} \mid \boldsymbol{M}^{N}) + \mathbf{E}_{L_{\#}^{N} \boldsymbol{\mu}^{N}} \left[ \log \frac{d\boldsymbol{M}^{N}}{dL_{\#}^{N} \boldsymbol{\pi}^{N}} \right]$$

$$\geq \mathbf{E}_{L_{\#}^{N} \boldsymbol{\mu}^{N}} \left[ \log \frac{d\boldsymbol{M}^{N}}{dL_{\#}^{N} \boldsymbol{\pi}^{N}} \right] . \tag{3.15}$$

For  $\nu \in \mathcal{P}_N(\mathcal{X})$  we set  $\mathcal{T}_N(\nu) = \{ \boldsymbol{x} \in \mathcal{X}^N : L^N(\boldsymbol{x}) = \nu \}$ . Then by the definition of  $\boldsymbol{\pi}^N$  (3.2) and  $U(\nu)$  (2.6) we have

$$L_{\#}^{N}\boldsymbol{\pi}^{N}(\nu) = \frac{|\mathcal{T}_{N}(\nu)|}{\mathbf{Z}^{N}} \exp\left(-NU(\nu)\right).$$

From (3.15) we thus conclude

$$\frac{1}{N}\mathcal{H}(\boldsymbol{\mu}^{N} \mid \boldsymbol{\pi}^{N}) \ge -\frac{1}{N}\log|\mathcal{P}_{N}(\mathcal{X})| + \frac{1}{N}\log\boldsymbol{Z}^{N} 
-\frac{1}{N}\operatorname{E}_{L_{\#}\boldsymbol{\mu}^{N}}\left[\log|\mathcal{T}_{N}|\right] + \operatorname{E}_{L_{\#}\boldsymbol{\mu}^{N}}[U].$$
(3.16)

The cardinality of  $\mathcal{P}_N(\mathcal{X})$  is given by  $\binom{N+d-1}{N} \leq N^{d-1}/d!$  and hence

$$\log |\mathcal{P}_N(\mathcal{X})| \le (d-1)\log N. \tag{3.17}$$

Moreover, by Stirling's formula (cf. Lemma A.1), it follows that for any  $\nu \in \mathcal{P}_N(\mathcal{X})$ 

$$-\frac{1}{N}\log|\mathcal{T}_N(\nu)| = -\frac{1}{N}\log\frac{N!}{\prod_{x\in\mathcal{X}}(N\nu(x))!}$$
$$= \sum_{x\in\mathcal{X}}\nu(x)\log\nu(x) + O\left(\frac{\log N}{N}\right). \tag{3.18}$$

Furthermore, we have

$$Z^N = \sum_{\nu \in \mathcal{P}_N(\mathcal{X})} e^{-NU(\nu)} |\mathcal{T}_N(\nu)|.$$

Hence, by using Sanov's and Varadhan's theorem [20] on the asymptotic evaluation of exponential integrals, it easily follows that

$$\lim_{N \to \infty} \frac{1}{N} \log \mathbf{Z}^N = -\inf_{\nu \in \mathcal{P}(\mathcal{X})} \mathcal{F}(\nu) =: -\mathcal{F}_0.$$
 (3.19)

Combing now (3.16) with (3.17)-(3.19) finishes the proof.

The other ingredient of the proof of Theorem 3.6 consists in proving the convergence of the metric derivatives (3.13) and slopes (3.14).

**Proposition 3.8** (Convergence of metric derivative and slopes). Let  $\mathbf{c}^N$  be an element of  $\mathrm{AC}([0,T],\mathcal{P}(\mathcal{X}^N))$ , and choose  $\boldsymbol{\psi}^N:[0,T]\to\mathcal{P}(\mathcal{X}^N)$  such that  $(\mathbf{c}^N,\boldsymbol{\psi}^N)$  solves the continuity equation. Furthermore, assume that

$$\boldsymbol{c}^{N} \stackrel{\tau}{\to} \mathbb{C},$$

with some measurable  $\mathbb{C}:[0,T]\to\mathcal{P}(\mathcal{P}(\mathcal{X}))$ , and that

$$\liminf_{N\to\infty}\int_0^T\frac{1}{N}\mathcal{A}^N(\boldsymbol{c}^N(t),\boldsymbol{\psi}^N(t))dt<\infty.$$

Then  $\mathbb{C} \in AC([0,T], \mathcal{P}(\mathcal{P}(\mathcal{X})))$ , and there exists  $\Psi : [0,T] \to \mathcal{P}(\mathcal{P}(\mathcal{X}))$ , for which  $(\mathbb{C}, \Psi)$  satisfy the continuity equation and for which we have

$$\liminf_{N\to\infty} \int_0^T \frac{1}{N} \mathcal{A}^N(\boldsymbol{c}^N(t), \boldsymbol{\psi}^N(t)) dt \ge \int_0^T \mathbb{A}(\mathbb{C}(t), \Psi(t)) \ dt$$

and

$$\liminf_{N\to\infty} \int_0^T \frac{1}{N} \mathcal{I}^N \left( \boldsymbol{c}^N(t) \right) dt \ge \int_0^T \mathbb{I} \left( \mathbb{C}(t) \right) dt.$$

*Proof.* Let us summarize consequences of the assumption  $\mathbf{c}^N \stackrel{\tau}{\to} \mathbb{C}$ . By Definition 3.5, this means  $L^N_\# \mathbf{c}^N(t) \rightharpoonup \mathbb{C}(t)$  for all  $t \in [0,T]$ . For  $x,y \in \mathcal{X}$  we define two auxiliary measures  $\mathbb{F}^{N;x,y;1}(t), \mathbb{F}^{N;x,y;2}(t) \in \mathcal{P}(\mathcal{P}(\mathcal{X}) \times \mathcal{P}(\mathcal{X}))$  by setting

$$\begin{split} \mathbb{F}^{N;x,y;1}(t,\nu,\mu) &:= \delta_{\nu^{N;x,y}}(\mu) \nu_x Q_{xy}^N(\nu) L_\#^N \pmb{c}^N(t,\nu) \\ \mathbb{F}^{N;x,y;2}(t,\nu,\mu) &:= \delta_{\mu^{N;y,x}}(\nu) \mu_y Q_{yx}^N(\mu) L_\#^N \pmb{c}^N(t,\mu), \end{split}$$

where  $\nu^{N;x,y} := \nu - \frac{\delta_x - \delta_y}{N}$ . Then, we have  $\mathbb{F}^{N;x,y;1}(t,\nu,\mu) = \mathbb{F}^{N;y,x;2}(t,\mu,\nu)$ . Due to (3.5) from Lemma 3.2 it holds

$$\mathbb{F}^{N;x,y;1}(t,\nu,\mu) \rightharpoonup \delta_{\nu}(\mu)\nu_x Q_{xy}(\nu)\mathbb{C}(t,\nu) \tag{3.20}$$

$$\mathbb{F}^{N;x,y;2}(t,\nu,\mu) \rightharpoonup \delta_{\mu}(\nu)\mu_x Q_{xy}(\mu)\mathbb{C}(t,\mu). \tag{3.21}$$

In the sequel, we will decompose the sum over all possible jumps of the particle system in different ways

$$\sum_{\boldsymbol{x},\boldsymbol{y}} f(\boldsymbol{x},\boldsymbol{y}) = \sum_{\substack{\nu,\mu \in \mathcal{P}_N(\mathcal{X})}} \sum_{\substack{\boldsymbol{x}:L^N\boldsymbol{x} = \nu\\\boldsymbol{y}:L^N\boldsymbol{y} = \mu}} f(\boldsymbol{x},\boldsymbol{y}) = \sum_{\substack{x,y \in \mathcal{X}}} \sum_{\substack{\nu \in \mathcal{P}_N(\mathcal{X})}} \sum_{i=1}^N \sum_{\substack{\boldsymbol{x}:L^N\boldsymbol{x} = \nu\\x_i = x}} f(\boldsymbol{x},\boldsymbol{x}^{i;y}).$$

where  $\mathbf{x}^{i;y} = \mathbf{x} - (x_i - y)\mathbf{e}^i$  and  $f: \mathcal{X}^N \times \mathcal{X}^N \to \mathbf{R}$  with  $f(\mathbf{x}, \mathbf{y}) = 0$  unless  $\mathbf{y} = \mathbf{x}^{i;y}$  for some  $i \in \{1, ..., N\}$  and  $y \in \mathcal{X}$ . We define the following vector field on  $\mathcal{P}(\mathcal{X}) \times \mathcal{P}(\mathcal{X})$ 

$$\begin{split} \mathbf{v}^{N;x,y}(t,\nu,\mu) &:= \frac{1}{2N} \delta_{\nu^{N;x,y}}(\mu) \sum_{\substack{\boldsymbol{x}:L^N \boldsymbol{x} = \nu \\ \boldsymbol{y}:L^N \boldsymbol{y} = \mu}} \left( \boldsymbol{\psi}_{\boldsymbol{y}}^N(t) - \boldsymbol{\psi}_{\boldsymbol{x}}^N(t) \right) \boldsymbol{w}_{\boldsymbol{x}\boldsymbol{y}}^N \left( \boldsymbol{c}^N(t) \right) \\ &= \frac{1}{2N} \delta_{\nu^{N;x,y}}(\mu) \sum_{i=1}^N \sum_{\substack{\boldsymbol{x}:L^N \boldsymbol{x} = \nu \\ \boldsymbol{x}_i = x}} \left( \boldsymbol{\psi}_{\boldsymbol{x}^{i;y}}^N(t) - \boldsymbol{\psi}_{\boldsymbol{x}}^N(t) \right) \boldsymbol{w}_{\boldsymbol{x}\boldsymbol{x}^{i;y}}^N \left( \boldsymbol{c}^N(t) \right), \end{split}$$

where  $\boldsymbol{w}_{\boldsymbol{x}\boldsymbol{y}}^{N}\left(\boldsymbol{c}^{N}(t)\right)$  is defined in (3.8). From the definition of  $\boldsymbol{v}^{N;x,y}(t)$  and the Cauchy-Schwarz inequality, it follows that for  $\nu, \mu \in \mathcal{P}(\mathcal{X})$  with  $\mu = \nu^{N;x,y}$  for some  $x, y \in \mathcal{X}$ 

$$egin{aligned} \left|\mathbf{v}^{N;x,y}(t,
u,\mu)
ight| &\leq \left(rac{1}{2N}\sum_{\substack{oldsymbol{x}:L^Noldsymbol{x}=
u\ oldsymbol{y}:L^Noldsymbol{y}=\mu}} \left(oldsymbol{\psi}_{oldsymbol{y}}^N(t) - oldsymbol{\psi}_{oldsymbol{x}}^N(t)
ight)^2oldsymbol{w}_{oldsymbol{x}oldsymbol{y}}^N\left(oldsymbol{c}^N(t)
ight)
ight)^{rac{1}{2}} imes \ \left(rac{1}{2N}\sum_{i=1}^N\sum_{oldsymbol{x}:L^Noldsymbol{x}=
u}oldsymbol{w}_{oldsymbol{x}oldsymbol{x}^N;y}\left(oldsymbol{c}^N(t)
ight)
ight)^{rac{1}{2}}. \end{aligned}$$

By using the identity

$$\frac{1}{N} \sum_{i=1}^{N} \sum_{\substack{\boldsymbol{x}: L^{N}(\boldsymbol{x}) = \nu \\ x_{i} = x}} \boldsymbol{c}_{\boldsymbol{x}}^{N}(t) = \sum_{\boldsymbol{x}} \delta_{\nu}(L^{N}(\boldsymbol{x})) \ L_{x}^{N}(\boldsymbol{x}) \ \boldsymbol{c}_{\boldsymbol{x}}^{N}(t) = L_{\#}^{N} \boldsymbol{c}^{N}(t, \nu) \ \nu_{x},$$

and the fact that the logarithmic mean is jointly concave and 1-homogeneous, we can conclude

$$\frac{1}{N} \sum_{i=1}^{N} \sum_{\substack{\boldsymbol{x}: L^{N} \boldsymbol{x} = \nu \\ x_{i} = x}} \boldsymbol{w}_{\boldsymbol{x} \boldsymbol{x}^{i;y}}^{N} \left( \boldsymbol{c}^{N}(t) \right) \leq \Lambda \left( \mathbb{F}^{N;x,y;1}(t,\nu,\nu^{N;x,y}), \mathbb{F}^{N;x,y;2}(t,\nu,\nu^{N;x,y}) \right),$$

which first shows that  $v^{N;x,y}(t) \ll \Lambda(\mathbb{F}^{N;x,y;1}(t),\mathbb{F}^{N;x,y;2}(t))$  as product measure on  $\mathcal{P}(\mathcal{X}) \times \mathcal{P}(\mathcal{X})$ . Moreover, by summation and integration over any Borel  $I \subset [0,T]$  we get

$$\int_{I} \sum_{\nu,\mu \in \mathcal{P}_{N}(\mathcal{X})} \left| \mathbf{v}^{N;x,y}(t,\nu,\mu) \right| dt \leq \left( \sqrt{T} \int_{0}^{T} \frac{1}{N} \mathcal{A}(\boldsymbol{c}^{N}(t), \boldsymbol{\psi}(t)) dt \right)^{\frac{1}{2}} \times \left( \frac{|I|}{2} \sum_{\nu,\mu \in \mathcal{P}_{N}(\mathcal{X})} \sup_{t \in I} \Lambda \left( \mathbb{\Gamma}^{N;x,y;1}(t,\nu,\mu), \mathbb{\Gamma}^{N;x,y;2}(t,\nu,\mu) \right) \right)^{\frac{1}{2}}.$$
(3.22)

The second sum is uniformly bounded in N, since  $\mathcal{X}$  is finite and by Lemma 3.2  $Q^N \to Q$  uniformly with Q continuous in the first argument on the compact space  $\mathcal{P}(\mathcal{X})$ . Now, from (3.22), we conclude that for some subsequence and all  $x, y \in \mathcal{X}$  we have  $\mathbf{v}^{N;x,y} \rightharpoonup \mathbf{v}^{x,y}$  with  $\mathbf{v}^{x,y}$  a Borel measure on  $[0,T] \times \mathcal{P}(\mathcal{X}) \times \mathcal{P}(\mathcal{X})$ . Using Jensen's inequality applied

to the 1-homogeneous jointly convex function  $\mathbf{R} \times \mathbf{R}^2_+ \ni (v, a, b) \mapsto \frac{v^2}{\Lambda(a, b)}$ , we get

$$\begin{split} & \int_0^T \frac{1}{N} \mathcal{A} \left( \boldsymbol{c}^N(t), \boldsymbol{\psi}^N(t) \right) \ dt \\ & = \int_0^T \frac{1}{2} \sum_{x,y} \sum_{\nu \in \mathcal{P}_N(\mathcal{X})} \frac{1}{N} \sum_{i=1}^N \sum_{\boldsymbol{x}: L^N \boldsymbol{x} = \nu} \frac{\left( \left( \boldsymbol{\psi}^N_{\boldsymbol{x}^{i;y}}(t) - \boldsymbol{\psi}^N_{\boldsymbol{x}}(t) \right) \boldsymbol{w}^N_{\boldsymbol{x} \boldsymbol{x}^{i;y}}(\boldsymbol{c}^N(t)) \right)^2}{\boldsymbol{w}^N_{\boldsymbol{x} \boldsymbol{x}^{i;y}}(\boldsymbol{c}^N(t))} \ dt \\ & \geq \int_0^T \frac{1}{2} \sum_{x,y} \sum_{\nu \in \mathcal{P}_N(\mathcal{X})} \frac{\left( \boldsymbol{v}^N(t, \nu, \boldsymbol{\nu}^{N;x,y}) \right)^2}{\Lambda(\mathbb{F}^{N;x,y;1}(t, \nu, \boldsymbol{\nu}^{N;x,y}), \mathbb{F}^{N;x,y;2}(t, \nu, \boldsymbol{\nu}^{N;x,y}))} \ dt \ . \end{split}$$

The last term can be written as

$$\frac{1}{2} \sum_{x,y} F(\mathbf{v}^{N;x,y}, \mathbb{F}^{N;x,y;1}, \mathbb{F}^{N;x,y;2}) \ ,$$

where the functional F on triples of measure on  $[0,T] \times \mathcal{P}(\mathcal{X})^2$  is defined via

$$F(\mathbf{v}, \mathbb{F}^1, \mathbb{F}^2) := \int_0^T \iint_{\mathcal{P}(\mathcal{X})^2} \alpha \bigg( \frac{d\mathbf{v}}{d\sigma}, \Lambda \bigg( \frac{d\mathbb{F}^1}{d\sigma}, \frac{d\mathbb{F}^2}{d\sigma} \bigg) \bigg) \ d\sigma \ dt \ ,$$

with  $\alpha$  the function defined in (2.13) and  $\sigma$  is any measure on  $[0,T] \times \mathcal{P}(\mathcal{X})^2$  such that  $\mathbb{V}, \mathbb{T}^1, \mathbb{T}^2 \ll \sigma$ . The definition does not depend on the choice of  $\sigma$  by the 1-homogeneity of  $\alpha$  and  $\Lambda$ . Then, by a general result on lower semicontinuity of integral functionals [8, Thm. 3.4.3] we can conclude, that

$$\liminf_{N\to\infty} F(\mathbf{v}^{N;x,y}, \mathbb{\Gamma}^{N;x,y;1}, \mathbb{\Gamma}^{N;x,y;2}) \ge F(\mathbf{v}^{x,y}, \mathbb{\Gamma}^{x,y;1}, \mathbb{\Gamma}^{x,y;2}) \; .$$

In particular, this implies

$$dv^{x,y} \ll \Lambda\left(\frac{d\mathbb{F}^{x,y;1}}{d\sigma}, \frac{d\mathbb{F}^{x,y;2}}{d\sigma}\right) d\sigma,$$

which by (3.20) and (3.21) is given in terms of

$$\Lambda\left(\frac{d\mathbb{F}^{x,y;1}}{d\sigma}(t,\nu,\mu),\frac{d\mathbb{F}^{x,y;2}}{d\sigma}(t,\nu,\mu)\right)d\sigma = \delta_{\nu}(\mu)\Lambda(\nu_{x}Q_{xy}(\nu),\nu_{y}Q_{yx}(\nu))\mathbb{C}(t,d\nu)dt.$$

Therefore, with the notation of Proposition 2.17, we obtain the statement

$$\lim_{N \to \infty} \inf \int_0^T \frac{1}{N} \mathcal{A} \left( \boldsymbol{c}^N(t), \boldsymbol{\psi}(t)^N \right) dt$$

$$\geq \frac{1}{2} \sum_{x,y} \int_0^T \int_{\mathcal{P}(\mathcal{X})} \frac{(\mathbb{V}_{xy}(t,\nu))^2}{\Lambda(\nu_x Q_{xy}(\nu), \nu_y Q_{yx}(\nu))} \mathbb{C}(t, d\nu) dt$$

$$= \int_0^T \vec{\mathbb{A}}(\mathbb{C}(t), \mathbb{V}(t)) dt \quad \text{with} \quad \mathbb{V}_{xy}(t, \nu) := \frac{d \mathbf{v}^{x,y}}{d \mathbb{C}(t) dt}.$$

From the convergence of the vector field  $\mathbf{v}^{N;x,y} \rightharpoonup \mathbf{v}^{x,y}$  it is straightforward to check that  $(\mathbb{C}, \mathbb{V}) \in \vec{\mathbb{CE}}_T(\mathbb{C}_0, \mathbb{C}_T)$  and hence by the conclusion of Proposition 2.17, there exists

 $\Psi: [0,T] \times \mathcal{P}(\mathcal{X}) \to \mathbf{R}^{\mathcal{X}}$  such that

$$\liminf_{N\to\infty}\frac{1}{N}\int_0^T \mathcal{A}\left(\boldsymbol{c}^N(t),\boldsymbol{\psi}^N(t)\right)dt \geq \int_0^T \vec{\mathbb{A}}(\mathbb{C}(t),\mathbb{V}(t))\ dt \geq \int_0^T \mathbb{A}(\mathbb{C}(t),\mathbb{\Psi}(t))\ dt,$$

which concludes the first part.

The liminf estimate of the Fisher information follows by a similar but simpler argument. The convex 1-homogeneous function  $\lambda(a,b)=(a-b)(\log a-\log b)$  allows to rewrite

$$\begin{split} &\frac{1}{N} \mathcal{I}^{N}\left(\boldsymbol{c}^{N}(t)\right) dt \\ &= \frac{1}{2} \sum_{x,y} \sum_{\nu \in \mathcal{P}_{N}(\mathcal{X})} \frac{1}{N} \sum_{i=1}^{N} \sum_{\substack{\boldsymbol{x}: L^{N} \boldsymbol{x} = \nu \\ x_{i} = x}} \lambda \left(\boldsymbol{c}_{\boldsymbol{x}}^{N}(t) \boldsymbol{Q}_{\boldsymbol{x} \boldsymbol{x}^{i;y}}^{N}(\boldsymbol{c}^{N}(t)), \boldsymbol{c}_{\boldsymbol{x}^{i;y}}^{N}(t) \boldsymbol{Q}_{\boldsymbol{x}^{i;y} \boldsymbol{x}}^{N}(\boldsymbol{c}^{N}(t))\right) \\ &\geq \frac{1}{2} \sum_{x,y} \iint_{\mathcal{P}(\mathcal{X})^{2}} \lambda \left(\frac{d\mathbb{F}^{N;x,y;1}(t)}{d\sigma}, \frac{d\mathbb{F}^{N;x,y;2}(t)}{d\sigma}\right) d\sigma. \end{split}$$

Then, the result follows by an application of [8, Thm. 3.4.3].

In order to apply Theorem 3.6, we still need to prove that a sequence of N-particle dynamics starting from nice initial conditions is tight.

**Lemma 3.9.** Let  $\mathbf{X}^N$  be the continuous Markov jump process with generator (3.6), then the sequence of laws of empirical measures is tight in the Skorokhod topology, i.e. it holds for any  $x \in \mathcal{X}$  and  $\varepsilon > 0$ 

$$\lim_{\delta \to 0} \limsup_{N \to \infty} \mathbb{P} \left[ \sup_{|t-s| \le \delta} \left| L_x^N(\mathbf{X}^N(t)) - L_x^N(\mathbf{X}^N(s)) \right| > \epsilon \right] = 0.$$
 (3.23)

The proof follows standard arguments for tightness of empirical measures of sequences of interacting particle systems. Since there is no original argument here, the exposition shall be brief, and we refer to [29] for more details, in a more general context of interacting particle systems. For example, see the first step in the proof of Theorem 2.1 in [29] for those arguments in the context of the simple exclusion process on a discrete torus.

*Proof.* The process

$$\begin{split} M_x^N(t) &= L_x^N(\mathbf{X}^N(t)) - L_x^N(\mathbf{X}^N(0)) \\ &- \int_0^t \sum_{y \neq x} L_x^N(\mathbf{X}^N(s)) \, Q_{xy}^N(L^N(\mathbf{X}^N(s))) - L_y^N(\mathbf{X}^N(s)) \, Q_{yx}^N(L^N(\mathbf{X}^N(s))) \, ds \end{split}$$

is a martingale. Since the rates are bounded,

$$\left| \int_{s}^{t} \sum_{y \neq x} L_{x}^{N}(\mathbf{X}^{N}(r)) Q_{xy}^{N}(L^{N}(\mathbf{X}^{N}(r))) - L_{y}^{N}(\mathbf{X}^{N}(r)) Q_{yx}^{N}(L^{N}(\mathbf{X}^{N}(r))) dr \right| \leq C|t-s|$$

and therefore, to prove (3.23), it is enough to show that

$$\mathbb{P}\left[\sup_{|t-s| \le \delta} \left| M_x^N(t) - M_x^N(s) \right| > \epsilon \right]$$

is small. To do so, we shall estimate the quadratic variation of the martingale. It is given by

$$\begin{split} \langle \boldsymbol{M}^{N} \rangle(t) &= \frac{1}{N^{2}} \sum_{y \neq x} \left( (NL_{x}^{N}(\mathbf{X}^{N}(t)) - 1)^{2} - (NL_{x}^{N}(\mathbf{X}^{N}(t)))^{2} \right) Q_{xy}^{N}(L^{N}(\mathbf{X}^{N}(s))) + \\ & \left( (NL_{x}^{N}(\mathbf{X}^{N}(t)) + 1\right)^{2} - \left( NL_{x}^{N}(\mathbf{X}^{N}(t))^{2} \right) Q_{yx}^{N}(L^{N}(\mathbf{X}^{N}(s))) + \\ & \frac{2}{N^{2}} \sum_{y \neq x} NL_{x}^{N}(\mathbf{X}^{N}(t)) \left( L_{x}^{N}(\mathbf{X}^{N}(s)) Q_{xy}^{N} \left( L^{N}(\mathbf{X}^{N}(s)) \right) - L_{y}^{N}(\mathbf{X}^{N}(s)) Q_{yx}^{N} \left( L^{N}(\mathbf{X}^{N}(s)) \right) \right). \end{split}$$

Using the boundedness of the rates, it is straightforward to see that

$$|\langle M^N \rangle(t)| \le CN^{-1}$$
.

The quadratic variation of the martingale vanishes. From Doob's martingale inequality, we deduce that for any  $\epsilon>0$ 

$$\mathbb{P}\left[\sup_{0\leq s\leq t\leq T}|M_x^N(t)-M_x^N(s)|>\epsilon\right]\longrightarrow 0.$$

Tightness of the sequence of processes follows.

We can now combine the work done so far to obtain our main result.

**Theorem 3.10** (Convergence of the particle system to the mean field equation). Let  $\mathbf{c}^N$  be a solution to (3.7). Moreover assume its initial distribution to be well-prepared

$$\frac{1}{N}\mathcal{F}^N(\mathbf{c}^N(0)) \to \mathbb{F}(\mathbb{C}(0)) - \mathcal{F}_0 \qquad with \qquad L_\#^N \mathbf{c}^N(0) \rightharpoonup \mathbb{C}(0) \qquad as \ N \to \infty.$$

Then it holds

$$L^N_{\#} \mathbf{c}^N(t) \rightharpoonup \mathbb{C}(t)$$
 for all  $t \in (0, \infty)$ ,

with  $\mathbb{C}$  a weak solution to (2.20) and moreover

$$\frac{1}{N} \mathcal{F}^{N}(\mathbf{c}^{N}(t)) \to \mathbb{F}(\mathbb{C}(t)) - \mathcal{F}_{0} \quad \text{for all } t \in (0, \infty).$$
 (3.24)

Proof. Fix T > 0. By the tightness Lemma 3.9, we have that the sequence of empirical measures  $L_{\#}^N \mathbf{c}^N : [0,T] \to \mathcal{P}(\mathcal{P}(\mathcal{X}))$  is tight w.r.t. the Skorokhod topology [4, Theorem 13.2]. Hence, there exist a measurable curve  $\mathbb{C} : [0,T] \to \mathcal{P}(\mathcal{P}(\mathcal{X}))$  such that up to a subsequence  $L_{\#}^N \mathbf{c}^N(t)$  weakly converges to  $\mathbb{C}(t)$  for all  $t \geq 0$ . By the Propositions 3.7 and 3.8, we get from Theorem 3.6 that (3.24) holds and  $\mathbb{C}$  is curve of maximal slope for the functional  $\mathbb{J}$ . By Proposition 2.21, it is characterized as weak solution to (2.20). By Lemma 3.2 the limiting rate matrix Q is Lipschitz on  $\mathcal{P}(\mathcal{X})$  providing uniqueness of the Liouville equation (2.20). Hence, the convergence actually holds for the full sequence.  $\square$ 

Corollary 3.11. In the setting of Theorem 3.10 assume in addition that

$$L_{\#}^{N} \mathbf{c}^{N}(0) \rightharpoonup \delta_{c(0)}$$
 for some  $c(0) \in \mathcal{P}(\mathcal{X})$ .

Then it holds

$$L^N_{\#} \mathbf{c}^N(t) \rightharpoonup \delta_{c(t)}$$
 for all  $t \in (0, \infty)$ ,

with c a solution to (1.1) and moreover

$$\frac{1}{N} \mathcal{F}^N(\mathbf{c}^N(t)) \to \mathcal{F}(c(t)) - \mathcal{F}_0 \qquad \text{for all } t \in (0, \infty) \ .$$

*Proof.* The proof is a direct application of Theorem 3.10 and a variance estimate for the particle system (Lemma B.1).  $\Box$ 

### 4. Properties of the metric W

In this section, we give the proof of Propostion 2.10 stating that W defines a distance on  $\mathcal{P}(\mathcal{X})$  and that the resulting metric space is seperable, complete and geodesic. The proof will be accomplished by a sequence of lemmas giving various estimates on and properties of W. Some work is needed in particular to show finiteness of W.

**Lemma 4.1.** For  $\mu, \nu$ , and T > 0 we have

$$\mathcal{W}(\mu,\nu) = \inf \left\{ \int_0^T \sqrt{\mathcal{A}(c(t),\psi(t))} \ dt : (c,\psi) \in CE_T(\mu,\nu) \right\}.$$

*Proof.* This follows from a standard reparametrization argument. See for instance [18, Thm. 5.4] for details in a similar situation.  $\Box$ 

From the previous lemma we easily deduce the triangular inequality

$$W(\mu, \eta) \le W(\mu, \nu) + W(\nu, \eta) \quad \forall \mu, \nu, \eta \in \mathcal{P}(\mathcal{X}) , \tag{4.1}$$

by concatenating two curves  $(c, \psi) \in \mathrm{CE}_T(\mu, \nu)$  and  $(c', \psi') \in \mathrm{CE}_T(\nu, \eta)$  to form a curve in  $\mathrm{CE}_{2T}(\mu, \eta)$ .

The next sequence of lemmas puts W in relation with the total variation distance on  $\mathcal{P}(\mathcal{X})$ . To proceed, we define similarly to [31], for every  $\mu \in \mathcal{P}(\mathcal{X})$ , the matrix

$$B_{xy}(\mu) := \begin{cases} \sum_{z \neq x} w_{xz}(\mu), & x = y, \\ -w_{xy}(\mu), & x \neq y \end{cases}$$

Now (2.16) can be rewritten as

$$\mathcal{W}^{2}(\mu,\nu) = \inf \left\{ \int_{0}^{1} \langle B(c(t))\psi(t), \psi(t) \rangle \ dt : (c,\psi) \in \mathrm{CE}_{1}(\mu,\nu) \right\},\,$$

where  $\langle \psi, \phi \rangle = \sum_{x \in \mathcal{X}} \psi_x \phi_x$  is the usual inner product on  $\mathbf{R}^{\mathcal{X}}$ .

**Lemma 4.2.** For any  $\mu, \nu \in \mathcal{P}(\mathcal{X})$  holds

$$\mathcal{W}(\mu,\nu) \ge \frac{1}{\sqrt{2}} \|\mu - \nu\|.$$

Moreover, for every a > 0, there exists a constant  $C_a$ , such that for all  $\mu, \nu \in \mathcal{P}^a(\mathcal{X})$ 

$$\mathcal{W}(\mu,\nu) \le C_a \|\mu - \nu\| .$$

*Proof.* The proof of the lower bound on W can be obtained very similar to [22, Proposition 2.9].

Let us show the upper bound. Following Lemma A.1. in [31], we notice that for  $\mu \in \mathcal{P}^a(\mathcal{X})$ , the map  $\psi \mapsto B(\mu)\psi$  has an image of dimension d-1. In addition, the dimension of the space  $\{a_x : \sum_{x \in \mathcal{X}} a_x = 0\}$  is d-1, therefore the map is surjective. From the above

we get that the matrix  $B(\mu)$  restricts to an isomorphism  $\tilde{B}(\mu)$ , on the d-1 dimensional space  $\{a_x : \sum a_x = 0\}$ . Now, since the mapping  $\mathcal{P}^a(\mathcal{X}) \ni \mu \to ||\tilde{B}^{-1}(\mu)||$ , is continuous with respect to the euclidean metric, we have an upper bound  $\frac{1}{c}$  by compactness. Also  $\mathcal{P}^a(\mathcal{X}) \ni \mu \to ||\tilde{B}(\mu)||$  has an upper bound C as a result of all entries in  $B(\mu)$  being uniformly bounded. From that we get

$$c\|\psi\| \le \|B(\mu)\psi\| \le C\|\psi\|, \forall \mu \in \mathcal{P}^a(\mathcal{X})$$

for some suitable positive constants.

Similarly to the proof of Lemma 3.19 in [31], for  $t \in [0,1]$ , we set  $c(t) = (1-t)\mu + t\nu$  and note that c(t) lies in  $\mathcal{P}^a(\mathcal{X})$ , since it is a convex set. Since  $\dot{c}(t) = \nu - \mu \in \operatorname{Ran} B(c(t))$ , there exists a unique element  $\psi(t)$  for which we have  $\dot{c}(t) = B(c(t))\psi(t)$ , and  $\|\psi(t)\| \leq \frac{1}{c}\|\mu - \nu\|$ .

From that we get

$$W^{2}(\mu, \nu) \leq \int_{0}^{1} \langle B(c(t))\psi(t), \psi(t) \rangle dt \leq \frac{1}{c^{2}} C \|\mu - \nu\|^{2}.$$

**Lemma 4.3.** For every  $\mu \in \mathcal{P}(\mathcal{X})$ ,  $\epsilon > 0$  there exists  $\delta > 0$  such that  $\mathcal{W}(\mu, \nu) < \epsilon$ , for every  $\nu \in \mathcal{P}(\mathcal{X})$ , with  $\|\mu - \nu\| \leq \delta$ 

The proof of this lemma is similar to the proof [31, Theorem 3.12] and uses comparison of W to corresponding quantity on the two point space  $\mathcal{X} = \{a, b\}$ . However, significantly more care is needed in the present setting to implement this argument. The reason being that the set of pairs of points x, y with  $Q_{x,y}(\mu) > 0$  now depends on  $\mu$ .

*Proof.* Let  $\epsilon > 0$  and  $\mu \in \mathcal{P}(\mathcal{X})$  be fixed. Since  $\mathcal{X}$  is finite, it holds with  $E_{\mu}$  defined in (2.11)

$$\inf \{Q_{xy}(\mu) : (x,y) \in E_{\mu}\} = a > 0.$$

Let  $B_r(\mu) = \{ \nu \in \mathcal{P}(\mathcal{X}) : \|\nu - \mu\| < r \}$  denote a r-neighborhood around  $\mu$ . Since,  $Q(\mu)$  is continuous in  $\mu$ , there exists for  $\eta > 0$  a  $\delta_1 > 0$  s.t.

$$\forall \nu \in B_{\delta_1}(\mu) \text{ holds } |Q(\nu) - Q(\mu)|_{L^{\infty}(\mathcal{X} \times \mathcal{X})} \leq \eta.$$

Especially, it holds by choosing  $\eta \leq a/2$  that  $E_{\mu} \subseteq E_{\nu}$  and in addition

$$\inf \{Q_{xy}(\nu) : (x,y) \in E_{\mu}, \nu \in B_{\delta_1}(\mu)\} \ge a/2.$$

For the next argument, observe that by the concavity of the logarithmic mean and a first order Taylor expansion holds for  $a, b, s, t, \eta > 0$ 

$$\Lambda((s+\eta)a, (t+\eta)b) \leq \Lambda(sa, tb) + \eta \left(\partial_s \Lambda(sa, tb) + \partial_t \Lambda(sa, tb)\right)$$
$$= \Lambda(sa, tb) + \eta \left(a\Lambda_1(sa, tb) + b\Lambda_2(sa, tb)\right),$$

where  $\Lambda_i$  is the *i*-th partial derivative of  $\Lambda$ . Therefore we can estimate for  $\nu \in B_{\delta}(\mu)$ .

$$\begin{split} &\Lambda(Q_{xy}(\nu)\nu(x),Q_{yx}(\nu)\nu(y)) - \Lambda(Q_{xy}(\mu)\nu(x),Q_{yx}(\mu)\nu(y)) \\ &\leq \Lambda((Q_{xy}(\mu)+\eta)\nu(x),(Q_{yx}(\mu)+\eta)\nu(y)) - \Lambda(Q_{xy}(\mu)\nu(x),Q_{yx}(\mu)\nu(y)) \\ &\leq \eta \Big(\nu(x)\Lambda_1(Q_{xy}(\mu)\nu(x),Q_{yx}(\mu)\nu(y)) + \nu(y)\Lambda_2(Q_{xy}(\mu)\nu(x),Q_{yx}(\mu)\nu(y))\Big) \\ &\leq \frac{2\eta}{a} \Big(Q_{xy}(\mu)\nu(x)\Lambda_1\big(Q_{xy}(\mu)\nu(x),Q_{yx}(\mu)\nu(y)\big) \\ &\qquad \qquad + Q_{yx}(\mu)\nu(y)\Lambda_2\big(Q_{xy}(\mu)\nu(x),Q_{yx}(\mu)\nu(y)\big)\Big) \\ &= \frac{2\eta}{a}\Lambda\big(Q_{xy}(\mu)\nu(x),Q_{yx}(\mu)\nu(y)\big) \leq \Lambda\big(Q_{xy}(\mu)\nu(x),Q_{yx}(\mu)\nu(y)\big), \end{split}$$

Moreover, the last identity follows directly from the one-homogeneity of the logarithmic mean. Furthermore, we used  $\eta \leq \frac{a}{2}$  to obtain the last estimate. Repeating the argument for the other direction we get

$$\frac{1}{2}\Lambda(Q_{xy}(\mu)\nu(x), Q_{yx}(\mu)\nu(y)) \leq \Lambda(Q_{xy}(\nu)\nu(x), Q_{yx}(\nu)\nu(y)) 
\leq 2\Lambda(Q_{xy}(\mu)\nu(x), Q_{yx}(\mu)\nu(y))$$
(4.2)

Now, let c be an absolutely continuous curve with respect to  $W_{Q(\mu)}$ , where  $W_{Q(\mu)}$  is the distance that corresponds to the linear Markov process with fixed rates  $Q(\mu)$ , and lives inside the ball  $B_{\delta_1}(\mu)$ , then it is also absolutely continuous w.r.t. W, and if  $\psi$  solves the continuity equation for c, with respect to the rates  $Q(\mu)$ , then there exists a  $\tilde{\psi}$ , that solves the continuity equation with respect to the variable rates Q(c(t)) and

$$\int_0^1 \mathcal{A}(c(t), \tilde{\psi}(t))dt \le 2 \int_0^1 \mathcal{A}_{Q(\mu)}(c(t), \psi(t))dt, \tag{4.3}$$

where  $\mathcal{A}_{Q(\mu)}$  is the action with fixed rate kernel  $Q(\mu)$ .

Indeed let  $\psi$  be a solution for the continuity equation for c with respect to the fixed rates  $Q(\mu)$ , i.e.

$$\dot{c}_x(t) = \sum_y (\psi_y(t) - \psi_x(t)) \Lambda(c_x(t)Q_{xy}(\mu), c_y(t)Q_{yx}(\mu)).$$

For  $(x,y) \in E_{\mu}$ , and  $t \in [0,1]$  we define

$$\tilde{v}_{xy}(t) := (\psi_y(t) - \psi_x(t))\Lambda(c_x(t)Q_{xy}(\mu), c_y(t)Q_{yx}(\mu)).$$

Then, it is easy to verify that  $(c, \tilde{v}) \in \vec{\mathrm{CE}}(c(0), c(1))$  (cf. Definition 2.5) and we can estimate

$$\int_{0}^{1} \vec{\mathcal{A}}(c(t), \tilde{v}(t)) dt = \int_{0}^{1} \frac{1}{2} \sum_{x,y} \alpha(\tilde{v}_{xy}(t), \Lambda(c_{x}(t)Q_{xy}(c(t)), c_{y}(t)Q_{yx}(c(t)))) dt 
= \int_{0}^{1} \frac{1}{2} \sum_{x,y} (\psi_{y}(t) - \psi_{x}(t))^{2} \frac{\Lambda(c_{x}(t)Q_{xy}(\mu), c_{y}(t)Q_{yx}(\mu))}{\Lambda(c_{x}(t)Q_{xy}(c(t)), c_{y}(t)Q_{yx}(c(t)))} 
\times \Lambda(c_{x}(t)Q_{xy}(\mu), c_{y}(t)Q_{yx}(\mu)) dt 
\stackrel{(4.2)}{\leq} \int_{0}^{1} \frac{1}{2} \sum_{x,y} (\psi_{y}(t) - \psi_{x}(t))^{2} 2\Lambda(c_{x}(t)Q_{xy}(\mu), c_{y}(t)Q_{yx}(\mu)) dt.$$

Now, the existence of  $\tilde{\psi}$  is a straightforward application of Lemma 2.7.

Having established (4.3), the final result will follow by a comparison with the two-point space for the Wasserstein distance with fixed rate kernel  $Q(\mu)$ .

For  $\nu \in B_{\delta}(\mu)$ , we can find a sequence of at most (d-1) measures  $\mu^i \in B_{\delta}(\mu)$ , such that  $\mu^0 = \mu$  and  $\mu^K = \nu$  and

supp 
$$(\mu^i - \mu^{i-1}) = \{x_i, y_i\} \in E_{\mu}$$
 for  $i = 1, ..., K$ .

Indeed we can use the following matching procedure: Find a pair (i,j) with  $\mu_i \neq \nu_i$  and  $\mu_j \neq \nu_j$ . Set  $h = \min\{|\mu_i - \nu_i|, |\mu_j - \nu_j|\}$ . Then define  $\mu_i^1 := \mu_i \pm h$  and  $\mu_j^1 := \mu_j \mp h$  with signs chosen as the sign of  $\nu_i - \mu_i$ . After this step at least (d-1)-coordinates of  $\mu^1$  and  $\nu$  agree. This procedure finishes after at most d-1 steps, because the defect mass of the last pair will match. Therewith, we can compare with the two-point space [31, Lemma 3.14]

$$\mathcal{W}_{Q(\mu)}(\mu^{i-1}, \mu^{i}) \leq \frac{1}{\sqrt{2p_{x_{i}y_{i}}}} \left| \int_{1-2\mu_{x_{i}}^{i}}^{1-2\mu_{x_{i}}^{i-1}} \sqrt{\frac{\operatorname{arctanh} r}{r}} dr \right| \leq \frac{\delta_{1}}{2},$$

with  $p_{x_iy_i} = Q_{x_iy_i}(\mu)\pi_{x_i}(\mu)$ . The last estimate follows from the fact, that the function  $\sqrt{\frac{\operatorname{arctanh} r}{r}}dr$ , is integrable in [-1,1]. Therefore, we can find a  $\delta \leq \delta_1$  such that for any a,b with  $|a-b| \leq \delta$ , we have  $\int_{1-2a}^{1-2b} \sqrt{\frac{\operatorname{arctanh} r}{r}}dr \leq \frac{\delta_1}{2}\min\{1,\sqrt{2p_{x_iy_i}}\}$ . Finally, by Lemma 4.2, we can infer that any curve has Euclidean length smaller than its action value. We can conclude that the  $\mathcal{A}_{Q(\mu)}$ -minimizing curve between any  $\mu^{i-1},\mu^i$ , stays inside the ball  $B_{\delta_1}(\mu)$ , from which we can further conclude that

$$\mathcal{W}(\mu^{i-1}, \mu^i) \le 2\mathcal{W}_{Q(\mu)}(\mu^{i-1}, \mu^i) \le 2\frac{\delta_1}{2} = \delta_1$$

By an application of the triangular inequality (4.1), we get  $W(\mu, \nu) \leq (d-1)\delta_1$ , and the proof concludes if we pick  $\delta$  such that  $(d-1)\delta_1 \leq \epsilon$ .

**Lemma 4.4.** For  $\mu_k, \mu \in \mathcal{P}(\mathcal{X})$ , we have

$$\mathcal{W}(\mu_k, \mu) \to 0$$
 iff  $\|\mu_k - \mu\| \to 0$ .

Moreover, the space  $\mathcal{P}(\mathcal{X})$ , along with the metric  $\mathcal{W}$ , is a complete space.

*Proof.* The proof is a direct application of Lemmas 4.2 and 4.3.

**Theorem 4.5** (Compactness of curves of finite action). Let  $\{(c^k, v^k)\}_k$ , with

$$(c^k, v^k) \in \overrightarrow{CE}_T(c^k(0), c^k(T)),$$

be a sequence of weak solutions to the continuity equation with uniformly bounded action

$$\sup_{k \in \mathbb{N}} \left\{ \int_0^T \vec{\mathcal{A}}(c^k(t), v^k(t)) dt \right\} \le C < \infty. \tag{4.4}$$

Then, there exists a subsequence and a limit (c, v), such that  $c^k$  converges uniformly to c in [0, T],  $(c, v) \in \vec{\operatorname{CE}}_T(c(0), c(T))$  and for the action we have

$$\liminf_{k \to \infty} \int_0^T \vec{\mathcal{A}}(c^k(t), v^k(t)) dt \ge \int_0^T \vec{\mathcal{A}}(c(t), v(t)) dt. \tag{4.5}$$

*Proof.* Let  $x, y \in \mathcal{X}$  and  $(c^k, v^k)$  be given as in the statement. Using the Cauchy-Schwarz inequality, we see that for any Borel  $I \subset [0, T]$  we have the a priori estimate on  $v^k$ 

$$\int_{I} \frac{1}{2} \sum_{x,y} \left| v_{xy}^{k}(t) \right| dt \leq \int_{0}^{T} \left( \vec{\mathcal{A}} \left( c^{k}(t), v^{k}(t) \right) \right)^{\frac{1}{2}} \left( \frac{1}{2} \sum_{x,y} w_{xy} \left( c(t) \right) \right)^{\frac{1}{2}} dt$$

$$\leq \sqrt{CT} \sqrt{C_{w} |I|},$$

with  $w_{xy}(c(t))$  from (2.7). Since Q is continuous on  $\mathcal{P}(\mathcal{X})$  by Definition 2.3,

$$\sup_{\nu \in \mathcal{P}(\mathcal{X})} \frac{1}{2} \sum_{x,y} w_{xy} (\nu) = C_w < \infty.$$

Together with the assumption (4.4), the whole r.h.s. is uniformly bounded in k. Therefore, for a subsequence holds  $v_{xy}^k \rightharpoonup v_{xy}$  as Borel measure on [0,T] and all  $x,y \in \mathcal{X}$ . Now, we choose a sequence of smooth test functions  $\varphi^{\varepsilon}$  in (2.12), which converge to the indicator of the interval  $[t_1,t_2]$  as  $\varepsilon \to 0$ . Therewith and using the above a priori estimate on  $v^k$ , we deduce

$$\left| c_x^k(t_2) - c_x^k(t_1) \right| \le \int_{t_1}^{t_2} \frac{1}{2} \sum_{y \in \mathcal{X}} \left( \left| v_{xy}^k(t) \right| + \left| v_{yx}^k(t) \right| \right) dt \le \sqrt{CC_w} \sqrt{|t_2 - t_1|}.$$

Hence,  $c^k$  is equi-continuous and therefore converges (upto a further subsequence) to some continuous curve c. This, already implies that we can pass to the limit in (2.12) and obtain that  $(c, v) \in CE_T$ .

Moreover, we can deduce since  $\nu \mapsto Q(\nu)$  is continuous for all  $x, y \in \mathcal{X}$  also  $c_{1;x,y}^k := c_x^k(t)Q_{xy}(c^k(t)) \to c_x(t)Q_{xy}(c(t)) =: c_{1;x,y}(t)$  and analogue with  $c_{2;x,y}^k := c_y^k(t)Q_{yx}(c^k(t))$ . We rewrite the action (2.14) as

$$\vec{\mathcal{A}}(c^k(t),v^k(t)) = \frac{1}{2} \sum_{x,y} \alpha \Big( v^k_{x,y}(t), \Lambda \Big( c^k_{1;x,y}(t), c^k_{2;x,y}(t) \Big) \Big)$$

The conclusion (4.5) follows now from [8, Thm. 3.4.3] by noting that  $(v, c_1, c_2) \mapsto \alpha(v, \Lambda(c_1, c_2))$  is l.s.c., jointly convex and 1-homogeneous and hence

$$\liminf_{k} \int_{0}^{T} \vec{\mathcal{A}}(c^{k}(t), v^{k}(t)) dt \ge \int_{0}^{T} \frac{1}{2} \sum_{x,y} \alpha(v_{x,y}(t), \Lambda(c_{1;x,y}(t), c_{2;x,y}(t))) dt$$

$$= \int_{0}^{T} \vec{\mathcal{A}}(c(t), v(t)) dt. \qquad \Box$$

We can now give the proof of Proposition 2.10:

Proof of Proposition 2.10. Symmetry of  $\mathcal{W}$  is obvious, the coincidence axiom follows from Lemma 4.2 and the triangular inequality from Lemma 4.1 as indicated above. The finiteness of  $\mathcal{W}$  comes by using Lemmas 4.2, 4.3 and the triangular inequality. Thus  $\mathcal{W}$  defines a metric. Completeness and separability follow directly from Lemmas 4.4 and 4.2. By the direct method of the calculus of variations and the compactness results Proposition 4.5, we obtain for any  $\mu, \nu \in \mathcal{P}(\mathcal{X})$  a curve  $(\gamma_t)_{t \in [0,1]}$  with minimal action connecting them, i.e.  $\mathcal{W}(\mu, \nu) = \int_0^1 \mathcal{A}(\gamma_t, \psi_t) dt = \int_0^1 |\gamma_t'|^2 dt$ , where in the last equality we used Proposition 2.11. From this, it is easy to see that  $\gamma$  is a constant speed geodesic.  $\square$ 

APPENDIX A. STIRLING FORMULA WITH EXPLICIT ERROR ESTIMATE

**Lemma A.1.** Let  $\nu \in \mathcal{P}_N(\mathcal{X})$ , then it holds

$$-\frac{\log(N+1)}{N} \le -\frac{1}{N}\log\frac{N!}{\prod_{x\in\mathcal{X}}(N\nu(x))!} - \sum_{x\in\mathcal{X}}\nu(x)\log\nu(x) \le \frac{|\mathcal{X}|\log N}{N} + \frac{1}{N}.$$

*Proof.* We write

$$-\log \frac{N!}{\prod_{x \in \mathcal{X}} (N\nu(x))!} = \sum_{x \in \mathcal{X}} \sum_{k=1}^{N\nu(x)} \log k - \sum_{k=1}^{N} \log k$$

$$\geq \sum_{x \in \mathcal{X}} \int_{1}^{N\nu(x)} \log y \, dy - \int_{1}^{N+1} \log(y) \, dy$$

$$= \sum_{x \in \mathcal{X}} \left( N\nu(x) \left( \log N\nu(x) - 1 \right) - 1 \right) - (N+1) \left( \log(N+1) - 1 \right) - 1$$

$$= N \sum_{x \in \mathcal{X}} \nu(x) \log \nu(x) + NR_{N},$$

where the remainder  $R_N$  can be estimated as follows

$$R_N = \frac{|\mathcal{X}|}{N} + \log \frac{N}{N+1} - \frac{\log(N+1)}{N} \ge -\frac{\log(N+1)}{N}$$

for  $|\mathcal{X}| \geq 2$  and  $N \geq 1$ . The other bound can be obtained by shifting the integration bounds appropriately in the above estimate.

# APPENDIX B. VARIANCE ESTIMATE FOR THE PARTICLE SYSTEM

**Lemma B.1.** For the N-Particle process  $X^N$  with generator 3.6 holds for some C > 0 and all  $t \in [0,T]$  with  $T < \infty$ 

$$\forall x \in \mathcal{X}: \operatorname{var}\left(L_x^N(\boldsymbol{X}^N(t))\right) \leq e^{Ct}\left(\operatorname{var}\left(L_x^N(\boldsymbol{X}^N(0))\right) + O(N^{-1})\right).$$

*Proof.* We denote with  $N_x(t) = NL_x^N(\mathbf{X}^N(t))$  the empirical process of the particle number at site x. The empirical density process of particles at site x is then  $N_x(t)/N = L_x^N(\mathbf{X}^N(t))$ . Therewith, we have

$$\frac{d}{dt}\operatorname{var}(N_{x}(t)) = \mathbb{E}[\mathcal{L}^{N}N_{x}^{2}(t)] - 2\mathbb{E}[N_{x}(t)]\mathbb{E}[\mathcal{L}^{N}N_{x}(t)]$$

$$= \mathbb{E}\left[\sum_{y}N_{x}(t)(N_{x}^{2}(t) - (N_{x}(t) - 1)^{2})Q_{xy}^{N}(L^{N}(\mathbf{X}^{N}(t)))\right]$$

$$+ \sum_{y}N_{y}(t)(N_{x}^{2}(t) - (N_{x}(t) + 1)^{2})Q_{yx}^{N}(L^{N}(\mathbf{X}^{N}(t)))\right]$$

$$- 2\mathbb{E}[N_{x}(t)]\mathbb{E}\left[\sum_{y}N_{x}(t)Q_{xy}^{N}(L^{N}(\mathbf{X}^{N}(t))) - N_{y}(t)Q_{yx}^{N}(L^{N}(\mathbf{X}^{N}(t)))\right]$$

$$= 2\mathbb{E}[N_{x}(t)^{2}Q_{xy}(L^{N}(\mathbf{X}^{N}(t)))] - 2\mathbb{E}\left[\sum_{y}N_{x}(t)N_{y}(t)Q_{yx}^{N}(L^{N}(\mathbf{X}^{N}(t)))\right]$$

$$- 2\mathbb{E}[N_{x}(t)]\mathbb{E}\left[\sum_{y}N_{x}(t)Q_{xy}^{N}(L^{N}(\mathbf{X}^{N}(t))) - N_{y}(t)Q_{yx}^{N}(L^{N}(\mathbf{X}^{N}(t)))\right] + O(N)$$

$$\leq C\operatorname{var}(N_{x}(t)) + C\sum_{y\neq x}\operatorname{var}(N_{x}(t))^{1/2}\operatorname{var}N_{y}(t)^{1/2}$$

$$\leq C\sum_{y}\operatorname{var}(N_{y}(t)) + O(N).$$

In theses computations, we used the fact that  $Q^N$  is uniformly bounded and that the state space is finite.

Hence

$$\frac{d}{dt}\operatorname{var}(N_x(t)/N) \le C\sum_{y}\operatorname{var}(N_y(t)/N) + O(N^{-1})$$

and therefore, using Gronwall's Lemma, as soon as the sum of initial variances goes to zero when N goes to infinity, it also goes to zero at any positive time, and uniformly on bounded time intervals.

### Acknowledgments

This work was done while the authors were enjoying the hospitality of the Hausdorff Research Institute for Mathematics during the Junior Trimester Program on Optimal Transport, whose support is gratefully acknowledged. We would like to thank Hong Duong for discussions on this topic. M.F. gratefully acknowledges funding from NSF FRG grant DMS-1361185 and GdR MOMAS. M.E. and A.S. acknowledge support by the German Research Foundation through the Collaborative Research Center 1060 *The Mathematics of Emergent Effects*.

#### References

- [1] S. Adams, N. Dirr, M. A. Peletier and J. Zimmer, From a large-deviations principle to the Wasser-stein gradient flow: a new micro-macro passage, Comm. Math. Phys., 307 (2011), 791–815.
- [2] L. Ambrosio, N. Gigli and G. Savaré, Gradient Flows in Metric Spaces and in the Space of Probability Measures, 2nd edition, Lectures in Mathematics ETH Zürich, Birkhäuser Verlag, Basel, 2008.
- [3] L. Ambrosio, G. Savaré and L. Zambotti, Existence and stability for Fokker-Planck equations with log-concave reference measure, Probab. Theory Related Fields, 145 (2009), 517–564.
- [4] P. Billingsley, *Probability and Measure*, 2nd edition, Wiley Series in Probability and Mathematical Statistics: Probability and Mathematical Statistics, John Wiley & Sons, Inc., New York, 1999.
- [5] F. Bolley, A. Guillin and C. Villani, Quantitative concentration inequalities for empirical measures on non-compact spaces, Probab. Theory Related Fields, 137 (2007), 541–593.
- [6] A. Budhiraja, P. Dupuis, M. Fischer and K. Ramanan, Limits of relative entropies associated with weakly interacting particle systems, Electron. J. Probab., 20 (2015), 22pp.
- [7] A. Budhiraja, P. Dupuis, M. Fischer and K. Ramanan, Local stability of Kolmogorov forward equations for finite state nonlinear Markov processes, Electron. J. Probab., 20 (2015), 30pp.
- [8] G. Buttazzo, Semicontinuity, Relaxation and Integral Representation in the Calculus of Variations, vol. 207 of Pitman Research Notes in Mathematics Series, Longman Scientific & Technical, Harlow; copublished in the United States with John Wiley & Sons, Inc., New York, 1989.
- [9] J. A. Carrillo, R. J. McCann and C. Villani, Kinetic equilibration rates for granular media and related equations: entropy dissipation and mass transportation estimates, Rev. Mat. Iberoamericana, 19 (2003), 971–1018.
- [10] J. A. Carrillo, R. J. McCann and C. Villani, Contractions in the 2-Wasserstein length space and thermalization of granular media, Arch. Ration. Mech. Anal., 179 (2006), 217–263.
- [11] P. Cattiaux, A. Guillin and F. Malrieu, Probabilistic approach for granular media equations in the non-uniformly convex case, Probab. Theory Related Fields, 140 (2008), 19–40.
- [12] P. Dai Pra and F. den Hollander, McKean-Vlasov limit for interacting random processes in random media, J. Statist. Phys., 84 (1996), 735–772.
- [13] S. Daneri and G. Savaré, Lecture notes on gradient flows and optimal transport, in Optimal Transportation, vol. 413 of London Math. Soc. Lecture Note Ser., Cambridge Univ. Press, Cambridge, 2014, 100–144.
- [14] D. A. Dawson and J. Gärtner, Large deviations from the McKean-Vlasov limit for weakly interacting diffusions, Stochastics, 20 (1987), 247–308.
- [15] E. De Giorgi, A. Marino and M. Tosques, Problems of evolution in metric spaces and maximal decreasing curve, Atti Accad. Naz. Lincei Rend. Cl. Sci. Fis. Mat. Natur. (8), 68 (1980), 180–187.
- [16] N. Dirr, V. Laschos and J. Zimmer, Upscaling from particle models to entropic gradient flows, J. Math. Phys., 53 (2012), 063704, 9 pp.
- [17] R. Dobrushin, Vlasov equations, Functional Analysis and Its Applications, 13 (1979), 48–58,96.
- [18] J. Dolbeault, B. Nazaret and G. Savaré, A new class of transport distances between measures, Calc. Var. Partial Differential Equations, 34 (2009), 193–231.
- [19] M. H. Duong, V. Laschos and M. Renger, Wasserstein gradient flows from large deviations of manyparticle limits, ESAIM Control Optim. Calc. Var., 19 (2013), 1166–1188.
- [20] P. Dupuis and R. S. Ellis, A Weak Convergence Approach to the Theory of Large Deviations, Wiley Series in Probability and Statistics: Probability and Statistics, John Wiley & Sons, Inc., New York, 1997, A Wiley-Interscience Publication.
- [21] M. Erbar, Gradient flows of the entropy for jump processes, Ann. Inst. H. Poincaré Probab. Statist., 50 (2014), 920–945.
- [22] M. Erbar and J. Maas, Ricci curvature of finite Markov chains via convexity of the entropy, Arch. Ration. Mech. Anal., 206 (2012), 997–1038.

- [23] M. Erbar and J. Maas, Gradient flow structures for discrete porous medium equations, Discrete Contin. Dyn. Syst., 34 (2014), 1355–1374.
- [24] M. Erbar, J. Maas and M. Renger, From large deviations to Wasserstein gradient flows in multiple dimensions, Electron. Commun. Probab., 20 (2015), 1–12.
- [25] [10.1016/j.matpur.2016.03.018] M. Fathi, A gradient flow approach to large deviations for diffusion processes, J. Math. Pures Appl., (2016).
- [26] M. Fathi and M. Simon, The gradient flow approach to hydrodynamic limits for the simple exclusion process, In P. Gonçalves and A. J. Soares, From Particle Systems to Partial Differential Equations III: Particle Systems and PDEs III, Braga, Portugal, December 2014, 167–184, Springer International Publishing, Cham, 2016.
- [27] N. Gigli and J. Maas, Gromov-Hausdorff convergence of discrete transportation metrics, SIAM J. Math. Anal., 45 (2013), 879–899.
- [28] R. Jordan, D. Kinderlehrer and F. Otto, The variational formulation of the Fokker-Planck equation, SIAM J. Math. Anal., 29 (1998), 1–17.
- [29] C. Kipnis and C. Landim, Scaling Limits of Interacting Particle Systems, vol. 320 of Grundlehren der Mathematischen Wissenschaften, Springer-Verlag, Berlin, 1999.
- [30] D. A. Levin, M. J. Luczak and Y. Peres, Glauber dynamics for the mean-field Ising model: Cut-off, critical power law, and metastability, Probab. Theory Related Fields, 146 (2010), 223–265.
- [31] J. Maas, Gradient flows of the entropy for finite Markov chains, J. Funct. Anal., 261 (2011), 2250–2292.
- [32] F. Malrieu, Convergence to equilibrium for granular media equations and their Euler schemes, Ann. Appl. Probab., 13 (2003), 540–560.
- [33] [10.1063/1.1699114] N. Metropolis, A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller and E. Teller, Equations of state calculations by fast computing machines, J. Chem. Phys., 21 (1953), 1087–1091.
- [34] A. Mielke, Geodesic convexity of the relative entropy in reversible Markov chains, Calc. Var. Partial Differential Equations, 48 (2013), 1–31.
- [35] A. Mielke, On evolutionary Γ-convergence for gradient systems, Springer International Publishing, Cham, 3 (2016), 187–249.
- [36] K. Oelschläger, A martingale approach to the law of large numbers for weakly interacting stochastic processes, Ann. Probab., 12 (1984), 458–479.
- [37] F. Otto, The geometry of dissipative evolution equations: The porous medium equation, Comm. Partial Differential Equations, 26 (2001), 101–174.
- [38] E. Sandier and S. Serfaty, Gamma-convergence of gradient flows with applications to Ginzburg-Landau, Comm. Pure Appl. Math., 57 (2004), 1627–1672.
- [39] A. Schlichting, Macroscopic limits of the Becker-Döring equations via gradient flows, arXiv: 1607.08735.
- [40] S. Serfaty, Gamma-convergence of gradient flows on Hilbert and metric spaces and applications, Discrete Contin. Dyn. Syst., 31 (2011), 1427–1451.
- [41] A.-S. Sznitman, Topics in Propagation of Chaos, in École d'Été de Probabilités de Saint-Flour XIX—1989, vol. 1464 of Lecture Notes in Math., Springer, Berlin, 1991, 165–251.

University of Bonn, Germany

 $E ext{-}mail\ address: erbar@iam.uni-bonn.de}$ 

UNIVERSITY OF CALIFORNIA, BERKELEY E-mail address: maxf@berkeley.edu

Weierstrass Institute

 $E ext{-}mail\ address: Vaios.laschos@wias-berlin.de}$ 

University of Bonn, Germany

 $E ext{-}mail\ address: schlichting@iam.uni-bonn.de}$