

A Background on the Schrödinger bridge problem

The following facts can be found for instance in the lecture notes [33, Sec. 4]. Consider $\mu, \nu \in \mathcal{P}(\mathcal{X})$ and the entropy-regularized optimal transport problem defining $T_\tau(\mu, \nu)$ in Eq. (6) with a cost function $c \in \mathcal{C}^\infty(\mathcal{X} \times \mathcal{X})$. This problem admits a unique solution γ^* and admits a dual formulation

$$T_\tau(\mu, \nu) = \max_{\varphi \in L^1(\mu), \psi \in L^1(\nu)} \int \varphi d\mu + \int \psi d\nu + \tau \left(1 - \int e^{(\varphi(x) + \psi(y) - c(x, y))/\tau} d\mu(x) d\nu(y) \right). \quad (15)$$

The dual problem admits a unique solution in $L^1(\mu) \times L^1(\nu)$ up to the transformation $(\varphi + \kappa, \psi - \kappa)$ for $\kappa \in \mathbb{R}$. Moreover, we have that $T_\tau(\mu, \nu) = \int \varphi d\mu + \int \psi d\nu$ and the primal-dual relation

$$\frac{d\gamma^*}{d\mu \otimes \nu}(x, y) = g(x, y) := e^{(\varphi(x) + \psi(y) - c(x, y))/\tau}, \quad \mu \otimes \nu \text{ a.e.}$$

Also, the potentials satisfy the following equations, which are the first-order optimality conditions:

$$\begin{cases} \varphi(x) = -\tau \log \left(\int e^{(\psi(y) - c(x, y))/\tau} d\nu(y) \right) \\ \psi(y) = -\tau \log \left(\int e^{(\varphi(x) - c(x, y))/\tau} d\mu(x) \right) \end{cases} \quad (16)$$

These equations are a priori only satisfied μ (resp. ν) almost everywhere, but they can be used to extend φ and ψ as continuous (in fact infinitely differentiable) functions over \mathcal{X} , which satisfy these equations everywhere and are unique in $\mathcal{C}^\infty(\mathcal{X})$, up to the additive invariance mentioned above. Throughout the paper, we refer to such functions (φ, ψ) as the *Schrödinger potentials* (they are also referred as EOT potentials in [33]). Let us conclude this section with two observations.

- (i) One can bound the oscillation of φ as follows

$$\text{osc}(\varphi) := \sup_x \varphi(x) - \inf_x \varphi(x) \leq \sup_{x, y} c(x, y) - \inf_{x, y} c(x, y) = \text{osc}(c). \quad (17)$$

This is obtained by upper bounding c inside the integral which leads to $\varphi(x) \leq \sup_{x, y} c(x, y) - \tau \log \left(\int e^{(\psi(y))/\tau} d\nu(y) \right)$. Subtracting the analogous lower bound, we observe that the log terms cancel and we get the bound on the oscillation.

- (ii) One can differentiate (16) to see that we have for $x \in \mathcal{X}$

$$\nabla \varphi(x) = \int \nabla_x c(x, y) d\gamma^*(y|x) = \mathbf{E}[\nabla_x c(X, Y) | X = x] \quad (18)$$

with $\text{Law}(X, Y) = \gamma^*$ and $\gamma^*(dy|x) = \int g(x, y) \nu(dy)$ is the conditional distribution of Y given X .

B Proof of Theorem 3.1

Let us first recall the statement of Thm. 3.1.

Theorem B.1 (Representer Theorem). *Let $\text{Fit} : \mathcal{P}(\mathcal{X})^T \rightarrow \mathbb{R}$ be any function.*

- (i) *If \mathcal{F} (Eq.(5)) admits a minimizer R^* then $(R_{t_1}^*, \dots, R_{t_T}^*)$ is a minimizer for F (Eq.(8)).*
(ii) *Conversely, if F admits a minimizer $\mu^* \in \mathcal{P}(\mathcal{X})^T$ then a minimizer for \mathcal{F} is built as*

$$R(\cdot) = \int_{\mathcal{X}^T} W^\tau(\cdot | x_1, \dots, x_T) dR_{t_1, \dots, t_T}(x_1, \dots, x_T)$$

where $W^\tau(\cdot | x_1, \dots, x_T)$ is the law of W^τ conditioned on passing through x_1, \dots, x_T at times t_1, \dots, t_T respectively and R_{t_1, \dots, t_T} is the composition of the transport plans $\gamma_{i, i+1}$ which are optimal in the definition of $T_{\tau_i}(\mu^{*(i)}, \mu^{*(i+1)})$, for $i = 1, \dots, T$.

This theorem is a direct consequence of the following lemma, which is similar in spirit to [1, Prop. B.1], but the terms involved are different.

Lemma B.2. *There exists $C > 0$ such that, for any $R \in \mathcal{P}(\Omega)$ and t_1, \dots, t_T a collection of instants, it holds*

$$\begin{aligned} H(R | W^\tau) &\stackrel{(\dagger)}{\geq} H(R_{t_1, \dots, t_T} | W_{t_1, \dots, t_T}^\tau) \\ &\stackrel{(\star)}{\geq} \sum_{i=1}^{T-1} H(R_{t_i, t_{i+1}} | p_{\tau_i}(R_{t_i} \otimes R_{t_{i+1}})) + \sum_{i=1}^T H(R_{t_i}) + C. \end{aligned}$$

The first inequality (\dagger) becomes an equality if and only if

$$R(\cdot) = \int_{\mathcal{X}^T} W(\cdot | x_1, \dots, x_T) dR_{t_1, \dots, t_T}(x_1, \dots, x_T)$$

where $W^\tau(\cdot | x_1, \dots, x_T)$ is the law of W^τ conditioned on passing through x_1, \dots, x_T at times t_1, \dots, t_T respectively. In addition, the second inequality (\star) becomes an equality if and only if R is Markovian.

Proof. The first inequality (\dagger) and the equality case follows from the additivity property of the relative entropy under conditioning [16, Eq. (A9)], namely that it holds

$$H(R | W^\tau) = H(R_{t_1, \dots, t_T} | W_{t_1, \dots, t_T}^\tau) + \int H(R(\cdot | x_1, \dots, x_T) | W^\tau(\cdot | x_1, \dots, x_T)) dR_{t_1, \dots, t_T}(x_1, \dots, x_T)$$

where the second term vanishes if and only if the conditional distributions $R(\cdot | x_1, \dots, x_T)$ are Brownian bridges, for R_{t_1, \dots, t_T} almost every (x_1, \dots, x_T) . For the second inequality (\star) , [20, Sec. 3.4] states that

$$H(R_{t_1, \dots, t_T} | W_{t_1, \dots, t_T}^\tau) \geq \sum_{i=1}^{T-1} H(R_{t_i, t_{i+1}} | W_{t_i, t_{i+1}}^\tau) - \sum_{i=2}^{T-1} H(R_{t_i} | W_{t_i}^\tau) =: E$$

with equality if and only if R_{t_1, \dots, t_T} is Markovian. This is the formula of [1, Prop. B.1], but this expression is unsuitable for our purposes and we need to further reorganise the terms in E .

Without loss of generality, let us assume that R_{t_i} are absolutely continuous with density $\frac{dR_{t_i}}{dx}(x) = r_i(x)$ (if this is not the case, both sides of the inequality are infinite) and let $V_{\mathcal{X}}$ be the Lebesgue volume of \mathcal{X} . On the one hand, since $W_{t_i}^\tau$ is the normalized volume measure on \mathcal{X} , it holds

$$H(R_{t_i} | W_{t_i}^\tau) = H(R_{t_i}) + \log(V_{\mathcal{X}}).$$

On the other hand, letting $\tau_i = \tau(t_{i+1} - t_i)$ it holds

$$W_{t_i, t_{i+1}}^\tau(dx, dy) = V_{\mathcal{X}}^{-1} p_{\tau_i}(x, y) dx dy$$

by definition of the transition probability density p of the Brownian motion on \mathcal{X} . It follows that for any $\mu, \nu \in \mathcal{P}(\mathcal{X})$ with finite differential entropy and $\gamma \in \Pi(\mu, \nu)$ it holds

$$\begin{aligned} H(\gamma | W_{t_i, t_{i+1}}^\tau) &= \int \log\left(\frac{d\gamma}{dx \otimes dy} \frac{V_{\mathcal{X}}}{p_{\tau_i}}\right) d\gamma(x, y) \\ &= \log(V_{\mathcal{X}}) + \int \log\left(\frac{d\gamma}{p_{\tau_i} d\mu \otimes \nu} \frac{d\mu}{dx} \frac{d\nu}{dy}\right) d\gamma \\ &= \log(V_{\mathcal{X}}) + H(\gamma | p_{\tau_i} \mu \otimes \nu) + H(\mu) + H(\nu) \end{aligned}$$

where we used the fact that $\gamma \in \Pi(\mu, \nu)$ to simplify the two last terms (see [34, Lem. 1.6] for more details on the change of reference measure in regularized optimal transport). Putting everything together, and using the fact that $R_{t_i, t_{i+1}} \in \Pi(R_{t_i}, R_{t_{i+1}})$ we get

$$\begin{aligned} E &= \log(V_{\mathcal{X}}) + \sum_{i=1}^{T-1} H(R_{t_i, t_{i+1}} | p_{\tau_i} R_{t_i} \otimes R_{t_{i+1}}) + \sum_{i=1}^{T-1} H(R_{t_i}) + \sum_{i=2}^T H(R_{t_i}) - \sum_{i=2}^{T-1} H(R_{t_i}) \\ &= \log(V_{\mathcal{X}}) + \sum_{i=1}^{T-1} H(R_{t_i, t_{i+1}} | p_{\tau_i} R_{t_i} \otimes R_{t_{i+1}}) + \sum_{i=1}^T H(R_{t_i}). \end{aligned}$$

which proves the formula. \square

Proof of Thm. 3.1. Clearly, a minimizer $R^* \in \mathcal{P}(\Omega)$ of $\mathcal{F}(R) = \text{Fit}(R_1, \dots, R_T) + \tau H(R|W^\tau)$ is of the form given by Lem. B.2. Let $\mu^{(i)} = R_{t_i}^*$ be its marginals and $\gamma^{(i)} = R_{t_i, t_{i+1}}^*$ which clearly satisfies $\gamma^{(i)} \in \Pi(\mu^{(i)}, \mu^{(i+1)})$. It holds, with $C = \log(V_{\mathcal{X}})$,

$$\begin{aligned} \mathcal{F}(R^*) &= \text{Fit}(\mu^{(1)}, \dots, \mu^{(T)}) + \tau \sum_{i=1}^{T-1} H(\gamma^{(i)} | p_{\tau_i} \mu^{(i)} \otimes \mu^{(i+1)}) + \tau H(\mu) + C \\ &\geq \text{Fit}(\mu^{(1)}, \dots, \mu^{(T)}) + \frac{\tau}{\tau_i} \sum_{i=1}^{T-1} T_{\tau_i}(\mu^{(i)}, \mu^{(i+1)}) + \tau H(\mu) + C \end{aligned}$$

where the last equality holds if and only if $R_{t_i, t_{i+1}}^* = \gamma^{(i)}$ is optimal in the definition of $T_{\tau_i}(\mu^{(i)}, \mu^{(i+1)})$. The claim follows. \square

C Proof of Theorem 3.3

Let us recall the statement of Theorem 3.3 and prove it, by an application of a convergence result proved independently in [14] and [15].

Theorem C.1 (Convergence). *Let $\mu_0 \in \mathcal{P}(\mathcal{X})^T$ be such that $F(\mu_0) < \infty$. Then for $\epsilon \geq 0$, there exists a unique solution $(\mu_s)_{s \geq 0}$ to the MFL dynamics (11). For $\epsilon > 0$, \mathcal{X} the d -torus and moreover assuming that μ_0 has a bounded absolute log-density, it holds*

$$F_\epsilon(\mu_s) - \min F_\epsilon \leq e^{-Cs} (F_\epsilon(\mu_0) - \min F_\epsilon).$$

where $C = \beta e^{-\alpha/\epsilon}$ for some $\alpha, \beta > 0$ independent of μ_0 and ϵ . Moreover, taking a smooth time-dependent ϵ_s that decays asymptotically as $\tilde{\alpha}/\log(s)$ for some $\tilde{\alpha} > \alpha$, it holds $F_0(\mu_s) - F_0(\mu^*) \lesssim \log(\log(s))/\log(s)$ and μ_s converges weakly to μ^* .

Proof. We verify the assumptions of [15], noticing that their proof can be adapted without difficulty to our context of families of T probability measures with $T \geq 1$. In that reference, the objective function is assumed of the form $G + aH$ for some $a > 0$ where H is the entropy and G satisfies certain properties. Here we will split the objective function differently between the well-posedness and the convergence result, so that G satisfies the required properties for each result.

For the well-posedness of the dynamics, we split the objective function F_ϵ as $G + (\tau + \epsilon)H$. [15, Assumption 1], about the stability and regularity of the first-variation V , is guaranteed Prop. C.2 below (the stability and regularity of the component $\delta \text{Fit}/\delta \mu$ is immediate) and leads to the well-posedness of the dynamics without requiring $\epsilon > 0$ (in fact we could even take $\epsilon \in [-\tau, +\infty[$ and have well-posedness).

For the global convergence guarantee we split the objective function F_ϵ as $(G + \tau H) + \epsilon H$ because we need the fact that $F_0 = G + \tau H$ is convex. [15, Assumption 2], which requires convexity of F_0 and existence of a minimizer for F_ϵ , is satisfied thanks to Prop. 3.2. For the uniform Log-Sobolev Inequality (LSI) ([15, Assumption 3]), we first remark that [35, Thm. 7.3] states that the uniform distribution over \mathcal{X} satisfies LSI with a constant ρ_0 that only depends on D the diameter of \mathcal{X} .

Now, the i -th component of the first-variation of F_0 is given by $V^{(i)}[\mu] + \tau \log(\mu_i)$. By the expression of Prop. 3.2, the oscillation $\text{osc}(V^{(i)}[\mu])$ of $V^{(i)}[\mu]$ (i.e. the difference between its maximum and minimum value over \mathcal{X}) is bounded (independently of ϵ and μ_0). Indeed, the gradient formula for $\delta \text{Fit}_\sigma/\delta \mu^{(i)}$ is nonnegative and bounded by $e^{D^2/(2\sigma^2)}$ and by App A Eq. (17), the Schrödinger potential $\varphi_{i,i+1}$ has an oscillation bounded by $\sup_{x,y} c_{\tau_i}(x,y) - \inf_{x,y} c_{\tau_i}(x,y)$ which is $D^2/2$ when $c(x,y) = \frac{1}{2}\|y-x\|^2$ or a less explicit constant that depends only on the domain \mathcal{X} for the cost c_{τ_i} . Combining these estimates with Lemma C.3 (which requires \mathcal{X} to be the d -torus), we get that the oscillation of $V^{(i)}[\mu] + \tau \log(\mu^{(i)})$ is bounded by some $\kappa > 0$ independent of ϵ and μ_0 (under the assumption that the log-density of μ_0 is absolutely bounded by $A > 0$).

It follows, by Holley and Strook perturbation criterion [36] that the density proportional to $e^{-(V^{(i)}[\mu] + \tau \log(\mu^{(i)}))/\epsilon}$ satisfies a LSI with constant $\rho \geq \alpha e^{-\beta/\epsilon}$ for some α, β independent of s, ϵ, μ_0 . Then [15, Thm. 3.2] guarantees the exponential convergence with rate e^{-Cs} with $C = 2\epsilon\rho$.

Finally, the convergence result with simulated annealing is a direct application of [15, Thm. 4.1] which proves the convergence rate in $\log(\log s)/\log s$. Since in addition the sublevel sets of F_0 are weakly compact and the minimizer μ^* is unique, standard considerations imply that μ_s converges weakly to μ^* . \square

Let us report a stability result concerning the Schrödinger potentials, which is a consequence of a more general result in [25] and is used in the proof above.

Proposition C.2. *Assume that $c \in \mathcal{C}^2(\mathcal{X} \times \mathcal{X})$. There exists $C > 0$ such that for all $\mu, \mu', \nu, \nu' \in \mathcal{P}(\mathcal{X})$, it holds*

$$\|\nabla\varphi - \nabla\varphi'\|_\infty + \|\nabla\psi - \nabla\psi'\|_\infty \leq C(W_2(\mu, \mu') + W_2(\nu, \nu')).$$

where (φ, ψ) (resp. (φ', ψ')) are the Schrödinger potentials (see App. A) associated to the pair of measures (μ, ν) (resp. (μ', ν')) and W_2 is the 2-Wasserstein distance.

Lemma C.3. *Let $(\mu_s)_{s \geq 0}$ be the solution of the Mean-Field Langevin Dynamics for $\epsilon \geq 0$ and assume that the absolute log-density of $\mu_0^{(i)}$ is bounded by $A > 0$ for each $i \in \{1, \dots, T\}$. Let \mathcal{X} be the d -torus. Then there exists $A' > 0$ (independent of ϵ) such that the absolute log-density of $\mu_s^{(i)}$ is bounded by A' , for all $s \geq 0$ and $i \in \{1, \dots, T\}$.*

Proof. The i -th component of the stochastic process associated to the Mean-field Langevin dynamics solves $dX_s^{(i)} = -\nabla V^{(i)}[\mu_s](X_s^{(i)})ds + \sqrt{2(\tau + \epsilon)}dB_s^{(i)}$ where $(B_s^{(i)})$ is a Brownian motion independent for each i and $\mu_s = (\text{Law}(X_s^{(1)}), \dots, \text{Law}(X_s^{(T)}))$ (the boundary reflection term is absent since we are considering the torus). Let $S > 0$ be an arbitrary time and let $P^{(i)} \in \mathcal{P}(\mathcal{C}([0, S]; \mathcal{X}))$ be the law of $X^{(i)}$ over the time interval $[0, S]$. By Girsanov's formula (see e.g. [1, Sec. 4.2]), it holds

$$\frac{dP^{(i)}}{dW^{(\tau+\epsilon)}}(X) = \frac{dP_0^{(i)}}{d\text{vol}}(X_0) \exp \left(\frac{2}{\tau + \epsilon} \left(\int_0^S -\nabla V^{(i)}[\mu_s](X_s) dX_s - \frac{1}{2} \int_0^S \|\nabla V^{(i)}[\mu_s](X_s)\|^2 ds \right) \right)$$

where vol is the uniform distribution over \mathcal{X} , and by Ito's formula

$$\begin{aligned} & \int_0^S -\nabla V^{(i)}[\mu_s](X_s) dX_s \\ &= V^{(i)}[\mu_0](X_0) - V^{(i)}[\mu_S](X_S) + \int_0^S \left((\partial_s V^{(i)}[\mu_s])(X_s) + (\tau + \epsilon) \Delta V^{(i)}[\mu_s](X_s) \right) ds. \end{aligned}$$

Since, for S fixed, all the quantities in the exponential are uniformly bounded (in particular, for the term involving $\partial_s V^{(i)}[\mu_s]$ this follows from Prop. C.2 which implies $\int_0^S \|\partial_s V^{(i)}[\mu_s]\|_\infty ds \leq C \sum_{i=1}^T \int_0^S (\int_{\mathcal{X}} \|v_s(x)\|^2 d\mu_s^{(i)}(x))^{1/2} ds \leq C \sqrt{ST} (\sum_{i=1}^T \int_0^S \int_{\mathcal{X}} \|v_s(x)\|^2 d\mu_s^{(i)}(x))^{1/2} \leq C \sqrt{ST} (F_\epsilon(\mu_0) - F_\epsilon(\mu_S))^{1/2}$ where v_s is the Wasserstein derivative of $(\mu_s^{(i)})_s$ and using facts from gradient flows theory [37]). This shows that the transition kernel of the process $X^{(i)}$ is upper and lower bounded by the heat kernel over \mathcal{X} . In particular, the density of $\mu_s^{(i)}$ is upper and lower bounded by positive constants over $[0, S]$. For times S' larger than S , we similarly have that $X_{S'}^{(i)}$ is obtained from $X_{S'-S}^{(i)}$ by a Markov process which is comparable to the heat diffusion on \mathcal{X} and thus its log-density is absolutely bounded. We conclude that the upper and lower bounds on $\log \mu_s^{(i)}$ are uniform in time. \square

D Proof of Proposition 3.2

Proposition D.1. *The function G is convex separately in each of its input (but not jointly), weakly continuous and its first-variation is given for $\mu \in \mathcal{P}(\mathcal{X})^T$ and $i \in [T]$ by*

$$V^{(i)}[\mu] = \frac{\delta \text{Fit}}{\delta \mu^{(i)}}[\mu] + \frac{\varphi_{i,i+1}}{t_{i+1} - t_i} + \frac{\psi_{i,i-1}}{t_i - t_{i-1}}, \quad \frac{\delta \text{Fit}}{\delta \mu^{(i)}}[\mu] : x \mapsto -\frac{\Delta t_i}{\lambda} \int \frac{g_\sigma(x-y)}{(g_\sigma * \mu^{(i)})(y)} d\hat{\mu}_{t_i}(y)$$

where $(\varphi_{i,j}, \psi_{i,j}) \in \mathcal{C}^\infty(\mathcal{X})$ are the Schrödinger potentials for $T_{\tau_i}(\boldsymbol{\mu}^{(i)}, \boldsymbol{\mu}^{(j)})$, with the convention that the corresponding term vanishes when it involves $\psi_{1,0}$ or $\varphi_{T,T+1}$. The function F is jointly convex and admits a unique minimizer $\boldsymbol{\mu}^*$, which has an absolutely continuous density (again denoted by $\boldsymbol{\mu}^*$) characterized by

$$(\boldsymbol{\mu}^*)^{(i)} \propto e^{-V^{(i)}[\boldsymbol{\mu}^*]/\tau}, \quad \text{for } i \in [T].$$

Proof. The properties and formulas for G and its first-variation follow from those for T_τ which are well-known (see [38, Prop. 2]) and for F it which are immediate. In particular, T_τ is convex separately in each of its variables as a supremum of linear forms but not jointly because of the product measure in Eq. (15). The joint convexity of F can be seen by a change of variable in Eq. (15): if μ admits a density ρ_μ with respect to Lebesgue, letting $\tilde{\varphi} = \varphi + \tau \log(\rho_\mu)$ it holds

$$\begin{aligned} T_\tau(\mu, \nu) &= \max_{\varphi \in L^1(\mu), \psi \in L^1(\nu)} \int \varphi d\mu + \int \psi d\nu + \tau \left(1 - \int e^{(\varphi(x) + \psi(y) - c(x,y))/\tau} d\mu(x) d\nu(y) \right). \\ &= \max_{\tilde{\varphi} \in L^1(\mu), \psi \in L^1(\nu)} \int \tilde{\varphi} d\mu + \int \psi d\nu + \tau \left(1 - \int e^{(\tilde{\varphi}(x) + \psi(y) - c(x,y))/\tau} dx d\nu(y) \right) - \tau H(\mu) \end{aligned}$$

Thus, as a supremum of linear forms $(\mu, \nu) \mapsto T_\tau(\mu, \nu) + \tau H(\mu)$ is a convex function. Now F is a sum of $T - 1$ terms of this form and convex functions, and is thus convex. For the uniqueness of the minimizer of F , note that \mathcal{F} is strictly convex and thus admits at most one minimizer. Since by Thm. 3.1, any minimizer of F can be mapped to a unique minimizer of \mathcal{F} , it follows that F has at most one minimizer. Its existence is guaranteed by the direct method in the calculus of variation, because F is weakly lower-semicontinuous and the set $\mathcal{P}(\mathcal{X})^T$ is weakly compact. Finally, the characterization of the minimizer can be formally deduced by writing the first order optimality condition $V^{(i)}[\boldsymbol{\mu}^*] + \tau \log(\boldsymbol{\mu}^{*(i)}) = 0$. The rigorous argument, which is standard, can be found in a similar context e.g. in [13, Lem. 10.4]. \square

E Solving the Schrödinger Bridge Problems with Sinkhorn's algorithm

At each iteration, in order to compute $V[\hat{\mu}[k]]$, one needs to compute the Schrödinger potentials $(\varphi_{i,i+1}, \psi_{i,i+1})$ associated to the $T - 1$ Schrödinger bridges problems $T_\tau(\hat{\mu}^{(i)}, \hat{\mu}^{(i+1)})$ (see Prop. 3.2). Among the various algorithms that can solve this problem [39], let us focus our discussion on the well-studied Sinkhorn's algorithm, which is alternate block maximization on the dual of Eq. (6).

Given two discrete probability measures $\hat{\mu}_m = \sum_{i=1}^m p_i \delta_{x_i}$ and $\hat{\nu}_m = \sum_{i=1}^m q_i \delta_{y_i}$ we define the cost matrix with entries $c_{i,j} = c(x_i, y_j)$ (which we approximate with $\frac{1}{2} \|x_i - y_j\|^2$ using Varadhan's formula). The iterates $u[\ell], v[\ell] \in \mathbb{R}^m$, $\ell \geq 1$ of Sinkhorn's algorithm are defined as :

$$u_i[\ell] = -\tau \log \left(\sum_{j=1}^m e^{(v_j[\ell-1] - c_{i,j})/\tau} q_j \right) \quad \text{and} \quad v_j[\ell] = -\tau \log \left(\sum_{i=1}^m e^{(u_i[\ell] - c_{i,j})/\tau} p_i \right).$$

This algorithm converges in value at a rate $O(1/(\tau k))$, see [40, 41]. In practice, Sinkhorn's iterations could be further sped up with non-linear acceleration methods [42]. Upon convergence, one can recover the Schrödinger potential $\varphi, \psi \in \mathcal{C}^\infty(\mathcal{X})$ via the formula

$$\varphi(x) = -\tau \log \left(\sum_{j=1}^m e^{(v_j^* - c(x, y_j))/\tau} q_j \right), \quad \psi(y) = -\tau \log \left(\sum_{i=1}^m e^{(u_i^* - c(x_i, y))/\tau} p_i \right). \quad (19)$$

Moreover the minimizer in Eq. (6), needed to recover P^* by Thm. 3.1, is given by $\gamma = \sum_{i,j} e^{(u_i^* + v_j^* - c_{i,j})/\tau} p_i q_j \delta_{(x_i, y_j)}$. Those consideration suggest two methods to implement Eq. (13):

- (i) estimate ∇G_m with automatic differentiation by backpropagating through a fixed number of Sinkhorn's iterations, or
- (ii) use the formula for ∇V given in Prop. 3.2 and plug-in the potentials of Eq. (19).

These alternatives are discussed in [43, Sec. 5.3] where (φ, ψ) are computed as a subroutine. There, the conclusion is that option (ii) is slightly more efficient, and this is the method we implemented.

F Simulated Annealing

A standard heuristic to accelerate diffusion-based algorithms with a temperature parameter τ is the so-called *simulated annealing* method [44], which consists in starting from a large value τ_0 and slowly decreasing it towards the desired value while the algorithm runs. In our context, a larger value for τ accelerates the convergence of both Sinkhorn’s algorithm and the MFL dynamics.

While Theorem 3.3 guarantees convergence for a temperature that decays sufficiently slowly, for practical purposes we opt for a faster rate. We used simulated annealing with geometric decay $\tau_t = \max\{cr^t, \tau_f\}$ for $r \in]0, 1[$ and $c > 0$ in order to quickly reach the desired temperature $\tau_f > 0$. Empirically, we find that also allowing the step size η and the (squared) data-fitting bandwidth σ^2 to scale with the temperature leads to further acceleration of the MFL dynamics, especially early on in the optimisation. We illustrate this for the example of Section 4.1 in Figure 5. In this experiment, particles of the MFL are started from $\mathcal{N}(0, 1.0)$. Consequently, some particles are distant from the data and require MFL to be run for a large number of iterations to converge. On the other hand, simulated annealing in (τ, σ, η) for the first 500 iterations with $\tau_0 = 5\tau$, followed by 2,000 iterations without annealing allows for the MFL dynamics to reach convergence much faster. The brief annealing phase allows for particles to quickly move away from their initial distribution towards what is a refined initial condition for a subsequent optimisation without annealing.

As can be seen in Figure 5(b), the result of MFL with annealing is slightly less noisy. To prevent noise at high temperature from causing particles in the MFL dynamics to stray far away from the observed data, we add a confining potential to the objective (7)

$$\text{Confine}_\sigma(R, \hat{\mu}) = - \int dR(y) \log \left[\int e^{\frac{-1}{2\sigma^2} \|x-y\|_2^2} d\hat{\mu}(x) \right],$$

where we have written $R = T^{-1} \sum_{i=1}^T R_{t_i}$ and $\hat{\mu} = T^{-1} \sum_{i=1}^T \hat{\mu}_{t_i}$ to be the mixtures (over time) respectively of the reconstructed and observed marginals. This has the effect of penalizing particles in the MFL dynamics which are far away from any observation. In the annealing examples, we used $\sigma = 5.0$ for the confining potential bandwidth.

We also investigated the setting where the temperature τ was fixed but annealing was carried out in the additional entropy term $\varepsilon \rightarrow 0$, as is analyzed in Theorem 3.3. As we can see in Figure 5, this actually results in a MFL dynamics that converges slower than MFL without annealing due to the additional level of noise injected. Interestingly, these results suggest that the lack of strong convexity when $\varepsilon = 0$ may not affect convergence in practice.

G Simulation details

All numerical experiments were run using a CPU-based implementation of MFL dynamics. A copy of the code to reproduce the figures in this article is available at <https://github.com/zsteve/mfl>.

In practice, we observed that the particles of the discretized MFL dynamics remain at a small bounded distance from the support of the observations throughout; we thus did not explicitly enforce the reflecting or periodic boundary.

G.1 Additional details for Section 4.1

We consider a bifurcating stochastic process (see Figure 1) in ambient space $\mathcal{X} \subset \mathbb{R}^{10}$, following (1) with $\tau = 1/4$ and time-dependent potential $\Psi(t, x) = \frac{1}{2}(x_1 - 1.5)^2(x_1 + 1.5)^2 + 10(x_2 + t)^2 + 10 \sum_{k=3}^{10} x_k^2$. Starting from an initial distribution $X_0 \sim \mathcal{N}(0, 0.1^2)$, particles are simulated over $t \in [0, 1.25]$ and were independently sampled at 10 evenly spaced timepoints $t_i, 1 \leq i \leq 10$.

The discretized MFL dynamics described in Section 3.4 was applied to each simulated dataset with $m = 100$ particles per timepoint and using the data fitting functional (3) with bandwidth $\sigma = 0.5$. Particles in the MFL dynamics were started from $\mathcal{N}(0, 0.1^2)$ and evolved following (13) with $\eta = 0.1$ for 2, 500 iterations. We repeated this for the following parameter values:

- λ (MFL): 0.0125, 0.025, 0.05, 0.1, 0.2,
- λ (gWOT): 0.000625, 0.00125, 0.0025, 0.005, 0.01,

- ε_{DF} (gWOT): 0.01,
- N (both): 1, 2, 4, 8, 16, 32, 64.

The approximate ground truth is taken to be a simulated dataset with 500 particles per timepoint. That is, we computed $(T^{-1} \sum_{i=1}^T D^2(\mu_{t_i}, R_{t_i}))^{1/2}$ where μ_{t_i} and R_{t_i} are respectively the ground truth and reconstructed marginals at time t_i , and the squared Energy Distance [29] between two measures (α, β) is defined to be

$$D^2(\alpha, \beta) = 2\mathbb{E}_{X \sim \alpha, Y \sim \beta} \|X - Y\| - \mathbb{E}_{X, X' \sim \alpha} \|X - X'\| - \mathbb{E}_{Y, Y' \sim \beta} \|Y - Y'\|. \quad (20)$$

G.2 Additional details for Section 4.2

The ambient space \mathcal{X} is taken as in [1], and dynamics follow (1) with $\tau = 1$ and $\Psi(t, x) = 1.25\|x - x^{(0)}\|_2^2\|x - x^{(1)}\|_2^2 + 10 \sum_{k=3}^{10} x_k^2$, where $x^{(0)} = [1.4, 1.4, 0, \dots, 0]$, $x^{(1)} = [-1.25, -1.25, 0, \dots, 0]$. We prescribe a growth rate $g(t, x) = 10(\tanh(2x_0) + 1)/2$ so that particles grow faster in the region $x_0 > 0$. Particles are started from $X_0 \sim \mathcal{N}(0, 0.1^2)$, and 10 timepoints with 50 particles each were sampled on the interval $t \in [0, 0.5]$. We fit MFL dynamics both with and without the modification for branching described in Section 4.2 with $\lambda = 0.025$ and $\rho = +\infty$ (since the growth rate is known exactly). Other hyperparameters were taken to be the same as in Section 4.1.

H Reprogramming dataset pre-processing details

For each set of subsampled snapshots, expression matrices were first centered and projected into 10 PCA dimensions. Similarly as in [1], a set of scaling factors were computed:

- Pairwise scaling factors: $\sigma_{\text{scale}}^2 = \mathbb{E}_{\hat{\mu}_{t_i}, \hat{\mu}_{t_{i+1}}} \left[\frac{\|X_{t_{i+1}} - X_{t_i}\|_2^2}{2} \right]$.
- Per-timepoint scaling factors: $\eta_{\text{scale}}^2 = \mathbb{E}_{\hat{\mu}_{t_i}, \hat{\mu}_{t_i}} \left[\frac{\|X_{t_i} - Y_{t_i}\|_2^2}{2} \right]$.

The cost function for each transport term in (6) was divided by σ_{scale}^2 so as to be of order one. Similarly, the pairwise squared distances in the data-fitting functional (3) were divided by η_{scale}^2 . The timepoints were mapped to $0 = t_1 \leq \dots \leq t_T = 1$ and the value of τ was chosen such that the effective transport regularization level (τ_i in (6)) was 0.1 for transport over 0.5-day interval. We applied MFL dynamics for $\lambda \in \{0.0125, 0.025, 0.05, 0.1, 0.2\}$ and other parameters as in Section 4.1, and gWOT for $\lambda \in \{0.000625, 0.00125, 0.0025, 0.005, 0.01\}$. Of these parameter values, we found that $\lambda = 0.025$ performed best for MFL in terms of Energy Distance to the full dataset (projected onto the previously calculated principal components) and similarly $\lambda = 0.01$ for gWOT.

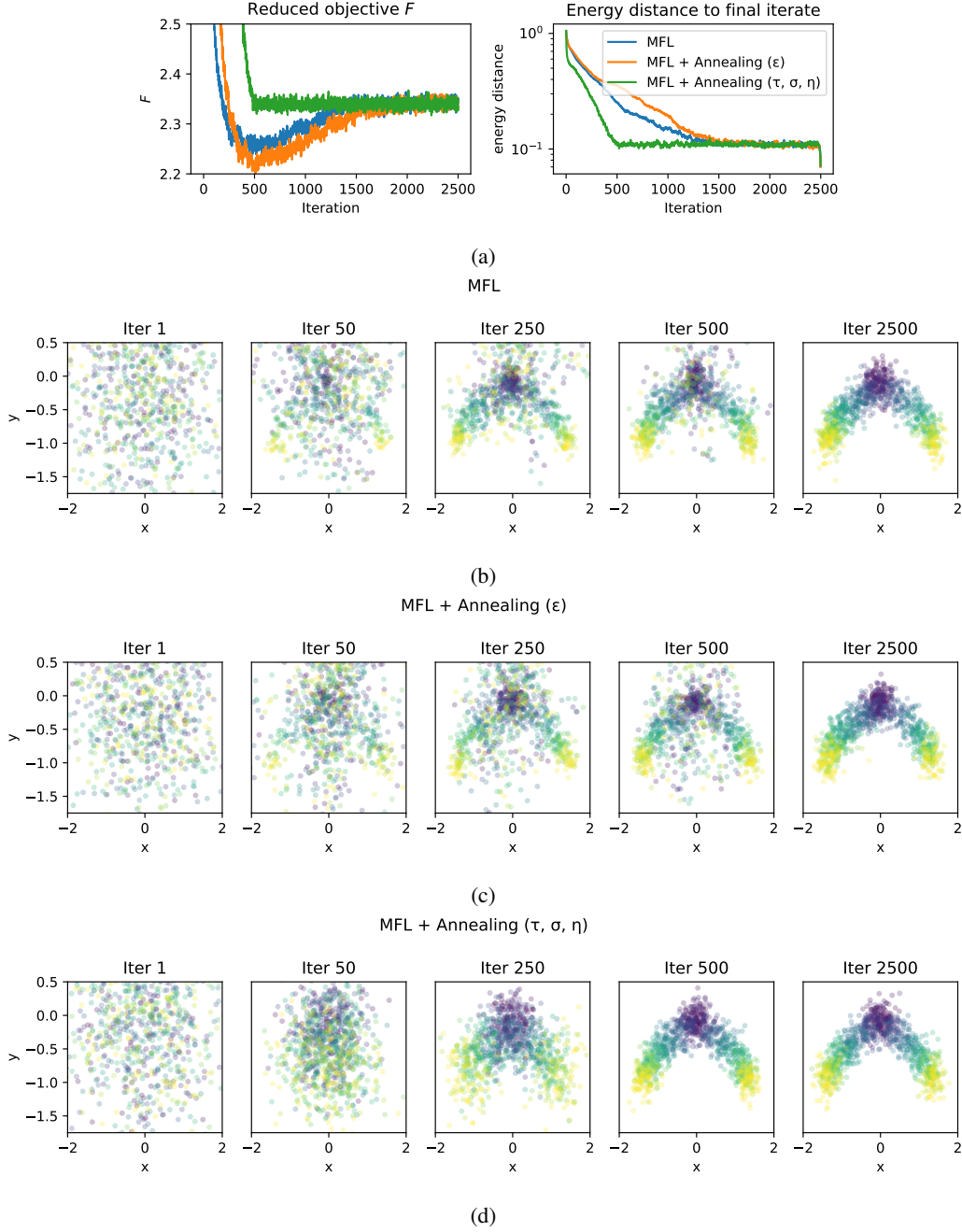


Figure 5: (a) Plot of the (i) reduced objective F (8) (compared to Fig. 2 the y -axis is zoomed in), and (ii) energy distance to final iterate, over 2500 iterations for MFL dynamics without annealing, with annealing in ϵ in the first 1000 iterations, and with annealing jointly in τ, σ, η over the first 500 iterations. (b, c, d) MFL iterates shown in the first 2 dimensions, without annealing, with annealing in ϵ , and with annealing in τ, σ, η .