

# On the Mean-field Analysis of Normalized Steepest Descent via Linear Minimization Oracles

**Yongtao Wu**  
EPFL

YONGTAO.WU@EPFL.CH

**Fanghui Liu**  
Shanghai Jiao Tong University

FANGHUI.LIU@SJTU.EDU.CN

**Taiji Suzuki**  
RIKEN AIP; The University of Tokyo

TAIJI@MIST.I.U-TOKYO.AC.JP

**Volkan Cevher**  
EPFL

VOLKAN.CEVHER@EPFL.CH

## Abstract

In this paper, we develop a mean-field analysis of normalized steepest descent for two-layer neural networks in Wasserstein space. We first study the induced gradient flow in continuous time, establishing global convergence guarantees in both Euclidean and Wasserstein spaces. Our analysis reveals a fundamental distinction between normalized and unnormalized dynamics: while the latter exhibits linear convergence, the normalized flow reaches the optimum in finite time. We further extend the framework to finite-particle and discrete-time settings. We introduce LMO-driven Langevin dynamics and develop adaptive LMO particle schemes, establishing non-asymptotic convergence and stationarity guarantees. Collectively, our results provide a theoretical foundation for normalized steepest descent, a class of optimization methods that has recently become popular for training neural networks.

## 1. Introduction

The growing cost of training large language models (LLMs) has renewed interest in optimization methods beyond standard Euclidean geometry. Adaptive optimizers such as Adam [23] and AdamW [28] rely on coordinate-wise second-moment estimates under the  $\ell_2$  norm, which only partially capture the geometry of large-scale neural network losses [48]. This has motivated a family of normalized steepest descent methods based on non-Euclidean geometries, including spectral-norm-based optimizers such as Muon and its variants [2, 22, 26, 35].

Despite their empirical success, the theoretical foundations of normalized steepest descent remain limited. A useful framework for analyzing optimization dynamics is the mean-field limit of two-layer neural networks [13, 14, 31, 32, 38, 41, 42], where training is described by an evolution over probability measures  $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ :

$$\min_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} \mathcal{F}(\mu) := F(\mu) + \lambda \text{Ent}(\mu), \quad \lambda \geq 0. \quad (1)$$

For  $\lambda > 0$ , the associated Wasserstein gradient flow corresponds to mean-field Langevin dynamics (MFLD), whose convergence properties are well understood in both continuous and discrete settings [11, 12, 21, 30, 32, 43]. See more detail on related work in Sec. B.

However, existing analyses exclusively study *unnormalized* steepest descent under  $\ell_2$  geometry. For normalized steepest descent, it remains unclear how to formulate the corresponding Wasserstein gradient flow and analyze its finite-particle approximation. This raises the following question:

*Can we establish a convergence theory for normalized steepest descent in Wasserstein space?*

In this work, we study two-layer neural networks in the mean-field regime and develop a convergence analysis of normalized steepest descent via linear minimization oracles (LMOs) [17, 25].

**Definition 1 (LMO)** *Let  $\|\cdot\|$  be a norm on  $\mathbb{R}^d$  with dual norm  $\|\cdot\|_*$ . We define:*

$$\text{lmo}(\mathbf{g}; \tau) := \arg \min_{\mathbf{y}: \|\mathbf{y}\| \leq \|\mathbf{g}\|_*} \langle \mathbf{g}, \mathbf{y} \rangle, \quad \forall \mathbf{g} \in \mathbb{R}^d, \tau \in \{0, 1\}.$$

By LMO, we can unify normalized and unnormalized steepest descent (with different radii) under the same framework. For instance, given a function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  and a sequence of step sizes  $\{\eta_k\}_{k \geq 1}$ , the LMO update is:  $\mathbf{x}_{k+1} = \mathbf{x}_k + \eta_k \text{lmo}(\nabla f(\mathbf{x}_k); \tau)$ . It allows for a general optimization framework, e.g.,  $\tau = 0$  corresponding to normalized steepest descent, and  $\tau = 1$  corresponding to (unnormalized) steepest descent; see more examples in Tab. 2.

## 2. Notation

We denote by  $\mathcal{P}_2(\mathbb{R}^d)$  the set of all probability measures on  $\mathbb{R}^d$  with finite second moment. The minimizer of  $\mathcal{F}$  is denoted by  $\mu^*$ . We denote by  $\mathcal{F}^*$  (resp.  $F^*$ ) the optimal value when  $\lambda > 0$  (resp.  $\lambda = 0$ ). For a functional  $F: \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$ , we denote its first variation with respect to  $\mu$  by  $\frac{\delta F}{\delta \mu}(\mu)$ , which is assumed to exist and to be continuously differentiable throughout the paper. We use  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  as the objective function when optimizing in Euclidean space while  $F$  and  $\mathcal{F}$  are the functionals  $\mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$ . We denote by  $\|\cdot\|$  a general primal norm, and by  $\|\cdot\|_*$  its dual norm on  $\mathbb{R}^d$ . We denote the negative Shannon entropy as  $\text{Ent}(\mu) := \int_{\mathbb{R}^d} \log\left(\frac{d\mu}{dx}(x)\right) d\mu(x)$ . For  $\mu \ll \nu$ , let  $H(\mu|\nu)$  denote the Kullback–Leibler divergence. We often identify absolutely continuous probability measures with their densities with respect to the Lebesgue measure for notational simplicity. The 2-Wasserstein distance between  $\mu, \nu \in \mathcal{P}_2(\mathbb{R}^d)$  is defined as:  $W_2(\mu, \nu) := (\inf_{\pi \in \Pi(\mu, \nu)} \int_{\mathbb{R}^d \times \mathbb{R}^d} \|\mathbf{x} - \mathbf{y}\|^2 d\pi(\mathbf{x}, \mathbf{y}))^{1/2}$ , where  $\Pi(\mu, \nu)$  denotes the set of all couplings of  $\mu$  and  $\nu$ . Let  $\mu \in \mathcal{P}_2(\mathbb{R}^d)$  and let  $f: \mathbb{R}^d \rightarrow \mathbb{R}^d$  be measurable, given a norm  $\|\cdot\|$ , we define the space of  $p$ -integrable vector fields as  $L^p(\mu; \|\cdot\|) := \left\{ f: \mathbb{R}^d \rightarrow \mathbb{R}^d \mid \|f\|_{L^p(\mu; \|\cdot\|)} < \infty \right\}$ , where  $\|f\|_{L^p(\mu; \|\cdot\|)} := \left( \int_{\mathbb{R}^d} \|f(x)\|^p d\mu(x) \right)^{1/p}$ ,  $1 \leq p < \infty$ . We denote by  $C^1(\mathbb{R}^d)$  the space of continuously differentiable functions on  $\mathbb{R}^d$ . We use  $\delta_x$  to denote the Dirac measure at  $x$ .

## 3. LMO gradient flow

In this section, we study the analysis of LMO gradient flow to optimization over probability measures based on the objective defined in Eq. (1), which is the extension of our result for optimization over Euclidean space (Thm. 7 in Sec. C of the appendix). We first introduce the following assumptions.

**Assumption 1** *The functional  $F$  is convex and  $\mathcal{F}$  admits a minimizer  $\mu^*$ .*

**Remark 2** *The same assumption has been made in Chizat [12] for analyzing two-layer networks. Proposition 2.5 of Hu et al. [21] guarantees the existence and uniqueness of the minimizer of  $\mathcal{F}$  under mild conditions. Hence Asm 1 can be verified.*

**Assumption 2 (Smoothness in Wasserstein space)** *We assume that the first variation of  $F$  is  $L$ -Lipschitz, i.e., there exists a constant  $L > 0$  such that for all  $\mu, \nu \in \mathcal{P}_2(\mathbb{R}^d)$  and  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$ :*

$$\left\| \nabla \frac{\delta F}{\delta \mu}(\mu)(\mathbf{x}) - \nabla \frac{\delta F}{\delta \mu}(\nu)(\mathbf{x}') \right\|_* \leq \frac{L}{2} (W_2(\mu, \nu) + \|\mathbf{x} - \mathbf{x}'\|),$$

**Remark 3** *Note that Asm 2 with the  $\ell_2$  norm is commonly used in the literature [12, 43]. We extend it to a general norm  $\|\cdot\|$  with the associated dual norm  $\|\cdot\|_*$ . Moreover, for two-layer mean-field neural networks, Asm 2 is satisfied under mild conditions on the neuron and the loss, as shown in Prop. 1 of the appendix.*

In the analysis of MFLD for neural networks, log-Sobolev inequalities play a crucial role in establishing linear convergence rates [12, 32] and it can be verified for two-layer mean-field neural network (Proposition 5.1 in Chizat [12]). In this work, we consider a generalized log-Sobolev assumption formulated with respect to the dual norm of a general norm, i.e.,  $\|\cdot\|_*$  [1, 5]. When  $\|\cdot\|$  is the Euclidean  $\ell_2$  norm, this assumption reduces to the standard log-Sobolev inequality.

**Assumption 3 (Generalized Log-Sobolev inequality)** *Fix a reference measure  $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ . Let  $\nu$  be the Gibbs measure on  $\mathbb{R}^d$  with density  $d\nu(x) \propto \exp\left(-\frac{1}{\lambda} \frac{\delta F(\mu)}{\delta \mu}(x)\right) dx$ . Assume there exists  $\zeta > 0$  such that for every probability measure  $\varphi \ll \nu$ , it holds  $H(\varphi \| \nu) \leq \frac{1}{2\zeta} \int \left\| \nabla \log \frac{d\varphi}{d\nu}(x) \right\|_*^2 d\varphi(x)$ .*

To deliver our analysis, we need the following assumption to avoid ‘‘ill-conditioned’’ measures.

**Assumption 4** *Given probability measures  $\{\mu_t\}$  along the trajectory of interest, there exists a constant  $\kappa \in (0, 1]$  such that, the associated Wasserstein gradient field  $g_{\mu_t}(\mathbf{x}) := \nabla \frac{\delta \mathcal{F}}{\delta \mu}(\mu_t)(\mathbf{x})$  satisfies the uniform lower bound  $\|g_{\mu_t}\|_{L^1(\mu_t; \|\cdot\|_*)} \geq \kappa \|g_{\mu_t}\|_{L^2(\mu_t; \|\cdot\|_*)}$ .*

**Remark:** Asm 4 admits a similar form of reverse Hölder’s inequality [34], see Lemma 1 in Sec. D.2 for further discussion. Asm 4 rules out the case where  $\|g_{\mu}(\mathbf{x})\|_*$  is highly concentrated on a  $\mu$ -negligible region (e.g., spiky gradients). Here we give a simple sufficient condition to Asm 4 as below, with proof deferred to Sec. D.1.

**Condition 1 (Sufficient condition of Asm 4)** *Assume that for all admissible probability measures  $\mu$  along the trajectory, there exists a constant  $G_{\max} < \infty$  such that  $\|g_{\mu}(\mathbf{x})\|_* \leq G_{\max}$  holds  $\mu$ -a.s. (classical bounded gradient assumption), and there exist constants  $m_1 > 0$ ,  $m_2 \in (0, 1]$  such that  $\mu\left(\left\{\mathbf{x} : \|g_{\mu}(\mathbf{x})\|_* \geq m_1\right\}\right) \geq m_2$ . Then Asm 4 holds with  $\kappa = \frac{m_2 m_1}{G_{\max}}$ .*

Note that the PDE associated with the MFLD as presented in Eq. (7) can be written as follows:

$$\partial_t \mu_t = \nabla \cdot (\mu_t \mathbf{v}_t), \quad \mathbf{v}_t := \nabla \frac{\delta F}{\delta \mu}(\mu_t) + \lambda \nabla \log \mu_t. \quad (2)$$

We are now ready to state our main theorem with an analysis on an LMO analogue of the dynamics in Eq. (2). The proof can be found in Sec. D.4.

---

**Algorithm 1: Adaptive LMO Particle Descent**


---

**Input:**  $N, \eta, S_{\text{init}}, \tau, K$ , initial particles  $(\mathbf{x}_0^i)_{i=1}^N \sim \mu_0$ , and  $s \in \{\text{global}, \text{per-particle}\}$   
**Output:**  $\mu_K = \frac{1}{N} \sum_{i=1}^N \delta_{\mathbf{x}_K^i}$   
 $S_{-1} \leftarrow S_{\text{init}}$   
**if**  $s = \text{per-particle}$  **then**  $S_{-1}^i \leftarrow S_{\text{init}}$  for all  $i \in [N]$   
**for**  $k = 0$  **to**  $K - 1$  **do**  
      $\mu_k \leftarrow \frac{1}{N} \sum_{i=1}^N \delta_{\mathbf{x}_k^i}$   
     **for**  $i = 1$  **to**  $N$  **do**  
          $\mathbf{g}_k^i \leftarrow \nabla \frac{\delta F}{\delta \mu}(\mu_k)(\mathbf{x}_k^i)$ ,  $\mathbf{v}_k^i \leftarrow \text{lmo}(\mathbf{g}_k^i; \tau)$   
         **if**  $s = \text{per-particle}$  **then**  $S_k^i \leftarrow S_{k-1}^i + \|\mathbf{v}_k^i\|^2$ ,  $\eta_k^i \leftarrow \eta / \sqrt{S_k^i}$   
     **end**  
     **if**  $s = \text{global}$  **then**  $S_k \leftarrow S_{k-1} + \frac{1}{N} \sum_{i=1}^N \|\mathbf{v}_k^i\|^2$ ,  $\eta_k \leftarrow \eta / \sqrt{S_k}$ ,  $\eta_k^i \leftarrow \eta_k$  for all  $i \in [N]$   
      $\mathbf{x}_{k+1}^i \leftarrow \mathbf{x}_k^i + \eta_k^i \mathbf{v}_k^i$  for all  $i \in [N]$   
**end**

---

**Theorem 4** Under Asm. 1 to 4, consider the following LMO-based continuity equation:

$$\partial_t \mu_t = \nabla \cdot (\mu_t \mathbf{v}_t), \quad \mathbf{v}_t := -\text{lmo} \left( \nabla \frac{\delta F}{\delta \mu}(\mu_t) + \lambda \nabla \log \mu_t; 0 \right). \quad (3)$$

Then for all  $t \geq 0$ , we can obtain:  $\mathcal{F}(\mu_t) - \mathcal{F}(\mu^*) \leq \left[ \max \left\{ \sqrt{\mathcal{F}(\mu_0) - \mathcal{F}(\mu^*)} - \frac{\kappa}{2} \sqrt{2\zeta \lambda} t, 0 \right\} \right]^2$ .

In particular,  $\mathcal{F}(\mu_t) = \mathcal{F}(\mu^*)$  for all  $t \geq t^*$ , with  $t^* = \frac{\sqrt{2}}{\kappa \sqrt{\zeta \lambda}} \sqrt{\mathcal{F}(\mu_0) - \mathcal{F}(\mu^*)}$ .

The result above can be considered as the extension of Thm. 7 (Euclidean space) to Wasserstein space, where finite time convergence is achieved instead of linear convergence.

## 4. Discrete time analysis

The result in the previous section establishes convergence for the continuous-time LMO flow in Wasserstein space with an infinite number of particles. In practice, however, one must work with a finite number of particles and a discrete-time dynamics. In this section, we study the convergence of the algorithm in discrete time with finite particles.

### 4.1. Discrete time analysis: adaptive LMO particle descent

Here, we work in a general non-convex setting without the requirement of Asm 3 and seek a stationarity guarantee. We consider the case without entropy regularization ( $\lambda = 0$ ) for simplicity. A technical issue is that an optimal step size relies on the knowledge of the smoothness constant  $L$ . Motivated by AdaGrad [16], we consider adaptive step sizes that automatically scale with the accumulated gradient, with the following two cases (see Alg. 1): **global** adaptive step size for all particles; **per-particle** adaptive step size for each particle. Below, we provide the analysis for the first choice. The proof and additional result for the second choice can be found in Sec. F.6.

**Theorem 5** Consider Alg. 1 with a global adaptive step size for normalized steepest descent ( $\tau = 0$ ). Let  $\{\mathbf{g}_k^i\}_{k \geq 0, i \in [N]}$ ,  $\{\mu_k\}_{k \geq 0}$  be produced by Alg. 1. Define  $\Delta_k := F(\mu_k) - F^*$ . Under Asm 2:

$$\min_{0 \leq k < K} \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_* \leq \frac{\Delta_0 \sqrt{S_{init}}}{K\eta} + \frac{L\eta}{\sqrt{K}} + \frac{2}{\sqrt{K}\eta} \left( \Delta_0 + \frac{L}{2} \eta^2 \log \left( \frac{S_{init} + K}{S_{init}} \right) \right).$$

#### 4.2. Discrete time analysis: MFLD

Note that in the finite-particle setting, the Kullback–Leibler (KL) term is not well-defined, since the negative entropy is not meaningful for a discrete empirical measure. Therefore, we instead consider the joint distribution of  $N$  particles. Specifically, let  $\mu^{(N)} \in \mathcal{P}^{(N)}$  denote the law of the particle system  $\mathbf{X} = (\mathbf{x}^i)_{i=1}^N$ , where  $\mathcal{P}^{(N)}$  is the space of probability measures on  $(\mathbb{R}^{d \times N}, \mathcal{B}(\mathbb{R}^{d \times N}))$ .

The discrete-time Euler–Maruyama scheme is:

$$\mathbf{x}_{k+1}^i = \mathbf{x}_k^i - \eta_k \mathbf{v}_k^i + \sqrt{2\lambda \eta_k} \boldsymbol{\xi}_k^i, \quad \boldsymbol{\xi}_k^i \sim \mathcal{N}(0, \mathbf{I}_d), \quad (4)$$

with step size  $\eta_k > 0$ . The term  $\mathbf{v}_k^i$  represents the drift direction. In previous works [43],  $\mathbf{v}_k^i$  has been analyzed as the full gradient  $\nabla \frac{\delta F}{\delta \mu}(\mu_k)(\mathbf{x}_k^i)$ , as a stochastic approximation computed on random mini-batches, or as a variance-reduced estimator. Here, for the theoretical analysis in discrete time, we will focus on the following:

$$\mathbf{v}_k^i = -\text{lmo} \left( \nabla \frac{\delta F}{\delta \mu}(\mu_k)(\mathbf{x}_k^i) + \lambda \nabla_i \log \mu_k^{(N)}(\mathbf{X}_k); 0 \right) - \lambda \nabla_i \log \mu_k^{(N)}(\mathbf{X}_k). \quad (5)$$

Similarly, when  $\tau = 1$  with  $\|\cdot\|$  being  $\ell_2$  norm,  $\mathbf{v}_k^i$  reduced to the full gradient  $\nabla \frac{\delta F}{\delta \mu}(\mu_k)(\mathbf{x}_k^i)$ . The motivation for such a choice is that we can rewrite the velocity field in the LMO Wasserstein gradient flow Eq. (3) as follows by adding and subtracting the score term:

$$\mathbf{v}_t = -\text{lmo} \left( \nabla \frac{\delta F}{\delta \mu}(\mu_t) + \lambda \nabla \log \mu_t; 0 \right) - \lambda \nabla \log \mu_t + \lambda \nabla \log \mu_t.$$

Therefore, Eq. (5) can be considered as the finite-particle Euler–Maruyama scheme for Eq. (3), with the last term  $\lambda \nabla \log \mu_t$  corresponding to the noise in Eq. (4). In Sec. E, we present a convergence analysis of LMO noisy particle descent and establish a discrete-time analogue of the finite-time convergence results derived for the continuous flows in Thm. 4 and 7. The precise statement and proof are provided in Thm. 10 and Sec. E.

## 5. Conclusion

This paper establishes convergence guarantees for two-layer neural networks using LMO Wasserstein gradient flows. This paper also provides a convergence rate for a practical finite-particle scheme as well as adaptive particle methods in a discrete time setting. Interesting open directions include extending to deep neural networks, weakening the uniform log-Sobolev and the non-highly concentrated gradient assumptions, and extending the theory to mirror update. Furthermore, developing a convergence analysis for adaptive particle-based LMO dynamics under both smoothness and log-Sobolev assumptions remains largely unexplored, particularly for entropy-regularized objectives.

## References

- [1] Radosław Adamczak, Witold Bednorz, and Paweł Wolff. Moment estimates implied by modified log-sobolev inequalities. *ESAIM: Probability and Statistics*, 21:467–494, 2017.
- [2] Kwangjun Ahn, Byron Xu, Natalie Abreu, Ying Fan, Gagik Magakyan, Pratyusha Sharma, Zheng Zhan, and John Langford. Dion: Distributed orthonormalized updates. *arXiv preprint arXiv:2504.05295*, 2025.
- [3] Amit Attia and Tomer Koren. Sgd with adagrad stepsizes: Full adaptivity with high probability to unknown parameters, unbounded gradients and affine variance. In *International Conference on Machine Learning*, pages 1147–1171. PMLR, 2023.
- [4] Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar. signsgd: Compressed optimisation for non-convex problems. In *International conference on machine learning*, pages 560–569. PMLR, 2018.
- [5] Sergey G Bobkov and Michel Ledoux. From brunn-minkowski to brascamp-lieb and to logarithmic sobolev inequalities. *Geometric and Functional Analysis*, 10(5):1028–1052, 2000.
- [6] Clément Bonet, Théo Uscidda, Adam David, Pierre-Cyril Aubin-Frankowski, and Anna Korba. Mirror and preconditioned gradient descent in wasserstein space. *Advances in Neural Information Processing Systems*, 37:25311–25374, 2024.
- [7] Siwan Boufadène and François-Xavier Vialard. On the global convergence of wasserstein gradient flow of the coulomb discrepancy. *SIAM Journal on Mathematical Analysis*, 57(4):4556–4587, 2025.
- [8] David Carlson, Volkan Cevher, and Lawrence Carin. Stochastic Spectral Descent for Restricted Boltzmann Machines. In *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics*, 2015.
- [9] David Carlson, Ya-Ping Hsieh, Edo Collins, Lawrence Carin, and Volkan Cevher. Stochastic spectral descent for discrete graphical models. *IEEE Journal of Selected Topics in Signal Processing*, 10(2):296–311, 2015.
- [10] David E Carlson, Edo Collins, Ya-Ping Hsieh, Lawrence Carin, and Volkan Cevher. Preconditioned spectral descent for deep learning. *Advances in neural information processing systems*, 28, 2015.
- [11] Fan Chen, Zhenjie Ren, and Songbo Wang. Uniform-in-time propagation of chaos for mean field langevin dynamics. In *Annales de l’Institut Henri Poincaré (B) Probabilités et statistiques*, volume 61, pages 2357–2404. Institut Henri Poincaré, 2025.
- [12] Lénaïc Chizat. Mean-field langevin dynamics : Exponential convergence and annealing. *Transactions on Machine Learning Research*, 2022. ISSN 2835-8856.
- [13] Lenaïc Chizat. Sparse optimization on measures with over-parameterized gradient descent. *Mathematical Programming*, 194(1):487–532, 2022.

- [14] Lenaïc Chizat and Francis Bach. On the global convergence of gradient descent for over-parameterized models using optimal transport. *Advances in neural information processing systems*, 31, 2018.
- [15] Lénaïc Chizat, Maria Colombo, and Xavier Fernández-Real. Convergence of drift-diffusion pdes arising as wasserstein gradient flows of convex functions. *arXiv preprint arXiv:2507.12385*, 2025.
- [16] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
- [17] Marguerite Frank and Philip Wolfe. An algorithm for quadratic programming. *Naval research logistics quarterly*, 3(1-2):95–110, 1956.
- [18] Kaja Gruntkowska, Alexander Gaponov, Zhirayr Tovmasyan, and Peter Richtárik. Error feedback for muon and friends. *arXiv preprint arXiv:2510.00643*, 2025.
- [19] Anming Gu and Juno Kim. Mirror mean-field langevin dynamics. *arXiv preprint arXiv:2505.02621*, 2025.
- [20] Elad Hazan, Kfir Levy, and Shai Shalev-Shwartz. Beyond convexity: Stochastic quasi-convex optimization. *Advances in neural information processing systems*, 28, 2015.
- [21] Kaitong Hu, Zhenjie Ren, David Šiška, and Łukasz Szpruch. Mean-field langevin dynamics and energy landscape of neural networks. In *Annales de l’Institut Henri Poincaré (B) Probabilités et statistiques*, pages 2043–2065. Institut Henri Poincaré, 2021.
- [22] Keller Jordan, Yuchen Jin, Vlado Boza, Jiacheng You, Franz Cesista, Laker Newhouse, and Jeremy Bernstein. Muon: An optimizer for hidden layers in neural networks, 2024. URL <https://kellerjordan.github.io/posts/muon/>.
- [23] Diederik P Kingma. Adam: A method for stochastic optimization. *International Conference on Learning Representations*, 2015.
- [24] Marcel Kreuter. Sobolev spaces of vector-valued functions. *Ulm University Faculty of Mathematics and Economics*, 2015.
- [25] Simon Lacoste-Julien and Martin Jaggi. On the global linear convergence of frank-wolfe optimization variants. *Advances in neural information processing systems*, 28, 2015.
- [26] Zichong Li, Liming Liu, Chen Liang, Weizhu Chen, and Tuo Zhao. Normuon: Making muon more efficient and scalable. *arXiv preprint arXiv:2510.05491*, 2025.
- [27] Fanghui Liu, Luca Viano, and Volkan Cevher. Understanding deep neural function approximation in reinforcement learning via epsilon-greedy exploration. *Advances in Neural Information Processing Systems*, 35:5093–5108, 2022.
- [28] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.

- [29] Yurii Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media, 2013.
- [30] Atsushi Nitanda. Improved particle approximation error for mean field neural networks. *Advances in Neural Information Processing Systems*, 37:113823–113845, 2024.
- [31] Atsushi Nitanda and Taiji Suzuki. Stochastic particle gradient descent for infinite ensembles. *arXiv preprint arXiv:1712.05438*, 2017.
- [32] Atsushi Nitanda, Denny Wu, and Taiji Suzuki. Convex analysis of the mean field langevin dynamics. In *International Conference on Artificial Intelligence and Statistics*, pages 9741–9757. PMLR, 2022.
- [33] Atsushi Nitanda, Anzelle Lee, Damian Tan Xing Kai, Mizuki Sakaguchi, and Taiji Suzuki. Propagation of chaos for mean-field langevin dynamics and its application to model ensemble. In *Forty-second International Conference on Machine Learning*, 2025. URL <https://openreview.net/forum?id=SjT6JK4KZv>.
- [34] Jorge Alberto Paz Moyado, Yamilet del Carmen Quintana Mato, José Manuel Rodríguez García, and José María Sigarreta Almira. New reverse hölder-type inequalities and applications. 2023.
- [35] Thomas Pethick, Wanyun Xie, Kimon Antonakopoulos, Zhenyu Zhu, Antonio Silveti-Falls, and Volkan Cevher. Training deep learning models with norm-constrained LMOs. In *Forty-second International Conference on Machine Learning*, 2025.
- [36] Boris T Polyak. Gradient methods for the minimisation of functionals. *USSR Computational Mathematics and Mathematical Physics*, 3(4):864–878, 1963.
- [37] Grant Rotskoff, Samy Jelassi, Joan Bruna, and Eric Vanden-Eijnden. Global convergence of neuron birth-death dynamics. In *International Conference on Machine Learning*, 2019.
- [38] Grant M Rotskoff and Eric Vanden-Eijnden. Neural networks as interacting particle systems: Asymptotic convexity of the loss landscape and universal scaling of the approximation error. *stat*, 1050:22, 2018.
- [39] Filippo Santambrogio. {Euclidean, metric, and Wasserstein} gradient flows: an overview. *Bulletin of Mathematical Sciences*, 7(1):87–154, 2017.
- [40] Louis Sharrock and Christopher Nemeth. Tuning-free sampling via optimization on the space of probability measures. *arXiv preprint arXiv:2510.25315*, 2025.
- [41] Justin Sirignano and Konstantinos Spiliopoulos. Mean field analysis of neural networks: A law of large numbers. *SIAM Journal on Applied Mathematics*, 80(2):725–752, 2020.
- [42] Mei Song, Andrea Montanari, and P Nguyen. A mean field view of the landscape of two-layers neural networks. *Proceedings of the National Academy of Sciences*, 115(33):E7665–E7671, 2018.
- [43] Taiji Suzuki, Denny Wu, and Atsushi Nitanda. Convergence of mean-field langevin dynamics: time-space discretization, stochastic gradient, and variance reduction. *Advances in Neural Information Processing Systems*, 36:15545–15577, 2023.

- [44] Alain-Sol Sznitman. Topics in propagation of chaos. In *Ecole d'été de probabilités de Saint-Flour XIX—1989*, pages 165–251. Springer, 2006.
- [45] Gerald Teschl. *Ordinary differential equations and dynamical systems*, volume 140. American Mathematical Soc., 2012.
- [46] Jan Vondrák. Submodularity and curvature: The optimal algorithm (combinatorial optimization and discrete algorithms). *23:253–266*, 2010.
- [47] Bohan Wang, Huishuai Zhang, Zhiming Ma, and Wei Chen. Convergence of adagrad for non-convex objectives: Simple proofs and relaxed assumptions. In *The Thirty Sixth Annual Conference on Learning Theory*, pages 161–190. PMLR, 2023.
- [48] Yushun Zhang, Congliang Chen, Tian Ding, Ziniu Li, Ruoyu Sun, and Zhiqian Luo. Why transformers need adam: A hessian perspective. *Advances in neural information processing systems*, 37:131786–131823, 2024.
- [49] Shuailong Zhu and Xiaohui Chen. Convergence analysis of the wasserstein proximal algorithm beyond geodesic convexity. *arXiv preprint arXiv:2501.14993*, 2025.

## Acknowledgments

We thank the workshop committee and reviewers for their work. This work was supported under project ID # 37 as part of the Swiss AI Initiative, through a grant from the ETH Domain and computational resources provided by the Swiss National Supercomputing Centre (CSCS) under the Alps infrastructure. This work was funded by the Swiss National Science Foundation (SNSF) under grant number 2000-1-240094. Research was sponsored by the Army Research Office and was accomplished under Grant Number W911NF-24-1-0048. TS was partially supported by JSPS KAKENHI (24K02905) and JST CREST (JPMJCR2015). This research is supported by the National Research Foundation, Singapore and the Ministry of Digital Development and Information under the AI Visiting Professorship Programme (award number AIVP-2024-004). Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not reflect the views of National Research Foundation, Singapore and the Ministry of Digital Development and Information.

## Contents of the Appendix

The Appendix is organized as follows:

- In Sec. [A](#), we summarize the main symbols and notation used throughout the paper.
- In Sec. [B](#), we provide further detail on related work.
- In Sec. [C](#), we provide theorems and proofs on LMO gradient flow in Euclidean space.
- In Sec. [D](#) and [E](#), we provide the theoretical analysis and the detailed proofs for the LMO-based mean-field Langevin dynamics, including the proof for continuous time in Sec. [D](#), and discrete setting in Sec. [E](#).
- In Sec. [F](#), we present the full technical proofs for our results on the optimization in measure space in the noiseless setting ( $\lambda = 0$ ) and an additional theorem on adaptive LMO descent, and the corresponding proof.
- The intuition on the adaptive step size design is given at Sec. [G](#).
- In Sec. [H](#), we study the annealed dynamics in the mean-field Langevin dynamics setting.
- We provide additional experiments with varying dimension, and comparison between global adaptive stepsize and per-particle adaptive stepsize in Sec. [I](#).
- We discuss the broader impact of this work in Sec. [J](#).

## Appendix A. Notation

A detailed summary of the notation can be found at Tab. [1](#).

The definition of the first variation of a functional  $F : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$  is as follows:

**Definition 6** For a functional  $F : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$ , its first variation  $\frac{\delta F}{\delta \mu}(\mu)(\mathbf{x})$  is defined by

$$\lim_{\varepsilon \rightarrow 0} \frac{F((1 - \varepsilon)\mu + \varepsilon\nu) - F(\mu)}{\varepsilon} = \int \frac{\delta F}{\delta \mu}(\mu)(\mathbf{x}) d(\nu - \mu)(\mathbf{x}),$$

for all  $\nu \in \mathcal{P}_2(\mathbb{R}^d)$ .

## Appendix B. Further detail on related work

### B.1. Steepest descent

Steepest descent was originally studied for unconstrained smooth optimization problems in  $\mathbb{R}^d$ . Early work by Carlson et al. [[8](#), [9](#), [10](#)] showed that the smoothness assumption of the objective function implies a local majorization bound that naturally induces a notion of steepest descent in possibly non-Euclidean geometries. This perspective yields unnormalized steepest descent directions. Spectral descent arises as a special case when the underlying norm is the spectral norm. Building on this idea, Jordan et al. [[22](#)] proposed Muon, which applies a normalized steepest descent update by orthogonalizing the momentum under a spectral norm constraint. More recently, Pethick et al. [[35](#)]

introduced Scion, providing a general framework for training neural networks using norm-constrained LMOs, including guidelines for norm choices, proper scalings across layers, and convergence guarantees. Several further modifications and extensions of Muon have since been proposed [2, 18]. On the other hand, Vondrák [46] analyzed a continuous greedy algorithm based on linear optimization over a matroid polytope, whereas we study normalized descent dynamics induced by norm-based oracles in Euclidean and Wasserstein spaces.

## B.2. Wasserstein gradient flows & MFLD

A growing body of work considers the training of overparameterized models at the level of probability measures, by viewing infinite-width two layer neural networks as minimizing a functional  $F : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$  in Wasserstein space, see Eq. (1). [14, 31, 42]. A neuron with parameter  $\mathbf{x}$  produces output  $h(\mathbf{x}, \mathbf{z})$  for input  $\mathbf{z}$ , and an infinitely wide network is represented by a probability distribution  $\mu$  over parameters, with predicted output  $f_\mu(\mathbf{z}) = \int h(\mathbf{x}, \mathbf{z}) d\mu(\mathbf{x})$ . Given a data distribution  $\mathcal{D}$  and a loss function  $\ell$ , the energy functional is defined as the population risk  $F(\mu) = \mathbb{E}_{(\mathbf{z}, y) \sim \mathcal{D}}[\ell(\int h(\mathbf{x}, \mathbf{z}) d\mu(\mathbf{x}), y)]$ . A finite-width network with parameters  $(\mathbf{x}^i)_{i=1}^N$  corresponds to the empirical measure  $\mu_N = \frac{1}{N} \sum_{i=1}^N \delta_{\mathbf{x}^i}$ , in which case the predicted output becomes  $f_{\mu_N}(\mathbf{z}) = \frac{1}{N} \sum_{i=1}^N h(\mathbf{x}^i, \mathbf{z})$ . Thus,  $\mu_N$  serves as a finite-particle approximation of  $\mu$ .

Previous works show that gradient descent on the finite network converges to a Wasserstein gradient flow in the mean-field limit, and use this perspective to obtain global convergence guarantees for mean-field neural networks [14, 37, 42]. Quantitative analysis with particle gradient descent is given in Chizat [13]. The extension to the mirror descent algorithm beyond gradient descent is presented in Bonet et al. [6]. In the entropy-regularized setting, the associated Wasserstein gradient flow becomes the MFLD. To minimize  $\mathcal{F}$ , a common approach is to consider the MFLD:

$$d\mathbf{x}_t = -\nabla \frac{\delta F}{\delta \mu}(\mu_t)(\mathbf{x}_t) dt + \sqrt{2\lambda} d\mathbf{w}_t, \mu_t = \text{Law}(\mathbf{x}_t), \quad (6)$$

where  $\mathbf{w}_t$  is a standard Brownian motion in  $\mathbb{R}^d$ . The law of a solution to MFLD solves the following PDE:

$$\partial_t \mu_t = \nabla \cdot \left( \mu_t \nabla \frac{\delta F}{\delta \mu}(\mu_t) \right) + \lambda \Delta \mu_t. \quad (7)$$

The convergence of MFLD to the optimum is proved in Hu et al. [21], Song et al. [42]. Later, under transport smoothness and uniform log-Sobolev conditions, one can obtain linear convergence of  $\mu_t$  toward the minimizer  $\mu^*$  of  $\mathcal{F}$  [12, 32].

To relate this mean-field limit to a practical finite model, one considers noisy particle gradient descent. For a finite particle system, given  $\mathbf{X} = (\mathbf{x}^i)_{i=1}^N \in \mathbb{R}^{Nd}$ , denote by its associated empirical measure as  $\mu_{\mathbf{X}} = \frac{1}{N} \sum_{i=1}^N \delta_{\mathbf{x}^i}$ , then a single iteration is:

$$\mathbf{x}_{k+1}^i = \mathbf{x}_k^i - \eta \nabla \frac{\delta F}{\delta \mu}(\mu_{\mathbf{X}_k})(\mathbf{x}_k^i) + \sqrt{2\lambda\eta} \boldsymbol{\xi}_k^i, i \in [N], \quad (8)$$

where  $\boldsymbol{\xi}_k^i \sim \mathcal{N}(0, I_d)$  are i.i.d. normal random variables. Passing to the limit  $\eta \rightarrow 0$  yields the stochastic differential equation:

$$d\mathbf{x}_t^i = -\nabla \frac{\delta F}{\delta \mu}(\mu_{\mathbf{X}_t})(\mathbf{x}_t^i) dt + \sqrt{2\lambda} d\mathbf{w}_t^i, i \in [N]. \quad (9)$$

When the number of particles grows to infinity, the empirical distribution  $\mu_{\mathbf{X}_t}$  converges to a deterministic measure  $\mu_t$  satisfying the MFLD in Eq. (6). In this case, the trajectories of particles are given by i.i.d. samples from Eq. (6), referred to as propagation of chaos [44].

Recent works establish propagation-of-chaos guarantees for mean-field Langevin dynamics, showing that the finite-particle system converges to the ideal mean-field limit with quantitative error bounds [11, 43]. Refined analysis is given by Nitanda [30], Nitanda et al. [33], which removes the dependence on the log-Sobolev inequality constants and shows that the objectives scale as  $\mathcal{O}(1/N)$  uniformly in time. Gu and Kim [19] further extends the idea of MFLD to mirror descent dynamics.

However, all previous results rely on the unnormalized Wasserstein gradient-flow structure, which does not cover the recent normalized steepest descent dynamics, so the existing convergence and propagation-of-chaos theories do not directly apply. In addition, Sharrock and Nemeth [40] studies adaptive step size for Wasserstein gradient descent in the geodesically convex setting, whereas our focus is on normalized and unnormalized steepest descent dynamics for nonconvex objectives.

Note that Eq. (9) has an equivalent vector form:

$$d\mathbf{X}_t = -N\nabla F(\mu_{\mathbf{X}_t})dt + \sqrt{2\lambda}d\mathbf{W}_t,$$

where  $\mathbf{W}_t$  is a standard Brownian motion in  $\mathbb{R}^{Nd}$ . Classical results imply that  $\mu_t^{(N)} = \text{Law}(\mathbf{X}_t)$  converges to the invariant Gibbs measure  $\frac{d\mu_*^{(N)}}{d\mathbf{X}}(\mathbf{X}) \propto \exp(-\frac{N}{\lambda}F(\mu_{\mathbf{X}}))$ . This stationary distribution is the unique minimizer of the finite-particle free energy

$$\mathcal{F}^{(N)}(\mu^{(N)}) = N \mathbb{E}_{\mathbf{X} \sim \mu^{(N)}}[F(\mu_{\mathbf{X}})] + \lambda \text{Ent}(\mu^{(N)}),$$

where  $\mu^{(N)} \in \mathcal{P}_2(\mathbb{R}^{dN})$ . A fundamental question is how well  $\mu_*^{(N)}$  approximates the ideal law  $(\mu^*)^{\otimes N}$ . Recent propagation-of-chaos results for MFLD [11, 43] prove

$$\frac{1}{N} \left( \mathcal{F}^{(N)}(\mu_*^{(N)}) - \mathcal{F}(\mu^*) \right) = \mathcal{O}\left(\frac{\lambda}{\zeta N}\right),$$

where  $\zeta$  is the constant in the log-Sobolev inequality. The resulting bounds depend inversely on log-Sobolev constants, which can deteriorate exponentially in the regularization parameter, leading to pessimistic particle complexity estimates. To mitigate this issue, Nitanda [30], Nitanda et al. [33] develop refined propagation-of-chaos and particle-bias bounds for mean-field neural networks, showing that the particle approximation error in the objective scales as  $\mathcal{O}(1/N)$  uniformly in time and removing the exponential sensitivity to log-Sobolev constants.

## Appendix C. Analysis of LMO gradient flow in Euclidean space

### C.1. LMO gradient flow in Euclidean space

We consider the unconstrained optimization problem  $\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$ , where  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  satisfies  $f \in C^1(\mathbb{R}^d)$ . Next, we state the following assumption, which is a generalized version of the strongly convex assumption with  $\ell_2$  norm in optimization literature [29].

**Assumption 5** *We assume the objective  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is  $\alpha$ -strongly convex w.r.t. the norm  $\|\cdot\|$ , i.e., there exists  $\alpha > 0$  such that:*

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\alpha}{2} \|\mathbf{y} - \mathbf{x}\|^2, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d.$$

Note that the norm  $\|\cdot\|$  can be  $\ell_p$  ( $1 \leq p \leq \infty$ ) or, more generally, any norm of interest.

Below, we define the gradient flow with LMO, which generalizes the classical gradient flow by replacing the direction  $-\nabla f(\mathbf{x}_t)$  with an LMO-induced direction.

$$d\mathbf{x}_t = \text{lmo}(\nabla f(\mathbf{x}_t); \tau) dt, \quad (10)$$

where  $\mathbf{x}_0 \in \mathbb{R}^d$  is the initialization of the optimization trajectory. In the case of using  $\ell_2$  norm with  $\tau = 1$  in LMO, classical result [39] shows that under Asm 5, one gets:  $f(\mathbf{x}_t) - f^* = (f(\mathbf{x}_0) - f^*)e^{-2\alpha t}$ . As a comparison, we provide the following theorem in the case of  $\tau = 0$  (i.e., normalized steepest descent).

**Theorem 7 (Convergence of normalized LMO gradient flow)** *Consider the gradient flow in Eq. (10) with any norm and  $\tau = 0$  in LMO, under Asm 5 and assume the gradient flow solution exists, then we have the following convergence result:*

$$f(\mathbf{x}_t) - f^* \leq \left[ \max\left\{ \sqrt{f(\mathbf{x}_0) - f^*} - \sqrt{\frac{\alpha}{2}t}, 0 \right\} \right]^2.$$

The proof can be found in Sec. C.2. This shows that the trajectory reaches the optimum in finite time,  $f(\mathbf{x}_t) = f^*$ ,  $\forall t \geq t^* := \sqrt{\frac{2}{\alpha}(f(\mathbf{x}_0) - f^*)}$ . Note that the assumption on the existence of the gradient flow solution can be verified by checking the local Lipschitz continuity of the corresponding LMO vector field and then applying the Picard–Lindelöf theorem [45]. Refer to Sec. C.3 for a concrete proof for the normalized gradient flow case.

## C.2. Proof of Thm. 7

**Proof** [Proof of Thm. 7] First, we prove that strong convexity under a general norm  $\|\cdot\|$  can still get a PL inequality: Strong convexity at  $\mathbf{x}$  with  $\mathbf{y} = \mathbf{x}^*$  gives

$$f^* \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}^* - \mathbf{x} \rangle + \frac{\alpha}{2} \|\mathbf{x}^* - \mathbf{x}\|^2.$$

Rearranging, using the Cauchy–Schwarz inequality for dual norms, and Young’s inequality:

$$f(\mathbf{x}) - f^* \leq \|\nabla f(\mathbf{x})\|_* \|\mathbf{x} - \mathbf{x}^*\| - \frac{\alpha}{2} \|\mathbf{x} - \mathbf{x}^*\|^2 \leq \frac{1}{2\alpha} \|\nabla f(\mathbf{x})\|_*^2 + \frac{\alpha}{2} \|\mathbf{x} - \mathbf{x}^*\|^2 - \frac{\alpha}{2} \|\mathbf{x} - \mathbf{x}^*\|^2.$$

As a result, we have:

$$f(\mathbf{x}) - f^* \leq \frac{1}{2\alpha} \|\nabla f(\mathbf{x})\|_*^2. \quad (11)$$

Next, we will prove the convergence. Along any differentiable trajectory  $\mathbf{x}(t)$ ,

$$\frac{d}{dt} f(\mathbf{x}(t)) = \langle \nabla f(\mathbf{x}(t)), \dot{\mathbf{x}}(t) \rangle = \langle \nabla f(\mathbf{x}(t)), \text{lmo}(\nabla f(\mathbf{x}(t)), 0) \rangle = -\|\nabla f(\mathbf{x}(t))\|_*. \quad (12)$$

where we use the equality  $\langle \mathbf{g}, \text{lmo}(\mathbf{g}, 0) \rangle = \min_{\|\mathbf{y}\| \leq 1} \langle \mathbf{g}, \mathbf{y} \rangle = -\|\mathbf{g}\|_*$ .

Define  $G(t) := \sqrt{f(\mathbf{x}(t)) - f^*} \geq 0$ . Combining Eq. (11) and Eq. (12) gives

$$\frac{d}{dt} f(\mathbf{x}(t)) = \frac{d}{dt} G(t)^2 = 2G(t) \frac{d}{dt} G(t) \leq -\sqrt{2\alpha} G(t).$$

Therefore, when  $G(t) > 0$ , we can divide both sides by  $2G(t)$  and obtain  $\frac{d}{dt}G(t) \leq -\sqrt{\frac{\alpha}{2}}$ . Integrating both sides gives

$$G(t) \leq G(0) - \sqrt{\frac{\alpha}{2}}t.$$

Hence  $G(t)$  decreases linearly until it possibly reaches 0. Once  $G(t) = 0$ , i.e.,  $f(\mathbf{x}(t)) = f^*$ , we have  $\nabla f(\mathbf{x}(t)) = 0$ , and from Eq. (12).

$$\frac{d}{dt}f(\mathbf{x}(t)) = -\|\nabla f(\mathbf{x}(t))\|_* = 0,$$

so  $G(t)$  remains constant at 0 thereafter.

Therefore, we get

$$f(\mathbf{x}(t)) - f^* \leq \left[ \max \left\{ \sqrt{f(\mathbf{x}(0)) - f^*} - \sqrt{\frac{\alpha}{2}}t, 0 \right\} \right]^2.$$

■

### C.3. Uniqueness and existence of normalized gradient flow

**Theorem 8 (Uniqueness and existence of normalized gradient flow)** *Under Asm 5, consider setting  $\tau = 0$  with  $\ell_2$  norm, the LMO gradient flow reduces to the following normalized gradient flow*

$$\dot{\mathbf{x}}(t) = -\frac{\nabla f(\mathbf{x}(t))}{\|\nabla f(\mathbf{x}(t))\|_2}, \quad \mathbf{x}(0) = \mathbf{x}_0, \quad (13)$$

with the convention that  $\dot{\mathbf{x}}(t) = 0$  whenever  $\nabla f(\mathbf{x}(t)) = 0$ . And we assume  $\nabla f(\mathbf{x}_0) \neq 0$  for a nontrivial trajectory. Then we have i) *Local existence and uniqueness. There exists  $T > 0$  and a unique trajectory  $\mathbf{x} : [0, T] \rightarrow \mathbb{R}^d$  satisfying Eq. (13) with  $\mathbf{x}(0) = \mathbf{x}_0$ . Moreover, for sufficiently small  $T$ , we have  $\nabla f(\mathbf{x}(t)) \neq 0$  for all  $t \in [0, T]$ , so uniqueness holds on  $[0, T]$ .*

ii) *Global existence. There exists a global absolutely continuous solution  $x : [0, \infty) \rightarrow \mathbb{R}^d$  to Eq. (13), and we have  $\|\mathbf{x}(t) - \mathbf{x}_0\| \leq t, \quad \forall t \geq 0$ . Moreover, if a trajectory reaches  $\mathbf{x}^*$  with  $\nabla f(\mathbf{x}^*) = 0$ , then multiple continuations of the solution may exist beyond that time.*

**Proof** First we show the local Lipschitzness. Define  $\phi : \mathbb{R}^d \setminus \{0\} \rightarrow \mathbb{R}^d$  by  $\phi(\mathbf{g}) := \frac{\mathbf{g}}{\|\mathbf{g}\|_2}$ . We show that if  $\|\mathbf{g}\|_2 \geq \varepsilon > 0$  and  $\|\mathbf{h}\|_2 \geq \varepsilon > 0$ , then

$$\|\phi(\mathbf{g}) - \phi(\mathbf{h})\|_2 \leq \frac{2}{\varepsilon} \|\mathbf{g} - \mathbf{h}\|_2. \quad (14)$$

Indeed, we have

$$\begin{aligned} \|\phi(\mathbf{g}) - \phi(\mathbf{h})\|_2 &= \left\| \frac{\mathbf{g}}{\|\mathbf{g}\|_2} - \frac{\mathbf{h}}{\|\mathbf{h}\|_2} \right\|_2 \leq \left\| \frac{\mathbf{g}}{\|\mathbf{g}\|_2} - \frac{\mathbf{g}}{\|\mathbf{h}\|_2} \right\|_2 + \left\| \frac{\mathbf{g}}{\|\mathbf{h}\|_2} - \frac{\mathbf{h}}{\|\mathbf{h}\|_2} \right\|_2 \\ &= \frac{|\|\mathbf{h}\|_2 - \|\mathbf{g}\|_2|}{\|\mathbf{h}\|_2} + \frac{\|\mathbf{g} - \mathbf{h}\|_2}{\|\mathbf{h}\|_2} \leq \frac{\|\mathbf{g} - \mathbf{h}\|_2}{\|\mathbf{h}\|_2} + \frac{\|\mathbf{g} - \mathbf{h}\|_2}{\|\mathbf{h}\|_2} \leq \frac{2}{\varepsilon} \|\mathbf{g} - \mathbf{h}\|_2. \end{aligned} \quad (15)$$

Fix  $\hat{\mathbf{x}}$  with  $\nabla F(\hat{\mathbf{x}}) \neq 0$  and set  $\varepsilon := \frac{1}{2}\|\nabla F(\hat{\mathbf{x}})\|_2 > 0$ . By continuity of  $\nabla F$ , there exists  $r > 0$  such that

$$\|\nabla f(\mathbf{x})\|_2 \geq \varepsilon, \quad \forall \mathbf{x} \in B(\hat{\mathbf{x}}, r). \quad (16)$$

For any  $\mathbf{x}, \mathbf{y} \in B(\hat{\mathbf{x}}, r)$ , apply Eq. (14) with  $g = \nabla f(\mathbf{x})$  and  $h = \nabla f(\mathbf{y})$ :

$$\left\| \frac{\nabla f(\mathbf{x})}{\|\nabla f(\mathbf{x})\|} - \frac{\nabla f(\mathbf{y})}{\|\nabla f(\mathbf{y})\|} \right\|_2 \leq \frac{2}{\varepsilon} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq \frac{2L}{\varepsilon} \|\mathbf{x} - \mathbf{y}\|_2,$$

where the last inequality uses the global Lipschitz property of  $\nabla F$  based on Asm 5.A. This proves that if  $\nabla F(\hat{\mathbf{x}}) \neq 0$ , then there exists a neighborhood  $U$  of  $\hat{\mathbf{x}}$  and a constant  $L_U < \infty$  such that

$$\left\| \frac{\nabla f(\mathbf{x})}{\|\nabla f(\mathbf{x})\|_2} - \frac{\nabla f(\mathbf{y})}{\|\nabla f(\mathbf{y})\|_2} \right\| \leq L_U \|\mathbf{x} - \mathbf{y}\|_2, \quad \forall \mathbf{x}, \mathbf{y} \in U. \quad (17)$$

Next, we are ready to show the local existence and uniqueness. Given  $\nabla f(\mathbf{x}_0) \neq 0$ , in its neighborhood, we have shown local Lipschitz property. By the Picard–Lindelöf theorem [45], there exists  $T > 0$  and a unique solution  $\mathbf{x} : [0, T] \rightarrow \mathbb{R}^d$  to Eq. (13) with  $\mathbf{x}(0) = \mathbf{x}_0$ . For sufficiently small  $T$ , the trajectory remains in the region where Eq. (16) holds, so  $\nabla f(\mathbf{x}(t)) \neq 0$  for all  $t \in [0, T]$ .

Next, clearly, we have  $\left\| \frac{\nabla f(\mathbf{x}(t))}{\|\nabla f(\mathbf{x}(t))\|_2} \right\|_2 = 1$ . Hence, any absolutely continuous solution satisfies

$$\|\mathbf{x}(t) - \mathbf{x}(0)\|_2 \leq \int_0^t \left\| \frac{\nabla f(\mathbf{x}(s))}{\|\nabla f(\mathbf{x}(s))\|_2} \right\|_2 ds = t,$$

Because trajectories cannot escape to infinity in finite time, standard continuation arguments imply that any local solution can be extended globally for all  $t \geq 0$ .

Lastly, we discuss the possible non-uniqueness at stationary points. If  $\nabla f(\mathbf{x}^*) = 0$ , then by convention we set  $\dot{\mathbf{x}}(t) = \mathbf{0}$  at  $\mathbf{x}^*$ . However, the right-hand side of Eq. (13) is not locally Lipschitz in any neighborhood of  $\mathbf{x}^*$ , since the normalization map  $\mathbf{g} \mapsto \mathbf{g}/\|\mathbf{g}\|$  becomes singular at  $\mathbf{g} = \mathbf{0}$ . Consequently, the Picard–Lindelöf theorem cannot be applied at such points, and multiple absolutely continuous solutions may originate from the same initial condition  $\mathbf{x}(0) = \mathbf{x}^*$ . ■

## Appendix D. Proof for LMO-based MFLD in continuous time

### D.1. Proof of Condition 1

**Proof** By the mass condition:

$$\|g_\mu\|_{L^1(\mu; \|\cdot\|_*)} = \int \|g_\mu(\mathbf{x})\|_* d\mu(\mathbf{x}) \geq \int_{\{\|g_\mu(\mathbf{x})\|_* \geq m_1\}} \|g_\mu(\mathbf{x})\|_* d\mu(\mathbf{x}) \geq m_1 m_2.$$

Using the uniform bound  $\|g_\mu(\mathbf{x})\|_* \leq G_{\max}$   $\mu$ -a.s.,

$$\|g_\mu\|_{L^2(\mu; \|\cdot\|_*)}^2 \leq \int G_{\max}^2 d\mu(\mathbf{x}) = G_{\max}^2.$$

Therefore,  $\|g_\mu\|_{L^1(\mu; \|\cdot\|_*)} \geq m_1 m_2 \geq \frac{m_1 m_2}{G_{\max}} \|g_\mu\|_{L^2(\mu; \|\cdot\|_*)}$ . This proves Asm 4 with  $\kappa = m_1 m_2 / G_{\max}$ .

The empirical case follows by replacing integrals with empirical averages:  $\int f(\mathbf{x}) d\hat{\mu}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N f(\mathbf{x}_i)$ . ■

## D.2. Connection of Asm 4 to reverse Hölder inequality

**Lemma 1 (Reverse form of Hölder's inequality [34])** Let  $1 < p \leq 2$  and let  $q = \frac{p}{p-1}$  be its conjugate exponent. For any  $g \in L^q(\mu)$  and  $h \in L^p(\mu)$  with  $\|h\|_{L^p(\mu)}, \|g\|_{L^q(\mu)} > 0$ , then we have:

$$\|hg\|_{L^1(\mu)} \geq \|h\|_{L^p(\mu)} \|g\|_{L^q(\mu)} \left( 1 - \frac{1}{p} \left\| \frac{|h|^{p/2}}{\|h\|_{L^p(\mu)}^{p/2}} - \frac{|g|^{q/2}}{\|g\|_{L^q(\mu)}^{q/2}} \right\|_{L^2(\mu)}^2 \right).$$

**Remark 9** Asm 4 can be viewed as a simplified consequence of the reverse form of Hölder's inequality (Lemma 1). Indeed, in the special case  $p = q = 2$  and  $h = 1$ , Lemma 1 yields the quantitative bound  $\|g\|_{L^1(\mu)} \geq \|g\|_{L^2(\mu)} \left( 1 - \frac{1}{2} \int \left( 1 - \frac{|g(\mathbf{x})|}{\|g\|_{L^2(\mu)}} \right)^2 d\mu(\mathbf{x}) \right)$ . Hence, whenever  $\left( 1 - \frac{|g(\mathbf{x})|}{\|g\|_{L^2(\mu)}} \right)$  is smaller, then Asm 4 holds. This type of non-degeneracy condition is conceptually related to the coverage or exploration inequalities used in reinforcement learning analysis (e.g., Lemma 3 in Liu et al. [27]).

## D.3. Smoothness for two-layer mean-field neural networks

**Proposition 1 (Smoothness for two-layer mean-field neural networks)** Consider the two-layer mean-field neural network  $f_\mu(\mathbf{z}) = \int h(\mathbf{x}, \mathbf{z}) d\mu(\mathbf{x})$ ,  $F(\mu) = \mathbb{E}_{(\mathbf{z}, y) \sim \mathcal{D}} [\ell(f_\mu(\mathbf{z}), y)]$ . Assume that: (i) the neuron  $h(\mathbf{x}, \mathbf{z})$  is differentiable in  $\mathbf{x}$ ; (ii)  $\nabla_{\mathbf{x}} h(\mathbf{x}, \mathbf{z})$  is bounded in  $\|\cdot\|_*$  and Lipschitz continuous with respect to  $\|\cdot\|$  and  $\|\cdot\|_*$  (e.g.,  $h(\mathbf{x}, \mathbf{z}) = \sigma(\mathbf{x}^\top \mathbf{z})$  with  $\sigma$  being tanh or sigmoid activation); (iii) the loss  $\ell(f, y)$  has Lipschitz continuous first derivative in  $f$  (e.g., logistic loss; squared loss under bounded outputs). Then Asm 2 holds.

**Proof** We first compute the first variation:

$$\frac{\delta F}{\delta \mu}(\mu)(\mathbf{x}) = \mathbb{E}_{(\mathbf{z}, y)} \left[ \partial_f \ell(f_\mu(\mathbf{z}), y) h(\mathbf{x}, \mathbf{z}) \right].$$

Taking gradient with respect to  $\mathbf{x}$  gives

$$\nabla \frac{\delta F}{\delta \mu}(\mu)(\mathbf{x}) = \mathbb{E}_{(\mathbf{z}, y)} \left[ \partial_f \ell(f_\mu(\mathbf{z}), y) \nabla_{\mathbf{x}} h(\mathbf{x}, \mathbf{z}) \right].$$

Let  $\mu, \nu \in \mathcal{P}_2(\mathbb{R}^d)$  and  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$ . Using the convexity of the dual norm,

$$\left\| \nabla \frac{\delta F}{\delta \mu}(\mu)(\mathbf{x}) - \nabla \frac{\delta F}{\delta \mu}(\nu)(\mathbf{x}') \right\|_* \leq \mathbb{E}_{(\mathbf{z}, y)} [A_1(\mathbf{z}, y) + A_2(\mathbf{z}, y)],$$

where we define

$$\begin{aligned} A_1(\mathbf{z}, y) &:= \left| \partial_f \ell(f_\mu(\mathbf{z}), y) - \partial_f \ell(f_\nu(\mathbf{z}), y) \right| \cdot \|\nabla_{\mathbf{x}} h(\mathbf{x}, \mathbf{z})\|_*, \\ A_2(\mathbf{z}, y) &:= \left| \partial_f \ell(f_\nu(\mathbf{z}), y) \right| \cdot \|\nabla_{\mathbf{x}} h(\mathbf{x}, \mathbf{z}) - \nabla_{\mathbf{x}} h(\mathbf{x}', \mathbf{z})\|_*. \end{aligned}$$

By assumption, there exist constants  $B_1, B_2, H_1, H_2 > 0$  such that for all  $\mathbf{z}, y, \mathbf{x}, \mathbf{x}'$ :

$$\begin{aligned} |\partial_f \ell(f, y)| &\leq B_1, \\ |\partial_f \ell(f_1, y) - \partial_f \ell(f_2, y)| &\leq B_2 |f_1 - f_2|, \\ \|\nabla_{\mathbf{x}} h(\mathbf{x}, \mathbf{z})\|_* &\leq H_1, \\ \|\nabla_{\mathbf{x}} h(\mathbf{x}, \mathbf{z}) - \nabla_{\mathbf{x}} h(\mathbf{x}', \mathbf{z})\|_* &\leq H_2 \|\mathbf{x} - \mathbf{x}'\|. \end{aligned}$$

Therefore,

$$A_1(\mathbf{z}, y) \leq B_2 H_1 |f_\mu(\mathbf{z}) - f_\nu(\mathbf{z})|, \quad A_2(\mathbf{z}, y) \leq B_1 H_2 \|\mathbf{x} - \mathbf{x}'\|.$$

Taking expectation gives

$$\left\| \nabla \frac{\delta F}{\delta \mu}(\mu)(\mathbf{x}) - \nabla \frac{\delta F}{\delta \mu}(\nu)(\mathbf{x}') \right\|_* \leq B_2 H_1 \mathbb{E}_{\mathbf{z}} |f_\mu(\mathbf{z}) - f_\nu(\mathbf{z})| + B_1 H_2 \|\mathbf{x} - \mathbf{x}'\|.$$

It remains to bound  $\mathbb{E}_{\mathbf{z}} |f_\mu(\mathbf{z}) - f_\nu(\mathbf{z})|$ . For any coupling  $\pi$  between  $\mu$  and  $\nu$ ,

$$\begin{aligned} |f_\mu(\mathbf{z}) - f_\nu(\mathbf{z})| &= \left| \int h(\mathbf{x}, \mathbf{z}) d\mu(\mathbf{x}) - \int h(\mathbf{x}', \mathbf{z}) d\nu(\mathbf{x}') \right| \\ &\leq \int |h(\mathbf{x}, \mathbf{z}) - h(\mathbf{x}', \mathbf{z})| d\pi(\mathbf{x}, \mathbf{x}'). \end{aligned}$$

Since  $\|\nabla_{\mathbf{x}} h(\mathbf{x}, \mathbf{z})\|_* \leq H_1$ , the map  $\mathbf{x} \mapsto h(\mathbf{x}, \mathbf{z})$  is  $H_1$ -Lipschitz w.r.t.  $\|\cdot\|$ , hence

$$|h(\mathbf{x}, \mathbf{z}) - h(\mathbf{x}', \mathbf{z})| \leq H_1 \|\mathbf{x} - \mathbf{x}'\|.$$

Therefore,

$$|f_\mu(\mathbf{z}) - f_\nu(\mathbf{z})| \leq H_1 \int \|\mathbf{x} - \mathbf{x}'\| d\pi(\mathbf{x}, \mathbf{x}'),$$

and taking the infimum over all couplings gives

$$|f_\mu(\mathbf{z}) - f_\nu(\mathbf{z})| \leq H_1 W_1(\mu, \nu) \leq H_1 W_2(\mu, \nu).$$

Taking expectation over  $\mathbf{z}$  yields

$$\mathbb{E}_{\mathbf{z}} |f_\mu(\mathbf{z}) - f_\nu(\mathbf{z})| \leq H_1 W_2(\mu, \nu).$$

Substituting back, we obtain

$$\begin{aligned} \left\| \nabla \frac{\delta F}{\delta \mu}(\mu)(\mathbf{x}) - \nabla \frac{\delta F}{\delta \mu}(\nu)(\mathbf{x}') \right\|_* &\leq B_2 H_1^2 W_2(\mu, \nu) + B_1 H_2 \|\mathbf{x} - \mathbf{x}'\|. \\ &\leq \frac{L}{2} (W_2(\mu, \nu) + \|\mathbf{x} - \mathbf{x}'\|), \end{aligned} \tag{18}$$

with

$$L = 2 \max\{B_2 H_1^2, B_1 H_2\}.$$

This proves the claim. ■

#### D.4. Proof of Thm. 4

**Proof** We first present the following lemma that is crucial in our analysis.

**Lemma 2 (Entropy Sandwich [12])** Under Asm. 2 and 3 and assuming that  $F(\mu)$  is convex w.r.t  $\mu$ , let  $\mu^*$  be the unique minimizer of  $\mathcal{F}$ . For any  $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ , given the definition of  $\nu$  in Asm 3, the following inequality holds:

$$\mathcal{F}(\mu) - \mathcal{F}(\mu^*) \leq \lambda H(\mu \| \nu). \quad (19)$$

Now we start our proof. Define the suboptimality  $\mathcal{E}(t) := \mathcal{F}(\mu_t) - \mathcal{F}(\mu^*) \geq 0$ . Since

$$\partial_t \mu_t + \nabla \cdot (\mu_t v_t) = 0,$$

starting from the chain rule in Wasserstein space and the definition of  $v_t$  gives

$$\begin{aligned} \frac{d}{dt} \mathcal{E}(t) &= \frac{d}{dt} \mathcal{F}(\mu_t) \\ &= \int \langle \nabla_{\mathbf{x}} \frac{\delta \mathcal{F}}{\delta \mu}(\mu_t)(\mathbf{x}), v_t(\mathbf{x}) \rangle d\mu_t(\mathbf{x}) \\ &= \int \langle \nabla_{\mathbf{x}} \frac{\delta \mathcal{F}}{\delta \mu}(\mu_t)(\mathbf{x}), \text{lmo}(\nabla_{\mathbf{x}} \frac{\delta \mathcal{F}}{\delta \mu}(\mu_t)(\mathbf{x}); 0) \rangle d\mu_t(\mathbf{x}) \\ &= - \int \left\| \nabla_{\mathbf{x}} \frac{\delta \mathcal{F}}{\delta \mu}(\mu_t)(\mathbf{x}) \right\|_* d\mu_t(\mathbf{x}) \\ &\leq -\kappa \left( \int \left\| \nabla_{\mathbf{x}} \frac{\delta \mathcal{F}}{\delta \mu}(\mu_t)(\mathbf{x}) \right\|_*^2 d\mu_t(\mathbf{x}) \right)^{1/2}, \end{aligned} \quad (20)$$

where in the second last step, we use Lemma 7 and the last step follows from Asm 4.

According to the ‘‘entropy sandwich’’ lemma (Lemma 2) and the definition of  $v_t$ ,

$$\mathcal{E}(t) \leq \lambda H(\mu_t \| \nu_t), \quad \nu_t \propto \exp\left(-\frac{1}{\lambda} \frac{\delta F}{\delta \mu}(\mu_t)\right), \quad (21)$$

we have:

$$\nabla_{\mathbf{x}} \log \frac{\mu_t}{\nu_t} = \nabla_{\mathbf{x}} \log \mu_t + \frac{1}{\lambda} \nabla_{\mathbf{x}} \frac{\delta F}{\delta \mu}(\mu_t) = \frac{1}{\lambda} \nabla_{\mathbf{x}} \frac{\delta \mathcal{F}}{\delta \mu}(\mu_t). \quad (22)$$

Next, by the uniform log-Sobolev inequality (Asm 3), we get

$$\frac{\mathcal{E}(t)}{\lambda} \leq H(\mu_t \| \nu_t) \leq \frac{1}{2\zeta \lambda^2} \int \left\| \nabla_{\mathbf{x}} \frac{\delta \mathcal{F}}{\delta \mu}(\mu_t)(\mathbf{x}) \right\|_*^2 d\mu_t(\mathbf{x}). \quad (23)$$

From Eq. (20) and Eq. (23), we obtain

$$\frac{d}{dt} \mathcal{E}(t) \leq -\kappa \left( \int \left\| \nabla_{\mathbf{x}} \frac{\delta \mathcal{F}}{\delta \mu}(\mu_t)(\mathbf{x}) \right\|_*^2 d\mu_t(\mathbf{x}) \right)^{1/2} \leq -\kappa (2\zeta \lambda \mathcal{E}(t))^{1/2}. \quad (24)$$

Let  $G(t) := \sqrt{\mathcal{E}(t)}$ . Then  $\frac{d}{dt} G(t) \leq -\frac{\kappa}{2} \sqrt{2\zeta \lambda}$  so that  $G(t) \leq \max\{G(0) - \frac{\kappa}{2} \sqrt{2\zeta \lambda} t, 0\}$ . Taking the square completes the proof.  $\blacksquare$

---

**Algorithm 2: LMO Noisy Particle Descent (LMO-NPD)**


---

**Data:** number of particles  $N$ ; step sizes  $\{\eta_k\}$ ; diffusion level  $\lambda > 0$ ; iterations  $K$

**Result:** empirical measure  $\mu_K^{(N)} = \frac{1}{N} \sum_{i=1}^N \delta_{\mathbf{x}_K^i}$

Initialize particles  $\mathbf{x}_0^i \sim \mu_0$  for  $i = 1, \dots, N$ ;

**for**  $k = 0$  **to**  $K - 1$  **do**

$\mu_k^{(N)} \leftarrow \frac{1}{N} \sum_{i=1}^N \delta_{\mathbf{x}_k^i}$ ;

**for**  $i = 1$  **to**  $N$  **do**

Evaluate mean-field force  $\mathbf{v}_k^i$  according to Eq. (87) or Eq. (5);

Sample Gaussian noise  $\boldsymbol{\xi}_k^i \sim \mathcal{N}(0, I_d)$ ;

$\mathbf{x}_{k+1}^i \leftarrow \mathbf{x}_k^i - \eta_k \mathbf{v}_k^i + \sqrt{2\lambda\eta_k} \boldsymbol{\xi}_k^i$ ;

**end**

**end**

---

## Appendix E. Discrete time analysis of LMO noisy particle descent

In this section, we consider the empirical risk minimization problem of mean-field neural networks. Throughout, the norm  $\|\cdot\|$  and  $\|\cdot\|_*$  on  $\mathbb{R}^{dN}$  are understood as the sum of the corresponding norms over  $N$  blocks. Let  $h(\mathbf{x}, \cdot) : \mathcal{Z} \rightarrow \mathbb{R}$  be the neuron with parameter  $\mathbf{x} \in \mathbb{R}^d$ , where  $\mathcal{Z}$  is the data space, and let  $F(\mu) = \frac{1}{n} \sum_{j=1}^N \ell(\mathbb{E}_{\mathbf{x} \sim \mu} [h(\mathbf{x}, \mathbf{z}_j)], y_j)$ , where  $(\mathbf{z}_i, y_i)_{i=1}^N$  are training data.

We then define the following objective on  $\mathcal{P}^{(N)}$ :

$$\mathcal{F}^N(\mu^{(N)}) = N \mathbb{E}_{\mathbf{X} \sim \mu^{(N)}} [F(\mu_{\mathbf{X}})] + \lambda \text{Ent}(\mu^{(N)}).$$

Below, we introduce the assumption on the loss and the neuron, which has been used in Nitanda [30].

### Assumption 6 (From Nitanda [30])

- The loss  $\ell(\cdot, y)$  is convex and  $\hat{L}$ -smooth.
- There exists  $\hat{R} > 0$  such that, for all  $\mathbf{z} \in \mathcal{Z}$ , we have  $\mathbb{E}_{\mathbf{x} \sim \mu^*} [|h(\mathbf{x}, \mathbf{z})|^2] \leq \hat{R}^2$ .

Below, we provide an assumption that is similar to Asm 2 and similar to the bounded gradient assumption.

**Assumption 7** *There exists a constant  $M_1 > 0$  such that  $\|\mathbf{v}_k^i\|_2 \leq M_1, \forall k \geq 0, i \in [N]$ . Fix a norm  $\|\cdot\|$  on  $\mathbb{R}^d$  with dual  $\|\cdot\|_*$ . There exists a constant  $M_2 > 0$  s.t.  $\forall \mu, \nu$  on  $(\mathbb{R}^d)^N$  and  $\forall \mathbf{Y}, \mathbf{X} \in (\mathbb{R}^d)^N$ :*

$$\|\nabla \log \mu(\mathbf{Y}) - \nabla \log \nu(\mathbf{X})\|_* \leq M_2 \left( W_2(\mu, \nu) + \|\mathbf{Y} - \mathbf{X}\| \right).$$

Now we are ready to provide the convergence theorem for Alg. 2 with proof deferred to Sec. E.1.

**Theorem 10 (Convergence of LMO-NPD)** *Let  $\{\mu_k^{(N)}\}_{k \geq 0}$  be the measure generated by Alg. 2 with drift selected by Eq. (87). Define:  $C_1 = \sup_{\mathbf{a} \in \mathbb{R}^{dN} \setminus \{0\}} \frac{\|\mathbf{a}\|_2^2}{\|\mathbf{a}\|_*^2}$ . and  $C_2 := \frac{8}{\kappa \sqrt{2\zeta\lambda}} (LN +$*

$\lambda M_2 N + \lambda M_2 \sqrt{N}$ )  $\sqrt{2\eta(M_1^2 + 2\lambda \mathbb{E}\|\xi\|^2)}$ . Consider a fixed step size  $\eta$  that satisfies:

$$0 < \eta \leq \min \left\{ 1, \frac{\kappa^2}{32 \lambda^2 M_2^2 C_1 N (M_1^2 + 2\lambda \mathbb{E}\|\xi\|^2)} \right\}.$$

Then, under Asm. 1 to 4, 6 and 7, for all  $k \geq 0$ ,

$$\begin{aligned} & \frac{1}{N} \mathcal{F}^{(N)}(\mu_k^{(N)}) - \mathcal{F}(\mu^*) \\ & \leq \frac{\hat{L}\hat{R}^2}{2N} + \frac{1}{N} \left[ \max \left\{ \sqrt{\mathcal{F}^{(N)}(\mu_0^{(N)}) - \mathcal{F}^{(N)}(\mu_*^{(N)})} - \frac{\kappa}{8} \sqrt{2\zeta\lambda k\eta}, C_2 \right\} \right]^2. \end{aligned}$$

### E.1. Proof for Thm. 10

#### E.1.1. AUXILIARY LEMMAS

**Lemma 3 (Theorem 1 in Nitanda [30])** Let  $\mu^*$  be the minimizer of  $\mathcal{F}$  and  $\mu_*^{(N)}$  the minimizer of  $\mathcal{F}^{(N)}$ , and denote by  $(\mu^*)^{\otimes N}$  the  $N$ -fold product of  $\mu^*$ . Under Asm. 1, 2 and 6, we have:

$$\frac{1}{N} \mathcal{F}^{(N)}(\mu_*^{(N)}) - \mathcal{F}(\mu^*) \leq \frac{\hat{L}\hat{R}^2}{2N}.$$

**Lemma 4 (LMO mismatch inequality on the unit ball)** Let  $(\mathbb{R}^d, \|\cdot\|)$  be a normed space with dual norm  $\|\cdot\|_*$ . Then for all  $\mathbf{h}, \mathbf{g} \in \mathbb{R}^d$ ,

$$0 \leq -\langle \mathbf{h}, \text{lmo}(\mathbf{h}; 0) - \text{lmo}(\mathbf{g}; 0) \rangle \leq 2 \|\mathbf{h} - \mathbf{g}\|_*.$$

**Proof** By optimality of  $\text{lmo}(\mathbf{h}; 0)$  for  $\mathbf{h}$  over  $\{\mathbf{s} : \|\mathbf{s}\| \leq 1\}$ ,

$$\langle \mathbf{h}, \text{lmo}(\mathbf{h}; 0) \rangle \leq \langle \mathbf{h}, \text{lmo}(\mathbf{g}; 0) \rangle,$$

hence

$$-\langle \mathbf{h}, \text{lmo}(\mathbf{h}; 0) - \text{lmo}(\mathbf{g}; 0) \rangle \geq 0.$$

By optimality of  $\text{lmo}(\mathbf{g}; 0)$  for  $\mathbf{g}$  over  $\{\mathbf{s} : \|\mathbf{s}\| \leq 1\}$ ,

$$\langle \mathbf{g}, \text{lmo}(\mathbf{g}; 0) \rangle \leq \langle \mathbf{g}, \text{lmo}(\mathbf{h}; 0) \rangle,$$

so

$$\langle \mathbf{g}, \text{lmo}(\mathbf{h}; 0) - \text{lmo}(\mathbf{g}; 0) \rangle \geq 0.$$

Therefore,

$$\begin{aligned} & -\langle \mathbf{h}, \text{lmo}(\mathbf{h}; 0) - \text{lmo}(\mathbf{g}; 0) \rangle \\ & = -\langle \mathbf{h} - \mathbf{g}, \text{lmo}(\mathbf{h}; 0) - \text{lmo}(\mathbf{g}; 0) \rangle - \langle \mathbf{g}, \text{lmo}(\mathbf{h}; 0) - \text{lmo}(\mathbf{g}; 0) \rangle \\ & \leq -\langle \mathbf{h} - \mathbf{g}, \text{lmo}(\mathbf{h}; 0) - \text{lmo}(\mathbf{g}; 0) \rangle \\ & \leq \|\mathbf{h} - \mathbf{g}\|_* \|\text{lmo}(\mathbf{h}; 0) - \text{lmo}(\mathbf{g}; 0)\|, \end{aligned}$$

where the last line uses Hölder's inequality. Finally, by the triangle inequality,

$$\|\text{lmo}(\mathbf{h}; 0) - \text{lmo}(\mathbf{g}; 0)\| \leq \|\text{lmo}(\mathbf{h}; 0)\| + \|\text{lmo}(\mathbf{g}; 0)\| \leq 2,$$

and hence

$$-\langle \mathbf{h}, \text{lmo}(\mathbf{h}; 0) - \text{lmo}(\mathbf{g}; 0) \rangle \leq 2 \|\mathbf{h} - \mathbf{g}\|_*.$$

■

**Lemma 5** Under the same setting and assumption as in Thm. 10, we have

$$\frac{d}{dt} \mathcal{F}^{(N)}(\nu_t) = - \int \nu_t(\mathbf{Y}) \|G_t(\mathbf{Y})\|_2^2 d\mathbf{Y} - \int \nu_t(\mathbf{Y}) G_t(\mathbf{Y}) \delta_t(\mathbf{Y}) d\mathbf{Y}.$$

where

$$\delta_t(\mathbf{Y}) := \mathbb{E}[-\text{lmo}(G_0(\mathbf{Y}_0); 0) - \lambda \nabla \log \nu_0(\mathbf{Y}_0) | \mathbf{Y}_t = \mathbf{Y}] - N \nabla F(\mathbf{Y}). \quad (25)$$

$$G_t(\mathbf{Y}) := \nabla \frac{\delta \mathcal{F}^{(N)}(\nu_t)}{\delta \mu^{(N)}}(\mathbf{Y}) = N \nabla F(\mathbf{Y}) + \lambda \nabla \log \nu_t(\mathbf{Y}) = \lambda \nabla \log \frac{\nu_t}{\mu_*^{(N)}}(\mathbf{Y}). \quad (26)$$

**Proof** The proof follows the idea in Nitanda [30]. We study the discrete update

$$\mathbf{x}_{k+1}^i = \mathbf{x}_k^i - \eta \mathbf{v}_k^i + \sqrt{2\lambda\eta} \boldsymbol{\xi}_k^i, \quad i \in \{1, \dots, N\},$$

where  $\boldsymbol{\xi}_k^i \sim \mathcal{N}(0, I_d)$  are i.i.d. Gaussian noises. Let  $\mathbf{X}_k = (\mathbf{x}_k^1, \dots, \mathbf{x}_k^N)$  and  $\mu_k^{(N)} = \text{Law}(\mathbf{X}_k)$ .

To analyze one step, we introduce the one-step interpolation process

$$d\mathbf{y}_t^i = -\mathbf{v}_k^i dt + \sqrt{2\lambda} d\mathbf{w}_t^i, \quad t \in [0, \eta],$$

with  $\mathbf{Y}_0 = \mathbf{X}_k$ . Its explicit solution is

$$\mathbf{y}_t^i = \mathbf{y}_0^i - t \mathbf{v}_k^i + \sqrt{2\lambda t} \boldsymbol{\xi}^i, \quad \boldsymbol{\xi}^i \sim \mathcal{N}(0, I_d),$$

so that  $\mathbf{Y}_\eta = \mathbf{X}_{k+1}$ . We denote  $\nu_t := \text{Law}(\mathbf{Y}_t)$ , hence  $\nu_0 = \mu_k^{(N)}$  and  $\nu_\eta = \mu_{k+1}^{(N)}$ .

For notational simplicity we identify measures with their densities with respect to the Lebesgue measure and write  $\mu_*^{(N)}(\mathbf{Y})$  for the density of the target Gibbs measure. Let  $\nu_{0t}(\mathbf{Y}_0, \mathbf{Y})$  be the joint law of  $(\mathbf{Y}_0, \mathbf{Y}_t)$ , and  $\nu_{t|0}, \nu_{0|t}$  the corresponding conditional distributions. Then

$$\nu_{0t}(\mathbf{Y}_0, \mathbf{Y}) = \nu_0(\mathbf{Y}_0) \nu_{t|0}(\mathbf{Y} | \mathbf{Y}_0) = \nu_t(\mathbf{Y}) \nu_{0|t}(\mathbf{Y}_0 | \mathbf{Y}).$$

The conditional density satisfies

$$\frac{\partial}{\partial t} \nu_{t|0}(\mathbf{Y} | \mathbf{Y}_0) = \nabla \cdot (\nu_{t|0}(\mathbf{Y} | \mathbf{Y}_0) \mathbf{V}_k) + \lambda \Delta \nu_{t|0}(\mathbf{Y} | \mathbf{Y}_0),$$

where  $\mathbf{V}_k := (\mathbf{v}_k^1, \dots, \mathbf{v}_k^N)$  is constant in  $t$ , i.e.,

$$\begin{aligned} \mathbf{V}_k(\mathbf{Y}_0) &= -\text{lmo}(N \nabla F(\mathbf{Y}_0) + \lambda \nabla \log \nu_0(\mathbf{Y}_0); 0) - \lambda \nabla \log \nu_0(\mathbf{Y}_0) \\ &= -\text{lmo}(G_0(\mathbf{Y}_0); 0) - \lambda \nabla \log \nu_0(\mathbf{Y}_0) \end{aligned} \quad (27)$$

where the LMO acts blockwise over the  $N$  particles. And the last inequality is because that  $\mu_*^{(N)}$  is with density proportional to  $\exp(-\frac{N}{\lambda}F(\mu_{\mathbf{Y}}))$ , we have  $\lambda \nabla \log \mu_*^{(N)}(\mathbf{Y}) = -N \nabla F(\mathbf{Y})$ , so we have (see Chen et al. [11] for the proof):

$$\nabla \frac{\delta \mathcal{F}^{(N)}(\nu_t)}{\delta \mu^{(N)}}(\mathbf{Y}) = N \nabla F(\mathbf{Y}) + \lambda \nabla \log \nu_t(\mathbf{Y}) = \lambda \nabla \log \frac{\nu_t}{\mu_*^{(N)}}(\mathbf{Y}). \quad (28)$$

Integrating against  $\nu_0(\mathbf{Y}_0)$  and using  $\nu_{0t} = \nu_0 \nu_{t|0} = \nu_t \nu_{0|t}$ , we get

$$\begin{aligned} \frac{\partial}{\partial t} \nu_t(\mathbf{Y}) &= \int \frac{\partial}{\partial t} \nu_{t|0}(\mathbf{Y} | \mathbf{Y}_0) \nu_0(\mathbf{Y}_0) d\mathbf{Y}_0 \\ &= \int \left( \nabla \cdot (\nu_{0t}(\mathbf{Y}_0, \mathbf{Y}) \mathbf{V}_k) + \lambda \Delta \nu_{0t}(\mathbf{Y}_0, \mathbf{Y}) \right) d\mathbf{Y}_0 \\ &= \nabla \cdot \left( \nu_t(\mathbf{Y}) \underbrace{\int \nu_{0|t}(\mathbf{Y}_0 | \mathbf{Y}) \mathbf{V}_k(\mathbf{Y}_0) d\mathbf{Y}_0}_{= \mathbb{E}[\mathbf{V}_k(\mathbf{Y}_0) | \mathbf{Y}_t = \mathbf{Y}]} \right) + \lambda \Delta \nu_t(\mathbf{Y}). \end{aligned}$$

By adding and subtracting yields

$$\frac{\partial}{\partial t} \nu_t(\mathbf{Y}) = \lambda \nabla \cdot \left( \nu_t(\mathbf{Y}) \nabla \log \frac{\nu_t}{\mu_*^{(N)}}(\mathbf{Y}) \right) + \nabla \cdot \left( \nu_t(\mathbf{Y}) \delta_t(\mathbf{Y}) \right),$$

where we define

$$\begin{aligned} \delta_t(\mathbf{Y}) &:= \mathbb{E}[\mathbf{V}_k(\mathbf{Y}_0) | \mathbf{Y}_t = \mathbf{Y}] - N \nabla F(\mathbf{Y}) \\ &= \mathbb{E}[-\text{Imo}(G_0(\mathbf{Y}_0); 0) - \lambda \nabla \log \nu_0(\mathbf{Y}_0) | \mathbf{Y}_t = \mathbf{Y}] - N \nabla F(\mathbf{Y}) \end{aligned} \quad (29)$$

Differentiating along  $t \mapsto \nu_t$  and integrating by parts, we obtain:

$$\begin{aligned} &\frac{d}{dt} \mathcal{F}^{(N)}(\nu_t) \\ &= \int \frac{\delta \mathcal{F}^{(N)}(\nu_t)}{\delta \mu^{(N)}}(\mathbf{Y}) \frac{\partial \nu_t}{\partial t}(\mathbf{Y}) d\mathbf{Y} \\ &= \lambda \int \frac{\delta \mathcal{F}^{(N)}(\nu_t)}{\delta \mu^{(N)}}(\mathbf{Y}) \nabla \cdot \left( \nu_t(\mathbf{Y}) \nabla \log \frac{\nu_t}{\mu_*^{(N)}}(\mathbf{Y}) \right) d\mathbf{Y} + \int \frac{\delta \mathcal{F}^{(N)}(\nu_t)}{\delta \mu^{(N)}}(\mathbf{Y}) \nabla \cdot \left( \nu_t(\mathbf{Y}) \delta_t(\mathbf{Y}) \right) d\mathbf{Y} \\ &= -\lambda \int \nu_t(\mathbf{Y}) \nabla \frac{\delta \mathcal{F}^{(N)}(\nu_t)}{\delta \mu^{(N)}}(\mathbf{Y})^\top \nabla \log \frac{\nu_t}{\mu_*^{(N)}}(\mathbf{Y}) d\mathbf{Y} - \int \nu_t(\mathbf{Y}) \nabla \frac{\delta \mathcal{F}^{(N)}(\nu_t)}{\delta \mu^{(N)}}(\mathbf{Y})^\top \delta_t(\mathbf{Y}) d\mathbf{Y}. \end{aligned}$$

And the proof is finished by using Eq. (25).  $\blacksquare$

**Lemma 6** Under the same setting and assumptions as in Thm. 10, fix a norm  $\|\cdot\|$  on  $\mathbb{R}^d$ . For all  $t \in [0, \eta]$  we have

$$\mathbb{E}_{\nu_{0t}} \left[ \sum_{i=1}^N \|\mathbf{y}_t^i - \mathbf{y}_0^i\|^2 \right] \leq 2Nt^2 M_1^2 + 4\lambda t N \mathbb{E} \|\boldsymbol{\xi}\|^2, \quad (30)$$

where  $\boldsymbol{\xi} \sim \mathcal{N}(0, I_d)$ . Moreover,

$$\mathbb{E}_{\nu_{0t}} \left[ \sum_{i=1}^N \|\mathbf{y}_t^i - \mathbf{y}_0^i\| \right] \leq N \sqrt{2t^2 M_1^2 + 4\lambda t \mathbb{E} \|\boldsymbol{\xi}\|^2}. \quad (31)$$

**Proof** Recall the interpolation SDE for  $t \in [0, \eta]$ :

$$d\mathbf{y}_t^i = -\mathbf{v}_k^i dt + \sqrt{2\lambda} d\mathbf{w}_t^i, \quad \mathbf{y}_0^i = \mathbf{x}_k^i,$$

whose explicit solution is

$$\mathbf{y}_t^i = \mathbf{y}_0^i - t\mathbf{v}_k^i + \sqrt{2\lambda t} \boldsymbol{\xi}^i, \quad \boldsymbol{\xi}^i \sim \mathcal{N}(0, \mathbf{I}_d) \text{ i.i.d.}$$

Hence

$$\mathbf{y}_t^i - \mathbf{y}_0^i = -t\mathbf{v}_k^i + \sqrt{2\lambda t} \boldsymbol{\xi}^i.$$

Using the elementary inequality (valid for any norm)

$$\|\mathbf{a} + \mathbf{b}\|^2 \leq 2\|\mathbf{a}\|^2 + 2\|\mathbf{b}\|^2,$$

we obtain

$$\begin{aligned} \|\mathbf{y}_t^i - \mathbf{y}_0^i\|^2 &\leq 2t^2\|\mathbf{v}_k^i\|^2 + 2\|\sqrt{2\lambda t} \boldsymbol{\xi}^i\|^2 \\ &= 2t^2\|\mathbf{v}_k^i\|^2 + 4\lambda t \|\boldsymbol{\xi}^i\|^2. \end{aligned}$$

Summing over  $i = 1, \dots, N$ ,

$$\sum_{i=1}^N \|\mathbf{y}_t^i - \mathbf{y}_0^i\|^2 \leq 2t^2 \sum_{i=1}^N \|\mathbf{v}_k^i\|^2 + 4\lambda t \sum_{i=1}^N \|\boldsymbol{\xi}^i\|^2.$$

Taking expectation and using Asm 7 (i.e.,  $\|\mathbf{v}_k^i\| \leq M_1$ ), as well as i.i.d. of  $\boldsymbol{\xi}^i$ , we get

$$\begin{aligned} \mathbb{E}_{\nu_{0t}} \left[ \sum_{i=1}^N \|\mathbf{y}_t^i - \mathbf{y}_0^i\|^2 \right] &\leq 2t^2 \mathbb{E} \left[ \sum_{i=1}^N \|\mathbf{v}_k^i\|^2 \right] + 4\lambda t \mathbb{E} \left[ \sum_{i=1}^N \|\boldsymbol{\xi}^i\|^2 \right] \\ &\leq 2Nt^2 M_1^2 + 4\lambda t N \mathbb{E} \|\boldsymbol{\xi}\|^2, \end{aligned}$$

which proves Eq. (30).

For the first-moment bound, by Jensen's inequality and Cauchy–Schwarz  $(\sum_{i=1}^N a_i)^2 \leq N \sum_{i=1}^N a_i^2$  for  $a_i \geq 0$ , we have

$$\begin{aligned} \mathbb{E}_{\nu_{0t}} \left[ \sum_{i=1}^N \|\mathbf{y}_t^i - \mathbf{y}_0^i\| \right] &\leq \left( \mathbb{E}_{\nu_{0t}} \left[ \left( \sum_{i=1}^N \|\mathbf{y}_t^i - \mathbf{y}_0^i\| \right)^2 \right] \right)^{1/2} \\ &\leq \left( N \mathbb{E}_{\nu_{0t}} \left[ \sum_{i=1}^N \|\mathbf{y}_t^i - \mathbf{y}_0^i\|^2 \right] \right)^{1/2} \\ &\leq \left( N (2Nt^2 M_1^2 + 4\lambda t N \mathbb{E} \|\boldsymbol{\xi}\|^2) \right)^{1/2} \\ &= N \sqrt{2t^2 M_1^2 + 4\lambda t \mathbb{E} \|\boldsymbol{\xi}\|^2}, \end{aligned}$$

where we used the previous bound in the second-to-last inequality. This proves Eq. (31). ■

## E.1.2. PROOF OF THE MAIN THEOREM

**Proof** We use the interpolation scheme and notation of Lemma 5. Specifically, Lemma 5 gives

$$\frac{d}{dt} \mathcal{F}^{(N)}(\nu_t) = - \int \nu_t(\mathbf{Y}) \|G_t(\mathbf{Y})\|_2^2 d\mathbf{Y} - \int \nu_t(\mathbf{Y}) G_t(\mathbf{Y}) \delta_t(\mathbf{Y}) d\mathbf{Y}. \quad (32)$$

where

$$\delta_t(\mathbf{Y}) := \mathbb{E}_{\mathbf{Y}_0}[-\text{lmo}(G_0(\mathbf{Y}_0)) - \lambda \nabla \log \nu_0(\mathbf{Y}_0) \mid \mathbf{Y}_t = \mathbf{Y}] - N \nabla F(\mathbf{Y}).$$

Define

$$\begin{aligned} r_t(\mathbf{Y}) &:= r_t^{\text{lmo}}(\mathbf{Y}) + r_t^{\text{sc}}(\mathbf{Y}), \\ r_t^{\text{lmo}}(\mathbf{Y}) &:= \text{lmo}(G_t(\mathbf{Y}); 0) - \mathbb{E}_{\mathbf{Y}_0}[\text{lmo}(G_0(\mathbf{Y}_0); 0) \mid \mathbf{Y}_t = \mathbf{Y}], \\ r_t^{\text{sc}}(\mathbf{Y}) &:= \lambda \nabla \log \nu_t(\mathbf{Y}) - \mathbb{E}_{\mathbf{Y}_0}[\lambda \nabla \log \nu_0(\mathbf{Y}_0) \mid \mathbf{Y}_t = \mathbf{Y}], \end{aligned}$$

Continuing from Eq. (32), we obtain

$$\begin{aligned} \frac{d}{dt} \mathcal{F}^{(N)}(\nu_t) &= - \int \nu_t \|G_t\|_2^2 - \int \nu_t G_t^\top \delta_t \\ &= - \int \nu_t \|G_t\|_2^2 - \int \nu_t G_t^\top (-\text{lmo}(G_t; 0) - G_t + r_t) \\ &= - \int \nu_t \|G_t\|_2^2 + \int \nu_t \langle G_t, \text{lmo}(G_t; 0) \rangle + \int \nu_t \|G_t\|_2^2 - \int \nu_t \langle G_t, r_t \rangle \\ &= \int \nu_t \langle G_t, \text{lmo}(G_t; 0) \rangle - \int \nu_t \langle G_t, r_t^{\text{lmo}} \rangle - \int \nu_t \langle G_t, r_t^{\text{sc}} \rangle \\ &= - \int \nu_t \|G_t\|_* - \int \nu_t \langle G_t, r_t^{\text{lmo}} \rangle - \int \nu_t \langle G_t, r_t^{\text{sc}} \rangle, \end{aligned}$$

where we abbreviate  $\nu_t(\mathbf{Y}) d\mathbf{Y}$  by  $\nu_t$  and use Lemma 7. We now bound the term  $-\int \nu_t \langle G_t, r_t^{\text{lmo}} \rangle$ . Using the joint distribution  $\nu_{0t}$  of  $(\mathbf{Y}_0, \mathbf{Y})$  and Lemma 4, we obtain

$$\begin{aligned} - \int \nu_t \langle G_t, r_t^{\text{lmo}} \rangle &= - \mathbb{E}_{\mathbf{Y} \sim \nu_t} \left[ \langle G_t(\mathbf{Y}), \text{lmo}(G_t(\mathbf{Y}); 0) - \mathbb{E}_{\mathbf{Y}_0 \mid \mathbf{Y}} \text{lmo}(G_0(\mathbf{Y}_0); 0) \rangle \right] \\ &= - \mathbb{E}_{(\mathbf{Y}_0, \mathbf{Y}) \sim \nu_{0t}} \left[ \langle G_t(\mathbf{Y}), \text{lmo}(G_t(\mathbf{Y}); 0) - \text{lmo}(G_0(\mathbf{Y}_0); 0) \rangle \right] \\ &\leq 2 \mathbb{E}_{(\mathbf{Y}_0, \mathbf{Y}) \sim \nu_{0t}} [\|G_t(\mathbf{Y}) - G_0(\mathbf{Y}_0)\|_*]. \end{aligned} \quad (33)$$

We now bound  $\|G_t(\mathbf{Y}) - G_0(\mathbf{Y}_0)\|_*$ . Recall

$$\begin{aligned} G_t(\mathbf{Y}) &= N \nabla F(\mathbf{Y}) + \lambda \nabla \log \nu_t(\mathbf{Y}), \\ G_0(\mathbf{Y}_0) &= N \nabla F(\mathbf{Y}_0) + \lambda \nabla \log \nu_0(\mathbf{Y}_0). \end{aligned} \quad (34)$$

By the triangle inequality in the dual norm,

$$\|G_t(\mathbf{Y}) - G_0(\mathbf{Y}_0)\|_* \leq N \|\nabla F(\mathbf{Y}) - \nabla F(\mathbf{Y}_0)\|_* + \lambda \|\nabla \log \nu_t(\mathbf{Y}) - \nabla \log \nu_0(\mathbf{Y}_0)\|_*. \quad (35)$$

Since  $\nabla_{\mathbf{y}^i} F(\mathbf{Y}) = \frac{1}{N} \nabla_{\frac{\delta F}{\delta \mu}}(\mu_{\mathbf{Y}})(\mathbf{y}^i)$  (Eq. (28)) and Asm 2, we get

$$\begin{aligned} N \|\nabla F(\mathbf{Y}) - \nabla F(\mathbf{Y}_0)\|_* &\leq N \cdot \frac{L}{2N} \sum_{i=1}^N (W_2(\mu_{\mathbf{Y}}, \mu_{\mathbf{Y}_0}) + \|\mathbf{y}^i - \mathbf{y}_0^i\|) \\ &\leq N \cdot \frac{L}{2} \frac{1}{\sqrt{N}} \left( \sum_{i=1}^N \|\mathbf{y}^i - \mathbf{y}_0^i\|^2 \right)^{1/2} + \frac{L}{2} \sum_{i=1}^N \|\mathbf{y}^i - \mathbf{y}_0^i\|, \end{aligned} \quad (36)$$

where in the last step we used the empirical coupling  $\pi := \frac{1}{N} \sum_{i=1}^N \delta_{(\mathbf{y}_0^i, \mathbf{y}^i)} \in \Pi(\mu_{\mathbf{Y}_0}, \mu_{\mathbf{Y}})$ .

By Asm 7, we also have:

$$\begin{aligned} \|\nabla \log \nu_t(\mathbf{Y}) - \nabla \log \nu_0(\mathbf{Y}_0)\|_* &\leq M_2 \left( W_2(\nu_t, \nu_0) + \sum_{i=1}^N \|\mathbf{y}^i - \mathbf{y}_0^i\| \right) \\ &\leq M_2 \left( (\mathbb{E}_{\nu_{0t}} \|\mathbf{Y} - \mathbf{Y}_0\|^2)^{1/2} + \sum_{i=1}^N \|\mathbf{y}^i - \mathbf{y}_0^i\| \right) \\ &= M_2 \left( (\mathbb{E}_{\nu_{0t}} \left[ \sum_{i=1}^N \|\mathbf{y}^i - \mathbf{y}_0^i\|^2 \right])^{1/2} + \sum_{i=1}^N \|\mathbf{y}^i - \mathbf{y}_0^i\| \right) \end{aligned} \quad (37)$$

where we choose the coupling  $\nu_{0t} = \text{Law}(\mathbf{Y}_0, \mathbf{Y})$  in the second to the last step and in the last step is by the definition mentioned in the notation section that the norm  $\|\cdot\|$  on  $\mathbb{R}^{dN}$  is understood as the sum of the corresponding norms over the  $N$  blocks.

Substituting Eq. (36) and Eq. (37) into Eq. (35) and taking expectation with respect to  $\nu_{0t}$ , we obtain

$$\begin{aligned} &\mathbb{E}_{\nu_{0t}} \|G_t(\mathbf{Y}) - G_0(\mathbf{Y}_0)\|_* \\ &\leq N \cdot \frac{L}{2} \cdot \frac{1}{\sqrt{N}} \mathbb{E}_{\nu_{0t}} \left( \sum_{i=1}^N \|\mathbf{y}^i - \mathbf{y}_0^i\|^2 \right)^{1/2} + \frac{L}{2} \mathbb{E}_{\nu_{0t}} \sum_{i=1}^N \|\mathbf{y}^i - \mathbf{y}_0^i\| \\ &\quad + \lambda M_2 \mathbb{E}_{\nu_{0t}} \sum_{i=1}^N \|\mathbf{y}^i - \mathbf{y}_0^i\| + \lambda M_2 \left( \mathbb{E}_{\nu_{0t}} \sum_{i=1}^N \|\mathbf{y}^i - \mathbf{y}_0^i\|^2 \right)^{1/2} \\ &\leq N \cdot \frac{L}{2} \cdot \frac{1}{\sqrt{N}} \left( \mathbb{E}_{\nu_{0t}} \sum_{i=1}^N \|\mathbf{y}^i - \mathbf{y}_0^i\|^2 \right)^{1/2} + \frac{L}{2} \mathbb{E}_{\nu_{0t}} \sum_{i=1}^N \|\mathbf{y}^i - \mathbf{y}_0^i\| \\ &\quad + \lambda M_2 \mathbb{E}_{\nu_{0t}} \sum_{i=1}^N \|\mathbf{y}^i - \mathbf{y}_0^i\| + \lambda M_2 \left( \mathbb{E}_{\nu_{0t}} \sum_{i=1}^N \|\mathbf{y}^i - \mathbf{y}_0^i\|^2 \right)^{1/2} \\ &\leq N \cdot \frac{L}{2} \cdot \frac{1}{\sqrt{N}} \left( 2Nt^2 M_1^2 + 4\lambda t N \mathbb{E} \|\boldsymbol{\xi}\|^2 \right)^{1/2} + \frac{L}{2} N \sqrt{2t^2 M_1^2 + 4\lambda t \mathbb{E} \|\boldsymbol{\xi}\|^2} \\ &\quad + \lambda M_2 N \sqrt{2t^2 M_1^2 + 4\lambda t \mathbb{E} \|\boldsymbol{\xi}\|^2} + \lambda M_2 \left( 2Nt^2 M_1^2 + 4\lambda t N \mathbb{E} \|\boldsymbol{\xi}\|^2 \right)^{1/2} \\ &= \left( LN + \lambda M_2 N + \lambda M_2 \sqrt{N} \right) \sqrt{2t^2 M_1^2 + 4\lambda t \mathbb{E} \|\boldsymbol{\xi}\|^2}. \end{aligned} \quad (38)$$

where we used Jensen's inequality to bound  $\mathbb{E}_{\nu_{0t}} (\sum_{i=1}^N \|\mathbf{y}^i - \mathbf{y}_0^i\|^2)^{1/2} \leq (\mathbb{E}_{\nu_{0t}} \sum_{i=1}^N \|\mathbf{y}^i - \mathbf{y}_0^i\|^2)^{1/2}$ , and Lemma 6 in the last inequality.

We now bound the term involving  $r_t^{\text{sc}}$ . By the tower property and Jensen's inequality, we have

$$\begin{aligned}
 & \int \nu_t(\mathbf{Y}) \|r_t^{\text{sc}}(\mathbf{Y})\|_*^2 d\mathbf{Y} \\
 &= \mathbb{E}_{\nu_t} \left\| \mathbb{E}[\lambda(\nabla \log \nu_t(\mathbf{Y}) - \nabla \log \nu_0(\mathbf{Y}_0)) \mid \mathbf{Y}] \right\|_*^2 \\
 &\leq \lambda^2 \mathbb{E}_{(\mathbf{Y}_0, \mathbf{Y}) \sim \nu_{0t}} \left[ \|\nabla \log \nu_t(\mathbf{Y}) - \nabla \log \nu_0(\mathbf{Y}_0)\|_*^2 \right] \quad (\text{by Jensen}) \\
 &\leq 2\lambda^2 M_2^2 \mathbb{E}_{\nu_{0t}} \|\mathbf{Y} - \mathbf{Y}_0\|^2 + 2\lambda^2 M_2^2 W_2^2(\nu_t, \nu_0) \quad (\text{by Asm 7 and } (a+b)^2 \leq 2a^2 + 2b^2) \\
 &\leq 4\lambda^2 M_2^2 \mathbb{E}_{\nu_{0t}} \|\mathbf{Y} - \mathbf{Y}_0\|^2 \quad (\text{using the same coupling as before}) \\
 &\leq 4\lambda^2 M_2^2 (2Nt^2 M_1^2 + 4\lambda t N \mathbb{E}\|\boldsymbol{\xi}\|^2) \quad (\text{by Lemma 6}).
 \end{aligned} \tag{39}$$

Next, using Hölder's inequality (Proposition 2.10 in Kreuter [24]), we obtain

$$\begin{aligned}
 & - \int \nu_t \langle G_t, r_t^{\text{sc}} \rangle \\
 &\leq \left( \int \nu_t \|G_t\|_*^2 \right)^{1/2} \left( \int \nu_t \|r_t^{\text{sc}}\|_*^2 \right)^{1/2} \\
 &\leq \sqrt{C} \left( \int \nu_t \|G_t\|_*^2 \right)^{1/2} \left( \int \nu_t \|r_t^{\text{sc}}\|_*^2 \right)^{1/2} \quad (\text{by norm equivalence } \|\mathbf{a}\| \leq \sqrt{C}\|\mathbf{a}\|_*) \\
 &\leq \left( \int \nu_t \|G_t\|_*^2 \right)^{1/2} \cdot \sqrt{C} \cdot 2\lambda M_2 (2Nt^2 M_1^2 + 4\lambda t N \mathbb{E}\|\boldsymbol{\xi}\|^2)^{1/2} \quad (\text{by Eq. (39)}).
 \end{aligned} \tag{40}$$

Thus, define

$$\begin{aligned}
 \beta_t &:= 2\lambda M_2 \sqrt{CN} \sqrt{2t^2 M_1^2 + 4\lambda t \mathbb{E}\|\boldsymbol{\xi}\|^2}. \\
 a_t &:= 2 \left( LN + \lambda M_2 N + \lambda M_2 \sqrt{N} \right) \sqrt{2t^2 M_1^2 + 4\lambda t \mathbb{E}\|\boldsymbol{\xi}\|^2}.
 \end{aligned} \tag{41}$$

Putting together the bounds for the three terms in Eq. (32) we obtain, for all  $t \in [0, \eta]$ ,

$$\begin{aligned}
 \frac{d}{dt} \mathcal{F}^{(N)}(\nu_t) - \mathcal{F}^{(N)}(\mu_*^{(N)}) &\leq - \int \nu_t \|G_t\|_* + \beta_t \left( \int \nu_t \|G_t\|_*^2 \right)^{1/2} + a_t \\
 &\leq -(\kappa - \beta_t) \left( \int \nu_t \|G_t\|_*^2 \right)^{1/2} + a_t \\
 &\leq -(\kappa - \beta_t) \sqrt{2\zeta\lambda} \sqrt{\mathcal{F}^{(N)}(\nu_t) - \mathcal{F}^{(N)}(\mu_*^{(N)})} + a_t.
 \end{aligned} \tag{42}$$

where we use Asm 4 in the second step, log-sobolev inequality Asm 3 and Lemma 2 in the third step when  $\beta_t \leq \kappa/2$  (which will be satisfied by the following argument).

Next, we choose  $\eta$  so that  $\beta_t \leq \kappa/2$  on  $t \in [0, \eta]$ . Since  $t \mapsto 2t^2 M_1^2 + 4\lambda t \mathbb{E}\|\boldsymbol{\xi}\|^2$  is increasing on  $[0, \infty)$ , we have  $\sup_{t \in [0, \eta]} \beta_t = \beta_\eta$ . For  $0 < \eta \leq 1$ , we bound

$$2t^2 M_1^2 + 4\lambda t \mathbb{E}\|\boldsymbol{\xi}\|^2 \leq 2\eta M_1^2 + 4\lambda \eta \mathbb{E}\|\boldsymbol{\xi}\|^2 = 2\eta(M_1^2 + 2\lambda \mathbb{E}\|\boldsymbol{\xi}\|^2),$$

and therefore, by Eq. (41),

$$\beta_\eta \leq 2\lambda M_2 \sqrt{CN} \sqrt{2\eta(M_1^2 + 2\lambda \mathbb{E}\|\xi\|^2)}. \quad (43)$$

Hence it suffices to choose  $\eta > 0$  such that

$$2\lambda M_2 \sqrt{CN} \sqrt{2\eta(M_1^2 + 2\lambda \mathbb{E}\|\xi\|^2)} \leq \frac{\kappa}{2},$$

Define

$$\eta_0 := \min \left\{ 1, \frac{\kappa^2}{32\lambda^2 M_2^2 C N (M_1^2 + 2\lambda \mathbb{E}\|\xi\|^2)} \right\}. \quad (44)$$

Then for all  $t \in [0, \eta]$  whenever  $0 < \eta \leq \eta_0$ , we have  $\beta_t \leq \kappa/2$  and thus  $\kappa - \beta_t \geq \kappa/2$ . Consequently, Eq. (42) simplifies to

$$\frac{d}{dt} \mathcal{F}^{(N)}(\nu_t) - \mathcal{F}^{(N)}(\mu_*^{(N)}) \leq -a \sqrt{\mathcal{F}^{(N)}(\nu_t) - \mathcal{F}^{(N)}(\mu_*^{(N)})} + a_t, \quad a := \frac{\kappa}{2} \sqrt{2\zeta\lambda} \quad (45)$$

Lastly, we convert this differential inequality into a per-iteration recursion. Fix an iteration  $k$  and consider the interpolation on  $t \in [0, \eta]$ . For notational convenience, set

$$\mathcal{E}(t) := \mathcal{F}^{(N)}(\nu_t) - \mathcal{F}^{(N)}(\mu_*^{(N)}), \quad \bar{a}_\eta := \sup_{t \in [0, \eta]} a_t.$$

From the previous step, we have for all  $t \in [0, \eta]$ ,

$$\mathcal{E}'(t) \leq -a \sqrt{\mathcal{E}(t)} + a_t, \quad a := \frac{\kappa}{2} \sqrt{2\zeta\lambda}. \quad (46)$$

Define  $y(t) := \sqrt{\mathcal{E}(t)}$ . Whenever  $y(t) > 0$ ,

$$2y(t)y'(t) = \mathcal{E}'(t) \leq -ay(t) + a_t \leq -ay(t) + \bar{a}_\eta.$$

Hence, for all  $t$  with  $y(t) > 0$ ,

$$y'(t) \leq -\frac{a}{2} + \frac{\bar{a}_\eta}{2y(t)}. \quad (47)$$

Set the threshold  $R := \frac{2\bar{a}_\eta}{a}$ . We distinguish two cases.

*Case 1:*  $y(t) \geq R$  for all  $t \in [0, \eta]$ . Then Eq. (47) gives

$$y'(t) \leq -\frac{a}{2} + \frac{a}{4} = -\frac{a}{4}, \quad t \in [0, \eta].$$

Integrating over  $[0, \eta]$  yields

$$y(\eta) \leq y(0) - \frac{a}{4}\eta.$$

*Case 2:*  $y(t_0) < R$  for some  $t_0 \in [0, \eta]$ . Let  $t_0$  be the first such time (if  $y(0) < R$ , take  $t_0 = 0$ ). We claim that  $y(t) \leq R$  for all  $t \in [t_0, \eta]$ . Indeed, suppose for contradiction that there exists  $t_1 \in (t_0, \eta]$  with  $y(t_1) > R$ . By continuity, define

$$s := \max\{t \in [t_0, t_1] : y(t) = R\}.$$

Then  $y(s) = R$  and  $y(t) > R$  for all  $t \in (s, t_1]$ . But whenever  $y(t) \geq R$ , Eq. (47) and the definition of  $R$  imply

$$y'(t) \leq -\frac{a}{2} + \frac{\bar{a}_\eta}{2y(t)} \leq -\frac{a}{2} + \frac{\bar{a}_\eta}{2R} = -\frac{a}{2} + \frac{a}{4} = -\frac{a}{4} < 0,$$

so  $y$  is strictly decreasing on  $(s, t_1]$ , contradicting  $y(t) > R$  there. Therefore  $y(t) \leq R$  on  $[t_0, \eta]$ , and in particular  $y(\eta) \leq R$ .

Combining the two cases, we obtain the one-step bound

$$y(\eta) \leq \max \left\{ y(0) - \frac{a}{4}\eta, \frac{2\bar{a}_\eta}{a} \right\}.$$

Recalling  $y(t) = \sqrt{\mathcal{E}(t)}$ , we get

$$\sqrt{\mathcal{F}^{(N)}(\mu_{k+1}^{(N)}) - \mathcal{F}^{(N)}(\mu_*^{(N)})} \leq \max \left\{ \sqrt{\mathcal{F}^{(N)}(\mu_k^{(N)}) - \mathcal{F}^{(N)}(\mu_*^{(N)})} - \frac{a}{4}\eta, \frac{2\bar{a}_\eta}{a} \right\}. \quad (48)$$

Iterating Eq. (48) yields, for all  $k \geq 0$ ,

$$\mathcal{F}^{(N)}(\mu_k^{(N)}) - \mathcal{F}^{(N)}(\mu_*^{(N)}) \leq \left[ \max \left\{ \sqrt{\mathcal{E}_0^{(N)}} - \frac{a}{4}k\eta, \frac{2\bar{a}_\eta}{a} \right\} \right]^2. \quad (49)$$

Finally, add and subtract  $\mathcal{F}^{(N)}(\mu_*^{(N)})$ :

$$\frac{1}{N} \mathcal{F}^{(N)}(\mu_k^{(N)}) - \mathcal{F}(\mu^*) = \frac{1}{N} \left( \mathcal{F}^{(N)}(\mu_k^{(N)}) - \mathcal{F}^{(N)}(\mu_*^{(N)}) \right) + \left( \frac{1}{N} \mathcal{F}^{(N)}(\mu_*^{(N)}) - \mathcal{F}(\mu^*) \right). \quad (50)$$

The finite-particle bias term is controlled by Lemma 3:

$$\frac{1}{N} \mathcal{F}^{(N)}(\mu_*^{(N)}) - \mathcal{F}(\mu^*) \leq \frac{\hat{L}\hat{R}^2}{2N}. \quad (51)$$

Combining Eq. (49)–Eq. (51) gives, for all  $k \geq 0$ ,

$$\begin{aligned} & \frac{1}{N} \mathcal{F}^{(N)}(\mu_k^{(N)}) - \mathcal{F}(\mu^*) \\ & \leq \frac{1}{N} \left[ \max \left\{ \sqrt{\mathcal{F}^{(N)}(\mu_0^{(N)}) - \mathcal{F}^{(N)}(\mu_*^{(N)})} - \frac{a}{4}k\eta, \frac{2\bar{a}_\eta}{a} \right\} \right]^2 + \frac{\hat{L}\hat{R}^2}{2N}, \quad a = \frac{\kappa}{2} \sqrt{2\zeta\lambda}. \end{aligned} \quad (52)$$

■

## Appendix F. LMO gradient flow in measure space in noiseless setting ( $\lambda = 0$ )

**Theorem 11 (Convergence of normalized LMO Wasserstein flow)** *Consider the following LMO Wasserstein flow:*

$$\partial_t \mu_t = \nabla \cdot (\mu_t v_t), \quad v_t := -\text{lmo} \left( \nabla \frac{\delta F}{\delta \mu}(\mu_t)(\mathbf{x}); 0 \right),$$

we assume the solution  $(\mu_t)_{t \geq 0}$  exists. Then under Asm. 4 and 8, for all  $t \geq 0$ ,

$$F(\mu_t) - F^* \leq \left[ \max \left\{ \sqrt{F(\mu_0) - F^*} - \kappa \sqrt{\frac{\alpha}{2}} t, 0 \right\} \right]^2.$$

In particular,  $F(\mu_t) = F^*$  for all  $t \geq t^*$ , where  $t^* = \frac{1}{\kappa} \sqrt{\frac{2}{\alpha}} \sqrt{F(\mu_0) - F^*}$ .

### F.1. Proof of Thm. 11

**Proof** Let  $\Delta(t) := F(\mu_t) - F^* \geq 0$ ,  $g_{\mu_t}(x) := \nabla_x \frac{\delta F}{\delta \mu}(\mu_t)(x)$ . By the chain rule for first variations along continuity equations,

$$\frac{d}{dt} F(\mu_t) = \int \langle g_{\mu_t}(x), v_t(x) \rangle d\mu_t(x).$$

Plugging  $v_t(x)$  and using Lemma 7, we get

$$\Delta'(t) = \frac{d}{dt} F(\mu_t) = -\|g_{\mu_t}\|_{L^1(\mu_t; \|\cdot\|_*)}. \quad (53)$$

By the non-degeneracy assumption Asm 4,

$$\|g_{\mu_t}\|_{L^1(\mu_t; \|\cdot\|_*)} \geq \kappa \|g_{\mu_t}\|_{L^2(\mu_t; \|\cdot\|_*)}. \quad (54)$$

Finally, by the transport PL inequality (Asm 8),

$$\|g_{\mu_t}\|_{L^2(\mu_t; \|\cdot\|_*)} \geq \sqrt{2\alpha \Delta(t)}. \quad (55)$$

Combining Eq. (53)–Eq. (55) yields the differential inequality

$$\Delta'(t) \leq -\kappa \sqrt{2\alpha} \sqrt{\Delta(t)}.$$

Let  $y(t) := \sqrt{\Delta(t)}$ . Wherever  $y(t) > 0$ , we have  $y'(t) = \Delta'(t)/(2\sqrt{\Delta(t)}) \leq -\kappa\sqrt{\alpha/2}$ . Integrating gives

$$y(t) \leq \max \left\{ y(0) - \kappa\sqrt{\alpha/2}t, 0 \right\},$$

and squaring yields the convergence rate given in Thm. 11. The expression for  $t^*$  followed by setting the bracketed term to zero.  $\blacksquare$

### F.2. Proof of Lemma 12

**Proof** Fix  $\mu$  and  $T$ , and define the displacement interpolation

$$\mu_t := (\text{id} + tv)_{\#}\mu, \quad t \in [0, 1],$$

where  $v(\mathbf{x}) := T(\mathbf{x}) - \mathbf{x}$ . Let  $h(t) := F(\mu_t)$ . Lemma A.2 in Chizat [12] provides the derivative formula

$$h'(t) = \int_{\mathbb{R}^d} \left\langle \nabla_{\mathbf{x}} \frac{\delta F}{\delta \mu}(\mu_t)(\mathbf{x} + tv(\mathbf{x})), v(\mathbf{x}) \right\rangle d\mu(\mathbf{x}). \quad (56)$$

Fix  $s, t \in [0, 1]$  and define  $G_{t,s}(\mathbf{x}) := \nabla_{\mathbf{x}} \frac{\delta F}{\delta \mu}(\mu_t)(\mathbf{x} + tv(\mathbf{x})) - \nabla_{\mathbf{x}} \frac{\delta F}{\delta \mu}(\mu_s)(\mathbf{x} + sv(\mathbf{x}))$ . Then, subtracting Eq. (56) at times  $t$  and  $s$ , we obtain

$$h'(t) - h'(s) = \int \langle G_{t,s}(\mathbf{x}), v(\mathbf{x}) \rangle d\mu(\mathbf{x}).$$

By  $|\langle u, v \rangle| \leq \|u\|_* \|v\|$  and Hölder's inequality applied to the scalar functions  $\|G_{t,s}(\cdot)\|_*$  and  $\|v(\cdot)\|$  (Proposition 2.10 in Kreuter [24]), we obtain

$$|h'(t) - h'(s)| \leq \|G_{t,s}\|_{L^2(\mu; \|\cdot\|_*)} \|v\|_{L^2(\mu; \|\cdot\|)}. \quad (57)$$

Next we bound  $\|G_{t,s}\|_{L^2(\mu; \|\cdot\|_*)}$ . By Asm 2:

$$\begin{aligned} \|G_{t,s}(\mathbf{x})\|_* &= \left\| \nabla_{\mathbf{x}} \frac{\delta F}{\delta \mu}(\mu_t)(\mathbf{x} + tv(\mathbf{x})) - \nabla_{\mathbf{x}} \frac{\delta F}{\delta \mu}(\mu_s)(\mathbf{x} + sv(\mathbf{x})) \right\|_* \\ &\leq \frac{L}{2} (W_2(\mu_t, \mu_s) + |t - s| \|v(\mathbf{x})\|). \end{aligned} \quad (58)$$

Recall that  $W_2$  is defined using the same norm  $\|\cdot\|$ : i.e.,

$$W_2^2(\mu_t, \mu_s) = \inf_{\pi \in \Pi(\mu_t, \mu_s)} \int \|\mathbf{x} - \mathbf{y}\|^2 d\pi(\mathbf{x}, \mathbf{y}).$$

Consider the measurable map  $\Phi_{t,s}(x) := (x + tv(x), x + sv(x))$  and define  $\pi_{t,s} := (\Phi_{t,s})_{\#}\mu$ . Then, for any Borel set  $A \subset \mathbb{R}^d$ , we have

$$\pi_{t,s}(A \times \mathbb{R}^d) = \mu(\{x : x + tv(x) \in A\}) = (\text{id} + tv)_{\#}\mu(A) = \mu_t(A),$$

and similarly,

$$\pi_{t,s}(\mathbb{R}^d \times A) = \mu(\{x : x + sv(x) \in A\}) = (\text{id} + sv)_{\#}\mu(A) = \mu_s(A).$$

Hence  $\pi_{t,s} \in \Pi(\mu_t, \mu_s)$ .

By the definition of the Wasserstein distance:

$$W_2^2(\mu_t, \mu_s) \leq \int \left\| (\mathbf{x} + tv(\mathbf{x})) - (\mathbf{x} + sv(\mathbf{x})) \right\|^2 d\mu(\mathbf{x}) = |t - s|^2 \int \|v(\mathbf{x})\|^2 d\mu(\mathbf{x}).$$

Therefore, we get:

$$\|G_{t,s}(\mathbf{x})\|_* \leq \frac{L}{2} |t - s| (\|v\|_{L^2(\mu; \|\cdot\|)} + \|v(\mathbf{x})\|).$$

Squaring and using  $(a + b)^2 \leq 2a^2 + 2b^2$ :

$$\|G_{t,s}(\mathbf{x})\|_*^2 \leq \frac{L^2}{2} |t - s|^2 \left( \|v\|_{L^2(\mu; \|\cdot\|)}^2 + \|v(\mathbf{x})\|^2 \right).$$

Integrating over  $\mu$  yields:

$$\begin{aligned} \|G_{t,s}\|_{L^2(\mu; \|\cdot\|_*)}^2 &= \int \|G_{t,s}(\mathbf{x})\|_*^2 d\mu(\mathbf{x}) \\ &\leq \frac{L^2}{2} |t - s|^2 \left( \|v\|_{L^2(\mu; \|\cdot\|)}^2 + \int \|v(\mathbf{x})\|^2 d\mu(\mathbf{x}) \right) \\ &= L^2 |t - s|^2 \|v\|_{L^2(\mu; \|\cdot\|)}^2. \end{aligned}$$

Plugging back Eq. (57), we obtain:

$$|h'(t) - h'(s)| \leq L |t - s| \|v\|_{L^2(\mu; \|\cdot\|)}^2. \quad (59)$$

Taking  $s = 0$  in Eq. (59) and integrating over  $t \in [0, 1]$ :

$$h(1) = h(0) + \int_0^1 h'(t) dt \leq h(0) + h'(0) + \int_0^1 |h'(t) - h'(0)| dt \leq h(0) + h'(0) + \frac{L}{2} \|v\|_{L^2(\mu; \|\cdot\|)}^2.$$

Finally, by definition  $h(0) = F(\mu)$  and  $h(1) = F(T_{\#}\mu)$ , and Eq. (56), substituting these identities into the previous inequality gives exactly Eq. (85), completing the proof.  $\blacksquare$

### F.3. Auxiliary lemma

**Lemma 7 (LMO identities)** Let  $\|\cdot\|$  be a norm on  $\mathbb{R}^d$  with dual norm  $\|\cdot\|_*$ , and let  $\text{lmo}(\cdot; \tau)$  be defined as in Thm. 1. Then, for every  $\mathbf{g} \in \mathbb{R}^d$ ,

$$\langle \mathbf{g}, \text{lmo}(\mathbf{g}; \tau) \rangle = -\|\mathbf{g}\|_*^{\tau+1}, \quad \|\text{lmo}(\mathbf{g}; \tau)\| = \|\mathbf{g}\|_*^\tau.$$

**Proof** If  $\mathbf{g} = 0$ , then  $\|\mathbf{g}\|_* = 0$  and any  $\mathbf{y}$  with  $\|\mathbf{y}\| \leq 0$  must be  $\mathbf{y} = 0$ , so  $\text{lmo}(0; \tau) = 0$  and both identities hold trivially. Hence we may assume  $\mathbf{g} \neq 0$ .

By Definition 1, we have

$$\text{lmo}(\mathbf{g}; \tau) = -\|\mathbf{g}\|_*^\tau s(\mathbf{g}), \quad \text{where } s(\mathbf{g}) \in \arg \max_{\|\mathbf{u}\| \leq 1} \langle \mathbf{g}, \mathbf{u} \rangle \text{ and } \langle \mathbf{g}, s(\mathbf{g}) \rangle = \|\mathbf{g}\|_*.$$

Since the maximum of a linear functional over the unit ball is attained on the unit sphere, we can choose  $s(\mathbf{g})$  with  $\|s(\mathbf{g})\| = 1$ .

Then

$$\langle \mathbf{g}, \text{lmo}(\mathbf{g}; \tau) \rangle = \langle \mathbf{g}, -\|\mathbf{g}\|_*^\tau s(\mathbf{g}) \rangle = -\|\mathbf{g}\|_*^\tau \langle \mathbf{g}, s(\mathbf{g}) \rangle = -\|\mathbf{g}\|_*^\tau \|\mathbf{g}\|_* = -\|\mathbf{g}\|_*^{\tau+1},$$

and

$$\|\text{lmo}(\mathbf{g}; \tau)\| = \|\|\mathbf{g}\|_*^\tau s(\mathbf{g})\| = \|\mathbf{g}\|_*^\tau \|s(\mathbf{g})\| = \|\mathbf{g}\|_*^\tau.$$

This proves the claim. ■

**Lemma 8 (Upper bound for  $\widehat{\Delta}_K$  in Eq. (73))** The term  $\widehat{\Delta}_K$  defined in Eq. (73) can be bounded as follows:

$$\widehat{\Delta}_K = \max_{0 \leq t < K} \Delta_t \leq \Delta_0 + \frac{L}{2} \eta^2 \log\left(\frac{S_{K-1}}{S_{\text{init}}}\right). \quad (60)$$

**Proof** From the smoothness inequality (Eq. (70), dropping the negative term  $-\eta_k \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_*^{1+\tau}$ ) we have

$$F(\mu_{k+1}) - F(\mu_k) \leq \frac{L}{2} \eta_k^2 \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_*^{2\tau}.$$

Summing over  $s = 0, \dots, k-1$  yields

$$F(\mu_k) - F(\mu_0) \leq \frac{L}{2} \sum_{s=0}^{k-1} \eta_s^2 \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_s^i\|_*^{2\tau} \leq \frac{L}{2} \eta^2 \log\left(\frac{S_{k-1}}{S_{\text{init}}}\right) \leq \frac{L}{2} \eta^2 \log\left(\frac{S_{K-1}}{S_{\text{init}}}\right),$$

for all  $k < K$ , where we use Lemma 9. Hence

$$\Delta_k = F(\mu_k) - F^* \leq \Delta_0 + \frac{L}{2} \eta^2 \log\left(\frac{S_{K-1}}{S_{\text{init}}}\right), \quad 0 \leq k < K,$$

and therefore

$$\widehat{\Delta}_K = \max_{0 \leq k < K} \Delta_k \leq \Delta_0 + \frac{L}{2} \eta^2 \log\left(\frac{S_{K-1}}{S_{\text{init}}}\right). \quad (61)$$

■

**Lemma 9 (From Wang et al. [47] (Lemma 10))** Let  $\{a_k\}_{k=0}^{\infty}$  be a sequence of non-negative real numbers with  $a_0 > 0$ . Then the following inequality holds:

$$\sum_{k=1}^K \frac{a_k}{\sum_{s=0}^k a_s} \leq \ln \left( \sum_{k=0}^K a_k \right) - \ln a_0.$$

**Lemma 10 (From Attia and Koren [3] (Lemma 4))** Let  $\{a_k\}_{k=0}^{\infty}$  be a sequence of non-negative real numbers. Then the following inequality holds:

$$\sum_{k=1}^K \frac{a_k}{\sqrt{\sum_{s=0}^k a_s}} \leq 2 \sqrt{\sum_{k=1}^K a_k}$$

**Lemma 11 (Polynomial growth of the accumulator  $S_{K-1}$  for  $\tau = 1$ )** Under the same setting and assumption as in Thm. 5, We have:

$$\log \left( \frac{S_{K-1}}{S_{\text{init}}} \right) \leq \log \left( 1 + \frac{4KL(F(\mu_0) - F^*) + \frac{1}{2}(L(\sqrt{N} + 1)\eta)^2 K^3}{S_{\text{init}}} \right). \quad (62)$$

**Proof** We work in the  $\tau = 1$  case. Recall that

$$\mathbf{g}_k^i := \nabla_{\mathbf{x}} \frac{\delta F}{\delta \mu}(\mu_k)(X_k^i), \quad \mathbf{v}_k^i := \text{lmo}(\mathbf{g}_k^i; 1), \quad \|\mathbf{v}_k^i\| = \|\mathbf{g}_k^i\|_*.$$

By using Asm 2, we have

$$\|\mathbf{g}_k^i - \mathbf{g}_{k-1}^i\|_* \leq \frac{L}{2} \left( W_2(\mu_k, \mu_{k-1}) + \|\mathbf{x}_k^i - \mathbf{x}_{k-1}^i\| \right).$$

Using the update  $\mathbf{x}_k^i = \mathbf{x}_{k-1}^i + \eta_{k-1} \mathbf{v}_{k-1}^i$  and the coupling that pairs each  $\mathbf{x}_{k-1}^j$  with  $\mathbf{x}_k^j$ , we obtain

$$\|\mathbf{x}_k^i - \mathbf{x}_{k-1}^i\| = \eta_{k-1} \|\mathbf{v}_{k-1}^i\|, \quad W_2(\mu_k, \mu_{k-1}) \leq \eta_{k-1} \sqrt{B_{k-1}},$$

where  $B_{k-1} = \frac{1}{N} \sum_{j=1}^N \|\mathbf{v}_{k-1}^j\|^2$ . Hence

$$\begin{aligned} \|\mathbf{g}_k^i\|_* &\leq \|\mathbf{g}_{k-1}^i\|_* + \frac{L}{2} \eta_{k-1} (\sqrt{B_{k-1}} + \|\mathbf{v}_{k-1}^i\|) \\ &\leq \|\mathbf{g}_{k-1}^i\|_* + \frac{L}{2} \eta_{k-1} (\sqrt{B_{k-1}} + \sqrt{\sum_{j=1}^N \|\mathbf{v}_{k-1}^j\|^2}) \\ &= \|\mathbf{g}_{k-1}^i\|_* + \frac{L}{2} \eta_{k-1} (\sqrt{N} + 1) \sqrt{B_{k-1}}. \end{aligned} \quad (63)$$

By the definition of the global adaptive stepsizes,

$$\eta_{k-1} = \frac{\eta}{\sqrt{S_{k-1}}}, \quad S_{k-1} = S_{\text{init}} + \sum_{s=0}^{k-1} B_s,$$

so  $\eta_{k-1}\sqrt{B_{k-1}} \leq \eta$ . Consequently,

$$\|\mathbf{g}_k^i\|_* \leq \|\mathbf{g}_{k-1}^i\|_* + \frac{L}{2}(\sqrt{N} + 1)\eta, \quad (64)$$

Iterating Eq. (64) gives, for all  $k \geq 1$ ,

$$\|\mathbf{g}_k^i\|_* \leq \|\mathbf{g}_0^i\|_* + k \frac{L}{2}(\sqrt{N} + 1)\eta. \quad (65)$$

Squaring Eq. (65) and using  $(a + b)^2 \leq 2a^2 + 2b^2$  yields

$$\|\mathbf{g}_k^i\|_*^2 \leq 2\|\mathbf{g}_0^i\|_*^2 + \frac{1}{2}(L(\sqrt{N} + 1)\eta)^2 k^2.$$

Thus, we have

$$B_k = \frac{1}{N} \sum_{j=1}^N \|v_k^j\|^2 = \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_*^2 \leq 2B_0 + \frac{1}{2}(L(\sqrt{N} + 1)\eta)^2 k^2, \quad B_0 := \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_0^i\|_*^2.$$

Therefore

$$\begin{aligned} S_{K-1} &= S_{\text{init}} + \sum_{k=0}^{K-1} B_k \\ &\leq S_{\text{init}} + \sum_{k=0}^{K-1} \left(2B_0 + \frac{1}{2}(L(\sqrt{N} + 1)\eta)^2 k^2\right) \\ &= S_{\text{init}} + 2KB_0 + \frac{1}{2}(L(\sqrt{N} + 1)\eta)^2 \sum_{k=0}^{K-1} k^2 \\ &\leq S_{\text{init}} + 2KB_0 + \frac{1}{2}(L(\sqrt{N} + 1)\eta)^2 K^3. \end{aligned} \quad (66)$$

Next, we will bound the term  $B_0 = \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_0^i\|_*^2$ . Continuing from Eq. (70) but with a fix step size  $\hat{\eta}$ , and  $k = 0$ , we have

$$F(\mu_1) - F(\mu_0) \leq -\hat{\eta} \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_0^i\|_*^2 + \frac{L}{2} \hat{\eta}^2 \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_0^i\|_*^2. \quad (67)$$

By setting  $\hat{\eta} = \frac{1}{L}$ , we get

$$B_0 \leq 2L(F(\mu_0) - F^*)$$

Finally, dividing Eq. (66) by  $S_{\text{init}}$ , applying the monotonicity of log yields Eq. (62).  $\blacksquare$

#### F.4. Global adaptive algorithm with unnormalized LMO particle descent

**Theorem 12** Consider the same setting as in Thm. 5 with the case  $\tau = 1$ , we have the following convergence guarantee:

$$\begin{aligned} \min_{0 \leq k < K} \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_* &\leq \frac{\sqrt{2\Delta_0} \sqrt[4]{S_{init}}}{\sqrt{K}\eta} + \frac{\sqrt{2}L\eta}{\sqrt{K}} + \frac{2\sqrt{2}}{\eta\sqrt{K}}\Delta_0 \\ &+ \frac{2\sqrt{2}}{\eta\sqrt{K}} \left( \frac{L}{2} \eta^2 \log \left( 1 + \frac{4KL\Delta_0 + \frac{1}{2}(L(\sqrt{N}+1)\eta)^2 K^3}{S_{init}} \right) \right). \end{aligned}$$

**Remark:** The proof can be found at Sec. F.5 As a comparison: i) When  $\tau = 1$ , the right-hand side contains a logarithmic dependence on the number of particles  $N$ , whereas this dependence disappears when  $\tau = 0$ . ii) When  $\tau = 0$ , we have  $\frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_*^{2\tau} = 1$  by Lemma 7. Thus, the step size  $\eta_k = \frac{\eta}{\sqrt{k}}$  does not depend on the past gradient. iii) The first term in the right-hand side of the bound is  $\mathcal{O}(K^{-1/2})$  for  $\tau = 1$  and  $\mathcal{O}(K^{-1})$  for  $\tau = 0$ .

#### F.5. Proof of Thm. 5 and Thm. 12

##### F.5.1. PROOF

**Proof** [Proof of Thm. 5] At iteration  $k$ , the LMO-APGD update can be written as a pushforward

$$\mu_{k+1} = (T_k)_\# \mu_k, \quad T_k(\mathbf{x}) = \mathbf{x} + \eta_k \mathbf{v}_k(\mathbf{x}),$$

where  $\mathbf{v}_k(\mathbf{x}) = \mathbf{v}_k^i$  whenever  $\mathbf{x} = \mathbf{x}_k^i$ . Define  $g_k(\mathbf{x}) := \nabla_{\mathbf{x}} \frac{\delta F}{\delta \mu}(\mu_k)(\mathbf{x})$  and  $g_k(\mathbf{x}) = \mathbf{g}_k^i$  whenever  $\mathbf{x} = \mathbf{x}_k^i$ . By Asm 2, we have Eq. (85), by applying with  $T_k$ , we obtain

$$F(\mu_{k+1}) - F(\mu_k) \leq \int \langle g_k(\mathbf{x}), T_k(\mathbf{x}) - \mathbf{x} \rangle d\mu_k(\mathbf{x}) + \frac{L}{2} \int \|T_k(\mathbf{x}) - \mathbf{x}\|^2 d\mu_k(\mathbf{x}) \quad (68)$$

$$= \int \langle g_k(\mathbf{x}), \eta_k \mathbf{v}_k(\mathbf{x}) \rangle d\mu_k(\mathbf{x}) + \frac{L}{2} \int \|\eta_k \mathbf{v}_k(\mathbf{x})\|^2 d\mu_k(\mathbf{x}). \quad (69)$$

Since  $\mu_k$  is the empirical measure  $\mu_k = \frac{1}{N} \sum_{i=1}^N \delta_{\mathbf{x}_k^i}$ , we have

$$\int \langle g_k(\mathbf{x}), \eta_k \mathbf{v}_k(\mathbf{x}) \rangle d\mu_k(\mathbf{x}) = \eta_k \frac{1}{N} \sum_{i=1}^N \langle \mathbf{g}_k^i, \mathbf{v}_k^i \rangle,$$

$$\int \|\eta_k \mathbf{v}_k(\mathbf{x})\|^2 d\mu_k(\mathbf{x}) = \eta_k^2 \frac{1}{N} \sum_{i=1}^N \|\mathbf{v}_k^i\|^2.$$

By the LMO identities (Lemma 7) for  $\tau \in \{0, 1\}$ , we have

$$\langle \mathbf{g}_k^i, \mathbf{v}_k^i \rangle = -\|\mathbf{g}_k^i\|_*^{1+\tau}, \quad \|\mathbf{v}_k^i\|^2 = \|\mathbf{g}_k^i\|_*^{2\tau},$$

and therefore

$$F(\mu_{k+1}) - F(\mu_k) \leq -\eta_k \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_*^{1+\tau} + \frac{L}{2} \eta_k^2 \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_*^{2\tau}. \quad (70)$$

Rearranging yields

$$\eta_k \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_*^{1+\tau} \leq F(\mu_k) - F(\mu_{k+1}) + \frac{L}{2} \eta_k^2 \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_*^{2\tau}. \quad (71)$$

Divide both sides of Eq. (71) by  $K\eta_k$  and then summing up over  $k = 0, \dots, K-1$ :

$$\frac{1}{K} \sum_{k=0}^{K-1} \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_*^{1+\tau} \leq \underbrace{\frac{1}{K} \sum_{k=0}^{K-1} \frac{F(\mu_k) - F(\mu_{k+1})}{\eta_k}}_{(\star)} + \frac{L}{2K} \sum_{k=0}^{K-1} \eta_k \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_*^{2\tau}. \quad (72)$$

We now bound the term  $(\star)$ . Define

$$\Delta_k := F(\mu_k) - F^*, \quad \widehat{\Delta}_K := \max_{0 \leq j < K} \Delta_j, \quad (73)$$

We set  $\Delta_{-1} := \Delta_0$  and  $\eta_{-1} := \frac{\eta}{\sqrt{S_{\text{init}}}}$  for convenience. Then

$$\begin{aligned} (\star) &= \frac{1}{K} \sum_{k=0}^{K-1} \frac{\Delta_k - \Delta_{k+1}}{\eta_k} \\ &= \frac{1}{K} \sum_{k=-1}^{K-1} \frac{\Delta_k - \Delta_{k+1}}{\eta_k} \quad (\text{the } k = -1 \text{ term vanishes since } \Delta_{-1} = \Delta_0) \\ &= \frac{1}{K} \left( \frac{\Delta_{-1}}{\eta_{-1}} + \sum_{k=0}^{K-1} \Delta_k \left( \frac{1}{\eta_k} - \frac{1}{\eta_{k-1}} \right) - \frac{\Delta_K}{\eta_{K-1}} \right) \\ &\leq \frac{\Delta_0}{K\eta_{-1}} + \frac{\widehat{\Delta}_K}{K} \sum_{k=0}^{K-1} \left( \frac{1}{\eta_k} - \frac{1}{\eta_{k-1}} \right), \end{aligned}$$

where we used  $\Delta_k \leq \widehat{\Delta}_K$ , and  $\eta_k$  is a non-increasing sequence, and dropped the negative term.

Recall that  $\eta_k = \eta/\sqrt{S_k}$  and by Lemma 7 we have

$$S_k = S_{k-1} + B_k, \quad B_k := \frac{1}{N} \sum_{i=1}^N \|\mathbf{v}_k^i\|^2 = \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_*^{2\tau}.$$

Then

$$\frac{1}{\eta_k} - \frac{1}{\eta_{k-1}} = \frac{\sqrt{S_k} - \sqrt{S_{k-1}}}{\eta} = \frac{S_k - S_{k-1}}{\eta(\sqrt{S_k} + \sqrt{S_{k-1}})} \leq \frac{S_k - S_{k-1}}{\eta\sqrt{S_k}} = \frac{B_k}{\eta\sqrt{S_k}} = \frac{\eta_k B_k}{\eta^2}.$$

Hence

$$(\star) \leq \frac{\Delta_0}{K\eta_{-1}} + \frac{\widehat{\Delta}_K}{K\eta^2} \sum_{k=0}^{K-1} \eta_k B_k.$$

Substituting into Eq. (72), we obtain

$$\frac{1}{K} \sum_{k=0}^{K-1} \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_*^{1+\tau} \leq \frac{\Delta_0}{K\eta_{-1}} + \frac{\widehat{\Delta}_K}{K\eta^2} \sum_{k=0}^{K-1} \eta_k B_k + \frac{L}{2K} \sum_{k=0}^{K-1} \eta_k B_k. \quad (74)$$

By Lemma 10, we have

$$\sum_{k=0}^{K-1} \eta_k B_k = \sum_{k=0}^{K-1} \frac{\eta B_k}{\sqrt{S_k}} = \sum_{k=0}^{K-1} \frac{\eta B_k}{\sqrt{\sum_{s=0}^k B_s + S_{\text{init}}}} \leq 2\eta \sqrt{\sum_{k=0}^{K-1} B_k}. \quad (75)$$

Substituting Eq. (75) into Eq. (74), we obtain

$$\min_{0 \leq k < K} \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_*^{1+\tau} \leq \frac{\Delta_0}{K\eta_{-1}} + \left( \frac{2\widehat{\Delta}_K}{K\eta} + \frac{L\eta}{K} \right) \sqrt{\sum_{k=0}^{K-1} B_k}. \quad (76)$$

**Case  $\tau = 0$ .** In this case, for each  $i$ ,  $\|\mathbf{v}_k^i\| \in \{0, 1\}$ . Indeed, when  $\tau = 0$  the LMO solves

$$\mathbf{v}_k^i = \text{lmo}(\mathbf{g}_k^i; 0) = \arg \min_{\|\mathbf{y}\| \leq 1} \langle \mathbf{g}_k^i, \mathbf{y} \rangle,$$

which returns a unit vector in the direction of  $-\mathbf{g}_k^i$  whenever  $\mathbf{g}_k^i \neq 0$ , and returns the zero vector only in the degenerate case  $\mathbf{g}_k^i = 0$ . Consequently,  $0 \leq B_k = \frac{1}{N} \sum_{i=1}^N \|\mathbf{v}_k^i\|^2 \leq 1$ .

Thus, the right of Eq. (76) can be further bounded by combining with the bound for  $\widehat{\Delta}_K$  in Lemma 8 (using  $S_K \leq S_{\text{init}} + K$ ):

$$\min_{0 \leq k < K} \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_* \leq \frac{\Delta_0}{K\eta_{-1}} + \frac{2}{K\eta} \left( \Delta_0 + \frac{L}{2} \eta^2 \log\left(\frac{S_{\text{init}} + K}{S_{\text{init}}}\right) \right) \sqrt{K} + \frac{L\eta}{K} \sqrt{K},$$

**Case  $\tau = 1$ .** We now bound the  $\sqrt{\sum_{k=0}^{K-1} B_k}$  terms in Eq. (76) using Young's inequality ( $ab \leq \frac{a^2}{4} + b^2$ ). Specifically,

$$\begin{aligned} \frac{2\widehat{\Delta}_K}{K\eta} \sqrt{\sum_{k=0}^{K-1} B_k} &\leq \frac{\sum_{k=0}^{K-1} B_k}{4K} + \frac{4\widehat{\Delta}_K^2}{\eta^2 K}, \\ \frac{L\eta}{K} \sqrt{\sum_{k=0}^{K-1} B_k} &\leq \frac{\sum_{k=0}^{K-1} B_k}{4K} + \frac{L^2\eta^2}{K}, \end{aligned}$$

and thus

$$\frac{1}{K} \sum_{k=0}^{K-1} \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_*^2 \leq \frac{\Delta_0}{K\eta_{-1}} + \frac{\sum_{k=0}^{K-1} B_k}{2K} + \frac{4\widehat{\Delta}_K^2}{\eta^2 K} + \frac{L^2\eta^2}{K}. \quad (77)$$

Substituting the bound for  $\widehat{\Delta}_K$  in Lemma 8 gives

$$\frac{1}{K} \sum_{k=0}^{K-1} \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_*^2 \quad (78)$$

$$\begin{aligned} &\leq \frac{\Delta_0}{K\eta_{-1}} + \frac{\sum_{k=0}^{K-1} B_k}{2K} + \frac{4}{\eta^2 K} \left( \Delta_0 + \frac{L}{2} \eta^2 \log\left(\frac{S_{K-1}}{S_{\text{init}}}\right) \right)^2 + \frac{L^2\eta^2}{K} \\ &= \frac{\Delta_0}{K\eta_{-1}} + \frac{1}{2} \frac{1}{K} \sum_{k=0}^{K-1} \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_*^2 + \frac{4}{\eta^2 K} \left( \Delta_0 + \frac{L}{2} \eta^2 \log\left(\frac{S_{K-1}}{S_{\text{init}}}\right) \right)^2 + \frac{L^2\eta^2}{K}. \quad (79) \end{aligned}$$

Lastly, by rearranging the term and using the bound for  $\log\left(\frac{S_{K-1}}{S_{\text{init}}}\right)$  in Lemma 11, using

$$\min_{0 \leq k < K} \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_* \leq \sqrt{\min_{0 \leq k < K} \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_*^2}.$$

and

$$\sqrt{a+b+c} \leq \sqrt{a} + \sqrt{b} + \sqrt{c}, \quad \text{for } a, b, c \geq 0,$$

we finishes the proof.  $\blacksquare$

### F.6. Per-particle adaptive LMO particle descent

In this subsection, we consider a per-particle AdaGrad variant of the LMO update, where each particle maintains its own accumulator and stepsize. We give a non-convex convergence bound in the Wasserstein setting for the case  $\tau = 0$ . The framework of the proof generally follows the idea in Attia and Koren [3].

**Theorem 13 (Non-convex analysis of per-particle AdaGrad LMO descent,  $\tau = 0$ )** *Under Asm 2, let  $(\mu_k)_{k \geq 0}$  be the sequence of empirical measures produced by Algorithm 1 with  $\tau = 0$ , i.e.,*

$$\mathbf{v}_k^i = \text{lmo}(\mathbf{g}_k^i; 0), \quad \mathbf{g}_k^i := \nabla_{\mathbf{x}} \frac{\delta F}{\delta \mu}(\mu_k)(X_k^i),$$

and per-particle stepsizes  $\eta_k^i = \eta / \sqrt{S_k^i}$  where  $S_k^i = S_{\text{init}} + \sum_{s=0}^k \|v_s^i\|^2$  and  $S_{\text{init}} > 0$ . Define the suboptimality  $\Delta_k := F(\mu_k) - F^*$ . Then, for any  $K \geq 1$ ,

$$\min_{0 \leq k < K} \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_* \leq \frac{\sqrt{S_{\text{init}} + K}}{K\eta} \left( \Delta_0 + \frac{L}{2} \eta^2 \log\left(\frac{S_{\text{init}} + K}{S_{\text{init}}}\right) \right). \quad (80)$$

**Proof** At iteration  $k$ , the update can be written as a pushforward

$$\mu_{k+1} = (T_k)_{\#} \mu_k, \quad T_k(\mathbf{x}) = \mathbf{x} + \eta_k(\mathbf{x}) \mathbf{v}_k(\mathbf{x}),$$

where  $\mathbf{v}_k(\mathbf{x}) = \mathbf{v}_k^i$  and  $\eta_k(\mathbf{x}) = \eta_k^i$  whenever  $\mathbf{x} = X_k^i$ , and

$$g_k(\mathbf{x}) := \nabla_{\mathbf{x}} \frac{\delta F}{\delta \mu}(\mu_k)(\mathbf{x}).$$

By Asm 2 (Eq. (85)), applied with  $T_k$ , we obtain

$$\begin{aligned} F(\mu_{k+1}) - F(\mu_k) &\leq \int \langle g_k(\mathbf{x}), T_k(\mathbf{x}) - \mathbf{x} \rangle d\mu_k(\mathbf{x}) + \frac{L}{2} \int \|T_k(\mathbf{x}) - \mathbf{x}\|^2 d\mu_k(\mathbf{x}) \\ &= \int \langle g_k(\mathbf{x}), \eta_k(\mathbf{x}) \mathbf{v}_k(\mathbf{x}) \rangle d\mu_k(\mathbf{x}) + \frac{L}{2} \int \|\eta_k(\mathbf{x}) \mathbf{v}_k(\mathbf{x})\|^2 d\mu_k(\mathbf{x}). \end{aligned}$$

Since  $\mu_k$  is empirical,  $\mu_k = \frac{1}{N} \sum_{i=1}^N \delta_{X_k^i}$ , we can write

$$F(\mu_{k+1}) - F(\mu_k) \leq \frac{1}{N} \sum_{i=1}^N \left( \eta_k^i \langle \mathbf{g}_k^i, \mathbf{v}_k^i \rangle + \frac{L}{2} (\eta_k^i)^2 \|\mathbf{v}_k^i\|^2 \right).$$

For  $\tau = 0$ , the LMO identities (Lemma 7) give

$$\langle \mathbf{g}_k^i, \mathbf{v}_k^i \rangle = -\|\mathbf{g}_k^i\|_*, \quad \|\mathbf{v}_k^i\|^2 \in \{0, 1\},$$

where  $\|\mathbf{v}_k^i\|^2 = 1$  whenever  $\mathbf{g}_k^i \neq 0$ , and  $\|\mathbf{v}_k^i\|^2 = 0$  when  $\mathbf{g}_k^i = 0$ . Hence,

$$F(\mu_{k+1}) - F(\mu_k) \leq -\frac{1}{N} \sum_{i=1}^N \eta_k^i \|\mathbf{g}_k^i\|_* + \frac{L}{2N} \sum_{i=1}^N (\eta_k^i)^2 \|\mathbf{v}_k^i\|^2. \quad (81)$$

Rearranging and using  $\Delta_k = F(\mu_k) - F^*$  gives

$$\frac{1}{N} \sum_{i=1}^N \eta_k^i \|\mathbf{g}_k^i\|_* \leq \Delta_k - \Delta_{k+1} + \frac{L}{2N} \sum_{i=1}^N (\eta_k^i)^2 \|\mathbf{v}_k^i\|^2. \quad (82)$$

Summing Eq. (82) over  $k = 0, \dots, K-1$  yields

$$\begin{aligned} \sum_{k=0}^{K-1} \frac{1}{N} \sum_{i=1}^N \eta_k^i \|\mathbf{g}_k^i\|_* &\leq \Delta_0 - \Delta_K + \frac{L}{2N} \sum_{k=0}^{K-1} \sum_{i=1}^N (\eta_k^i)^2 \|\mathbf{v}_k^i\|^2 \\ &\leq \Delta_0 + \frac{L}{2N} \sum_{k=0}^{K-1} \sum_{i=1}^N (\eta_k^i)^2 \|\mathbf{v}_k^i\|^2. \end{aligned} \quad (83)$$

We now bound the double sum on the right-hand side. Fix a particle  $i$  and set

$$a_k^i := \|\mathbf{v}_k^i\|^2, \quad S_k^i = S_{\text{init}} + \sum_{s=0}^k a_s^i, \quad \eta_k^i = \frac{\eta}{\sqrt{S_k^i}}.$$

Then

$$\sum_{k=0}^{K-1} (\eta_k^i)^2 \|\mathbf{v}_k^i\|^2 = \eta^2 \sum_{k=0}^{K-1} \frac{a_k^i}{S_k^i}.$$

By Lemma 9 applied to the sequence  $\{a_k^i\}_{k \geq 0}$  with initial offset  $S_{\text{init}}$ , we have

$$\sum_{k=0}^{K-1} \frac{a_k^i}{S_{\text{init}} + \sum_{s=0}^k a_s^i} \leq \log \left( S_{\text{init}} + \sum_{s=0}^{K-1} a_s^i \right) - \log(S_{\text{init}}) = \log \left( \frac{S_K^i}{S_{\text{init}}} \right).$$

Thus

$$\sum_{k=0}^{K-1} (\eta_k^i)^2 \|\mathbf{v}_k^i\|^2 \leq \eta^2 \log \left( \frac{S_K^i}{S_{\text{init}}} \right).$$

Summing over  $i$  and using Jensen's inequality ( $\frac{1}{N} \sum_i \log u_i \leq \log(\frac{1}{N} \sum_i u_i)$ ), we obtain

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N \sum_{k=0}^{K-1} (\eta_k^i)^2 \|\mathbf{v}_k^i\|^2 &\leq \eta^2 \cdot \frac{1}{N} \sum_{i=1}^N \log \left( \frac{S_K^i}{S_{\text{init}}} \right) \\ &\leq \eta^2 \log \left( \frac{1}{N} \sum_{i=1}^N \frac{S_K^i}{S_{\text{init}}} \right) = \eta^2 \log \left( \frac{\bar{S}_K}{S_{\text{init}}} \right), \end{aligned}$$

where  $\bar{S}_K := \frac{1}{N} \sum_{i=1}^N S_K^i$ . Moreover, since  $a_k^i = \|\mathbf{v}_k^i\|^2 \leq 1$  for all  $k$ , we have

$$S_K^i = S_{\text{init}} + \sum_{s=0}^{K-1} a_s^i \leq S_{\text{init}} + K, \quad \text{so} \quad \bar{S}_K \leq S_{\text{init}} + K.$$

Therefore

$$\frac{1}{N} \sum_{i=1}^N \sum_{k=0}^{K-1} (\eta_k^i)^2 \|\mathbf{v}_k^i\|^2 \leq \eta^2 \log\left(\frac{S_{\text{init}} + K}{S_{\text{init}}}\right).$$

Plugging this into Eq. (83) gives

$$\sum_{k=0}^{K-1} \frac{1}{N} \sum_{i=1}^N \eta_k^i \|\mathbf{g}_k^i\|_* \leq \Delta_0 + \frac{L}{2} \eta^2 \log\left(\frac{S_{\text{init}} + K}{S_{\text{init}}}\right). \quad (84)$$

Next, we lower bound the left-hand side by using the pointwise lower bound on the stepsizes. For every  $i$  and  $k$ ,

$$S_k^i \leq S_{\text{init}} + \sum_{s=0}^{K-1} a_s^i \leq S_{\text{init}} + K, \quad \Rightarrow \quad \eta_k^i = \frac{\eta}{\sqrt{S_k^i}} \geq \frac{\eta}{\sqrt{S_{\text{init}} + K}}.$$

Hence

$$\frac{\eta}{\sqrt{S_{\text{init}} + K}} \sum_{k=0}^{K-1} \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_* \leq \sum_{k=0}^{K-1} \frac{1}{N} \sum_{i=1}^N \eta_k^i \|\mathbf{g}_k^i\|_*.$$

Combining this with Eq. (84), we obtain

$$\sum_{k=0}^{K-1} \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_* \leq \frac{\sqrt{S_{\text{init}} + K}}{\eta} \left( \Delta_0 + \frac{L}{2} \eta^2 \log\left(\frac{S_{\text{init}} + K}{S_{\text{init}}}\right) \right).$$

Dividing both sides by  $K$  gives

$$\min_{0 \leq k < K} \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_* \leq \frac{1}{K} \sum_{k=0}^{K-1} \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_* \leq \frac{\sqrt{S_{\text{init}} + K}}{K\eta} \left( \Delta_0 + \frac{L}{2} \eta^2 \log\left(\frac{S_{\text{init}} + K}{S_{\text{init}}}\right) \right).$$

This completes the proof. ■

## Appendix G. Intuition on the adaptive step size design

Given Asm 2, we can derive the following lemma.

**Lemma 12 (Upper bound of  $F$  by Asm 2)** Under Asm 2, for every  $\mu \in \mathcal{P}_2(\mathbb{R}^d)$  and every measurable transport map  $T \in L^2(\mu; \|\cdot\|)$ , we have:

$$F(T_{\#}\mu) \leq F(\mu) + \int \left\langle \nabla_{\frac{\delta F}{\delta \mu}}(\mu)(\mathbf{x}), T(\mathbf{x}) - \mathbf{x} \right\rangle d\mu(\mathbf{x}) + \frac{L}{2} \int \|T(\mathbf{x}) - \mathbf{x}\|^2 d\mu(\mathbf{x}). \quad (85)$$

**Remark:** Eq. (85) is an analogue of the descent lemma in Euclidean space, and the proof can be found in Sec. F.2.

To understand why the adaptive step size appears naturally, let us start by freezing the step sizes and examining the basic descent mechanism. The analysis above is based on the Wasserstein-type descent lemma Eq. (85), plugging the LMO update  $T_k(\mathbf{x}) = \mathbf{x} + \eta \mathbf{v}_k(\mathbf{x})$  gives the one-step inequality:

$$F(\mu_{k+1}) - F(\mu_k) \leq \frac{\eta}{N} \sum_{i=1}^N \langle \mathbf{g}_k^i, \mathbf{v}_k^i \rangle + \frac{L\eta^2}{2N} \sum_{i=1}^N \|\mathbf{v}_k^i\|^2 = -\frac{\eta}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_*^{1+\tau} + \frac{L\eta^2}{2N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_*^{2\tau},$$

where the last inequality is by Lemma 7. Summing  $k = 0, \dots, K-1$ , rearranging, and dividing by  $\eta$  yields

$$\sum_{k=0}^{K-1} \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_*^{1+\tau} \leq \frac{\Delta_0}{\eta} + \frac{L\eta}{2} \sum_{k=0}^{K-1} \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_*^{2\tau}. \quad (86)$$

The right-hand side is minimized at the optimal stepsize:

$$\eta^* = \sqrt{\frac{2\Delta_0}{L \sum_{k=0}^{K-1} \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_*^{2\tau}}}.$$

Plugging  $\eta^*$  back into Eq. (86) gives the following rates for  $\tau = 0$ :  $\min_k \frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_k^i\|_* \leq \sqrt{\frac{2\Delta_0 L}{K}}$ . Thus, under the standard quadratic Wasserstein-type descent lemma Eq. (85), we can achieve the  $\mathcal{O}(K^{-1/2})$  stationarity rate. The adaptive step size designed in Alg. 1 is precisely a computable surrogate for the ideal but non-implementable  $\eta^*$ .

## Appendix H. Annealing dynamics

We consider the following LMO noisy particle descent rule, which augments the standard noisy particle gradient descent in Eq. (8) with an LMO oracle:

$$\mathbf{x}_{k+1}^i = \mathbf{x}_k^i + \eta \text{lmo} \left( \nabla_{\delta\mu} \frac{\delta F}{\delta \mu}(\mu_{\mathbf{x}_k})(\mathbf{x}_k^i); 0 \right) + \sqrt{2\lambda\eta} \boldsymbol{\xi}_k^i, \quad (87)$$

for  $i \in [N]$ , where  $\boldsymbol{\xi}_k^i \sim \mathcal{N}(0, I_d)$  are i.i.d. Gaussian random variables. Passing formally to the limits  $\eta \rightarrow 0$  and  $N \rightarrow \infty$  motivates the study of the following nonlinear Fokker–Planck equation:

$$\begin{aligned} \partial_t \mu_t &= \nabla \cdot (\mu_t \mathbf{v}_t), \\ \mathbf{v}_t &:= -\text{lmo} \left( \nabla_{\delta\mu} \frac{\delta F}{\delta \mu}(\mu_t); 0 \right) + \lambda_t \nabla \log \mu_t. \end{aligned} \quad (88)$$

Compared with Eq. (3), the key difference is that the first variation term is now passed solely through an LMO operator and the entropy regularization is allowed to vary over time. We emphasize that this is an annealing dynamics, where the regularization parameter  $\lambda_t$  decays to zero as  $t \rightarrow \infty$ , following the annealing framework for MFLD introduced in Chizat [12]. For our analysis, we introduce a Polyak-Lojasiewicz(PL)-type assumptions in Wasserstein space.

**Assumption 8 (Wasserstein PL inequality)** *Let  $\mu^* \in \mathcal{P}_2(\mathbb{R}^d)$  be a minimizer of  $F$ , and denote  $F^* := F(\mu^*)$ . Define  $g_\mu(x) := \nabla_x \frac{\delta F}{\delta \mu}(\mu)(x)$ . We assume that there exists  $\alpha > 0$  such that for every  $\mu$ ,  $\|g_\mu\|_- L^2(\mu; \|\cdot\|_*)^2 \geq 2\alpha(F(\mu) - F^*)$ .*

**Remark 14** Note that such assumption have been studied in prior work [7, 15].

We also need the following regularity assumption along the trajectory.

**Assumption 9 (Divergence control along the trajectory)** Consider the dynamics Eq. (88), denote by  $g_\mu(\mathbf{x}) := \nabla \frac{\delta F}{\delta \mu}(\mu)(\mathbf{x})$ . we assume there exists  $C_{\text{div}} < \infty$  such that  $\forall t \geq 0$ ,  $\left| \int (\nabla \cdot g_{\mu_t}) d\mu_t \right| \leq C_{\text{div}}$ .

Now we have the following convergence guarantee for such annealing LMO dynamics.

**Theorem 15** Under Asm. 4, 8 and 9, consider the flow Eq. (88) with  $\lambda_t = \frac{\gamma}{\log(e+t)}$  for some  $\gamma > 0$ , and assume the solution exists, denoted by  $\{\mu_t\}_{t \geq 0}$ . Then there exists  $t_0 \geq 0$  such that for all  $t \geq t_0$ , we have:

$$F(\mu_t) - F(\mu^*) = O(\log^{-2} t).$$

**Proof** [Proof of Thm. 15] Denote by  $\mathcal{E}(t) := F(\mu_t) - F(\mu^*)$ . Since  $\mu_t$  solves the continuity equation  $\partial_t \mu_t = -\nabla \cdot (\mu_t v_t)$  with  $v_t = s_t - \lambda_t \nabla \log \mu_t$  and  $s_t(\mathbf{x}) := \text{lmo}(g_{\mu_t}(\mathbf{x}); 0)$ , the chain rule gives

$$\frac{d}{dt} F(\mu_t) = \int_{\mathbb{R}^d} \langle g_t(\mathbf{x}), v_t(\mathbf{x}) \rangle d\mu_t(\mathbf{x}), \quad g_t := g_{\mu_t}. \quad (89)$$

By the LMO property on the unit ball,  $\langle g_t(\mathbf{x}), s_t(\mathbf{x}) \rangle = -\|g_t(\mathbf{x})\|_*$  for  $\mu_t$ -a.e.  $\mathbf{x}$ , hence

$$\frac{d}{dt} F(\mu_t) = - \int_{\mathbb{R}^d} \|g_t(\mathbf{x})\|_* d\mu_t(\mathbf{x}) - \lambda_t \int_{\mathbb{R}^d} \langle g_t(\mathbf{x}), \nabla \log \mu_t(\mathbf{x}) \rangle d\mu_t(\mathbf{x}). \quad (90)$$

Next, we use integration by parts for the cross term. Write  $d\mu_t(\mathbf{x}) = \rho_t(\mathbf{x}) d\mathbf{x}$ . Then  $\nabla \log \mu_t = \nabla \rho_t / \rho_t$  and

$$\int \langle g_t, \nabla \log \mu_t \rangle d\mu_t = \int \langle g_t, \nabla \rho_t \rangle d\mathbf{x} = - \int (\nabla \cdot g_t) \rho_t d\mathbf{x} = - \int (\nabla \cdot g_t) d\mu_t,$$

Therefore

$$\frac{d}{dt} F(\mu_t) = -\|g_t\|_{L^1(\mu_t; \|\cdot\|_*)} + \lambda_t \int_{\mathbb{R}^d} (\nabla \cdot g_t)(\mathbf{x}) d\mu_t(\mathbf{x}). \quad (91)$$

By Assumption 9,

$$\frac{d}{dt} F(\mu_t) \leq -\|g_t\|_{L^1(\mu_t; \|\cdot\|_*)} + \lambda_t C_{\text{div}}. \quad (92)$$

Now, by Asm 4, Asm 8 :

$$\|g_t\|_{L^1(\mu_t; \|\cdot\|_*)} \geq \kappa \|g_t\|_{L^2(\mu_t; \|\cdot\|_*)} \geq \kappa \sqrt{2\alpha \mathcal{E}(t)}.$$

Plugging the previous bound to Eq. (92) gives

$$\mathcal{E}'(t) \leq -\kappa \sqrt{2\alpha \mathcal{E}(t)} + \lambda_t C_{\text{div}},$$

Now specialize to the schedule  $\lambda_t = \gamma / \log(e+t)$  and define

$$b(t) := \frac{C_{\text{div}} \gamma}{\log(e+t)}, \quad y(t) := \sqrt{\mathcal{E}(t)}.$$

Since  $\mathcal{E}$  is nonnegative and absolutely continuous,  $y$  is absolutely continuous and differentiable a.e. on  $\{y > 0\}$ , with  $\mathcal{E}'(t) = 2y(t)y'(t)$ . Thus, for a.e.  $t$  with  $y(t) > 0$ ,

$$2y(t)y'(t) \leq -\kappa\sqrt{2\alpha}y(t) + b(t). \quad (93)$$

Dividing by  $2y(t)$  yields

$$y'(t) \leq -\frac{\kappa\sqrt{2\alpha}}{2} + \frac{b(t)}{2y(t)}. \quad (94)$$

Define the barrier

$$h(t) := \frac{2}{\kappa\sqrt{2\alpha}}b(t) = \frac{2C_{\text{div}\gamma}}{\kappa\sqrt{2\alpha}\log(e+t)}.$$

Whenever  $y(t) \geq h(t)$ , we have  $b(t)/(2y(t)) \leq \kappa\sqrt{2\alpha}/4$ , and Eq. (94) implies

$$y'(t) \leq -\frac{\kappa\sqrt{2\alpha}}{4} \quad \text{for a.e. such } t. \quad (95)$$

Fix any  $T \geq 0$ . If  $y(t) \geq h(t)$  for all  $t \geq T$ , then integrating Eq. (95) yields

$$y(t) \leq y(T) - \frac{\kappa\sqrt{2\alpha}}{4}(t-T) \xrightarrow[t \rightarrow \infty]{} -\infty,$$

a contradiction. Hence, for every  $T \geq 0$  there exists  $t \geq T$  such that

$$y(t) \leq h(t). \quad (96)$$

Since  $h$  is  $C^1$  and strictly decreasing,

$$h'(t) = -\frac{2C_{\text{div}\gamma}}{\kappa\sqrt{2\alpha}(e+t)\log^2(e+t)}.$$

As  $h'(t) \uparrow 0$  when  $t \rightarrow \infty$ , there exists  $t_* \geq 0$  such that

$$h'(t) \geq -\frac{\kappa\sqrt{2\alpha}}{4} \quad \forall t \geq t_*. \quad (97)$$

Choose  $t_0 \geq t_*$  such that  $y(t_0) \leq h(t_0)$ . We claim that

$$y(t) \leq h(t) \quad \forall t \geq t_0.$$

Otherwise, define

$$\hat{t} := \inf\{t \geq t_0 : y(t) > h(t)\}.$$

Then  $y(\hat{t}) = h(\hat{t})$  and for  $t > \hat{t}$  sufficiently close to  $\hat{t}$ ,  $y(t) \geq h(t)$ . Integrating Eq. (95) and using Eq. (97) gives

$$y(t) - y(\hat{t}) \leq -\frac{\kappa\sqrt{2\alpha}}{4}(t - \hat{t}), \quad h(t) - h(\hat{t}) \geq -\frac{\kappa\sqrt{2\alpha}}{4}(t - \hat{t}),$$

which implies  $y(t) - h(t) \leq 0$ , a contradiction. Hence the barrier is invariant.

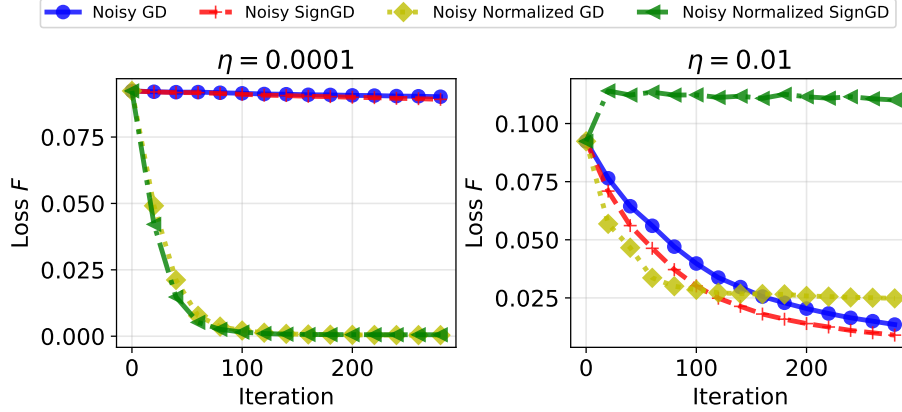


Figure 1: Results on the teacher–student regression task for the mean field Langevin dynamics with LMO noisy particle descent.

Therefore, for all  $t \geq t_0$ ,

$$\sqrt{\mathcal{E}(t)} = y(t) \leq h(t) = \frac{2C_{\text{div}}\gamma}{\kappa\sqrt{2\alpha} \log(e+t)},$$

and squaring yields

$$\mathcal{E}(t) \leq \frac{4C_{\text{div}}^2}{2\alpha\kappa^2} \lambda_t^2 \quad \forall t \geq t_0.$$

Thus,  $F(\mu_t) - F(\mu^*) = O(\log^{-2} t)$ . ■

## Appendix I. Experiment

### I.1. Main experiment

We consider a teacher–student regression task following the setup in Zhu and Chen [49]. Given  $M$  input samples  $\{z_j\}_{j=1}^M \subset \mathbb{R}^d$  drawn i.i.d. from the standard Gaussian distribution, the labels are generated by a teacher model  $y_j = \sin(\hat{\mathbf{x}}^\top z_j)$ ,  $j = 1, \dots, M$ , where  $\hat{\mathbf{x}} \in \mathbb{R}^d$  is a fixed random vector drawn from the standard Gaussian distribution. The student model is a two-layer neural network with  $N$  particles with tanh activation and mean-field scaling:  $h(\mathbf{x}_i, \mathbf{z}) = \tanh(\mathbf{x}_i^\top \mathbf{z})$ , where  $\mathbf{x}_i \in \mathbb{R}^d$  denotes the weight vector of the  $i$ -th neuron and  $\mathbf{X} = (\mathbf{x}_i)_{i=1}^N \in \mathbb{R}^{N \times d}$ . We consider the entropy regularized objective in Eq. (1), with  $F(\mu_{\mathbf{X}}) = \frac{1}{M} \sum_{j=1}^M \left( \frac{1}{N} \sum_{i=1}^N h(\mathbf{x}_i, \mathbf{z}_j) - y_j \right)^2$  and  $\lambda = 0.01$ . We compare different LMO noisy particle descent methods based on Eq. (87), including noisy GD, noisy SignGD, noisy normalized GD, and noisy normalized SignGD. We choose both a smaller and a larger step size:  $\eta = 10^{-4}$  and  $\eta = 10^{-2}$ . The convergence results are shown in Fig. 1, where each curve represents the mean over three random seeds. For a very small step size, *noisy GD and noisy SignGD essentially do not make progress, while normalized methods converge rapidly*. For a larger step size, unnormalized methods converge faster compared with normalized methods. This behavior is consistent with our previous theoretical findings, which show that the normalized LMO

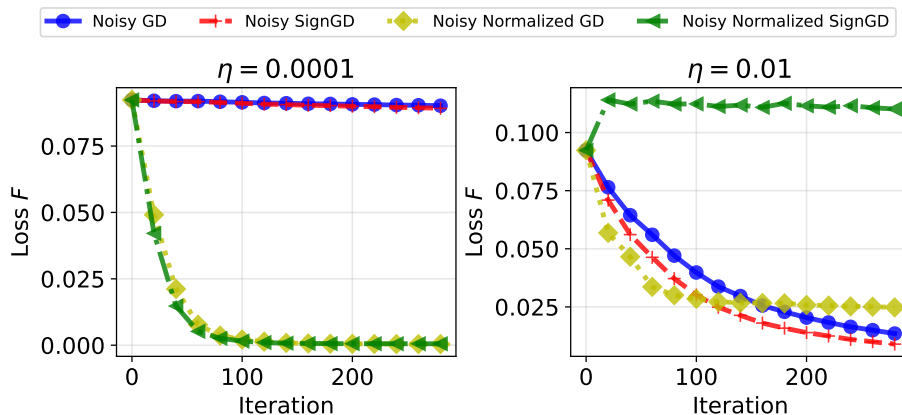


Figure 2: Results on the teacher–student regression task for the mean field Langevin dynamics with LMO noisy particle descent. The dimension of the neuron  $d = 100$ .

with  $\tau = 0$  enjoys faster convergence than the unnormalized one in the gradient flow setting, i.e., with small step size.

Additional experiments on adaptive step-size schemes, comparing global and per-particle variants under different initialization settings, are provided in Sec. I of the appendix. These results are consistent with our main findings, showing that normalized steepest descent converges faster in the small step-size regime.

## I.2. Additional experiment

In Fig. 1 of the previous section, the result is under  $d = 2$ . In this section, we set  $d = 100$  under the same teacher–student regression task, the additional result is given in Fig. 2. We can still see that for a very small step size, noisy GD and noisy SignGD essentially make much less progress compared to normalized methods. And such a gap becomes smaller as the step size increases.

Next, we run experiments regarding the adaptive stepsize scheme developed in Sec. 4.1. Under the same teacher–student regression setting as in Sec. I.1, we conducted additional experiments comparing global adaptive and per-particle adaptive step sizes for both normalized and unnormalized methods. The result is provided in Fig. 3. We have the following observations: Firstly, the results are consistent with our original findings: normalized methods converge faster in the small step-size regime than the unnormalized method.

The difference between global adaptive and per-particle adaptive step sizes is extremely small in this setting. Since all particles are initialized in the same way, their trajectories tend to remain similar, leading to comparable per-particle quantities. Combined with a relatively large number of particles ( $m = 100$ ) in this experiment, this can result in very similar accumulators and hence nearly identical effective step sizes and convergence behavior. Our original paper only shows that both achieve the same convergence rate instead of distinguishing between these two variants.

To further investigate potential differences, we design an additional experiment to amplify particle difference. Specifically, we reduce the number of particles from  $N = 100$  to  $N = 3$ , and replace the original standard Gaussian initialization with a log-normal scaling initialization of the

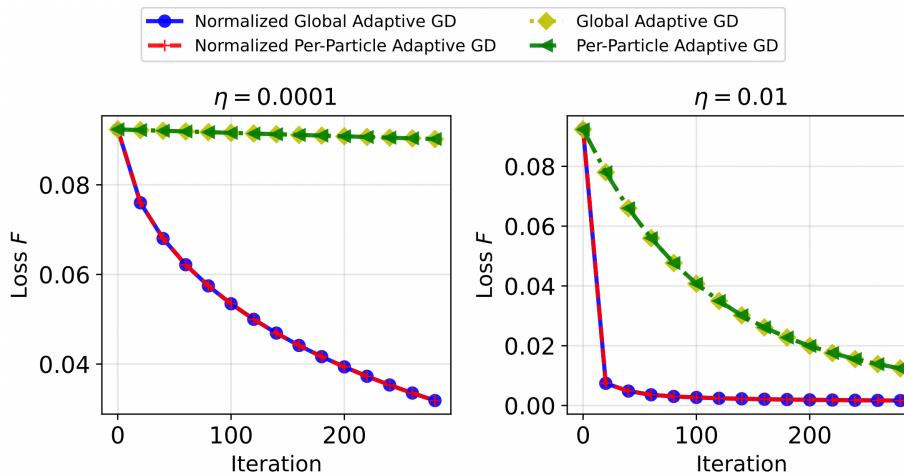


Figure 3: Results on the teacher–student regression task for comparing global adaptive stepsize and per-particle adaptive stepsize.

form  $x_i^0 = \alpha_i \tilde{x}_i$ ,  $\tilde{x}_i \sim \mathcal{N}(0, I_d)$ ,  $\alpha_i = \exp(\xi_i)$ ,  $\xi_i \sim \mathcal{N}(0, 1)$ , independently across particles  $i = 1, \dots, N$ . This increases the variability of their magnitudes.

The results are shown in Fig. 4. Under the normalized scheme, both global and per-particle adaptive methods remain essentially identical as in previous experiment, since the step size does not depend on the magnitude of past gradients.

However, for the unnormalized methods, we observe a noticeable (though still moderate) difference, where the global adaptive method performs slightly better than the per-particle adaptive one. While we do not provide a formal theoretical analysis for this phenomenon, an explanation is that when  $N$  is small and particle-wise gradient magnitudes differ significantly under different initialization, the global adaptive scheme averages these variations, leading to a more stable step size, whereas the per-particle scheme adapts each particle independently and may introduce higher variance in the updates.

## Appendix J. Broader Impact

This paper provides a theoretical analysis of optimization methods for language model training, with a focus on normalized steepest descent dynamics in Wasserstein space and their particle-based implementations. The primary goal of this work is to improve the theoretical understanding of large-scale non-convex model training, which may indirectly contribute to reducing computational costs and energy consumption in modern machine learning systems. From a methodological perspective, our results provide principled guidance for the analysis of geometry-aware optimization algorithms in Wasserstein space, and do not constitute or evaluate applications in any specific domain. We do not foresee direct negative societal or ethical consequences arising uniquely from this work, and this work does not raise new ethical concerns.

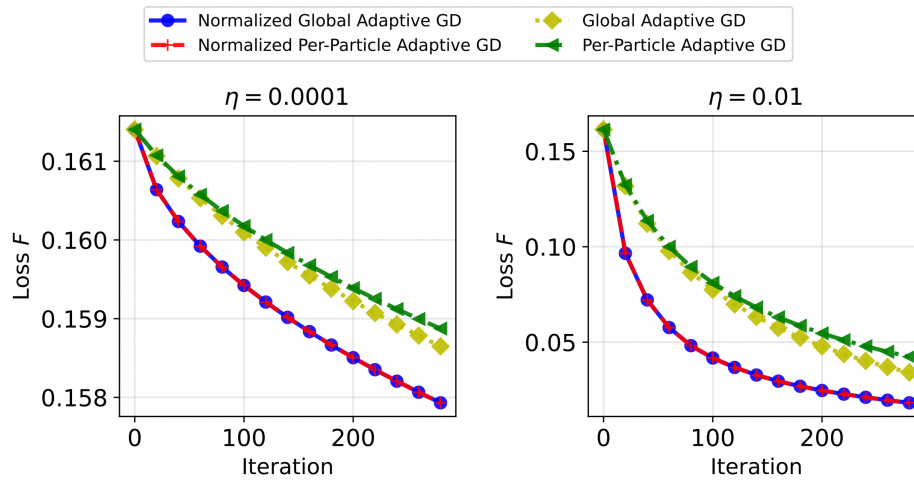


Figure 4: Results on the teacher–student regression task for comparing global adaptive stepsize and per-particle adaptive stepsize with log-normal scaling initialization.

Table 1: Core symbols and constants used in this paper.

Symbol	Type	Description
<b>Optimization on <math>\mathbb{R}^d</math></b>		
$\mathbf{x}_k, \mathbf{x}(t)$	$\mathbb{R}^d$	Parameter at iteration $k$ or time $t$
$f(\mathbf{x})$	$\mathbb{R}^d \rightarrow \mathbb{R}$	Objective in Euclidean space
$f^*$	$\mathbb{R}$	Optimal value $f^* := \min_{\mathbf{x}} f(\mathbf{x})$
$\nabla f(\mathbf{x})$	$\mathbb{R}^d$	Gradient of $f$ at $\mathbf{x}$
$\ \cdot\ , \ \cdot\ _*$	–	A norm on $\mathbb{R}^d$ and its dual: $\ \mathbf{y}\ _* := \sup_{\ \mathbf{x}\  \leq 1} \langle \mathbf{x}, \mathbf{y} \rangle$
$\text{lmo}(\mathbf{g}; \tau)$	$\mathbb{R}^d$	LMO direction (any minimizer), Def. 1
$\eta_k$	$> 0$	Stepsize at iteration $k$
$\Delta_k$	$\mathbb{R}_{\geq 0}$	Suboptimality: $\Delta_k := f(\mathbf{x}_k) - F^*$
<b>Optimization in measure space and particle systems</b>		
$\mathcal{P}_2(\mathbb{R}^d)$	–	Probability measures on $\mathbb{R}^d$ with finite second moment
$\mu, \mu_t$	$\mathcal{P}_2(\mathbb{R}^d)$	Measure (generic / continuous time)
$F(\mu)$	$\mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$	Energy functional on measures
$\frac{\delta F}{\delta \mu}(\mu)(\mathbf{x})$	$\mathbb{R}$	First variation of $F$ at $\mu$ evaluated at $\mathbf{x}$
$g_\mu(\mathbf{x})$	$\mathbb{R}^d$	Wasserstein gradient field: $g_\mu(\mathbf{x}) := \nabla_{\mathbf{x}} \frac{\delta F}{\delta \mu}(\mu)(\mathbf{x})$
$\mathcal{F}(\mu)$	$\mathbb{R}$	Objective in measure space: $\mathcal{F}(\mu) := F(\mu) + \lambda \text{Ent}(\mu)$
$\text{Ent}(\mu)$	$\mathbb{R}$	(Negative) Shannon entropy: if $\mu$ has density $\rho$ , $\text{Ent}(\mu) = \int \rho(\mathbf{x}) \log \rho(\mathbf{x}) d\mathbf{x}$
$W_2(\mu, \nu)$	$\mathbb{R}_{\geq 0}$	2-Wasserstein distance between $\mu, \nu$
$\mathbf{X}_k = (\mathbf{x}_k^i)_{i=1}^N$	$(\mathbb{R}^d)^N$	Collection of particles at iteration $k$
$\mu_{\mathbf{X}_k}$	$\mathcal{P}_2(\mathbb{R}^d)$	Empirical measure: $\mu_{\mathbf{X}_k} := \frac{1}{N} \sum_{i=1}^N \delta_{\mathbf{x}_k^i}$ (and often $\mu_k := \mu_{\mathbf{X}_k}$ )
$\mu^{(N)}, \mu_t^{(N)}, \mu_k^{(N)}$	$\mathcal{P}_2(\mathbb{R}^{dN})$	Law of the $N$ -particle system on $\mathbb{R}^{dN}$ (continuous / discrete time)
$\mathcal{F}^{(N)}(\mu^{(N)})$	$\mathbb{R}$	Finite-particle free energy on $\mathcal{P}_2(\mathbb{R}^{dN})$
$\mu_*^{(N)}$	$\mathcal{P}_2(\mathbb{R}^{dN})$	Minimizer of $\mathcal{F}^{(N)}$ (finite-particle optimum)
$(\mu^*)^{\otimes N}$	$\mathcal{P}_2(\mathbb{R}^{dN})$	$N$ -fold product of the mean-field minimizer $\mu^*$
$\xi, \mathbf{w}$	$\mathbb{R}^d$	Gaussian noise / Brownian motion
$N$	$\mathbb{N}$	Number of particles
$\mathbf{g}_k^i$	$\mathbb{R}^d$	Mean-field gradient at particle: $\mathbf{g}_k^i := \nabla_{\mathbf{x}} \frac{\delta F}{\delta \mu}(\mu_{\mathbf{X}_k})(\mathbf{x}_k^i)$
$\mathbf{v}_k^i$	$\mathbb{R}^d$	Drift / update direction at particle $i$ (cf. Eq. (5))
<b>Constants in assumptions and rates</b>		
$L$	$> 0$	Smoothness / transport-smoothness constant (Assumps. 5.A, 2)
$\alpha$	$> 0$	Strong convexity constant (Assumps. 5.B)
$\kappa$	$(0, 1]$	Non-degeneracy constant (Assump. 4)
$\zeta$	$> 0$	Log-Sobolev constant (Assump. 3)
$C_1$	$> 0$	Norm-comparison constant on $\mathbb{R}^{dN}$ used in Thm. 10: $C_1 := \sup_{\mathbf{a} \neq 0} \frac{\ \mathbf{a}\ ^2}{\ \mathbf{a}\ _*^2}$
$M_1$	$> 0$	Uniform drift bound (Assump. 7)
$M_2$	$> 0$	Score Lipschitz constant (Assump. 7)

Table 2: Examples of LMO and the updates in Euclidean space.<sup>1</sup>

Instance	Norm	$\tau$	$\text{lmo}(\mathbf{g}; \tau)$	Update
<b>Unnormalized steepest descent</b>				
GD [36] / Adam [23]	$\ell_2$	1	$-\mathbf{g}$	$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta_k \nabla f(\mathbf{x}_k)$
SignGD [9]	$\ell_\infty$	1	$-\ \mathbf{g}\ _1 \text{sign}(\mathbf{g})$	$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta_k \ \nabla f(\mathbf{x}_k)\ _1 \text{sign}(\nabla f(\mathbf{x}_k))$
Spectral descent [9]	Spectral	1	$-\sum_i \sigma_i \mathbf{u}\mathbf{v}^\top$	$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta_k \sum_i \sigma_i \mathbf{u}\mathbf{v}^\top$
<b>Normalized steepest descent</b>				
Normalized GD [20]	$\ell_2$	0	$-\frac{\mathbf{g}}{\ \mathbf{g}\ _2}$	$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta_k \frac{\nabla f(\mathbf{x}_k)}{\ \nabla f(\mathbf{x}_k)\ _2}$
Normalized SignGD [4]	$\ell_\infty$	0	$-\text{sign}(\mathbf{g})$	$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta_k \text{sign}(\nabla f(\mathbf{x}_k))$
Muon [22] / Scion [35]	Spectral	0	$-\mathbf{u}\mathbf{v}^\top$	$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta_k \mathbf{u}\mathbf{v}^\top$

<sup>1</sup> In the spectral norm case, the gradient  $\mathbf{g} = \nabla f(\mathbf{x}_k)$  is a matrix with singular value decomposition  $\nabla f(\mathbf{x}_k) = \mathbf{u} \text{diag}(\sigma_1, \dots, \sigma_r) \mathbf{v}^\top$ , and the dual norm is the nuclear norm  $\|\nabla f(\mathbf{x}_k)\|_* = \sum_i \sigma_i$ .