
Data Generation without Function Estimation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Estimating the score function—or other population-density-dependent functions
2 —is a fundamental component of most generative models. However, such function
3 estimation is computationally and statistically challenging. *Can we avoid function*
4 *estimation for data generation?* We propose an **estimation-free** generative method:
5 *A set of points* whose locations are *deterministically* updated with (inverse) *gradient*
6 *descent* can transport a uniform distribution to arbitrary data distribution, in the
7 mean field regime, **without function estimation, training neural networks, and**
8 **even noise injection**. The proposed method is built upon recent advances in
9 the physics of interacting particles. Leveraging recent advances in mathematical
10 physics, we prove that the proposed method samples from the true underlying data
11 distribution in the asymptotic regime, without making any structural assumptions
12 on the distribution.

13 1 Introduction

14 Given i.i.d. samples from an unknown distribution, how to generate a new sample from the distribu-
15 tion? Existing generative models often rely on estimating functions depending on population density
16 such as the score function. Such an estimation is both statistically and computationally challenging.
17 Computationally, estimating the score function even for a simple Gaussian mixture model with
18 maximum likelihood estimation of parameters is NP-hard (Arora et al., 2001). Statistically, score
19 estimation suffers the curse of dimensionality (Wibisono et al., 2024). This raises a fundamental
20 question: can we avoid function estimation for data generation? Ultimately, what is typically available
21 during training is an empirical distribution over i.i.d. samples, for which the score function is not even
22 defined. Can we generate new samples directly from this discrete empirical distribution, avoiding
23 (score) function estimation?

24 We demonstrate that data generation can be achieved without function estimation. Our proposed
25 approach is a novel generative method operating entirely on the empirical distribution of a finite
26 set of training samples. We construct an interactive system that iteratively updates the positions
27 of these data points in two phases: (1) applying standard gradient descent to shape data points (2)
28 performing an inverse gradient descent step to generate new samples. This method builds upon recent
29 advances in the theoretical study of systems of interacting particles (Duerinckx and Serfaty, 2020;
30 Frank and Matzke, 2025). Leveraging this rich literature, we prove that the proposed estimation-free
31 method can transport a uniform measure over a finite ball/sphere to an arbitrary data distribution in
32 the mean-field regime, where the number of samples tends to infinity. Finally, we experimentally
33 validate estimation-free data generation using a finite number of points.

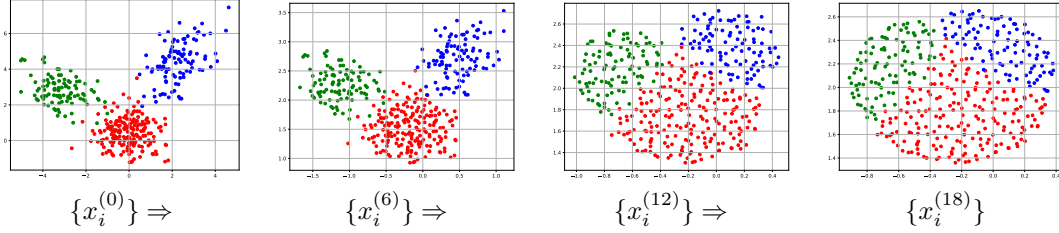


Figure 1: **Convergence of the empirical distribution with optimization.** Scatter plot of gradient descent iterates $x_i^{(k)}$, defined in (3), initialized at $x_i^{(0)}$ drawn from a Gaussian mixture distribution. Different mixture components (modes) are distinguished by color. As k evolves, the points become uniformly distributed on the circle.

34 2 Method

35 Given the support of an empirical distribution over points $x_1, \dots, x_n \in \mathbb{R}^d$, consider the following
 36 optimization problem:

$$x_1^*, \dots, x_n^* := \arg \min_{x_1, \dots, x_n \in \mathbb{R}^d} \left(E_n^{(\epsilon)}(x_1, \dots, x_n) := \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j \neq i} W_\epsilon^{(s)}(x_i - x_j) \right), \quad (1)$$

37 where $W_\epsilon^{(s)}(x-y) = \frac{\|x-y\|^2}{2} + \frac{1}{s(\|x-y\|^2 + \epsilon)^{s/2}}$ is called a "attractive-repulsive power-law interaction
 38 potential" (Balagué et al., 2013; Shu, 2025). The squared-norm term encourages attraction between
 39 particles, while the inverse-norm term induces repulsion, preventing collapse¹. These competing
 40 forces lead to a striking structure in the distribution of the optimized points x_1^*, \dots, x_n^* in the
 41 asymptotic regime as $n \rightarrow \infty$, where the limiting distribution is characterized by

$$\arg \min_{p \in \Omega} \left(E(p) := \int W^{(s)}(x-y) p(x) p(y) dx dy \right), \Omega := \{\text{probability measures over } \mathbb{R}^d\}, \quad (2)$$

42 where $W^{(s)} := W_0^{(s)}$. The minimizer is unique up to translation and is uniform over

$$\begin{cases} \text{a ball (Carrillo and Shu, 2023)} & \text{if } s = d - 2, \\ \text{a sphere (Frank et al., 2025)} & \text{if } 2 < d \text{ and } -2 \leq s < d - 4, \end{cases}$$

43 with a radius that is finite and computable in closed form. By optimizing the locations x_1, \dots, x_n ,
 44 one can asymptotically transport any empirical distribution to a simple uniform distribution over a
 45 compact ball or sphere. This optimization can be performed via standard gradient descent (GD):

$$\forall i: \quad x_i^{(k+1)} = x_i^{(k)} - \gamma \nabla_{x_i} E_n^{(\epsilon)}(x_1^{(k)}, \dots, x_n^{(k)}). \quad (3)$$

46 Figure 1 plots the time evolution of points $\{x_i^{(k)}\}$ indexed by k starting from the initial points
 47 $x_i^{(0)}$ drawn i.i.d. from Gaussian mixture. We observe that the distribution becomes uniform as k
 48 increases. Although the distribution is uniform, it contains very interesting information about the
 49 initial Gaussian mixture distribution. For example, the ratio of area of different colors are proportional
 50 to the number of points in each cluster. This structure indicates that it may be possible to recover
 51 original data distribution from the final state. Whereas in the initial state $x_i^{(0)}$ the data points are
 52 i.i.d., evolving according to (3) induces strong correlations and entanglement among them under the
 53 gradient descent dynamics. Thus, in the limit $k \rightarrow \infty$, although the marginal distribution of each
 54 data point becomes uniform, their joint distribution still conveys potentially useful information about
 55 the original distribution.

56 Even more interesting is the inverse of gradient descent: assuming GD converges to the minimizer
 57 of E_n , inverse GD maps points uniformly distributed on a sphere (or disk) to an arbitrary empirical
 58 distribution supported on a finite set of points. The Proposition 11 in the Appendix guarantees that
 59 this inverse exists and is computable under mild assumptions. Indeed, the inverse process can be

¹For $s = 0$, the repulsive term is replaced by the logarithm of the norm, as given by $s \rightarrow 0$ (Shu, 2025).

60 efficiently computed using gradient descent on a convex function, without requiring the score function
 61 or any other statistical information about the population density. Combining gradient descent (GD)
 62 and its inverse can deterministically transport an empirical distribution to a uniform distribution and
 63 vice versa. This insight enables a three-step sampling method that operates on a finite set of points:
 64 (1) apply GD to the n available samples from the target distribution; (2) add a new point drawn
 65 uniformly at random from a sphere or ball; and (3) apply inverse GD. We call this three step method
 66 **estimation-free sampling (EFS)**.

67 (1) *Forward optimization*: Apply gradient descent to $E_n^{(\epsilon)}$ (see (3)) starting from i.i.d. samples
 68 $x_1^{(0)}, \dots, x_n^{(0)}$ drawn from the target distribution. In the next section, we show that this optimization
 69 acts as a transport map, pushing the empirical distribution toward a uniform distribution in the
 70 asymptotic regime.

71 (2) *Augmentation*: Add a new point $y^{(k)}$, sampled from a uniform distribution, to the set $x_0^{(k)}, \dots, x_n^{(k)}$.
 72 To generate $y^{(k)}$, we first estimate the center and radius of the ball or sphere enclosing the
 73 points $x_i^{(k)}$. $y^{(k)}$ can also be generated by interpolating between two points in the set, i.e.,
 74 $y^{(k)} = (1 - t)x_i^{(k)} + tx_j^{(k)}$ for some $i \neq j$ and $t \in (0, 1)$.

75 (3) *Backward optimization*: Compute the inverse GD for the newly generated point $y^{(k)}$ by optimizing
 76 a convex function (defined in Proposition 11), using efficient gradient descent with a constant step
 77 size. While the forward process involves n points, the backward step involves $n + 1$ points, consisting
 78 of the n points from the forward step and the newly generated point from step (2). Find Algorithm
 79 details in Appendix A.

80 3 Theory

81 We demonstrate that the proposed estimation-free sampling (EFS) method can generate new samples
 82 from a distribution in the asymptotic regime. As established in Proposition 11, backward step is
 83 convex problem with well-established theoretical guarantees. In contrast, the forward optimization
 84 is non-convex. Nevertheless, its global convergence can be analyzed using recent advances in
 85 Wasserstein gradient flows on interactive-repulsive energy (Duerinckx and Serfaty, 2020; Carrillo
 86 and Shu, 2023). As the potential function is shift-invariant, *all statements hold up to translation*, and
 87 we refrain from repeating "up to translation" for simplicity.

88 3.1 Forward Transport to Uniform Distribution

89 While our ultimate goal is to analyze the backward step of EFS, which generates new samples, the
 90 forward step also plays a critical role in data generation. Here, we show that the forward step of EFS
 91 *transports* the data distribution to a uniform distribution².

92 Consider a continuous-time model where the samples $x_i^{(k)}$ updated by gradient descent, modeled by
 93 the following ordinary differential equation (ODE):

$$\frac{dx_i}{dt} = -\frac{1}{n} \sum_{j \neq i} \nabla W^{(s)}(x_i - x_j), \quad (4)$$

94 ODEs have been widely used to describe the limiting behavior of gradient descent with an infinitesi-
 95 mally small step size (Su et al., 2016; Zhang et al., 2021a,b; Chizat and Bach, 2018). Analyzing the
 96 above system becomes challenging due to the coupling between the variables x_i , as the complexity
 97 increases with the number of particles n . A common approach to address this challenge is to analyze
 98 the evolution of the empirical distribution of the particles at a macroscopic level, rather than tracking
 99 their individual trajectories. Define the empirical distribution as

$$\mu_t^{(n)} := \frac{1}{n} \sum_{i=1}^n \delta_{x_i(t)}, \quad (5)$$

100 where δ_x denotes the Dirac measure centered at x . Consider the following PDE:

$$\frac{d\mu}{dt} = \operatorname{div} \left(\mu_t(x) \int \nabla W^{(s)}(x - y) \mu_t(y) dy \right), \quad (6)$$

²A map f transports measure μ to ν if the distribution of $f(x)$ is ν when $x \sim \mu$ (Villani et al., 2008).

101 where div denotes the divergence operator acting on the associated vector field. Duerinckx and
 102 Serfaty (2020) prove $\mu_t^{(n)}$ converges in the weak sense to μ_t , the solution of the PDE above.

103 **Theorem 1** (Duerinckx and Serfaty (2020)). *If $\mu_0^{(n)}$ converges to a regular measure μ_0 in Wasserstein
 104 distance, then $\mu_t^{(n)}$ converges weakly to μ_t , provided $d - 2 \leq s < d$.*

105 The macroscopic measure μ_t is significantly easier to analyze than the microscopic trajectories $x_i(t)$.
 106 While the energy E_n is generally non-convex in $x_i(t)$, the limiting energy E obeys linear interpolation
 107 convexity in μ (Shu, 2025). This convexity can be exploited to characterize the asymptotic behavior
 108 of μ_t as $t \rightarrow \infty$ (Shu, 2025).

109 **Theorem 2** (Frank and Matzke (2025); Carrillo and Shu (2023)). *The steady state of
 110 μ_t is the global minimizer of E (up to translation) which is the uniform measure over:
 111 $\begin{cases} \text{a ball (Carrillo and Shu, 2023)} & \text{if } d - 2 \leq s < d, \\ \text{a sphere (Frank et al., 2025)} & \text{if } 2 < d \text{ and } -2 \leq s < d - 4, \end{cases}$ with finite radius.*

112 Notably, proving that a steady state (i.e., a local minimizer) is in fact a global minimizer required
 113 decades of mathematical development, beginning with the foundational work of Frostman (1935),
 114 and has only recently been fully resolved in certain parameter regimes (Frank and Matzke, 2025).
 115 These advances allow us to analyze the asymptotic dynamics of the forward step in EFS.

116 **Corollary 3.** *If $\mu_0^{(n)}$ converges to μ_0 in Wasserstein-2 distance, then the measure $\mu_t^{(n)}$ converges in
 117 the weak sense to the uniform distribution over a ball of finite radius as $n \rightarrow \infty$ and $t \rightarrow \infty$, for
 118 $s = d - 2$.*

119 3.2 Backward Transport to Target Distribution

120 It is straightforward to verify that the backward optimization can exactly recover the original training
 121 data $x_k^{(i)}$ by initializing EFS with $y_k = x_k^{(i)}$. However, our goal is not to reconstruct existing data
 122 points, but to generate new samples from the underlying data distribution. We show that this form of
 123 generalization is achievable in the asymptotic regime.

124 Consider the continuous-time formulation of the backward optimization process:

$$\frac{dy^{(n)}}{d\bar{t}} = \frac{1}{n} \sum_{i=1}^n \nabla W^{(s)} \left(y_t^{(n)} - x_i(t) \right), \quad \text{where } x_i(t) \text{ is defined in (4)}. \quad (7)$$

125 Here, the differential $d\bar{t}$ indicates that the process is reversed time (Anderson, 1982). This dynamics
 126 corresponds to EFS in the limit of an infinitesimally small step size γ and a large terminal time
 127 $T \rightarrow \infty$. We prove this dynamics transport μ_t back to μ_0 in the asymptotic regime as $n \rightarrow \infty$. In
 128 other words, the distribution of $y_0^{(n)}$ converges to μ_0 as $n \rightarrow \infty$.

129 **Theorem 4.** *Assume $\mu_0^{(n)}$ converges to μ_0 in Wasserstein-2 distance. Suppose that $y_t^{(n)}$ is a random
 130 variable with law μ_t , where μ_t is the solution to the continuity equation (6). Then, the distribution
 131 of $y_0^{(n)}$ —obtained from the reversed-time ODE (7)—converges to μ_0 as $n \rightarrow \infty$, for continuous
 132 measures μ_0 and μ_t , and for $s = d - 2$.*

133 To generate new samples, we assume access to i.i.d. samples from the target distribution. Since the
 134 empirical measure $\mu_0^{(n)}$ converges to μ_0 in Wasserstein distance (Villani et al., 2008) which ensures
 135 the assumption holds in the last theorem. According to Corollary 3, the measure μ_t converges to a
 136 uniform distribution as $t \rightarrow \infty$. Consequently, sampling from a nearly uniform distribution allows
 137 for effective recovery of the initial data distribution μ_0 .

138 Thus, a set of points whose locations are updated via (inverse) gradient descent, can provably
 139 generate new samples from an arbitrary distribution μ_0 —without function estimation, training neural
 140 networks, or injecting noise. Avoiding explicit function estimation enables us to establish asymptotic
 141 guarantees for EFS without any structural assumptions on the data distribution. In comparison,
 142 methods such as Langevin dynamics depend on the log-concavity of the distribution and its explicit
 143 form $\mu_0 = \frac{e^{-U(x)}}{\int e^{-U(x)} dx}$. Instead, EFS requires only i.i.d. samples. For more details on related works,
 144 please see Section B in the Appendix.

145 **A Note on Proofs and Experimental Results**

146 Due to space constraints, we defer the proofs, experimental results, related work, and discussion of
147 limitations to the Appendix.

148 **References**

- 149 Sanjeev Arora, Ravi Kannan, and Ankur Moitra. Learning mixtures of gaussians is np-hard. In
150 *Proceedings 34th Annual ACM Symposium on Theory of Computing (STOC)*. ACM, 2001.
- 151 Andre Wibisono, Yihong Wu, and Kaylee Yingxi Yang. Optimal score estimation via empirical bayes
152 smoothing. In *The Thirty Seventh Annual Conference on Learning Theory*, pages 4958–4991.
153 PMLR, 2024.
- 154 Mitia Duerinckx and Sylvia Serfaty. Mean field limit for coulomb-type flows. *Duke Mathematical*
155 *Journal*, 169(15):2887–2935, 2020.
- 156 Rupert L Frank and Ryan W Matzke. Minimizers for an aggregation model with attractive–repulsive
157 interaction. *Archive for Rational Mechanics and Analysis*, 249(2):15, 2025.
- 158 Daniel Balagué, José A Carrillo, Thomas Laurent, and Gaël Raoul. Nonlocal interactions by
159 repulsive–attractive potentials: radial ins/stability. *Physica D: Nonlinear Phenomena*, 260:5–25,
160 2013.
- 161 Ruiwen Shu. A family of explicit minimizers for interaction energies. *arXiv preprint*
162 *arXiv:2501.14666*, 2025.
- 163 José A Carrillo and Ruiwen Shu. From radial symmetry to fractal behavior of aggregation equilibria
164 for repulsive–attractive potentials. *Calculus of Variations and Partial Differential Equations*, 62
165 (1):28, 2023.
- 166 Rupert L Frank, Rupert L Frank, Rupert L Frank, and Ryan W Matzke. Minimizers for an aggregation
167 model with attractive–repulsive interaction. *Archive for Rational Mechanics and Analysis*, 249(2):
168 15, 2025.
- 169 Cédric Villani et al. *Optimal transport: old and new*, volume 338. Springer, 2008.
- 170 Weijie Su, Stephen Boyd, and Emmanuel J Candes. A differential equation for modeling nesterov’s
171 accelerated gradient method: Theory and insights. *Journal of Machine Learning Research*, 17
172 (153):1–43, 2016.
- 173 Peiyuan Zhang, Antonio Orvieto, and Hadi Daneshmand. Rethinking the variational interpretation
174 of accelerated optimization methods. *Advances in Neural Information Processing Systems*, 34:
175 14396–14406, 2021a.
- 176 Peiyuan Zhang, Antonio Orvieto, Hadi Daneshmand, Thomas Hofmann, and Roy S Smith. Revisiting
177 the role of euler numerical integration on acceleration and stability in convex optimization. In
178 *International Conference on Artificial Intelligence and Statistics*, pages 3979–3987. PMLR, 2021b.
- 179 Lenaic Chizat and Francis Bach. On the global convergence of gradient descent for over-parameterized
180 models using optimal transport. *Advances in neural information processing systems*, 31, 2018.
- 181 O Frostman. Potentiel d’équilibre et theorie des ensembles. *thesis*, 1935.
- 182 Brian DO Anderson. Reverse-time diffusion equation models. *Stochastic Processes and their*
183 *Applications*, 12(3):313–326, 1982.
- 184 Keith Ball et al. An elementary introduction to modern convex geometry. *Flavors of geometry*, 31
185 (1-58):26, 1997.
- 186 George EP Box and Mervin E Muller. A note on the generation of random normal deviates. *The*
187 *Annals of Mathematical Statistics*, 29(2):610–611, 1958.
- 188 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
189 Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural*
190 *Information Processing Systems*, volume 27, 2014.
- 191 Danilo Jimenez Rezende and Shakir Mohamed. Variational inference with normalizing flows. In
192 *International Conference on Machine Learning*, pages 1530–1538. PMLR, 2015.

- 193 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
194 Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint*
195 *arXiv:2011.13456*, 2020.
- 196 Luigi Ambrosio, Nicola Gigli, and Giuseppe Savaré. *Gradient flows: in metric spaces and in the*
197 *space of probability measures*. Springer Science & Business Media, 2008.
- 198 Richard Jordan, David Kinderlehrer, and Felix Otto. The variational formulation of the fokker–planck
199 equation. *SIAM journal on mathematical analysis*, 29(1):1–17, 1998.
- 200 Andrei Kolmogoroff. Über die analytischen methoden in der wahrscheinlichkeitsrechnung. *Mathe-*
201 *matische Annalen*, 104:415–458, 1931.
- 202 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised
203 learning using nonequilibrium thermodynamics. In *International conference on machine learning*,
204 pages 2256–2265. pmlr, 2015.
- 205 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*
206 *neural information processing systems*, 33:6840–6851, 2020.
- 207 Min Jae Song. Cryptographic hardness of score estimation. In *The Thirty-eighth Annual Conference*
208 *on Neural Information Processing Systems*, 2024.
- 209 Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks.
210 In *International conference on machine learning*. PMLR, 2017.
- 211 Michael Arbel, Anna Korba, Adil Salim, and Arthur Gretton. Maximum mean discrepancy gradient
212 flow. *Advances in Neural Information Processing Systems*, 32, 2019.
- 213 Youssef Mroueh, Tom Sercu, and Anant Raj. Sobolev descent. In *The 22nd International Conference*
214 *on Artificial Intelligence and Statistics*, pages 2976–2985. PMLR, 2019.
- 215 Petr Mokrov, Alexander Korotin, Lingxiao Li, Aude Genevay, Justin M Solomon, and Evgeny
216 Burnaev. Large-scale wasserstein gradient flows. *Advances in Neural Information Processing*
217 *Systems*, 34:15243–15256, 2021.
- 218 David Alvarez-Melis, Yair Schiff, and Youssef Mroueh. Optimizing functionals on the space of
219 probabilities with input convex neural network. In *Annual Conference on Neural Information*
220 *Processing Systems*, 2021.
- 221 Jiaojiao Fan, Qinsheng Zhang, Amirhossein Taghvaei, and Yongxin Chen. Variational wasserstein
222 gradient flow. In *International Conference on Machine Learning*, pages 6185–6215. PMLR, 2022.
- 223 Clément Bonet, Nicolas Courty, François Septier, and Lucas Drumetz. Efficient gradient flows in
224 sliced-wasserstein space. *Transactions on Machine Learning Research Journal*, 2022.
- 225 Charlotte Bunne, Laetitia Papaxanthos, Andreas Krause, and Marco Cuturi. Proximal optimal
226 transport modeling of population dynamics. In *International Conference on Artificial Intelligence*
227 *and Statistics*, pages 6511–6528. PMLR, 2022.
- 228 Fabian Altekrüger, Johannes Hertrich, and Gabriele Steidl. Neural wasserstein gradient flows for
229 discrepancies with riesz kernels. In *International Conference on Machine Learning*, pages 664–690.
230 PMLR, 2023.
- 231 Brandon Amos, Lei Xu, and J Zico Kolter. Input convex neural networks. In *International conference*
232 *on machine learning*, pages 146–155. PMLR, 2017.
- 233 Atsushi Nitanda and Taiji Suzuki. Stochastic particle gradient descent for infinite ensembles. *arXiv*
234 *preprint arXiv:1712.05438*, 2017.
- 235 Hadi Daneshmand, Jason D Lee, and Chi Jin. Efficient displacement convex optimization with
236 particle gradient descent. In *International Conference on Machine Learning*, pages 6836–6854.
237 PMLR, 2023.

- 238 Yilun Xu, Ziming Liu, Max Tegmark, and Tommi Jaakkola. Poisson flow generative models.
239 *Advances in Neural Information Processing Systems*, 35:16782–16795, 2022.
- 240 Yilun Xu, Ziming Liu, Yonglong Tian, Shangyuan Tong, Max Tegmark, and Tommi Jaakkola.
241 Pfgm++: Unlocking the potential of physics-inspired generative models. In *International Confer-*
242 *ence on Machine Learning*, pages 38566–38591. PMLR, 2023.
- 243 Alexander Kolesov, Manukhov Stepan, Vladimir V Palyulin, and Alexander Korotin. Field matching:
244 an electrostatic paradigm to generate and transfer data. *arXiv preprint arXiv:2502.02367*, 2025.
- 245 Ruiwen Shu and Jiuya Wang. Generalized erdős-turán inequalities and stability of energy minimizers.
246 *arXiv e-prints*, pages arXiv–2110, 2021a.
- 247 Filippo Santambrogio. {Euclidean, metric, and Wasserstein} gradient flows: an overview. *Bulletin of*
248 *Mathematical Sciences*, 2017.
- 249 Ruiwen Shu and Jiuya Wang. The sharp erdős-turán inequality. *arXiv preprint arXiv:2109.11006*,
250 2021b.
- 251 Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. A
252 kernel two-sample test. *The Journal of Machine Learning Research*, 13(1):723–773, 2012.
- 253 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
254 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward
255 Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner,
256 Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance
257 deep learning library. In *Advances in Neural Information Processing Systems 32 (NeurIPS 2019)*,
258 pages 8024–8035, 2019. URL [https://papers.nips.cc/paper_files/paper/2019/file/
259 bdbca288fee7f92f2bfa9f7012727740-Paper.pdf](https://papers.nips.cc/paper_files/paper/2019/file/bdbca288fee7f92f2bfa9f7012727740-Paper.pdf).
- 260 Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to
261 document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 262 Jonathan Masci, Ueli Meier, Dan Ciresan, and Jürgen Schmidhuber. Stacked convolutional auto-
263 encoders for hierarchical feature extraction. In *International Conference on Artificial Neural*
264 *Networks*, pages 52–59. Springer, 2011.

Appendix

266 A Algorithm

267 For completeness, we present the EFS algorithm in Algorithm 1 and its forward and backward
 268 optimization subroutines in Algorithm 2 and Algorithm 3, respectively. Remarkably, even when the
 269 distribution is uniform over a high-dimensional ball, the samples are effectively drawn uniformly
 270 from the unit sphere, which is denoted by \mathbb{S}^{d-1} , as almost all the volume concentrates near the
 271 boundary (Ball et al., 1997).

Algorithm 1: Estimation Free Sampling (EFS)

Input: I.I.D samples $\{x_1^{(0)}, \dots, x_n^{(0)}\}$
 1 **Parameters:** Step size γ , number of forward iterations k , learning rate of backward (proximal
 step) β , number of iterations for each backward proximal step T , potential parameters s and ϵ
 /* Forward the training data to the sphere by gradient descent */
 2 **Set** $\{x_0^{(j)}, \dots, x_n^{(j)}\}_{j=1}^k \leftarrow \text{Forward}(\{x_1^{(0)}, \dots, x_n^{(0)}\}, \gamma, k, s, \epsilon)$
 272 /* Draw a new random point uniformly from the sphere */
 3 **Set** $c = \frac{1}{n} \sum_{i=1}^n x_i^{(k)}, r = \frac{1}{n} \sum_{i=1}^n \|c - x_i^{(k)}\|$
 4 **Draw** $\nu \sim \mathbb{S}^{d-1}(c, r)$ // Sphere with center c and radius r
 /* Backward the new data point to the original space */
 5 **Set** $y^0 \leftarrow \text{Backward}(\{x_1^{(j)}, \dots, x_n^{(j)}\}_{j=1}^k, y^{(k)}, \gamma, k, \beta, T, s, \epsilon)$
Output: Generated sample $y^{(0)}$

Algorithm 2: Forward Optimization (Forward)

Input: Training data $\{x_1^{(0)}, \dots, x_n^{(0)}\}$
 1 **Parameters:** step size γ , number of forward iterations k , potential parameters s and ϵ
 2 **for** $j \leftarrow 0$ **to** $k - 1$ **do**
 3 **for** $i \leftarrow 1$ **to** n **do**
 4 $\Delta_i = \frac{1}{n-1} \sum_{a \in [n], a \neq i} \nabla W_\epsilon^{(s)}(x_i^{(j)} - x_a^{(j)})$
 5 $x_i^{(j+1)} = x_i^{(j)} - \gamma \Delta_i$
 6 **end**
 7 **end**
Output: $\{x_0^{(j)}, \dots, x_n^{(j)}\}_{j=1}^k$

Algorithm 3: Backward Optimization (Backward)

Input: Data $\{x_1^{(j)}, \dots, x_n^{(j)}\}_{j=1}^k$, new sample $y^{(k)}$
 1 **Parameters:** Step size γ , number of forward iterations k , learning rate of backward (proximal
 step) β , number of iterations for each backward proximal step T , potential parameters s and ϵ
 2 **Set:** $j = k$
 3 **while** $j \geq 0$ **do**
 4 **Set:** $v_0 = y^{(j)}$
 273 5 **for** $t \leftarrow 0$ **to** T **do**
 6 $\nabla = \frac{\gamma}{n} \sum_{i \in [n]} \nabla W_\epsilon^{(s)}(v_t - x_i^{(j)})$
 7 $\Delta = v_t - y^{(j)} - \nabla$
 8 $v_{t+1} = v_t - \beta \Delta$
 9 **end**
 10 $j = j - 1$ and $y^{(j)} = v_T$
 11 **end**
Output: $y^{(0)}$

274 **B Related works**

275 **Data generation with estimation.** Generative models traditionally rely on transporting simple to
 276 complex distributions. A classical example is the Box–Muller transform: given independent random
 277 variables $\theta \sim \text{Uniform}[0, 2\pi]$ and $r \sim \text{Exponential}$, the random vector $v = r(\cos \theta, \sin \theta)$ follows a
 278 standard two-dimensional Normal distribution (Box and Muller, 1958). This illustrates a fundamental
 279 idea: sampling from a complex distribution can be achieved via a nonlinear transformation of samples
 280 drawn from a simpler distribution.

281 Modern generative models adopt this principle using function approximation: Given latent variable
 282 x sampled from a Gaussian distribution, the goal is to find a parametric function f_θ such that the
 283 distribution of $f_\theta(x)$ approximates the data distribution. In generative adversarial networks (GANs),
 284 f_θ is implemented by a neural network and trained via adversarial objectives (Goodfellow et al., 2014).
 285 In normalizing flows, f_θ is a sequence of invertible and differentiable transformations optimized
 286 via maximum likelihood (Rezende and Mohamed, 2015). Score-based diffusion models, gradually
 287 transform data using estimated score functions of diffusion density (Song et al., 2020). In contrast to
 288 these methods, our proposed approach computes a non-linear transformation using a set of interacting
 289 points, thereby avoiding the estimation of transport map f_θ .

290 **Variational perspective towards generative models.** Despite its distinct mechanism, EFS shares
 291 conceptual parallels with diffusion models. Both can be interpreted as discretizations of different
 292 "*Wasserstein gradient flows*" (Ambrosio et al., 2008). This connection arises from a common
 293 variational principle: many simple distributions from which i.i.d. samples are easily drawn, such as
 294 the Gaussian distribution or the uniform distribution on a sphere, can be viewed as equilibrium states
 295 of physical systems with minimum energy. For example, Gaussian distribution is the minimizer of
 296 the negative entropy functional over the space of probability measures (Jordan et al., 1998), while the
 297 uniform distribution on the unit sphere minimizes a Riesz-type interaction energy (Frank and Matzke,
 298 2025). More generally, such distributions can be described as

$$p^* = \arg \min_{p \in \Omega} F(p), \quad \Omega := \{\text{probability measures over } \mathbb{R}^d\}, \quad (8)$$

299 where F is an appropriate energy functional.

300 This variational formulation provides a principled mechanism to transport any distribution to the
 301 minimum energy state by gradient descent on F in the space of probability densities. Diffusion models,
 302 optimize $F(p) = \int \log(p(x))p(x)dx$ using Wasserstein gradient flow (Jordan et al., 1998). More
 303 interesting observation is that the reverse of gradient flow which can transport Gaussian distribution
 304 to any distribution. For entropy, inverse Wasserstein gradient flow is implemented by the Kolmogorov
 305 backward equation (Kolmogoroff, 1931), enabling transformation from a Gaussian to arbitrary target
 306 distributions. This key result forms the basis of modern generative diffusion models (Sohl-Dickstein
 307 et al., 2015; Ho et al., 2020; Song et al., 2020), where the reverse-time stochastic process includes a
 308 drift term proportional to the score function of intermediate densities (Anderson, 1982). However,
 309 estimating this score function for arbitrary distributions remains statistically and computationally
 310 challenging. In this vein, Wibisono et al. (2024) recently proved that estimating the score function of
 311 sub-Gaussian distributions with Lipschitz-continuous scores suffers from a curse of dimensionality
 312 in the required sample complexity. Song (2024) further demonstrated that, under lattice-based
 313 cryptographic hardness assumptions, score estimation remains computationally intractable even when
 314 the sample complexity is polynomial in the relevant parameters. To circumvent these challenges, our
 315 method replaces the entropy-based potential with a Riesz-type interaction energy and optimizes it
 316 directly over discrete empirical measures.

317 Inspired by Formulation (8), one may attempt to optimize functionals of the probability measure p
 318 that quantify its distance to a target distribution p^* . For example, the celebrated work of Arjovsky
 319 et al. (2017) introduce the Wasserstein GAN, which optimizes the Wasserstein-1 distance between the
 320 model distribution p and the target distribution p^* by solving its Kantorovich–Rubinstein dual via an
 321 adversarial training objective over both the generator (which parameterizes p) and the critic (which
 322 approximates the dual). More closely related to gradient flows on the space of probability measures,
 323 Arbel et al. (2019) consider the Wasserstein gradient flow of the Maximum Mean Discrepancy
 324 (MMD) functional and provide particle-based implementations. Additionally, Mroueh et al. (2019)
 325 propose Sobolev Descent, deriving a particle algorithm from the Sobolev integral probability metric
 326 that transports samples along smooth descent paths. More recently, there is been a surge in using

327 the seminal proximal point method of Jordan, Kinderlehrer, and Otto (JKO) (Jordan et al., 1998), to
 328 discretize the backward Wasserstein gradient flows over the chosen energy function F (Mokrov
 329 et al., 2021; Alvarez-Melis et al., 2021; Fan et al., 2022; Bonet et al., 2022; Bunne et al., 2022;
 330 Altekruiger et al., 2023). Central to almost all of these works, is the estimation of the transport
 331 map in JKO proximal step with input convex neural networks (Amos et al., 2017). In contrast to
 332 these works, our approach does not require parameterizing the transport map, the functionals that
 333 solve the variational formulations, or the underlying ODE/PDE induced by the continuity equation.
 334 Instead, it generates data via direct (inverse) gradient descent on a finite set of points without function
 335 estimation.

336 **Particle gradient descent.** Nitanda and Suzuki (2017) introduce "particle gradient descent" for
 337 optimizing a given energy function F over sparse measures as

$$\min_{x_1, \dots, x_n} F \left(\frac{1}{n} \sum_{i=1}^n \delta(x_i) \right), \quad \delta(x) : \text{the Dirac measure at } x. \quad (9)$$

338 Particle gradient descent is the gradient descent optimizing the location of x_1, \dots, x_n . Single layer
 339 neural networks can be viewed as particle gradient descent on a specific energy functions. Motivated
 340 by this connection, Chizat and Bach (2018) study the connection between gradient descent on
 341 particles and gradient flow in the space of probability measure equipped with Wasserstein-2 metric
 342 in asymptotic regime $n \rightarrow \infty$. Daneshmand et al. (2023) establishes a non-asymptotic convergence
 343 analysis for particle gradient descent on displacement convex functions. While primary focus of these
 344 studies is on analyzing single-layer neural networks, we leverage (inverse) particle gradient descent
 345 to develop an estimation-free generative method.

346 **Physics-inspired methods.** Xu et al. (2022) developed a generative model that maps samples from
 347 a uniform distribution over an infinite-radius hemisphere to an arbitrary target distribution. Their
 348 model relies on estimating the *Poisson vector field* parameterized by neural networks. This idea,
 349 primarily inspired by physical systems—especially electrostatic theory—led to further developments
 350 in subsequent works (Xu et al., 2023; Kolesov et al., 2025).

351 To cope with the challenges of sampling from an infinite-radius hemisphere, Xu et al. (2022) simulate
 352 the corresponding ODE by perturbing the training data and then estimate the negative normalized field
 353 from these perturbed samples. This contrasts with our approach, in which the primitive distribution is
 354 uniform on a finite-dimensional compact manifold (the sphere), and is therefore easy to sample from.
 355 Although we also transport samples from a uniform distribution to an arbitrary data distribution,
 356 our transport mechanism does not rely on function estimation; rather, it is purely based on (inverse)
 357 gradient descent over a finite set of points.

358 C Discussions

359 We demonstrate, both theoretically and experimentally, that it is possible to transport a uniform
 360 distribution to a target distribution without estimating a score function—using only i.i.d. samples
 361 from the target. Our method, Estimation-Free Sampling (EFS), avoids both noise injection and
 362 function estimation. Instead, it introduces interactions between samples by optimizing an energy
 363 functional, thereby shaping the empirical distribution of samples. EFS lies at the intersection of three
 364 foundational fields—mathematical optimization, potential theory, and generative AI—offering rich
 365 opportunities for interdisciplinary research.

366 **Optimization for Sampling.** While most generative models are built upon stochastic differential
 367 equations, EFS is purely a *deterministic optimization method*, opening new avenues for the optimiza-
 368 tion community to use powerful optimization techniques—such as accelerated methods, higher-order
 369 optimization methods, and efficient stochastic techniques—to generate data.

370 Notably, EFS involves both convex and non-convex optimization. Linking the non-convex gradient
 371 descent dynamics to a well-studied Wasserstein gradient flow, we establish asymptotic convergence
 372 guarantees. However, understanding the behavior of gradient descent in non-asymptotic regimes
 373 remains an open challenge. We believe the introduced asymptotic results provides a solid foundation
 374 for future non-asymptotic analyses.

375 **Potential Theory.** EFS induces attractive-repulsive interactions between training samples using
376 a well-known potential function studied extensively in fractional potential theory (Shu and Wang,
377 2021a; Frank et al., 2025; Duerinckx and Serfaty, 2020). By linking generative modeling to this rich
378 theoretical framework, we gain access to powerful analytical tools for understanding and designing
379 generative models. In particular, our analysis leverages results on Wasserstein gradient flows of
380 attractive-repulsive energies (Duerinckx and Serfaty, 2020), as well as variational analysis of these
381 energy functionals (Carrillo and Shu, 2023; Frank and Matzke, 2025). This bridge between theory
382 and practice opens new directions for developing principled and reliable generative models with
383 theoretical guarantees.

384 Potential theory can design new interaction potentials tailored to data generation and practical
385 applications. As noted in our experiments, one major challenge was the choice of the power s in the
386 potential function, which led to numerical instability in high dimensions. Investigating alternative
387 potentials that avoid such issues—particularly those that do not involve dimension-dependent blow-
388 up—may help resolve computational challenges of EFS.

389 **Generative AI.** EFS opens up several promising directions for future research, particularly in
390 scaling to large-scale machine learning applications. A key practical challenge lies in the quadratic
391 time complexity of the forward optimization step with respect to the number of training samples. We
392 conjecture that stochastic optimization techniques could mitigate this issue and significantly reduce
393 computational cost.

394 D Proofs

395 **Notations.** $\mathbb{S}^{d-1}(c, r)$ denotes sphere with center $c \in \mathbb{R}^d$ and radius r . Ω denotes the set of
396 probability measure over \mathbb{R}^d with Wasserstein-2 metric, which is denoted by $\mathcal{W}_2(\mu, \nu)$ where
397 $\mu, \nu \in \Omega$. Given vector function $v : \mathbb{R}^d \rightarrow \mathbb{R}^d$, its *divergence* is defined as $\operatorname{div}(v) = \sum_{i=1}^d \frac{dv_i}{dx_i}$.
398 Suppose $E : \Omega \rightarrow \mathbb{R}$. Then, the *first variation* of E with respect to μ is denoted by $\frac{dE}{d\mu}$ (Santambrogio,
399 2017). Function $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$ transport $\mu \in \Omega$ to $\nu \in \Omega$ if for a random vector $x \in \mathbb{R}^d$ drawn from
400 μ , $f(x)$ has density ν . Recall $W_\epsilon^{(s)}$ is the potential function defined as

$$W_\epsilon^{(s)}(x - y) = \frac{\|x - y\|^2}{2} + \frac{1}{s(\|x - y\|^2 + \epsilon)^{s/2}}. \quad (10)$$

401 For simplicity, we use the shorthand notation $W^{(s)} = W_0^{(s)}$ and $W^{(s)}$ by W . The sign $*$ denotes the
402 standard convolution as

$$f * g(y) = \int f(y - x)g(x)dx \quad (11)$$

403 $\|f\|_{L^2}$ is the L_2 functional norm defined as

$$\|f\|_{L^2}^2 = \int \|f(x)\|^2 dx \quad (12)$$

404 **Weak/Strong convergence.** A sequence of measures $\mu^{(n)} \in \Omega$ is said to converge weakly to μ if,
405 for all bounded continuous test functions $f : \mathbb{R}^d \rightarrow \mathbb{R}$, the integrals $\int f(x)\mu^{(n)}(x) dx$ converge to
406 $\int f(x)\mu(x) dx$. A stronger notion is convergence in the $L^2(\nu)$ norm, which requires that

$$\lim_{n \rightarrow \infty} \int \|\mu^{(n)}(x) - \mu(x)\|^2 \nu(x) dx = 0.$$

407 **MMD.** Define $\operatorname{MMD} : \Omega \times \Omega \rightarrow \mathbb{R}_+$ as

$$\operatorname{MMD}^2(\mu, \nu) = \int K(x - y)(\mu(x) - \nu(x))dx(\mu(y) - \nu(y))dy, \quad \text{with } K(\Delta) := \frac{1}{s\|\Delta\|^s} \quad (13)$$

408 where MMD denotes the Maximum Mean Discrepancy between two measures μ and ν .

409 Suppose that \widehat{K} and $\widehat{\Delta}$ denote the Fourier transforms of the functions $K(x)$ and $\Delta(x) := \mu(x) - \nu(x)$,
 410 respectively. Viewing MMD as a convolution and Plancherel's theorem allow us to write MMD
 411 as (Shu and Wang, 2021b)

$$\text{MMD}^2(\mu, \nu) = \int \widehat{K}(w) |\widehat{\Delta}(w)|^2 dw, \quad (14)$$

412 where $\widehat{K}(w) = C \|w\|^{-d+s}$ with a constant C depending only on d and s (Frank and Matzke, 2025).
 413 Substituting the Fourier transform into the equation above establishes that $\text{MMD}(\mu, \nu) = 0$ implies
 414 $\mu = \nu$. Notably, it is sufficient for the Fourier transform \widehat{K} to be strictly positive almost everywhere
 415 to guarantee that K is a *universal kernel* (Gretton et al., 2012).

416 To avoid potential confusion, note that the MMD associated with the kernel K defined above is not a
 417 metric. While machine learning often relies on positive semi-definite kernels to ensure that MMD
 418 defines a valid metric, the kernel K is not positive semi-definite. As a result, the corresponding MMD
 419 does not obey the triangle inequality.

420 **Gradient dominance of E .** We establish an important property of the energy function E . Define
 421 functional $F : \Omega \rightarrow L_2$ as

$$F(\mu)(y) = \nabla \frac{dE}{d\mu} = \int W(y-x) \mu(x) dx. \quad (15)$$

422 Remarkably, μ is a steady state of the Wasserstein gradient flow on E if $F(\mu) = 0$. The next theorem
 423 represent MMD using F .

424 **Theorem 5.** *Suppose that $\mu \in \Omega$ and $\nu \in \Omega$ have the same first-moment, then*

$$\|F(\mu) - F(\nu)\|_{L_2}^2 = \text{MMD}^2(\mu, \nu)$$

425 *holds for $s = d - 2$.*

426 An application of the above theorem recovers the result of Carrillo and Shu (2023): all steady states
 427 μ satisfying $F(\mu) = 0$ are global minimizers of E (up to translation). More precisely, the theorem
 428 implies that for all μ, ν such that $F(\mu) = F(\nu) = 0$, the following holds:

$$0 = \|F(\mu) - F(\nu)\|_{L_2}^2 = \text{MMD}^2(\mu, \nu). \quad (16)$$

429 Beyond recovering existing results, the theorem will also be used to analyze the backward optimization
 430 for EFS in Section D.1.

431 *Proof.* Let $\widehat{f}(w)$ denote the Fourier transform of function $f(x)$. Define $\widehat{\Delta} := \widehat{\mu} - \widehat{\nu}$ which is
 432 equivalent to the Fourier transform of $\Delta := \mu - \nu$. As discussed in Section D, MMD can be written
 433 as

$$\text{MMD}^2(\nu, \mu) = \int \widehat{K}(w) |\widehat{\Delta}(w)|^2 dw \quad (17)$$

434 According to the definition,

$$\nabla W(z) = \nabla K(z) + \text{idem}(z), \quad \text{idem}(z) := z \quad (18)$$

435 Replacing the above formula into the definition of F obtains

$$\|F(\mu) - F(\nu)\|_{L_2}^2 = \|(\nabla K + \text{idem}) * \mu - (\nabla K + \text{idem}) * \nu\|_{L_2}^2, \quad (19)$$

436 where $*$ denotes convolution defined in (11) and L_2 is the norm 2 for functions defined in (12). It is
 437 easy to check that

$$(\text{idem} * \mu)(y) = y \underbrace{\int \mu(x) dx}_{=1} - \int x \mu(x) dx \quad (20)$$

$$= y \underbrace{\int \nu(x) dx}_{=1} - \underbrace{\int x \mu(x) dx}_{\mathbb{E}_\mu[x]} \quad (21)$$

$$= y \int \nu(x) dx - \underbrace{\int x \nu(x) dx}_{\mathbb{E}_\nu[x]} = (\text{idem} * \nu)(y). \quad (22)$$

438 The last equation holds because the first moments of μ and ν are equal. This observation allows us to
 439 significantly simplify the expression for the quantity of interest as follows:

$$\|F(\mu) - F(\nu)\|_{L_2}^2 = \|\nabla K * \mu - \nabla K * \nu\|_{L_2}^2 \quad (23)$$

440 Recall two fundamental properties of Fourier transform as: $\begin{cases} \widehat{f * g} = \widehat{f} \widehat{g} \\ \widehat{\partial f} = iw \widehat{f} \end{cases}$. These properties
 441 together with Parseval's theorem yield

$$\|F(\mu) - F(\nu)\|_{L_2}^2 = \int \|w\|^2 \left(\widehat{K}(w)\right)^2 |\widehat{\Delta}(w)|^2 dw \quad (24)$$

442 For $s = d - 2$, it is easy to check that

$$\|w\|^2 (\widehat{K}(w))^2 = \widehat{K}(w) \quad (25)$$

443 holds given $\widehat{K}(w) = C\|w\|^{-2}$ (Frank and Matzke, 2025). Combining all the results completes the
 444 proof:

$$\|F(\mu) - F(\nu)\|_{L_2}^2 \stackrel{(24)}{=} \int \|w\|^2 \left(\widehat{W}(w)\right)^2 |\widehat{\Delta}(w)|^2 dw \quad (26)$$

$$\stackrel{(25)}{=} \int \widehat{W}(w) |\widehat{\Delta}(w)|^2 dw \quad (27)$$

$$\stackrel{(17)}{=} \text{MMD}^2(\mu, \nu) \quad (28)$$

445 \square

446 **A convergence result for ODEs.** Consider the following two differential equations

$$\begin{cases} \frac{dy}{dt} = F_t(y_t) \\ \frac{dy^{(n)}}{dt} = F_t^{(n)}(y^{(n)}) \end{cases}, \quad y_0 = y_0^{(n)} \quad (29)$$

447 where $F_t, F_t^{(n)} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ and $y_t, y_t^{(n)} \in \mathbb{R}^d$. If $F_t^{(n)}$ converges to F_t as $n \rightarrow \infty$, then $y_t^{(n)}$
 448 converges to y_t in interval $t \in [0, T)$.

449 **Lemma 6.** *Define*

$$\epsilon_n := \int \|F_t(y) - F_t^{(n)}(y)\|^2 dy.$$

450 *Suppose $F_t(y)$ is almost surely L -Lipschitz in y . If $\lim_{n \rightarrow \infty} \epsilon_n = 0$, then $\int_0^T \|y_t - y_t^{(n)}\|^2 dt$ converges
 451 to 0 as $n \rightarrow \infty$.*

452 *Proof.* We start by writing down ODEs in the integral form and consider their difference,

$$\begin{aligned} \|y^{(n)}(t) - y(t)\| &= \left\| \int_0^t \left[F_s^{(n)}(y^{(n)}) - F_s(y) \right] ds \right\| \\ &= \left\| \int_0^t \left[F_s^{(n)}(y^{(n)}) - F_s(y^{(n)}) \right] ds + \int_0^t \left[F_s(y^{(n)}) - F_s(y) \right] ds \right\| \\ &\leq t \sup_y \left\| F_s^{(n)}(y) - F_s(y) \right\| + L \int_0^t \|y^{(n)}(s) - y(s)\| ds \\ &\leq t\epsilon_n + L \int_0^t \|y^{(n)}(s) - y(s)\| ds \end{aligned} \quad (30)$$

453 The last inequality was followed by convergence of $F_s^{(n)}$ to F_s , and boundedness of L^∞ norm by L^2 ,
 454 the first term is converging to zero, and bounded by $\epsilon_n \xrightarrow{n \rightarrow \infty} 0$.

455 Define $\Delta_n(t) := \sup_{s \leq t} \|y^{(n)}(s) - y(s)\|$. By taking supremum from (30),

$$\Delta_n(t) \leq t\epsilon_n + L \int_0^t \Delta_n(s) ds.$$

456 In turn, applying Grönwall's inequality yields,

$$\Delta_n(t) \leq t\epsilon_n e^{Lt} \xrightarrow{n \rightarrow \infty} 0.$$

457 Thus,

$$\left\| y^{(n)}(s) - y(s) \right\| \xrightarrow{n \rightarrow \infty} 0,$$

458 uniformly for all $s \leq t$ and proof is completed. \square

459 **Continuity equation and transporting distributions.** Let $v_t : \mathbb{R}^d \rightarrow \mathbb{R}^d$ be a general vector field.
460 $\mu_t \in \Omega$ is a solution of continuity equation associated with v_t if it obeys

$$\frac{d\mu_t}{dt} = -\operatorname{div}(\mu_t v_t) \quad (31)$$

461 Given the vector field, we define the following ODE

$$\frac{dy_t}{dt} = v_t(y_t) \quad (32)$$

462 The above ODE transports μ_0 to μ_t as stated in the following lemma.

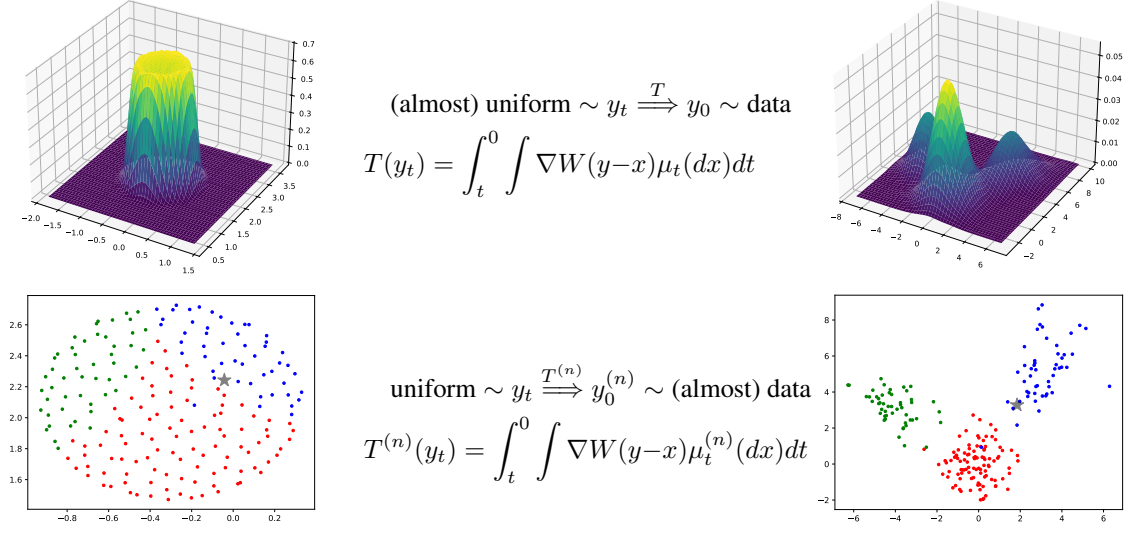
463 **Lemma 7** (Lemma 8.1.6. of (Ambrosio et al., 2008)). *Suppose that the vector field v_t obeys the*
464 *following 3 conditions*

465 (1) $\int |v_t(x)| d\mu_t(x) dx < \infty.$

466 (2) *For every compact subset $B \subset \mathbb{R}^d$, $\sup_{x \in B} |v_t(x)| < \infty.$*

467 (3) $v_t(x)$ is Lipschitz

468 *If y_0 is a random variable with distribution μ_0 , then y_t is a random variable with distribution μ_t for*
469 *a finite t .*



Proof sketch for Theorem 4. Given the solution to the continuity equation (6), we define the map T above, which provably transports μ_t to data distribution, as stated in Lemma 8. Recall that Theorem 2 states that μ_t converges to the uniform distribution, from which new samples can be drawn. T is not implementable as it requires μ_t , the solution of the continuity equation (6). We show that the empirical distribution $\mu_t^{(n)}$, defined over the point set $\{x_1(t), \dots, x_n(t)\}$, yields a transport map $T^{(n)}$ (as defined above) that converges to T as $n \rightarrow \infty$. The \star symbol in the second row indicates y_t (left) and $y_0^{(n)}$ (right). Colored points represent $\{x_1(t), \dots, x_n(t)\}$ on the left and $\{x_1(0), \dots, x_n(0)\}$ on the right.

471 **Recap.** Recall that $\mu_t^{(n)}$ defined as $\mu_t^{(n)} = \frac{1}{n} \sum_{i=1}^n \delta_{x_i(t)}$ where δ is the Dirac measure and $x_i(t)$
472 obeys

$$\frac{dx_i}{dt} = -\frac{1}{n} \sum_{j \neq i} \nabla W^{(s)}(x_i - x_j). \quad (33)$$

473 Using $x_i(t)$, we define $y_t^{(n)}$ which obeys

$$\frac{dy^{(n)}}{d\bar{t}} = F_t^{(n)}(y_t^{(n)}), \quad F_t^{(n)}(y) := \int \nabla W(y-x) \mu_t^{(n)}(x) dx, \quad (34)$$

474 where $d\bar{t}$ indicates that the process is reversed time (Anderson, 1982).

475 **Reversed-time transport.** Recall μ_t as the solution of continuity equation for Wasserstein-2
476 gradient flow:

$$\frac{d\mu}{dt} = \text{div} \left(\mu(x) \int \nabla W^{(s)}(x-y) \mu_t(y) dy \right), \quad (35)$$

477 Notably, $F_t(y) := \int \nabla W(y-x) \mu_t(x) dx$ represents the vector field associated with the above
478 dynamics. Given this vector field, define the following ODE

$$\frac{dy}{dt} = -F_t(y_t), \quad \text{where } F_t(y) := \int \nabla W(y-x) \mu_t(x) dx \quad (36)$$

479 The above ODE transports μ_t backward to μ_0 as stated in the following lemma.

480 **Lemma 8.** Let y_0 be a random vector drawn from μ_0 . Then, the distribution of y_t is μ_0 .

481 The above result together with a simple change of variables obtains the inverse transport from μ_t to
 482 μ_0 .

483 **Corollary 9.** Consider the following reversed-time dynamics:

$$\frac{dy_t}{dt} = F_t(y_t) \quad (37)$$

484 The above dynamics transports μ_t to μ_0 as long as μ_t is not a steady state of (35).

485 *Proof.* According to definition

$$\lim_{\epsilon \rightarrow 0} \frac{y_{t+\epsilon} - y_t}{\epsilon} = -F_t(y_t) \implies \frac{dy}{dt} = \lim_{\epsilon \rightarrow 0} \frac{y_{t-\epsilon} - y_t}{\epsilon} = \lim_{\epsilon \rightarrow 0} F_{t-\epsilon}(y_{t-\epsilon}) \quad (38)$$

486 where $F_{t-\epsilon}(y) = \int \nabla W(y-x) \mu_{t-\epsilon}(x) dx$. Lemma 10 implies F_t is Lipschitz. Thus, invoking
 487 Lemma 6 concludes the statement. \square

488 While the flow y_t transports μ_t to the data distribution, its construction relies on access to μ_t , which
 489 is defined as the solution to a PDE. Since μ_t is not feasible to compute, we cannot implement the
 490 reverse-time ODE (36) to recover y_0 . However, Theorem 4 restated below ensures that EFS can
 491 reconstruct the backward dynamic as $n \rightarrow \infty$.

492 **Theorem 4.** Assume $\mu_0^{(n)}$ converges to μ_0 in Wasserstein-2 distance. Suppose that $y_t^{(n)}$ is a random
 493 variable with law μ_t , where μ_t is the solution to the continuity equation (6). Then, the distribution of
 494 $y_0^{(n)}$ —obtained from the reversed-time ODE (7),

$$\frac{dy^{(n)}}{dt} = \frac{1}{n} \sum_{i=1}^n \nabla W^{(s)}(y^{(n)} - x_i)$$

495 converges to μ_0 as $n \rightarrow \infty$, for regular measures μ_0 and μ_t , and for $s = d - 2 > 0$.

496 *Proof.* We prove that $y_0^{(n)}$ converges to y_0 , whose distribution is μ_0 , according to Lemma 8. To
 497 establish this convergence, we use Theorem 1 of (Duerinckx and Serfaty, 2020), which shows that
 498 $\mu_t^{(n)}$ converges to μ_t in MMD as $n \rightarrow \infty$, namely

$$\lim_{n \rightarrow \infty} \text{MMD}(\mu_t, \mu_t^{(n)}) = 0, \quad \text{for } \beta > 0. \quad (39)$$

499 As proven in the next lemma, $F_t(y)$ is Lipschitz. This Lipschitz continuity, together with the result
 500 above, ensures that the conditions in Lemma 6 are satisfied. Thus, $y_0^{(n)}$ converges to y_0 , as guaranteed
 501 by the lemma. \square

502 **Lemma 10.** Recall vector function $F_t(y) = \int \nabla W^{(d-2)}(y-x) \mu_t(x) dx$ where μ_t is the solution
 503 of PDE (35) starting from μ_0 . Then, $F_t(y)$ is L -Lipschitz for $L = c_1 + c_2(E(\mu_0) + 1)^{\frac{2d+2}{d}}$ where
 504 constants c_1 and s_2 depend only on d .

505 *Proof.* The first step is to cast Lipschitzness to the boundedness of the following integral

$$\int \frac{1}{\|y-x\|^d} \mu_t(x) dx \leq L \implies F_t \text{ is } L\text{-Lipschitz} \quad (40)$$

506 Then, we use properties of Wasserstein gradient flow to prove

$$\int \frac{1}{\|y-x\|^d} \mu_t(x) dx \leq c_1 + c_2 E(\mu_0), \quad (41)$$

507 where c_1 and c_2 are constants that only depend on d . For the similarity, we introduce notation $a \lesssim b$
 508 if there are constants c_1 and c_2 that only depends on d such that $a \leq c_1 + c_2 b$

509 *Step 1. A Sufficient Condition for Lipschitz-ness.* According to the definition, we have

$$\|F_t(y) - F_t(y')\|_2 = \left\| \int \nabla K(y-x)\mu_t(dx) - \int \nabla K(y'-x)\mu_t(dx) \right\|_2 \quad (42)$$

$$\leq \|y - y'\| \int_0^1 \nabla^2 K(\gamma y + (1-\gamma)y' - x)\mu_t(dx) d\gamma \quad (43)$$

$$\leq \|y - y'\|_2 \max_y \left\| \int \nabla^2 K(y-x)\mu_t(dx) \right\|_2 \quad (44)$$

510 To prove that the right side is bounded, we bound the norm of $\nabla^2 K$ as

$$\left\| \int \nabla^2 K(y-x)\mu_t(dx) \right\| = \left\| \int -\frac{I_d}{\|y-x\|^d} + (d)\frac{(y-x)(y-x)^\top}{\|y-x\|^{d+2}}\mu_t(dx) \right\| \quad (45)$$

$$\leq (d) \int \frac{1}{\|y-x\|^d} \mu_t(x) dx \quad (46)$$

511 Thus, the above integral bounds the Lipschitz constant of F_t .

512 *Step II. Gradient flow property.* Since μ_t is a gradient flow, the following holds

$$E(\mu_t) \leq E(\mu_0) \quad (47)$$

513 We establish an important consequence of the above inequality, namely the following bound

$$\sup_y \left| \int \frac{1}{\|y-x\|^{d-1}} \mu_t(x) dx \right| \lesssim (1 + E(\mu_0))^2 \quad (48)$$

514 holds. Using Jensen's inequality, we have

$$\int \frac{1}{\|x-y\|} d\mu_t(x) \leq \left(\int \frac{1}{\|x-y\|^{d-1}} d\mu_t(x) \right)^{\frac{1}{d-1}} \quad (49)$$

515 Combining the last two bounds obtains

$$\int \frac{1}{\|y-x\|^d} \mu_t(x) dx \leq \left| \int \frac{1}{\|y-x\|^{d-1}} \mu_t(x) dx \right| \left| \int \frac{1}{\|y-x\|} \mu_t(x) dx \right| \lesssim (1 + E(\mu_0))^{\frac{2d+2}{d}} \quad (50)$$

516 The above equation concludes that F_t is Lipschitz according to step I. To complete the proof, we
517 need to prove equation (48).

518 We first establish a consequence of $E(\mu_t) \leq E(\mu_0)$ for μ_t . Let $\widehat{\mu}(w)$ denote Fourier transform of
519 $\mu_t(x)$. It is easy to check that

$$|\widehat{\mu}(w)| \leq \int |\mu_t(x)| dx = 1 \quad (51)$$

520 Without loss of generality, we can assume that $\int x\mu_t(x) dx = 0$. An application of Jensen's
521 inequality yields

$$\int \|x\|^2 d\mu_t(x) = \int \|x - \int y\mu_t(y) dy\|^2 d\mu_t(x) \quad (52)$$

$$\leq \int \|x-y\|^2 d\mu_t(x) d\mu_t(y) dx dy \quad (53)$$

$$\lesssim E(\mu_t) \lesssim E(\mu_0) \quad (54)$$

522 According properties of Fourier transform, we have

$$\left| \frac{\partial \widehat{\mu}}{\partial w_i \partial w_j} \right| = \int |x_i x_j| \mu_t(x) dx \quad (55)$$

$$\leq \frac{1}{2} \left(\int x_i^2 \mu_t(x) dx + \int x_j^2 \mu_t(x) dx \right) \quad (56)$$

$$\leq \int \|x\|^2 \mu_t(x) \quad (57)$$

$$\stackrel{(54)}{\leq} E(\mu_0) \quad (58)$$

523 Gagliardo–Nirenberg interpolation inequality yields

$$\|\widehat{\mu}\|_{L^1} \lesssim (\|\partial^2 \widehat{\mu}\|_{L^\infty} + 1)(\|\widehat{\mu}\|_{L^\infty} + 1) \quad (59)$$

$$\stackrel{(58)}{\lesssim} (1 + E(\mu_0))(\|\widehat{\mu}\|_{L^\infty} + 1) \quad (60)$$

$$\stackrel{(51)}{\lesssim} (1 + E(\mu_0))^2 \quad (61)$$

524 We will use the above bound to complete the proof.

525 Recall the Fourier transform of the radial function $\|x\|^{-s}$ is $C\|w\|^{-d+s}$ for $0 < s < d$ where C
 526 is a constant depending on d and s (Frank and Matzke, 2025). Using the Parseval's theorem, (47)
 527 translates to the following inequality in the complex (Fourier) domain

$$C \int \frac{|\widehat{\mu}(w)|^2}{\|w\|^2} dw = \int \frac{1}{\|y-x\|^{d-2}} d\mu_t(x) d\mu_t(y) \leq (d-2)E(\mu_0) \quad (62)$$

528 holds. Similarly, we can write (48) as

$$\int \frac{1}{\|w\|} |\widehat{\mu}(w)| dw \lesssim (1 + E(\mu_0))^2 \quad (63)$$

529 Define set $A = \{w \mid \|w\| < 1\}$, a straightforward application of Cauchy-Schwarz obtains

$$\left(\int_A \|w\|^{-1} |\widehat{\mu}(w)| dw \right)^2 \leq \int_A \|w\|^{-2} |\widehat{\mu}(w)|^2 dw \int_A dw \quad (64)$$

$$\stackrel{(62)}{\lesssim} E(\mu_0) \quad (65)$$

530 To establish (63), we need to bound the above integral taken over the complement of A denoted by
 531 A^c :

$$\int_{A^c} \frac{1}{\|w\|} |\widehat{\mu}(w)| dw \leq \int |\widehat{\mu}(w)| dw \stackrel{(61)}{\lesssim} (1 + E(\mu_0))^2 \quad (66)$$

532 Combining the last two inequality concludes (63), hence the proof is complete.

533

□

534 D.2 Proof of the Auxiliary Lemma 8

535 **Lemma 8.** *Let y_0 be a random vector drawn from μ_0 such that $E(\mu_0) < \infty$. Consider the following*
 536 *ODE*

$$\frac{dy}{dt} = -F_t(y_t), \quad F_t(y) := \int \nabla W(y-x) \mu_t(x) dx dt$$

537 *Then, the distribution of y_t is μ_t .*

538 *Proof.* If the vector field F_t obeys 3 conditions in Lemma 7, then invoking the lemma concludes the
 539 proof. It remains to validate necessary conditions for the vector field. Since μ_t is a gradient flow, it
 540 obeys

$$\int \frac{1}{(d-2)\|x-y\|^{d-2}} d\mu_t(x) d\mu_t(y) + \frac{1}{2} \int \|x-y\|^2 d\mu_t(x) d\mu_t(y) = E(\mu_t) \leq E(\mu_0) \quad (67)$$

541 **Condition 1:** $\int \|F_t(x)\| \mu_t(x) dx < \infty$. According to the definition, we have

$$\|F_t(x)\| \mu_t(x) dx = \int \frac{x-y}{\|x-y\|^d} + (x-y) d\mu_t(y) \|\mu_t(x) dx \quad (68)$$

$$\leq \int \frac{1}{\|x-y\|^{d-1}} d\mu_t(x) d\mu_t(y) + \int \|x-y\| d\mu_t(x) d\mu_t(y) \quad (69)$$

542 We bound each term in the above upper-bound. The first term can be bounded using Jensen's
 543 inequality as

$$\int \|x - y\| d\mu_t(x) d\mu_t(y) \leq \sqrt{\int \|x - y\|^2 d\mu_t(x) d\mu_t(y)} \stackrel{(67)}{\leq} \sqrt{E(\mu_0)} \quad (70)$$

544 Similarly, we get

$$\int \frac{1}{\|x - y\|} d\mu_t(x) d\mu_t(y) \leq \left(\int \|x - y\|^2 d\mu_t(x) d\mu_t(y) \right)^{1/(d-2)} \stackrel{(67)}{\leq} (E(\mu_0))^{1/(d-2)} \quad (71)$$

545 The above inequality yields

$$\int \frac{1}{\|x - y\|^{d-1}} d\mu_t(x) d\mu_t(y) \leq \left(\int \frac{1}{\|x - y\|^{d-2}} d\mu_t(x) d\mu_t(y) \right) \int \frac{1}{\|x - y\|} d\mu_t(x) d\mu_t(y) \quad (72)$$

$$\leq E(\mu_0)^{(d-1)/(d-2)} \quad (73)$$

546 **Condition 2.** $\|F_t(x)\| < \infty$ for a bounded x . Since all statements hold up to translation, we can
 547 assume $\int y \mu_t(y) dy = 0$. This allows us to simplify the expression for the vector field as

$$\|F_t(x)\| = \left\| \frac{x - y}{\|x - y\|^d} d\mu_t(y) + \underbrace{x - \int y \mu_t(y)}_{=0} \right\| \quad (74)$$

$$\leq \int \frac{1}{\|x - y\|^{d-1}} d\mu_t(y) + \|x\| \quad (75)$$

$$\stackrel{(48)}{\leq} c_1 + c_2(1 + E(\mu_0)^2) + \|x\|, \quad (76)$$

548 where c_1 and c_2 are constants independent from t .

549 **Condition 3. Lipschitzness.** Lemma 10 ensures F_t is Lipschitz. □

550 **E Inverse Gradient Descent**

551 **Proposition 11.** *Consider the following proximal optimization problem*

$$(y_1^*, y_2^*, \dots, y_n^*) = \arg \min_{y_1, y_2, \dots, y_n \in \mathbb{R}^d} \frac{1}{2} \sum_{i=1}^d \|y_i - x_i^{(k)}\|^2 - \gamma E_n^{(\epsilon)}(y_1, y_2, \dots, y_n) \quad (77)$$

552 *The above optimization is convex with solution $y_i^* = x_i^{(k-1)}$ for all $i \in [n]$, as long as the learning*
 553 *rate γ is sufficiently small.*

554 *Proof.* Let the objective of the proximal step (77) defined as,

$$H(y_1, \dots, y_n) := \frac{1}{2} \sum_{i=1}^n \|y_i - x_i^{(k)}\|^2 - \gamma E_n^{(\epsilon)}(y_1, \dots, y_n). \quad (78)$$

555 And for ease of notation, let $Y = (y_1, \dots, y_n) \in (\mathbb{R}^d)^n$. It is easy to observe that,

$$\nabla^2 H(Y) = I_{n \times d} - \gamma \nabla^2 E_n^{(\epsilon)}(Y).$$

556 Thus, the optimization problem (11) is convex as long as we can prove an upper bound on the spectral
 557 norm of $\nabla^2 E_n^{(\epsilon)}(Y)$ and choose a small learning rate for the inverse gradient map γ . From this point
 558 onward, we drop the subscript/superscript ϵ to streamline the notation. Thus W is defined as

$$W(z) = -\frac{m}{(\|z\|^2 + \epsilon)^{m/2}} + \frac{1}{2}\|z\|^2$$

559 with direct calculations we have,

$$\nabla W(z) = (1 + m^2(r^2 + \epsilon)^{-\frac{m}{2}-1})z,$$

560 where $r := \|z\|$. Again with differentiation,

$$\nabla^2 W(z) = (1 + m^2(r^2 + \epsilon)^{-\frac{m}{2}-1})I_d - m^2(m+2)(r^2 + \epsilon)^{-\frac{m}{2}-2}zz^\top.$$

561 Hence, for every unit vector $\nu \in \mathbb{S}^d$,

$$\begin{aligned} \nu^\top \nabla^2 W(z) \nu &\leq 1 + m^2(r^2 + \epsilon)^{-\frac{m}{2}-1} \\ &\leq 1 + m^2 \epsilon^{-\frac{m}{2}-1}. \end{aligned} \quad (79)$$

562 Recall that,

$$E_n(Y) = \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j \neq i} W(y_i - y_j).$$

563 Taking derivatives, for each i we have,

$$\nabla_{y_i} E_n(Y) = \frac{1}{n(n-1)} \sum_{j \neq i} \nabla W(y_i - y_j).$$

564 Next, for $j \neq i$,

$$\nabla_{y_i, y_j}^2 E_n(Y) = \frac{-1}{n(n-1)} \nabla^2 W(y_i - y_j)$$

565 and,

$$\nabla_{y_i}^2 E_n(Y) = \frac{1}{n(n-1)} \sum_{j \neq i} \nabla^2 W(y_i - y_j).$$

566 Taking an arbitrary block vector $\nu \in (\mathbb{R}^d)^n$,

$$\begin{aligned} \langle \nu, \nabla^2 E_n(Y) \nu \rangle &= \frac{1}{n(n-1)} \left(\sum_{i=1}^n \nu_i^\top \left(\sum_{j \neq i} \nabla^2 W(y_i - y_j) \right) \nu_i + \sum_{i \neq j} \nu_i^\top (-\nabla^2 W(y_i - y_j)) \nu_i \right) \\ &= \frac{1}{n(n-1)} \left(\frac{1}{2} \sum_{i \neq j} (\nu_i - \nu_j)^\top \nabla^2 W(y_i - y_j) (\nu_i - \nu_j) \right). \end{aligned}$$

567 In turn, by spectral bound on each block in (79),

$$\begin{aligned} \langle \nu, \nabla^2 E_n(Y) \nu \rangle &\leq \frac{1}{2n(n-1)} \left(\sum_{i \neq j} ((1 + m^2(\|y_i - y_j\|^2 + \epsilon)^{-\frac{m}{2}-1}) \|\nu_i - \nu_j\|^2) \right) \\ &\leq \frac{1}{2n(n-1)} \left(\sum_{i \neq j} (1 + m^2 \epsilon^{-\frac{m}{2}-1}) \|\nu_i - \nu_j\|^2 \right) \tag{80} \\ &= \frac{(1 + m^2 \epsilon^{-\frac{m}{2}-1})}{2n(n-1)} \left(\sum_{i \neq j} \|\nu_i - \nu_j\|^2 \right) \\ &\leq \frac{(1 + m^2 \epsilon^{-\frac{m}{2}-1})}{2n(n-1)} \left(2n \sum_i \|\nu_i\|^2 \right) \\ &= \frac{1 + m^2 \epsilon^{-\frac{m}{2}-1}}{n-1} \|\nu\|^2, \end{aligned}$$

568 where in (80) we used the (79).

569 Hence, for the operator norm of $\nabla^2 E_n(Y)$ we infer that,

$$\|\nabla^2 E_n(Y)\|_2 \leq \frac{1 + m^2 \epsilon^{-\frac{m}{2}-1}}{n-1}.$$

570 Thus, by choosing any sufficiently small $\gamma \in (0, \frac{n-1}{1+m^2 \epsilon^{-\frac{m}{2}-1}})$,

$$\nabla^2 H(Y) \succ 0,$$

571 and objective of the proximal step (78) is strongly convex and hence is uniquely minimized.

572 Considering the optimal solution $(y_1^*, y_2^*, \dots, y_n^*)$ to (78), by first order optimality condition,

$$\nabla_{y_i} H(y_1^*, y_2^*, \dots, y_n^*) = y_i^* - x_i^{(k)} - \gamma \nabla_{y_i} E_n(y_1^*, y_2^*, \dots, y_n^*) = 0.$$

573 Thus,

$$x_i^{(k)} = y_i^* - \gamma \nabla_{y_i} E_n(y_1^*, y_2^*, \dots, y_n^*),$$

574 and hence for all $i \in [n]$,

$$y_i^* = x_i^{(k-1)},$$

575 and proof is concluded. □

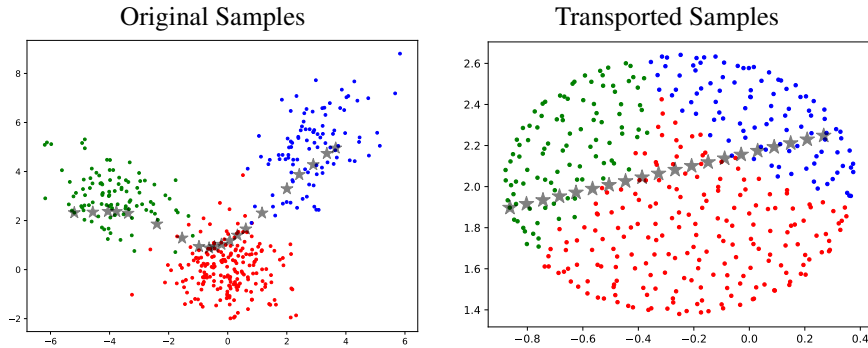


Figure 2: **Interpolation in Latent Space.** Colored points are samples drawn from a Gaussian mixture model (left) and their positions after forward optimization (right). The \ast s are the generated sample y in EFS (right), and their location after backward optimization in EFS (left). The straight line in the latent space (right) is transformed into a curved trajectory (left) that aligns with the data distribution. Colors correspond to different mixture components.

576 F Experiments

577 While we have proven EFS can transport an almost uniform distribution to an arbitrary target
 578 distribution, our result is asymptotic—it holds in the limit as $n \rightarrow \infty$. This asymptotic regime greatly
 579 simplifies the theoretical analysis. However, practical applications are inherently constrained to finite
 580 n . We bridge this gap with experiments.

581 As the first estimation-free generative algorithm, EFS needs further research for a comprehensive
 582 benchmarking. Our experiments underscore both the practical challenges and the promising potential
 583 of this novel generative method. The code is implemented in PyTorch (Paszke et al., 2019) and was
 584 executed on a single NVIDIA RTX 6000 GPU with 48GB memory. The total execution time is under
 585 15 minutes.

586 **Mixture of Gaussians: Capturing Data Geometry.** EFS effectively recovers the underlying data
 587 distribution even after transforming into a set of points uniformly distributed on a circle. Figure 2
 588 displays the original samples drawn from a Gaussian mixture (left) and their transformed positions
 589 after the EFS forward optimization step (right), where they are mapped onto the circle. New samples
 590 are then generated along a straight line on circle. Backward optimization maps these linear samples
 591 onto a curved trajectory that closely follows the geometry of the original data distribution. The
 592 backward process adaptively distorts pairwise distances to avoid low-density regions. Furthermore,
 593 we observe that EFS creatively generate unseen new samples.

594 **Generating samples from Swiss roll.** We repeat the experiments on the Swiss roll dataset, similar
 595 to those on the Gaussian mixture in Figures 1, and 2. Table 1 summarizes the parameter choices
 596 for this dataset. We evaluate: (i) convergence to a uniform distribution via forward optimization,
 597 (ii) generation of new samples via backward optimization, and (iii) preservation of data geometry
 598 through interpolation.

599 (i) *Forward optimization:* In Figure 3, we observe that the distribution of data points becomes
 600 uniform over the unit ball. Recall that Corollary 3 proves this result in the asymptotic regime
 601 as $n \rightarrow \infty$. Figure 3 shows that this asymptotic result provides a good approximation even
 602 when $n = 500$.

603 (ii) *Backward optimization:* We assert that EFS can generate new samples from the Swiss roll
 604 dataset using a given set of i.i.d. samples. While Theorem 4 establishes that EFS draws
 605 samples from the data distribution in the mean-field regime, Figure 4 supports this result in
 606 the non-asymptotic setting.

607 (iii) *Interpolation:* Figure 5 illustrates how backward optimization can generate new samples
 608 from the Swiss roll dataset via interpolation. Similar to the Gaussian mixture case shown in

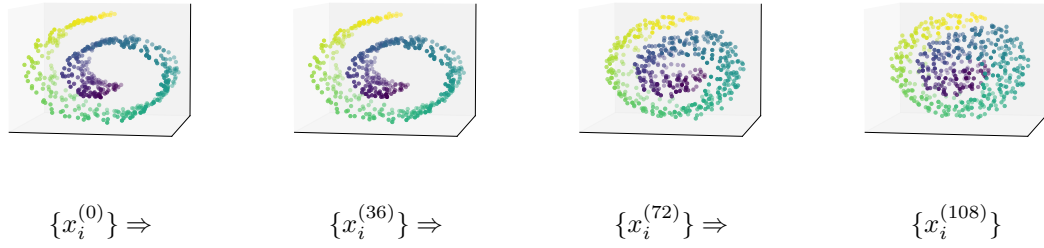


Figure 3: **Forward convergence for Swiss roll.** Scatter plot of gradient descent iterates $x_i^{(k)}$, defined in (3), initialized at $x_i^{(0)}$ drawn from a Swiss roll with noise level 0.2. As k evolves, the points become uniformly distributed akin to 2d example of Gaussian mixture in Figure 1. Notably, Corollary 3 establishes this convergence in the mean-field regime as $n \rightarrow \infty$.

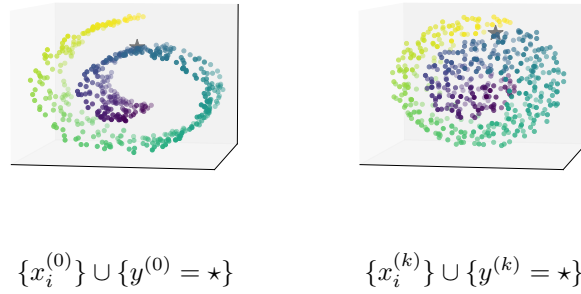


Figure 4: **Sampling:** The black \star is the generated sample $y^{(0)}$; Left: $x_i^{(0)} \stackrel{i.i.d.}{\sim}$ Swiss roll with noise level 0.2; Right: $x_i^{(k)}$ obtained by forward optimization. Colors: mixture components. The location of the new sample is updated by Backward Optimization in EFS

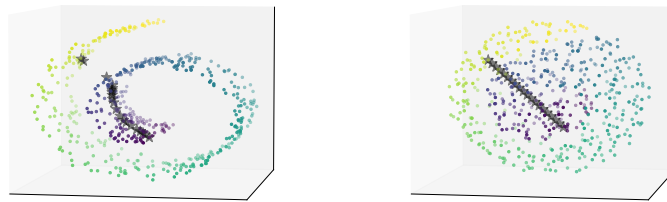


Figure 5: **Interpolation in Latent Space.** Colored points are samples drawn from a Swiss roll (left) and their positions after forward optimization (right). The \star s are the generated sample $y^{(0)}$ in EFS (right), and their location after applying backward step of EFS (left). The straight line in the latent space (right) is transformed into a curved trajectory (left) that aligns with the data distribution. We observe that the algorithm avoids generating data from no density regions.

609
610

Figure 2, we observe that EFS respects the data density and avoids generating samples in low-density regions.

611 **MNIST: Generating New Samples.** EFS can generate new MNIST images (LeCun et al., 1998).
 612 Theorem 4 requires high exponents in the range $d - 2 \leq m < d$. For MNIST with $d = 784$,
 613 this exponent leads to numerical overflow due to limited floating-point precision. To avoid this
 614 issue, we reduce the data dimension before applying EFS. Specifically, we use a convolutional
 615 autoencoder (Masci et al., 2011) to compress the data into a 15-dimensional latent space (see
 616 Appendix for architectural and training details). We then apply encoding \rightarrow EFS \rightarrow decoding to
 617 generate new samples. Figure 6 shows generated samples alongside their nearest neighbors (in
 618 Euclidean distance) from decoded original samples $x_1^{(0)}, \dots, x_n^{(0)}$. The generated digits exhibit
 619 stylistic variations distinct from their closest training examples, with minimum Euclidean distance of
 620 0.7508. To emphasize subtle differences, we include an animated image that alternates between the
 621 generated sample and its nearest neighbor (open with Adobe Reader). See Table 1 for details.

Generated Images	
Closest Image	
Animation (see with Adobe)	

Figure 6: **Generative MNIST.** *First row:* Images generated by EFS. These are obtained by applying forward and backward optimization in the latent space of an autoencoder, followed by decoding into the image space. *Second row:* The closest decoded training samples $\{x_i^{(0)}\}_{i=1}^n$ for each generated image, determined by Euclidean distance. *Third row:* An animation illustrating subtle differences between the generated images and their nearest neighbors (use Adobe Reader to see). *Note:* The minimum Euclidean distance between each generated image and its nearest neighbor is 0.75.

622 **MNIST: Interpolation Comparison.** Recall that in the previous experiment, we used an autoen-
 623 coder for dimensionality reduction. Notably, autoencoders can generate samples with interpolation in
 624 the latent space. To highlight differences, we compare the quality of the autoencoder samples with
 625 those produced by EFS.

626 For EFS, we generate new samples by linearly interpolating between uniformly distributed samples
 627 after forward optimization as $y^{(k)} = (1 - t)x_i^{(k)} + tx_j^{(k)}$. Then, we apply the backward optimization
 628 to $y^{(k)}$. Similarly, autoencoder samples are obtain by decoding $(1 - t)x_i^{(0)} + tx_j^{(0)}$ where $x_i^{(0)}$
 629 are encoded samples in the latent space. Figure 7 shows generate samples with these two different
 630 strategies. Observe interpolation in the autoencoder’s latent space often leads to unrealistic outputs
 631 that do not resemble valid digits. As illustrated in the cartoon of Figure 7, such interpolation may
 632 cross regions of low data density. In stark contrast, EFS effectively avoids generating samples from
 633 these low-density regions, as similarly observed for Gaussian mixtures in Figure 6.

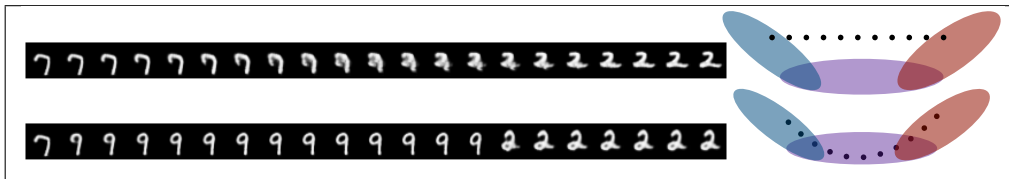


Figure 7: **Interpolation Comparison.** In the interpolation between images 7 and 2, the intermediate autoencoder outputs do not resemble valid digits (1st row), whereas our method produces a transition via $7 \rightarrow 9 \rightarrow 2$ (2nd row). On the right, we visualize the interpolation paths cross low density regions for autoencoder, while moves along the underlying digit clusters for EFS. Moreover, for EFS, it can be observed that different images are generated for each digit, illustrating that interpolation paths evolve smoothly within clusters.

634 F.1 Autoencoder details

635 As previously mentioned, the experiments in Section F employ an autoencoder. In this section, we
 636 provide further details on the architecture and training procedure of the autoencoder. The model is

	γ	k	T	β	ϵ	s	n
Gaussian mixtures	0.1	31	300	0.1	0.001	1	400
MNIST	0.05	120	300	0.01	0.001	$d - 2$	15000
Swiss roll	0.05	120	300	0.1	0.001	$d - 2$	500

Table 1: Parameters of EFS

637 a standard convolutional autoencoder consisting of two encoding layers followed by two decoding
638 layers. The structure of the encoder is summarized in Listing 1. Importantly, the encoder serves as a
639 dimensionality reduction mechanism. Without this reduction, computing $W^{(s)}$, as required in EFS,
640 becomes numerically infeasible.

641 F.2 Training Protocol

642 We trained the autoencoder on the MNIST dataset using the Adam optimizer with a learning rate of
643 10^{-3} , a batch size of 250, and for 120 epochs. The training loss, measured by mean squared error
644 (MSE), decreased to 0.008 after 120 epochs. Training was conducted on an NVIDIA RTX 6000
645 GPU with 48 GB of memory. The implementation was done in PyTorch (Paszke et al., 2019) (see the
646 notebook `neurips-figures4-5.ipynb` for details). No preprocessing was applied to the data.

```

647 Autoencoder(
648   (encoder): Sequential(
649     (0): Conv2d(1, 16, kernel_size=(3, 3), stride=(2, 2), padding=(1,
650       1))
651     (1): ReLU(inplace=True)
652     (2): Conv2d(16, 32, kernel_size=(3, 3), stride=(2, 2), padding=(1,
653       1))
654     (3): ReLU(inplace=True)
655     (4): Conv2d(32, 15, kernel_size=(7, 7), stride=(1, 1))
656   )
657   (decoder): Sequential(
658     (0): ConvTranspose2d(15, 32, kernel_size=(7, 7), stride=(1, 1))
659     (1): ReLU(inplace=True)
660     (2): ConvTranspose2d(32, 16, kernel_size=(3, 3), stride=(2, 2),
661       padding=(1, 1), output_padding=(1, 1))
662     (3): ReLU(inplace=True)
663     (4): ConvTranspose2d(16, 1, kernel_size=(3, 3), stride=(2, 2),
664       padding=(1, 1), output_padding=(1, 1))
665     (5): Sigmoid()
666   )
667 )
668 )

```

Listing 1: Autoencoder