
On Fitting Flow Models with Large Sinkhorn Couplings

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Flow models transform data gradually from one modality (e.g. noise) onto another
2 (e.g. images) following a time-dependent velocity field, trained to fit segments
3 connecting pairs of source and target points. When the pairing between source and
4 target points is given, training flow models boils down to a supervised regression
5 problem. When no such pairing exists, as is the case when generating data from
6 noise, training flows is much harder. A popular approach lies in picking source and
7 target points independently. This can, however, lead to velocity fields that are slow
8 to train but also costly to integrate at inference time. *In theory*, one would greatly
9 benefit from training flow models by sampling pairs from an optimal transport
10 (OT) coupling, since this would lead to an optimal flow in the sense of [Benamou
11 and Brenier](#). *In practice*, recent works have proposed to sample *mini-batches* of
12 n source and n target points and reorder them using an OT solver to form *better*
13 pairs. These works have advocated using batches of size $n \approx 256$, and considered
14 OT solvers that return couplings that are either sharp (using e.g. the Hungarian
15 algorithm) or blurred (using the [Sinkhorn](#) algorithm). We follow in the footsteps
16 of these works by exploring the benefits of increasing this mini-batch size n by
17 three to four orders of magnitude, and look more carefully on the effect of the
18 entropic regularization ε used in [Sinkhorn](#). Our analysis is facilitated by new
19 scale invariant quantities to report the sharpness of a coupling, while our sharded
20 computations across multiple GPUs and nodes allow scaling up n . We show that
21 in both synthetic and image generation tasks, flow models greatly benefit when
22 fitted with large [Sinkhorn](#) couplings, with a low entropic regularization ε .

23 1 Introduction

24 Finding a map that can transform a source into a target measure is a task at the core of generative
25 modeling and unpaired modality translation. Following the widespread popularity of GAN formu-
26 lations [[Goodfellow et al., 2014](#)], the field has greatly benefited from a gradual, time-dependent
27 parameterization of these transformations as normalizing flows [[Rezende and Mohamed, 2015](#)] and
28 neural ODEs [[Chen et al., 2018](#)]. Such flow models are now commonly estimated using flow match-
29 ing [[Lipman et al., 2024](#)]. While a velocity formulation substantially increases the expressivity of
30 generative models, this results on the other hand in a higher cost at inference time due to the ad-
31 ditional burden of running an ODE solver. Indeed, a common drawback of Neural-ODE solvers is
32 that they require potentially many steps, and therefore many passes through the flow network,
33 to generate data. In principle, to mitigate this problem, the gold standard for such continuous-time
34 transformations is given by the solution of the [Benamou and Brenier](#) dynamical optimal transport
35 (OT) problem, which should be equivalent, if trained perfectly, to a 1-step generation achieved by
36 the [Monge](#) map formulation [[Santambrogio, 2015](#), §1.3]. In practice, while the mathematics [[Villani,
37 2003](#)] of optimal transport have contributed to the understanding of these methods [[Liu et al., 2022](#)],
38 the jury seems to be still out on ruling whether tools from the computational OT toolbox [[Peyré and](#)

39 [Cuturi, 2019](#)], which is typically used to compute large scale couplings on data [[Klein et al., 2025](#)],
 40 can decisively help with the estimation of flows in high-dimensional / high-sample sizes regimes.

41 **Stochastic interpolants.** The flow matching (FM) framework [[Lipman et al., 2024](#)], introduced
 42 in concurrent seminal papers [[Peluchetti, 2022](#), [Lipman et al., 2023](#), [Albergo and Vanden-Eijnden,](#)
 43 [2023](#), [Neklyudov et al., 2023](#)] proposes to estimate a flow model by leveraging a pre-defined interpo-
 44 lation μ_t between source μ_0 and target μ_1 measures — the stochastic interpolant following the termi-
 45 nology of [Albergo and Vanden-Eijnden](#). That interpolation is the crucial ingredient used to fit a pa-
 46 rameterized velocity field with a regression loss. In practice, such an interpolation can be formed by
 47 sampling $X_0 \sim \mu_0$ independently of $X_1 \sim \mu_1$ and defining μ_t as the law of $X_t := (1-t)X_0 + tX_1$.
 48 One can then fit a parameterized time-dependent velocity field $\mathbf{v}_\theta(t, \mathbf{x})$ that minimizes the expecta-
 49 tion of $\|X_1 - X_0 - \mathbf{v}_\theta(X_T, T)\|^2$ w.r.t. X_0, X_1 and T a random time variable in $[0, 1]$. This
 50 procedure (hereafter abbreviated as Independent-FM, I-FM) has been immensely successful, but
 51 can suffer from high variance, and its loss can never be zero [[Liu, 2022](#)]. Furthermore, minimizing
 52 its loss cannot recover an optimal transport path: the effect of this can be measured by noticing a
 53 high curvature when integrating the ODE needed to form an output from an input sample point \mathbf{x}_0 .

54 **From I-FM to Batch-OT FM.** To fit exactly the OT framework, ideally one would choose μ_t to
 55 be the [McCann](#) interpolation between μ_0 and μ_1 , which would be $\mu_t := ((1-t)\text{Id} + tT^*)_{\#}\mu_0$,
 56 where T^* is the [Monge](#) map connecting μ_0 to μ_1 . Unfortunately, this insight is irrelevant; knowing
 57 T^* means that no flow needs to be trained. Adopting a more practical perspective, [Pooladian et al.](#)
 58 [[2023](#)] and [Tong et al. \[2023\]](#) proposed to modify I-FM and select pairs of source and target points
 59 using discrete OT solvers. Concretely, they sample mini-batches $\mathbf{x}_0^1, \dots, \mathbf{x}_0^n$ from μ_0 and $\mathbf{x}_1^1, \dots, \mathbf{x}_1^n$
 60 from μ_1 ; compute an $n \times n$ OT coupling matrix; sample pairs of indices (i_ℓ, j_ℓ) from that bistochastic
 61 matrix, and feed the flow model with pairs $\mathbf{x}_0^{i_\ell}, \mathbf{x}_1^{j_\ell}$. This approach, referred to as Batch OT-FM in
 62 the literature, was recently used and adapted in [Tian et al. \[2024\]](#), [Generale et al. \[2024\]](#), [Klein et al.](#)
 63 [[2023](#)], [Davtyan et al. \[2025\]](#), [Kim et al. \[2024\]](#). Despite their appeal, these modifications have not
 64 yet been widely adopted. The consensus stated recently by [Lipman et al. \[2024\]](#) seems to be still
 65 that *"the most popular class of affine probability paths is instantiated by the independent coupling"*.

66 **Can mini-batch OT really help?** We try to answer this question by noticing first that the evalua-
 67 tions carried out in all of the references cited above use batch sizes of $2^8 = 256$ points, more rarely
 68 $2^{10} = 1024$, upper bounded by $2^{12} = 4096$ for [Kim et al. \[2024\]](#). We believe that for many of
 69 these works this might be due to a reliance on the Hungarian algorithm [[Kuhn, 1955](#)] whose $O(n^3)$
 70 complexity is prohibitive for large n . We also notice that, while these works also consider entropic
 71 OT (EOT) [[Cuturi, 2013](#)], they stick to a single ε regularization value in their evaluations (e.g. 0.2
 72 [Kim et al. \[2024\]](#)). We go back to the drawing board in this paper, and study whether batch OT-FM
 73 can reliably work, and if so at which regimes of mini-batch size n , regularization ε , and for which
 74 data dimensions d . Our contributions are:

- 75 • We provide a general lower bound on the statistical hardness of achieving optimal transport while
 76 using couplings supported on noise and data pairs, establishing the need to move to larger n .
- 77 • Rather than drawing an artificial line between Batch-OT (in Hungarian or EOT form) and I-FM,
 78 we leverage the fact that *all* of these approaches can be interpolated using EOT: Hungarian cor-
 79 responds to $\varepsilon \rightarrow 0$ while I-FM is recovered with $\varepsilon \rightarrow \infty$. I-FM is therefore a particular case of
 80 Batch-OT with infinite regularization, which can be continuously moved towards batch-OT.
- 81 • We modify the [Sinkhorn](#) algorithm when used with the squared-Euclidean cost: we drop norms
 82 and only use negative dot-product. This improves stability and still returns the correct solution.
- 83 • We define a renormalized entropy for couplings, to pin them efficiently on a scale of 0 (bijective
 84 assignment induced by a permutation, e.g. that returned by the Hungarian algorithm) to 1 (inde-
 85 pendent coupling). This quantity is useful because, unlike transport cost or entropy regularization
 86 ε , it is bounded in $[0, 1]$ and does not depend on the data dimension d or coupling size $n \times n$.
- 87 • We explore in our experiments substantially different regimes for n and ε . We vary the mini-batch
 88 size from $n = 2^{11} = 2048$ to $n = 2^{21} = 2,097,152$ and consider an adaptive grid to set ε that
 89 results in [Sinkhorn](#) couplings whose normalized entropy is distributed within $[0, 1]$.

90 2 Background Material on Optimal Transport and Flow Matching

91 Let $\mathcal{P}_2(\mathbb{R}^d)$ denote the space of probability measures over \mathbb{R}^d with finite second moment. Let
 92 $\mu, \nu \in \mathcal{P}_2(\mathbb{R}^d)$, and let $\Gamma(\mu, \nu)$ be the set of joint probability measures in $\mathcal{P}_2(\mathbb{R}^d \times \mathbb{R}^d)$ with left-

93 marginal μ and right-marginal ν . The OT problem in its [Kantorovich](#) formulation is:

$$W_2(\mu, \nu)^2 := \inf_{\pi \in \Gamma(\mu, \nu)} \iint \frac{1}{2} \|x - y\|^2 d\pi(x, y). \quad (1)$$

94 A minimizer of (1) is called an *OT coupling measure*, denoted π^* . If μ was a noise source and ν a
 95 data target measure, π^* would be the perfect coupling to sample pairs of noise and data to learn flow
 96 models: sample $\mathbf{x}_0, \mathbf{x}_1 \sim \pi^*$ and ensure the flow models bring \mathbf{x}_0 to \mathbf{x}_1 . Such optimal couplings π^*
 97 are in fact induced by *pushforward maps*: when paired optimally, a point \mathbf{x}_0 can only be associated
 98 with a $\mathbf{x}_1 = T(\mathbf{x}_0)$, where $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is the [Monge](#) optimal transport map, defined as follows:

$$T^*(\mu, \nu) := \arg \min_{T: T_{\#}\mu = \nu} \int \frac{1}{2} \|\mathbf{x} - T(\mathbf{x})\|^2 d\mu(\mathbf{x}) \quad (2)$$

99 where the push-forward constraint $T_{\#}\mu = \nu$ means that for $X \sim \mu$ one has $T(X) \sim \nu$. [Monge](#) OT
 100 maps have been characterized by [Brenier](#) in great detail:

101 **Theorem 1** ([\[Brenier, 1991\]](#)). *If $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ has an absolutely continuous density then (2) is solved*
 102 *by a map T^* of the form $T^* = \nabla u$, where $u : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex. Moreover if u is a convex potential*
 103 *that is such that $\nabla u_{\#}\mu = \nu$ then ∇u solves (2).*

104 As a result of Theorem 1, one can choose an arbitrary convex potential u , a starting measure μ , and
 105 define a synthetic task to train flow matching models between $\mu_0 := \mu$ and $\mu_1 := \nabla u_{\#}\mu$, for which
 106 a ground truth coupling π^* is known. Inspired by [\[Korotin et al., 2021\]](#) who considered the same
 107 result to benchmark [Monge](#) [1781] map solvers, we use this setting in § 4.1 to benchmark batch-OT.

108 **Entropic OT.** Entropic regularization [\[Cuturi, 2013\]](#) has become the most popular approach to
 109 estimate a finite sample analog of π^* using samples $(\mathbf{x}_1, \dots, \mathbf{x}_n)$ and $(\mathbf{y}_1, \dots, \mathbf{y}_n)$. Using a regu-
 110 larization strength $\varepsilon > 0$, a cost matrix $\mathbf{C} := [\frac{1}{2}\|\mathbf{x}_i - \mathbf{y}_j\|^2]_{ij}$ between these samples, the entropic
 111 OT (EOT) problem can be presented in primal and dual forms as:

$$\min_{\mathbf{P} \in \mathbb{R}_+^{n \times n}, \mathbf{P}\mathbf{1}_n = \mathbf{P}^T\mathbf{1}_n = \mathbf{1}_n/n} \langle \mathbf{P}, \mathbf{C} \rangle - \varepsilon H(\mathbf{P}), \quad \max_{\mathbf{f}, \mathbf{g} \in \mathbb{R}^n} \frac{1}{n} \langle \mathbf{f} + \mathbf{g}, \mathbf{1}_n \rangle - \varepsilon \langle \exp\left(\frac{\mathbf{f} \oplus \mathbf{g} - \mathbf{C}}{\varepsilon}\right), \mathbf{1}_n \times \mathbf{1}_n \rangle, \quad (3)$$

112 where $H(\mathbf{P}) = -\langle \mathbf{P}, \log(\mathbf{P}) \rangle$ is the discrete entropy functional.

113 The optimal solutions to (3) are usually found
 114 with the [Sinkhorn](#) algorithm, as presented in
 115 Algorithm 1, where for a matrix \mathbf{S} we write
 116 $\min_{\varepsilon}(\mathbf{S}) := [-\varepsilon \log(\mathbf{1}^\top e^{-\mathbf{S}/\varepsilon})]_i$, and \oplus is
 117 the tensor sum of two vectors, i.e. $(\mathbf{f} \oplus \mathbf{g})_{ij} :=$
 118 $\mathbf{f}_i + \mathbf{g}_j$. The optimal dual variables (3) $(\mathbf{f}^\varepsilon, \mathbf{g}^\varepsilon)$
 119 can then be used to instantiate a valid coupling
 120 matrix $\mathbf{P}^\varepsilon = \exp((\mathbf{f}^\varepsilon \oplus \mathbf{g}^\varepsilon - \mathbf{C})/\varepsilon)$, which
 121 approximately solves the finite-sample counter-
 122 part of (1). An important remark is that as
 123 $\varepsilon \rightarrow 0$, the solution \mathbf{P}^ε converges to the opti-
 124 mal transport matrix solving (1), while $\mathbf{P}^\varepsilon \rightarrow$
 125 $\frac{1}{n^2} \mathbf{1}_n \times \mathbf{1}_n$ as $\varepsilon \rightarrow \infty$. These two limiting points
 126 coincide with the *optimal assignment* matrix (or
 127 optimal permutation as returned e.g. by the
 128 Hungarian algorithm [\[Kuhn, 1955\]](#)), and the
 129 uniform independent coupling used in I-FM.

130 FM methods use a stochastic interpolant μ_t
 131 with law $X_t := (1-t)X_0 + tX_1$, to minimize
 132 the expectation of a squared-norm regression
 133 loss $\min_{\theta} \mathbb{E}_{T, X_0, X_1} \|X_1 - X_0 - \mathbf{v}_\theta(X_T, T)\|^2$

134 where $X_0 \sim \mu_0, X_1 \sim \mu_1$ and T a random variable in $[0, 1]$. In I-FM, this interpolant is implemented
 135 by taking independent batches of samples $\mathbf{x}_0^1, \dots, \mathbf{x}_0^n$ from $\mu_0, \mathbf{x}_1^1, \dots, \mathbf{x}_1^n$ from μ_1 , and t_1, \dots, t_n
 136 time values sampled in $[0, 1]$, to form the loss values $\|\mathbf{x}_1^k - \mathbf{x}_0^k - \mathbf{v}_\theta((1-t_k)\mathbf{x}_0^k + t_k\mathbf{x}_1^k, t_k)\|^2$. In
 137 the formalism of [Pooladian et al. \[2023\]](#) and [Tong et al. \[2023\]](#), the same samples $\mathbf{x}_0^1, \dots, \mathbf{x}_0^n$ and
 138 $\mathbf{x}_1^1, \dots, \mathbf{x}_1^n$ are first fed into a discrete optimal matching solver. This outputs a bistochastic coupling

Algorithm 1 SINK($\mathbf{X} \in \mathbb{R}^{n \times d}, \mathbf{Y} \in \mathbb{R}^{n \times d}, \varepsilon, \tau$)

```

1:  $\mathbf{f}, \mathbf{g} \leftarrow \mathbf{0}_n, \mathbf{0}_n$ .
2:  $\mathbf{C} \leftarrow [\frac{1}{2}\|\mathbf{x}_i - \mathbf{y}_j\|^2]_{ij}, i \leq n, j \leq n$ 
3: while  $\|\exp\left(\frac{\mathbf{f} \oplus \mathbf{g} - \mathbf{C}}{\varepsilon}\right) \mathbf{1}_n - \frac{1}{n} \mathbf{1}_n\|_1 > \tau$  do
4:    $\mathbf{f} \leftarrow \varepsilon \log \frac{1}{n} \mathbf{1}_n + \min_{\varepsilon}(\mathbf{C} - \mathbf{f} \oplus \mathbf{g}) + \mathbf{f}$ 
5:    $\mathbf{g} \leftarrow \varepsilon \log \frac{1}{n} \mathbf{1}_n + \min_{\varepsilon}(\mathbf{C}^\top - \mathbf{g} \oplus \mathbf{f}) + \mathbf{g}$ 
6: end while
7: return  $\mathbf{f}, \mathbf{g}, \mathbf{P} = \exp((\mathbf{f} \oplus \mathbf{g} - \mathbf{C})/\varepsilon)$ 

```

Algorithm 2 FM 1-Step($\mu_0, \mu_1, n, \text{OT-SOLVE}$)

```

1:  $\mathbf{X}_0 = (\mathbf{x}_0^1, \dots, \mathbf{x}_0^n) \sim \mu_0$ 
2:  $\mathbf{X}_1 = (\mathbf{x}_1^1, \dots, \mathbf{x}_1^n) \sim \mu_1$ 
3:  $\mathbf{P} \leftarrow \text{OT-SOLVE}(\mathbf{X}_0, \mathbf{X}_1)$  or  $\mathbf{I}_n/n$ 
4:  $(i_1, j_1), \dots, (i_n, j_n) \sim \mathbf{P}$ 
5:  $t_1, \dots, t_n \leftarrow \text{TIMESAMPLER}$ 
6:  $\tilde{\mathbf{x}}^k \leftarrow (1-t_k)\mathbf{x}_0^{i_k} + t_k\mathbf{x}_1^{j_k}$ , for  $k \leq n$ 
7:  $\mathcal{L}(\theta) = \sum_k \|\mathbf{x}_1^{j_k} - \mathbf{x}_0^{i_k} - \mathbf{v}_\theta(\tilde{\mathbf{x}}^k, t_k)\|^2$ 
8:  $\theta \leftarrow \text{GRADIENT-UPDATE}(\nabla \mathcal{L}(\theta))$ 

```

139 matrix $\mathbf{P} \in \mathbb{R}^{n \times n}$ which is then used to *re-shuffle* the n pairs originally provided to be better cou-
 140 pled, and which should help the velocity field fit straighter trajectories, with less training steps. The
 141 procedure is summarized in Algorithm 2 and adapted to our setup and notations. The choice \mathbf{I}_n/n
 142 corresponds to I-FM, as it would return the original untouched pairs $(\mathbf{x}_0^k, \mathbf{x}_1^k)$. Equivalently, I-FM
 143 would also be recovered if the coupling was the independent coupling $\mathbf{1}_{n \times n}/n^2$, up to the difference
 144 on carrying out stratified sampling (which would result in each noise/image observed once per mini-
 145 batch) or sampling with replacement. More recently, Davtyan et al. [2025] have proposed to keep
 146 a memory of that matching effort across mini-batches, by updating a large (of the size of the entire
 147 dataset) assignment permutation between noise and full-batch data that is locally refreshed with the
 148 output of the Hungarian method run on a small batch.

149 3 Prepping Sinkhorn for Large Batch Size and Dimension.

150 **On Using Large Batch Size and Selecting $\varepsilon > 0$.** The motivation to use larger batch sizes for
 151 Batch-OT lies in the fundamental bias introduced by using small batches in light of the OT curse
 152 of dimensionality [Chewi et al., 2024, Fatras et al., 2019], which cannot be traded off with more
 153 iterations on the flow matching loss. Specifically, we provide the following lower bound that char-
 154 acterizes the statistical hardness of optimal transport, and defer its proof to the Appendix A.1.

155 **Proposition 2.** *Suppose the support of μ_1 has intrinsic dimension r , formalized in Assumption 5.*
 156 *Define the coupling $X_0, X_1 \sim \pi_n$ as follows: first draw $\mathbf{X}_0 \sim \mu_0^{\otimes n}$ and $\mathbf{X}_1 \sim \mu_1^{\otimes n}$, then sample*
 157 *$X_0, X_1 \sim \hat{\pi}_n(\mathbf{X}_0, \mathbf{X}_1)$ for any coupling rule $\hat{\pi}_n$ supported on $\mathbf{X}_0, \mathbf{X}_1$. Then, for any $\mathbf{x}_0 \in \mathbb{R}^d$,*

$$\text{Var}_{X_0, X_1 \sim \pi_n}(X_1 | X_0 = \mathbf{x}_0) \geq cn^{-2/r},$$

158 *where $c > 0$ is a constant depending only on C and r of Assumption 5.*

159 Note that the above proposition covers the case of using couplings that are supported on batches
 160 of noise and data, as in Algorithm 1. When μ_0 admits a density, the conditional variance under
 161 exact OT would be zero. Thus, Proposition 2 shows the curse of dimensionality in learning optimal
 162 transport *with any high-dimensional data distribution* μ_1 , which is in contrast to minimax lower
 163 bounds (e.g. Chewi et al. [2024, Theorem 2.15]) that only show the hardness for *some* unknown pair
 164 of distributions. This generality is at the expense of limiting the (stochastic) coupling to be supported
 165 on $(\mathbf{X}_0, \mathbf{X}_1)$, which is the relevant setting for flow matching. This curse of dimensionality becomes
 166 milder under the *manifold hypothesis* where $r \ll d$, but still advocates for the use of large n .

167 The necessity of varying ε is that this regularization can offset the bias between a regularized empir-
 168 ical OT matrix and its coupling measure counterpart, with favorable sample complexity [Genevay
 169 et al., 2018, Mena and Niles-Weed, 2019, Rigollet and Stromme, 2025].

170 **Automatic Rescaling of ε .** A practical problem arising when running the Sinkhorn algorithm lies
 171 in choosing the ε parameter. As described earlier, while \mathbf{P}^ε does follow a path from the optimal
 172 permutation (i.e., returned by the Hungarian algorithm) to the independent coupling, as ε varies
 173 from 0 to ∞ , what matters in practice is to pick relevant values in between these two extremes. To
 174 avoid using a fixed grid that risks becoming irrelevant as we vary n and d , we revisit the strategy
 175 originally used in [Cuturi, 2013] to divide the cost matrix \mathbf{C} by its mean, median or maximal value,
 176 as implemented for instance in [Flamary et al., 2021]. While needed to avoid underflow when
 177 instantiating a kernel matrix $\mathbf{K} = e^{-\mathbf{C}/\varepsilon}$, that strategy is not relevant when using the log-sum-exp
 178 operator in our implementation (as advocated in [Peyré and Cuturi, 2019, Remark 4.23]), since the
 179 \min_ε in our implementation is *invariant* to a constant shift in \mathbf{C} , whereas mean, median and max
 180 statistics are not. We propose instead to use the *standard deviation* (STD) of the cost matrix. Indeed,
 181 the dispersion of costs around their mean has more relevance as a scale than the mean of these costs
 182 itself. The STD can be computed in nd^2 time / memory, without having to instantiate the cost matrix.
 183 When this memory cost increase from d to d^2 is too high, we subsample $n = 2^{14} = 16,384$ points.
 184 In what follows, we always pass the ε value to the Sinkhorn algorithm 1 as $\tilde{\varepsilon} := \text{std}(\mathbf{C}) \times \varepsilon$, where
 185 ε is now a scale-free quantity selected in a logarithmic grid within $[0.001, 1.0]$.

Scale-Free Renormalized Coupling Entropy. While useful to keep computations stable across
 runs, the rescaling of ε still does not provide a clear idea of whether a computed coupling \mathbf{P}^ε from
 n to n points is sharp (close to an optimal permutation) or blurred (closer to what I-FM would
 use). While a distance to the independent coupling can be computed, that to the optimal Hungarian
 permutation cannot, of course, be derived without computing it beforehand which would incur a
 prohibitive cost. Instead, we resort to a fundamental information inequality used in [Cuturi, 2013]:

if \mathbf{P} is a valid coupling between two marginal probability vectors \mathbf{a}, \mathbf{b} , then one has $\frac{1}{2}(H(\mathbf{a}) + H(\mathbf{b})) \leq H(\mathbf{P}) \leq H(\mathbf{a}) + H(\mathbf{b})$. As a result, for any ε , we define the *renormalized entropy* \mathcal{E} of a coupling of \mathbf{a}, \mathbf{b} :

$$\mathcal{E}(\mathbf{P}) := \frac{2H(\mathbf{P})}{H(\mathbf{a}) + H(\mathbf{b})} - 1 \in (0, 1].$$

186 When $\mathbf{a} = \mathbf{b} = \mathbf{1}_n/n$, as considered in this work, this simplifies to $\mathcal{E}(\mathbf{P}) := H(\mathbf{P})/\log n - 1$.
 187 Independently of the size n and ε , $\mathcal{E}(\mathbf{P}^\varepsilon)$ provides a simple measure of the proximity of \mathbf{P}^ε to an
 188 optimal assignment matrix (as \mathcal{E} gets closer to 0) or to the independent coupling (as \mathcal{E} reaches 1).
 189 As a result we report $\mathcal{E}(\mathbf{P}^\varepsilon)$ rather than ε in our plots (or to be more accurate, the *average* of $\mathcal{E}(\mathbf{P}^\varepsilon)$
 190 computed over multiple mini-batches). Figures 7 and 9 in the appendix are indexed by ε instead.

191 **From Squared Euclidean Costs to Dot-products.** Using the notation $T^*(\mu, \nu)$ introduced in (2),
 192 we notice an equivariance property of **Monge** maps. For $\mathbf{s} \in \mathbb{R}^d$ and $r \in \mathbb{R}_+$ we write $L_{r,\mathbf{s}}$ for the
 193 dilation and translation map $L_{r,\mathbf{s}}(\mathbf{x}) = r\mathbf{x} + \mathbf{s}$. Naturally, $L_{r,\mathbf{s}}^{-1}(\mathbf{x}) = (\mathbf{x} - \mathbf{s})/r = L_{1/r, -\mathbf{s}/r}(\mathbf{x})$, but
 194 also $L_{r,\mathbf{s}} = \nabla w_{r,\mathbf{s}}$ where $w_{r,\mathbf{s}}(\mathbf{x}) := \frac{r}{2}\|\mathbf{x}\|^2 - \mathbf{s}^T \mathbf{x}$ is convex. This is summarized in the following
 195 result, for which we provide a short proof in the appendix.

196 **Lemma 3.** *The Monge map $T(\mu, \nu)$ is equivariant w.r.t to dilation and translation maps, as*
 197 $T^*((L_{r,\mathbf{s}})_\# \mu, (L_{r',\mathbf{s}'})_\# \nu) = L_{r',\mathbf{s}'} \circ T^*(\mu, \nu) \circ L_{r,\mathbf{s}}^{-1}$.

198 In practice this equivariance means that, when focusing on permutation matrices (which can be
 199 seen as the discrete counterparts of **Monge** maps), one is free to rescale and shift either point cloud.
 200 This remark has a practical implication when running **Sinkhorn** as well. When using the squared-
 201 Euclidean distance matrix, the cost matrix is a sum of a correlation term with two rank-1 norm terms,
 202 $\mathbf{C} = -\mathbf{X}^T \mathbf{Y} + \frac{1}{2}(\boldsymbol{\xi} \mathbf{1}_n^T + \mathbf{1}_n \boldsymbol{\gamma}^T)$ where $\boldsymbol{\xi}$ and $\boldsymbol{\gamma}$ are the vectors composed of the n squared norms of
 203 vectors in \mathbf{X} and \mathbf{Y} . Yet, due to the constraints $\mathbf{P} \mathbf{1}_n = \mathbf{a}$, $\mathbf{P}^T \mathbf{1}_n = \mathbf{b}$, any modification to the cost
 204 matrix of the form $\tilde{\mathbf{C}} = \mathbf{C} - \mathbf{c} \mathbf{1}_n^T - \mathbf{1}_n \mathbf{d}^T$, where $\mathbf{c}, \mathbf{d} \in \mathbb{R}^n$ only shifts the (3) objective by a constant:
 205 $\langle \mathbf{P}, \tilde{\mathbf{C}} \rangle = \langle \mathbf{P}, \mathbf{C} \rangle - \frac{1}{n} \mathbf{1}_n^T \mathbf{c} - \frac{1}{n} \mathbf{1}_n^T \mathbf{d}$. In practice, this means that norms only perturb **Sinkhorn** without
 206 altering the optimal coupling, and one should focus on the negative correlation matrix $\mathbf{C} := -\mathbf{X}^T \mathbf{Y}$,
 207 replacing Line 2 in Algorithm 1. We do observe substantial stability gains of these properly rescaled
 208 costs when comparing two point clouds (see Appendix A.2).

209 **Warm-starting Sinkhorn.** Solving the EOT problem (3) from scratch for each new batch of noise-
 210 data pairs $(\mathbf{X}_0, \mathbf{X}_1)$ is generally unnecessarily costly, since the solution is discarded each time a new
 211 batch is drawn. For large batch sizes, we propose to use the OT solution to i th batch $(\mathbf{X}_0^{(i)}, \mathbf{X}_1^{(i)})$ by
 212 warm-starting Sinkhorn for the $(i+1)$ th batch $(\mathbf{X}_0^{(i+1)}, \mathbf{X}_1^{(i+1)})$. Let $(\mathbf{f}^*, \mathbf{g}^*)$ be the optimal dual
 213 potentials for a given batch (\mathbf{X}, \mathbf{Y}) . These potentials can be extended to the continuous domain:

$$\mathbf{f}(\mathbf{x}) = \varepsilon \log \frac{1}{n} + \min_{\varepsilon} (\mathbf{C}(\mathbf{x}, \mathbf{y}_j) - \mathbf{g}_j), \quad \mathbf{g}(\mathbf{y}) = \varepsilon \log \frac{1}{n} + \min_{\varepsilon} (\mathbf{C}(\mathbf{x}_i, \mathbf{y}) - \mathbf{f}_i).$$

214 For a new batch $(\mathbf{X}', \mathbf{Y}')$, we use the above formula to initialize the potentials $(\mathbf{f}', \mathbf{g}')$, i.e.
 215 $(\mathbf{f}', \mathbf{g}') \leftarrow (\mathbf{f}(\mathbf{x}'_i)_i, \mathbf{g}(\mathbf{y}'_j)_j)$. Since (3) is strictly convex, the choice of initialization has no influ-
 216 ence on the solution. In practice, we find that warm-starting Sinkhorn substantially reduces the
 217 number of iterations required and the overall runtime of OTFM.

218 **Principal Component Analysis (PCA).** With the dot-product cost we can further use PCA to reduce
 219 the dimensionality and significantly speed up Sinkhorn computation. Let \mathbf{x} and \mathbf{y} represent noise
 220 and data samples respectively, and let $\mathbf{A} \in \mathbb{R}^{k \times d}$ denote the projection matrix whose rows contain
 221 top- k PCA directions. The PCA reconstruction of \mathbf{y} is $\mathbf{A}^T \mathbf{A} \mathbf{y}$, and $\mathbf{x}^T \mathbf{y} \approx \mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{y} = \bar{\mathbf{x}}^T \bar{\mathbf{y}}$,
 222 where $\bar{\mathbf{x}}$ and $\bar{\mathbf{y}}$ are the projection of \mathbf{x} and \mathbf{y} onto the PCA subspace. For large n , we compute the
 223 cost matrix on the fly per Sinkhorn iteration and avoid materializing the entire matrix at once, hence
 224 this reduction occurs per iteration. In our experiments, we can achieve $\sim 10x$ speedup in Sinkhorn
 225 computation from PCA, without sacrificing generation quality; see Appendix A.3 for details.

226 **Scaling Up Sinkhorn to Millions of High-Dimensional Points.** When guiding flow matching with
 227 batch-OT as presented in Algorithm 2, our ambition is to vary n and ε so that the coupling \mathbf{P}^ε used
 228 to sample indices can be both large ($n \approx 10^6$) and sharp if needed, i.e. with an ε that can be brought
 229 to arbitrarily low levels so that $\mathcal{E}(\mathbf{P}^\varepsilon) \approx 0$. To that end, we leverage the OTT-JAX implementation
 230 of the **Sinkhorn** algorithm [Cuturi et al., 2022], which can be natively sharded across multi-GPUs,
 231 or more generally multiple nodes of GPU machines equipped with efficient interconnect. In that
 232 approach, inspired by the earlier mono-GPU implementation of [Feydy, 2020], all n points from

ImageNet-32					ImageNet-64				
NFE \rightarrow $n \downarrow$	4	8	16	Adaptive 115 ± 1	NFE \rightarrow $n \downarrow$	4	8	16	Adaptive 269 ± 1
I-FM	66.4	24.3	12.1	5.55	I-FM	80.1	37.0	19.5	9.32
2048	38.2	16.8	10.0	5.89	4096	50.3	25.0	15.8	9.39
65536	33.1	15.1	9.28	4.88	32768	48.8	24.6	15.7	9.08
524288	31.5	14.8	9.19	4.85	131072	46.9	23.9	15.4	8.99

Table 1: FID for models trained across different OT batch sizes. We use the best checkpoint (w.r.t FID at Dopri5) for each model, restricting results to the setting where the relative epsilon value $\varepsilon = 0.1$ for ease of presentation (more detailed results can be seen in the plots of Figure 3).

233 source and target are sharded across GPUs and nodes (we have used either 1 or 2 nodes of 8 GPUs
234 each, either NVIDIA H100 or A100). A crucial point in our implementation is that the cost matrix
235 $\mathbf{C} = -\mathbf{X}\mathbf{Y}^T$ (following remark above) is never instantiated globally. Instead, it is recomputed at
236 each \min_ε operation in Lines 4 and 5 of Algorithm 1 locally, for these shards. All sharded results
237 are then gathered to recover \mathbf{f}, \mathbf{g} newly assigned after that iteration. When outputted, we use \mathbf{f}^ε
238 and \mathbf{g}^ε and, analogously, never instantiate the full \mathbf{P}^ε matrix (this would be impossible at sizes $n \approx 10^6$
239 we consider) but instead, materialize it blockwise to do stratified index sampling corresponding to
240 Line 4 in Algorithm 2. We use the Gumbel-softmax trick to vectorize the categorical sampling of
241 each of these lines to select, for each line index i , the corresponding column j_i .

242 4 Experiments

243 We revisit the application of Algorithm 2 using the modifications to the Sinkhorn algorithm outlined
244 in Section 3 to consider various benchmark tasks for which I-FM has been used. We consider syn-
245 thetic tasks in which the ground-truth Monge map is known, and benchmark unconditioned image
246 generation using CIFAR-10 [Krizhevsky et al., 2009], and the 32×32 and 64×64 downsampled
247 variants [Chrabaszcz et al., 2017] of the ImageNet dataset [Deng et al., 2009].

248 4.1 Synthetic Benchmark Tasks, $d = 32 \sim 256$

249 We consider in this section synthetic benchmarks of dimensionality ($d = 32 \sim 256$). We favor
250 this synthetic setting over other data sources with similar dimensions (e.g. single-cell data [Bunne
251 et al., 2024]) in order to have access to the ground-truth reconstruction loss, which helps elucidate
252 the impact of OT batch size n and ε . Details on the FM training setup are provided in Appendix A.4.

253 **Piecewise Affine Brenier Map.** The source is a stan-
254 dard Gaussian and the target is obtained by mapping
255 it through the gradient of a potential, itself a (con-
256 vex) piecewise quadratic function obtained using the
257 pointwise maximum of k rank-deficient parabolas:

$$u(\mathbf{x}) := \max_{i \leq k} \frac{1}{2} \|\mathbf{x}\|^2 + \frac{1}{2} \|\mathbf{A}_i(\mathbf{x} - \mathbf{m}_i)\|^2 - \|\mathbf{A}_i \mathbf{m}_i\|^2, \quad (4)$$

258 where $\mathbf{A}_i \sim \text{Wishart}(\frac{d}{2}, I_d)$, $\mathbf{m}_i \sim \mathcal{N}(0, 3I_d)$, $c_i \sim$
259 $\mathcal{N}(0, 1)$ and all means are centered around zero after
260 sampling. In practice, this yields a transport map
261 of the form $\nabla u(\mathbf{x}) = \mathbf{x} + \mathbf{A}_{i^*}(\mathbf{x} - \mathbf{m}_{i^*})$ where
262 i^* is the potential selected for that particular \mathbf{x} (i.e.
263 the argmax in (4)). The correction $-\|\mathbf{A}_i \mathbf{m}_i\|^2$ is
264 designed to ensure that these potentials are sampled
265 equally even when \mathbf{m}_i is sampled far from 0. The
266 number of potentials k is set to $d/16$. Examples of this map are shown in Appendix A.5 for dimen-
267 sion 128. We consider this setting in dimensions $d = 32, 64, 128, 256$.

268 **Korotin et al. Benchmark.** The source is a predefined Gaussian mixture and the ground-truth OT
269 map is a pre-trained ICNN. We consider this benchmark in various dimensions $d = 32 \sim 256$. Due
270 to space constraints, we refer the reader to Appendix A.6 for results in this case.

271 **Results.** From Figure 2 we find that increasing n is generally impactful and beneficial for all metrics.
272 The interest of decreasing ε , while beneficial in smaller dimensions, can be less pronounced in higher

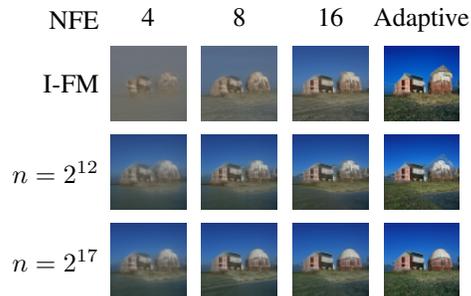


Figure 1: Samples generated from models trained on **ImageNet-64**. n denotes the total OT batch size. We use $\varepsilon = 0.1$ and the Euler solver (Dopri5 for adaptive with NFE ≈ 270). More samples in Figure 14.

273 dimensions. Indeed, we find that renormalized entropies around ≈ 0.1 should be advocated, if one
 274 considers the computational effort needed to get these samples, pictured at the bottom of each figure.

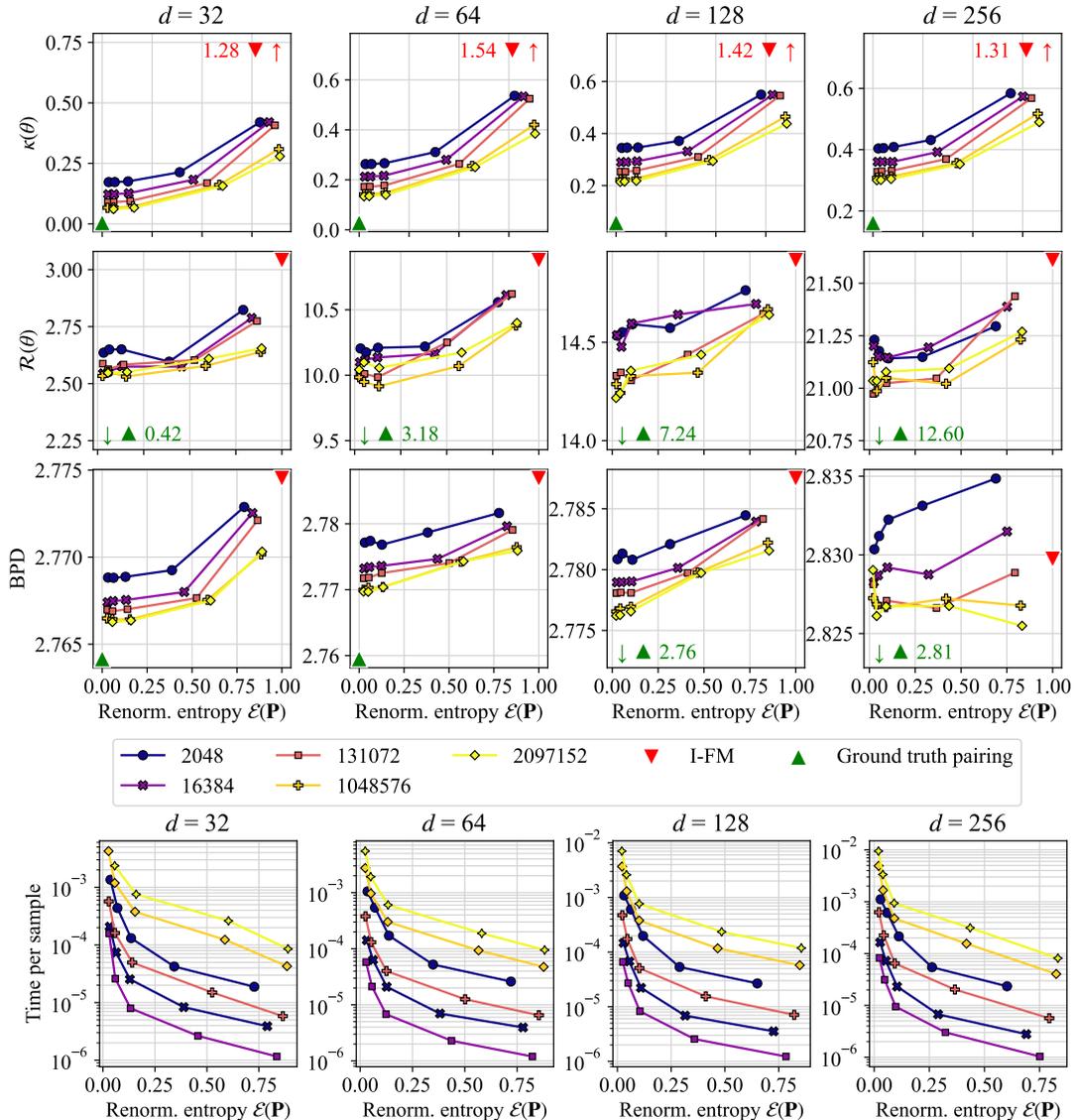


Figure 2: Results on the **piecewise affine OT Map benchmark**. The three top rows present (in that order) curvature, reconstruction and BPD metrics. Below, we provide compute times associated with running the **Sinkhorn** algorithm as a per-example cost. This per-example cost is the total time needed to run **Sinkhorn** to get $n \times n$ coupling divided by n . That cost would be 0 when using I-FM. We observe across all dimensions improvements of all metrics.

275 **4.2 Unconditioned Image Generation, $d = 3072 \sim 12288$.**

276 As done originally in [Lipman et al., 2023], we consider unconditional generation of the CIFAR-
 277 10, ImageNet-32 and ImageNet-64 datasets. Details on network parameterization and training are
 278 provided in Appendix A.7, and we defer results for CIFAR-10 generation to Appendix A.8 due to
 279 space constraints. ImageNet-32 and ImageNet-64 generation results are shown in Figures 3 and 4.
 280 Compared to results reported in [Tong et al., 2023] we observe slightly better FID scores (about
 281 0.1 when using the Dopri5 solver for instance) for I-FM. Compared to CIFAR-10, these datasets
 282 are more suitable for our large OT batch sizes as they contain significantly more samples, and we
 283 continue to observe the benefits of larger batch size and proper choice of renormalized entropy.

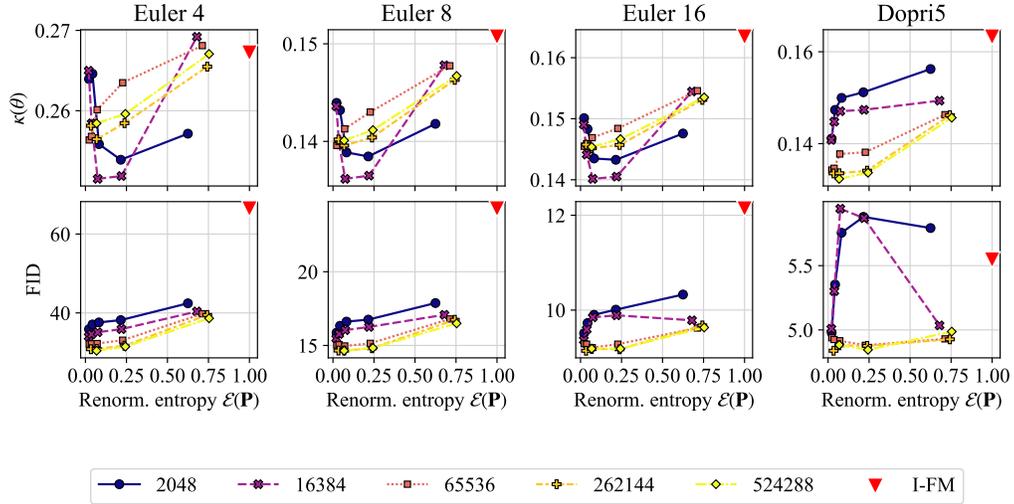


Figure 3: **ImageNet-32** experiment metrics. We observe that both FID and curvature are smaller when using larger OT batch size, and smaller renormalized entropy tends to result in better metrics.

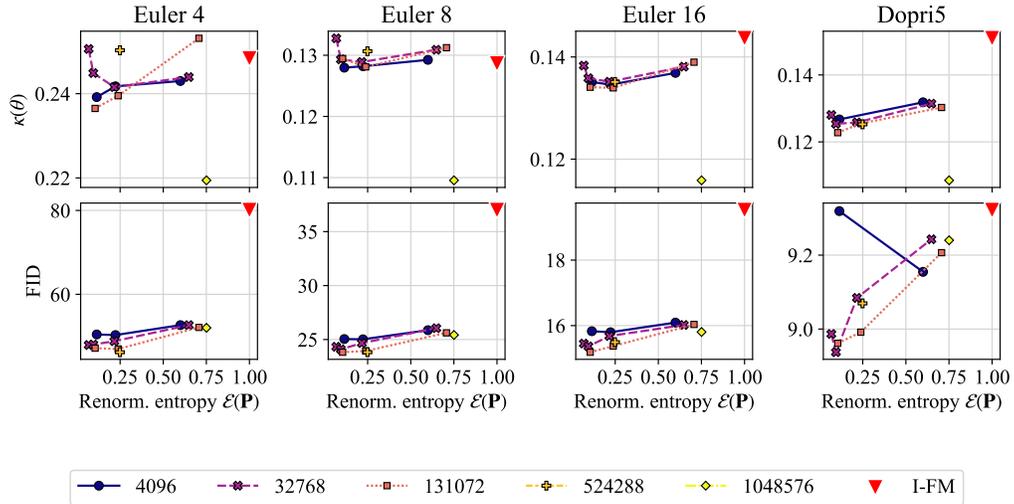


Figure 4: **ImageNet-64** results: Curvature and FID obtained with Euler integration with varying number of steps, as well as Dopri5 integration.

284 **Conclusion**

285 Our experiments suggest that guiding flow models with large scale Sinkhorn couplings can prove
 286 beneficial for downstream performance. We have tested this hypothesis by computing and sampling
 287 from both crisp and blurry $n \times n$ Sinkhorn coupling matrices for sizes n in the millions of points,
 288 placing them on an intuitive scale from 0 (close to using an optimal permutation as returned e.g.
 289 by the Hungarian algorithm) to 1 (equivalent to the independent sampling approach popularized
 290 by Lipman et al. [2023]). This involved efficient multi-GPU parallelization, realizing scales which,
 291 to our knowledge, were never achieved previously in the literature. Although the scale of these
 292 computations may seem large, they are still relatively cheap compared to the price one has to pay
 293 to optimize the FM loss, and, additionally, are completely independent from model training. As a
 294 result, they should be carried out prior to any training. While we have not explored the possibility of
 295 launching multiple jobs with them (to ablate, e.g., for other fundamental aspects of model training
 296 such as learning rates), we leave a more careful tuning of these training runs for future work. We
 297 claim that paying this relatively small price to log and sample paired indices obtained from large
 298 scale couplings results for mid-sized problems in great returns in the form of faster training and
 299 faster inference, thanks to the straightness of the flows learned with the batch-OT procedure.

300 References

- 301 Michael S Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic inter-
302 polants. In *11th International Conference on Learning Representations, ICLR 2023*, 2023.
- 303 Jean-David Benamou and Yann Brenier. A computational fluid mechanics solution to the monge-
304 kantorovich mass transfer problem. *Numerische Mathematik*, 84(3):375–393, 2000.
- 305 Yann Brenier. Polar factorization and monotone rearrangement of vector-valued functions. *Communi-
306 cations on Pure and Applied Mathematics*, 44(4), 1991. doi: 10.1002/cpa.3160440402.
- 307 Charlotte Bunne, Geoffrey Schiebinger, Andreas Krause, Aviv Regev, and Marco Cuturi. Optimal
308 transport for single-cell and spatial omics. *Nature Reviews Methods Primers*, 4(1):58, 2024.
- 309 Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary
310 differential equations. *Advances in neural information processing systems*, 31, 2018.
- 311 Sinho Chewi, Jonathan Niles-Weed, and Philippe Rigollet. Statistical optimal transport. *arXiv
312 preprint arXiv:2407.18163*, 2024.
- 313 Patryk Chrabaszcz, Ilya Loshchilov, and Frank Hutter. A downsampled variant of imagenet as an
314 alternative to the cifar datasets. *arXiv preprint arXiv:1707.08819*, 2017.
- 315 Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In *Advances in
316 neural information processing systems*, pages 2292–2300, 2013.
- 317 Marco Cuturi, Laetitia Meng-Papaxanthos, Yingtao Tian, Charlotte Bunne, Geoff Davis, and Olivier
318 Teboul. Optimal transport tools (ott): A jax toolbox for all things wasserstein. *arXiv preprint
319 arXiv:2201.12324*, 2022.
- 320 Aram Davtyan, Leello Tadesse Dadi, Volkan Cevher, and Paolo Favaro. Faster inference of flow-
321 based generative models via improved data-noise coupling. In *The Thirteenth International Con-
322 ference on Learning Representations*, 2025.
- 323 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hier-
324 archical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*,
325 pages 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.
- 326 Kilian Fatras, Younes Zine, Rémi Flamary, Rémi Gribonval, and Nicolas Courty. Learning with
327 minibatch wasserstein: asymptotic and gradient properties. *arXiv preprint arXiv:1910.04091*,
328 2019.
- 329 Jean Feydy. *Analyse de données géométriques, au delà des convolutions*. PhD thesis, Université
330 Paris-Saclay, 2020.
- 331 Rémi Flamary, Nicolas Courty, Alexandre Gramfort, Mokhtar Z. Alaya, Aurélie Boisbunon, Stanis-
332 las Chambon, Laetitia Chapel, Adrien Corenflos, Kilian Fatras, Nemo Fournier, Léo Gautheron,
333 Nathalie T.H. Gayraud, Hicham Janati, Alain Rakotomamonjy, Ievgen Redko, Antoine Rolet,
334 Antony Schutz, Vivien Seguy, Danica J. Sutherland, Romain Tavenard, Alexander Tong, and
335 Titouan Vayer. Pot: Python optimal transport. *Journal of Machine Learning Research*, 22(78):
336 1–8, 2021.
- 337 Adam P Generale, Andreas E Robertson, and Surya R Kalidindi. Conditional variable flow matching:
338 Transforming conditional densities with amortized conditional optimal transport. *arXiv preprint
339 arXiv:2411.08314*, 2024.
- 340 Aude Genevay, Lénaïc Chizat, Francis Bach, Marco Cuturi, and Gabriel Peyré. Sample complexity
341 of sinkhorn divergences. *arXiv preprint arXiv:1810.02733*, 2018.
- 342 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
343 Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural infor-
344 mation processing systems*, pages 2672–2680, 2014.
- 345 Leonid Kantorovich. On the transfer of masses (in russian). *Doklady Akademii Nauk*, 37(2), 1942.

- 346 Beomsu Kim, Yu-Guan Hsieh, Michal Klein, Marco Cuturi, Jong Chul Ye, Bahjat Kawar, and
347 James Thornton. Simple reflow: Improved techniques for fast flow models. *arXiv preprint*
348 *arXiv:2410.07815*, 2024.
- 349 Dominik Klein, Giovanni Palla, Marius Lange, Michal Klein, Zoe Piran, Manuel Gander, Laetitia
350 Meng-Papaxanthos, Michael Sterr, Lama Saber, Changying Jing, et al. Mapping cells through
351 time and space with moscot. *Nature*, pages 1–11, 2025.
- 352 Leon Klein, Andreas Krämer, and Frank Noé. Equivariant flow matching. *Advances in Neural*
353 *Information Processing Systems*, 36:59886–59910, 2023.
- 354 Alexander Korotin, Lingxiao Li, Aude Genevay, Justin Solomon, Alexander Filippov, and Evgeny
355 Burnaev. Do Neural Optimal Transport Solvers Work? A Continuous Wasserstein-2 Benchmark.
356 In *NeurIPS*, 2021.
- 357 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny im-
358 ages.(2009), 2009.
- 359 Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics*
360 *quarterly*, 2(1-2):83–97, 1955.
- 361 Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow
362 matching for generative modeling. In *The Eleventh International Conference on Learning Repre-*
363 *sentations*, 2023. URL <https://openreview.net/forum?id=PqvMRDCJT9t>.
- 364 Yaron Lipman, Marton Havasi, Peter Holderrieth, Neta Shaul, Matt Le, Brian Karrer, Ricky TQ
365 Chen, David Lopez-Paz, Heli Ben-Hamu, and Itai Gat. Flow matching guide and code. *arXiv*
366 *preprint arXiv:2412.06264*, 2024.
- 367 Qiang Liu. Rectified flow: A marginal preserving approach to optimal transport. *arXiv preprint*
368 *arXiv:2209.14577*, 2022.
- 369 Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and
370 transfer data with rectified flow. 2022.
- 371 Robert J McCann. A convexity principle for interacting gases. *Advances in mathematics*, 128(1):
372 153–179, 1997.
- 373 Gonzalo Mena and Jonathan Niles-Weed. Statistical bounds for entropic optimal transport: sample
374 complexity and the central limit theorem. *Advances in neural information processing systems*, 32,
375 2019.
- 376 Gaspard Monge. Mémoire sur la théorie des déblais et des remblais. *Histoire de l’Académie Royale*
377 *des Sciences*, 1781.
- 378 Kirill Neklyudov, Rob Brekelmans, Daniel Severo, and Alireza Makhzani. Action matching: Learn-
379 ing stochastic dynamics from samples. In *International conference on machine learning*, pages
380 25858–25889. PMLR, 2023.
- 381 Stefano Peluchetti. Non-denoising forward-time diffusions, 2022. URL <https://openreview.net/forum?id=oVfIKuhqfC>.
- 383 Gabriel Peyré and Marco Cuturi. Computational optimal transport. *Foundations and Trends in*
384 *Machine Learning*, 11(5-6), 2019. ISSN 1935-8245.
- 385 Aram-Alexandre Pooladian, Heli Ben-Hamu, Carles Domingo-Enrich, Brandon Amos, Yaron Lip-
386 man, and Ricky TQ Chen. Multisample flow matching: Straightening flows with minibatch cou-
387 plings. In *International Conference on Machine Learning*, pages 28100–28127. PMLR, 2023.
- 388 Danilo Jimenez Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *Pro-*
389 *ceedings of the 32nd International Conference on International Conference on Machine Learning-*
390 *Volume 37*, pages 1530–1538, 2015.
- 391 Philippe Rigollet and Austin J Stromme. On the sample complexity of entropic optimal transport.
392 *The Annals of Statistics*, 53(1):61–90, 2025.

- 393 Filippo Santambrogio. *Optimal transport for applied mathematicians*. Springer, 2015.
- 394 Richard Sinkhorn. A relationship between arbitrary positive matrices and doubly stochastic matrices.
395 *Ann. Math. Statist.*, 35:876–879, 1964.
- 396 Qingwen Tian, Yuxin Xu, Yixuan Yang, Zhen Wang, Ziqi Liu, Pengju Yan, and Xiaolin Li. Equiflow:
397 Equivariant conditional flow matching with optimal transport for 3d molecular conformation pre-
398 diction. *arXiv preprint arXiv:2412.11082*, 2024.
- 399 Alexander Tong, Nikolay Malkin, Kilian Fatras, Lazar Atanackovic, Yanlei Zhang, Guillaume
400 Huguet, Guy Wolf, and Yoshua Bengio. Simulation-free schrödinger bridges via score and
401 flow matching. *arXiv preprint arXiv:2307.03672*, 2023.
- 402 Alexander Tong, Kilian FATRAS, Nikolay Malkin, Guillaume Huguet, Yanlei Zhang, Jarrid Rector-
403 Brooks, Guy Wolf, and Yoshua Bengio. Improving and generalizing flow-based generative models
404 with minibatch optimal transport. *Transactions on Machine Learning Research*, 2024.
- 405 Cédric Villani. *Topics in optimal transportation*. Number 58. American Mathematical Soc., 2003.

406 **A Appendix**

407 **A.1 The Necessity of Large OT Batch Size**

408 Here, we formalize the assumptions in and provide the proof of Proposition 2.

409 Assuming that μ_0 admits a density, if we were to couple X_0 and X_1 through optimal transport,
 410 by Theorem 1 we would have $\text{Var}(X_1 | X_0) = 0$ a.s. over X_0 , where variance is the sum of co-
 411 ordinate variances. In general, any coupling that provides $\text{Var}(X_1 | X_0) = 0$ allows for one-step
 412 generation, simply by performing least-squares regression to learn $\mathbb{E}[X_1 | X_0]$. Therefore, we adopt
 413 $\text{Var}(X_1 | X_0)$ as a measure of success of a coupling.

414 Recall from 2 that to obtain a pair of samples X_0, X_1 for training, we first draw n i.i.d. samples
 415 $\mathbf{X}_0 \sim \mu_0^{\otimes n}$ and $\mathbf{X}_1 \sim \mu_1^{\otimes n}$. Then, we sample $X_0, X_1 \sim \hat{\pi}_n(\mathbf{X}_0, \mathbf{X}_1)$, where $\hat{\pi}_n$ denotes the
 416 discrete optimal (entropic) transport solution between the uniform distribution on \mathbf{X}_0 and \mathbf{X}_1 . We
 417 only require $\hat{\pi}_n(\mathbf{X}_0, \mathbf{X}_1)$ to be supported on $\mathbf{X}_0 \times \mathbf{X}_1$, as formalized by the following assumption.

418 **Assumption 4.** $\hat{\pi}_n(\mathbf{X}_0, \mathbf{X}_1)$ is supported on \mathbf{X}_0 and \mathbf{X}_1 , more precisely,

$$\hat{\pi}_n(\mathbf{X}_0, \mathbf{X}_1) = \sum_{i,j} P_{ij}(\mathbf{X}_0, \mathbf{X}_1) \delta_{(\mathbf{x}_0^{(i)}, \mathbf{x}_1^{(j)})},$$

419 where $(P_{ij}(\mathbf{X}_0, \mathbf{X}_1))_{ij}$ is some bistochastic matrix, equivariant under permutations of \mathbf{X}_0 and \mathbf{X}_1 ,
 420 and δ denotes the Dirac measure.

421 To capture the intrinsic dimension of data, we can impose the following assumption on μ_1 .

Assumption 5. For X and X' drawn independently from the data distribution μ_1 , we have

$$\mathbb{P}[\|X - X'\| \leq t] \leq Ct^r,$$

422 for all $t > 0$ and some $C, r > 0$.

423 Note that the volume of an r -dimensional ball of radius t is proportional to t^r . Therefore, r in the
 424 above assumption roughly captures the intrinsic dimension of data, typically assumed to be much
 425 less than the ambient dimension, i.e. $r \ll d$.

426 We are now ready to present the proof of Proposition 2, which we repeat here for ease of reference.

427 **Proposition 6.** *Suppose $\hat{\pi}_n$ is any coupling rule that satisfies Assumption 4, and that μ_1 satisfies*
 428 *Assumption 5. Define the coupling $X_0, X_1 \sim \pi_n$ as follows: first draw $\mathbf{X}_0 \sim \mu_0^{\otimes n}$ and $\mathbf{X}_1 \sim \mu_1^{\otimes n}$,*
 429 *then sample $X_0, X_1 \sim \hat{\pi}_n(\mathbf{X}_0, \mathbf{X}_1)$. Then, for any $\mathbf{x}_0 \in \mathbb{R}^d$, we have*

$$\text{Var}_{X_0, X_1 \sim \pi_n}(X_1 | X_0 = \mathbf{x}_0) \geq cn^{-2/r},$$

430 where $c > 0$ is a constant depending only on C and r .

431 To prove Proposition 6, we use the fact that

$$\text{Var}(X_1 | X_0 = \mathbf{x}_0) = \frac{1}{2} \mathbb{E}[\|X_1 - X'_1\|^2 | X_0 = \mathbf{x}_0],$$

432 where X_1 and X'_1 are drawn independently from $\pi_n(\cdot | X_0 = \mathbf{x}_0)$. However, X_1 and X'_1 essentially
 433 come from different batches \mathbf{X}_1 and \mathbf{X}'_1 , and their only dependence is through being coupled with
 434 $X_0 = \mathbf{x}_0$. We can remove this dependence by lower bounding the variance by the minimum distance
 435 between two batches of i.i.d. samples \mathbf{X}_1 and \mathbf{X}'_1 . This is performed by the following lemma.

436 **Lemma 7.** *Let π_n be as defined in Proposition 6. Then, for any $\mathbf{x}_0 \in \mathbb{R}^d$, we have*

$$\text{Var}_{X_0, X_1 \sim \pi_n}(X_1 | X_0 = \mathbf{x}_0) \geq \frac{1}{2} \mathbb{E}_{\mathbf{X}_1, \mathbf{X}'_1 \sim \mu_1^{\otimes n} \otimes \mu_1^{\otimes n}}[D(\mathbf{X}_1, \mathbf{X}'_1)],$$

437 where $D(\mathbf{X}_1, \mathbf{X}'_1) := \min_{\mathbf{x}_1 \in \mathbf{X}_1, \mathbf{x}'_1 \in \mathbf{X}'_1} \|\mathbf{x}_1 - \mathbf{x}'_1\|^2$.

438 *Proof.* To draw $X_1, X'_1 \sim \pi_n(\cdot | X_0 = \mathbf{x}_0) \otimes \pi_n(\cdot | X_0 = \mathbf{x}_0)$ we can draw $\mathbf{X}_0, \mathbf{X}'_0 \sim \mu_0^{\otimes n} \otimes \mu_0^{\otimes n}$
 439 and $\mathbf{X}_1, \mathbf{X}'_1 \sim \mu_1^{\otimes n} \otimes \mu_1^{\otimes n}$. We then replace the first sample in \mathbf{X}_0 and \mathbf{X}'_0 with \mathbf{x}_0 to condition on
 440 $X_0 = \mathbf{x}_0$ (in particular, we rely on the equivariance of $\hat{\pi}_n$), and denote them by $\mathbf{X}_0(\mathbf{x}_0)$ and $\mathbf{X}'_0(\mathbf{x}_0)$.

441 We can then write

$$\begin{aligned}
\text{Var}(X_1 | X_0 = \mathbf{x}_0) &= \frac{1}{2} \mathbb{E}[\|X_1 - X'_1\|^2 | X_0 = \mathbf{x}_0] \\
&= \frac{1}{2} \mathbb{E}[\mathbb{E}[\|X_1 - X'_1\|^2 | \mathbf{X}_0(X_0), \mathbf{X}'_0(X_0), \mathbf{X}_1, \mathbf{X}'_1] | X_0 = \mathbf{x}_0] \\
&= \frac{1}{2} \mathbb{E} \left[\sum_{\mathbf{x}_1 \in \mathbf{X}_1, \mathbf{x}'_1 \in \mathbf{X}'_1} \hat{\pi}_n(\mathbf{X}_0(\mathbf{x}_0), \mathbf{X}_1)(\mathbf{x}_1 | \mathbf{x}_0) \hat{\pi}_n(\mathbf{X}'_0(\mathbf{x}_0), \mathbf{X}'_1)(\mathbf{x}'_1 | \mathbf{x}_0) \|\mathbf{x}_1 - \mathbf{x}'_1\|^2 \right] \\
&\geq \frac{1}{2} \mathbb{E} \left[D(\mathbf{X}_1, \mathbf{X}'_1) \sum_{\mathbf{x}_1 \in \mathbf{X}_1, \mathbf{x}'_1 \in \mathbf{X}'_1} \hat{\pi}_n(\mathbf{X}_0(\mathbf{x}_0), \mathbf{X}_1)(\mathbf{x}_1 | \mathbf{x}_0) \hat{\pi}_n(\mathbf{X}'_0(\mathbf{x}_0), \mathbf{X}'_1)(\mathbf{x}'_1 | \mathbf{x}_0) \right] \\
&\geq \frac{1}{2} \mathbb{E}[D(\mathbf{X}_1, \mathbf{X}'_1)],
\end{aligned}$$

442 which finishes the proof. \square

443 Using the above lemma, to prove Proposition 6, we only need to estimate the expected distance
444 between two batches of samples from μ_1 .

445 *Proof of Proposition 6.* Let $\mathbf{X}_1, \mathbf{X}'_1$ be independent batches of n i.i.d. samples from μ_1 . We use
446 expand our notation by letting $D(\mathbf{X}_1, \mathbf{X}'_1) := \min_{\mathbf{x}_1 \in \mathbf{X}_1} \|\mathbf{x}_1 - \mathbf{x}'_1\|$ be the distance between a single
447 sample and a batch. By the Markov inequality, for any $t > 0$ we have

$$\begin{aligned}
\mathbb{E}[D(\mathbf{X}_1, \mathbf{X}'_1)] &\geq t \mathbb{P}[D(\mathbf{X}_1, \mathbf{X}'_1) \geq t] \\
&= t \mathbb{E} \left[\mathbb{P} \left[\bigcap_{\mathbf{x}'_1 \in \mathbf{X}'_1} \{D(\mathbf{X}_1, \mathbf{x}'_1) \geq t\} \mid \mathbf{X}_1 \right] \right] \\
&= t \mathbb{E} \left[\mathbb{P}[D(\mathbf{X}_1, X'_1) \geq t \mid \mathbf{X}_1]^n \right] && \text{(Independence)} \\
&\geq t \mathbb{P}[D(\mathbf{X}_1, X'_1) \geq t]^n && \text{(Jensen's Inequality)} \\
&= t \mathbb{E} \left[\mathbb{P} \left[D(\mathbf{X}_1, X'_1) \geq t \mid X'_1 \right]^n \right] \\
&= t \mathbb{E} \left[\mathbb{P} \left[\bigcap_{\mathbf{x}_1 \in \mathbf{X}_1} \{\|\mathbf{x}_1 - X'_1\| \geq t\} \mid X'_1 \right]^n \right] \\
&= t \mathbb{E} \left[\mathbb{P}[\|X_1 - X'_1\| \geq t \mid X'_1]^n \right] \\
&\geq t \mathbb{P}[\|X_1 - X'_1\| \geq t]^n && \text{(Jensen's Inequality)} \\
&\geq t(1 - Ct^r)^{n^2} && \text{(Assumption 5)}
\end{aligned}$$

448 Choosing $t = (2Cn^2)^{-1/r}$ and using the inequality $(1 - 1/(2x))^x \geq 1/2$ for all $x \geq 1$ yields
449 $\mathbb{E}[D(\mathbf{X}_1, \mathbf{X}'_1)] \geq (2Cn^2)^{-1/r}/2$, which completes the proof. \square

450 *Proof of Lemma 3.* Following Brenier's theorem, let u be a convex potential such that $T^*(\mu, \nu) =$
451 ∇u . Set $F := L_{r',s'} \circ \nabla u \circ L_{r,s}^{-1}$. Then F is the composition of the gradients of 3 convex func-
452 tions. Because the Jacobians of $L_{r,s}$ and $L_{r',s'}^{-1}$ are respectively $r\mathbf{I}_d$ and \mathbf{I}_d/r , they commute with the
453 Hessian of u . Therefore the Jacobian of F is symmetric, positive definite, and F is the gradient of a
454 convex potential that pushes $(L_{r,s})\#\mu$ to $(L_{r',s'})\#\nu$, and is therefore their Monge map by Brenier's
455 theorem. \square

456 A.2 Using the negative dot-product cost rather than squared-Euclidean in Sinkhorn

457 As we mention in the main text, entropically regularized optimal transport plan for the squared
458 Euclidean cost can be equivalently recast using exclusively the negative scalar product $\langle \mathbf{x}, \mathbf{y} \rangle \mapsto$
459 $-\langle \mathbf{x}, \mathbf{y} \rangle$ between source and target, and not on any absolute measure of scale. To see this, consider
460 an affine map $\bar{\mathbf{x}} = \alpha \mathbf{x} + \beta$ with $\alpha > 0$. Then:

$$\langle \bar{\mathbf{x}}, \mathbf{y} \rangle = \alpha \langle \mathbf{x}, \mathbf{y} \rangle + \langle \beta, \mathbf{y} \rangle.$$

461 The second term is a rank-1 term that will be absorbed by the optimal dual potentials (see (3)) and
 462 the factor α amounts to a rescaling of the entropic regularization level ε . In particular, when there is
 463 only translation, then $\alpha = 1$ and the transport plans are identical for the same ε .

464 Therefore for input data $\mathbf{X} \in \mathbb{R}^{n \times d}$, $\mathbf{Y} \in \mathbb{R}^{m \times d}$ the Sinkhorn transport plan depends only on the
 465 dot-product cost $-\mathbf{X}\mathbf{Y}^T$. We argue that it is always more natural to use the dot-product cost than the
 466 full squared-Euclidean cost, and we find that in practice using directly the dot-product can improve
 467 the numerical conditioning of the Sinkhorn algorithm. This is because we drop terms arising from
 468 the squared norm, which can be very large. This becomes especially important for single-precision
 469 floating point computations, as is the case for the large scale GPU applications we consider.

470 We illustrate this in Figure 5: we sample $N = 8192$ points $\{\mathbf{x}_i\}_{i=1}^N$ in dimension $d = 128$ from
 471 the Gaussian example described in Section 4.1 and map them through the piecewise affine Brenier
 472 map, i.e. $\mathbf{y}_i = T(\mathbf{x}_i)$. We then introduce a translation, $\bar{\mathbf{y}}_i = \mathbf{y}_i + 5$. We use the Sinkhorn algo-
 473 rithm (Algorithm 1) with either the dot-product or squared Euclidean cost to compute the transport
 474 plan and record the number of iterations taken, and our computations are carried out on GPU with
 475 single-precision arithmetic. Even though we already use log-domain computation tricks to prevent
 476 under/overflow, we find that for small ε , Sinkhorn with squared Euclidean cost begins to suffer from
 477 numerical issues and fails to converge within the iteration limit of 50,000.

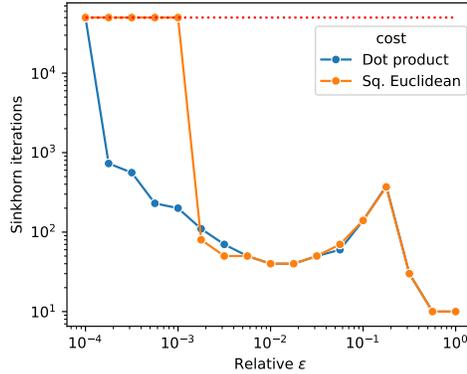


Figure 5: Number of Sinkhorn iterations against relative ε for Gaussian piecewise affine OT exam-
 ple.

478 A.3 Sinkhorn Speedup from PCA

479 Table 2 presents the average wall-clock time of running Sinkhorn on ImageNet-64 with a batch
 480 size of 131072 and $\varepsilon = 0.1$. As can be seen, we can reduce dimension by a factor of almost 25,
 481 which reduces time by a factor of 10, while having no significant impact on the quality of generated
 482 images measured by FID. Moreover, the normalized entropy demonstrates that the coupling obtained
 483 from the reduced-dimensional cost matrix has the same sharpness as the original coupling, which is
 484 expected since PCA will mostly preserve the dot product cost, resulting in similar couplings.

485 A.4 Velocity Field Parameterization and Training: Synthetic Benchmarks

486 The velocity fields are parameterized as MLPs with 5 hidden layers, each of size 512 when $d =$
 487 32, 64 and 1024 when $d = 128, 256$. Time in $[0, 1]$ is encoded using $d/8$ Fourier encodings. All
 488 models are trained with unpaired batches: the sampling in Line 1 of Algorithm 2 is done as $\bar{\mathbf{X}}_0 \sim \mu$
 489 while for Line 2, $\mathbf{X}_1 := T_0(\mathbf{X}'_0)$ where \mathbf{X}'_0 is a new sample from μ and T is applied to each of
 490 the n points described in \mathbf{X}'_0 . All models are trained for 8192 steps, with effective batch sizes of
 491 2048 samples (256 per GPU) to average a gradient, a learning rate of 10^{-3} (we tested with 10^{-2}
 492 or 10^{-4} , the former was unstable while the latter was less efficient on a subset of runs). The model
 493 marked as \blacktriangle in the plots is a flow model trained with *perfect* supervision, i.e. given *ground-truth*
 494 *paired samples* $\mathbf{X}_0 \sim \mu$ and $\mathbf{X}_1 := T_0(\mathbf{X}_0)$, provided in the correct order. I-FM is marked as \blacktriangledown .
 495 For all other runs, we vary ε (reporting renormalized entropy $\mathcal{E}(\mathbf{P}^\varepsilon)$) and the total batch size n used
 496 to compute couplings, somewhere between 2048 and 2,097,152. These runs are carried out on a
 497 single node with 8 GPUs, and therefore the data is sharded in blocks of size $n/8$ when running the
 498 Sinkhorn algorithm.

	$k = 500$	$k = 1000$	$k = 3000$	$k = 12288$ (full dimension)
<i>Sinkhorn time</i>	<i>1.45s</i>	<i>1.82s</i>	<i>4.05s</i>	<i>14.1s</i>
FID@NFE=4	48.4	48.1	47.0	47.3
FID@NFE=8	24.7	24.4	24.0	24.2
FID@NFE=16	16.0	15.8	15.8	15.8
FID@Dopri5 (Adaptive)	9.17	9.33	9.46	9.51
Renormalized Entropy	0.247	0.239	0.232	0.236

Table 2: Sinkhorn runtime per batch and FID for different solvers and different PCA dimension k . The model is trained on ImageNet-64 with OT batch size = 131072 and $\varepsilon = 0.1$. The difference with the result of Table 1 is due to a different attention implementation in the architecture.

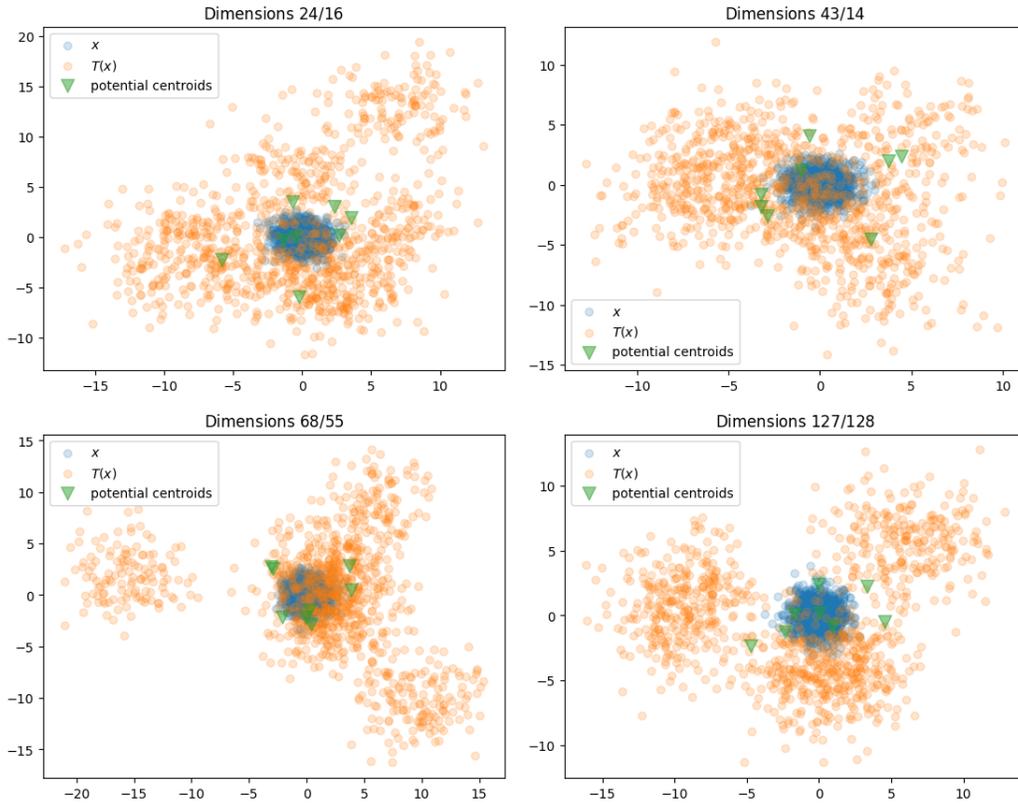


Figure 6: Example of the maps generated in our piecewise affine benchmark task. In these plots $d = 128$ and there are therefore $128/16 = 8$ quadratic potentials sampled around 0. These 2D plots illustrate the action of the same 128 dimensional map, pictured using 2D projections over pairs chosen in $[1, \dots, 128]$.

499 A.5 Gaussian Transported with a Piecewise Affine Ground-Truth OT Map

500 We present in Figure 6 examples of our piecewise affine OT map generation, corresponding to results
501 presented more widely in Figures 2 and 7.

502 A.6 Korotin et al. Benchmark Examples

503 The reader may find examples of the Korotin et al. benchmark in their paper, App. A.1, Figure 6.

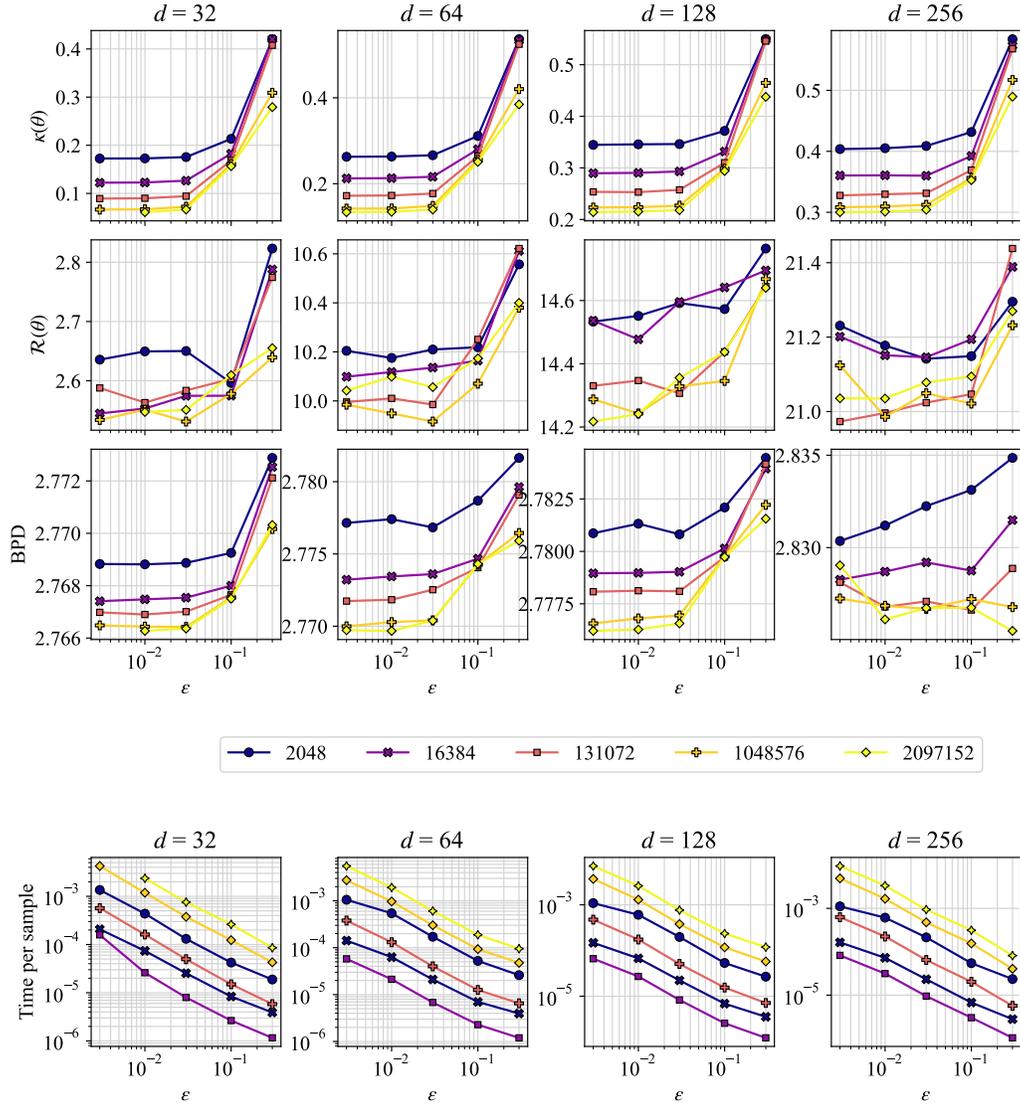


Figure 7: Plots corresponding to Figure 2 in main paper, on piecewise affine synthetic benchmark, using directly the relative epsilon parameter as the x-axis (log-scale), instead of re-normalized entropy.

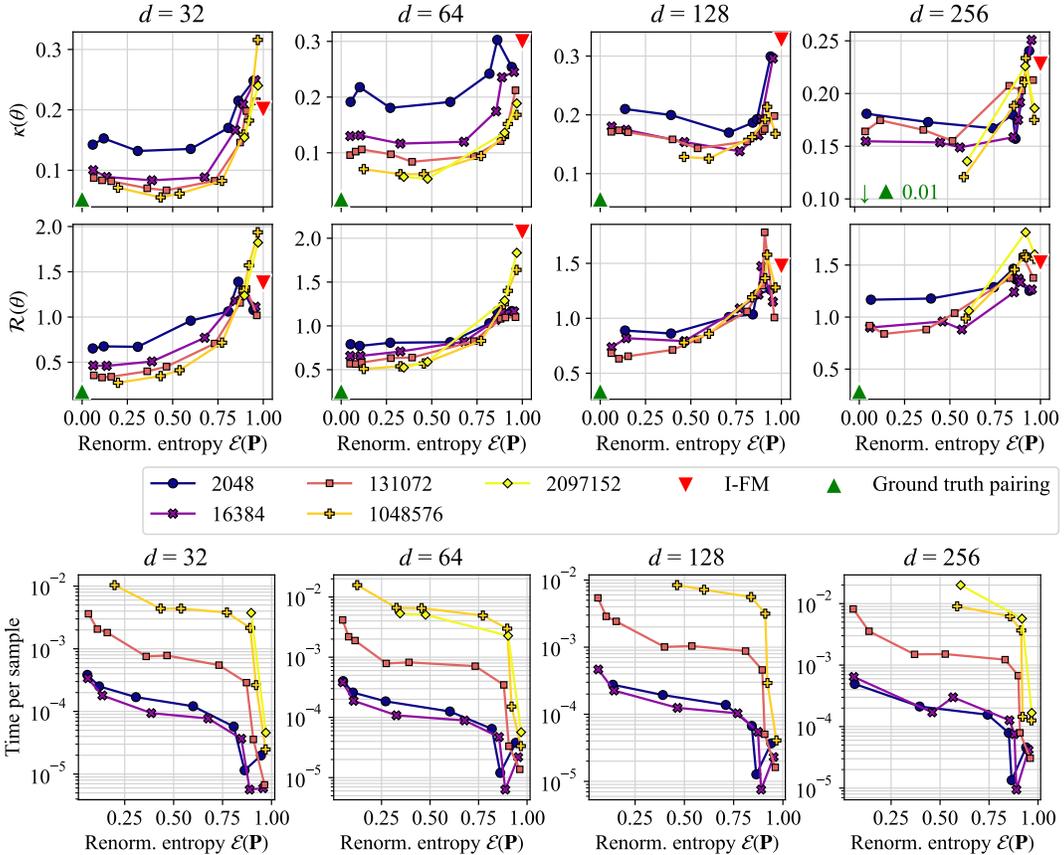


Figure 8: Results on the **Korotin benchmark**. As with Figure 2, we compute curvature and reconstruction metrics, and compute times below. Some of the runs for largest OT batch sizes n are provided in the supplementary. These runs suggest that to train OT models in these dimensions increasing n is overall beneficial across the board.

504 A.7 Velocity Field Parameterization and Training: Image Generation

505 We use the network parameterization given in [Tong et al., 2024] for CIFAR-10 and those given
 506 in [Pooladian et al., 2023] for ImageNet 32 and ImageNet 64. We follow their recommendations
 507 on setting learning rates, batch sizes (to average gradients) as well as total number of iterations:
 508 we train respectively for 400k, 438k and 957k using effective batch sizes advocated in their paper,
 509 respectively 16×8 , 128×8 and 50×16 .

510 A.8 CIFAR-10 Results

511 Results are presented in Figure 10. Compared to results reported in [Tong et al., 2023] we observe
 512 slightly better FID scores (about 0.1) for both I-FM and Batch OT-FM. Note that the size of the
 513 dataset itself (50k, 100k when including random flipping as we do) is comparable (if not slightly
 514 lower) to our largest batch size $n = 131,072$, meaning some images are duplicated. Overall, the
 515 results show the benefit of relatively larger batch sizes and suitably small ε , that is more pronounced
 516 at lower NFE.

517 We show generated images in Figure 12. We see general quantitative and qualitative improvements
 518 for larger OT batch size and smaller renormalized entropy. However, these improvements are not
 519 as significant as our observation for the more complex down-sampled ImageNet datasets in Appen-
 520 dices A.9 and A.10, likely due to the fact that the dataset size is much smaller. We also plot BPD
 521 as a function of renormalized entropy for CIFAR-10, ImageNet-32, and ImageNet-64, in Figure 11.

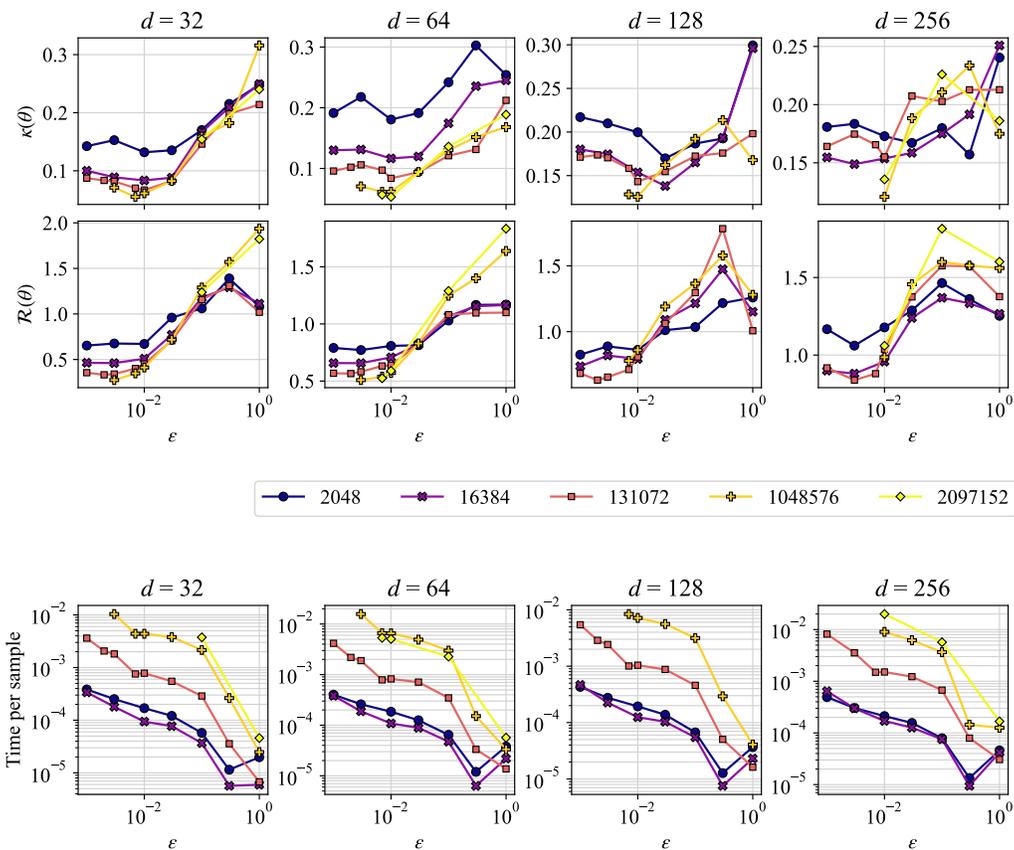


Figure 9: Plots for the [Korotin et al.](#) benchmark, shown initially in Figure 8, using the relative epsilon ϵ parameter directly in the x-axis, in logarithmic scale.

522 **A.9 ImageNet-32 Detailed Results**

523 Figure 13 shows generated images using I-FM and OT-FM with different batch sizes and different
 524 ODE solvers. As expected, the greatest improvements in the quality of images occur with smaller
 525 number of integration steps, which demonstrates the benefit of OT-FM for reducing inference cost.

526 **A.10 ImageNet-64 Detailed Results**

527 We also perform experiments on the 64×64 downsampled ImageNet dataset, where we observe an
 528 even bigger gap between I-FM and OT-FM with large batch size both in terms of metrics (Figure 4)
 529 and in terms of qualitative results (Figure 14). This observation implies that with a proper choice of
 530 entropy and batch size, OT-FM is a promising approach to reduce inference cost and generate higher
 531 quality high-resolution images.

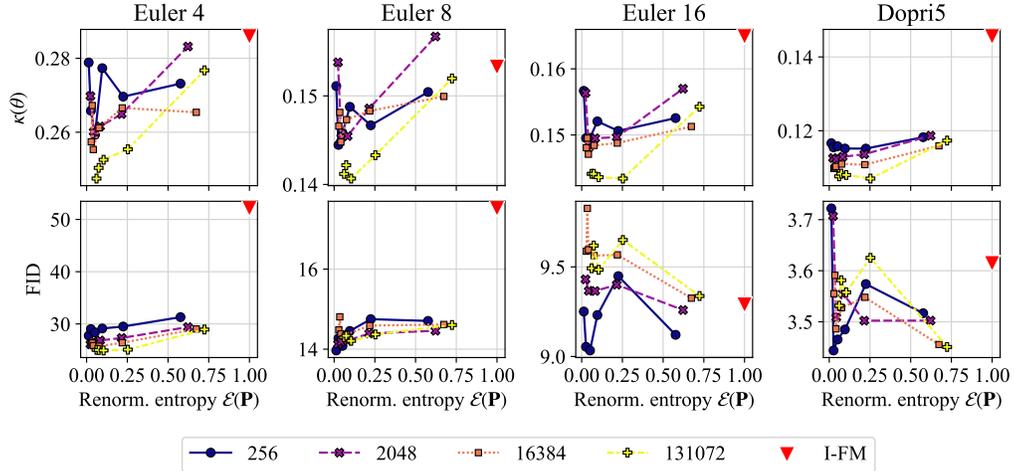


Figure 10: Experiment metrics for **CIFAR-10** image generation. We evaluate the trained models using the Euler solver with three different number of steps, and with the Dopri5 solver and adaptive steps. The plots demonstrate the benefits of a larger OT batch size to achieve significantly smaller curvature, and moderately smaller FID at low number of integration steps. CIFAR-10 is not necessarily the best setup to evaluate the performance of OT based FM, since the number of points is relatively low (the batch sizes we consider involve in fact resampling *data*). Our experiments also suggest that in this setting, lower renormalized entropy generally benefits the performance.

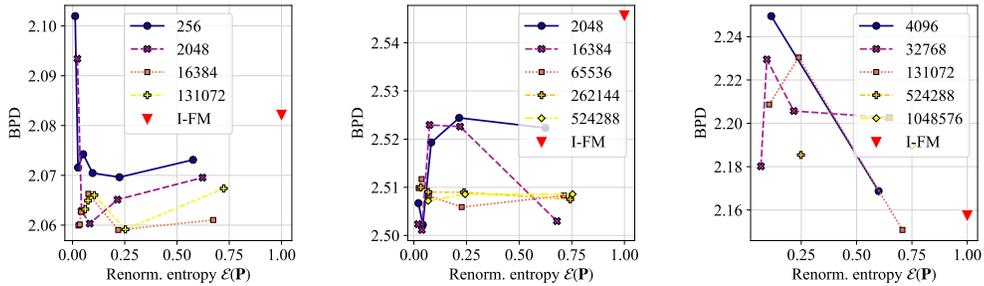


Figure 11: BPD for **CIFAR-10** (Left), **ImageNet-32** (Middle) and **ImageNet-64** (Right). The BPDs are computed using Dopri5 integration, evaluated on 50 times steps, and computed using 8 vectors for the Hutchinson trace estimator. As a consequence of its high number of function evaluations, the Dopri5 solver relies less on straightness of the flows. Therefore, we do not observe a significant difference across batch sizes.

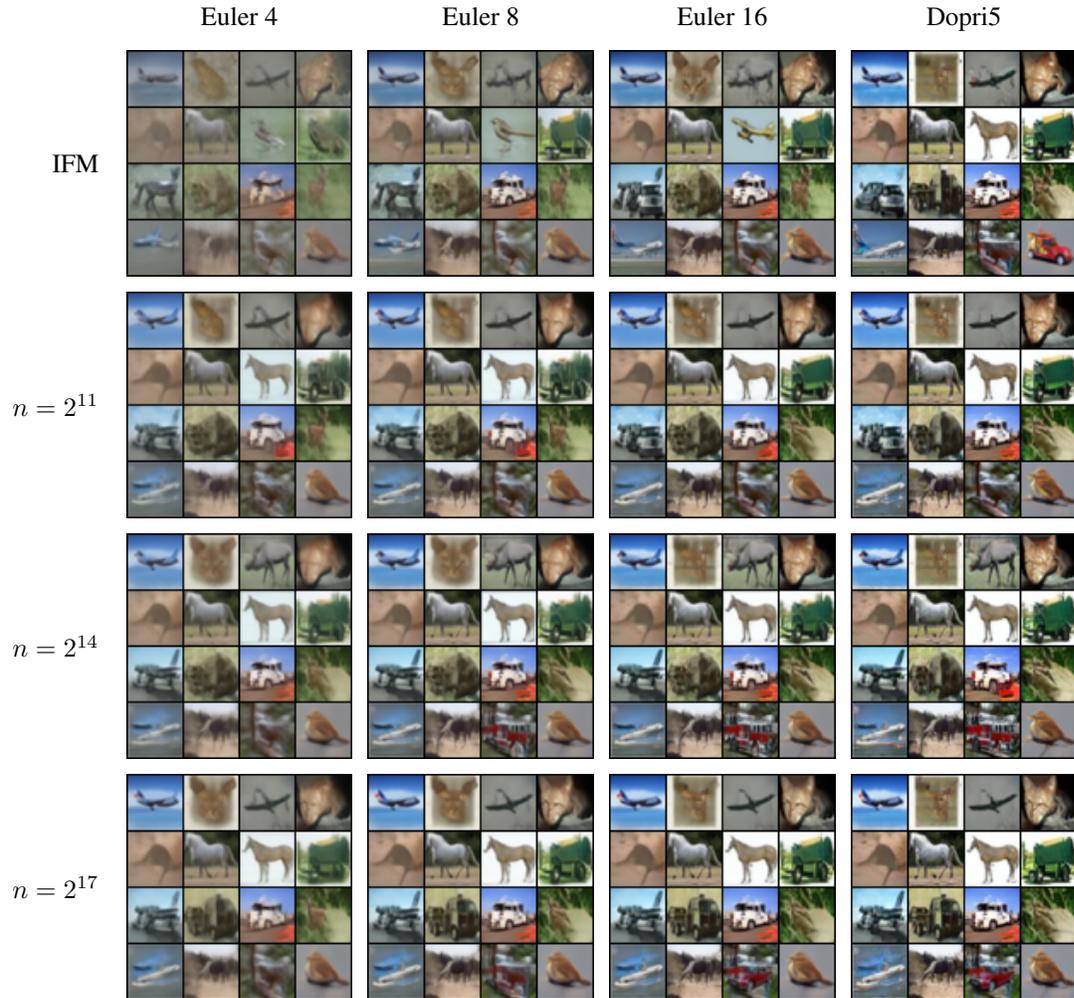


Figure 12: Non-curated images generated from models trained on **CIFAR-10**. The number following Euler denotes NFE, while Dopri5 uses an adaptive number of evaluations. n denotes the total batch size for the Sinkhorn algorithm. We use OT-FM models trained with $\varepsilon = 0.01$.

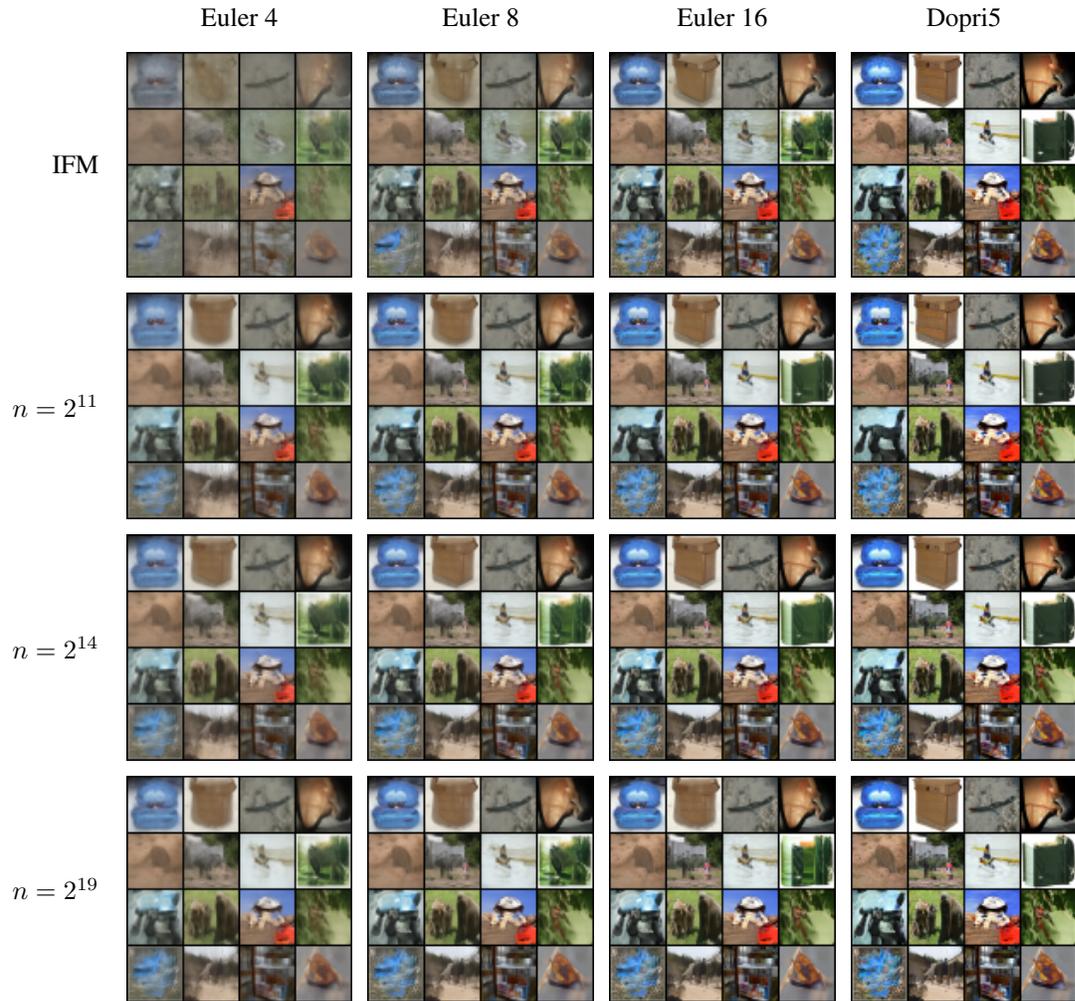


Figure 13: Non-curated images generated from models trained on **ImageNet-32**. n denotes the total batch size for the Sinkhorn algorithm. We use OT-FM models trained with $\varepsilon = 0.1$.

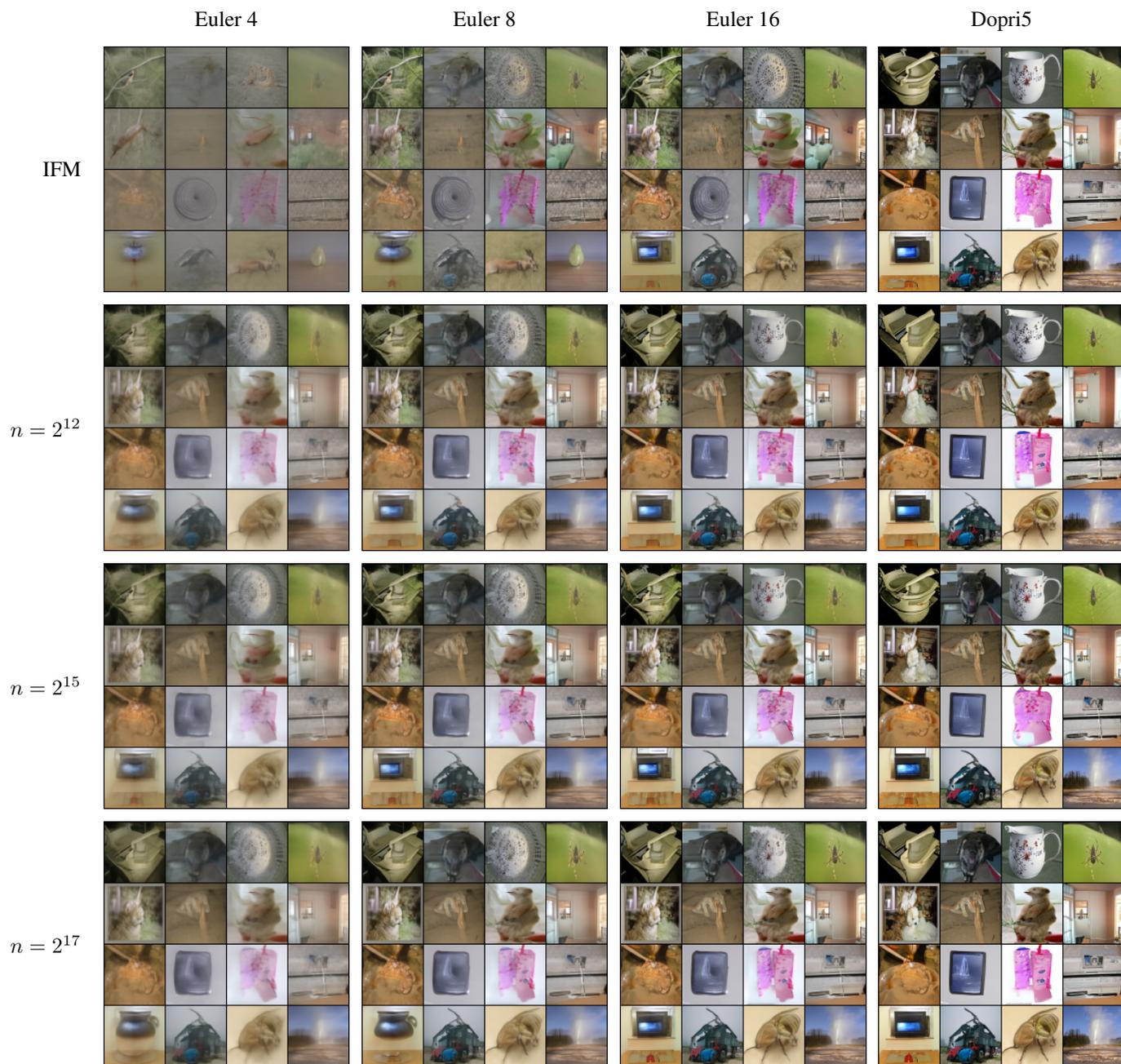


Figure 14: Non-curated images generated from models trained on **ImageNet-64**. n denotes the total batch size for the Sinkhorn algorithm. We use models trained with a varying trained with $\varepsilon = 0.1$.