
Meta Optimal Transport

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We study the use of amortized optimization to predict optimal transport (OT) maps
2 from the input measures, which we call *Meta OT*. This helps repeatedly solve sim-
3 ilar OT problems between different measures by leveraging the knowledge and in-
4 formation present from past problems to rapidly predict and solve new problems.
5 Otherwise, standard methods ignore the knowledge of the past solutions and sub-
6 optimally re-solve each problem from scratch. We instantiate Meta OT models in
7 discrete and continuous (Wasserstein-2) settings between images, spherical data,
8 and color palettes and use them to improve the computational time of standard OT
9 solvers by multiple orders of magnitude.

10 1 Introduction

11 Optimal transportation [Villani, 2009, Ambrosio, 2003, Santambrogio, 2015, Peyré et al., 2019,
12 Merigot and Thibert, 2021] is thriving in domains including economics [Galichon, 2016], rein-
13 forcement learning [Dadashi et al., 2021, Fickinger et al., 2021], style transfer [Kolkin et al., 2019],
14 generative modeling [Arjovsky et al., 2017, Seguy et al., 2018, Huang et al., 2020, Rout et al., 2021],
15 geometry [Solomon et al., 2015, Cohen et al., 2021], domain adaptation [Courty et al., 2017, Redko
16 et al., 2019], signal processing [Kolouri et al., 2017], fairness [Jiang et al., 2020], and cell repro-
17 gramming [Schiebinger et al., 2019]. A core component in these settings is to couple two measures
18 (α, β) supported on domains $(\mathcal{X}, \mathcal{Y})$ by solving a transport optimization problem such as the *primal*
19 *Kantorovich problem*, which is defined by:

$$\pi^*(\alpha, \beta, c) \in \arg \min_{\pi \in \mathcal{U}(\alpha, \beta)} \int_{\mathcal{X} \times \mathcal{Y}} c(x, y) d\pi(x, y), \quad (1)$$

20 where the *optimal coupling* π^* is a joint distribution over the product space, $\mathcal{U}(\alpha, \beta)$ is the set of
21 admissible couplings between α and β , and $c : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ is the *ground cost*, that represents a
22 notion of distance between elements in \mathcal{X} and elements in \mathcal{Y} .

23 **Challenges.** Unfortunately, solving eq. (1) *once* is computationally expensive between general mea-
24 sures and computationally cheaper alternatives are an active research topic: *Entropic optimal trans-*
25 *port* [Cuturi, 2013] smooths the transport problem with an entropy penalty, and *sliced distances*
26 [Kolouri et al., 2016, 2018, 2019, Deshpande et al., 2019] solve OT between 1-dimensional projec-
27 tions of the measures, where eq. (1) can be solved easily.

28 Furthermore, when an optimal transport method is deployed in practice, eq. (1) is not just solved
29 a single time, but is *repeatedly* solved for new scenarios between different input measures (α, β) .
30 For example, the measures could be representations of images we care about optimally transporting
31 between and in deployment we would receive a stream of new images to couple. Repeatedly solving
32 optimal transport problems also comes up in the context of comparing seismic signals [Engquist
33 and Froese, 2013] and in single-cell perturbations [Bunne et al., 2021, 2022b,a]. Standard optimal
34 transport solvers deployed in this setting would re-solve the optimization problems from scratch, but
35 this ignores the shared structure and information present between different coupling problems.

36 **Overview and outline.** We study the use of amortized optimization and machine learning methods
 37 to rapidly solve multiple optimal transport problems and predict the solution from the input measures
 38 (α, β) . This setting involves learning a *meta* model to predict the solution to the optimal transport
 39 problem, which we will refer to as *Meta Optimal Transport*. We learn Meta OT models to predict
 40 the solutions to optimal transport problems and significantly improve the computational time and
 41 number of iterations needed to solve eq. (1) between discrete (sect. 3.1) and continuous (sect. 3.2)
 42 measures. The paper is organized as follows: sect. 2 recalls the main concepts needed for the rest
 43 of the paper, in particular the formulations of the entropy regularized and unregularized optimal
 44 transport problems and the basic notions of amortized optimization; sect. 3 presents the Meta OT
 45 models and algorithms; and sect. 4 empirically demonstrates the effectiveness of Meta OT.

46 **Settings that are not Meta OT.** Meta OT is not useful in OT settings that do *not* involve *repeatedly*
 47 solving OT problems over a fixed distribution, including 1) standard generative modeling settings,
 48 such as Arjovsky et al. [2017] that estimate the OT distance between the data and model distri-
 49 butions, and 2) the out-of-sample setting of Seguy et al. [2018], Perrot et al. [2016] that couple
 50 measures and then extrapolate the map to larger measures containing the original measures.

51 2 Preliminaries and background

52 2.1 Dual optimal transport solvers

53 We review foundations of optimal transportation, following the notation of Peyré et al. [2019] in
 54 most places. The discrete setting often favors the entropic regularized version since it can be com-
 55 puted efficiently and in a parallelized way using the Sinkhorn algorithm. On the other hand, the
 56 continuous setting is often solved from samples using convex potentials. While the primal Kan-
 57 torovich formulation in eq. (1) provides an intuitive problem description, optimal transport problems
 58 are rarely solved directly in this form due to the high-dimensionality of the couplings π and the diffi-
 59 culty of satisfying the coupling constraints $\mathcal{U}(\alpha, \beta)$. Instead, most computational OT solvers use the
 60 *dual* of eq. (1), which we build our Meta OT solvers on top of in discrete and continuous settings.

61 2.1.1 Entropic OT between discrete measures with the Sinkhorn algorithm

62 Let $\alpha := \sum_{i=1}^m a_i \delta_{x_i}$ and $\beta := \sum_{i=1}^n b_i \delta_{y_i}$ be
 63 *discrete* measures, where δ_z is a Dirac at point
 64 z and $a \in \Delta_{m-1}$ and $b \in \Delta_{n-1}$ are in the
 65 *probability simplex* defined by

$$\Delta_{k-1} := \{x \in \mathbb{R}^k : x \geq 0 \text{ and } \sum_i x_i = 1\}. \quad (2)$$

Algorithm 1 Sinkhorn($\alpha, \beta, c, \epsilon, f_0 = 0$)

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

$$g_i \leftarrow \epsilon \log b - \epsilon \log (K^\top \exp\{f_{i-1}/\epsilon\})$$

$$f_i \leftarrow \epsilon \log a - \epsilon \log (K \exp\{g_i/\epsilon\})$$

end for

Compute P_N from f_N, g_N using eq. (6)

return $P_N \approx P^*$

66 **Discrete OT.** In the discrete setting, eq. (1) simplifies to the *linear program*

$$P^*(\alpha, \beta, c) \in \arg \min_{P \in U(a,b)} \langle C, P \rangle \quad U(a, b) := \{P \in \mathbb{R}_+^{n \times m} : P1_m = a, \quad P^\top 1_n = b\} \quad (3)$$

67 where P is a *coupling matrix*, $P^*(\alpha, \beta)$ is the *optimal* coupling, and the *cost* can be discretized as a
 68 matrix $C \in \mathbb{R}^{m \times n}$ with entries $C_{i,j} := c(x_i, y_j)$, and $\langle C, P \rangle := \sum_{i,j} C_{i,j} P_{i,j}$,

69 **Entropic OT.** The linear program above can be regularized adding the entropy of the coupling to
 70 smooth the objective as in Cominetti and Martín [1994], Cuturi [2013], resulting in:

$$P^*(\alpha, \beta, c, \epsilon) \in \arg \min_{P \in U(a,b)} \langle C, P \rangle - \epsilon H(P) \quad (4)$$

71 where $H(P) := -\sum_{i,j} P_{i,j} (\log(P_{i,j}) - 1)$ is the discrete entropy of a coupling matrix P .

72 **Entropic OT dual.** As presented in Peyré et al. [2019, Prop. 4.4], the dual of eq. (4) is

$$f^*, g^* \in \arg \max_{f \in \mathbb{R}^n, g \in \mathbb{R}^m} \langle f, a \rangle + \langle g, b \rangle - \epsilon \langle \exp\{f/\epsilon\}, K \exp\{g/\epsilon\} \rangle, \quad K_{i,j} := \exp\{-C_{i,j}/\epsilon\}, \quad (5)$$

73 where $K \in \mathbb{R}^{m \times n}$ is the *Gibbs kernel* and the *dual variables* or *potentials* $f \in \mathbb{R}^n$ and $g \in \mathbb{R}^m$ are
 74 associated, respectively, with the marginal constraints $P1_m = a$ and $P^\top 1_n = b$. The optimal duals
 75 depend on the problem, e.g. $f^*(\alpha, \beta, c, \epsilon)$, but we omit this dependence for notational simplicity.

76 **Recovering the primal solution from the duals.** Given optimal duals f^*, g^* that solve eq. (5) the
 77 optimal coupling P^* to the primal problem in eq. (4) can be obtained by

$$P_{i,j}^*(\alpha, \beta, c, \epsilon) := \exp\{f_i^*/\epsilon\} K_{i,j} \exp\{g_j^*/\epsilon\} \quad (K \text{ is defined in eq. (5)}) \quad (6)$$

78 **The Sinkhorn algorithm.** Algorithm 1 summarizes the log-space version, which takes closed-form
 79 block coordinate ascent updates on eq. (5) obtained from the first-order optimality conditions [Peyré
 80 et al., 2019, Remark 4.21]. We will use it to fine-tune predictions made by our Meta OT models.

81 **Computing the error.** Standard implementations of the Sinkhorn algorithm, such as Flamary et al.
 82 [2021], Cuturi et al. [2022], measure the error of a candidate dual solution (f, g) by computing the
 83 deviation from the marginal constraints, which we will also use in comparing our solution quality:

$$\text{err}(f, g; \alpha, \beta, c) := \|P1_m - a\|_1 + \|P^\top 1_n - b\|_1 \quad (\text{compute } P \text{ from eq. (6)}) \quad (7)$$

84 **Mapping between the duals.** The first-order optimality conditions of eq. (5) also provide an equiv-
 85 alence between the optimal dual potentials that we will make use of:

$$g(f; b, c) := \epsilon \log b - \epsilon \log (K^\top \exp\{f/\epsilon\}). \quad (8)$$

86 2.1.2 Wasserstein-2 OT between continuous (Euclidean) measures with dual potentials

87 Let α and β be continuous measures in Euclidean
 88 space $\mathcal{X} = \mathcal{Y} = \mathbb{R}^d$ (with α absolutely contin-
 89 uous with respect to the Lebesgue measure) and
 90 the ground cost be the squared Euclidean distance
 91 $c(x, y) := \|x - y\|_2^2$. Then the minimum of eq. (1)
 92 defines the square of the Wasserstein-2 distance:

$$W_2^2(\alpha, \beta) := \min_{\pi \in \mathcal{U}(\alpha, \beta)} \int_{\mathcal{X} \times \mathcal{Y}} \|x - y\|_2^2 d\pi(x, y) = \min_T \int_{\mathcal{X}} \|x - T(x)\|_2^2 d\alpha(x), \quad (9)$$

93 where T is a *transport map* pushing α to β , i.e. $T\# \alpha = \beta$ with the *pushforward operator* defined
 94 by $T\# \alpha(B) := \alpha(T^{-1}(B))$ for any measurable set B .

95 **Convex dual potentials.** The primal form in eq. (9) is difficult to solve, as in the discrete setting, due
 96 to the difficulty of representing the coupling and satisfying the constraints. Makkuva et al. [2020],
 97 Taghvaei and Jalali [2019], Korotin et al. [2019, 2021b, 2022] propose to instead solve the dual:

$$\psi^*(\cdot; \alpha, \beta) \in \arg \min_{\psi \in \text{convex}} \int_{\mathcal{X}} \psi(x) d\alpha(x) + \int_{\mathcal{Y}} \bar{\psi}(y) d\beta(y), \quad (10)$$

98 where ψ is a convex function referred to as a *convex potential*, and $\bar{\psi}(y) := \max_{x \in \mathcal{X}} \langle x, y \rangle - \psi(x)$ is
 99 the *Legendre-Fenchel transform* or *convex conjugate* of ψ [Fenchel, 1949, Rockafellar, 2015]. The
 100 potential ψ is often approximated with an input-convex neural network (ICNN) [Amos et al., 2017].

101 **Recovering the primal solution from the dual.** Given an optimal dual ψ^* for eq. (10), Brenier
 102 [1991] remarkably shows that an optimal map T^* for eq. (9) can be obtained with differentiation:

$$T^*(x) = \nabla_x \psi^*(x). \quad (11)$$

103 **Wasserstein-2 Generative Networks (W2GNs).** Korotin et al. [2019] model ψ_φ and $\bar{\psi}_\varphi$ in eq. (10)
 104 with two separate ICNNs parameterized by φ . The separate model for $\bar{\psi}_\varphi$ is useful because the
 105 conjugate operation in eq. (10) becomes computationally expensive. They optimize the loss:

$$\mathcal{L}(\varphi) := \underbrace{\mathbb{E}_{x \sim \alpha} [\psi_\varphi(x)] + \mathbb{E}_{y \sim \beta} [\langle \nabla \bar{\psi}_\varphi(y), y \rangle - \psi_\varphi(\nabla \bar{\psi}_\varphi(y))]}_{\text{Cyclic monotone correlations (dual objective)}} + \underbrace{\gamma \mathbb{E}_{y \sim \beta} \|\nabla \psi_\varphi \circ \nabla \bar{\psi}_\varphi(y) - y\|_2^2}_{\text{Cycle-consistency regularizer}} \quad (12)$$

106 where φ is a detached copy of the parameters and γ is a hyper-parameter. The first term are the
 107 *cyclic monotone correlations* [Chartrand et al., 2009, Taghvaei and Jalali, 2019], that optimize the
 108 dual objective in eq. (10), and the second term provides *cycle consistency* [Zhu et al., 2017] to
 109 estimate the conjugate $\bar{\psi}$. Algorithm 2 shows how \mathcal{L} is typically optimized using samples from the
 110 measures, which we use to fine-tune Meta OT predictions.

Algorithm 2 W2GN(α, β, φ_0)

for iteration $i = 1$ to N **do**
 Sample from (α, β) and estimate $\mathcal{L}(\varphi_{i-1})$
 Update φ_i with approximation to $\nabla_\varphi \mathcal{L}(\varphi_{i-1})$
end for
return $T_N(\cdot) := \nabla_x \psi_{\varphi_N}(\cdot) \approx T^*(\cdot)$

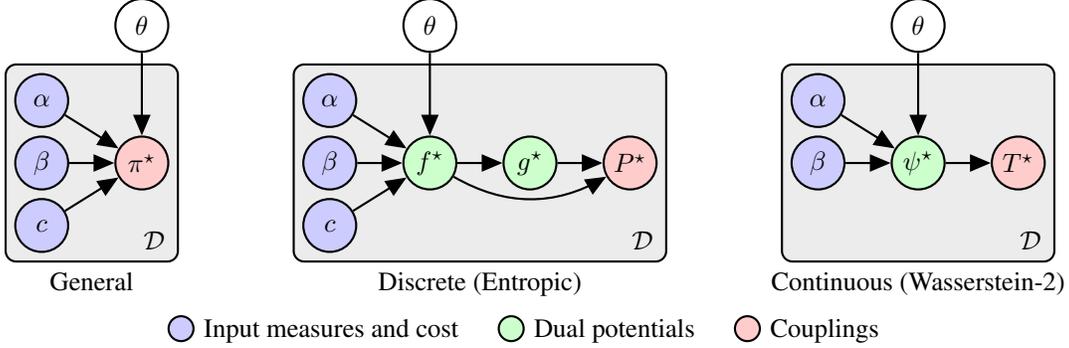


Figure 1: Meta OT uses objective-based amortization for optimal transport. In the general formulation, the *parameters* θ capture shared structure in the *optimal couplings* π^* between multiple input measures and costs over some *distribution* \mathcal{D} . In practice, we learn this shared structure over the *dual potentials* which map back to the coupling: f^* in discrete settings and ψ^* in continuous ones.

111 2.2 Amortized optimization and learning to optimize

112 Our paper is an application of amortized optimization methods that predict the solutions of opti-
 113 mization problems, as surveyed in, e.g., Chen et al. [2021], Amos [2022]. We use the basic setup
 114 from Amos [2022], which considers unconstrained continuous optimization problems of the form

$$z^*(\phi) \in \arg \min_z J(z; \phi), \quad (13)$$

115 where J is the objective, $z \in \mathcal{Z}$ is the *domain*, and $\phi \in \Phi$ is some *context* or *parameterization*. In
 116 other words, the context conditions the objective but is not optimized over. Given a *distribution over*
 117 *contexts* $\mathcal{P}(\phi)$, we learn a model \hat{z}_θ parameterized by θ to approximate eq. (13), i.e. $\hat{z}_\theta(\phi) \approx z^*(\phi)$.
 118 J will be differentiable for us, so we optimize the parameters using *objective-based learning* with

$$\min_{\theta} \mathbb{E}_{\phi \sim \mathcal{P}(\phi)} J(\hat{z}_\theta(\phi); \phi), \quad (14)$$

119 which does *not* require ground-truth solutions z^* and can be optimized with a gradient-based solver.
 120 While we focus on optimizing eq. (14) because we do not assume easy access to ground-truth solu-
 121 tions $z^*(\phi)$, one alternative is *regression-based learning* if the solutions are easily available:

$$\min_{\theta} \mathbb{E}_{\phi \sim \mathcal{P}(\phi)} \|z^*(\phi) - \hat{z}_\theta(\phi)\|_2^2. \quad (15)$$

122 3 Meta Optimal Transport

123 Figure 1 illustrates our key contribution of connecting objective-based amortization in eq. (14) to
 124 optimal transport. We consider solving *multiple* OT problems and learning shared structure and
 125 correlations between them. We denote a joint *meta-distribution* over the input measures and costs
 126 with $\mathcal{D}(\alpha, \beta, c)$, which we call *meta* to distinguish it from the measures α, β .

127 In general, we could introduce a model that directly predicts the primal solution to eq. (1), i.e.
 128 $\pi_\theta(\alpha, \beta, c) \approx \pi^*(\alpha, \beta, c)$ for $(\alpha, \beta, c) \sim \mathcal{D}$. This is difficult for the same reason why most compu-
 129 tational methods do not operate directly in the primal space: the optimal coupling is often a high-
 130 dimensional joint distribution with non-trivial marginal constraints. We instead turn to predicting
 131 the dual variables used by today’s solvers.

132 3.1 Meta OT between discrete measures

133 We build on standard methods for entropic OT reviewed in sect. 2.1.1 between discrete measures
 134 $\alpha := \sum_{i=1}^m a_i \delta_{x_i}$ and $\beta := \sum_{i=1}^n b_i \delta_{x_i}$ with $a \in \Delta_{m-1}$ and $b \in \Delta_{n-1}$ coupled using a cost c . In the
 135 Meta OT setting, the measures and cost are the contexts for amortization and sampled from a *meta-*
 136 *distribution*, i.e. $(\alpha, \beta, c) \sim \mathcal{D}(\alpha, \beta, c)$. For example, sects. 4.1 and 4.2 considers meta-distributions
 137 over the weights of the atoms, i.e. $(a, b) \sim \mathcal{D}$, where \mathcal{D} is a distribution over $\Delta_{m-1} \times \Delta_{n-1}$.

Algorithm 3 Training Meta OT

Initialize amortization model with θ_0
for iteration i **do**
 Sample $(\alpha, \beta, c) \sim \mathcal{D}$
 Predict duals \hat{f}_θ or $\hat{\varphi}_\theta$ on the sample
 Estimate the loss in eq. (17) or eq. (18)
 Update θ_{i+1} with a gradient step
end for

Algorithm 4 Fine-tuning with Sinkhorn

Predict duals $\hat{f}_\theta(\alpha, \beta, c)$
return Sinkhorn($\alpha, \beta, c, \epsilon, \hat{f}_\theta$)

Algorithm 5 Fine-tuning with W2GN

Predict dual ICNN parameters $\hat{\varphi}_\theta(\alpha, \beta, c)$
return W2GN($\alpha, \beta, c, T, \hat{\varphi}_\theta$)

138 **Amortization objective.** We will seek to predict the *optimal* potential. At optimality, the pair of
139 potentials are related to each other via eq. (8), i.e. $g(f; \alpha, \beta, c) := \epsilon \log b - \epsilon \log (K^\top \exp\{f/\epsilon\})$
140 where $K \in \mathbb{R}^{m \times n}$ is the *Gibbs kernel* from eq. (5). Hence, it is sufficient to predict one of the
141 potentials, e.g. f , and recover the other. We thus re-formulate eq. (5) to just optimize over f with

$$f^*(\alpha, \beta, c, \epsilon) \in \arg \min_{f \in \mathbb{R}^n} J(f; \alpha, \beta, c), \quad (16)$$

142 where $-J(f; \alpha, \beta, c) := \langle f, a \rangle + \langle g, b \rangle - \epsilon \langle \exp\{f/\epsilon\}, K \exp\{g/\epsilon\} \rangle$ is the (negated) dual objective.
143 Even though most solvers optimize over f and g jointly as in eq. (16), amortizing over these would
144 likely need: 1) to have a higher capacity than a model just predicting f , and 2) to learn how f and g
145 are connected through eq. (8) while in eq. (16) we explicitly provide this knowledge.

146 **Amortization model.** We predict the solution to eq. (16) with $\hat{f}_\theta(\alpha, \beta, c)$ parameterized by θ ,
147 resulting in a computationally efficient approximation $\hat{f}_\theta \approx f^*$. Here we use the notation $\hat{f}_\theta(\alpha, \beta, c)$
148 to mean that the model \hat{f}_θ depends on *representations* of the input measures and cost. In our settings,
149 we define \hat{f}_θ as a fully-connected MLP mapping from the atoms of the measures to the duals.

150 **Amortization loss.** Applying objective-based amortization from eq. (14) to the dual in eq. (16)
151 completes our learning setup. Our model should best-optimize the expectation of the dual objective

$$\min_{\theta} \mathbb{E}_{(\alpha, \beta, c) \sim \mathcal{D}} J(\hat{f}_\theta(\alpha, \beta, c); \alpha, \beta, c), \quad (17)$$

152 which is appealing as it does not require ground-truth solutions f^* . Algorithm 3 shows a basic
153 training loop for eq. (17) using a gradient-based optimizer such as Adam [Kingma and Ba, 2014].

154 **Sinkhorn fine-tuning.** The dual prediction made by \hat{f}_θ with an associated \hat{g} can easily be input as
155 the initialization to a standard Sinkhorn solver as shown in algorithm 4. This allows us to deploy the
156 predicted potential with Sinkhorn to obtain the optimal potentials with only a few extra iterations.

157 **On accelerated solvers.** Here we have only considered fine-tuning the Meta OT prediction with
158 a log-Sinkhorn solver. Meta OT can also be combined with accelerated variants of entropic OT
159 solvers such as Thibault et al. [2017], Altschuler et al. [2017], Alaya et al. [2019], Lin et al. [2019]
160 that would otherwise solve every problem from scratch.

161 3.2 Meta OT between continuous measures (Wasserstein-2)

162 We take an analogous approach to predicting the Wasserstein-2 map between continuous measures
163 for Wasserstein-2 as reviewed in sect. 2.1.2. Here the measures α, β are supported in continuous
164 space $\mathcal{X} = \mathcal{Y} = \mathbb{R}^d$ and we focus on computing Wasserstein-2 couplings from instances sampled
165 from a *meta-distribution* $(\alpha, \beta) \sim \mathcal{D}(\alpha, \beta)$. The cost c is not included in \mathcal{D} as it remains fixed to the
166 squared Euclidean cost everywhere here.

167 One challenge here is that the optimal dual potential $\psi^*(\cdot; \alpha, \beta)$ in eq. (10) is a convex function and
168 not simply a finite-dimensional real vector. The dual potentials in this setting are approximated by,
169 e.g., an ICNN. We thus propose a *Meta ICNN* that predicts the *parameters* φ of an ICNN ψ_φ that
170 approximates the optimal dual potentials, which can be seen as a hypernetwork [Stanley et al., 2009,
171 Ha et al., 2016]. The dual prediction made by $\hat{\varphi}_\theta$ can easily be input as the initial value to a standard
172 W2GN solver as shown in algorithm 5. App. B discusses other modeling choices we considered:
173 we tried models based on MAML [Finn et al., 2017] and neural processes [Garnelo et al., 2018b,a].

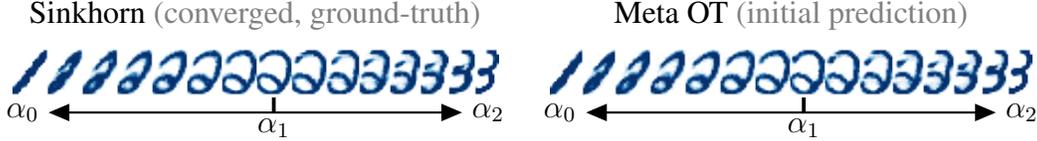


Figure 2: Interpolations between MNIST test digits using couplings obtained from (left) solving the problem with Sinkhorn, and (right) Meta OT model’s initial prediction, which is ≈ 100 times computationally cheaper and produces a nearly identical coupling.

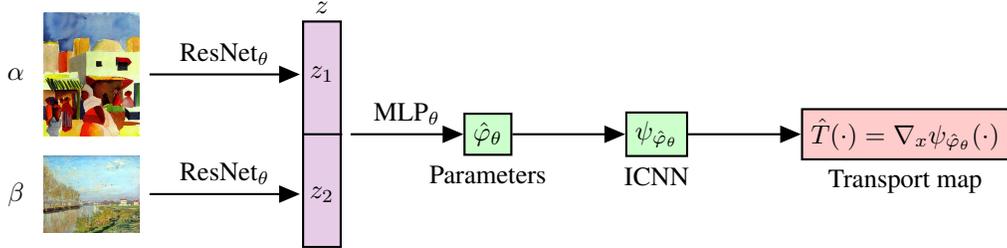


Figure 3: A Meta ICNN for image-based input measures. A shared ResNet processes the input measures α and β into latents z that are decoded with an MLP into the parameters φ of an ICNN dual potential ψ_φ . The derivative of the ICNN provides the transport map \hat{T} .

Table 1: Sinkhorn runtime (seconds) to reach a marginal error of 10^{-3} . Meta OT’s initial prediction takes $\approx 5 \cdot 10^{-5}$ seconds.

Initialization	MNIST	Spherical
Zeros	$7.7 \cdot 10^{-3} \pm 1.2 \cdot 10^{-3}$	$1.4 \pm 1.9 \cdot 10^{-1}$
Gaussian	$7.7 \cdot 10^{-3} \pm 1.4 \cdot 10^{-3}$	$1.1 \pm 2.0 \cdot 10^{-1}$
Meta OT	$3.9 \cdot 10^{-3} \pm 1.6 \cdot 10^{-3}$	$0.44 \pm 1.5 \cdot 10^{-1}$

Table 2: Color transfer runtimes and values.

	Iter	Runtime (s)	Dual Value
Meta OT + W2GN	None	$3.5 \cdot 10^{-3} \pm 2.7 \cdot 10^{-4}$	$0.90 \pm 6.08 \cdot 10^{-2}$
	1k	$0.93 \pm 2.27 \cdot 10^{-2}$	$1.0 \pm 2.57 \cdot 10^{-3}$
	2k	$1.84 \pm 3.78 \cdot 10^{-2}$	$1.0 \pm 5.30 \cdot 10^{-3}$
W2GN	1k	$0.90 \pm 1.62 \cdot 10^{-2}$	$0.96 \pm 2.62 \cdot 10^{-2}$
	2k	$1.81 \pm 3.05 \cdot 10^{-2}$	$0.99 \pm 1.14 \cdot 10^{-2}$

We report the mean and standard deviation across 10 test instances.

174 **Amortization objective.** We build on the W2GN formulation [Korotin et al., 2019] and seek parameters φ^* optimizing the dual ICNN potentials ψ_φ and $\overline{\psi_\varphi}$ with $\mathcal{L}(\varphi; \alpha, \beta)$ from eq. (12). We
 175 chose W2GN due to the stability, but could also easily use other losses optimizing ICNN potentials.
 176

177 **Amortization model: the Meta ICNN.** We predict the solution to eq. (12) with $\hat{\varphi}_\theta(\alpha, \beta)$ parameterized by θ , resulting in a computationally efficient approximation to the optimum $\hat{\varphi}_\theta \approx \varphi^*$.
 178 Figure 3 instantiates a convolutional Meta ICNN model using a ResNet-18 [He et al., 2016] architecture for coupling image-based measures. We again emphasize that α, β used with the model here
 179 are *representations* of measures, which in our cases are simply images.
 180
 181

182 **Amortization loss.** Applying objective-based amortization from eq. (14) to the W2GN loss in
 183 eq. (12) completes our learning setup. Our model should best-optimize the expectation of the loss:

$$\min_{\theta} \mathbb{E}_{(\alpha, \beta) \sim \mathcal{D}} \mathcal{L}(\hat{\varphi}_\theta(\alpha, \beta); \alpha, \beta). \quad (18)$$

184 As in the discrete setting, it does not require ground-truth solutions φ^* and we learn it with Adam.

185 4 Experiments

186 We demonstrate how Meta OT models improve the convergence of the state-of-the-art solvers in
 187 settings where solving multiple OT problems naturally arises. We implemented our code in JAX
 188 [Bradbury et al., 2018] as an extension to the the Optimal Transport Tools (OTT) package [Cuturi
 189 et al., 2022]. App. C covers further experimental and implementation details, and shows that all of
 190 our experiments take a few hours to run on our single Quadro GP100 GPU.

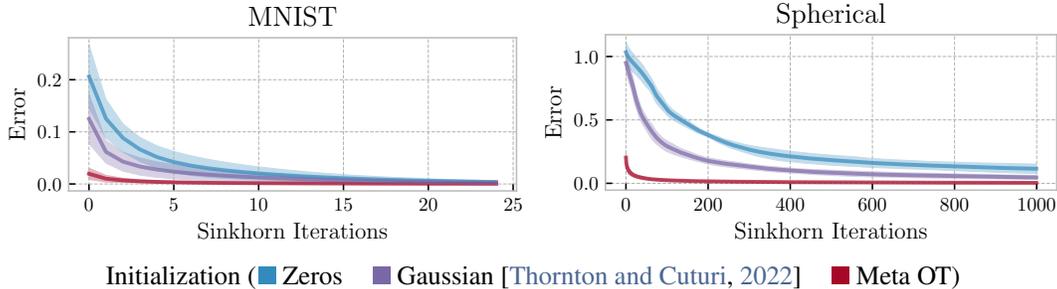


Figure 4: Meta OT successfully predicts warm-start initializations that significantly improve the convergence of Sinkhorn iterations on test data. The error is the marginal error defined in eq. (7).

191 4.1 Discrete OT between MNIST digits

192 Images provide a natural setting for Meta OT where the distribution over images provide the meta-
 193 distribution \mathcal{D} over OT problems. Given a pair of images α_0 and α_1 , each grayscale image is
 194 cast as a discrete measure in 2-dimensional space where the intensities define the probabilities of
 195 the atoms. The goal is to compute the optimal transport interpolation between the two measures
 196 as in, e.g., Peyré et al. [2019, §7]. Formally, this means computing the optimal coupling P^* by
 197 solving the entropic optimal transport problem between α_0 and α_1 and computing the interpolates
 198 as $\alpha_t = (t \text{proj}_y + (1-t) \text{proj}_x) \# P^*$, for $t \in [0, 1]$, where $\text{proj}_x(x, y) := x$ and $\text{proj}_y(x, y) = y$.
 199 We selected $\epsilon = 10^{-2}$ as app. A shows that it gives interpolations that are not too blurry or sharp.

200 Our Meta OT model \hat{f}_θ (sect. 3.1) is an MLP that predicts the transport map between pairs of MNIST
 201 digits. We train on every pair from the standard training dataset. Figure 2 shows that even without
 202 fine-tuning, Meta OT’s predicted Wasserstein interpolations between the measures are close to the
 203 ground-truth interpolations obtained from running the Sinkhorn algorithm to convergence. We then
 204 fine-tune Meta OT’s prediction with Sinkhorn as in algorithm 4. Figure 4 shows that the near-
 205 optimal predictions can be quickly refined in fewer iterations than running Sinkhorn with the default
 206 initialization, and table 1 shows the runtime required to reach the default threshold, which uses the
 207 default marginal error threshold of 10^{-3} . We compare our learned initialization to the standard zero
 208 initialization, as well as the Gaussian initialization proposed in Thornton and Cuturi [2022], which
 209 takes a continuous Gaussian approximation of the measures and initializes the potentials to be the
 210 known coupling between the Gaussians. This Gaussian initialization assumes the squared Euclidean
 211 cost, which is not the case in our spherical transport problem, but we find it is still helpful over the
 212 zero initialization.

213 4.2 Discrete OT for supply-demand transportation on spherical data

214 We next set up a synthetic transport problem between supply and demand locations where the supply
 215 and demands may change locations or quantities frequently, creating another Meta OT setting to be
 216 able to rapidly solve the new instances. We specifically consider measures living on the 2-sphere
 217 defined by $\mathcal{S}_2 := \{x \in \mathbb{R}^3 : \|x\| = 1\}$, i.e. $\mathcal{X} = \mathcal{Y} = \mathcal{S}_2$, with the transport cost given by the
 218 spherical distance $c(x, y) = \arccos(\langle x, y \rangle)$. We then randomly sample supply locations uniformly
 219 from Earth’s landmass and demand locations from Earth’s population density to induce a class of
 220 transport problems on the sphere obtained from the CC-licensed dataset from Doxsey-Whitfield et al.
 221 [2015]. Figure 5 shows that the predicted transport maps on test instances are close to the optimal
 222 maps obtained from Sinkhorn to convergence. Similar to the MNIST setting, fig. 4 and table 1 show
 223 improved convergence and runtime.

224 4.3 Continuous Wasserstein-2 color transfer

225 The problem of color transfer between two images consists in mapping the color palette of one image
 226 into the other one. The images are required to have the same number of channels, for example RGB
 227 images. The continuous formulation that we use from Korotin et al. [2019], takes i.e. $\mathcal{X} = \mathcal{Y} =$
 228 $[0, 1]^3$ with c being the squared Euclidean distance. We collected ≈ 200 public domain images from
 229 WikiArt and trained a Meta ICNN model from sect. 3.2 to predict the color transfer maps between

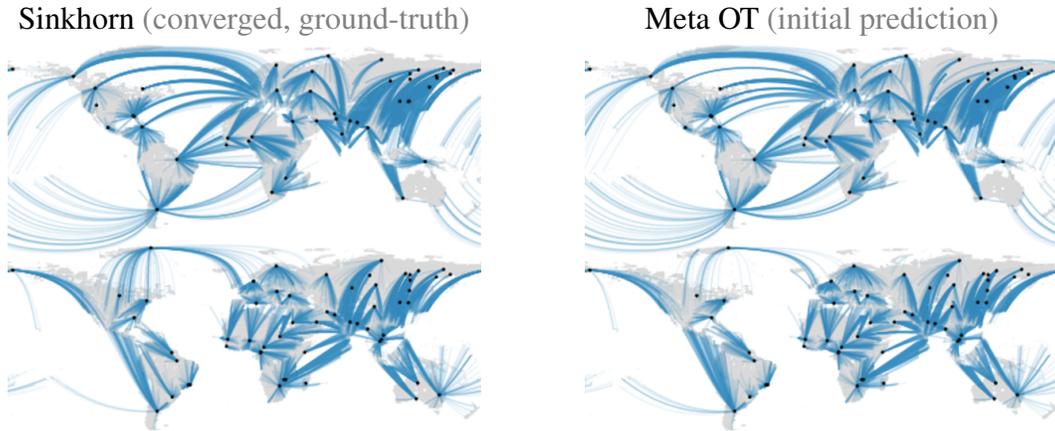


Figure 5: Test set coupling predictions of the spherical transport problem. Meta OT’s initial prediction is ≈ 37500 times faster than solving Sinkhorn to optimality. Supply locations are shown as black dots and the blue lines show the spherical transport maps T going to demand locations at the end. The sphere is visualized with the Mercator projection.

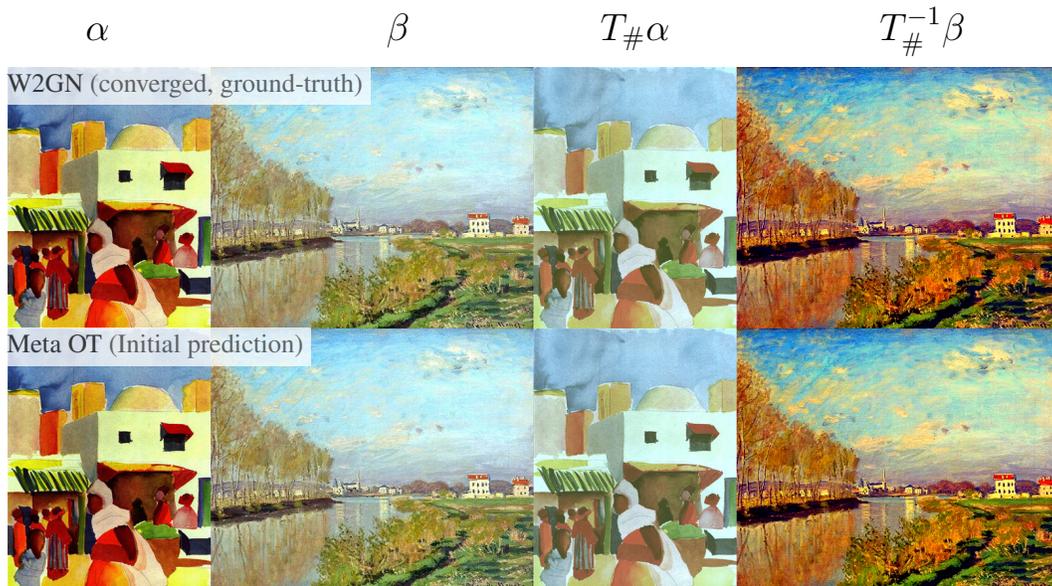


Figure 6: Color transfers with a Meta ICNN on test pairs of images. The objective is to optimally transport the continuous RGB measure of the first image α to the second β , producing an invertible transport map T . Meta OT’s prediction is ≈ 1000 times faster than training W2GN from scratch. The image generating α is *Market in Algiers* by August Macke (1914) and β is *Argenteuil, The Seine* by Claude Monet (1872), obtained from WikiArt.

230 every pair of them. Figure 6 shows the predictions on test pairs and fig. 7 shows the convergence in
 231 comparison to the standard W2GN learning. Table 2 reports runtimes and app. E shows additional
 232 results.

233 5 Related work

234 **Efficiently estimating OT maps.** To compute OT maps with fixed cost between pairs of measures
 235 efficiently, neural OT models [Korotin et al., 2019, Li et al., 2020, Korotin et al., 2021a, Mokrov
 236 et al., 2021, Korotin et al., 2021b] leverage ICNNs to estimate maps between continuous high-

237 dimensional measures given samples from these, and Litvinenko et al. [2021], Scetbon et al. [2021],
 238 Forrow et al. [2019], Sommerfeld et al. [2019], Scetbon et al. [2022], Muzellec and Cuturi [2019],
 239 Bonet et al. [2021] leverage structural assumptions on coupling and cost matrices to reduce the
 240 computational and memory complexity. In the meta-OT setting, we consider learning to rapidly
 241 compute OT mappings between new pairs measures. All these works can hence potentially benefit
 242 from an acceleration effect by leveraging amortization similarly.

243 **Embedding measures where OT distances are**
 244 **discriminative.** Effort has been invested in learning
 245 encodings/projections of measures through
 246 a nested optimization problem, which aims to
 247 find discriminative embeddings of the measures
 248 to be compared [Genevay et al., 2018, Deshpande
 249 et al., 2019, Nguyen and Ho, 2022]. While these
 250 works share an encoder and/or a projection across
 251 task with the aim of leveraging more discrimina-
 252 tive alignments (and hence an OT distance with a
 253 metric different from the Euclidean metric), our
 254 work differs in the sense that we find good initial-
 255 zations to solve the OT problem itself with fixed
 256 cost more efficiently across tasks.

257 **Optimal transport and amortization.** Few pre-
 258 vious works in the OT literature leverage amor-
 259 tization. Courty et al. [2018] learn a latent space in which the Wasserstein distance between the
 260 measure’s embeddings is equivalent to the Euclidean distance. Concurrent work [Nguyen and Ho,
 261 2022] amortizes the estimation of the optimal projection in the max-sliced objective, which differs
 262 from our work where we instead amortize the estimation of the optimal coupling directly. Also,
 263 Lacombe et al. [2021] learns to predict Wasserstein barycenters of pixel images by training a con-
 264 volutional networks that, given images as input, outputs their barycenters. Our work is hence a
 265 generalization of this pixel-based work to general measures – both discrete and continuous. A limita-
 266 tion of Lacombe et al. [2021] is that it does not provide alignments, as the amortization networks
 267 predicts the barycenter directly rather than individual couplings.

268 6 Conclusions, future directions, and limitations

269 We have presented foundations for modeling and learning to solve OT problems with Meta OT by
 270 using amortized optimization to predict optimal transport plans. This works best in applications that
 271 require solving multiple OT problems with shared structure. We instantiated it to speed up entropic
 272 regularized optimal transport and unregularized optimal transport with squared cost by multiple
 273 orders of magnitude. We envision extensions of the work in:

- 274 1. **Meta OT models.** While we mostly consider models based on hypernetworks, other meta-
 275 learning paradigms can be connected in. In the discrete setting, we only considered settings
 276 where the cost remains fixed, but the Meta OT model can also be conditioned on the cost
 277 by considering the entire cost matrix as an input (which may be too large for most models
 278 to handle), or considering a lower-dimensional parameterization of the cost that changes
 279 between the Meta OT problem instances.
- 280 2. **OT algorithms.** While we instantiated models on top of log-Sinkhorn and W2GN, Meta
 281 OT could be built on top of other methods.
- 282 3. **OT applications** that are computationally expensive and repeatedly solved, e.g. in multi-
 283 marginal and barycentric settings, or for Gromov-Wasserstein distances between metric-
 284 measure spaces.

285 **Limitations.** While we have illustrated successful applications of Meta OT, it is also important to
 286 understand the limitations: 1) **Meta OT does not make previously intractable problems tractable.**
 287 All of the baseline OT solvers we consider solve our problems within milliseconds or seconds. 2)
 288 **Out-of-distribution generalization.** Meta OT may not generate good predictions on instances that
 289 are not close to the training OT problems from the meta-distribution \mathcal{D} over the measures and cost.
 290 If the model makes a bad prediction, one fallback option is to re-solve the instance from scratch.

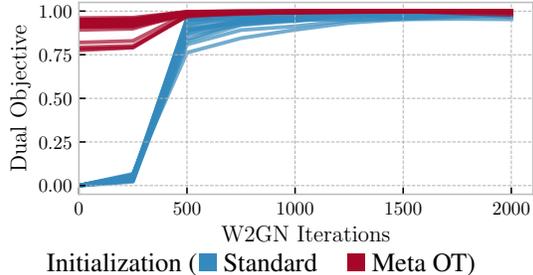


Figure 7: Convergence on color transfer test instances using W2GN. Meta ICNNs predicts warm-start initializations that significantly improve the (normalized) dual objective values.

291 References

- 292 Mokhtar Z Alaya, Maxime Berar, Gilles Gasso, and Alain Rakotomamonjy. Screening sinkhorn algorithm for
293 regularized optimal transport. *Advances in Neural Information Processing Systems*, 32, 2019.
- 294 Jason Altschuler, Jonathan Niles-Weed, and Philippe Rigollet. Near-linear time approximation algorithms for
295 optimal transport via sinkhorn iteration. *Advances in neural information processing systems*, 30, 2017.
- 296 Luigi Ambrosio. Lecture notes on optimal transport problems. In *Mathematical aspects of evolving interfaces*,
297 pages 1–52. Springer, 2003.
- 298 Brandon Amos. Tutorial on amortized optimization for learning to optimize over continuous domains. *arXiv*
299 *preprint arXiv:2202.00665*, 2022.
- 300 Brandon Amos, Lei Xu, and J Zico Kolter. Input convex neural networks. In *International Conference on*
301 *Machine Learning*, pages 146–155. PMLR, 2017.
- 302 Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *Inter-*
303 *national conference on machine learning*, pages 214–223. PMLR, 2017.
- 304 Clément Bonet, Titouan Vayer, Nicolas Courty, François Septier, and Lucas Drumetz. Subspace detours meet
305 gromov–wasserstein. *Algorithms*, 14(12):366, 2021.
- 306 James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George
307 Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable trans-
308 formations of Python+NumPy programs. *GitHub*, 2018. URL <http://github.com/google/jax>.
- 309 Yann Brenier. Polar factorization and monotone rearrangement of vector-valued functions. *Communications*
310 *on pure and applied mathematics*, 44(4):375–417, 1991.
- 311 Charlotte Bunne, Stefan G Stark, Gabriele Gut, Jacobo Sarabia del Castillo, Kjong-Van Lehmann, Lucas Pelk-
312 mans, Andreas Krause, and Gunnar Ratsch. Learning single-cell perturbation responses using neural optimal
313 transport. *bioRxiv*, 2021.
- 314 Charlotte Bunne, Andreas Krause, and Marco Cuturi. Supervised training of conditional monge maps. *arXiv*
315 *preprint arXiv:2206.14262*, 2022a.
- 316 Charlotte Bunne, Laetitia Papaxanthos, Andreas Krause, and Marco Cuturi. Proximal optimal transport mod-
317 eling of population dynamics. In *International Conference on Artificial Intelligence and Statistics*, pages
318 6511–6528. PMLR, 2022b.
- 319 Rick Chartrand, Brendt Wohlberg, Kevin Vixie, and Erik Bollt. A gradient descent solution to the monge-
320 kantorovich problem. *Applied Mathematical Sciences*, 3(22):1071–1080, 2009.
- 321 Tianlong Chen, Xiaohan Chen, Wuyang Chen, Howard Heaton, Jialin Liu, Zhangyang Wang, and Wotao Yin.
322 Learning to optimize: A primer and a benchmark. *arXiv preprint arXiv:2103.12828*, 2021.
- 323 Samuel Cohen, Brandon Amos, and Yaron Lipman. Riemannian convex potential maps. In *International*
324 *Conference on Machine Learning*, pages 2028–2038. PMLR, 2021.
- 325 Roberto Cominetti and J San Martín. Asymptotic analysis of the exponential penalty trajectory in linear pro-
326 gramming. *Mathematical Programming*, 67(1):169–187, 1994.
- 327 Nicolas Courty, Rémi Flamary, Amaury Habrard, and Alain Rakotomamonjy. Joint distribution optimal trans-
328 portation for domain adaptation. *Advances in Neural Information Processing Systems*, 30, 2017.
- 329 Nicolas Courty, Rémi Flamary, and Mélanie Ducoffe. Learning wasserstein embeddings. In *International*
330 *Conference on Learning Representations*, 2018.
- 331 Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. *Advances in neural informa-*
332 *tion processing systems*, 26:2292–2300, 2013.
- 333 Marco Cuturi, Laetitia Meng-Papaxanthos, Yingtao Tian, Charlotte Bunne, Geoff Davis, and Olivier Teboul.
334 Optimal transport tools (ott): A jax toolbox for all things wasserstein. *arXiv preprint arXiv:2201.12324*,
335 2022.
- 336 Robert Dadashi, Léonard Hussenot, Matthieu Geist, and Olivier Pietquin. Primal wasserstein imitation learning.
337 In *ICLR 2021-Ninth International Conference on Learning Representations*, 2021.

- 338 Ishan Deshpande, Yuan-Ting Hu, Ruoyu Sun, Ayis Pyrros, Nasir Siddiqui, Sanmi Koyejo, Zhizhen Zhao, David
339 Forsyth, and Alexander G Schwing. Max-sliced wasserstein distance and its use for gans. In *Proceedings of
340 the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10648–10656, 2019.
- 341 Erin Doxsey-Whitfield, Kytt MacManus, Susana B Adamo, Linda Pistoiesi, John Squires, Olena Borkovska,
342 and Sandra R Baptista. Taking advantage of the improved availability of census data: a first look at the
343 gridded population of the world, version 4. *Papers in Applied Geography*, 1(3):226–234, 2015.
- 344 Bjorn Engquist and Brittany D Froese. Application of the wasserstein metric to seismic signals. *arXiv preprint
345 arXiv:1311.4581*, 2013.
- 346 Werner Fenchel. On conjugate convex functions. *Canadian Journal of Mathematics*, 1(1):73–77, 1949.
- 347 Arnaud Fickinger, Samuel Cohen, Stuart Russell, and Brandon Amos. Cross-domain imitation learning via
348 optimal transport. In *International Conference on Learning Representations*, 2021.
- 349 Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep
350 networks. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference
351 on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 1126–1135. PMLR,
352 06–11 Aug 2017. URL <https://proceedings.mlr.press/v70/finn17a.html>.
- 353 Rémi Flamary, Nicolas Courty, Alexandre Gramfort, Mokhtar Z Alaya, Aurélie Boisbunon, Stanislas Chambon,
354 Laetitia Chapel, Adrien Corenflos, Kilian Fatras, Nemo Fournier, et al. Pot: Python optimal transport.
355 *Journal of Machine Learning Research*, 22(78):1–8, 2021.
- 356 Aden Forrow, Jan-Christian Hütter, Mor Nitzan, Philippe Rigollet, Geoffrey Schiebinger, and Jonathan Weed.
357 Statistical optimal transport via factored couplings. In *The 22nd International Conference on Artificial
358 Intelligence and Statistics*, pages 2454–2465. PMLR, 2019.
- 359 Alfred Galichon. Optimal transport methods in economics. In *Optimal Transport Methods in Economics*.
360 Princeton University Press, 2016.
- 361 Marta Garnelo, Dan Rosenbaum, Christopher Maddison, Tiago Ramalho, David Saxton, Murray Shanahan,
362 Yee Whye Teh, Danilo Rezende, and SM Ali Eslami. Conditional neural processes. In *International Con-
363 ference on Machine Learning*, pages 1704–1713. PMLR, 2018a.
- 364 Marta Garnelo, Jonathan Schwarz, Dan Rosenbaum, Fabio Viola, Danilo J Rezende, SM Eslami, and Yee Whye
365 Teh. Neural processes. *arXiv preprint arXiv:1807.01622*, 2018b.
- 366 Aude Genevay, Gabriel Peyre, and Marco Cuturi. Learning generative models with sinkhorn divergences. In
367 Amos Storkey and Fernando Perez-Cruz, editors, *Proceedings of the Twenty-First International Conference
368 on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of Machine Learning Research*, pages
369 1608–1617. PMLR, 09–11 Apr 2018. URL [https://proceedings.mlr.press/v84/genevay18a.
370 html](https://proceedings.mlr.press/v84/genevay18a.html).
- 371 David Ha, Andrew Dai, and Quoc V Le. Hypernetworks. *arXiv preprint arXiv:1609.09106*, 2016.
- 372 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In
373 *European conference on computer vision*, pages 630–645. Springer, 2016.
- 374 Chin-Wei Huang, Ricky TQ Chen, Christos Tsirigotis, and Aaron Courville. Convex potential flows: Universal
375 probability distributions with optimal transport and convex optimization. In *International Conference on
376 Learning Representations*, 2020.
- 377 J. J. Hull. A database for handwritten text recognition research. *IEEE Transactions on Pattern Analysis and
378 Machine Intelligence*, 16(5):550–554, 1994. doi: 10.1109/34.291440.
- 379 Ray Jiang, Aldo Pacchiano, Tom Stepleton, Heinrich Jiang, and Silvia Chiappa. Wasserstein fair classification.
380 In *Uncertainty in Artificial Intelligence*, pages 862–872. PMLR, 2020.
- 381 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint
382 arXiv:1412.6980*, 2014.
- 383 Nicholas Kolkin, Jason Salavon, and Gregory Shakhnarovich. Style transfer by relaxed optimal transport and
384 self-similarity. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
385 pages 10051–10060, 2019.
- 386 Soheil Kolouri, Yang Zou, and Gustavo K Rohde. Sliced wasserstein kernels for probability distributions. In
387 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5258–5267, 2016.

- 388 Soheil Kolouri, Se Rim Park, Matthew Thorpe, Dejan Slepcev, and Gustavo K Rohde. Optimal mass transport:
389 Signal processing and machine-learning applications. *IEEE signal processing magazine*, 34(4):43–59, 2017.
- 390 Soheil Kolouri, Phillip E Pope, Charles E Martin, and Gustavo K Rohde. Sliced wasserstein auto-encoders. In
391 *International Conference on Learning Representations*, 2018.
- 392 Soheil Kolouri, Kimia Nadjahi, Umut Simsekli, Roland Badeau, and Gustavo K Rohde. Generalized sliced
393 wasserstein distances. *arXiv preprint arXiv:1902.00434*, 2019.
- 394 Alexander Korotin, Vage Egiazarian, Arip Asadulaev, Alexander Safin, and Evgeny Burnaev. Wasserstein-2
395 generative networks. *arXiv preprint arXiv:1909.13082*, 2019.
- 396 Alexander Korotin, Lingxiao Li, Aude Genevay, Justin M Solomon, Alexander Filippov, and Evgeny Burnaev.
397 Do neural optimal transport solvers work? a continuous wasserstein-2 benchmark. *Advances in Neural
398 Information Processing Systems*, 34:14593–14605, 2021a.
- 399 Alexander Korotin, Lingxiao Li, Justin Solomon, and Evgeny Burnaev. Continuous wasserstein-2 barycenter
400 estimation without minimax optimization. *arXiv preprint arXiv:2102.01752*, 2021b.
- 401 Alexander Korotin, Daniil Selikhanovych, and Evgeny Burnaev. Neural optimal transport. *arXiv preprint
402 arXiv:2201.12220*, 2022.
- 403 Julien Lacombe, Julie Digne, Nicolas Courty, and Nicolas Bonneel. Learning to generate wasserstein barycen-
404 ters, 2021. URL <https://openreview.net/forum?id=2ioNzsz61vw>.
- 405 Lingxiao Li, Aude Genevay, Mikhail Yurochkin, and Justin Solomon. Continuous regularized wasserstein
406 barycenters. *arXiv preprint arXiv:2008.12534*, 2020.
- 407 Tianyi Lin, Nhat Ho, and Michael I Jordan. On the acceleration of the sinkhorn and greenhorn algorithms for
408 optimal transport. *arXiv preprint arXiv:1906.01437*, 2019.
- 409 Alexander Litvinenko, Youssef Marzouk, Hermann G Matthies, Marco Scavino, and Alessio Spantini. Com-
410 puting f-divergences and distances of high-dimensional probability density functions—low-rank tensor ap-
411 proximations. *arXiv preprint arXiv:2111.07164*, 2021.
- 412 Ashok Makkuva, Amirhossein Taghvaei, Sewoong Oh, and Jason Lee. Optimal transport mapping via input
413 convex neural networks. In *International Conference on Machine Learning*, pages 6672–6681. PMLR, 2020.
- 414 Quentin Merigot and Boris Thibert. Optimal transport: discretization and algorithms. In *Handbook of Numer-
415 ical Analysis*, volume 22, pages 133–212. Elsevier, 2021.
- 416 Petr Mokrov, Alexander Korotin, Lingxiao Li, Aude Genevay, Justin M Solomon, and Evgeny Burnaev. Large-
417 scale wasserstein gradient flows. *Advances in Neural Information Processing Systems*, 34:15243–15256,
418 2021.
- 419 Boris Muzellec and Marco Cuturi. Subspace detours: Building transport plans that are optimal on subspace
420 projections. *Advances in Neural Information Processing Systems*, 32, 2019.
- 421 Khai Nguyen and Nhat Ho. Amortized projection optimization for sliced wasserstein generative models. *arXiv
422 preprint arXiv:2203.13417*, 2022.
- 423 Michaël Perrot, Nicolas Courty, Rémi Flamary, and Amaury Habrard. Mapping estimation for discrete optimal
424 transport. *Advances in Neural Information Processing Systems*, 29, 2016.
- 425 Gabriel Peyré, Marco Cuturi, et al. Computational optimal transport: With applications to data science. *Founda-
426 tions and Trends® in Machine Learning*, 11(5-6):355–607, 2019.
- 427 Ievgen Redko, Nicolas Courty, Rémi Flamary, and Devis Tuia. Optimal transport for multi-source domain
428 adaptation under target shift. In *The 22nd International Conference on Artificial Intelligence and Statistics*,
429 pages 849–858. PMLR, 2019.
- 430 Ralph Tyrell Rockafellar. Convex analysis. In *Convex analysis*. Princeton university press, 2015.
- 431 Litu Rout, Alexander Korotin, and Evgeny Burnaev. Generative modeling with optimal transport maps. In
432 *International Conference on Learning Representations*, 2021.
- 433 Andrei A Rusu, Dushyant Rao, Jakub Sygnowski, Oriol Vinyals, Razvan Pascanu, Simon Osindero, and Raia
434 Hadsell. Meta-learning with latent embedding optimization. In *International Conference on Learning Rep-
435 resentations*, 2018.

- 436 Filippo Santambrogio. Optimal transport for applied mathematicians. *Birkäuser, NY*, 55(58-63):94, 2015.
- 437 Meyer Scetbon, Marco Cuturi, and Gabriel Peyré. Low-rank sinkhorn factorization. In *International Conference on Machine Learning*, pages 9344–9354. PMLR, 2021.
- 439 Meyer Scetbon, Gabriel Peyré, and Marco Cuturi. Linear-time gromov wasserstein distances using low rank couplings and costs. In *International Conference on Machine Learning*, pages 19347–19365. PMLR, 2022.
- 441 Geoffrey Schiebinger, Jian Shu, Marcin Tabaka, Brian Cleary, Vidya Subramanian, Aryeh Solomon, Joshua Gould, Siyan Liu, Stacie Lin, Peter Berube, Lia Lee, Jenny Chen, Justin Brumbaugh, Philippe Rigollet, Konrad Hochedlinger, Rudolf Jaenisch, Aviv Regev, and Eric S. Lander. Optimal-transport analysis of single-cell gene expression identifies developmental trajectories in reprogramming. *Cell*, 176(4):928–943.e22, 2019. ISSN 0092-8674. doi: <https://doi.org/10.1016/j.cell.2019.01.006>. URL <https://www.sciencedirect.com/science/article/pii/S009286741930039X>.
- 447 Vivien Seguy, Bharath Bhushan Damodaran, Remi Flamary, Nicolas Courty, Antoine Rolet, and Mathieu Blondel. Large scale optimal transport and mapping estimation. In *International Conference on Learning Representations*, 2018.
- 450 Justin Solomon, Fernando De Goes, Gabriel Peyré, Marco Cuturi, Adrian Butscher, Andy Nguyen, Tao Du, and Leonidas Guibas. Convolutional wasserstein distances: Efficient optimal transportation on geometric domains. *ACM Transactions on Graphics (TOG)*, 34(4):1–11, 2015.
- 453 Max Sommerfeld, Jörn Schrieber, Yoav Zemel, and Axel Munk. Optimal transport: Fast probabilistic approximation with exact solvers. *J. Mach. Learn. Res.*, 20:105–1, 2019.
- 455 Kenneth O Stanley, David B D’Ambrosio, and Jason Gauci. A hypercube-based encoding for evolving large-scale neural networks. *Artificial life*, 15(2):185–212, 2009.
- 457 Amirhossein Taghvaei and Amin Jalali. 2-wasserstein approximation via restricted convex potentials with application to improved training for gans. *arXiv preprint arXiv:1902.07197*, 2019.
- 459 Matthew Tancik, Ben Mildenhall, Terrance Wang, Divi Schmidt, Pratul P Srinivasan, Jonathan T Barron, and Ren Ng. Learned initializations for optimizing coordinate-based neural representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2846–2855, 2021.
- 462 Alexis Thibault, Lenaïc Chizat, Charles Dossal, and Nicolas Papadakis. Overrelaxed sinkhorn-knopp algorithm for regularized optimal transport. *arXiv preprint arXiv:1711.01851*, 2017.
- 464 James Thornton and Marco Cuturi. Rethinking initialization of the sinkhorn algorithm. *arXiv preprint arXiv:2206.07630*, 2022.
- 466 Cédric Villani. *Optimal transport: old and new*, volume 338. Springer, 2009.
- 467 Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017.
- 470 Luisa Zintgraf, Kyriacos Shiarli, Vitaly Kurin, Katja Hofmann, and Shimon Whiteson. Fast context adaptation via meta-learning. In *International Conference on Machine Learning*, pages 7693–7702. PMLR, 2019.

472 **Checklist**

- 473 1. For all authors...
- 474 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
475 contributions and scope? [Yes] We hope so
- 476 (b) Did you describe the limitations of your work? [Yes] In sect. 6
- 477 (c) Did you discuss any potential negative societal impacts of your work? [No] We do not
478 immediately foresee any that our work would add that the broader optimal transport
479 field doesn't already have
- 480 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
481 them? [Yes]
- 482
- 483 2. If you are including theoretical results... (This is not a theory paper)
- 484 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 485 (b) Did you include complete proofs of all theoretical results? [N/A]
- 486 3. If you ran experiments...
- 487 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
488 mental results (either in the supplemental material or as a URL)? [Yes]
- 489 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
490 were chosen)? [Yes]
- 491 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
492 ments multiple times)? [Yes] We show results from multiple trials in most places
- 493 (d) Did you include the total amount of compute and the type of resources used (e.g., type
494 of GPUs, internal cluster, or cloud provider)? [Yes]
- 495 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 496 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 497 (b) Did you mention the license of the assets? [Yes]
- 498 (c) Did you include new assets either in the supplemental material or as a URL? [No]
- 499 (d) Did you discuss whether and how consent was obtained from people whose data
500 you're using/curating? [N/A]
- 501 (e) Did you discuss whether the data you are using/curating contains personally identifi-
502 able information or offensive content? [N/A]
- 503
- 504 5. If you used crowdsourcing or conducted research with human subjects...
- 505 (a) Did you include the full text of instructions given to participants and screenshots, if
506 applicable? [N/A]
- 507 (b) Did you describe any potential participant risks, with links to Institutional Review
508 Board (IRB) approvals, if applicable? [N/A]
- 509 (c) Did you include the estimated hourly wage paid to participants and the total amount
510 spent on participant compensation? [N/A]

511 **A Selecting ϵ for MNIST**

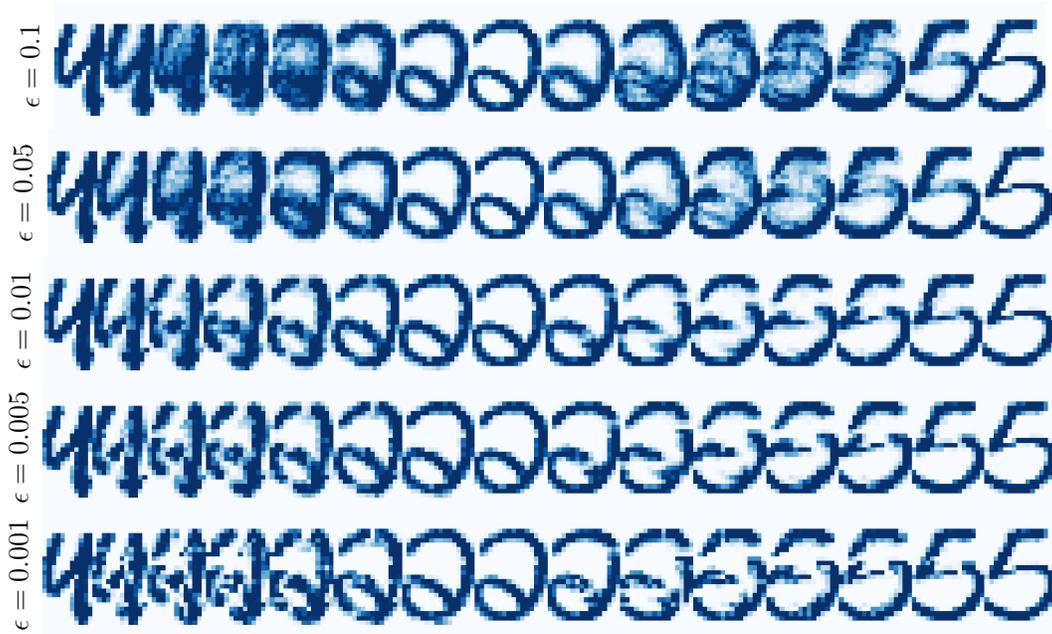


Figure 8: We selected $\epsilon = 10^{-2}$ for our MNIST coupling experiments as it results in transport maps that are not too blurry or sharp.

512 **B Other models for continuous OT**

513 While developing the hyper-network or Meta ICNN in [sect. 3.2](#) for predicting couplings between
 514 continuous measures, we considered alternative modeling formulations briefly documented in this
 515 section. We finalized only the hyper-network model because it is conceptually the most similar to
 516 predicting the optimal dual variables in the continuous setting and results in rapid predictions.

517 **B.1 Optimization-based meta-learning (MAML-inspired)**

518 The model-agnostic meta-learning setup proposed in MAML [[Finn et al., 2017](#)] could also be ap-
 519 plied in the Meta OT setting to learn an adaptable initial parameterization. In the continuous setting,
 520 one initial version would take a parameterized dual potential model $\psi_\varphi(x)$ and seek to learn an initial
 521 parameterization φ_0 so that optimizing a loss such as the W2GN loss \mathcal{L} from [eq. \(12\)](#) results in
 522 a minimal $\mathcal{L}(\varphi_K)$ after adapting the model for K steps. Formally, this would optimize:

$$\arg \min_{\varphi_0} \mathcal{L}(\varphi_K) \quad \text{where} \quad \varphi_{t+1} = \varphi_t - \nabla_{\varphi} \mathcal{L}(\varphi_t) \tag{19}$$

523 [Tancik et al. \[2021\]](#) explores similar learned initializations for coordinate-based neural implicit rep-
 524 resentations for 2D images, CT scan reconstruction, and 3d shape and scene recovery from 2D
 525 observations.

526 **Challenges for Meta OT.** The transport maps given by $T = \nabla \psi$ can significantly vary depending on
 527 the input measures α, β . We found it difficult to learn an initialization that can be rapidly adapted,
 528 and optimizing [eq. \(19\)](#) is more computationally expensive than [eq. \(18\)](#) as it requires unrolling
 529 through many evaluations of the transport loss \mathcal{L} . And, we found that *only* learning to predict
 530 the optimal parameters with [eq. \(18\)](#), conditional on the input measures, and then fine-tuning with
 531 W2GN to be stable.

532 **Advantages for Meta OT.** Exploring MAML-inspired methods could further incorporate the knowl-
 533 edge that the model’s prediction is going to be fine-tuned into the learning process. One promising

534 direction we did not try could be to integrate some of the ideas from LEO [Rusu et al., 2018] and
535 CAVIA [Zintgraf et al., 2019], which propose to learn a latent space for the parameters where the
536 initialization is also conditional on the input.

537 B.2 Neural process and conditional Monge maps

538 The (conditional) neural process models considered in Garnelo et al. [2018b,a] can also be adapted
539 for the Meta OT setting, and is similar to the model proposed in Bunne et al. [2022a]. In the
540 continuous setting, this would result in a dual potential that is also conditioned on a representation
541 of the input measures, e.g. $\psi_\varphi(x; z)$ where $z := f_\varphi^{\text{emb}}(\alpha, \beta)$ is a learned embedding of the input
542 measures that is learned with the parameters of ψ . This could be formulated as

$$\arg \min_{\varphi} \mathbb{E}_{(\alpha, \beta) \sim \mathcal{D}} \mathcal{L}(\varphi, f_\varphi^{\text{emb}}(\alpha, \beta)), \quad (20)$$

543 where \mathcal{L} modifies the model used in the loss eq. (12) to also be conditioned on the context extracted
544 from the measures.

545 **Challenges for Meta OT.** This raises the issue on best-formulating the model to be conditional on
546 the context. One way could be to append z to the input point x in the domain. Bunne et al. [2022a]
547 proposes to use the Partially Input-Convex Neural Network (PICNN) from [Amos et al., 2017] to
548 make the model convex with respect to x and not z .

549 **Advantages for Meta OT.** A large advantage is that the representation z of the measures α, β would
550 be significantly lower-dimensional than the parameters φ that our Meta OT models are predicting.

551 C Additional experimental and implementation details

552 We have attached the Jax source code necessary to run and reproduce all of the experiments in our
553 paper and will open-source all of it. Here is a basic overview of the files:

```
554 | meta_ot  Meta OT Python library code
    |   | conjugate.py  Exact conjugate solver for the continuous setting
    |   | data.py
    |   | models.py
    |   | utils.py
    | config  Hydra configuration for the experiments (containing hyper-parameters)
    | train_discrete.py  Train Meta OT models for discrete OT
    | train_color_single.py  Train a single ICNN with W2GN between 2 images (for debugging)
    | train_color_meta.py  Train a Meta ICNN with W2GN
    | plot_mnist.py  Visualize the MNIST couplings
    | plot_world_pair.py  Visualize the spherical couplings
    | eval_color.py  Evaluate the Meta ICNN in the continuous setting
    | eval_discrete.py  Evaluate the Meta ICNN for the discrete tasks
```

555 Connecting to the data is one difficulty in running the experiments. The easiest experiment to re-run
 556 is the MNIST one, which will automatically download the dataset:

```
557 ./train_discrete.py # Train the model, outputting to <exp_dir>  

  558 ./eval_discrete.py <exp_dir> # Evaluate the learned models  

  559 ./plot_mnist.py <exp_dir> # Produce further visualizations
```

562 **C.1 Hyper-parameters**

563 We briefly summarize the hyper-parameters we used for training, which we did not extensively tune.
 564 In the discrete setting, we use the same hyper-parameters for the MNIST and spherical settings.

Table 3: Discrete OT hyper-parameters.

Name	Value
Batch size	128
Number of training iterations	50000
MLP Hidden Sizes	[1024, 1024, 1024]
Adam learning rate	1e-3

Table 4: Continuous OT hyper-parameters.

Name	Value
Meta batch size (for α, β)	8
Inner batch size (to estimate \mathcal{L})	1024
Cycle loss weight (γ)	3.
Adam learning rate	1e-3
ℓ_2 weight penalty	1e-6
Max grad norm (for clipping)	1.
Number of training iterations	200000
Meta ICNN Encoder	ResNet18
Encoder output size (both measures)	256 \times 2
Meta ICNN Decoder Hidden Sizes	[512]

566 **C.2 Sinkhorn convergence times, varying thresholds**

567 In the main paper, [table 1](#) reports the runtime of Sinkhorn to reach a convergence threshold of the
 568 marginal error being below a tolerance of 10^{-3} , which is the default value used in many solvers.
 569 [app. C.2](#) report the results from sweeping over other thresholds and show that Meta OT’s initializa-
 570 tion is consistently able to help.

Table 5: Sinkhorn runtime to reach a thresholded marginal error on MNIST.

Initialization	Threshold= 10^{-2}	Threshold= 10^{-3}	Threshold= 10^{-4}	Threshold= 10^{-5}
Zeros	$4.5 \cdot 10^{-3} \pm 1.5 \cdot 10^{-3}$	$7.7 \cdot 10^{-3} \pm 1.2 \cdot 10^{-3}$	$1.1 \cdot 10^{-2} \pm 1.8 \cdot 10^{-3}$	$1.5 \cdot 10^{-2} \pm 2.3 \cdot 10^{-3}$
Gaussian	$4.1 \cdot 10^{-3} \pm 1.2 \cdot 10^{-3}$	$7.7 \cdot 10^{-3} \pm 1.4 \cdot 10^{-3}$	$1.1 \cdot 10^{-2} \pm 1.7 \cdot 10^{-3}$	$1.4 \cdot 10^{-2} \pm 2.4 \cdot 10^{-3}$
Meta OT	$2.3 \cdot 10^{-3} \pm 9.2 \cdot 10^{-6}$	$3.9 \cdot 10^{-3} \pm 1.6 \cdot 10^{-3}$	$6.7 \cdot 10^{-3} \pm 1.4 \cdot 10^{-3}$	$1.0 \cdot 10^{-2} \pm 2.4 \cdot 10^{-3}$

Table 6: Sinkhorn runtime to reach a thresholded marginal error on the spherical transport problem.

Initialization	Threshold= 10^{-2}	Threshold= 10^{-3}	Threshold= 10^{-4}	Threshold= 10^{-5}
Zeros	$8.8 \cdot 10^{-1} \pm 1.3 \cdot 10^{-1}$	$1.4 \pm 1.9 \cdot 10^{-1}$	$2.1 \pm 3.6 \cdot 10^{-1}$	$2.8 \pm 5.6 \cdot 10^{-1}$
Gaussian	$5.6 \cdot 10^{-1} \pm 9.9 \cdot 10^{-2}$	$1.1 \pm 2.0 \cdot 10^{-1}$	$1.7 \pm 3.5 \cdot 10^{-1}$	$2.4 \pm 5.4 \cdot 10^{-1}$
Meta OT	$7.8 \cdot 10^{-2} \pm 3.4 \cdot 10^{-2}$	$0.44 \pm 1.5 \cdot 10^{-1}$	$0.97 \pm 3.2 \cdot 10^{-1}$	$1.7 \pm 6.8 \cdot 10^{-1}$

571 **C.3 Experimental runtimes and convergence**

572 **App. C.3** shows the convergence during training of Meta OT models in the discrete and continuous
573 settings over 10 trials on our single Quadro GP100 GPU. The MNIST models are consistently trained
574 to optimality within 2 minutes (!) while the continuous model takes a few hours to train.

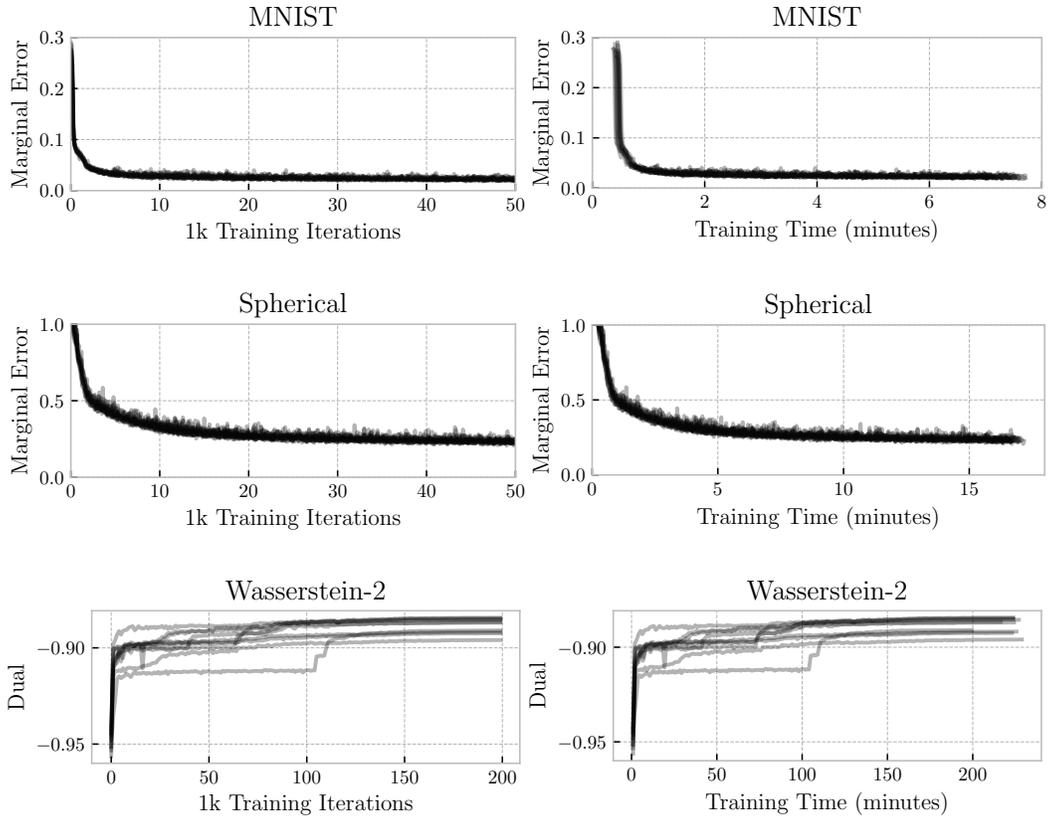


Figure 9: Convergence of Meta OT models during training, reported over iterations and wall-clock time. We run each experiment for 10 trials with different seeds and report each trial as a line.

575 **D Out-of-distribution generalization**

576 **App. D** tests the ability of Meta OT to predict potentials for out-of-distribution input data. We
 577 consider the pairwise training and evaluation on the following datasets: 1) MNIST; 2) USPS [Hull,
 578 1994] (upscaled to have the same size as the MNIST); 3) Google Doodles dataset* with classes Crab,
 579 Cat and Faces; 4) sparsified random uniform data in [0,1] where sparsity (zeroing values below 0.95)
 580 is used to mimic the sparse signal in black-and-white images. For each pair, eg, MNIST-USPS, we
 581 train on one dataset and use the other to predict the potentials. The comparison is done using the
 582 same metric as before, i.e., the deviation from the marginal constraints defined in eq. (7).

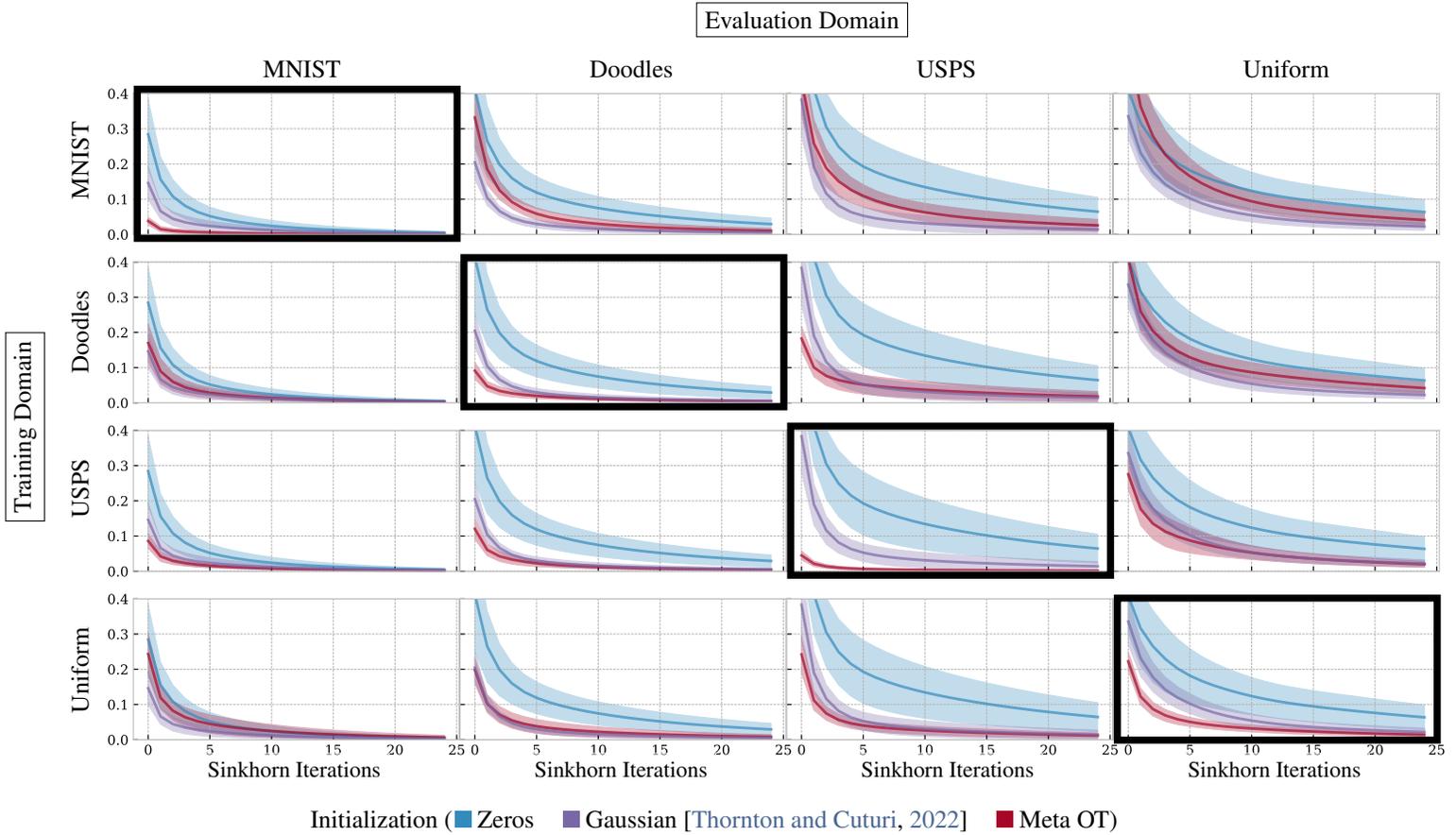


Figure 10: Cross-domain experiments.

*<https://quickdraw.withgoogle.com/data>

583 **E Additional color transfer results**

584 We next show additional color transfer results from the experiments in [sect. 4.3](#) on the following
585 public domain images from [WikiArt](#):

- 586 • Distant View of the Pyramids by Winston Churchill (1921)
- 587 • Charing Cross Bridge, Overcast Weather by Claude Monet (1900)
- 588 • Houses of Parliament by Claude Monet (1904)
- 589 • October Sundown, Newport by Childe Hassam (1901)
- 590 • Landscape with House at Ceret by Juan Gris (1913)
- 591 • Irises in Monet's Garden by Claude Monet (1900)
- 592 • Crystal Gradation by Paul Klee (1921)
- 593 • Senecio by Paul Klee (1922)
- 594 • Váza s květinami by Josef Capek (1914)
- 595 • Sower with Setting Sun by Vincent van Gogh (1888)
- 596 • Three Trees in Grey Weather by Claude Monet (1891)
- 597 • Vase with Daisies and Anemones by Vincent van Gogh (1887)

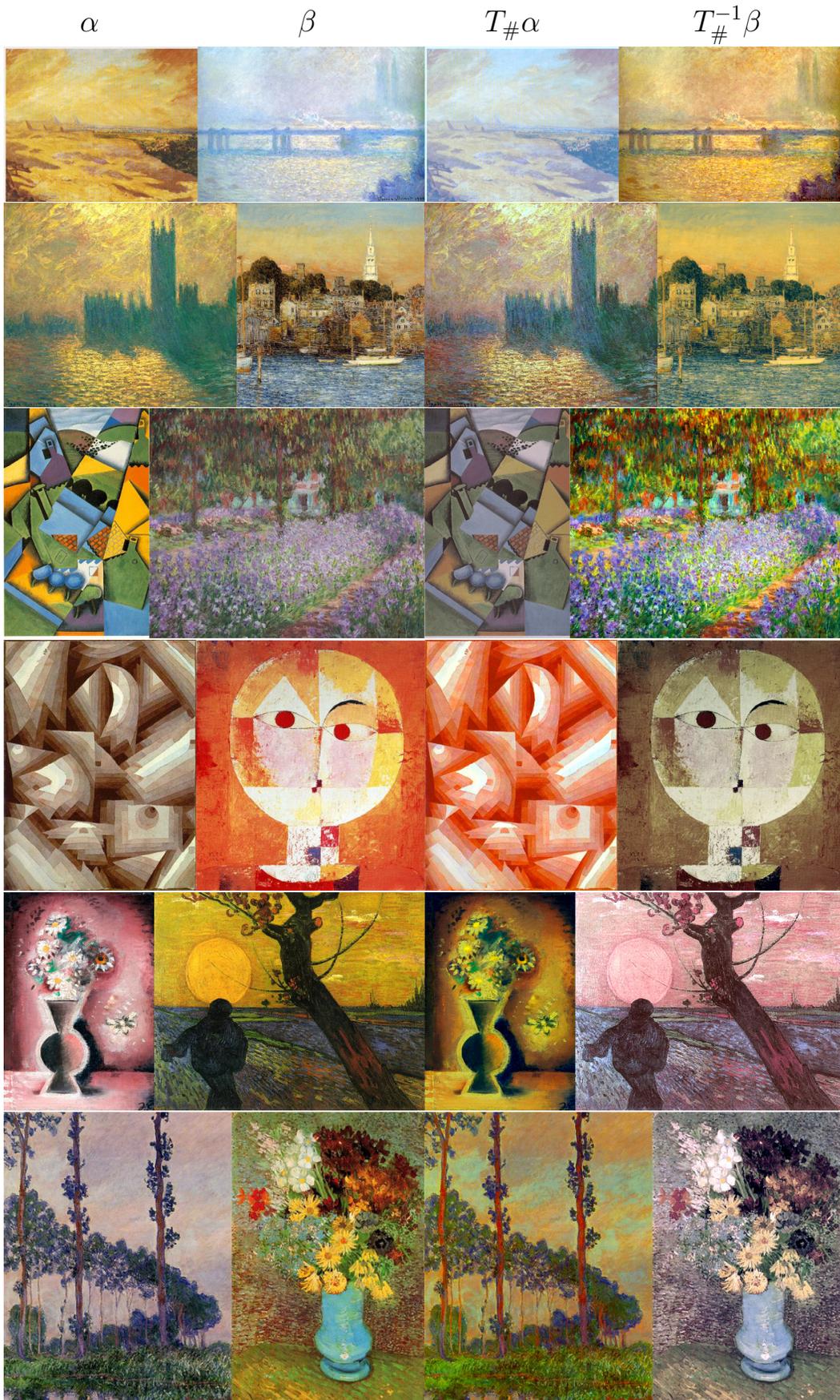


Figure 11: Meta ICNN (initial prediction). The sources are given in the beginning of app. E.



Figure 12: Meta ICNN + W2GN fine-tuning. The sources are given in the beginning of app. E.



Figure 13: W2GN (final). The sources are given in the beginning of app. E.