# A BICONVEX FORMULATION FOR TRANSPORT OF MIX-TURE MODELS

Anonymous authors
Paper under double-blind review

# **ABSTRACT**

Optimal transport (OT) provides a principled framework for mapping between probability distributions. Despite extensive progress in the field, OT remains computationally demanding, and the resulting transport plans are often difficult to interpret. Here, we propose Optimal Mixture Transport (OMT), an efficient algorithm that leverages mixture modeling and entropic regularization to yield interpretable transport plans. We show that transport between mixtures, in particular mixtures of Gaussians which are universal approximators in  $L^2$ , can be formulated as a biconvex optimization problem with a unique minimizer. This formulation not only reduces computational cost, but also provides component-level correspondences, offering insights into complex distributions. We demonstrate the practicality and effectiveness of OMT across a diverse collection of synthetic benchmarks and real-world datasets, including large-scale single-cell RNA sequencing measurements.

## 1 Introduction

Optimal Transport (OT) offers a powerful mathematical framework for comparing probability distributions and finding optimal mappings between them (Santambrogio, 2015). Its versatility has led to advances in diverse fields, including domain adaptation (Grave et al., 2019; Struckmeier et al., 2023; Chuang et al., 2023; Fernandes Montesuma et al., 2025), data integration and alignment (Demetci et al., 2022), and predicting cell fates (Tong et al., 2020; Bunne et al., 2023; 2024). At its core, OT seeks to find the most cost-effective way to transform one probability distribution into another, subject to constraints on the total mass being transported (Peyré et al., 2019; Villani et al., 2008).

A major challenge in OT has been its high computational cost. Cuturi (2013) introduced an entropy regularization term to the OT objective to obtain a strictly convex problem (EOT) and an elegant solution known as the Sinkhorn algorithm. However, even with EOT, sample-to-sample transportation remains limited by the curse of dimensionality and can be slow on large datasets (Genevay et al., 2018). To mitigate this, mini-batch strategies have been developed to approximate the transport plan by operating on subsets of the data (Fatras et al., 2021b;a). While computationally cheaper, these methods often yield suboptimal transport plans, as cost estimation from subsets can be inaccurate and satisfying the mass preservation constraint of balanced OT becomes difficult.

To improve transport accuracy over batches, one prominent class of methods approximates the OT path by non-parametric interpolation within the Wasserstein space. These techniques range from deterministic approaches, such as Progressive Optimal Transport (ProgOT) (Kassraie et al., 2024), to stochastic methods based on gradient flows and Schrödinger bridges (Albergo & Vanden-Eijnden, 2023; Albergo et al., 2024). Stochastic methods, often employing neural networks such as those based on gradient flows (Daniels et al., 2021) or Schrödinger bridges (Gushchin et al., 2023b;a), also typically necessitate inner iterations to achieve accurate transport maps. Such methods construct a sequence of intermediate distributions to bridge the source and target, often requiring numerous intermediate steps, significant memory overhead, and many inner iterations to converge to an accurate solution. Furthermore, simpler displacement strategies, like the McCann interpolation (McCann, 1997) used in ProgOT, do not always produce interpretable intermediate distributions. While regularization techniques, can improve the robustness of the transport map approximation (Buzun et al., 2024b), their performance is sensitive to the regularization parameters.

In contrast to non-parametric and relaxation-based approaches, an effective strategy for large-scale problems is to adopt a parametric model of the data, thereby simplifying the task. Following this direction, we propose Optimal Mixture Transport (OMT), an efficient and scalable framework for computing EOT between mixture models. While EOT is generally intractable for most parametric families, a closed-form solution exists for transport between two Gaussian distributions in both balanced and unbalanced settings (Janati et al., 2020). Building on this result, we tailor the framework to the Gaussian family, focusing on Gaussian Mixture Models (GMMs). This specialization is powerful as GMMs are universal function approximators capable of representing any sufficiently smooth density with arbitrary precision (Goodfellow et al., 2016). Our formulation recasts the transport problem as a uniquely solvable biconvex optimization, yielding a computationally efficient and theoretically grounded alternative for large-scale transport tasks. Furthermore, we observe that one limitation of many learned coupling transport maps or dual functions is their pronounced directional bias (source → target), which leads to performance degradation when inverted. In contrast, we show that OMT maps remain robust regardless of transport direction.

Our main contributions are summarized as follows:

- We propose a parametric EOT framework, called OMT, which operates by transporting sub-populations rather than individual samples.
- We show that the OMT formulation is strictly biconvex and, when solved as part of a global
  optimization algorithm, this subproblem converges to a unique solution in a single step.
- Building on closed-form results for entropic Gaussian transport, we propose a formulation of OMT within Gaussian distributions, as a flexible and expressive parametric family.
- Through experiments on synthetic and real-world datasets, we demonstrate that OMT consistently matches or surpasses the performance of state-of-the-art OT solvers, while requiring substantially less computation and memory.

# 2 BACKGROUND

**Optimal transport:** For  $\mathcal{X}, \mathcal{Y} \subset \mathbb{R}^d$ , probability measures  $\mu_0, \mu_1 \in \mathcal{P}(\mathbb{R}^d)$  and  $c : \mathcal{X} \times \mathcal{Y} \to \mathbb{R}$ , a cost function associated with transporting a unit of mass from a point in  $\mathcal{X}$  to a point in  $\mathcal{Y}$ , the minimum total cost of transport can be obtained as

$$\inf_{T \sharp \mu_0 = \mu_1} \int_{\mathcal{X}} c(\mathbf{x}, T(\mathbf{x})) \, d\mu_0(\mathbf{x}), \tag{1}$$

where  $T_{\sharp}\mu_0=\mu_1$  denotes the pushforward of  $\mu_0$  by T, defined as  $\mu_1(A)=\mu_0(T^{-1}(A)), \forall A\subset\mathcal{Y}$ , ensuring mass conservation. Problem (1) is known as the Monge problem (Peyré et al., 2019) which seeks a map  $T:\mathbb{R}^d\to\mathbb{R}^d$  referred to as the transport map between  $\mu_0$  and  $\mu_1$ . The Monge formulation is often problematic because the optimization is over a non-convex set of maps, and a deterministic map T may not exist. To bypass this, Kantorovich presented a relaxed formulation, which seeks a distribution  $\pi\in\mathbb{R}^d\times\mathbb{R}^d$  referred to as "coupling" of  $\mu_0$  and  $\mu_1$ , as

$$\inf_{\pi \in \prod(\mu_0, \mu_1)} \int_{\mathcal{X} \times \mathcal{Y}} c(\mathbf{x}, \mathbf{y}) \, d\pi(\mathbf{x}, \mathbf{y}), \tag{2}$$

where  $\prod (\mu_0, \mu_1) := \{ \pi \in \mathcal{P}(\mathcal{X} \times \mathcal{Y}) \mid P_{\mathcal{X}} : \pi \to \mu_0, P_{\mathcal{Y}} : \pi \to \mu_1 \}.$ 

When  $\mathcal{X} = \mathcal{Y}$  and  $c(\mathbf{x}, \mathbf{y}) = d(\mathbf{x}, \mathbf{y})^p$ ,  $p \ge 1$ , where d is a distance on  $\mathcal{X}$ , the Kantorovich problem is equivalent to the Wasserstein p-distance between probability measures,  $\mathcal{W}_p(\mu_0, \mu_1)$ . Specifically, for  $c(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2^2$ , the Kantorovich problem yields the squared Wasserstein-2 distance:

$$\mathcal{W}_2^2(\mu_0, \mu_1) = \inf_{\pi \in \Pi(\mu_0, \mu_1)} \int_{\mathcal{X} \times \mathcal{Y}} \|\mathbf{x} - \mathbf{y}\|_2^2 d\pi(\mathbf{x}, \mathbf{y}). \tag{3}$$

According to the Brenier Theorem (Santambrogio, 2015), in (3), if  $\mu_0$  and  $\mu_1$  are absolutely continuous with respect to the Lebesgue measure, there exists a unique optimal solution that can be expressed as  $\pi^* = (I_d \times T^*) \sharp \mu_0$ , where  $I_d$  stands for the identity map, and  $T^*$  is the unique minimizer of (1). Moreover, the unique optimal transport  $T(\mathbf{x}) = \nabla \phi(\mathbf{x})$ , where  $\phi$  is a convex function.

**Entropic optimal transport:** Entropic optimal transport introduces an entropic regularization term to (3), transforming the problem into a strictly convex optimization that can be efficiently solved using algorithms like the Sinkhorn-Knopp method (Cuturi, 2013; Janati et al., 2020). For a regularization parameter  $\varepsilon > 0$ , the entropic optimal transport cost is defined as:

$$d_{\varepsilon}(\mu_0, \mu_1) = \inf_{\pi \in \Pi(\mu_0, \mu_1)} \left\{ \int_{\mathcal{X} \times \mathcal{V}} \|\mathbf{x} - \mathbf{y}\|_2^2 d\pi(\mathbf{x}, \mathbf{y}) - 2\varepsilon H(\pi) \right\}. \tag{4}$$

Considering  $D_{KL}(P||Q) = \int dP \left( \log \frac{dP}{dQ} - 1 \right) + dQ$ , minimizing the objective in (4) is equivalent to minimizing

$$\min_{\pi \in \prod(\mu_0, \mu_1)} \int_{\mathcal{X} \times \mathcal{Y}} \|\mathbf{x} - \mathbf{y}\|_2^2 d\pi(\mathbf{x}, \mathbf{y}) + 2\varepsilon D_{KL}(\pi \| \mu_0 \otimes \mu_1).$$

**Transport between Gaussian measures:** When both mass measures are Gaussian distributions, i.e.,  $\mu_0 = \mathcal{N}(\mathbf{m}_0, \Sigma_0)$  and  $\mu_1 = \mathcal{N}(\mathbf{m}_1, \Sigma_1)$ , (3) is simplified as

$$W_2^2(\mathcal{N}(\mathbf{m}_0, \Sigma_0), \mathcal{N}(\mathbf{m}_1, \Sigma_1)) = \|\mathbf{m}_0 - \mathbf{m}_1\|_2^2 + tr\{\Sigma_0 + \Sigma_1 - 2\Gamma\},\tag{5}$$

where  $\Gamma = \left(\Sigma_0^{\frac{1}{2}}\Sigma_1\Sigma_0^{\frac{1}{2}}\right)^{\frac{1}{2}}$  (Bhatia et al., 2019; Janati et al., 2020). Furthermore, the optimal transport map  $T^*: \mathbb{R}^d \to \mathbb{R}^d$  that pushes  $\mu_0$  forward to  $\mu_1$  admits a closed-form as follows.

$$T^*(\mathbf{x}) = A(\mathbf{x} - \mathbf{m}_0) + \mathbf{m}_1,\tag{6}$$

where  $A = \Sigma_0^{-1} \# \Sigma_1 = \Sigma_0^{-\frac{1}{2}} \left( \Sigma_0^{\frac{1}{2}} \Sigma_1 \Sigma_0^{\frac{1}{2}} \right)^{\frac{1}{2}} \Sigma_0^{-\frac{1}{2}} = \Sigma_0^{-\frac{1}{2}} \Gamma \Sigma_0^{-\frac{1}{2}}$ , corresponding to the geometric mean of the precision and covariance matrices at source and target points, respectively.

Gaussian mixture models: Gaussian mixture models (GMMs) are widely used for density estimation due to several key advantages: (i) As a linear combination of Gaussian distributions, GMMs allow for analytical tractability and have favorable asymptotic properties. (ii) GMMs are universal approximators for continuous density functions: any smooth density can be approximated with arbitrary accuracy by a mixture of Gaussians with enough number of components (Titterington et al., 1985; Scott, 2015; Zeevi & Meir, 1997). (iii) Many real-world datasets are naturally organized into clusters with unimodal distributions, making GMMs particularly effective for modeling such structures. These motivate our exploration of the optimal transport problem for Gaussian mixtures.

Using GMMs for density approximation involves approximating density functions by a convex combination of "basis" densities (Zeevi & Meir, 1997). Consider the set of square-integrable density functions in  $\mathbb{R}^d$ , denoted as  $\mathcal{F}=\{f\mid f\in L^2(\mathbb{R}^d),\ f\geq 0,\ \int_{\mathbb{R}^d}f(\mathbf{x})d\mathbf{x}=1\}$ . We define the set of GMM densities with K components,  $\mathcal{G}_K$ , as:

$$\mathcal{G}_K = \left\{ f_K^{\boldsymbol{\theta}} \mid f_K^{\boldsymbol{\theta}}(\cdot) = \sum_{i=1}^K \alpha_i \phi(\cdot, \mu_i, \Sigma_i), \ \alpha_i > 0, \ \sum_{i=1}^K \alpha_i = 1 \right\},\tag{7}$$

where  $\boldsymbol{\theta} = \{\alpha_i, \boldsymbol{\mu}_i, \Sigma_i\}_{i=1}^K$  represents the collection of parameters for the K components.  $\mathcal{G}_K \subset \mathcal{F}$  and the universal approximation property implies that for any  $f \in \mathcal{F}$ ,  $\lim_{K \to \infty} \inf_{\theta} \mathcal{D}(f, f_K^{\boldsymbol{\theta}}) = 0$ , where  $\mathcal{D}$  denotes a distance (Titterington et al., 1985; Zeevi & Meir, 1997).

# 3 RELATED WORK

To address the limitations of Sinkhorn-based methods, researchers turned to deep learning, giving rise to Neural Optimal Transport (Neural OT) (Makkuva et al., 2020; Korotin et al., 2023), which uses neural networks to learn a continuous mapping between distributions, while enforcing theoretical constraints (Genevay et al., 2018; Buzun et al., 2024b). Another direction directly learns the transport map via neural networks, transforming samples from a source to a target distribution. This approach is widely used in domain adaptation and generative modeling, where models such as normalizing flows learn invertible maps from simple to complex distributions. A prominent example of Neural OT connects diffusion models with the theory of Schrödinger Bridges, a classic stochastic

transport problem. This establishes a learning framework for diffusion models equivalent to solving a Schrödinger Bridge problem, which can be viewed as a form of neural EOT (Gushchin et al., 2024).

Alongside neural methods, other strategies tackle scalability through iterative, mini-batch frameworks. PROGOT (Kassraie et al., 2024), constructs the transport map sequentially. While this approach can be parallelized and accelerated using modern frameworks like OTT-Jax (Cuturi et al., 2022), it requires substantial memory and computational resources. Similarly, stochastic and neural OT methods require extensive training with many samples, making the process time-consuming and computationally expensive.

Parametric OT simplifies the transport problem by assuming data distributions belong to a parametric family, which often yields computationally more efficient solutions. A prominent example involves Gaussian distributions, for which both  $W_2(\mu_0, \mu_1)$  and its entropically regularized version (Eq. 4) admit closed-form solutions Kassraie et al. (2024). Building on this, the parametric formulation has been extended to the more general case of Gaussian Mixture Models (GMMs). This body of work approximates (bounds) the Wasserstein distance between Gaussian components, proposed as the aggregated Wasserstein distance (Chen et al., 2019) or  $GW_2$  (Delon & Desolneux, 2020), by considering the transport between their individual components. This approach reduces the computational complexity from being dependent on the number of data points to the number of mixture components, offering a scalable solution for high-density data. However, existing studies have often been limited to simpler applications, such as simple 2D tasks and color transfer. A recent extension leverages  $GW_2$  for unsupervised domain adaptation, facilitating label transfer from a source domain to a target domain (Fernandes Montesuma et al., 2025).

A key challenge in using  $GW_2$  lies in optimizing the component weights, which reduces the task to a discrete OT problem, a computationally challenging paradigm that may lack a unique solution. In this work, we extend parametric OT to the entropic mixture transport setting. We show that this formulation is strictly biconvex, yielding a strictly biconvex formulation that guarantees uniqueness for both the transport plan over mixing weights and the individual component distributions. Moreover, building on the findings of Kassraie et al. (2024), we specialize the proposed formulation to overparameterized GMM, a regime in which a global convergence of Expectation-Maximization can be established (Xu et al., 2024), supporting practical viability across different tasks.

## 4 Transport problem for mixture models

Let  $\nu \in M_K(\mathbb{R}^d)$  denote a mixture model in  $\mathbb{R}^d$  with K components:

$$\nu = \sum_{i=0}^{K} \alpha_i \mu_i,\tag{8}$$

where  $\mu_i$  are probability measures and  $\sum_i \alpha_i = 1$ ,  $\alpha_i \geq 0$ ,  $\forall i$ .

**Definition 1** (Mixture transport coupling). Given two measures  $\nu_0 \in M_{K_0}(\mathbb{R}^d)$  and  $\nu_1 \in M_{K_1}(\mathbb{R}^d)$ , we define the mixture transport coupling as follows:

$$\pi_{\mathcal{M}}^* := \underset{\pi \in \Pi(\nu_0, \nu_1) \cap M_K(\mathbb{R}^{2d})}{\arg \min} \int_{\mathcal{X} \times \mathcal{Y}} \|\mathbf{x} - \mathbf{y}\|_2^2 d\pi(\mathbf{x}, \mathbf{y}), \tag{9}$$

where  $\prod (\nu_0, \nu_1) := \{ \pi \in \mathcal{P}(\mathcal{X} \times \mathcal{Y}) \mid P_{\mathcal{X}} : \pi \to \nu_0(\mathbf{x}), P_{\mathcal{Y}} : \pi \to \nu_1(\mathbf{y}) \}$  and  $K \leq K_0 K_1$ .

Therefore, the transport policy belongs to the mixture model family and can be expressed as

$$d\pi(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^{K} \omega_i dp_i(\mathbf{x}, \mathbf{y}), \text{ where } \forall i, \ p_i \in \mathcal{P}(\mathbb{R}^d).$$

Note that the trivial choice of  $dp_1 = \ldots = dp_K$ ,  $\omega_1 = \ldots = \omega_K = 1/K$  makes the solver in (9) equal to  $W_2(\nu_0, \nu_1)$  in (3). We next constrain that marginals of the components of the transport

coupling are the same as the components of the source and target functions:

$$\mathcal{D}_{\mathcal{M}}(\nu_{0}, \nu_{1}) = \min_{\Omega, P} \sum_{i,j} \omega_{ij} \int_{\mathcal{X} \times \mathcal{Y}} \|\mathbf{x} - \mathbf{y}\|_{2}^{2} dp_{ij}(\mathbf{x}, \mathbf{y}),$$
s.t. 
$$\mathbf{1}^{\mathbf{T}} \Omega = \boldsymbol{\alpha}_{0}, \quad \Omega^{T} \mathbf{1} = \boldsymbol{\alpha}_{1}$$

$$\forall i, \sum_{i} \int_{\mathcal{Y}} dp_{ij}(\mathbf{x}, \mathbf{y}) = \mu_{0_{i}}, \quad \forall j, \sum_{i} \int_{\mathcal{X}} dp_{ij}(\mathbf{x}, \mathbf{y}) = \mu_{1_{j}},$$
(10)

where  $\Omega = [w_{ij}]$  denotes the matrix of mixture weights. With this constraint,  $\mathcal{D}_{\mathcal{M}}(\nu_0, \nu_1) \geq \mathcal{W}_2(\nu_0, \nu_1)$  and equality is achieved in the  $K \to \infty$  limit when the source and target functions are over-parametrized by a dense mixture model family (e.g., GMM (Goodfellow et al., 2016)). The problem in (10) is similar to minimizing the aggregated Wasserstein distance, which was proposed for comparing hidden Markov models with Gaussian state conditional distributions (Chen et al., 2019).

## 4.1 REGULARIZED MIXTURE TRANSPORT

A common approach to ensure the uniqueness of solutions in optimal transport problems is to introduce an entropy regularization term, which makes the objective function strictly convex and improves the numerical stability of optimization. In the context of mixture transport optimization in (10), we adopt a similar approach by incorporating a weighted average entropy term as a regularizer.

**Definition 2** (Optimal Mixture Transport). We introduce two forms of regularization into the mixture transport problem (9): (i) a component-wise regularizer, and (ii) a mixing-matrix regularizer, controlled respectively by parameters  $\varepsilon_1, \varepsilon_2 > 0$ . The resulting problem is formulated as the following optimization:

$$\mathcal{D}_{OMT} := \min_{\omega_{ij}, dp_{ij}} \sum_{i,j}^{K} \omega_{ij} \left[ \int_{\mathcal{X} \times \mathcal{Y}} \|\mathbf{x} - \mathbf{y}\|_{2}^{2} dp_{ij}(\mathbf{x}, \mathbf{y}) - \varepsilon_{1} H(p_{ij}) \right] - \varepsilon_{2} H(\Omega) .$$

$$for \quad \Omega = [\omega_{ij}]_{K_{0} \times K_{1}} \in \mathcal{S}^{K-1}, \quad P = [p_{ij}]_{K_{0} \times K_{1}}, \quad where \quad p_{ij} \in \prod (\mu_{0_{i}}, \mu_{1_{j}})$$

$$(11)$$

Minimizing the objective in (11) is equivalent to minimizing  $\mathcal{L}_{\varepsilon_1,\varepsilon_2}(\Omega,P)$ , defined as follows.

$$\mathcal{L}_{\varepsilon_{1},\varepsilon_{2}}(\Omega, P) = \sum_{i,j}^{K} \omega_{ij} \left[ \int_{\mathcal{X} \times \mathcal{Y}} \|\mathbf{x} - \mathbf{y}\|_{2}^{2} dp_{ij}(\mathbf{x}, \mathbf{y}) + \varepsilon_{1} D_{KL}(p_{ij} \| \mu_{0_{i}} \otimes \mu_{1_{j}}) \right] + \varepsilon_{2} D_{KL}(\Omega \| \boldsymbol{\alpha}_{0} \otimes \boldsymbol{\alpha}_{1})$$
(12)

**Remark 1.** The problem in Eq. 11 is a generalization of entropic optimal transport in the sense that Eq. 11 collapses to entropic optimal transport when  $K_0 = K_1 = 1$ .

Therefore, we consider an optimization problem of the form

$$\min_{\Omega, P} \mathcal{L}_{\varepsilon_1, \varepsilon_2}(\Omega, P) \tag{13}$$

s.t. 
$$\mathbf{1}\Omega = \boldsymbol{\alpha}_0, \quad \Omega^T \mathbf{1} = \boldsymbol{\alpha}_1, \quad \int_{\mathcal{Y}} dp_{ij}(\mathbf{x}, \mathbf{y}) = \mu_{0_i}, \quad \int_{\mathcal{X}} dp_{ij}(\mathbf{x}, \mathbf{y}) = \mu_{1_j},$$
 (14)

where  $\Omega \in \mathcal{S}^{K-1}$  and  $P \in \mathcal{P}^K(\mathcal{X} \times \mathcal{Y})$ .

Eq. 13 no longer defines a convex program. However, as we show now in Lemma 2, the objective is biconvex. Moreover, while biconvex problems don't have unique solutions generally, Eq. 13 has a unique minimizer that can be obtained efficiently (Theorem 2 and Corollary 2 below).

**Lemma 1.** For any  $\varepsilon_1, \varepsilon_2 > 0$ ,  $\mathcal{L}_{\varepsilon_1, \varepsilon_2}(\Omega, P)$  is strictly biconvex.

Floudas & Visweswaran (1990) proposed the *Global Optimization Algorithm* (GOP) to solve constrained biconvex problems. It decomposes the optimization into disjoint blocks similar to the *Alternate Convex Search*(ACS) method and exploits the convex substructure of the problem by a primal-relaxed dual approach (Gorski et al., 2007). The GOP algorithm is guaranteed to terminate after a finite number of steps for an  $\epsilon$ -global optimum solution, for any  $\epsilon > 0$  (**Theorem 4.11**,

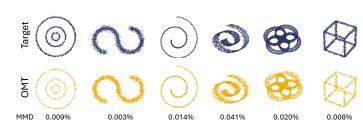


Figure 1: Transporting samples from a normal distribution to various target distributions with entropic OMT. **Top:** target point cloud distributions. **Bottom:** distributions generated by OMT after training with 10,000 samples. MMD between the target and OMT-generated samples expressed as a percentage.

**Corollary 4.12**, Ref. (Gorski et al., 2007)). As mentioned above, the uniqueness of this global optimum is not guaranteed in the general case. However, by exploiting the structure of Eq. 13, we show that its solution is unique and obtained in a single iteration:

**Theorem 1.** For the optimization problem defined in (13), the GOP algorithm converges to a unique solution in a single iteration.

## 4.2 REGULARIZED MIXTURE TRANSPORT FOR GMMS

If the probability measures  $\nu_0 \in G_{K_0}(\mathbb{R}^d)$  and  $\nu_1 \in G_{K_1}(\mathbb{R}^d)$  are defined as mixtures of Gaussian distributions, by using the results from the optimal mixture transport framework introduced in Section 4.1, we can compute an optimal transport plan between the two GMMs. Notably, the resulting optimal mixture transport between two GMMs can also be shown to be a GMM itself, thereby preserving the Gaussian structure in the transported distribution.

**Corollary 1.** Let  $\nu_0 \in G_{K_0}(\mathbb{R}^d)$  and  $\nu_1 \in G_{K_1}(\mathbb{R}^d)$  be two Gaussian mixture models (GMMs) in  $\mathbb{R}^d$  with  $K_0$  and  $K_1$  components, respectively. Then, the optimal mixture transport map between  $\nu_0$  and  $\nu_1$  is itself a Gaussian mixture model with K components, where  $K \leq K_0K_1$ .

# 5 EXPERIMENTS

**Synthetic datasets.** We conduct two sets of simulation experiments. The first set focuses on synthetic 2D tasks with multiple target distributions, designed to demonstrate the capability of the proposed optimal mixture transport strategy. As shown in Figure 1, OMT successfully recovers the target shapes across all cases. In these tasks, the source data is sampled from a normal distribution.

The second set of experiments evaluates our method on the W2-Benchmark tasks (Korotin et al., 2021), which are widely adopted in recent studies on both neural and non-neural OT. We compare the proposed OMT method against state-of-the-art approaches, including ExNOT (Buzun et al., 2024a), ENOT (Gushchin et al., 2024), PROGOT (Kassraie et al., 2024), as well as the classical entropic OT (EOT) solver. Figure 2 presents the comparative performance of OMT across three evaluation metrics: Sinkhorn divergence ( $D_\varepsilon$ ), mean squared error (MSE), and runtime. In all experiments, OMT was trained with  $K_s=3$ ,  $K_t=15$ ,  $\varepsilon_{1,2}=0.01$ . For dimensions d>64, we impose a diagonal structure on the covariance matrix instead of using the full covariance. Appendix B reports the transport cost ( $T_c$ ) and total memory usage for each method. As shown in Figures 2 and 6(Appendix), OMT consistently outperforms EOT, ENOT, and, in most cases, PROGOT. It also outperforms ExNOT at higher latent dimensions. Note that methods like PROGOT and EOT are sample-based solvers, whereas OMT, similar to neural OT, solves the continuous transport problem at the distribution level. Despite this difference, OMT still performs reasonably well on sample-to-sample metrics such as MSE. Considering all metrics along with transport costs, OMT achieves strong overall performance while using substantially less resources, as reflected by shorter runtimes and smaller memory usage.

To investigate the stability of OMT under noise, we conduct an ablation study: noise is added to the source data during training, while the original clean data is used for evaluation. We consider two types of perturbations: white noise, controlled by  $\sigma$  and dropout noise with probability p. Figure 3 demonstrates that OMT produces the most robust OT plans under both noise models, considering the relative change in MSE in response to input perturbations. Overall, OMT consistently delivers strong performance across metrics, often matching or exceeding existing baselines, highlighting both the robustness and competitiveness of our approach.

**Single-cell RNA Sequencing Data.** OT has emerged as a powerful tool in computational biology, with applications such as aligning cell populations across conditions and inferring their trajectories

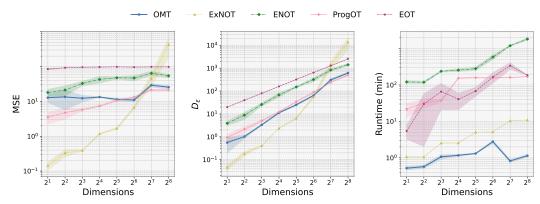


Figure 2: Comparison of OMT against baseline methods on the Wasserstein-2 benchmark tasks in Korotin et al. (2021). The reported plots are averaged over both forward and backward directions. All results are evaluated on the test set with 10,000 samples and averaged across five random initializations. MSE captures fidelity in sample-to-sample transportation whereas  $D_\varepsilon$  is more suitable for transportation between distributions. The runtime is measured on allocated nodes of a cluster, each equipped with one NVIDIA A100 GPU, 4 Intel Xeon Gold 6330N CPU cores, and 128 GB of RAM. The reported time corresponds to the optimization of the transport plan, and the metric calculations are excluded

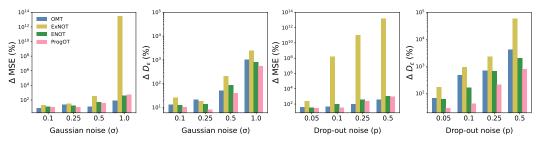


Figure 3: Stability of OT solvers under noise in the W2 benchmark task. The performance of OMT together with baseline methods is evaluated under two types of perturbations applied to the source data. The two left panels illustrate the effect of additive Gaussian noise with increasing standard deviations, while the two right panels show the effect of dropout noise. Reported values indicate performance changes relative to the noise-free case, evaluated on the test set with 10,000 samples and averaged over five random initializations.

over time (Tong et al., 2020; Bunne et al., 2023; 2024). Here, we focus on single-cell RNA sequencing (scRNA-seq) as our primary real-world application. This technology generates high-dimensional molecular profiles by measuring the expression of thousands of genes at single-cell resolution. We consider three scRNA-seq datasets: one human dataset, sci-Plex (Srivatsan et al., 2020), and two 10x Genomics mouse brain datasets, one collected during development (Gao et al., 2024) and the other during aging (Jin et al., 2025). The sci-Plex data serves as a common benchmark for assessing OT performance on real-world biological data (Cuturi et al., 2023; Janati et al., 2020). Consistent with previous work, our analysis focuses on a subset of this dataset comprising three cell lines (A549, K562, and MCF7) exposed to five different cancer treatments for 24 hours. Following the preprocessing steps recommended in Cuturi et al. (2023), the final dataset contains 77, 920 cells and

Table 1: Average  $D_{\varepsilon} \downarrow$  values for the forward and backward OMT mappings compared to PROGOT on the human scRNA-seq dataset (Srivatsan et al., 2020). The results correspond to  $d_{PCA} = 16$  and are reported as the mean over 5 randomly initialized runs, with standard deviations. Additional results for other dimensions and computational costs are provided in Appendix C.

	Belinostat	Dacinostat	Givinostat	Quisinostat	Hesperadin
EOT	$17.4 \pm 0.01$	$18.4 \pm 0.01$	$17.5 \pm 0.01$	$17.6 \pm 0.03$	$17.5 \pm 0.01$
ProgOT	$8.43 \pm 0.01$	$8.84 \pm 0.03$	$8.82 \pm 0.04$	$9.52 \pm 0.01$	$8.06 \pm 0.01$
OMT	$7.91 \pm 0.06$	$8.12 \pm 0.04$	$\boldsymbol{8.75 \pm 0.14}$	$\boldsymbol{8.72 \pm 0.10}$	$8.00 \pm 0.30$

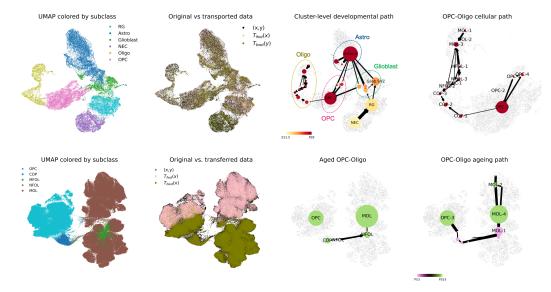


Figure 4: OPC—Oligo trajectories across the mouse lifespan. Top row: developmental dataset from the mouse visual cortex. From left to right: (1) UMAP projection showing distinct neural cell subclasses. (2) The alignment between the original measured data and the transferred data using OMT. (3) The inferred global developmental trajectory at the cluster level, tracing paths from early progenitors like neuroepithelial cells (NEC) and radial glia (RG). (4) The specific cellular pathway detailing the differentiation from OPC to oligodendrocytes. Bottom row: mouse aging dataset. From left to right: (1) UMAP projection of cell subtypes within the oligodendrocyte lineage. (2) The alignment of original and transferred data distributions. (3) The network graph illustrating the stages of the myelination cycle in aged mice. (4) The inferred aging pathway.

34,636 genes. Similarly, for the transport analysis, we perform dimensionality reduction using PCA, retaining the same number of PCs as in Kassraie et al. (2024). Table 1 summarizes the results for OMT against EOT and ProgOT, which is the top-performing baseline. The results indicate that our OMT model outperforms ProgOT across all treatment conditions. For this comparison, ProgOT is configured with the recommended scheduling parameters from its original publication, with K=4.

We note that, here, the data subset for each task is relatively small ( $\sim 10^4$  cells). While this scale is computationally feasible for sample-based approaches like PRoGOT, it does not represent the large-scale datasets in modern single-cell studies. To extend our analysis beyond small-scale data, we apply OMT to larger scRNA-seq datasets from the mouse brain, encompassing the entire lifespan from development to aging. For brain development, we use data from the visual cortex spanning a wide period from embryonic days to postnatal days (E11.5-P28) (Gao et al., 2024). For aging, we consider data from Jin et al. (2025) collected from 108 mice, span six brain regions at two timepoints: adult (P53–69) and aged (P540–553). Our analysis focuses on the cellular dynamics of the oligodendrocyte lineage, including oligodendrocyte precursor cells (OPCs) and mature oligodendrocytes (Oligos). These glial cells, which are responsible for myelinating axons to facilitate neural communication, exhibit significant heterogeneity in their lifespan and function, making them a suitable candidate for studying time-dependent cellular transitions (Marques et al., 2016; Jin et al., 2025).

After preprocessing (Appendix C), the data includes 32,998 cells and 9,900 highly variable genes (HVGs) from the developmental data, alongside 253,468 cells and 9,359 HVGs from the ageing dataset. We utilized a VAE model to learn a compressed representation of the cells. The OMT model was then trained on these low-dimensional embeddings ( $d_z=10$ ). OMT is applied across 11 consecutive time pairs between E11.5 and P28 for the developmental data, and between adult and aged time points for ageing data. Figure 4 summarizes the analysis of the mouse datasets. The UMAP plots show that the cell population transported by the model, whether forward or backward in time, closely mirrors the empirical cell distribution at the target timepoints. This demonstrates the model's ability to learn the global distribution across cell subclasses. The right panels of the figure illustrate the clear developmental and aging trajectories revealed by our OMT model. The transport map reveals the known developmental pathway, beginning with neuroepithelial cells (NECs) that mature into radial glia (RG). These cells subsequently differentiate into glioblasts (Gliob), which are the common progenitors for both the astrocyte (Astro) and the OPC-Oligo lineages. This temporal



Figure 5: Performance of OMT for unpaired image-to-image translation on the MNIST and CIFAR-10 datasets. For each dataset, the top row shows original samples from the source distribution,  $x \sim \nu_0$ , and the bottom row shows the corresponding transported images  $T_{00}^{\nu_0 \to \nu_1}$ .

progression is visually represented by a color gradient, transitioning from yellow (E11.5) to dark red (P28). Focusing on the OPC-Oligo lineage, the rightmost column provides a detailed view of this population during development (top) and aging (bottom). It highlights the specific cellular maturation sequence from oligodendrocyte precursor cells (OPCs) to committed oligodendrocyte precursors (COPs), newly formed oligodendrocytes (NFOLs), myelin-forming oligodendrocytes (MFOLs), and finally, mature oligodendrocytes (MOLs).

**Image Datasets.** To further demonstrate the applicability of the the proposed OMT framework beyond tabular data, we apply it to an unpaired image-to-image translation task using two benchmark datasets: MNIST (LeCun, 1998) and CIFAR-10 (Krizhevsky et al., 2009). In MNIST, the task involves translating images of one digit into another (e.g., learning transport maps such as  $T:1\to 7$ ). Similarly, in CIFAR-10, the goal is to translate images from one semantic class (e.g., airplane) into another (e.g., bird). Although OMT can in principle be applied directly to raw image data, the resulting mappings are not semantically meaningful and fail to capture class-level translations. To enable this, we first train an autoencoder on the entire dataset, covering all classes, to obtain compact and semantically meaningful low-dimensional embeddings. Within this latent space, the OMT is then applied to learn optimal transport maps across different classes.

Figures 5 and 11 illustrate representative examples of these class-to-class translations for test images in MNIST and CIFAR-10, respectively. Quantitative evaluation of the generated translations is reported in Table 2 using the widely adopted Fréchet Inception Distance (FID). These experiments highlight that, OMT can be effectively extended to image-based applications as well. For context and to benchmark our performance against established OT based approaches, we also report the FID scores for WGAN (Arjovsky et al., 2017) and WGAN-GP (Gulrajani et al., 2017). The results show that OMT performs in a similar range to WGAN on CIFAR-10, while outperforming both WGAN and WGAN-GP on the MNIST dataset. See Appendix D for further implementation details, including the autoencoder architectures used for dimensionality reduction and the hyperparameters for OMT.

	MNIST	CIFAR-10
WGAN	$6.7 \pm 0.4$	55.2
WGAN-GP	$7.43 \pm 0.3$	39.4
OMT	$1.2 \pm 0.1$	$56.13 \pm 1.5$

Table 2: FID ↓ values for unpaired image translation on the MNIST (grayscale) and CIFAR-10 (color) datasets. Reported values for WGAN and WGAN-GP are taken from previous studies (Choi et al., 2023; Rout et al., 2022; Qian et al., 2021). Results for OMT are computed over 10 random initializations.

## 6 CONCLUSION

In this work, we introduced OMT, a new family of EOT solvers that enhance performance by moving beyond sample-to-sample transportation toward subpopulation-level transportation, leveraging mixture-model representations and the closed-form structure of Gaussian families. OMT achieves computational efficiency through a strictly biconvex formulation which, when embedded in a global optimization framework, ensures that each subproblem converges in a single step to a unique solution, thereby providing a stable and reliable estimator of the OT plan. Empirically, we showed OMT matches or exceeds the performance of state-of-the-art non-neural OT solvers, while remaining competitive with neural approaches, but with substantially lower computational and memory requirements. One promising direction for future work is extending OMT to the unbalanced OT setting, particularly for mixtures with unequal component masses. This is especially relevant for real-world applications such as single-cell RNA-seq analysis, where unbalanced transport naturally arises. A current limitation, however, lies in applying OMT to high-resolution image generation tasks, where straightforward extensions are not yet practical. Another interesting avenue would be to explore a neural extension of OMT, enabling scalable applications in such high-dimensional domains.

# REFERENCES

- Michael S Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic interpolants. *International conference on learning representations*, 2023.
- Michael S Albergo, Mark Goldstein, Nicholas M Boffi, Rajesh Ranganath, and Eric Vanden-Eijnden. Stochastic interpolants with data-dependent couplings. *International conference on machine learning*, 2024.
  - Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *International conference on machine learning*, pp. 214–223. PMLR, 2017.
    - Rajendra Bhatia, Tanvi Jain, and Yongdo Lim. On the bures—wasserstein distance between positive definite matrices. *Expositiones Mathematicae*, 37(2):165–191, 2019.
    - Charlotte Bunne, Stefan G Stark, Gabriele Gut, Jacobo Sarabia Del Castillo, Mitch Levesque, Kjong-Van Lehmann, Lucas Pelkmans, Andreas Krause, and Gunnar Rätsch. Learning single-cell perturbation responses using neural optimal transport. *Nature methods*, 20(11):1759–1768, 2023.
    - Charlotte Bunne, Geoffrey Schiebinger, Andreas Krause, Aviv Regev, and Marco Cuturi. Optimal transport for single-cell and spatial omics. *Nature Reviews Methods Primers*, 4(1):58, 2024.
    - Nazar Buzun, Maksim Bobrin, and Dmitry V Dylov. Enot: Expectile regularization for fast and accurate training of neural optimal transport. *arXiv preprint arXiv:2403.03777*, 2024a.
    - Nazar Buzun, Maksim Bobrin, and Dmitry V Dylov. Expectile regularization for fast and accurate training of neural optimal transport. *Advances in Neural Information Processing Systems*, 37: 119811–119837, 2024b.
    - Yukun Chen, Jianbo Ye, and Jia Li. Aggregated wasserstein distance and state registration for hidden markov models. *IEEE transactions on pattern analysis and machine intelligence*, 42(9):2133–2147, 2019.
    - Jaemoo Choi, Jaewoong Choi, and Myungjoo Kang. Generative modeling through the semi-dual formulation of unbalanced optimal transport. *Advances in Neural Information Processing Systems*, 36:42433–42455, 2023.
    - Ching-Yao Chuang, Stefanie Jegelka, and David Alvarez-Melis. Infoot: Information maximizing optimal transport. In *International Conference on Machine Learning*, pp. 6228–6242. PMLR, 2023.
    - Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. *Advances in neural information processing systems*, 26, 2013.
    - Marco Cuturi, Laetitia Meng-Papaxanthos, Yingtao Tian, Charlotte Bunne, Geoff Davis, and Olivier Teboul. Optimal transport tools (ott): A jax toolbox for all things wasserstein. *arXiv preprint arXiv:2201.12324*, 2022.
    - Marco Cuturi, Michal Klein, and Pierre Ablin. Monge, bregman and occam: Interpretable optimal transport in high-dimensions with feature-sparse maps. *International Conference on Machine Learning*, 2023.
    - Max Daniels, Tyler Maunu, and Paul Hand. Score-based generative neural networks for large-scale optimal transport. *Advances in neural information processing systems*, 34:12955–12965, 2021.
    - Julie Delon and Agnes Desolneux. A wasserstein-type distance in the space of gaussian mixture models. *SIAM Journal on Imaging Sciences*, 13(2):936–970, 2020.
    - Pinar Demetci, Rebecca Santorella, Björn Sandstede, William Stafford Noble, and Ritambhara Singh. Scot: single-cell multi-omics alignment with optimal transport. *Journal of computational biology*, 29(1):3–18, 2022.
    - Kilian Fatras, Thibault Séjourné, Rémi Flamary, and Nicolas Courty. Unbalanced minibatch optimal transport; applications to domain adaptation. In *International conference on machine learning*, pp. 3186–3197. PMLR, 2021a.

- Kilian Fatras, Younes Zine, Szymon Majewski, Rémi Flamary, Rémi Gribonval, and Nicolas Courty. Minibatch optimal transport distances; analysis and applications. *arXiv preprint arXiv:2101.01792*, 2021b.
  - Eduardo Fernandes Montesuma, Fred Maurice Ngolè Mboula, and Antoine Souloumiac. Optimal transport for domain adaptation through gaussian mixture models. *Transactions on Machine Learning Research*, 2025.
  - Christodoulos A Floudas and Vishy Visweswaran. Primal-relaxed dual global optimization approach. *Journal of Optimization Theory and Applications*, 78(2):187–225, 1993.
  - Christodoulos A Floudas and Viswanathan Visweswaran. A global optimization algorithm (gop) for certain classes of nonconvex nlps—i. theory. *Computers & chemical engineering*, 14(12): 1397–1417, 1990.
  - Yuan Gao, Cindy TJ van Velthoven, Changkyu Lee, Emma D Thomas, Darren Bertagnolli, Daniel Carey, Tamara Casper, Anish Bhaswanth Chakka, Rushil Chakrabarty, Michael Clark, et al. Continuous cell type diversification throughout the embryonic and postnatal mouse visual cortex development. *bioRxiv*, 2024.
  - Aude Genevay, Gabriel Peyré, and Marco Cuturi. Learning generative models with sinkhorn divergences. In *International Conference on Artificial Intelligence and Statistics*, pp. 1608–1617. PMLR, 2018.
  - Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT press Cambridge, 2016.
  - Jochen Gorski, Frank Pfeuffer, and Kathrin Klamroth. Biconvex sets and optimization with biconvex functions: a survey and extensions. *Mathematical methods of operations research*, 66(3):373–407, 2007.
  - Edouard Grave, Armand Joulin, and Quentin Berthet. Unsupervised alignment of embeddings with wasserstein procrustes. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 1880–1890. PMLR, 2019.
  - Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of wasserstein gans. *Advances in neural information processing systems*, 30, 2017.
  - Nikita Gushchin, Alexander Kolesov, Alexander Korotin, Dmitry P Vetrov, and Evgeny Burnaev. Entropic neural optimal transport via diffusion processes. *Advances in Neural Information Processing Systems*, 36:75517–75544, 2023a.
  - Nikita Gushchin, Alexander Kolesov, Petr Mokrov, Polina Karpikova, Andrei Spiridonov, Evgeny Burnaev, and Alexander Korotin. Building the bridge of schrödinger: A continuous entropic optimal transport benchmark. *Advances in Neural Information Processing Systems*, 36:18932–18963, 2023b.
  - Nikita Gushchin, Alexander Kolesov, Alexander Korotin, Dmitry P Vetrov, and Evgeny Burnaev. Entropic neural optimal transport via diffusion processes. *Advances in Neural Information Processing Systems*, 36, 2024.
  - Hicham Janati, Boris Muzellec, Gabriel Peyré, and Marco Cuturi. Entropic optimal transport between unbalanced gaussian measures has a closed form. *Advances in neural information processing systems*, 33:10468–10479, 2020.
  - Kelly Jin, Zizhen Yao, Cindy TJ van Velthoven, Eitan S Kaplan, Katie Glattfelder, Samuel T Barlow, Gabriella Boyer, Daniel Carey, Tamara Casper, Anish Bhaswanth Chakka, et al. Brain-wide cell-type-specific transcriptomic signatures of healthy ageing in mice. *Nature*, 638(8049):182–196, 2025.
  - Parnian Kassraie, Aram-Alexandre Pooladian, Michal Klein, James Thornton, Jonathan Niles-Weed, and Marco Cuturi. Progressive entropic optimal transport solvers. *Advances in Neural Information Processing Systems*, 37:19561–19590, 2024.

- Alexander Korotin, Lingxiao Li, Aude Genevay, Justin M Solomon, Alexander Filippov, and Evgeny Burnaev. Do neural optimal transport solvers work? a continuous wasserstein-2 benchmark. *Advances in neural information processing systems*, 34:14593–14605, 2021.
  - Alexander Korotin, Daniil Selikhanovych, and Evgeny Burnaev. Neural optimal transport. *International conference on learning representations*, 2023.
  - Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-100 and cifar-10 (canadian institute for advanced research). *URL http://www. cs. toronto. edu/kriz/cifar. html. MIT License*, 2009.
  - Yann LeCun. The mnist database of handwritten digits. http://yann. lecun. com/exdb/mnist/, 1998.
  - Ashok Makkuva, Amirhossein Taghvaei, Sewoong Oh, and Jason Lee. Optimal transport mapping via input convex neural networks. In *International Conference on Machine Learning*, pp. 6672–6681. PMLR, 2020.
  - Sueli Marques, Amit Zeisel, Simone Codeluppi, David Van Bruggen, Ana Mendanha Falcão, Lin Xiao, Huiliang Li, Martin Häring, Hannah Hochgerner, Roman A Romanov, et al. Oligodendrocyte heterogeneity in the mouse juvenile and adult central nervous system. *Science*, 352(6291):1326–1329, 2016.
  - Robert J McCann. A convexity principle for interacting gases. *Advances in mathematics*, 128(1): 153–179, 1997.
  - Gabriel Peyré, Marco Cuturi, et al. Computational optimal transport: With applications to data science. *Foundations and Trends*® *in Machine Learning*, 11(5-6):355–607, 2019.
  - Wenliang Qian, Yang Xu, Wangmeng Zuo, and Hui Li. Self sparse generative adversarial networks. *arXiv preprint arXiv:2101.10556*, 2021.
  - Litu Rout, Alexander Korotin, and Evgeny Burnaev. Generative modeling with optimal transport maps. *International conference on machine learning*, 2022.
  - Filippo Santambrogio. Optimal transport for applied mathematicians, volume 87. Springer, 2015.
  - David W Scott. *Multivariate density estimation: theory, practice, and visualization.* John Wiley & Sons, 2015.
  - Sanjay R Srivatsan, José L McFaline-Figueroa, Vijay Ramani, Lauren Saunders, Junyue Cao, Jonathan Packer, Hannah A Pliner, Dana L Jackson, Riza M Daza, Lena Christiansen, et al. Massively multiplex chemical transcriptomics at single-cell resolution. *Science*, 367(6473):45–51, 2020.
  - Oliver Struckmeier, Ievgen Redko, Anton Mallasto, Karol Arndt, Markus Heinonen, and Ville Kyrki. Learning representations that are closed-form monge mapping optimal with application to domain adaptation. *Transactions on Machine Learning Research*, 2023.
  - David Michael Titterington, Adrian FM Smith, and Udi E Makov. *Statistical analysis of finite mixture distributions*. John Wiley, 1985.
  - Alexander Tong, Jessie Huang, Guy Wolf, David Van Dijk, and Smita Krishnaswamy. Trajectorynet: A dynamic optimal transport network for modeling cellular dynamics. In *International conference on machine learning*, pp. 9526–9536. PMLR, 2020.
  - Cédric Villani et al. Optimal transport: old and new, volume 338. Springer, 2008.
  - Weihang Xu, Maryam Fazel, and Simon S Du. Toward global convergence of gradient em for over-paramterized gaussian mixture models. *Advances in Neural Information Processing Systems*, 37:10770–10800, 2024.
- Assaf J Zeevi and Ronny Meir. Density estimation through convex combinations of densities: approximation and estimation bounds. *Neural Networks*, 10(1):99–109, 1997.

**APPENDIX** 

### A Proofs

**Lemma 2.** For any  $\varepsilon_1, \varepsilon_2 > 0$ ,  $\mathcal{L}_{\varepsilon_1, \varepsilon_2}(\Omega, P)$  is strictly biconvex.

*Proof.* A function  $f: \mathcal{X} \times \mathcal{Y} \to \mathbb{R}$  is called *biconvex* if, for fixed  $x \in \mathcal{X}$ , the function f(x,y) is convex in y, and for fixed  $y \in \mathcal{Y}$ , it is convex in x. According to **Theorem 3.1** Gorski et al. (2007), f(x,y) is biconvex if and only if for all  $(x_1,y_1), (x_1,y_2), (x_2,y_1), (x_2,y_2) \in \mathcal{X} \times \mathcal{Y}$  and all  $\lambda, \tau \in [0,1]$ , the following inequality holds:

$$f(x_{\lambda}, y_{\tau}) \leq \lambda \tau f(x_1, y_1) + (1 - \lambda)\tau f(x_2, y_1) + \lambda (1 - \tau)f(x_1, y_2) + (1 - \lambda)(1 - \tau)f(x_2, y_2),$$

where 
$$(x_{\lambda}, y_{\tau}) := (\lambda x_1 + (1 - \lambda)x_2, \ \tau y_1 + (1 - \tau)y_2).$$

Then, for given  $\omega_{ij_{\lambda}} = \lambda \tilde{\omega}_{ij} + (1 - \lambda) \tilde{\tilde{\omega}}_{ij}$  and  $dp_{ij_{\tau}} = \tau \tilde{dp}_{ij} + (1 - \tau) \tilde{dp}_{ij}$ , the following inequality must hold.

$$\mathcal{L}_{\varepsilon_{1},\varepsilon_{2}}(\Omega_{\lambda},P_{\tau}) < \lambda \tau \mathcal{L}_{\varepsilon_{1},\varepsilon_{2}}(\tilde{\Omega},\tilde{P}) + (1-\lambda)\tau \mathcal{L}_{\varepsilon_{1},\varepsilon_{2}}(\tilde{\tilde{\Omega}},\tilde{P}) + \lambda(1-\tau)\mathcal{L}_{\varepsilon_{1},\varepsilon_{2}}(\tilde{\Omega},\tilde{\tilde{P}}) + (1-\lambda)(1-\tau)\mathcal{L}_{\varepsilon_{1},\varepsilon_{2}}(\tilde{\tilde{\omega}},\tilde{\tilde{P}})$$
(15)

$$\mathcal{L}_{\varepsilon_{1},\varepsilon_{2}}(\Omega_{\lambda}, P_{\tau}) = \sum_{i,j}^{K} \left(\lambda \tilde{\omega}_{ij} + (1-\lambda)\tilde{\omega}_{ij}\right) \left[ \int_{\mathcal{X} \times \mathcal{Y}} \|\mathbf{x} - \mathbf{y}\|_{2}^{2} \left(\tau d\tilde{p}_{ij}(\mathbf{x}, \mathbf{y}) + (1-\tau)d\tilde{\tilde{p}}_{ij}(\mathbf{x}, \mathbf{y})\right) \right] + \\
\varepsilon_{2} D_{KL}(\Omega_{\lambda} \|\boldsymbol{\alpha}_{0} \otimes \boldsymbol{\alpha}_{1}) + \sum_{i,j}^{K} \left(\lambda \tilde{\omega}_{ij} + (1-\lambda)\tilde{\tilde{\omega}}_{ij}\right) \varepsilon_{1} D_{KL}(p_{ij_{\tau}} \|\mu_{0_{i}} \otimes \mu_{1_{j}}) \\
= \lambda \tau \sum_{i,j} \tilde{\omega}_{ij} \int \|\mathbf{x} - \mathbf{y}\|_{2}^{2} d\tilde{p}_{ij}(\mathbf{x}, \mathbf{y}) + (1-\lambda)\tau \sum_{i,j} \tilde{\omega}_{ij} \int \|\mathbf{x} - \mathbf{y}\|_{2}^{2} d\tilde{p}_{ij}(\mathbf{x}, \mathbf{y}) + \\
\lambda (1-\tau) \sum_{i,j} \tilde{\omega}_{ij} \int \|\mathbf{x} - \mathbf{y}\|_{2}^{2} d\tilde{p}_{ij}(\mathbf{x}, \mathbf{y}) + (1-\lambda)(1-\tau) \sum_{i,j} \tilde{\omega}_{ij} \int \|\mathbf{x} - \mathbf{y}\|_{2}^{2} d\tilde{p}_{ij}(\mathbf{x}, \mathbf{y}) + \\
\varepsilon_{1} \sum_{i,j} \left(\lambda \tilde{\omega}_{ij} + (1-\lambda)\tilde{\tilde{\omega}}_{ij}\right) \underbrace{D_{KL}(p_{ij_{\tau}} \|\mu_{0_{i}} \otimes \mu_{1_{j}})}_{E} + \varepsilon_{2} \underbrace{D_{KL}(\Omega_{\lambda} \|\boldsymbol{\alpha}_{0} \otimes \boldsymbol{\alpha}_{1})}_{F}. \tag{16}$$

$$\lambda \tau \mathcal{L}_{\varepsilon_{1},\varepsilon_{2}}(\tilde{\Omega},\tilde{P}) + (1-\lambda)\tau \mathcal{L}_{\varepsilon_{1},\varepsilon_{2}}(\tilde{\Omega},\tilde{P}) + \lambda(1-\tau)\mathcal{L}_{\varepsilon_{1},\varepsilon_{2}}(\tilde{\Omega},\tilde{P}) + (1-\lambda)(1-\tau)\mathcal{L}_{\varepsilon_{1},\varepsilon_{2}}(\tilde{\omega},\tilde{P}) = \lambda\tau \sum_{i,j} \tilde{\omega}_{ij} \int \|\mathbf{x} - \mathbf{y}\|_{2}^{2} d\tilde{p}_{ij}(\mathbf{x},\mathbf{y}) + \varepsilon_{1} \sum_{i,j} \lambda \tilde{\omega}_{ij} \left(\tau D_{KL}(\tilde{p}_{ij_{\tau}} \|\mu_{0_{i}} \otimes \mu_{1_{j}})\right) + \varepsilon_{2}\tau \lambda D_{KL}(\tilde{\Omega} \|\boldsymbol{\alpha}_{0} \otimes \boldsymbol{\alpha}_{1}) + (1-\lambda)\tau \sum_{i,j} \tilde{\omega}_{ij} \int \|\mathbf{x} - \mathbf{y}\|_{2}^{2} d\tilde{p}_{ij}(\mathbf{x},\mathbf{y}) + \varepsilon_{2}\tau (1-\lambda)D_{KL}(\tilde{\Omega} \|\boldsymbol{\alpha}_{0} \otimes \boldsymbol{\alpha}_{1}) + \varepsilon_{2}\tau (1-\lambda)D_{KL}(\tilde{\Omega} \|\boldsymbol{\alpha}_{0} \otimes \boldsymbol{\alpha}_{1}) + \lambda(1-\tau)\sum_{i,j} \tilde{\omega}_{ij} \int \|\mathbf{x} - \mathbf{y}\|_{2}^{2} d\tilde{p}_{ij}(\mathbf{x},\mathbf{y}) + \varepsilon_{2}\tau (1-\lambda)D_{KL}(\tilde{\Omega} \|\boldsymbol{\alpha}_{0} \otimes \boldsymbol{\alpha}_{1}) + \varepsilon_{2}\tau (1-\tau)\sum_{i,j} \tilde{\omega}_{ij} \int \|\mathbf{x} - \mathbf{y}\|_{2}^{2} d\tilde{p}_{ij}(\mathbf{x},\mathbf{y}) + \varepsilon_{2}\tau (1-\tau)\sum_{i,j} \tilde{\omega}_{ij} \int \|\mathbf{x} - \mathbf{y}\|_{2}^{2} d\tilde{p}_{ij}(\mathbf{x},\mathbf{y}) + \varepsilon_{2}\tau (1-\tau) \sum_{i,j} \tilde{\omega}_{ij} \int \|\mathbf{x} - \mathbf{y}\|_{2}^{2} d\tilde{p}_{ij}(\mathbf{x},\mathbf{y}) + \varepsilon_{2}\tau (1-\tau) \sum_{i,j} \tilde{\omega}_{ij} \int \|\mathbf{x} - \mathbf{y}\|_{2}^{2} d\tilde{p}_{ij}(\mathbf{x},\mathbf{y}) + \varepsilon_{2}\tau (1-\tau) \sum_{i,j} \tilde{\omega}_{ij} \int \|\mathbf{x} - \mathbf{y}\|_{2}^{2} d\tilde{p}_{ij}(\mathbf{x},\mathbf{y}) + \varepsilon_{2}\tau (1-\tau) \sum_{i,j} \tilde{\omega}_{ij} \int \|\mathbf{x} - \mathbf{y}\|_{2}^{2} d\tilde{p}_{ij}(\mathbf{x},\mathbf{y}) + \varepsilon_{2}\tau (1-\tau) \sum_{i,j} \tilde{\omega}_{ij} \int \|\mathbf{x} - \mathbf{y}\|_{2}^{2} d\tilde{p}_{ij}(\mathbf{x},\mathbf{y}) + \varepsilon_{2}\tau (1-\tau) \sum_{i,j} \tilde{\omega}_{ij} \int \|\mathbf{x} - \mathbf{y}\|_{2}^{2} d\tilde{p}_{ij}(\mathbf{x},\mathbf{y}) + \varepsilon_{2}\tau (1-\tau) \sum_{i,j} \tilde{\omega}_{ij} \int \|\mathbf{x} - \mathbf{y}\|_{2}^{2} d\tilde{p}_{ij}(\mathbf{x},\mathbf{y}) + \varepsilon_{2}\tau (1-\tau) \sum_{i,j} \tilde{\omega}_{ij} \int \|\mathbf{x} - \mathbf{y}\|_{2}^{2} d\tilde{p}_{ij}(\mathbf{x},\mathbf{y}) + \varepsilon_{2}\tau (1-\tau) \sum_{i,j} \tilde{\omega}_{ij} \int \|\mathbf{x} - \mathbf{y}\|_{2}^{2} d\tilde{p}_{ij}(\mathbf{x},\mathbf{y}) + \varepsilon_{2}\tau (1-\tau) \sum_{i,j} \tilde{\omega}_{ij} \int \|\mathbf{x} - \mathbf{y}\|_{2}^{2} d\tilde{p}_{ij}(\mathbf{x},\mathbf{y}) + \varepsilon_{2}\tau (1-\tau) \sum_{i,j} \tilde{\omega}_{ij}(\mathbf{x},\mathbf{y}) + \varepsilon_{2}\tau (1-\tau) \sum_{i,j}$$

$$\varepsilon_{1} \sum_{i,j} \lambda \tilde{\omega}_{ij} \left( (1-\tau) D_{KL}(\tilde{p}_{ij_{\tau}} \| \mu_{0_{i}} \otimes \mu_{1_{j}}) \right) + \varepsilon_{2} (1-\tau) \lambda D_{KL}(\tilde{\Omega} \| \boldsymbol{\alpha}_{0} \otimes \boldsymbol{\alpha}_{1}) + \\
(1-\lambda)(1-\tau) \sum_{i,j} \tilde{\omega}_{ij} \int \|\mathbf{x} - \mathbf{y}\|_{2}^{2} \tilde{d} \tilde{p}_{ij}(\mathbf{x}, \mathbf{y}) + \varepsilon_{1} \sum_{i,j} (1-\lambda) \tilde{\omega}_{ij} \left( (1-\tau) D_{KL}(\tilde{p}_{ij} \| \mu_{0_{i}} \otimes \mu_{1_{j}}) \right) + \\
\underline{\qquad \qquad \qquad } \\
\varepsilon_{2} (1-\tau)(1-\lambda) D_{KL}(\tilde{\Omega} \| \boldsymbol{\alpha}_{0} \otimes \boldsymbol{\alpha}_{1}),$$

$$= A + B + C + D + \varepsilon_{2} \left( \underbrace{\lambda D_{KL}(\tilde{\Omega} \| \boldsymbol{\alpha}_{0} \otimes \boldsymbol{\alpha}_{1}) + (1 - \lambda) D_{KL}(\tilde{\tilde{\Omega}} \| \boldsymbol{\alpha}_{0} \otimes \boldsymbol{\alpha}_{1})}_{G} \right) +$$

$$\varepsilon_{1} \sum_{i,j} \left( \lambda \tilde{\omega}_{ij} + (1 - \lambda \tilde{\tilde{\omega}}_{ij}) \right) \left( \underbrace{\tau D_{KL}(\tilde{p}_{ij} \| \mu_{0_{i}} \otimes \mu_{1_{j}}) + (1 - \tau) D_{KL}(\tilde{\tilde{p}}_{ij} \| \mu_{0_{i}} \otimes \mu_{1_{j}})}_{H} \right).$$

$$(17)$$

For any fixed q,  $D_{KL}(p||q)$  is strictly convex in p. Consequently, we have E < H and F < G, which together imply that inequality (15) holds.

**Theorem 2.** For the optimization problem defined in (13), the GOP algorithm converges to a unique solution in a single iteration.

*Proof.* Consider the biconvex optimization problem defined as follows:

$$\min \{ \mathcal{L}_{\varepsilon_1 \varepsilon_2}(\Omega, P), \ (\Omega, P) \in \Lambda \}.$$

Let's begin by selecting an arbitrary initial point  $Z_0=(\Omega_0,P_0)\in\Lambda$  and set the iteration index s=0. Without loss of generality, we assume that  $\omega_{ij}^0>0$ , for all i,j. We then solve the following convex optimization problem with respect to P, keeping  $\Omega_s$  fixed.

$$\min_{P} \sum_{i,j} \omega_{ij}^{0} \left[ \int_{\mathcal{X} \times \mathcal{Y}} \|\mathbf{x} - \mathbf{y}\|_{2}^{2} dp_{ij}(\mathbf{x}, \mathbf{y}) + \varepsilon_{1} D_{KL}(p_{ij} \| \mu_{0_{i}} \otimes \mu_{1_{j}}) \right]$$
s.t. 
$$\int_{\mathcal{Y}} dP(\mathbf{x}, \mathbf{y}) = \boldsymbol{\mu}_{0}, \quad \int_{\mathcal{X}} dP(\mathbf{x}, \mathbf{y}) = \boldsymbol{\mu}_{1} \tag{18}$$

Since the objective function in (A) is convex in P for any fixed  $\Omega$  and Slater's condition is satisfied Floudas & Visweswaran (1993), strong duality holds. Accordingly, problem (A) admits the following strong dual formulation:

$$\max_{\Phi, \Psi} \min_{P} \sum_{i,j} \omega_{ij}^{0} \left[ \int_{\mathcal{X} \times \mathcal{Y}} \|\mathbf{x} - \mathbf{y}\|_{2}^{2} dp_{ij}(\mathbf{x}, \mathbf{y}) + \varepsilon_{1} D_{KL}(p_{ij} \| \mu_{0_{i}} \otimes \mu_{1_{j}}) \right] -$$
(19)

$$\sum_{i,j} \omega_{ij}^{0} \left[ \int_{\mathcal{X}} \varphi_{ij}(\mathbf{x}) \left( \int_{\mathcal{Y}} dp_{ij}(\mathbf{x}, \mathbf{y}) - d\mu_{0_{i}}(\mathbf{x}) \right) + \int_{\mathcal{Y}} \psi_{ij}(\mathbf{y}) \left( \int_{\mathcal{X}} dp_{ij}(\mathbf{x}, \mathbf{y}) - d\mu_{1_{j}}(\mathbf{y}) \right) \right]. \tag{20}$$

Here,  $\phi_{ij}, \psi_{ij} \geq 0$  are the Lagrange multipliers associated with the marginal constraints.

Denoting the inner minimization problem in (19) by  $\min_{P} f(P, \Phi, \Psi)$ , we seek the optimal policy that minimizes the loss in (19). To do so, we compute the functional derivative of the loss with respect to  $dp_{ij}(\mathbf{x}, \mathbf{y})$ .

$$\frac{df}{dp_{ij}(\mathbf{x}, \mathbf{y})} = \omega_{ij}^{0} \left( \|\mathbf{x} - \mathbf{y}\|_{2}^{2} + \varepsilon_{1} \log \frac{dp_{ij}(\mathbf{x}, \mathbf{y})}{d\mu_{0_{i}}(\mathbf{x}) d\mu_{1_{i}}(\mathbf{y})} - \varphi_{ij}(\mathbf{x}) - \psi_{ij}(\mathbf{y}) \right). \tag{21}$$

Since the objective is strictly convex in P, it admits a unique solution that is independent of  $\Omega$ .

$$dp_{ij}^*(\mathbf{x}, \mathbf{y}) = \exp\left(\frac{\varphi_{ij}(\mathbf{x}) + \psi_{ij}(\mathbf{y}) - \|\mathbf{x} - \mathbf{y}\|_2^2}{\varepsilon_1}\right) d\mu_{0_i}(\mathbf{x}) d\mu_{1_j}(\mathbf{y}). \tag{22}$$

Substituting  $dp_{ij}^*$  in (19), the dual form can be written as:

$$\max_{\Phi,\Psi} \sum_{i,j} \omega_{ij}^{0} \left[ \int_{\mathcal{X}} \varphi_{ij}(\mathbf{x}) d\mu_{0_{i}}(\mathbf{x}) - \int_{\mathcal{Y}} \psi_{ij}(\mathbf{y}) d\mu_{0_{j}}(\mathbf{y}) - \varepsilon_{1} \left( \int_{\mathcal{X} \times \mathcal{Y}} dp_{i,j}^{*}(\mathbf{x}, \mathbf{y}) - 1 \right) \right].$$
 (23)

To determine the optimal Lagrange multipliers that maximize the dual objective, we differentiate the loss function in (23), denoted as  $g(\Phi, \Psi, P^*)$  with respect to each multiplier as follows.

$$\frac{dg}{\varphi_{ij}(\mathbf{x})} = \omega_{ij}^{0} \left( d\mu_{0_i}(\mathbf{x}) - \int_{\mathcal{Y}} \exp\left(\frac{\varphi_{ij}(\mathbf{x}) + \psi_{ij}(\mathbf{y}) - \|\mathbf{x} - \mathbf{y}\|_{2}^{2}}{\varepsilon_{1}}\right) d\mu_{0_i}(\mathbf{x}) d\mu_{1_j}(\mathbf{y}) \right)$$
(24)

$$\frac{dg}{\psi_{ij}(\mathbf{x})} = \omega_{ij}^{0} \left( d\mu_{1_{j}}(\mathbf{x}) - \int_{\mathcal{X}} \exp\left(\frac{\varphi_{ij}(\mathbf{x}) + \psi_{ij}(\mathbf{y}) - \|\mathbf{x} - \mathbf{y}\|_{2}^{2}}{\varepsilon_{1}}\right) d\mu_{0_{i}}(\mathbf{x}) d\mu_{1_{j}}(\mathbf{y}) \right)$$
(25)

Assuming that, for all i, j, the measures  $\mu_{0_i}$  and  $\mu_{1_j}$  have finite second-order moments, the pair  $(\varphi_{ij}, \psi_{ij})$  is optimal if and only if the following conditions are satisfied.

$$\int_{\mathcal{Y}} \exp\left(\frac{\varphi_{ij}(\mathbf{x}) + \psi_{ij}(\mathbf{y}) - \|\mathbf{x} - \mathbf{y}\|_{2}^{2}}{\varepsilon_{1}}\right) d\mu_{1_{j}}(\mathbf{y}) = 1, \qquad \int_{\mathcal{X}} \exp\left(\frac{\varphi_{ij}(\mathbf{x}) + \psi_{ij}(\mathbf{y}) - \|\mathbf{x} - \mathbf{y}\|_{2}^{2}}{\varepsilon_{1}}\right) d\mu_{0_{i}}(\mathbf{x}) = 1,$$
(26)

which is equivalent to the following expressions for the optimal multipliers:

$$\varphi_{ij}(\mathbf{x}) = -\varepsilon_1 \log \int_{\mathcal{Y}} \exp\left(\frac{\psi_{ij}(\mathbf{y}) - \|\mathbf{x} - \mathbf{y}\|_2^2}{\varepsilon_1}\right) d\mu_{1_j}(\mathbf{y}),$$
 (27)

$$\psi_{ij}(\mathbf{y}) = -\varepsilon_1 \log \int_{\mathcal{X}} \exp \left( \frac{\varphi_{ij}(\mathbf{x}) - \|\mathbf{x} - \mathbf{y}\|_2^2}{\varepsilon_1} \right) d\mu_{0_i}(\mathbf{x}). \tag{28}$$

As observed, the optimal Lagrange multipliers are also independent of  $\Omega$ . Therefore, the unique optimal solution of (19) remains the same for any fixed  $\Omega_s$ , implying

$$\forall s, \quad \mathcal{L}_{\varepsilon_{1}\varepsilon_{2}}(\Omega_{s}, P^{*}) \leq \mathcal{L}_{\varepsilon_{1}\varepsilon_{2}}(\Omega_{s}, P),$$

$$\forall s \neq s', \quad \mathcal{L}_{\varepsilon_{1}\varepsilon_{2}}(\Omega_{s}, P^{*}) = \mathcal{L}_{\varepsilon_{1}\varepsilon_{2}}(\Omega_{s'}, P^{*}),$$

$$\forall s \neq s', \quad \mathcal{L}_{\varepsilon_{1}\varepsilon_{2}}(\Omega_{s}, P^{*}) = \mathcal{L}_{\varepsilon_{1}\varepsilon_{2}}(\Omega_{s'}, P_{s'}), \text{ iff } P_{s'} = P^{*}$$

$$(29)$$

Proceeding to the next step, we set s=1, which gives  $P_1=P^*$ . For a given  $\varepsilon_2>0$ , we solve the following strictly convex optimization problem with respect to  $\Omega$ , keeping P fixed.

$$\min_{\Omega} \sum_{i,j} \omega_{ij} \left[ \int_{\mathcal{X} \times \mathcal{Y}} \|\mathbf{x} - \mathbf{y}\|_{2}^{2} dp_{ij}^{*}(\mathbf{x}, \mathbf{y}) + \varepsilon_{1} D_{KL}(p_{ij}^{*} \| \mu_{0_{i}} \otimes \mu_{1_{j}}) \right] + \varepsilon_{2} D_{KL}(\Omega \| \boldsymbol{\alpha}_{0} \otimes \boldsymbol{\alpha}_{1}) \tag{30}$$

s.t. 
$$\mathbf{1}\Omega = \boldsymbol{\alpha}_0, \quad \Omega^T \mathbf{1} = \boldsymbol{\alpha}_1$$
 (31)

This formulation, similar to equation (A), admits the following dual form:

$$\max_{\boldsymbol{\lambda}, \boldsymbol{\tau}} \min_{\Omega} \sum_{i,j} \omega_{ij} \left[ \int_{\mathcal{X} \times \mathcal{Y}} \|\mathbf{x} - \mathbf{y}\|_{2}^{2} dp_{ij}^{*}(\mathbf{x}, \mathbf{y}) + \varepsilon_{1} D_{KL}(p_{ij}^{*} \|\mu_{0_{i}} \otimes \mu_{1_{j}}) \right] + \varepsilon_{2} D_{KL}(\Omega \|\boldsymbol{\alpha}_{0} \otimes \boldsymbol{\alpha}_{1}) -$$
(32)

$$\sum_{i} \lambda_{i} \left( \sum_{j} \omega_{ij} - \alpha_{0_{i}} \right) - \sum_{j} \tau_{j} \left( \sum_{i} \omega_{ij} - \alpha_{1_{j}} \right). \tag{33}$$

The inner minimization in (32), denoted  $f'(\Omega, \lambda, \tau)$ , admits a unique solution for the optimal weights. These weights can be derived by computing the functional derivative of the objective with respect to  $\omega_{ij}$ , yielding:

$$\frac{df'}{\omega_{ij}} = \underbrace{\int_{\mathcal{X}\times\mathcal{Y}} \|\mathbf{x} - \mathbf{y}\|_{2}^{2} dp_{ij}^{*}(\mathbf{x}, \mathbf{y}) + \varepsilon_{1} D_{KL}(p_{ij}^{*} \|\mu_{0_{i}} \otimes \mu_{1_{j}}) + \varepsilon_{2} \log \frac{\omega_{ij}}{\alpha_{0_{i}} \alpha_{1_{j}}} - \lambda_{i} - \tau_{f}(34)}_{\mathcal{L}_{p_{ij}^{*}}}$$

$$\omega_{ij}^{*} = \exp\left(\frac{\lambda_{i} + \tau_{j} - \mathcal{L}_{p_{ij}^{*}}}{\varepsilon_{2}}\right) \alpha_{0_{i}} \alpha_{1_{j}}$$
(35)

 To obtain the optimal Lagrangian multipliers that maximize the loss in (32), denoted  $\max_{\lambda, \tau} g'(\lambda, \tau, \Omega^*)$ , we compute the partial derivatives of g' with respect to each multiplier.

$$\frac{dg'}{\lambda_i} = \alpha_{0_i} - \sum_{j} \exp\left(\frac{\lambda_i + \tau_j - \mathcal{L}_{p_{ij}^*}}{\varepsilon_2}\right) \alpha_{0_i} \alpha_{1_j}$$
(36)

$$\frac{dg'}{\tau_j} = \alpha_{1_j} - \sum_{i} \exp\left(\frac{\lambda_i + \tau_j - \mathcal{L}_{p_{ij}^*}}{\varepsilon_2}\right) \alpha_{0_i} \alpha_{1_j}$$
(37)

Solving these yields the optimal multipliers as:

$$\lambda_{i} = -\varepsilon_{2} \log \sum_{j} \exp \left( \frac{\tau_{j} - \mathcal{L}_{p_{ij}^{*}}}{\varepsilon_{2}} \right) \alpha_{1_{j}}, \qquad \tau_{j} = -\varepsilon_{2} \log \sum_{i} \exp \left( \frac{\lambda_{i} - \mathcal{L}_{p_{ij}^{*}}}{\varepsilon_{2}} \right) \alpha_{0_{i}}$$
(38)

Since the optimal weight in (35) minimizes  $\mathcal{L}_{\varepsilon_1\varepsilon_2}(\Omega, P_1)$  uniquely for the fixed choice  $P_1 = P^*$ , it follows that:

$$\mathcal{L}_{\varepsilon_1 \varepsilon_2}(\Omega^*, P_1) \le \mathcal{L}_{\varepsilon_1 \varepsilon_2}(\Omega_0, P_1) \le \mathcal{L}_{\varepsilon_1 \varepsilon_2}(\Omega_0, P_0) \tag{39}$$

Now, lets update the optimization by setting  $\Omega_1 = \Omega^*$ , and advancing to step s = 2. According to (29), the next update satisfies:

$$\mathcal{L}_{\varepsilon_1 \varepsilon_2}(\Omega_1, P^*) = \min_{\mathcal{D}} \mathcal{L}_{\varepsilon_1 \varepsilon_2}(\Omega_1, P). \tag{40}$$

Since we find that  $P_2 = P_1 = P^*$  and consequently,  $\Omega_2 = \Omega_1 = \Omega^*$ , the stopping criterion of the overall alternating optimization procedure is met after just a single iteration.

**Corollary 2.** Let  $\nu_0 \in G_{K_0}(\mathbb{R}^d)$  and  $\nu_1 \in G_{K_1}(\mathbb{R}^d)$  be two Gaussian mixture models (GMMs) in  $\mathbb{R}^d$  with  $K_0$  and  $K_1$  components, respectively. Then, the optimal mixture transport map between  $\nu_0$  and  $\nu_1$  is itself a Gaussian mixture model with K components, where  $K \leq K_0K_1$ .

*Proof.* According to **Theorem** 2, the optimal mixture transport policy between each pair  $\mu_0^i(\mathbf{x}), \mu_1^j(\mathbf{y})$ , is independent of the weight variable and is given by:

$$dp_{ij}^*(\mathbf{x}, \mathbf{y}) = \exp\left(\frac{\varphi_{ij}(\mathbf{x}) + \psi_{ij}(\mathbf{y}) - \|\mathbf{x} - \mathbf{y}\|_2^2}{\varepsilon_1}\right) d\mu_{0_i}(\mathbf{x}) d\mu_{1_j}(\mathbf{y}),$$

where  $\phi_{ij}$  and  $\psi_{ij}$  are Lagrange multipliers defined by::

$$\varphi_{ij}(\mathbf{x}) = -\varepsilon_1 \log \int_{\mathcal{Y}} \exp\left(\frac{\psi_{ij}(\mathbf{y}) - \|\mathbf{x} - \mathbf{y}\|_2^2}{\varepsilon_1}\right) d\mu_{1_j}(\mathbf{y}),$$
  
$$\psi_{ij}(\mathbf{y}) = -\varepsilon_1 \log \int_{\mathcal{X}} \exp\left(\frac{\varphi_{ij}(\mathbf{x}) - \|\mathbf{x} - \mathbf{y}\|_2^2}{\varepsilon_1}\right) d\mu_{0_i}(\mathbf{x}).$$

For  $\mu_0^i(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{m}_{i_x}, \Sigma_{i_{xx}})$  and  $\mu_1^j(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{m}_{j_y}, \Sigma_{j_{yy}})$ , it was shown in Janati et al. (2020) that  $\varphi_{ij}$  and  $\psi_{ij}$  admit closed-form solutions in the form of quadratic functions as follows (**Proposition** 1 in Janati et al. (2020)).

$$\varphi_{ij}(\mathbf{x}) = -(\mathbf{x} - \mathbf{m}_{i_x})^T U_{ij}(\mathbf{x} - \mathbf{m}_{i_x}), \quad U_{ij} = \Sigma_{j_{yy}} \left(\Sigma_{ij}^{\varepsilon_1} + \varepsilon_1 \mathbf{I}_d\right)^{-1} - \mathbf{I}_d$$

$$\psi_{ij}(\mathbf{y}) = -(\mathbf{y} - \mathbf{m}_{j_y})^T V_{ij}(\mathbf{y} - \mathbf{m}_{j_y}), \quad V_{ij} = \left(\Sigma_{ij}^{\varepsilon_1} + \varepsilon_1 \mathbf{I}_d\right)^{-1} \Sigma_{i_{xx}} - \mathbf{I}_d$$
(41)

$$\text{where } \Sigma_{ij}^{\varepsilon_1} = \Sigma_{ixx}^{\frac{1}{2}} \Gamma_{ij}^{\varepsilon_1} \Sigma_{ixx}^{-\frac{1}{2}} - \frac{\varepsilon_1}{2} I_d \text{, and } \Gamma_{ij}^{\varepsilon_1} = (\Sigma_{ixx}^{\frac{1}{2}} \Sigma_{jyy} \Sigma_{ixx}^{\frac{1}{2}} + \frac{\varepsilon_1^2}{4} I_d)^{\frac{1}{2}}.$$

Accordingly, the closed-form unique solution for  $dp_{ij}^*$  can be obtained as:

$$p_{ij}^*(\mathbf{x}, \mathbf{y}) = \mathcal{N}\left(\begin{bmatrix}\mathbf{x}\\\mathbf{y}\end{bmatrix} | \begin{bmatrix}\mathbf{m}_{i_x}\\\mathbf{m}_{i_y}\end{bmatrix}, \begin{bmatrix}\Sigma_{i_{xx}} & \Sigma_{ij}^{\varepsilon_1}\\\Sigma_{ij}^{\varepsilon_1} & \Sigma_{i_{yy}}\end{bmatrix}\right)$$

Therefore, the optimal mixture transport policy is itself a GMM, given by:

$$\pi(\mathbf{x}, \mathbf{y}) = \sum_{i,j}^{K} \omega_{ij} p_{ij}(\mathbf{x}, \mathbf{y}) = \sum_{i,j} \omega_{ij} \mathcal{N}\left(\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} | \begin{bmatrix} \mathbf{m}_{i_x} \\ \mathbf{m}_{i_y} \end{bmatrix}, \begin{bmatrix} \Sigma_{i_{xx}} & \Sigma_{i_{xy}} \\ \Sigma_{i_{xy}}^T & \Sigma_{i_{yy}} \end{bmatrix}\right). \tag{42}$$

# B W2-BENCHMARK TASK

### B.1 EXPERIMENT

 For the continuous Wasserstein-2 benchmark task, we adapted the experimental setup from the publicly available repository of Korotin et al. (2021). We evaluated all models across a range of dimensions (d) with corresponding training sample sizes (n). For each configuration, performance was assessed on a separate test set of 10,000 samples. To ensure statistical robustness, every experiment was repeated five times with different random initializations. These same settings were also used for the ablation study on the impact of noise.

The specific dimension and sample size pairs were as follows:

- for  $d \in \{2, 4\}$  with n = 10,000,
- for  $d \in \{8, 16, 32\}$  with n = 20,000,
- for  $d \in \{8, 16, 32\}$  with n = 20,000,
- for  $d \in \{8, 16, 32, 64\}$  with n = 20,000,
- for  $d \in \{128, 256\}$  with n = 50, 000.

## **B.2** Training configurations

We compared our OMT solver against several state-of-the-art baselines. We implemented EOT, ExNOT, and ProgOT using their official versions in the OTT-JAX toolbox (Cuturi et al., 2022), following the recommended settings from the tutorials.

- EOT: The entropy regularization was set to  $\varepsilon = 0.1$ , with a maximum of  $10^6$  iterations.
- ProgOT: We used the recommended schedulers with K=4 steps.
- ExNOT: We employed the recommended network architecture, using a 5-layer MLP for each potential function, with 128 nodes per hidden layer, and trained for a maximum of 10<sup>5</sup> iterations.
- ENOT: We employed the original code released by the authors for this benchmark, using the same configuration as reported.
- OMT (proposed): We set the number of source components to K=5 and target components to K=15. Gaussian mixture models (GMMs) were fitted to the source and target data using Python's scikit-learn, employing a full covariance structure for  $d \le 64$ . The model was trained for a maximum of  $10^5$  iterations with  $\varepsilon=0.01$  for both entropy regularizers.

Table 3: Transportation costs for different methods.

Method	Transportation Cost
OMT	$\sum_{i,j}^{K} \omega_{ij}^* \int_{\mathcal{X} \times \mathcal{Y}} \ \mathbf{x} - \mathbf{y}\ _2^2 dp_{ij}^*(\mathbf{x}, \mathbf{y})$
ExNOT	$\int_{\mathcal{X}} f^*(\mathbf{x}) d\nu_0(\mathbf{x}) + \int_{\mathcal{Y}} g^*(\mathbf{y}) d\nu_1(\mathbf{y})$
ENOT	$\int_{\mathcal{X}} \ \mathbf{x} - T^*(\mathbf{x})\ _2^2 \ d\nu_0(\mathbf{x})$
ProgOT	$\int_{\mathcal{X}\times\mathcal{Y}} \ \mathbf{x} - \mathbf{y}\ _2^2 d\pi^*(\mathbf{x}, \mathbf{y})$
ЕОТ	$\int_{\mathcal{X}\times\mathcal{Y}} \ \mathbf{x} - \mathbf{y}\ _2^2 d\pi^*(\mathbf{x}, \mathbf{y})$

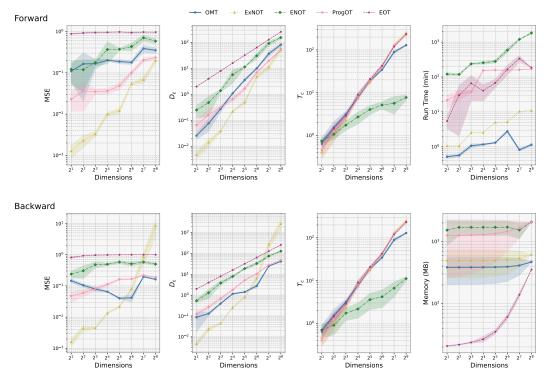


Figure 6: Comparison of OT solvers on forward and backward paths in Wasserstein-2 benchmark tasks Korotin et al. (2021). Results are computed on the test set with 10,000 samples and averaged over five random initializations. Transport cost  $T_c$  is defined in Table 3. Among the evaluated methods, PROGOT exhibits the highest computational cost among OT solvers.

# C SINGLE-CELL DATA ANALYSIS

## C.1 DATA AVAILABILITY

- sci-Plex3 data can be downloaded from NCBI GEO (#GSE139944).
- Mouse developmental data is available through Neuroscience Multi-omic Data Archive (NeMO), (RRID:SCR-016152). The 10x scRNA-seq dataset is available at https://assets.nemoarchive.org/dat-0oyried.
- Mouse ageing scRNA-seq is also available through NeMO, https://nemoarchive.org/, and can be accessed at https://assets.nemoarchive.org/dat-61kfys3.

## C.2 Preprocessing

For human scRNA-seq data, sci-Plex, we followed the same processing steps recommended in Cuturi et al. (2023). Genes which appear in less than 20cells, and cells with less that 20 gene expressed are excluded. Then we normalized gene expression, by first normalized to counts per million (CPM) and then transformed using the formula  $\log{(CPM+1)}$ . Then we whiten the data and apply PCA.

For the mouse scRNA-seq datasets, we preprocessed the raw count matrix. First, we performed library size normalization by converting counts to counts per million (CPM), followed by log-transformation. For feature selection, we chose a subset of highly variable genes combined with a list of known marker genes from the mouse brain atlas.

## C.3 OMT TRAINING

For each dataset, we first performed dimensionality reduction and then trained the OMT model on the resulting low-dimensional embeddings. The specific hyperparameters were tailored to each dataset.

sci-Plex Dataset. As previously described, we used PCA for dimensionality reduction. On the resulting PCA embeddings, we trained the OMT model with the number of source components set to  $K_s=3$  and target components to  $K_t=5$ . The entropy regularization parameter for both  $(\Omega,P)$  was set to 0.01.

Mouse Brain Datasets. For the mouse scRNA-seq data, we first trained a variational autoencoder (VAE) to learn a compressed cellular representation in a latent space of dimension  $d_z=10$ . We then trained the OMT model on these VAE embeddings. The number of components was set within a range of 5 to 25, with the specific value chosen based on the biological context; we typically used approximately twice the number of known cell types present at the analyzed timepoints.

## C.4 Additional results

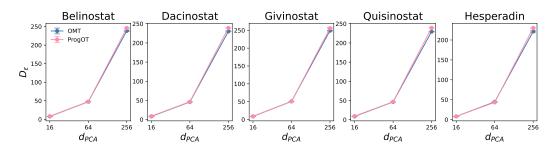


Figure 7: Average  $D_{\varepsilon} \downarrow$  values for the forward and backward OMT mappings compared to PROGOT sci-Plex dataset (Srivatsan et al., 2020). The results reported as the mean over 5 randomly initialized runs, with standard deviations.

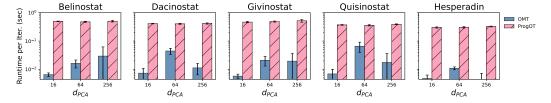


Figure 8: Comparison of the average runtime per iteration for OMT and ProgOT on the sci-Plex dataset as a function of latent space dimensionality  $d_{PCA}$ . The sharp decrease in OMT's runtime at d=256 arises from switching from a full covariance approximation to a diagonal structure. Results are shown as the mean  $\pm$  standard deviation over 5 randomly initialized runs.

#### D IMAGE TRANSLATION TASKS

### D.1 OMT TRAINING

Similar to our approach with scRNA-seq data, our method for image datasets involves a two-stage process. We first train a deep neural network to learn a low-dimensional representation of the images, and then train the OMT model on these resulting embeddings. For the MNIST dataset, we employed a convolutional VAE featuring a dual-decoder design, where each decoder reconstructs images for the source and target domains, respectively. The latent dimension for this network was set to  $d_z=10$ . The full architecture is detailed in Table 4.

For the CIFAR-10 dataset, we utilized the DoubleRessNet architecture, described in Table 5, with a latent space dimension of  $d_z=32$ . For all experiments on these datasets, the subsequent OMT model was trained using 10 components for both the source and target measure and  $\varepsilon=0.01$ . We found the number of components choice to be robust, as preliminary experiments with other values did not yield significant changes in the final results.

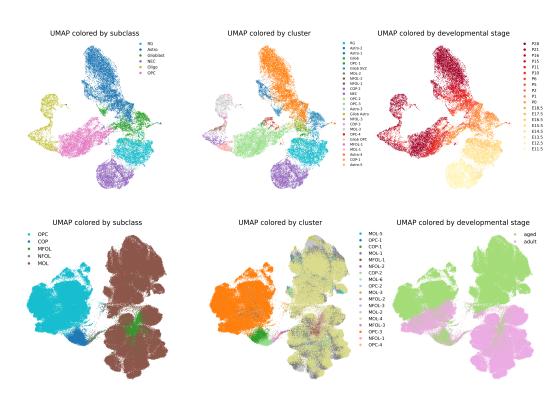


Figure 9: UMAP plots of the mouse scRNA-seq datasets. (Top) The developmental dataset from the visual cortex, including 32, 998 cells and 9, 900 HVGs (Gao et al., 2024). (Bottom) The aging dataset, consisting of 253, 468 cells and 9, 359 HVGs from six brain regions (Jin et al., 2025). For each dataset, the subfigures from left to right display the same embedding colored by cell subclass, cell type, and timepoint, respectively.

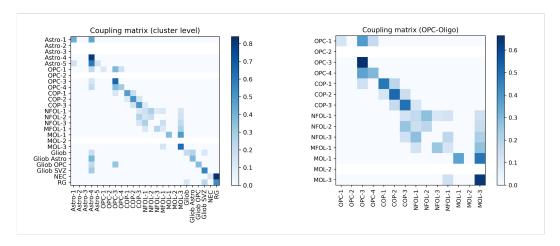


Figure 10: OMT-derived coupling matrices for the mouse developmental dataset. These heatmaps show the learned transport coupling by OMT. (Left) The cluster-level coupling matrix for all analyzed non-neuronal cell types. The strong diagonal indicates self-renewal or state maintenance at the developed stage, while off-diagonal values highlight developmental transitions, such as from glioblasts (Gliob) to astrocyte and oligodendrocyte lineages. (Right) A detailed view of the coupling matrix focused specifically on the OPC-Oligo lineage, illustrating the sequential maturation pathway from oligodendrocyte precursor cells (OPCs) to mature oligodendrocytes (MOLs). These findings should be considered in conjunction with the results reported in Figure 4.

Table 4: Architecture of the VAE for MNIST.

Layer	Configuration Details	Activation
-	Encoder	
Input Layer	2 x (Batch, 1, N, N) Image	-
Dropout		-
Conv2d (x2)	$1 \rightarrow 16$ channels, kernel=(5,5), stride=2	ReLU
Norm2d	16 features	-
Conv2d	$16 \rightarrow 32$ channels, kernel=(3,3), stride=2	ReLU
Norm2d	32 features	-
Conv2d	$32 \rightarrow 32$ channels, kernel=(3,3), stride=2	ReLU
Norm2d	32 features	-
Flatten	Reshapes feature map to (Batch, 128)	-
Linear	$128 \rightarrow 100 \text{ units}$	ReLU
	Decoders	
Linear	$d_z \rightarrow 100 \text{ units}$	ReLU
Linear	$100 \rightarrow 128$ units	ReLU
Unflatten	Reshapes to (Batch, 32, 2, 2)	-
ConvTranspose2d	$32 \rightarrow 32$ channels, kernel=(3,3), stride=2	ReLU
Norm2d	32 features	-
ConvTranspose2d	$32 \rightarrow 16$ channels, kernel=(5,5), stride=2	ReLU
Norm2d	16 features	-
ConvTranspose2d	$16 \rightarrow 1$ channel, kernel=(5,5), stride=2	ReLU
Norm2d	1 feature	-
ConvTranspose2d	$1 \rightarrow 1$ channel, kernel=(4,4)	ReLU

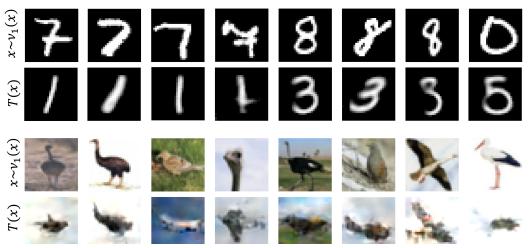


Figure 11: Additional results for unpaired image-to-image translation on the MNIST and CIFAR-10 datasets. For each dataset, the top row shows samples from the target distribution  $x \sim \nu_1$ , while the bottom row shows the corresponding transported images  $T_{\text{OMT}}^{\nu_1 \to \nu_0}$ , generated by OMT in the backward direction.

Table 5: Architecture of the DoubleRessNet for CIFAR-10.

Layer	Configuration Details	Activation		
Shared Encoder				
Input Layer	2 x (Batch, 3, H, W)	_		
Conv2d	$3 \rightarrow 32$ channels, kernel=9, stride=1	ReLU		
Norm2d	32 features	-		
Conv2d	$32 \rightarrow 64$ channels, kernel=3, stride=2	ReLU		
Norm2d	64 features	-		
Conv2d	$64 \rightarrow 128$ channels, kernel=3, stride=2	ReLU		
Norm2d	128 features	-		
Residual Block	4 stacked blocks (128 channels)	ReLU		
Flatten	Reshapes to (Batch, $H \times 64$ )	-		
Linear	$H \times 64 \rightarrow H \times 16$ units	ReLU		
Norm1d	$H \times 16$ features	-		
	Latent Space			
Linear	$H \times 16 \rightarrow H \times 4$ units	Tanh		
	Decoders			
Linear	$H \times 4 \rightarrow H \times 16$ units	-		
Norm1d	$H \times 16$ features	-		
Linear	$H \times 16 \rightarrow H \times 64$ units	ReLU		
Norm1d	$H \times 64$ features	-		
Unflatten	Reshapes to (Batch, $H$ , 8, 8)	-		
Upsample, Conv2d	$H \rightarrow 64$ channels, kernel=3, upsample=2	ReLU		
Norm2d	64 features	-		
Upsample, Conv2d	$64 \rightarrow 32$ channels, kernel=3, upsample=2	ReLU		
Norm2d	32 features	-		

Linear

Upsample, Conv2d  $32 \rightarrow 3$  channels, kernel=9, stride=1