
Unbalanced Optimal Transport meets Sliced-Wasserstein

Anonymous Authors¹

Abstract

Optimal transport (OT) has emerged as a powerful framework to compare probability measures, a fundamental task in many statistical and machine learning problems. Substantial advances have been made over the last decade in designing OT variants which are either computationally and statistically more efficient, or more robust to the measures/datasets to compare. Among them, sliced OT distances have been extensively used to mitigate optimal transport’s cubic algorithmic complexity and curse of dimensionality. In parallel, unbalanced OT was designed to allow comparisons of more general positive measures, while being more robust to outliers. In this paper, we propose to combine these two concepts, namely slicing and unbalanced OT, to develop a general framework for efficiently comparing positive measures. We propose two new loss functions based on the idea of slicing unbalanced OT, and study their induced topology and statistical properties. We then develop a fast Frank-Wolfe-type algorithm to compute these losses, and show that our methodology is modular as it encompasses and extends prior related work. We finally conduct an empirical analysis of our loss functions and methodology on both synthetic and real datasets, to illustrate their relevance and applicability.

1. Introduction

Positive measures are ubiquitous in various fields, including data sciences and machine learning (ML) where they commonly serve as data representations. A common example is the density fitting task, which arises in generative modeling (Arjovsky et al., 2017; De Bortoli et al., 2021): the observed samples can be represented as a discrete positive measure α and the goal is to find a parametric measure β_η

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the Workshop on New Frontiers in Learning, Control, and Dynamical Systems at the International Conference on Machine Learning (ICML). Do not distribute.

which fits the best α . This can be achieved by training a model that minimizes a loss function over η , usually defined as a distance between α and β_η . Therefore, it is important to choose a meaningful discrepancy with desirable statistical, robustness and computational properties. In particular, some settings require comparing arbitrary positive measures, *i.e.* measures whose total mass can have an arbitrary value, as opposed to probability distributions, whose total mass is equal to 1. In cell biology (Schiebinger et al., 2019), for example, measures are used to represent and compare gene expressions of cell populations, and the total mass represents the population size.

(Unbalanced) Optimal Transport. Optimal transport has been chosen as a loss function in various ML applications. OT defines a distance between two positive measures of same mass α and β (*i.e.* $m(\alpha) = m(\beta)$) by moving the mass of α toward the mass of β with least possible effort. The mass equality can nevertheless be hindering by imposing a normalization of α and β to enforce $m(\alpha) = m(\beta)$, which is potentially spurious and makes the problem less interpretable. In recent years, OT has then been extended to settings where measures have different masses, leading to the *unbalanced OT* (UOT) framework (Liero et al., 2018; Kondratyev et al., 2016; Chizat et al., 2018b). An appealing outcome of this new OT variant is its robustness to outliers which is achieved by discarding them before transporting α to β . UOT has been useful for many theoretical and practical applications, *e.g.* theory of deep learning (Chizat & Bach, 2018; Rotskoff et al., 2019), biology (Schiebinger et al., 2019; Demetci et al., 2022) and domain adaptation (FAtlas et al., 2021). We refer to (Séjourné et al., 2022a) for an extensive survey of UOT. Computing OT requires to solve a linear program whose complexity is in $\mathcal{O}(n^3 \log n)$. Besides, accurately estimating OT distances through empirical distributions is challenging as OT suffers from the curse of dimension (Dudley, 1969). A common workaround is to rely on OT variants with lower complexities and better statistical properties. Among the most popular, we can list entropic OT (Cuturi, 2013), minibatch OT (FAtlas et al., 2020) and sliced OT (Radon, 2005; Bonneel et al., 2015). In this paper, we will focus on the latter.

Slicing (U)OT and related work. Sliced OT leverages the OT 1D closed-form solution to define a new cost. It averages the OT cost between projections of (α, β) on 1D

subspaces of \mathbb{R}^d . For 1D data, the OT solution can be computed through a sort algorithm, leading to an appealing $\mathcal{O}(n \log(n))$ complexity (Peyré et al., 2019). Furthermore, it has been shown to lift useful topological and statistical properties of OT from 1-dimensional to multi-dimensional settings (Bayraktar & Guo, 2021; Nadjahi et al., 2020; Goldfeld & Greenewald, 2021). It therefore helps to mitigate the curse of dimensionality making SOT-based algorithms theoretically-grounded, statistically efficient and efficiently solvable even on large-scale settings. These appealing properties motivated the development of several variants and generalizations, e.g. to different types or distributions of projections (Kolouri et al., 2019; Deshpande et al., 2019; Nguyen et al., 2020; Ohana et al., 2023) and non-Euclidean data (Bonet et al., 2023a; 2022a; 2023b). The slicing operation has also been applied to partial OT (Bonneel & Coeurjolly, 2019; Bai et al., 2022; Sato et al., 2020), a particular case of UOT, in order to speed up comparisons of unnormalized measures at large scale. However, while (sliced) partial OT allows to compare measures with different masses, it assumes that each input measure is discrete and supported on points that all share the same mass (typically 1). In contrast, the Gaussian-Hellinger-Kantorovich (GHK) distance (Liero et al., 2018), another popular formulation of UOT, allows to compare measures with different masses *and* supported on points with varying masses, and has not been studied jointly with slicing.

Contributions. This paper presents the first general framework combining UOT and slicing. Our main contribution is the introduction of two novel sliced variants of UOT, respectively called *Sliced UOT* (SUOT) and *Unbalanced Sliced OT* (USOT). SUOT and USOT both leverage one-dimensional projections and the newly-proposed implementation of UOT in 1D (Séjourné et al., 2022b), but differ in the penalization used to relax the constraint on the equality of masses: USOT essentially performs a global reweighting of the inputs measures (α, β) , while SUOT reweights each projection of (α, β) . Our work builds upon the Frank-Wolfe-type method (Frank & Wolfe, 1956) recently proposed in (Séjourné et al., 2022b) to efficiently compute GHK between univariate measures, an instance of UOT which has not yet been combined with slicing. We derive the associated theoretical properties, along with the corresponding fast and GPU-friendly algorithms. We demonstrate its versatility and efficiency on challenging experiments, where slicing is considered on a non-Euclidean hyperbolic manifold, as a similarity measure for document classification, or for computing barycenters of geoclimatic data.

Outline. In Section 2, we provide background knowledge on UOT and sliced OT (SOT). In Section 3, we define our two new loss functions (SUOT and USOT) and prove their metric, topological, statistical and duality properties in wide generality. We then detail in Section 4 the numerical imple-

mentation of SUOT and USOT based on the Frank-Wolfe algorithm. We investigate their empirical performance on hyperbolic and geophysical data as well as document classification in Section 5.

2. Background

Unbalanced Optimal Transport. We denote by $\mathcal{M}_+(\mathbb{R}^d)$ the set of all positive Radon measures on \mathbb{R}^d . For any $\alpha \in \mathcal{M}_+(\mathbb{R}^d)$, $\text{supp}(\alpha)$ is the support of α and $m(\alpha) = \int_{\mathbb{R}^d} d\alpha(x)$ the mass of α . We recall the standard formulation of unbalanced OT (Liero et al., 2018), which uses φ -divergences for regularization.

Definition 2.1. (Unbalanced OT) Let $\alpha, \beta \in \mathcal{M}_+(\mathbb{R}^d)$. Let $\varphi : \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$ be an *entropy function*, i.e. φ is convex, lower semicontinuous, $\text{dom}(\varphi) \triangleq \{x \in \mathbb{R}, \varphi(x) < +\infty\} \subset [0, +\infty)$ and $\varphi(1) = 0$. Denote $\varphi'_\infty \triangleq \lim_{x \rightarrow +\infty} \varphi(x)/x$. The φ -divergence between α and β is defined as,

$$D_\varphi(\alpha|\beta) \triangleq \int_{\mathbb{R}^d} \varphi\left(\frac{d\alpha}{d\beta}(x)\right) d\beta(x) + \varphi'_\infty \int_{\mathbb{R}^d} d\alpha^\perp(x), \quad (1)$$

where α^\perp is defined as $\alpha = (d\alpha/d\beta)\beta + \alpha^\perp$. Given two entropy functions (φ_1, φ_2) and a cost $C_d : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$, the unbalanced OT problem between α and β reads

$$\text{UOT}(\alpha, \beta) \triangleq \inf_{\pi \in \mathcal{M}_+(\mathbb{R}^d \times \mathbb{R}^d)} \int C_d(x, y) d\pi(x, y) + D_{\varphi_1}(\pi_1|\alpha) + D_{\varphi_2}(\pi_2|\beta), \quad (2)$$

where (π_1, π_2) denote the marginal distributions of π .

When $\varphi_1 = \varphi_2$ and $\varphi_1(x) = 0$ for $x = 1$, $\varphi_1(x) = +\infty$ otherwise, (2) boils down to the Kantorovich formulation of OT (or *balanced OT*), which we denote by $\text{OT}(\alpha, \beta)$. Indeed, in that case, $D_{\varphi_1}(\pi_1|\alpha) = D_{\varphi_2}(\pi_2|\beta) = 0$ if $\pi_1 = \alpha$ and $\pi_2 = \beta$, $D_{\varphi_1}(\pi_1|\alpha) = D_{\varphi_2}(\pi_2|\beta) = +\infty$ otherwise.

Under suitable choices of entropy functions (φ_1, φ_2) , $\text{UOT}(\alpha, \beta)$ allows to compare α and β even when $m(\alpha) \neq m(\beta)$ and can discard outliers, which makes it more robust than $\text{OT}(\alpha, \beta)$. Two common choices are $\varphi(x) = \rho|x - 1|$ and $\varphi(x) = \rho(x \log(x) - x + 1)$, where $\rho > 0$ is a *characteristic radius* w.r.t. C_d . They respectively correspond to $D_\varphi = \rho\text{TV}$ (total variation distance (Chizat et al., 2018a)) and $D_\varphi = \rho\text{KL}$ (Kullback-Leibler divergence).

The UOT problem has been shown to admit an equivalent formulation obtained by deriving the dual of (2) and proving strong duality. Based on Proposition 2.2, computing $\text{UOT}(\alpha, \beta)$ consists in optimizing a pair of continuous functions (f, g) .

Proposition 2.2. (*Liero et al., 2018, Corollary 4.12*) The UOT problem (2) can equivalently be written as

$$\text{UOT}(\alpha, \beta) = \sup_{f \oplus g \leq C_d} \int \varphi_1^\circ(f(x)) d\alpha(x) + \int \varphi_2^\circ(g(y)) d\beta(y), \quad (3)$$

where for $i \in \{1, 2\}$, $\varphi_i^\circ(x) \triangleq -\varphi_i^*(-x)$ with $\varphi_i^*(x) \triangleq \sup_{y \geq 0} xy - \varphi_i(y)$ the Legendre transform of φ_i , and $f \oplus g \leq C_d$ means that for $(x, y) \sim \alpha \otimes \beta$, $f(x) + g(y) \leq C_d(x, y)$.

In this paper, we mainly focus on the *GHK setting*, both theoretically and computationally. It corresponds to (2) with $C_d(x, y) = \|x - y\|^2$, $D_{\varphi_i} = \rho_i \text{KL}$, leading to $\varphi_i^\circ(x) = \rho_i(1 - e^{-x/\rho_i})$. UOT(α, β) is known to be computationally intensive (Pham et al., 2020), thus motivating the development of methods that can scale to dimensions and sample sizes encountered in ML applications.

Sliced Optimal Transport. Among the many workarounds that have been proposed to overcome the OT computational bottleneck (Peyré et al., 2019), Sliced OT (Rabin et al., 2012) has attracted a lot of attention due to its computational benefits and theoretical guarantees. We define it below.

Definition 2.3 (Sliced OT). Let $\mathbb{S}^{d-1} \triangleq \{\theta \in \mathbb{R}^d : \|\theta\| = 1\}$ be the unit sphere in \mathbb{R}^d . For $\theta \in \mathbb{S}^{d-1}$, denote by $\theta^* : \mathbb{R}^d \rightarrow \mathbb{R}$ the linear map such that for $x \in \mathbb{R}^d$, $\theta^*(x) \triangleq \langle \theta, x \rangle$. Let σ be the uniform probability over \mathbb{S}^{d-1} . For $\alpha, \beta \in \mathcal{M}_+(\mathbb{R}^d)$, the *Sliced OT* problem reads

$$\text{SOT}(\alpha, \beta) \triangleq \int_{\mathbb{S}^{d-1}} \text{OT}(\theta_\#^* \alpha, \theta_\#^* \beta) d\sigma(\theta), \quad (4)$$

where for any measurable function f and $\xi \in \mathcal{M}_+(\mathbb{R}^d)$, $f_\# \xi$ is the *push-forward measure* of ξ by f , i.e. for any measurable set $A \subset \mathbb{R}$, $f_\# \xi(A) \triangleq \xi(f^{-1}(A))$, $f^{-1}(A) \triangleq \{x \in \mathbb{R}^d : f(x) \in A\}$.

Note that $\theta_\#^* \alpha, \theta_\#^* \beta$ are two measures supported on \mathbb{R} , therefore $\text{OT}(\theta_\#^* \mu, \theta_\#^* \nu)$ is defined in terms of a cost function $C_1 : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$. Since OT between univariate measures can be efficiently computed, SOT(α, β) can provide significant computational advantages over OT(α, β) in large-scale settings. In practice, if α and β are discrete measures supported on $\{x_i\}_{i=1}^n$ and $\{y_i\}_{i=1}^n$ respectively, the standard procedure for approximating SOT(α, β) consists in (i) sampling m i.i.d. samples $\{\theta_j\}_{j=1}^m$ from σ , (ii) computing $\text{OT}((\theta_j^*)_\# \alpha, (\theta_j^*)_\# \beta)$, $j = 1, \dots, m$. Computing OT between univariate discrete measures amounts to sorting (Peyré et al., 2019, Section 2.6), thus step (ii) involves $\mathcal{O}(n \log n)$ operations for each θ_j .

SOT(α, β) is defined in terms of the Kantorovich formulation of OT, hence inherits the following drawbacks: SOT(α, β) $< +\infty$ only when $m(\alpha) = m(\beta)$, and may not provide meaningful comparisons in presence of outliers. To

overcome such limitations, prior work have proposed sliced versions of partial OT (Bonneel & Coeurjolly, 2019; Bai et al., 2022), a particular instance of UOT. However, their contributions only apply to measures whose samples have constant mass. We generalize their line of work in the next section.

3. Sliced Unbalanced OT and Unbalanced Sliced OT: Theoretical Analysis

We propose two strategies to make unbalanced OT scalable, by leveraging sliced OT. We formulate two loss functions (Definition 3.1), then study their theoretical properties and discuss their implications.

Definition 3.1. Let $\alpha, \beta \in \mathcal{M}_+(\mathbb{R}^d)$. The **Sliced Unbalanced OT** loss (SUOT) and the **Unbalanced Sliced OT** loss (USOT) between α and β are defined as,

$$\text{SUOT}(\alpha, \beta) \triangleq \int_{\mathbb{S}^{d-1}} \text{UOT}(\theta_\#^* \alpha, \theta_\#^* \beta) d\sigma(\theta), \quad (5)$$

$$\text{USOT}(\alpha, \beta) \triangleq \inf_{(\pi_1, \pi_2) \in \mathcal{M}_+(\mathbb{R}^d) \times \mathcal{M}_+(\mathbb{R}^d)} \text{SOT}(\pi_1, \pi_2) + D_{\varphi_1}(\pi_1 | \alpha) + D_{\varphi_2}(\pi_2 | \beta). \quad (6)$$

SUOT(α, β) compares α and β by solving the UOT problem between $\theta_\#^* \alpha$ and $\theta_\#^* \beta$ for $\theta \sim \sigma$. Note that SUOT extends the sliced partial OT problem (Bonneel & Coeurjolly, 2019; Bai et al., 2022) (where $D_{\varphi_i} = \rho_i \text{TV}$) by allowing the use of arbitrary φ -divergences. On the other hand, USOT is a completely novel approach and stems from the following property on UOT (Liero et al., 2018, Equations (4.21)): $\text{UOT}(\alpha, \beta) = \inf_{(\pi_1, \pi_2) \in \mathcal{M}_+(\mathbb{R}^d)^2} \text{OT}(\pi_1, \pi_2) + D_{\varphi_1}(\pi_1 | \alpha) + D_{\varphi_2}(\pi_2 | \beta)$.

SUOT vs. USOT. As outlined in Definition 3.1, SUOT and USOT differ in how the transportation problem is penalized: SUOT(α, β) regularizes the marginals of π_θ for $\theta \sim \sigma$ where π_θ denotes the solution of $\text{UOT}(\theta_\#^* \alpha, \theta_\#^* \beta)$, while USOT(α, β) operates a geometric normalization directly on (α, β) . We illustrate this difference on the following practical setting: we consider $(\alpha, \beta) \in \mathcal{M}_+(\mathbb{R}^2)$ where α is polluted with some outliers, and we compute SUOT(α, β) and USOT(α, β). We plot the input measures and the sampled projections $\{\theta_k\}_k$ (Figure 1, left), the marginals of π_{θ_k} for SUOT and the marginals of $(\theta_k)_\#^* \pi$ for USOT (Figure 1, right). As expected, SUOT marginals change for each θ_k . We also observe that the source outliers have successfully been removed for any θ when using USOT, while they may still appear with SUOT (e.g. for $\theta = 120^\circ$): this is a direct consequence of the penalization terms D_{φ_i} in USOT, which operate on (α, β) rather than on their projections.

Theoretical analysis. In the rest of this section, we prove a set of theoretical properties of SUOT and USOT. All proofs

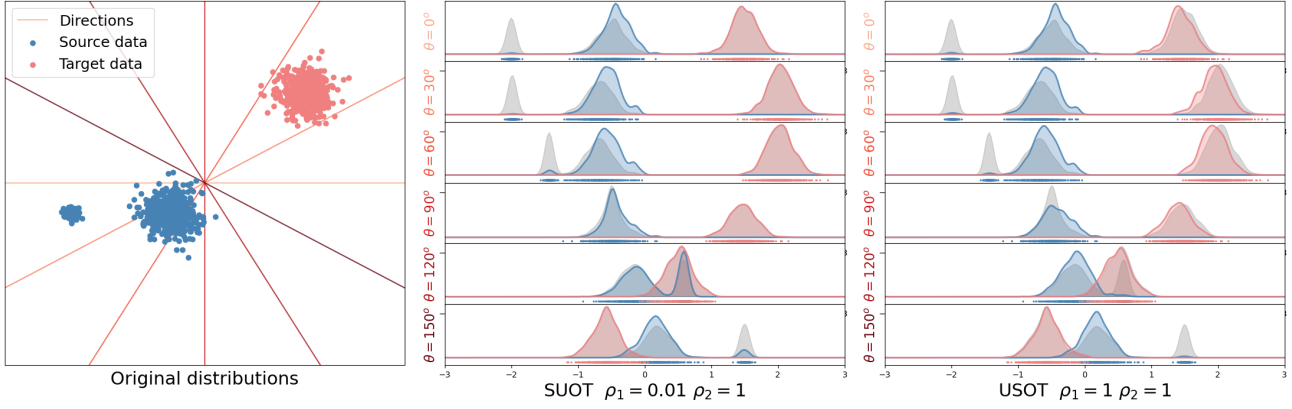


Figure 1: **Toy illustration** on the behaviors of SUOT and USOT. (left) Original 2D samples and slices used for illustration. KDE density estimations of the projected samples: grey, original distributions, colored, distributions reweighed by SUOT (center), and reweighed by USOT (right).

are provided in Appendix A. We first identify the conditions on the cost C_1 and entropies φ_1, φ_2 under which the infimum is attained in $\text{UOT}(\theta_{\sharp}^* \alpha, \theta_{\sharp}^* \beta)$ for $\theta \in \mathbb{S}^{d-1}$ and in $\text{USOT}(\alpha, \beta)$: the formal statement is given in Appendix A. We also show that these optimization problems are convex, both SUOT and USOT are jointly convex w.r.t. their input measures, and that strong duality holds (Theorem 3.7).

Next, we prove that both SUOT and USOT preserve some topological properties of UOT, starting with the metric axioms as stated in the next proposition.

Proposition 3.2. (Metric properties) (i) Suppose UOT is non-negative, symmetric and/or definite on $\mathcal{M}_+(\mathbb{R}) \times \mathcal{M}_+(\mathbb{R})$. Then, SUOT is respectively non-negative, symmetric and/or definite on $\mathcal{M}_+(\mathbb{R}^d) \times \mathcal{M}_+(\mathbb{R}^d)$. If there exists $p \in [1, +\infty)$ s.t. for any $(\alpha, \beta, \gamma) \in \mathcal{M}_+(\mathbb{R})$, $\text{UOT}^{1/p}(\alpha, \beta) \leq \text{UOT}^{1/p}(\alpha, \gamma) + \text{UOT}^{1/p}(\gamma, \beta)$, then $\text{SUOT}^{1/p}(\alpha, \beta) \leq \text{SUOT}^{1/p}(\alpha, \gamma) + \text{SUOT}^{1/p}(\gamma, \beta)$.

(ii) For $\alpha, \beta \in \mathcal{M}_+(\mathbb{R}^d)$, $\text{USOT}(\alpha, \beta) \geq 0$. If $\varphi_1 = \varphi_2$, USOT is symmetric. If $D_{\varphi_1}, D_{\varphi_2}$ are definite, so is USOT.

By Proposition 3.2(i), establishing the metric axioms of UOT between univariate measures (e.g., as detailed in (Séjourné et al., 2022a, Section 3.3.1)) suffices to prove the metric axioms of SUOT between multivariate measures. Since e.g. GHK (Liero et al., 2018, Theorem 7.25) is a metric for $p = 2$, then so is the associated SUOT.

In our next theorem, we show that SUOT, USOT and UOT are equivalent, under certain assumptions on the entropies (φ_1, φ_2) , cost functions, and input measures (α, β) .

Theorem 3.3. (Equivalence of SUOT, USOT, UOT) Let $X \subset \mathbb{R}^d$ be a compact set with radius R . Let $p \in [1, +\infty)$. Assume $C_1(x, y) = |x - y|^p$, $C_d(x, y) = \|x - y\|^p$, $D_{\varphi_1} =$

$D_{\varphi_2} = \rho \text{KL}$. Then, for $\alpha, \beta \in \mathcal{M}_+(X)$,

$$\text{SUOT}(\alpha, \beta) \leq \text{USOT}(\alpha, \beta) \leq \text{UOT}(\alpha, \beta), \text{ and} \quad (7)$$

$$\text{UOT}(\alpha, \beta) \leq c(m(\alpha), m(\beta), \rho, R) \text{SUOT}(\alpha, \beta)^{1/(d+1)}, \quad (8)$$

where $c(m(\alpha), m(\beta), \rho, R)$ is constant depending on $m(\alpha), m(\beta), \rho, R$, which is non-decreasing in $m(\alpha)$ and $m(\beta)$. Additionally, assume there exists $M > 0$ s.t. $m(\alpha) \leq M, m(\beta) \leq M$. Then, $c(m(\alpha), m(\beta), \rho, R)$ no longer depends on $m(\alpha), m(\beta)$, which proves the equivalence of SUOT, USOT and UOT.

Theorem 3.3 is an application of a more general result, which we derive in the appendix. In particular, we show that the first two inequalities in (7) hold under milder assumptions on φ_1, φ_2 and C_1, C_d . The equivalence of SUOT, USOT and UOT is useful to prove that SUOT and USOT metrize the weak* convergence when UOT does, e.g. in the GHK setting (Liero et al., 2018, Theorem 7.25). Before formally stating this result, we recall that a sequence of positive measures $(\alpha_n)_{n \in \mathbb{N}^*}$ converges weakly to $\alpha \in \mathcal{M}_+(\mathbb{R}^d)$ (denoted $\alpha_n \rightharpoonup \alpha$) if for any continuous $f : \mathbb{R}^d \rightarrow \mathbb{R}$, $\lim_{n \rightarrow +\infty} \int f d\alpha_n = \int f d\alpha$.

Theorem 3.4. (Weak* metrization) Assume $D_{\varphi_1} = D_{\varphi_2} = \rho \text{KL}$. Let $p \in [1, +\infty)$ and consider $C_1(x, y) = |x - y|^p$, $C_d(x, y) = \|x - y\|^p$. Let (α_n) be a sequence of measures in $\mathcal{M}_+(X)$ and $\alpha \in \mathcal{M}_+(X)$, where $X \subset \mathbb{R}^d$ is compact with radius $R > 0$. Then, $\alpha_n \rightharpoonup \alpha \Leftrightarrow \lim_{n \rightarrow +\infty} \text{SUOT}(\alpha_n, \alpha) = 0 \Leftrightarrow \lim_{n \rightarrow +\infty} \text{USOT}(\alpha_n, \alpha) = 0$.

The metrization of weak* convergence is an important property when comparing measures. For instance, it can be leveraged to justify the well-posedness of approximating an unbalanced Wasserstein gradient flow (Ambrosio et al., 2005) using SUOT, as done in (Bonet et al., 2022b; Candau-Tilh,

2020) for SOT. Unbalanced Wasserstein gradient flows have been a key tool in deep learning theory, e.g. to prove global convergence of 1-hidden layer neural networks (Chizat & Bach, 2018; Rotskoff et al., 2019).

We now specialize some metric and topological properties to sliced partial OT, a particular case of SUOT. Theorem 3.5 shows that our framework encompasses existing approaches and more importantly, helps complement their analysis (Bonneel & Coeurjolly, 2019; Bai et al., 2022).

Theorem 3.5. (Properties of Sliced Partial OT) Assume $C_1(x, y) = |x - y|$ and $D_{\varphi_1} = D_{\varphi_2} = \rho \text{TV}$. Then, USOT satisfies the triangle inequality. Additionally, for any $(\alpha, \beta) \in \mathcal{M}_+(\mathbb{X})$ where $\mathbb{X} \subset \mathbb{R}^d$ is compact with radius R , $\text{UOT}(\alpha, \beta) \leq c(\rho, R) \text{SUOT}(\alpha, \beta)^{1/(d+1)}$, and USOT and SUOT both metrize the weak* convergence.

We move on to the statistical properties and prove that SUOT offers important statistical benefits, as it lifts the *sample complexity* of UOT from one-dimensional setting to multi-dimensional ones. In what follows, for any $\alpha \in \mathcal{M}_+(\mathbb{R}^d)$, we use $\hat{\alpha}_n$ to denote the empirical approximation of α over $n \geq 1$ i.i.d. samples, i.e. $\hat{\alpha}_n = \frac{1}{n} \sum_{i=1}^n \delta_{Z_i}$, $Z_i \sim \alpha$.

Theorem 3.6. (Sample complexity) If for $\mu, \nu \in \mathcal{M}_+(\mathbb{R})$, $\mathbb{E}|\text{UOT}(\mu, \nu) - \text{UOT}(\hat{\mu}_n, \hat{\nu}_n)| \leq \kappa(n)$, then for $\alpha, \beta \in \mathcal{M}_+(\mathbb{R}^d)$, $\mathbb{E}|\text{SUOT}(\alpha, \beta) - \text{SUOT}(\hat{\alpha}_n, \hat{\beta}_n)| \leq \kappa(n)$.

If for $\mu, \nu \in \mathcal{M}_+(\mathbb{R})$, $\mathbb{E}|\text{UOT}(\mu, \hat{\mu}_n)| \leq \xi(n)$, then for $\alpha, \beta \in \mathcal{M}_+(\mathbb{R}^d)$, $\mathbb{E}|\text{SUOT}(\alpha, \hat{\alpha}_n)| \leq \xi(n)$.

Theorem 3.6 means that SUOT enjoys a *dimension-free* sample complexity, even when comparing multivariate measures: this advantage is recurrent of sliced divergences (Nadjahi et al., 2020) and further motivates their use on high-dimensional settings. The sample complexity rates $\kappa(n)$ or $\xi(n)$ can be deduced from the literature on UOT for univariate measures, for example we refer to (Vacher & Vialard, 2022) for the GHK setting. Establishing the statistical properties of USOT may require extending (Nietert et al., 2022): we leave this question for future work.

We conclude this section by deriving the dual formulations of SUOT, USOT and proving that strong duality holds. We will consider that σ is approximated with $\hat{\sigma}_K = \frac{1}{K} \sum_{k=1}^K \delta_{\theta_k}$, $\theta_k \sim \sigma$. This corresponds to the routine case in practice, as practitioners usually resort to a Monte Carlo approximation to estimate the expectation w.r.t. σ defining sliced OT.

Theorem 3.7. (Strong duality) For $i \in \{1, 2\}$, let φ_i be an entropy function s.t. $\text{dom}(\varphi_i^*) \cap \mathbb{R}_-$ is non-empty, and either $0 \in \text{dom}(\varphi_i)$ or $m(\alpha), m(\beta) \in \text{dom}(\varphi_i)$. Define $\mathcal{E} \triangleq \{\forall \theta \in \text{supp}(\sigma_K), f_\theta \oplus g_\theta \leq C_1\}$. Let $f_{\text{avg}} \triangleq \int_{\mathbb{S}^{d-1}} f_\theta d\hat{\sigma}_K(\theta)$, $g_{\text{avg}} \triangleq \int_{\mathbb{S}^{d-1}} g_\theta d\hat{\sigma}_K(\theta)$.

Then, SUOT (5) and USOT (6) can be equivalently written

for $\alpha, \beta \in \mathcal{M}_+(\mathbb{R}^d)$ as,

$$\begin{aligned} \text{SUOT}(\alpha, \beta) &= \sup_{(f_\theta), (g_\theta) \in \mathcal{E}} \int_{\mathbb{S}^{d-1}} \left(\int \varphi_1^\circ(f_\theta \circ \theta^*(x)) d\alpha(x) \right. \\ &\quad \left. + \int \varphi_2^\circ(g_\theta \circ \theta^*(y)) d\beta(y) \right) d\hat{\sigma}_K(\theta) \end{aligned} \quad (9)$$

$$\begin{aligned} \text{USOT}(\alpha, \beta) &= \sup_{(f_\theta), (g_\theta) \in \mathcal{E}} \int \varphi_1^\circ(f_{\text{avg}} \circ \theta^*(x)) d\alpha(x) \\ &\quad + \int \varphi_2^\circ(g_{\text{avg}} \circ \theta^*(y)) d\beta(y) \end{aligned} \quad (10)$$

We conjecture that strong duality also holds for σ Lebesgue over \mathbb{S}^{d-1} , and discuss this aspect in Appendix A. Theorem 3.7 has important practical implications, since it justifies the Frank-Wolfe-type algorithms that we develop in Section 4 to compute SUOT and USOT in practice.

4. Computing SUOT and USOT with Frank-Wolfe algorithms

We propose two algorithms by leveraging our strong duality result (Theorem 3.7) along with a Frank-Wolfe algorithm (FW, Frank & Wolfe (1956)) introduced in (Séjourné et al., 2022b) to optimize UOT dual (3). Our methods, summarized in Algorithms 1 and 2, can be applied for smooth $D_{\varphi_1}, D_{\varphi_2}$: this is satisfied for GHK (where $D_{\varphi_i} = \rho_i \text{KL}$), but not for sliced partial OT (where $D_{\varphi_i} = \rho_i \text{TV}$, Bai et al. (2022)). We refer to Appendix B for more technical details on our methodology and its theoretical justification.

FW is an iterative procedure which aims at maximizing a functional \mathcal{H} over a compact convex set \mathcal{E} , by maximizing a linear approximation $\nabla \mathcal{H}$: given iterate x^t , FW solves the linear oracle $r^{t+1} \in \arg \max_{r \in \mathcal{E}} \langle \nabla \mathcal{H}(x^t), r \rangle$ and performs a convex update $x^{t+1} = (1 - \gamma_{t+1})x^t + \gamma_{t+1}r^{t+1}$, with γ_{t+1} typically chosen as $\gamma_{t+1} = 2/(2 + t + 1)$. We call this step `FWStep` in our pseudo-code. When applied in (Séjourné et al., 2022b) to compute UOT(α, β) dual (3), `FWStep` updates (f_t, g_t) s.t. $f_t \oplus g_t \leq C_d$, and the linear oracle is the balanced dual of OT(α_t, β_t) where (α_t, β_t) are normalized versions of (α, β) . Updating (α_t, β_t) involves (f_t, g_t) and $\rho = (\rho_1, \rho_2)$: we refer to this routine as `NORM`($\alpha, \beta, f_t, g_t, \rho$) and report the closed-form updates in Appendix B. In other words, computing UOT amounts to solve a sequence of OT problems, which can efficiently be done for univariate measures (Séjourné et al., 2022b).

Analogously to UOT, and by Theorem 3.7, we propose to compute SUOT(α, β) and USOT(α, β) based on their dual forms. FW iterates consists in solving a sequence of sliced OT problems. We derive the updates for the `FWStep`

Algorithm 1 – SUOT

Input: $\alpha, \beta, F, (\theta_k)_{k=1}^K, \rho = (\rho_1, \rho_2)$
Output: SUOT(α, β), (f_θ, g_θ)
 $(f_\theta, g_\theta) \leftarrow (0, 0)$
for $t = 0, 1, \dots, F - 1$, **for** $\theta \in (\theta_k)_{k=1}^K$ **do**
 $(\alpha_\theta, \beta_\theta) \leftarrow \text{Norm}(\theta_\#^* \alpha, \theta_\#^* \beta, f_\theta, g_\theta, \rho)$
 $(r_\theta, s_\theta) \leftarrow \text{SlicedDual}(\alpha_\theta, \beta_\theta)$
 $(f_\theta, g_\theta) \leftarrow \text{FWStep}(f_\theta, g_\theta, r_\theta, s_\theta, \gamma t)$
end for
 Return SUOT(α, β), (f_θ, g_θ) as in (9)

Algorithm 2 – USOT

Input: $\alpha, \beta, F, (\theta_k)_{k=1}^K, \rho = (\rho_1, \rho_2)$
Output: USOT(α, β), (f_{avg}, g_{avg})
 $(f_\theta, g_\theta, f_{avg}, g_{avg}) \leftarrow (0, 0, 0, 0)$
for $t = 0, 1, \dots, F - 1$, **for** $\theta \in (\theta_k)_{k=1}^K$ **do**
 $(\pi_1, \pi_2) \leftarrow \text{Norm}(\alpha, \beta, f_{avg}, g_{avg}, \rho)$
 $(r_\theta, s_\theta) \leftarrow \text{SlicedDual}(\theta_\#^* \pi_1, \theta_\#^* \pi_2)$
 $r_{avg}, s_{avg} \leftarrow \text{AvgPot}(r_\theta), \text{AvgPot}(s_\theta)$
 $(f_{avg}, g_{avg}) \leftarrow \text{FWStep}(f_{avg}, g_{avg}, r_{avg}, s_{avg}, \gamma t)$
end for
 Return USOT(α, β), (f_{avg}, g_{avg}) as in (10)

tailored for SUOT and USOT in Appendix B, and re-use the aforementioned `Norm` routine. For USOT, we implement an additional routine called `AvgPot`((f_θ)) to compute $\int f_\theta d\hat{\sigma}_K(\theta)$ given the sliced potentials (f_θ).

A crucial difference is the need of SOT dual potentials (r_θ, s_θ) to call `Norm`. However, past implementations only return the loss SOT(α, β) for e.g. training models (Deshpande et al., 2019; Nguyen et al., 2020). Thus we designed two novel (GPU) implementations in PyTorch (Paszke et al., 2019) which return them. The first one leverages that the gradient of OT(α, β) w.r.t. (α, β) are optimal (f, g), which allows to backpropagate OT($\theta_\#^* \alpha, \theta_\#^* \beta$) w.r.t. (α, β) to obtain (r_θ, s_θ). The second implementation computes them in parallel on GPUs using their closed form, which to the best of our knowledge is a new sliced algorithm. We call `SlicedDual`($\theta_\#^* \alpha, \theta_\#^* \beta$) the step returning optimal (r_θ, s_θ) solving OT($\theta_\#^* \alpha, \theta_\#^* \beta$) for all θ . Both routines preserve the $O(N \log N)$ per slice time complexity and can be adapted to any SOT variant. Thus, our FW approach is modular in that one can reuse the SOT literature. We illustrate this by computing USOT between distributions in the hyperbolic Poincaré disk. (Figure 2).

Algorithmic complexity. FW algorithms and its variants have been widely studied theoretically. Computing `SlicedDual` has a complexity $O(KN \log N)$, where N is the number of samples, and K the number of projections of $\hat{\sigma}_K$. The overall complexity of SUOT and USOT is thus $O(FKN \log N)$, where F is the number of FW iterations needed to reach convergence. Our setting falls under the assumptions of (Lacoste-Julien & Jaggi, 2015, Theorem 8), thus ensuring fast convergence of our methods. We plot in Appendix B empirical evidence that a few iterations of FW ($F \leq 20$) suffice to reach numerical precision.

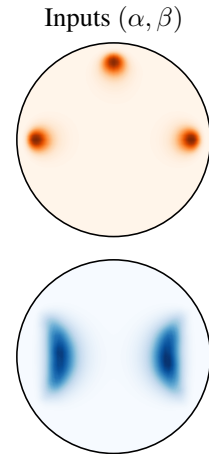
Outputting marginals of SUOT and USOT. The optimal primal marginals of UOT (therefore, SUOT and USOT) are geometric normalizations of inputs (α, β) with discarded outliers. Their computation involves the `Norm` routine, using optimal dual potentials. This is how we compute marginals in Figures 1, 2 and 4: see Appendix B.

Stochastic USOT. In practice, $\hat{\sigma}_K = \frac{1}{K} \sum_i^K \delta_{\theta_i}$ is fixed, and (f_{avg}, g_{avg}) are computed w.r.t. $\hat{\sigma}_K$. However, $\mathbb{E}_{\theta_k \sim \sigma}[\hat{\sigma}_K] = \sigma$. Thus, assuming Theorem 3.7 holds for σ , we have $\mathbb{E}_{\theta_k \sim \sigma}[f_{avg}(x)] = \int f_\theta(\theta^*(x)) d\sigma(\theta)$ if we sample a new $\hat{\sigma}_K$ at each FW step. This approach, which we refer to as, *Stochastic USOT*, should output a more accurate estimate of the USOT w.r.t. σ , but is more expensive: we need to sort projected data w.r.t new projections at each iteration. More importantly, for balanced OT ($\varphi^\circ(x) = x$), USOT = SOT and this idea remains valid for sliced OT. See Section 5 for applications.

5. Experiments

Comparing hyperbolic datasets. We display in Figure 2 the impact of the parameter $\rho = \rho_1 = \rho_2$ on the optimal marginals of USOT. To illustrate the modularity of our FW algorithm, our inputs are synthetic mixtures of Wrapped Normal Distribution on the 2-hyperbolic manifold \mathbb{H} (Nagano et al., 2019), so that the FW oracle is hyperbolic sliced OT (Bonet et al., 2022a). The parameter θ characterizes on \mathbb{H} any geodesic curve passing through the origin, and each sample is projected by taking the shortest path to such geodesics. Once projected on a geodesic curve, we sort data and compute SOT w.r.t. hyperbolic metric $d_{\mathbb{H}}$.

We display the 2-hyperbolic manifold on the Poincaré disc. The measure α (in red) is a mixture of 3 isotropic normal distributions, with a mode at the top of the disc playing the role of an outlier. The measure β is a mixture of two anisotropic normal distributions, whose means are close to two modes of α , but are slightly shifted at the disc’s center. We illustrate several take-home messages, stated in Section 3. First, the optimal



275
 276
 277
 278
 279
 280
 281
 282
 283
 284
 285
 286
 287
 288
 289
 290
 291
 292
 293
 294
 295
 296
 297
 298
 299
 300
 301
 302
 303
 304
 305
 306
 307
 308
 309
 310
 311
 312
 313
 314
 315
 316
 317
 318
 319
 320
 321
 322
 323
 324
 325
 326
 327
 328
 329

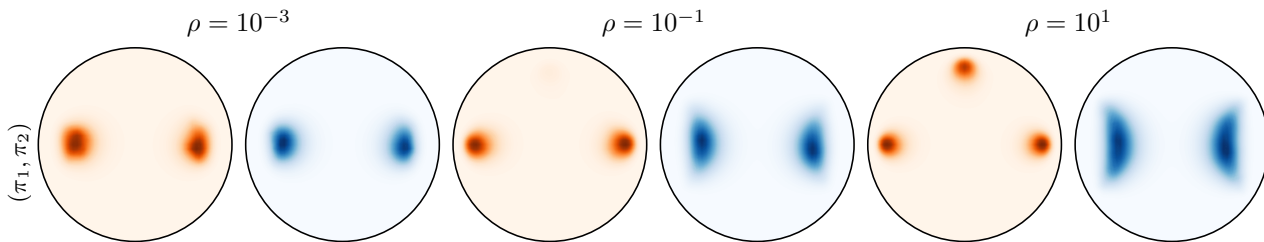


Figure 2: KDE estimation (kernel $e^{-d_{\text{sq}}^2/\sigma}$) of optimal (π_1, π_2) of USOT(α, β) when $D_{\varphi_i} = \rho \text{KL}$.

marginals (π_1, π_2) are renormalisation of (α, β) accounting for their geometry, which are able to remove outliers for properly tuned ρ . When ρ is large, $(\pi_1, \pi_2) \simeq (\alpha, \beta)$ and we retrieve SOT. When ρ is too small, outliers are removed, but we see a shift of the modes, so that modes of (π_1, π_2) are closer to each other, but do not exactly correspond to those of (α, β) . Second, note that such plot cannot be made with SUOT, since the optimal marginals depend on the projection θ (see Figure 1). Third, we are indeed able to reuse any variant of SOT existing in the literature.

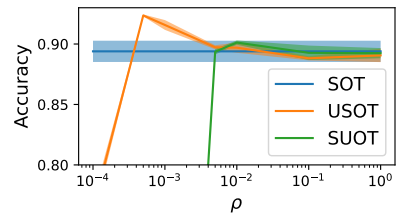
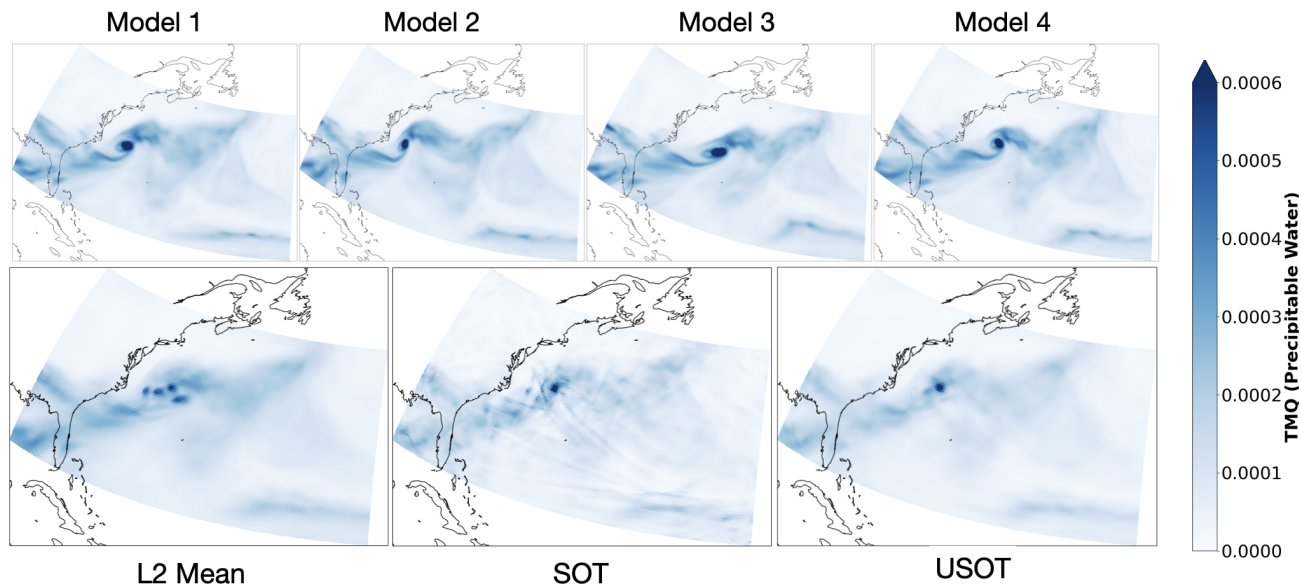
Document classification. To show the benefits of our proposed losses over SOT, we consider a document classification problem (Kusner et al., 2015). Documents are represented as distributions of words embedded with *word2vec* (Mikolov et al., 2013) in dimension $d = 300$. Let D_k be the k -th document and $x_1^k, \dots, x_{n_k}^k \in \mathbb{R}^d$ be the set of words in D_k . Then, $D_k = \sum_{i=1}^{n_k} w_i^k \delta_{x_i^k}$ where w_i^k is the frequency of x_i^k in D_k normalized s.t. $\sum_{i=1}^{n_k} w_i^k = 1$. Given a loss function L , the document classification task is solved by computing the matrix $(L(D_k, D_\ell))_{k,\ell}$, then using a k-nearest neighbor classifier. Since a word typically appears several times in a document, the measures are not uniform and sliced partial OT (Bonnel & Coeurjolly, 2019; Bai et al., 2022) cannot be used in this setting. The aim of this experiment is to show that by discarding possible outliers using a well chosen parameter ρ , USOT is able to outperform SOT and SUOT on this task. We consider BBCSport dataset (Kusner et al., 2015), Movies reviews (Pang et al., 2002) and the Goodreads dataset (Maharjan et al., 2017) on two tasks (genre and likability). We report in Table 1 the accuracy of SUOT, USOT and the stochastic USOT (SUSOT) compared with SOT, OT and UOT computed with the majorization minimization algorithm (Chapel et al., 2021) or approximated with the Sinkhorn algorithm (Pham et al., 2020). All the benchmark methods are computed using the POT library (Flamary et al., 2021). For sliced methods (SOT, SUOT, USOT and SUSOT), we average over 3 computations of the loss matrix and report the standard deviation in Table 1. The number of neighbors was selected via cross validation. The results in Table 1 are reported for ρ yielding the best accuracy, and we display an ablation of this parameter on the BBCSport dataset in Figure 3. We observe that when

ρ is tuned, USOT outperforms SOT, just as UOT outperforms OT. Note that OT and UOT cannot be used in large scale settings (typically large documents) as their complexity scale cubically. We report in Appendix C runtimes on the Goodreads dataset. In particular, computing the OT matrix took 3 times longer than computing the USOT matrix on GPU. Moreover, we were unable to run UOT using POT on the Movies and Goodreads datasets in a reasonable amount of time, due to their computational complexity.

Barycenter on geophysical data. OT barycenters have been an important topic of interest (Bonet et al., 2022b; Le et al., 2021). To compute barycenters under the USOT geometry on a fixed grid, we employ a mirror-descent strategy similar to (Cuturi & Doucet, 2014a, Algorithm (1)): see Appendix C. We showcase unbalanced sliced OT barycenter using climate model data. Ensembles of multiple models are commonly employed to reduce biases and evaluate uncertainties in climate projections (e.g. (Sanderson et al., 2015; Thao et al., 2022)). The commonly used Multi-Model Mean approach assumes models are centered around true values and averages the ensemble with equal or varying weights. However, spatial averaging may fail in capturing specific characteristics of the physical system at stake. We propose to use USOT barycenter here instead. We consider the ClimateNet dataset (Prabhat et al., 2021), and more specifically the TMQ (precipitable water) indicator. The ClimateNet dataset is a human-expert-labeled curated dataset that captures notably tropical cyclones (TCs). In order to simulate the output of several climate models, we take a specific instant (first date of 2011) and deform the data with the elastic deformation from TorchVision (Paszke et al., 2019), in an area located close to the eastern part of the U.S. We obtain 4 different TCs (Figure 4, first row). As expected, the classical L2 spatial mean (Figure 4, second row) reveals 4 different TCs centers/modes, which is undesirable. Since the total TMQ mass in the considered zone varies between the different models, a direct application of SOT is impossible, or requires a normalization of the mass that has undesired effect as can be seen on the second picture of the second row. Finally, we show the result of the USOT barycenter with $\rho_1 = 1e1$ (related to the data) and $\rho_2 = 1e4$ (related to the barycenter). As a result, the corresponding barycenter has

Table 1: Accuracy on document classification

	BBCSport	Movies	Goodreads genre	Goodreads like
OT	91.64	68.88	52.75	70.60
UOT	96.27	-	-	-
Sinkhorn UOT	93.64	63.8	42.55	66.06
SOT	89.39 \pm 0.76	66.95 \pm 0.45	50.09 \pm 0.51	65.60 \pm 0.20
SUOT	90.12 \pm 0.15	67.84 \pm 0.37	50.15 \pm 0.04	66.72 \pm 0.38
USOT	92.36 \pm 0.07	69.21 \pm 0.37	51.87 \pm 0.56	67.41 \pm 1.06
SUSOT	92.45 \pm 0.39	69.53 \pm 0.53	51.93 \pm 0.53	67.33 \pm 0.26

Figure 3: Ablation on BBCSport of the parameter ρ .Figure 4: **Barycenter of geophysical data.** (First row) Simulated output of 4 different climate models depicting different scenarios for the evolution of a tropical cyclone (Second row) Results of different averaging/aggregation strategies.

only one apparent mode which is the expected behaviour. The considered measures have a size of 100×200 , and we run the barycenter algorithm for 500 iterations (with $K = 64$ projections), which takes 3 minutes on a commodity GPU. UOT barycenters for this size of problems are untractable, and to the best of our knowledge, this is the first time such large scale unbalanced OT barycenters can be computed. This experiment encourages an in-depth analysis of the relevance of this aggregation strategy for climate modeling and related problems.

6. Conclusion and Discussion

We proposed two losses merging unbalanced and sliced OT, with theoretical guarantees and an efficient Frank-Wolfe algorithm which allows to reuse any sliced OT variant. We highlighted experimentally the performance improvement over SOT, and described novel applications of unbalanced OT barycenters of positive measures, with a new case study on geophysical data. These novel results and algorithms pave the way to numerous new applications of sliced vari-

ants of OT: we believe our contributions will motivate practitioners to further explore their use in ML applications, without having to pre-process probability measures.

An immediate drawback is the induced additional computational cost w.r.t. SOT. While our empirical results show that SUOT and USOT significantly outperform SOT, and though the complexity is sub-quadratic in the number of samples, our FW approach uses SOT as a subroutine, rendering it necessarily more expensive. Another practical burden comes from the introduction of hyperparameters (ρ_1, ρ_2), which requires cross-validation when possible. A future direction would be to derive efficient strategies to tune (ρ_1, ρ_2), maybe w.r.t. the applicative context, and complement possible interpretations of ρ as a “threshold” for the geometric information encoded by C_1, C_d . On the other hand, while OT between univariate measures defines a reproducing kernel and sliced OT takes advantage of this property (Kolouri et al., 2016; Carriere et al., 2017), some of our experiments suggest this no longer holds for UOT (therefore, for SUOT, USOT). This leaves as an open direction the design of OT-based kernel methods between arbitrary positive measures.

References

- Ambrosio, L., Gigli, N., and Savaré, G. *Gradient flows: in metric spaces and in the space of probability measures*. Springer Science & Business Media, 2005.
- Arjovsky, M., Chintala, S., and Bottou, L. Wasserstein generative adversarial networks. In *International conference on machine learning*, pp. 214–223. PMLR, 2017.
- Bai, Y., Schmitzer, B., Thorpe, M., and Kolouri, S. Sliced optimal partial transport. *arXiv preprint arXiv:2212.08049*, 2022.
- Bayraktar, E. and Guo, G. Strong equivalence between metrics of Wasserstein type. *Electronic Communications in Probability*, 26(none):1 – 13, 2021. doi: 10.1214/21-ECP383.
- Beck, A. and Teboulle, M. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Research Letters*, 31(3):167–175, 2003.
- Bogachev, V. I. and Ruas, M. A. S. *Measure theory*, volume 1. Springer, 2007.
- Bonet, C., Chapel, L., Drumetz, L., and Courty, N. Hyperbolic sliced-wasserstein via geodesic and horospherical projections. *arXiv preprint arXiv:2211.10066*, 2022a.
- Bonet, C., Courty, N., Septier, F., and Drumetz, L. Efficient gradient flows in sliced-wasserstein space. *Transactions on Machine Learning Research*, 2022b.
- Bonet, C., Berg, P., Courty, N., Septier, F., Drumetz, L., and Pham, M.-T. Spherical sliced-wasserstein. In *The Eleventh International Conference on Learning Representations*, 2023a.
- Bonet, C., Malézieux, B., Rakotomamonjy, A., Drumetz, L., Moreau, T., Kowalski, M., and Courty, N. Sliced-wasserstein on symmetric positive definite matrices for m/eeg signals. In *Proceedings of the 40th International Conference on Machine Learning*, 2023b.
- Bonneel, N. and Coeurjolly, D. Spot: sliced partial optimal transport. *ACM Transactions on Graphics (TOG)*, 38(4): 1–13, 2019.
- Bonneel, N., Rabin, J., Peyré, G., and Pfister, H. Sliced and radon wasserstein barycenters of measures. *Journal of Mathematical Imaging and Vision*, 51:22–45, 2015.
- Bonnotte, N. *Unidimensional and evolution methods for optimal transportation*. PhD thesis, Université Paris Sud-Paris XI; Scuola normale superiore (Pise, Italie), 2013.
- Candau-Tilh, J. Wasserstein and sliced-wasserstein distances. Master’s thesis, Université Pierre et Marie Curie, 2020.
- Carriere, M., Cuturi, M., and Oudot, S. Sliced wasserstein kernel for persistence diagrams. In *International conference on machine learning*, pp. 664–673. PMLR, 2017.
- Chapel, L., Flamary, R., Wu, H., Févotte, C., and Gasso, G. Unbalanced optimal transport through non-negative penalized linear regression. *Advances in Neural Information Processing Systems*, 34:23270–23282, 2021.
- Chizat, L. and Bach, F. On the global convergence of gradient descent for over-parameterized models using optimal transport. *Advances in neural information processing systems*, 31, 2018.
- Chizat, L., Peyré, G., Schmitzer, B., and Vialard, F.-X. An interpolating distance between optimal transport and fisher–rao metrics. *Foundations of Computational Mathematics*, 18:1–44, 2018a.
- Chizat, L., Peyré, G., Schmitzer, B., and Vialard, F.-X. Unbalanced optimal transport: Dynamic and kantorovich formulations. *Journal of Functional Analysis*, 274(11): 3090–3123, 2018b.
- Cuturi, M. Sinkhorn distances: Lightspeed computation of optimal transport. *Advances in neural information processing systems*, 26, 2013.
- Cuturi, M. and Doucet, A. Fast computation of wasserstein barycenters. In Xing, E. P. and Jebara, T. (eds.), *Proceedings of the 31st International Conference on Machine Learning*, volume 32 of *Proceedings of Machine Learning Research*, pp. 685–693, Beijing, China, 22–24 Jun 2014a. PMLR.
- Cuturi, M. and Doucet, A. Fast computation of wasserstein barycenters. In *International conference on machine learning*, pp. 685–693. PMLR, 2014b.
- De Bortoli, V., Thornton, J., Heng, J., and Doucet, A. Diffusion schrödinger bridge with applications to score-based generative modeling. *Advances in Neural Information Processing Systems*, 34:17695–17709, 2021.
- Demetci, P., Santorella, R., Chakravarthy, M., Sandstede, B., and Singh, R. Scotv2: Single-cell multiomic alignment with disproportionate cell-type representation. *Journal of Computational Biology*, 29(11):1213–1228, 2022.
- Deshpande, I., Hu, Y.-T., Sun, R., Pyrros, A., Siddiqui, N., Koyejo, S., Zhao, Z., Forsyth, D., and Schwing, A. G. Max-sliced wasserstein distance and its use for gans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10648–10656, 2019.
- Dudley, R. M. The speed of mean glivenko-cantelli convergence. *The Annals of Mathematical Statistics*, 40(1): 40–50, 1969.

- 495 Fatras, K., Zine, Y., Flamary, R., Gribonval, R., and Courty,
496 N. Learning with minibatch wasserstein : asymptotic
497 and gradient properties. In Chiappa, S. and Calandra,
498 R. (eds.), *Proceedings of the Twenty Third International
499 Conference on Artificial Intelligence and Statistics*,
500 volume 108 of *Proceedings of Machine Learning Re-
501 search*, pp. 2131–2141, Online, 26–28 Aug 2020. PMLR.
502 URL [http://proceedings.mlr.press/v108/
503 fatras20a.html](http://proceedings.mlr.press/v108/fatras20a.html).
- 504 Fatras, K., Sejourne, T., Flamary, R., and Courty, N. Un-
505 balanced minibatch optimal transport; applications to
506 domain adaptation. In Meila, M. and Zhang, T. (eds.),
507 *Proceedings of the 38th International Conference on Ma-
508 chine Learning*, volume 139 of *Proceedings of Machine
509 Learning Research*, pp. 3186–3197. PMLR, 18–24 Jul
510 2021. URL [http://proceedings.mlr.press/
511 v139/fatras21a.html](http://proceedings.mlr.press/v139/fatras21a.html).
- 513 Flamary, R., Courty, N., Gramfort, A., Alaya, M. Z., Bois-
514 bunon, A., Chambon, S., Chapel, L., Corenflos, A., Fatras,
515 K., Fournier, N., et al. Pot: Python optimal transport. *The
516 Journal of Machine Learning Research*, 22(1):3571–3578,
517 2021.
- 519 Frank, M. and Wolfe, P. An algorithm for quadratic pro-
520 gramming. *Naval research logistics quarterly*, 3(1-2):
521 95–110, 1956.
- 522 Genevay, A., Chizat, L., Bach, F., Cuturi, M., and Peyré,
523 G. Sample complexity of sinkhorn divergences. In *The
524 22nd international conference on artificial intelligence
525 and statistics*, pp. 1574–1583. PMLR, 2019.
- 527 Goldfeld, Z. and Greenewald, K. Sliced mutual
528 information: A scalable measure of statistical depen-
529 dence. In Ranzato, M., Beygelzimer, A., Dauphin,
530 Y., Liang, P., and Vaughan, J. W. (eds.), *Advances
531 in Neural Information Processing Systems*, vol-
532 ume 34, pp. 17567–17578. Curran Associates, Inc.,
533 2021. URL [https://proceedings.neurips.
534 cc/paper_files/paper/2021/file/
535 92c4661685bf6681f6a33b78ef729658-Paper.
536 pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/92c4661685bf6681f6a33b78ef729658-Paper.pdf).
- 537 Kolouri, S., Zou, Y., and Rohde, G. K. Sliced wasserstein
538 kernels for probability distributions. In *Proceedings of
539 the IEEE Conference on Computer Vision and Pattern
540 Recognition*, pp. 5258–5267, 2016.
- 542 Kolouri, S., Nadjahi, K., Simsekli, U., Badeau, R., and
543 Rohde, G. Generalized sliced wasserstein distances. *Ad-
544 vances in neural information processing systems*, 32,
545 2019.
- 547 Kondratyev, S., Monsaingeon, L., and Vorotnikov, D. A
548 fitness-driven cross-diffusion system from population dy-
549 namics as a gradient flow. *Journal of Differential Equa-
tions*, 261(5):2784–2808, 2016.
- Kusner, M., Sun, Y., Kolkin, N., and Weinberger, K. From
word embeddings to document distances. In *International
conference on machine learning*, pp. 957–966. PMLR,
2015.
- Lacoste-Julien, S. and Jaggi, M. On the global linear con-
vergence of frank-wolfe optimization variants. *Advances
in neural information processing systems*, 28, 2015.
- Le, K., Nguyen, H., Nguyen, Q. M., Pham, T., Bui, H.,
and Ho, N. On robust optimal transport: Computational
complexity and barycenter computation. In Ranzato, M.,
Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan,
J. W. (eds.), *Advances in Neural Information Processing
Systems*, volume 34, pp. 21947–21959, 2021.
- Liero, M., Mielke, A., and Savaré, G. Optimal entropy-
transport problems and a new hellinger–kantorovich dis-
tance between positive measures. *Inventiones mathemati-
cae*, 211(3):969–1117, 2018.
- Maharjan, S., Arevalo, J., Montes, M., González, F. A., and
Solorio, T. A multi-task approach to predict likability of
books. In *Proceedings of the 15th Conference of the Eu-
ropean Chapter of the Association for Computational Lin-
guistics: Volume 1, Long Papers*, pp. 1217–1227, 2017.
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., and
Dean, J. Distributed representations of words and phrases
and their compositionality. *Advances in neural informa-
tion processing systems*, 26, 2013.
- Nadjahi, K., Durmus, A., Chizat, L., Kolouri, S., Shahram-
pour, S., and Simsekli, U. Statistical and topological
properties of sliced probability divergences. *Advances
in Neural Information Processing Systems*, 33:20802–
20812, 2020.
- Nagano, Y., Yamaguchi, S., Fujita, Y., and Koyama, M.
A wrapped normal distribution on hyperbolic space for
gradient-based learning. In *International Conference on
Machine Learning*, pp. 4693–4702. PMLR, 2019.
- Nguyen, K., Ho, N., Pham, T., and Bui, H. Distributional
sliced-wasserstein and applications to generative model-
ing. *arXiv preprint arXiv:2002.07367*, 2020.
- Nietert, S., Goldfeld, Z., and Cummings, R. Outlier-robust
optimal transport: Duality, structure, and statistical analy-
sis. In *Proceedings of The 25th International Conference
on Artificial Intelligence and Statistics*. PMLR, 2022.
- Ohana, R., Nadjahi, K., Rakotomamonjy, A., and Ralaivola,
L. Shedding a pac-bayesian light on adaptive sliced-
wasserstein distances. In *Proceedings of the 40th Inter-
national Conference on Machine Learning*, 2023.

- 550 Pang, B., Lee, L., and Vaithyanathan, S. Thumbs up? senti-
551 ment classification using machine learning techniques.
552 In *Proceedings of EMNLP*, pp. 79–86, 2002.
553
- 554 Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J.,
555 Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga,
556 L., et al. Pytorch: An imperative style, high-performance
557 deep learning library. *Advances in neural information*
558 *processing systems*, 32, 2019.
559
- 560 Peyré, G., Cuturi, M., et al. Computational optimal trans-
561 port: With applications to data science. *Foundations and*
562 *Trends® in Machine Learning*, 11(5-6):355–607, 2019.
563
- 564 Pham, K., Le, K., Ho, N., Pham, T., and Bui, H. On un-
565 balanced optimal transport: An analysis of Sinkhorn
566 algorithm. In III, H. D. and Singh, A. (eds.), *Pro-*
567 *ceedings of the 37th International Conference on Ma-*
568 *chine Learning*, volume 119 of *Proceedings of Machine*
569 *Learning Research*, pp. 7673–7682. PMLR, 13–18 Jul
570 2020. URL [https://proceedings.mlr.press/](https://proceedings.mlr.press/v119/pham20a.html)
571 [v119/pham20a.html](https://proceedings.mlr.press/v119/pham20a.html).
572
- 573 Piccoli, B. and Rossi, F. Generalized wasserstein distance
574 and its application to transport equations with source.
575 *Archive for Rational Mechanics and Analysis*, 211:335–
576 358, 2014.
577
- 578 Prabhat, Kashinath, K., Mudigonda, M., Kim, S., Kapp-
579 Schwoerer, L., Graubner, A., Karaismailoglu, E., von
580 Kleist, L., Kurth, T., Greiner, A., Mahesh, A., Yang,
581 K., Lewis, C., Chen, J., Lou, A., Chandran, S., Toms,
582 B., Chapman, W., Dagon, K., Shields, C. A., O’Brien,
583 T., Wehner, M., and Collins, W. Climateset: an expert-
584 labeled open dataset and deep learning architecture for
585 enabling high-precision analyses of extreme weather.
586 *Geoscientific Model Development*, 14(1):107–124, 2021.
587 doi: 10.5194/gmd-14-107-2021. URL [https://gmd.](https://gmd.copernicus.org/articles/14/107/2021/)
588 [copernicus.org/articles/14/107/2021/](https://gmd.copernicus.org/articles/14/107/2021/).
589
- 590 Rabin, J., Peyré, G., Delon, J., and Bernot, M. Wasser-
591 stein barycenter and its application to texture mixing. In
592 *Scale Space and Variational Methods in Computer Vision:*
593 *Third International Conference, SSVM 2011, Ein-Gedi,*
594 *Israel, May 29–June 2, 2011, Revised Selected Papers 3*,
595 pp. 435–446. Springer, 2012.
596
- 597 Radon, J. 1.1 über die bestimmung von funktionen durch
598 ihre integralwerte längs gewisser mannigfaltigkeiten.
599 *Classic papers in modern diagnostic radiology*, 5(21):
600 124, 2005.
601
- 602 Rotskoff, G., Jelassi, S., Bruna, J., and Vanden-Eijnden,
603 E. Global convergence of neuron birth-death dynamics.
604 *arXiv preprint arXiv:1902.01843*, 2019.
- Sanderson, B. M., Knutti, R., and Caldwell, P. A represen-
tative democracy to reduce interdependency in a multi-
model ensemble. *Journal of Climate*, 28(13):5171–5194,
2015.
- Santambrogio, F. Optimal transport for applied mathemati-
cians. *Birkhäuser, NY*, 55(58-63):94, 2015.
- Sato, R., Yamada, M., and Kashima, H. Fast unbalanced op-
timal transport on a tree. *Advances in neural information*
processing systems, 33:19039–19051, 2020.
- Schiebinger, G., Shu, J., Tabaka, M., Cleary, B., Subrama-
nian, V., Solomon, A., Gould, J., Liu, S., Lin, S., Berube,
P., et al. Optimal-transport analysis of single-cell gene
expression identifies developmental trajectories in repro-
gramming. *Cell*, 176(4):928–943, 2019.
- Séjourné, T., Feydy, J., Vialard, F.-X., Trounev, A., and
Peyré, G. Sinkhorn divergences for unbalanced optimal
transport. *arXiv preprint arXiv:1910.12958*, 2019.
- Séjourné, T., Peyré, G., and Vialard, F.-X. Unbalanced
optimal transport, from theory to numerics. *arXiv preprint*
arXiv:2211.08775, 2022a.
- Séjourné, T., Vialard, F.-X., and Peyré, G. Faster unbal-
anced optimal transport: Translation invariant sinkhorn
and 1-d frank-wolfe. In *International Conference on Ar-*
tificial Intelligence and Statistics, pp. 4995–5021. PMLR,
2022b.
- Simons, S. *Minimax and monotonicity*. Springer, 2006.
- Thao, S., Garvik, M., Mariethoz, G., and Vrac, M. Com-
bining global climate models using graph cuts. *Cli-*
mate Dynamics, 59:2345–2361, 2022. URL <https://hal.science/hal-03620538>.
- Vacher, A. and Vialard, F.-X. Stability of semi-dual unbal-
anced optimal transport: fast statistical rates and conver-
gent algorithm. 2022.
- Xi, J. and Niles-Weed, J. Distributional convergence
of the sliced wasserstein process. *arXiv preprint*
arXiv:2206.00156, 2022.

A. Postponed proofs for Section 3

A.1. Existence of minimizers

We provide the formal statement and detailed proof on the existence of a solution for both SUOT and USOT, as mentioned in Section 3.

Proposition A.1. (Existence of minimizers) *Assume that C_1 is lower-semicontinuous and that either (i) $\varphi'_{1,\infty} = \varphi'_{2,\infty} = +\infty$, or (ii) C_1 has compact sublevels on $\mathbb{R} \times \mathbb{R}$ and $\varphi'_{1,\infty} + \varphi'_{2,\infty} + \inf C_1 > 0$. Then the solution of $\text{SUOT}(\alpha, \beta)$ and $\text{USOT}(\alpha, \beta)$ exist, i.e. the infimum in (5) and (6) is attained. More precisely, there exists (π_1, π_2) which attains the infimum for $\text{USOT}(\alpha, \beta)$ (see Equation (6)). Concerning $\text{SUOT}(\alpha, \beta)$, there exists for any $\theta \in \text{supp}(\sigma)$ a plan π_θ attaining the infimum in $\text{UOT}(\theta_\#^* \alpha, \theta_\#^* \beta)$ (see Equation (2)).*

Proof. We leverage (Liero et al., 2018, Theorem 3.3) to prove this proposition. In the setting of SUOT, if such assumptions (i) or (ii) are satisfied for (α, β) , then they also hold for $(\theta_\#^* \alpha, \theta_\#^* \beta)$ for any $\theta \in \mathbb{S}^{d-1}$. Hence, $\text{UOT}(\theta_\#^* \alpha, \theta_\#^* \beta)$ admits a solution π^θ .

Concerning USOT, note that one necessarily has $m(\pi_1) = m(\pi_2)$, otherwise $\text{SOT}(\pi_1, \pi_2) = +\infty$. From (Liero et al., 2018, Equation (3.10)), that for any admissible (π_1, π_2, π) , one has

$$\text{USOT}(\alpha, \beta) \geq m(\pi) \inf C_1 + m(\alpha) \varphi_1\left(\frac{m(\pi)}{m(\alpha)}\right) + m(\beta) \varphi_2\left(\frac{m(\pi)}{m(\beta)}\right).$$

In both settings the above bounds implies coercivity of the functional of USOT w.r.t. the masses of the measures (π_1, π_2, π) . Thus there exists $M > 0$ such that $m(\pi_1) = m(\pi_2) = m(\pi) < M$, otherwise $\text{USOT}(\alpha, \beta) = +\infty$. By the Banach-Alaoglu theorem, the set of bounded measures (π_1, π_2) is compact, and the set of plans π with such marginals is also compact because \mathbb{R}^d is Polish and C_1 is lower-semicontinuous (Santambrogio, 2015, Theorem 1.7). Because the functional of USOT is lower-semicontinuous in (π_1, π_2, π) and we can restrict optimization over a compact set, we have existence of minimizers for USOT by standard proofs of calculus of variations. \square

A.2. Metric properties: Proof of Proposition 3.2

Proof of Proposition 3.2. Metric properties of SUOT. Symmetry and non-negativity are immediate. Assume $\text{SUOT}(\alpha, \beta) = 0$. Since σ is the uniform distribution on \mathbb{S}^{d-1} , then for any $\theta \in \mathbb{S}^{d-1}$, $\text{UOT}(\theta_\#^* \alpha, \theta_\#^* \beta) = 0$, and since UOT is assumed to be definite, then $\theta_\#^* \alpha = \theta_\#^* \beta$. By (Bogachev & Ruas, 2007, Proposition 3.8.6), this implies that α and β have the same Fourier transform. By injectivity of the Fourier transform, we conclude that $\alpha = \beta$, hence SUOT is definite. The triangle inequality results from applying the Minkowski inequality then the triangle inequality for $\text{UOT}^{1/p}$ for $p \in [1, +\infty)$: for any $\alpha, \beta, \gamma \in \mathcal{M}_+(\mathbb{R}^d)$,

$$\begin{aligned} & \text{SUOT}^{1/p}(\alpha, \beta) \\ &= \left(\int_{\mathbb{S}^{d-1}} \text{UOT}(\theta_\#^* \alpha, \theta_\#^* \beta) d\sigma(\theta) \right)^{1/p} \\ &\leq \left(\int_{\mathbb{S}^{d-1}} [\text{UOT}^{1/p}(\theta_\#^* \alpha, \theta_\#^* \gamma) + \text{UOT}^{1/p}(\theta_\#^* \gamma, \theta_\#^* \beta)]^p d\sigma(\theta) \right)^{1/p} \\ &\leq \left(\int_{\mathbb{S}^{d-1}} [\text{UOT}^{1/p}(\theta_\#^* \alpha, \theta_\#^* \gamma)]^p d\sigma(\theta) \right)^{1/p} + \left(\int_{\mathbb{S}^{d-1}} [\text{UOT}^{1/p}(\theta_\#^* \gamma, \theta_\#^* \beta)]^p d\sigma(\theta) \right)^{1/p} \\ &= \text{SUOT}^{1/p}(\alpha, \gamma) + \text{SUOT}^{1/p}(\gamma, \beta). \end{aligned}$$

Metric properties of USOT. Let $(\alpha, \beta) \in \mathcal{M}_+(\mathbb{R}^d)$. Non-negativity is immediate, as USOT is defined as a program minimizing a sum of positive terms. SOT is symmetric, thus when $\varphi_1 = \varphi_2$, we obtain symmetry of the functional w.r.t. (α, β) . Assume D_φ is definite, i.e. $D_\varphi(\alpha|\beta) = 0$ implies $\alpha = \beta$. Assume now that $\text{USOT}(\alpha, \beta) = 0$, and denote by (π_1, π_2) the optimal marginals attaining the infimum in (6). $\text{USOT}(\alpha, \beta) = 0$ implies that $\text{SOT}(\pi_1, \pi_2) = 0$, $D_\varphi(\pi_1|\alpha) = 0$ and $D_\varphi(\pi_2|\beta) = 0$. These three terms are definite, which yields $\alpha = \pi_1 = \pi_2 = \beta$, hence the definiteness of USOT. \square

A.3. Comparison of SUOT, USOT, SOT, and proof of Theorem 3.3

In this section, we establish several bounds to compare SUOT, USOT and SOT on the space of compactly-supported measures. We provide the detailed derivations and auxiliary lemmas needed for the proofs. Note that Theorem 3.3 is a direct consequence from Theorems A.2 to A.4.

Theorem A.2. *Let X be a compact subset of \mathbb{R}^d with radius R and consider $\alpha, \beta \in \mathcal{M}_+(X)$. Then, $\text{SUOT}(\alpha, \beta) \leq \text{USOT}(\alpha, \beta)$.*

Proof. To show that $\text{SUOT}(\alpha, \beta) \leq \text{USOT}(\alpha, \beta)$, we use a sub-optimality argument. Let π be the solution $\text{USOT}(\alpha, \beta)$ and denote by (π_1, π_2) the marginals of π . For any $\theta \in \mathbb{S}^{d-1}$, denote by π_θ the solution of $\text{OT}(\theta_\#^* \pi_1, \theta_\#^* \pi_2)$. By definition of USOT, the marginals of π_θ are given by $(\theta_\#^* \pi_1, \theta_\#^* \pi_2)$. Since the sequence $(\pi_\theta)_\theta$ is suboptimal for the problem $\text{SUOT}(\alpha, \beta)$, one has

$$\text{SUOT}(\alpha, \beta) \leq \int_{\mathbb{S}^{d-1}} \left\{ \int \mathbf{C}_1 d\pi_\theta + D_{\varphi_1}(\theta_\#^* \pi_1 | \theta_\#^* \alpha) + D_{\varphi_2}(\theta_\#^* \pi_2 | \theta_\#^* \beta) \right\} d\sigma(\theta) \quad (11)$$

$$\leq \int_{\mathbb{S}^{d-1}} \int \mathbf{C}_1 d\pi_\theta d\sigma(\theta) + D_{\varphi_1}(\pi_1 | \alpha) + D_{\varphi_2}(\pi_2 | \beta) \quad (12)$$

$$= \text{USOT}(\alpha, \beta), \quad (13)$$

where the second inequality results from Lemma A.5, and the last equality follows from the definition of $\text{USOT}(\alpha, \beta)$. \square

Theorem A.3. *Let X be a compact subset of \mathbb{R}^d with radius R and consider $\alpha, \beta \in \mathcal{M}_+(X)$. Additionally, let $p \in [1, +\infty)$ and assume $\mathbf{C}_1(x, y) = |x - y|^p$ for $(x, y) \in \mathbb{R} \times \mathbb{R}$ and $\mathbf{C}_d(x, y) = \|x - y\|^p$ for $(x, y) \in \mathbb{R}^d \times \mathbb{R}^d$. Then, $\text{USOT}(\alpha, \beta) \leq \text{UOT}(\alpha, \beta)$.*

Proof. By (Bonnotte, 2013, Proposition 5.1.3), $\text{SOT}(\mu, \nu) \leq K \text{OT}(\mu, \nu)$ with $K \leq 1$. Let π be the solution of $\text{UOT}(\alpha, \beta)$ with marginals (π_1, π_2) . These marginals are sub-optimal for $\text{USOT}(\alpha, \beta)$, we have

$$\text{USOT}(\alpha, \beta) \leq \text{SOT}(\pi_1, \pi_2) + D_{\varphi_1}(\pi_1 | \alpha) + D_{\varphi_2}(\pi_2 | \beta), \quad (14)$$

$$\leq \text{OT}(\pi_1, \pi_2) + D_{\varphi_1}(\pi_1 | \alpha) + D_{\varphi_2}(\pi_2 | \beta), \quad (15)$$

$$= \text{UOT}(\alpha, \beta), \quad (16)$$

where the last equality is obtained because π is optimal in $\text{UOT}(\alpha, \beta)$. \square

Theorem A.4. *Let X be a compact subset of \mathbb{R}^d with radius R and consider $\alpha, \beta \in \mathcal{M}_+(X)$. Additionally, let $p \in [1, +\infty)$ and assume $\mathbf{C}_1(x, y) = |x - y|^p$ for $(x, y) \in \mathbb{R}$ and $\mathbf{C}_d(x, y) = \|x - y\|^p$ for $(x, y) \in \mathbb{R}^d$. Let $\rho > 0$ and assume $D_{\varphi_1} = D_{\varphi_2} = \rho \text{KL}$. Then, $\text{UOT}(\alpha, \beta) \leq c \text{SUOT}(\alpha, \beta)^{1/(d+1)}$, where $c = c(m(\alpha), m(\beta), \rho, R)$ is a non-decreasing function of $m(\alpha)$ and $m(\beta)$.*

Proof. We adapt the proof of (Bonnotte, 2013, Lemma 5.1.4), which establishes a bound between OT and SOT. The first step consists in bounding from above the distance between two regularized measures.

Let $\psi : \mathbb{R}^d \rightarrow \mathbb{R}_+$ be a smooth and radial function verifying $\text{supp}(\psi) \subseteq B_d(\mathbf{0}, 1)$ and $\int_{\mathbb{R}^d} \psi(x) d\text{Leb}(x) = 1$. Let $\psi_\lambda(x) = \lambda^{-d} \psi(x/\lambda) / \mathcal{A}(\mathbb{S}^{d-1})$ where $\mathcal{A}(\mathbb{S}^{d-1})$ is the surface area of \mathbb{S}^{d-1} , i.e. $\mathcal{A}(\mathbb{S}^{d-1}) = \frac{2\pi^{d/2}}{\Gamma(d/2)}$ with Γ the gamma function. For any function f defined on \mathbb{R}^s ($s \geq 1$), denote by $\mathcal{F}[f]$ the Fourier transform of f defined for $x \in \mathbb{R}^s$ as $\mathcal{F}[f](x) = \int_{\mathbb{R}^s} f(w) e^{-i\langle w, x \rangle} dw$. Let $\alpha_\lambda = \alpha * \varphi_\lambda$ and $\beta_\lambda = \beta * \varphi_\lambda$ where $*$ is the convolution operator. Let (f, g) such that $f \oplus g \leq \mathbf{C}_d$. By using the isometry properties of the Fourier transform and the definition of ψ_λ , then representing the variables with polar coordinates, we have

$$\int_{\mathbb{R}^d} \varphi^\circ(f(x)) d\alpha_\lambda(x) = \int_{\mathbb{R}^d} \mathcal{F}[\varphi^\circ \circ f](w) \mathcal{F}[\alpha](w) \mathcal{F}[\psi](\lambda w) dw \quad (17)$$

$$= \int_{\mathbb{S}^{d-1}} \int_0^{+\infty} \mathcal{F}[\varphi^\circ \circ f](r\theta) \mathcal{F}[\alpha](r\theta) \mathcal{F}[\psi](\lambda r) r^{d-1} dr d\sigma(\theta). \quad (18)$$

Since $\varphi^\circ \circ f$ is a real-valued function, $\mathcal{F}[\varphi^\circ \circ f]$ is an even function, then

$$\int_{\mathbb{R}^d} \varphi^\circ(f(x)) d\alpha_\lambda(x) \quad (19)$$

$$= \frac{1}{2} \int_{\mathbb{S}^{d-1}} \int_{\mathbb{R}} \mathcal{F}[\varphi^\circ \circ f](r\theta) \mathcal{F}[\alpha](r\theta) \mathcal{F}[\psi](\lambda r) |r|^{d-1} dr d\sigma(\theta) \quad (20)$$

$$= \frac{1}{2} \int_{\mathbb{S}^{d-1}} \int_{\mathbb{R}} \mathcal{F}[\varphi^\circ \circ f](r\theta) \mathcal{F}[\theta_\#^* \alpha](r) \mathcal{F}[\psi](\lambda r) |r|^{d-1} dr d\sigma(\theta) \quad (21)$$

$$= \frac{1}{2} \int_{\mathbb{S}^{d-1}} \int_{\mathbb{R}} \mathcal{F}[\varphi^\circ \circ f](r\theta) \left(\int_{-R}^R e^{-ir u} d\theta_\#^* \alpha(u) \right) \mathcal{F}[\psi](\lambda r) |r|^{d-1} dr d\sigma(\theta) \quad (22)$$

$$= \frac{1}{2} \int_{\mathbb{S}^{d-1}} \int_{\mathbb{R}} \left(\int_{\mathbb{R}^d} \int_{-R}^R \varphi^\circ(f(x)) e^{-ir(u+\langle \theta, x \rangle)} d\theta_\#^* \alpha(u) \right) \mathcal{F}[\psi](\lambda r) |r|^{d-1} dx dr d\sigma(\theta). \quad (23)$$

Equation (21) follows from the property of push-forward measures, (22) results from the definition of the Fourier transform and $u \in [-R, R]$, and (23) results from the definition of the Fourier transform and Fubini's theorem. By making a change of variables (x becomes $x - u\theta$), we obtain

$$\int_{\mathbb{R}^d} \varphi^\circ(f(x)) d\alpha_\lambda(x) \quad (24)$$

$$= \frac{1}{2} \int_{\mathbb{S}^{d-1}} \int_{\mathbb{R}} \int_{\mathbb{R}^d} \int_{-R}^R \varphi^\circ(f(x - u\theta)) e^{-ir\langle \theta, x \rangle} d\theta_\#^* \alpha(u) \mathcal{F}[\psi](\lambda r) |r|^{d-1} dx dr d\sigma(\theta) \quad (25)$$

$$= \frac{1}{2} \int_{\mathbb{S}^{d-1}} \int_{\mathbb{R}} \int_{B_d(\mathbf{0}, 2R+\lambda)} \int_{-R}^R \varphi^\circ(f(x - u\theta)) e^{-ir\langle \theta, x \rangle} d\theta_\#^* \alpha(u) \mathcal{F}[\psi](\lambda r) |r|^{d-1} dx dr d\sigma(\theta), \quad (26)$$

where (26) follows from the assumption that $\text{supp}(\alpha) \subseteq B_d(\mathbf{0}, R)$. Indeed, this implies that $\text{supp}(\alpha_\lambda) \subseteq B_d(\mathbf{0}, R + \lambda)$, thus the domain of $x \mapsto \varphi^\circ \circ f(x - u\theta)$ is contained in $B_d(\mathbf{0}, 2R + \lambda)$.

Similarly, one can show that

$$\int_{\mathbb{R}^d} \varphi^\circ(g(y)) d\beta_\lambda(y) \quad (27)$$

$$= \frac{1}{2} \int_{\mathbb{S}^{d-1}} \int_{\mathbb{R}} \int_{B_d(\mathbf{0}, 2R+\lambda)} \int_{-R}^R \varphi^\circ(g(y - u\theta)) e^{-ir\langle \theta, y \rangle} d\theta_\#^* \beta(u) \mathcal{F}[\psi](\lambda r) |r|^{d-1} dy dr d\sigma(\theta). \quad (28)$$

By (26) and (28), and applying Fubini's theorem, we obtain

$$\int_{\mathbb{R}^d} \varphi^\circ(f(x)) d\alpha_\lambda(x) + \int_{\mathbb{R}^d} \varphi^\circ(g(y)) d\beta_\lambda(y) \quad (29)$$

$$\leq \frac{1}{2} \int_{\mathbb{R}} \int_{B_d(\mathbf{0}, 2R+\lambda)} \int_{\mathbb{S}^{d-1}} \left\{ \int_{-R}^R \varphi^\circ(f(x - u\theta)) d\theta_\#^* \alpha(u) + \int_{-R}^R \varphi^\circ(g(x - u\theta)) d\theta_\#^* \beta(u) \right\} e^{-ir\langle \theta, x \rangle} \mathcal{F}[\psi](\lambda r) |r|^{d-1} d\sigma(\theta) dx dr \quad (30)$$

$$\leq c_1 (2R + \lambda)^d \int_{\mathbb{S}^{d-1}} \text{UOT}(\theta_\#^* \alpha, \theta_\#^* \beta) d\sigma(\theta) \int_{\mathbb{R}} \lambda^{-d} |\mathcal{F}[\psi](r)| |r|^{d-1} |dr| \quad (31)$$

$$\leq c_2 (2R + \lambda)^d \lambda^{-d} \text{SUOT}(\alpha, \beta) \quad (32)$$

where $c_1 > 0$ is independent from α and β , and $c_2 = c_1 \int_{\mathbb{R}} |\mathcal{F}[\psi](r)| |r|^{d-1} dr$. Equation (32) is obtained by taking the supremum of (30) over the set of potentials (\tilde{f}, \tilde{g}) such that for $u \in [-R, R]$, $\exists (x, \theta) \in B_d(\mathbf{0}, 2R + \lambda) \times \mathbb{S}^{d-1}$, $\tilde{f}(u) = f(x - u\theta)$, $\tilde{g}(u) = g(x - u\theta)$, which is included in the set of potentials (f', g') s.t. $f' : \mathbb{R} \rightarrow \mathbb{R}$, $g' : \mathbb{R} \rightarrow \mathbb{R}$ and $f' \oplus g' \leq C_1$.

We deduce from the dual formulation of UOT (3) and (32) that,

$$\text{UOT}(\alpha_\lambda, \beta_\lambda) \leq c_2(2R + \lambda)^d \lambda^{-d} \text{SUOT}(\alpha, \beta). \quad (33)$$

The last step of the proof consists in relating $\text{UOT}(\alpha_\lambda, \beta_\lambda)$ with $\text{UOT}(\alpha, \beta)$. For any (f, g) such that $f \oplus g \leq C_d$, we have

$$\int_{\mathbb{R}^d} \varphi^\circ(f(x)) d\alpha(x) + \int_{\mathbb{R}^d} \varphi^\circ(g(y)) d\beta(y) - \text{UOT}(\alpha_\lambda, \beta_\lambda) \quad (34)$$

$$\leq \int_{\mathbb{R}^d} \varphi^\circ(f(x)) d\alpha(x) + \int_{\mathbb{R}^d} \varphi^\circ(g(x)) d\beta(x) - \int_{\mathbb{R}^d} \varphi^\circ(f(x)) d\alpha_\lambda(x) - \int_{\mathbb{R}^d} \varphi^\circ(g(y)) d\beta_\lambda(y) \quad (35)$$

$$\leq \int_{\mathbb{R}^d} \{\varphi^\circ(f(x)) - \psi_\lambda * \varphi^\circ(f(x))\} d\alpha(x) + \int_{\mathbb{R}^d} \{\varphi^\circ(g(y)) - \psi_\lambda * \varphi^\circ(g(y))\} d\beta(y). \quad (36)$$

For $x \in \mathbb{R}^d$,

$$\varphi^\circ(f(x)) - \psi_\lambda * \varphi^\circ(f(x)) = \frac{\lambda^{-d}}{\mathcal{A}(\mathbb{S}^{d-1})} \int_{\mathbb{R}^d} (\varphi^\circ(f(x)) - \varphi^\circ(f(y))) \psi\left(\frac{x-y}{\lambda}\right) dy \quad (37)$$

$$\leq \frac{\lambda^{-d}}{\mathcal{A}(\mathbb{S}^{d-1})} \int_{\mathbb{R}^d} |\varphi^\circ(f(x)) - \varphi^\circ(f(y))| \psi\left(\frac{x-y}{\lambda}\right) dy, \quad (38)$$

Since $D_\varphi = \rho \text{KL}$, then for $z \in \mathbb{R}$, $\varphi^\circ(z) = \rho(1 - e^{-z/\rho})$, so for $(x, y) \in \mathbb{R}^d \times \mathbb{R}^d$,

$$\varphi^\circ(f(x)) - \varphi^\circ(f(y)) = \rho(e^{-f(y)/\rho} - e^{-f(x)/\rho}) \quad (39)$$

By Lemma A.8, the potentials (f, g) are bounded by constants depending on $m(\alpha), m(\beta)$, thus we can bound (39) as follows.

$$|\varphi^\circ(f(x)) - \varphi^\circ(f(y))| \leq \rho e^{-\lambda^*/\rho} (1 - e^{-R/\rho}), \quad (40)$$

with $\lambda^* \in [-R + \frac{\rho}{2} \log \frac{m(\alpha)}{m(\beta)}, \frac{R}{2} + \frac{\rho}{2} \log \frac{m(\alpha)}{m(\beta)}]$. We thus derive the following upper-bound on (38).

$$\varphi^\circ(f(x)) - \psi_\lambda * \varphi^\circ(f(x)) \leq \frac{\lambda^{-d}}{\mathcal{A}(\mathbb{S}^{d-1})} \rho e^{-\lambda^*/\rho} (1 - e^{-R/\rho}) \int_{\mathbb{R}^d} \psi\left(\frac{x-y}{\lambda}\right) dy \quad (41)$$

$$\leq \frac{\lambda^{-d+1}}{\mathcal{A}(\mathbb{S}^{d-1})} \rho e^{-\lambda^*/\rho} (1 - e^{-R/\rho}) \int_{\mathbb{R}^d} \frac{1}{\lambda} \psi\left(\frac{x-y}{\lambda}\right) dy \quad (42)$$

$$\leq \frac{\lambda^{-d+1}}{\mathcal{A}(\mathbb{S}^{d-1})} \sqrt{\frac{m(\beta)}{m(\alpha)}} \rho e^{R/\rho} (1 - e^{-R/\rho}) \int_{\mathbb{R}^d} \frac{1}{\lambda} \psi\left(\frac{x-y}{\lambda}\right) dy \quad (43)$$

By doing the change of variables $z = (y - x)/\lambda$ and using the fact that ψ is a radial function and $\int_{\mathbb{R}^d} \psi(z) d\text{Leb}(z) = 1$, we obtain $\int_{\mathbb{R}^d} \frac{1}{\lambda} \psi\left(\frac{x-y}{\lambda}\right) dy = 1$. Therefore,

$$\varphi^\circ(f(x)) - \psi_\lambda * \varphi^\circ(f(x)) \leq \frac{\lambda^{-d+1}}{\mathcal{A}(\mathbb{S}^{d-1})} \sqrt{\frac{m(\beta)}{m(\alpha)}} \rho e^{R/\rho} (1 - e^{-R/\rho}) \quad (44)$$

$$\leq \frac{\lambda}{\mathcal{A}(\mathbb{S}^{d-1})} \sqrt{\frac{m(\beta)}{m(\alpha)}} \rho e^{R/\rho} (1 - e^{-R/\rho}). \quad (45)$$

Similarly, using the bounds on g in Lemma A.8, one can show that

$$|\varphi^\circ(g(x)) - \varphi^\circ(g(y))| \leq \rho e^{\lambda^*/\rho} (e^{R/\rho} - e^{-R/\rho}) \leq \rho \sqrt{\frac{m(\alpha)}{m(\beta)}} e^{R/2\rho} (e^{R/\rho} - e^{-R/\rho}), \quad (46)$$

therefore,

$$\varphi^\circ(g(x)) - \psi_\lambda * \varphi^\circ(g(x)) \leq \frac{\lambda}{\mathcal{A}(\mathbb{S}^{d-1})} \sqrt{\frac{m(\alpha)}{m(\beta)}} \rho e^{R/2\rho} \left(e^{R/\rho} - e^{-R/\rho} \right). \quad (47)$$

We conclude that,

$$\int_{\mathbb{R}^d} \varphi^\circ(f(x)) d\alpha(x) + \int_{\mathbb{R}^d} \varphi^\circ(g(y)) d\beta(y) - \text{UOT}(\alpha_\lambda, \beta_\lambda) \quad (48)$$

$$\leq \frac{\lambda\rho}{\mathcal{A}(\mathbb{S}^{d-1})} \left\{ m(\alpha) e^{-\lambda^*/\rho} \left(1 - e^{-R/\rho} \right) + m(\beta) e^{\lambda^*/\rho} \left(e^{R/\rho} - e^{-R/\rho} \right) \right\} \quad (49)$$

$$\leq \frac{\lambda\rho}{\mathcal{A}(\mathbb{S}^{d-1})} \sqrt{m(\alpha)m(\beta)} \left\{ e^{R/\rho} \left(1 - e^{-R/\rho} \right) + e^{R/2\rho} \left(e^{R/\rho} - e^{-R/\rho} \right) \right\} \quad (50)$$

Taking the supremum on both sides over (f, g) such that $f \oplus g \leq C_d$ yields,

$$\text{UOT}(\alpha, \beta) - \text{UOT}(\alpha_\lambda, \beta_\lambda) \quad (51)$$

$$\leq \frac{\lambda\rho}{\mathcal{A}(\mathbb{S}^{d-1})} \left\{ m(\alpha) e^{-\lambda^*/\rho} \left(1 - e^{-R/\rho} \right) + m(\beta) e^{\lambda^*/\rho} \left(e^{R/\rho} - e^{-R/\rho} \right) \right\} \quad (52)$$

$$\leq \frac{\lambda\rho}{\mathcal{A}(\mathbb{S}^{d-1})} \sqrt{m(\alpha)m(\beta)} \left\{ e^{R/\rho} \left(1 - e^{-R/\rho} \right) + e^{R/2\rho} \left(e^{R/\rho} - e^{-R/\rho} \right) \right\}. \quad (53)$$

Finally, by combining (33) with the above inequality, we obtain

$$\text{UOT}(\alpha, \beta) \quad (54)$$

$$\leq \frac{\lambda\rho}{\mathcal{A}(\mathbb{S}^{d-1})} \sqrt{m(\alpha)m(\beta)} \left\{ e^{R/\rho} \left(1 - e^{-R/\rho} \right) + e^{R/2\rho} \left(e^{R/\rho} - e^{-R/\rho} \right) \right\} \quad (55)$$

$$+ c_2(2R + \lambda)^d \lambda^{-d} \text{SUOT}(\alpha, \beta) \quad (56)$$

$$\leq c\lambda(1 + (2R + \lambda)^d \lambda^{-(d+1)}) \text{SUOT}(\alpha, \beta), \quad (57)$$

where c is a constant satisfying $c \geq c_2$ and

$$c \geq \rho \sqrt{m(\alpha)m(\beta)} \left\{ e^{R/\rho} \left(1 - e^{-R/\rho} \right) + e^{R/2\rho} \left(e^{R/\rho} - e^{-R/\rho} \right) \right\} / \mathcal{A}(\mathbb{S}^{d-1}). \quad (58)$$

We conclude the proof by plugging $\lambda = R^{d/(d+1)} \text{SUOT}(\alpha, \beta)^{1/(d+1)}$ in (57) and using the fact that $\text{SUOT}(\alpha, \beta)$ is bounded from above: $\text{SUOT}(\alpha, \beta) \leq \rho(m(\alpha) + m(\beta))$ since on the one hand, π is suboptimal in (3) thus $\text{UOT}(\alpha, \beta) \leq \rho(m(\alpha) + m(\beta))$, and on the other hand, $m(\alpha) = m(\theta_\#^* \alpha)$ for any $\theta \in \mathbb{S}^{d-1}$. \square

Lemma A.5. For any $\theta \in \mathbb{S}^{d-1}$ and $\alpha, \beta \in \mathcal{M}_+(\mathbb{R}^d)$, $D_\varphi(\theta_\#^* \alpha | \theta_\#^* \beta) \leq D_\varphi(\alpha | \beta)$.

Proof. For $\alpha, \beta \in \mathcal{M}_+(\mathbb{R}^s)$ with $s \geq 1$, the dual characterization of φ -divergences reads (Liero et al., 2018, Theorem 2.7)

$$D_\varphi(\alpha | \beta) = \sup_{f \in \mathcal{E}(\mathbb{R}^s)} \int_{\mathbb{R}^s} \varphi^\circ(f(x)) d\beta(x) - \int_{\mathbb{R}^s} f(x) d\alpha(x),$$

where $\mathcal{E}(\mathbb{R}^s)$ denotes the space of lower semi-continuous functions from \mathbb{R}^s to $\mathbb{R} \cup \{+\infty\}$. Therefore, for any $\theta \in \mathbb{S}^{d-1}$ and $\alpha, \beta \in \mathcal{M}_+(\mathbb{R}^d)$,

$$D_\varphi(\theta_\#^* \alpha | \theta_\#^* \beta) = \sup_{f \in \mathcal{E}(\mathbb{R})} \int_{\mathbb{R}} \varphi^\circ(f(t)) d(\theta_\#^* \beta)(t) - \int_{\mathbb{R}} f(t) d(\theta_\#^* \alpha)(t) \quad (59)$$

$$= \sup_{g: \mathbb{R}^d \rightarrow \mathbb{R} \text{ s.t. } \exists f \in \mathcal{E}(\mathbb{R}), g = f \circ \theta^*} \int_{\mathbb{R}^d} \varphi^\circ(g(x)) d\beta(x) - \int_{\mathbb{R}^d} g(x) d\alpha(x) \quad (60)$$

where (60) results from the definition of push-forward measures. We conclude the proof by observing that the supremum in (60) is taken over a subset of $\mathcal{E}(\mathbb{R}^d)$. \square

Lemma A.6. (*Santambrogio, 2015, Proposition 1.11*) Let $p \in [1, +\infty)$ and assume $C_d(x, y) = \|x - y\|^p$. Let α, β with compact support, such that $C_d(x, y) \leq R^p$ for $(x, y) \in \text{supp}(\alpha) \times \text{supp}(\beta)$. Then without loss of generality the dual potentials (f, g) of $\text{UOT}(\alpha, \beta)$ satisfy $f(x) \in [0, R]$ and $g(y) \in [-R, R]$.

Lemma A.7. (*Séjourné et al., 2022b, Proposition 2*) Define the translation-invariant dual formulation

$$\text{UOT}(\alpha, \beta) = \sup_{f \oplus g \leq C_d} \sup_{\lambda \in \mathbb{R}} \int \varphi_1^\circ(f + \lambda) d\alpha + \int \varphi_2^\circ(g - \lambda) d\beta. \quad (61)$$

Let $\rho > 0$ and assume $D_{\varphi_1} = D_{\varphi_2} = \rho \text{KL}$. Take optimal potentials (f, g) in (61). Then optimal potentials in (3) are given by $(f + \lambda^*(f, g), g - \lambda^*(f, g))$, where the optimal translation λ^* reads

$$\lambda^*(f, g) \triangleq \frac{1}{2} \left[S_\rho^\beta(g) - S_\rho^\alpha(f) \right], \quad S_\rho^\alpha(f) \triangleq -\rho \log \int e^{-f/\rho} d\alpha,$$

and we call $S_\rho^\alpha(f)$ the soft-minimum of f . When $m(\alpha) = 1$ and $m \leq f(x) \leq M$, then $m \leq S_\rho^\alpha(f) \leq M$.

Lemma A.8. Assume (α, β) have compact support such that, for $(x, y) \in \text{supp}(\alpha) \times \text{supp}(\beta)$, $C(x, y) \leq R$. Then, without loss of generality, one can restrict the optimization of the dual formulation (3) of $\text{UOT}(\alpha, \beta)$ over the set of potentials satisfying for $(x, y) \in \text{supp}(\alpha) \times \text{supp}(\beta)$,

$$f(x) \in [\lambda^*, \lambda^* + R], \quad g(y) \in [-\lambda^* - R, -\lambda^* + R],$$

where $\lambda^* \in [-R + \frac{\rho}{2} \log \frac{m(\alpha)}{m(\beta)}, \frac{R}{2} + \frac{\rho}{2} \log \frac{m(\alpha)}{m(\beta)}]$. In particular, one has

$$f(x) \in [-R + \frac{\rho}{2} \log \frac{m(\alpha)}{m(\beta)}, \frac{3R}{2} + \frac{\rho}{2} \log \frac{m(\alpha)}{m(\beta)}], \quad g(y) \in [-\frac{3R}{2} - \frac{\rho}{2} \log \frac{m(\alpha)}{m(\beta)}, 2R - \frac{\rho}{2} \log \frac{m(\alpha)}{m(\beta)}]$$

Proof. Consider the translation-invariant dual formulation (61): if (f, g) are optimal, then for any $\lambda \in \mathbb{R}$, $(f + \lambda, g - \lambda)$ are also optimal. We leverage the structure of the dual constraint $f \oplus g \leq C_d$ with Lemma A.6. Since for $(x, y) \in \text{supp}(\alpha) \times \text{supp}(\beta)$, $C_d(x, y) \leq R$, then without loss of generality, $f(x) \in [0, R]$ and $g(y) \in [-R, R]$. The potentials (f, g) are optimal for the translation-invariant dual energy, and we need a bound for the original dual functional (3). To this end, we leverage Lemma A.7 to compute the optimal translation, such that $(f, g) = (f + \lambda^*(f, g), g - \lambda^*(f, g))$. Let $\bar{\alpha} = \alpha/m(\alpha)$ and $\bar{\beta} = \beta/m(\beta)$ be the normalized probability measures. The translation can be written as,

$$\lambda^*(f, g) = \frac{1}{2} \left[S_\rho^{\bar{\beta}}(g) - S_\rho^{\bar{\alpha}}(f) \right] + \frac{\rho}{2} \log \frac{m(\alpha)}{m(\beta)}, \quad (62)$$

where the functional S_ρ^α is defined in Lemma A.7. Since $\bar{\alpha}$ and $\bar{\beta}$ are probability measures, then by (Genevay et al., 2019, Proposition 1), $f(x) \in [0, R]$ and $g(x) \in [-R, R]$ respectively imply $S_\rho^{\bar{\alpha}}(f) \in [0, R]$ and $S_\rho^{\bar{\beta}}(g) \in [-R, R]$. Combining these bounds on $S_\rho^{\bar{\alpha}}(f)$, $S_\rho^{\bar{\beta}}(g)$ with the expression of $\lambda^*(f, g)$ (62) yields the desired bounds on the optimal potentials (f, g) of the dual formulation (3). \square

A.4. Metrizing weak* convergence: Proof of Theorem 3.4

Proof. Let (α_n) be a sequence of measures in $\mathcal{M}_+(X)$ and $\alpha \in \mathcal{M}_+(X)$, where $X \subset \mathbb{R}^d$ is compact with radius $R > 0$. First, we assume that $\alpha_n \rightharpoonup \alpha$. Then, by (Liero et al., 2018, Theorem 2.25), under our assumptions, $\alpha_n \rightharpoonup \alpha$ is equivalent to $\lim_{n \rightarrow +\infty} \text{UOT}(\alpha_n, \alpha) = 0$. This implies that $\lim_{n \rightarrow +\infty} \text{SUOT}(\alpha_n, \alpha) = 0$ and $\lim_{n \rightarrow +\infty} \text{USOT}(\alpha_n, \alpha) = 0$, since by Theorem 3.3 and non-negativity of SUOT (Proposition 3.2),

$$0 \leq \text{SUOT}(\alpha_n, \alpha) \leq \text{USOT}(\alpha_n, \alpha) \leq \text{UOT}(\alpha_n, \alpha).$$

Conversely, assume either that $\lim_{n \rightarrow +\infty} \text{SUOT}(\alpha_n, \alpha) = 0$ or $\lim_{n \rightarrow +\infty} \text{USOT}(\alpha_n, \alpha) = 0$. First assume there exists $M > 0$ such that for large enough $n \in \mathbb{N}^*$, $m(\alpha_n) \leq M$, then by Theorem 3.3, there exists $c > 0$ such that $\text{UOT}(\alpha_n, \alpha) \leq c(\text{SUOT}(\alpha_n, \alpha))^{1/(d+1)}$. Since c doesn't depend on the masses $(m(\alpha_n), m(\alpha))$, it does not depend on n . By Theorem 3.3, it yields metric equivalence between SUOT, USOT and UOT, thus $\lim_{n \rightarrow +\infty} \text{UOT}(\alpha_n, \alpha) = 0$. By (Liero et al., 2018, Theorem 2.25), we eventually obtain $\alpha_n \rightharpoonup \alpha$, which is the desired result.

The remaining step thus consists in proving that the sequence of masses $(m(\alpha_n))_{n \in \mathbb{N}^*}$ is indeed uniformly bounded by $M > 0$ for large enough n . Note that for any $(\alpha, \beta) \in \mathcal{M}_+(\mathbb{R}^d)$, one has $\text{UOT}(\alpha, \beta) \geq \rho(\sqrt{m(\alpha)} - \sqrt{m(\beta)})^2$. Indeed one has $\text{UOT}(\alpha, \beta) \geq \mathcal{D}(\lambda, -\lambda)$, where \mathcal{D} denotes the dual functional (3) and $\lambda = \frac{\rho}{2} \log \frac{m(\alpha)}{m(\beta)}$. Note that the pair $(\lambda, -\lambda)$ are feasible dual potentials for the constraint $f \oplus g \leq C_d$, because the cost C_d is positive in our setting. The property of push-forwards measures means that for any $\theta \in \mathbb{S}^{d-1}$, one has $m(\theta_\#^* \alpha) = m(\alpha)$. Therefore, we obtain the following bounds for n large enough.

$$\begin{aligned} \text{USOT}(\alpha_n, \alpha) &\geq \text{SUOT}(\alpha_n, \alpha) \geq \int_{\mathbb{S}^{d-1}} \rho \left(\sqrt{m(\theta_\#^* \alpha_n)} - \sqrt{m(\theta_\#^* \alpha)} \right)^2 d\sigma(\theta), \\ &= \rho(\sqrt{m(\alpha_n)} - \sqrt{m(\alpha)})^2. \end{aligned}$$

Hence, $\lim_{n \rightarrow +\infty} \text{SUOT}(\alpha_n, \alpha) = 0$ or $\lim_{n \rightarrow +\infty} \text{USOT}(\alpha_n, \alpha) = 0$ implies $\lim_{n \rightarrow +\infty} m(\alpha_n) = m(\alpha)$. In other terms the mass of sequence converges and is thus uniformly bounded for large enough n . Since we proved that $m(\alpha_n) < M$ and $m(\alpha)$ is finite, it ends the proof. \square

A.5. Application to sliced partial OT: Proof of Theorem 3.5

The proof of Theorem 3.5 relies on a formulation for SUOT and USOT when $D_{\varphi_1} = D_{\varphi_2} = \rho \text{TV}$, which we prove below. Equation (63) is proved in (Piccoli & Rossi, 2014), and can then be applied to SUOT. We include it for completeness. Equation (64) is our contribution and is specific to USOT.

Lemma A.9. *Let $\rho > 0$ and assume $D_{\varphi_1} = D_{\varphi_2} = \rho \text{TV}$ and $C_d(x, y) = \|x - y\|$. Then, for any $(\alpha, \beta) \in \mathcal{M}_+(\mathbb{R}^d)$,*

$$\text{UOT}(\alpha, \beta) = \sup_{f \in \mathcal{E}} \int f(x) d(\alpha - \beta)(x), \quad (63)$$

where

$$\mathcal{E} = \{f : \mathbb{R}^d \rightarrow \mathbb{R}, \|f\|_{\text{Lip}} \leq 1, \|f\|_\infty \leq \rho\},$$

and $\|f\|_\infty \triangleq \sup_{x \in \mathbb{R}^d} |f(x)|$ and $\|f\|_{\text{Lip}} \triangleq \sup_{(x,y) \in \mathbb{R}^d} \frac{|f(x) - f(y)|}{C_d(x,y)}$.

Furthermore, for $C_1(x, y) = |x - y|$ and an empirical approximation $\hat{\sigma}_N = \frac{1}{N} \sum_{i=1}^N \delta_{\theta_i}$ of σ , one has

$$\text{USOT}(\alpha, \beta) = \sup_{(f_\theta) \in \mathcal{E}} \int_{\mathbb{R}^d} \left(\int_{\mathbb{S}^{d-1}} f_\theta(\theta^*(x)) d\hat{\sigma}_N(\theta) \right) d(\alpha - \beta)(x), \quad (64)$$

where

$$\mathcal{E} = \{\forall \theta \in \text{supp}(\hat{\sigma}_N), f_\theta : \mathbb{R} \rightarrow \mathbb{R}, \|f_\theta\|_{\text{Lip}} \leq 1, \|\int_{\mathbb{S}^{d-1}} f_\theta \circ \theta^* d\hat{\sigma}_N(\theta)\|_\infty \leq \rho\},$$

and the Lipschitz norm here is defined w.r.t. C_1 as $\|f\|_{\text{Lip}} \triangleq \sup_{(x,y) \in \mathbb{R}^d} \frac{|f(x) - f(y)|}{C_1(x,y)}$

Proof. We start with the formulation of Equation 3 and Theorem 3.7. For USOT one has

$$\begin{aligned} \text{USOT}(\alpha, \beta) &= \sup_{f_\theta(\cdot) \oplus g_\theta(\cdot) \leq C_1} \int \varphi_1^\circ \left(\int_{\mathbb{S}^{d-1}} f_\theta(\theta^*(x)) d\sigma_N(\theta) \right) d\alpha(x) \\ &\quad + \int \varphi_2^\circ \left(\int_{\mathbb{S}^{d-1}} g_\theta(\theta^*(y)) d\sigma_N(\theta) \right) d\beta(y). \end{aligned}$$

When $D_\varphi = \rho \text{TV}$, the function φ° reads $\varphi^\circ(x) = x$ for $x \in [-\rho, \rho]$, $\varphi^\circ(x) = \rho$ when $x \geq \rho$, and $\varphi^\circ(x) = -\infty$ otherwise. Noting $f_{\text{avg}}(x) = \int_{\mathbb{S}^{d-1}} f_\theta(\theta^*(x)) d\sigma_N(\theta)$ and $g_{\text{avg}}(x) = \int_{\mathbb{S}^{d-1}} g_\theta(\theta^*(x)) d\sigma_N(\theta)$. This formula on φ° imposes $f_{\text{avg}}(x) \geq -\rho$ and $g_{\text{avg}}(x) \geq -\rho$. Furthermore, since we perform a supremum w.r.t. $(f_{\text{avg}}, g_{\text{avg}})$ where φ° attains a plateau, then without loss of generality, we can impose the constraint $f_{\text{avg}}(x) \leq \rho$ and $g_{\text{avg}}(x) \geq \rho$, as it will have no impact on the optimal dual functional value. Thus we have that $\|f_{\text{avg}}\|_\infty \leq \rho$ and $\|g_{\text{avg}}\|_\infty \leq \rho$. To obtain the Lipschitz

property, we use the constraint that $f_\theta(\cdot) \oplus g_\theta(\cdot) \leq C_1$ for any $\theta \in \text{supp}(\sigma_N)$, as well as (Santambrogio, 2015, Proposition 3.1). Thus by using c-transform for the cost $C_1(x, y) = |x - y|$, we can take w.l.o.g $f_\theta(\cdot) = -g_\theta(\cdot)$ with $f_\theta(\cdot)$ a 1-Lipschitz function. Thus w.l.o.g we can perform the supremum over $(f_\theta)_\theta \in \mathcal{E}$, and rephrase the functional as desired, since we have that $\varphi^\circ(f_{avg}) = f_{avg}$.

The proof for UOT is exactly the same, except that our inputs are (f, g) instead of (f_θ, g_θ) . \square

We can now prove Theorem 3.5.

Proof of Theorem 3.5. First we prove that in that setting USOT is a metric. Reusing Lemma A.9, we have that for any measures (α, β, γ)

$$\begin{aligned} \text{USOT}(\alpha, \gamma) &= \sup_{(f_\theta)_\theta \in \mathcal{E}} \int \left(\int_{\mathbb{S}^{d-1}} f_\theta(\theta^*(x)) d\sigma_N \right) d(\alpha - \gamma)(x) \\ &= \sup_{(f_\theta)_\theta \in \mathcal{E}} \int \left(\int_{\mathbb{S}^{d-1}} f_\theta(\theta^*(x)) d\sigma_N \right) d(\alpha - \beta + \beta - \gamma)(x) \\ &\leq \sup_{(f_\theta)_\theta \in \mathcal{E}} \int \left(\int_{\mathbb{S}^{d-1}} f_\theta(\theta^*(x)) d\sigma_N \right) d(\alpha - \beta)(x) \\ &\quad + \sup_{(f_\theta)_\theta \in \mathcal{E}} \int \left(\int_{\mathbb{S}^{d-1}} f_\theta(\theta^*(x)) d\sigma_N \right) d(\beta - \gamma)(x) \\ &= \text{USOT}(\alpha, \beta) + \text{USOT}(\beta, \gamma). \end{aligned}$$

Note that reusing Lemma A.9, we have that SUOT is a sliced integral probability metric over the space of bounded and Lipschitz functions. More precisely, we satisfy the assumptions of (Nadjahi et al., 2020, Theorem 3), so that one has $\text{UOT}(\alpha, \beta) \leq c(\rho, R)(\text{SUOT}(\alpha, \beta))^{1/(d+1)}$.

To prove that USOT and SUOT metrize the weak* convergence, the proof is very similar to that of Theorem 3.4 detailed above. Assuming that $\alpha_n \rightarrow \alpha$ implies $\text{SUOT}(\alpha_n, \alpha) \rightarrow 0$ and $\text{USOT}(\alpha_n, \alpha) \rightarrow 0$ is already proved in Appendix A.4. To prove the converse, the proof is also the same, i.e. we use the property that SUOT, USOT and UOT are equivalent metrics, which holds as we assumed that supports of (α, β) are compact in a ball of radius R . Note that since the bound $\text{UOT}(\alpha, \beta) \leq c(\rho, R)(\text{SUOT}(\alpha, \beta))^{1/(d+1)}$ holds independently of the measure's masses, we do not need to uniformly bound $m(\alpha_n)$, compared to the KL setting of Theorem 3.4. \square

A.6. Sample complexity: Proof of Theorem 3.6

Theorem 3.6 is obtained by adapting (Nadjahi et al., 2020, Theorems 4 and 5). We provide the detailed derivations below.

Proof of Theorem 3.6. Let α, β in $\mathcal{M}_+(\mathbb{R}^d)$ with respective empirical approximations $\hat{\alpha}_n, \hat{\beta}_n$ over n samples. By using the definition of SUOT, the triangle inequality and the assumed sample complexity of UOT for univariate measures, we show that

$$\mathbb{E} \left| \text{SUOT}(\alpha, \beta) - \text{SUOT}(\hat{\alpha}_n, \hat{\beta}_n) \right| \tag{65}$$

$$= \mathbb{E} \left| \int_{\mathbb{S}^{d-1}} \{ \text{UOT}(\theta_\#^* \alpha, \theta_\#^* \beta) - \text{UOT}(\theta_\#^* \hat{\alpha}_n, \theta_\#^* \hat{\beta}_n) \} d\sigma(\theta) \right| \tag{66}$$

$$\leq \mathbb{E} \left\{ \int_{\mathbb{S}^{d-1}} | \text{UOT}(\theta_\#^* \alpha, \theta_\#^* \beta) - \text{UOT}(\theta_\#^* \hat{\alpha}_n, \theta_\#^* \hat{\beta}_n) | d\sigma(\theta) \right\} \tag{67}$$

$$\leq \int_{\mathbb{S}^{d-1}} \mathbb{E} | \text{UOT}(\theta_\#^* \alpha, \theta_\#^* \beta) - \text{UOT}(\theta_\#^* \hat{\alpha}_n, \theta_\#^* \hat{\beta}_n) | d\sigma(\theta) \tag{68}$$

$$\leq \int_{\mathbb{S}^{d-1}} \kappa(n) d\sigma(\theta) = \kappa(n), \tag{69}$$

which completes the proof for the first setting.

Next, let $\alpha \in \mathcal{M}_+(\mathbb{R}^d)$ with corresponding empirical approximation $\hat{\alpha}_n$. Then, using the definition of SUOT, the triangle inequality (w.r.t. integral) and the assumed convergence rate in UOT,

$$\mathbb{E} |\text{SUOT}(\hat{\alpha}_n, \alpha)| \tag{70}$$

$$= \mathbb{E} \left| \int_{\mathbb{S}^{d-1}} \text{UOT}(\theta_{\#}^* \hat{\alpha}_n, \theta_{\#}^* \alpha) d\sigma(\theta) \right| \leq \mathbb{E} \left\{ \int_{\mathbb{S}^{d-1}} |\text{UOT}(\theta_{\#}^* \hat{\alpha}_n, \theta_{\#}^* \alpha)| d\sigma(\theta) \right\} \tag{71}$$

$$\leq \int_{\mathbb{S}^{d-1}} \mathbb{E} |\text{UOT}(\theta_{\#}^* \hat{\alpha}_n, \theta_{\#}^* \alpha)| d\sigma(\theta) \leq \int_{\mathbb{S}^{d-1}} \xi(n) d\sigma(\theta) = \xi(n). \tag{72}$$

Additionally, if we assume that $\text{UOT}^{1/p}$ satisfies non-negativity, symmetry and the triangle inequality on $\mathcal{M}_+(\mathbb{R}) \times \mathcal{M}_+(\mathbb{R})$, then by Proposition 3.2, $\text{SUOT}^{1/p}$ verifies these three metric properties on $\mathcal{M}_+(\mathbb{R}^d) \times \mathcal{M}_+(\mathbb{R}^d)$, and we can derive its sample complexity as follows. For any α, β in $\mathcal{M}_+(\mathbb{R}^d)$ with respective empirical approximations $\hat{\alpha}_n, \hat{\beta}_n$, applying the triangle inequality yields for $p \in [1, +\infty)$,

$$\left| \text{UOT}^{1/p}(\alpha, \beta) - \text{UOT}^{1/p}(\hat{\alpha}_n, \hat{\beta}_n) \right| \leq \text{UOT}^{1/p}(\hat{\alpha}_n, \alpha) + \text{UOT}^{1/p}(\hat{\beta}_n, \beta). \tag{73}$$

Taking the expectation of (73) with respect to $\hat{\alpha}_n, \hat{\beta}_n$ gives,

$$\mathbb{E} \left| \text{SUOT}^{1/p}(\alpha, \beta) - \text{SUOT}^{1/p}(\hat{\alpha}_n, \hat{\beta}_n) \right| \leq \mathbb{E} |\text{SUOT}^{1/p}(\hat{\alpha}_n, \alpha)| + \mathbb{E} |\text{SUOT}^{1/p}(\hat{\beta}_n, \beta)| \tag{74}$$

$$\leq \{\mathbb{E} |\text{SUOT}(\hat{\alpha}_n, \alpha)|\}^{1/p} + \{\mathbb{E} |\text{SUOT}(\hat{\beta}_n, \beta)|\}^{1/p} \tag{75}$$

$$\leq \xi(n)^{1/p} + \xi(n)^{1/p} = 2\xi(n)^{1/p}, \tag{76}$$

where (75) is immediate if $p = 1$, and results from applying Hölder's inequality on \mathbb{S}^{d-1} if $p > 1$, and (76) follows from (72). \square

A.7. Strong duality: Proof of Theorem 3.7

Proof of Theorem 3.7. Note that the result for SUOT is already proved in Lemma A.12. Thus we focus on the proof of duality for USOT. We start from the definition of USOT, reformulate it to apply the strong duality result of Proposition A.10 and obtain our reformulation. We first have that

$$\begin{aligned} \text{USOT}(\alpha, \beta) &= \inf_{(\pi_1, \pi_2) \in \mathcal{M}_+(\mathbb{R}^d)^2} \left\{ \text{SOT}(\pi_1, \pi_2) + D_{\varphi_1}(\pi_1 | \alpha) + D_{\varphi_2}(\pi_2 | \beta) \right\}, \\ &= \inf_{(\pi_1, \pi_2) \in \mathcal{M}_+(\mathbb{R}^d)^2} \left\{ \int_{\mathbb{S}^{d-1}} \left[\sup_{f_{\theta} \oplus g_{\theta} \leq C_1} \int f_{\theta} d(\theta_{\#}^* \pi_1) + \int g_{\theta} d(\theta_{\#}^* \pi_2) \right] d\hat{\sigma}_K(\theta) \right. \\ &\quad \left. + \sup_{\tilde{f} \in \mathcal{E}(\mathbb{R}^d)} \int \varphi_1^{\circ}(\tilde{f}(x)) d\alpha(x) - \int \tilde{f}(x) d\pi_1(x) \right. \\ &\quad \left. + \sup_{\tilde{g} \in \mathcal{E}(\mathbb{R}^d)} \int \varphi_2^{\circ}(\tilde{g}(y)) d\beta(y) - \int \tilde{g}(y) d\pi_2(y) \right\}, \\ &= \inf_{(\pi_1, \pi_2) \in \mathcal{M}_+(\mathbb{R}^d)^2} \left\{ \sup_{f_{\theta} \oplus g_{\theta} \leq C_1} \int_{\mathbb{S}^{d-1}} \left[\int f_{\theta} d(\theta_{\#}^* \pi_1) + \int g_{\theta} d(\theta_{\#}^* \pi_2) \right] d\hat{\sigma}_K(\theta) \right. \\ &\quad \left. + \sup_{\tilde{f} \in \mathcal{E}(\mathbb{R}^d)} \int \varphi_1^{\circ}(\tilde{f}(x)) d\alpha(x) - \int \tilde{f}(x) d\pi_1(x) \right. \\ &\quad \left. + \sup_{\tilde{g} \in \mathcal{E}(\mathbb{R}^d)} \int \varphi_2^{\circ}(\tilde{g}(y)) d\beta(y) - \int \tilde{g}(y) d\pi_2(y) \right\}, \end{aligned}$$

where $\mathcal{E}(\mathbb{R}^d)$ denotes a set of lower-semicontinuous functions, and the last equality holds thanks to Lemma A.11.

1100 We focus now on verifying that Proposition A.10 holds, so that we can swap the infimum and the supremum. Define the
 1101 functional

$$\begin{aligned}
 1102 \mathcal{L}((\pi_1, \pi_2), ((f_\theta)_\theta, (g_\theta)_\theta, \tilde{f}, \tilde{g})) &\triangleq \int_{\mathbb{S}^{d-1}} \left[\int f_\theta d(\theta_\#^* \pi_1) + \int g_\theta d(\theta_\#^* \pi_2) \right] d\hat{\sigma}_K(\theta) \\
 1103 &+ \int \varphi_1^\circ(\tilde{f}(x)) d\alpha(x) - \int \tilde{f}(x) d\pi_1(x) \\
 1104 &+ \int \varphi_2^\circ(\tilde{g}(y)) d\beta(y) - \int \tilde{g}(y) d\pi_2(y).
 \end{aligned}$$

1105
 1106
 1107
 1108
 1109 One has that,

- 1110 • For any $((f_\theta)_\theta, (g_\theta)_\theta, \tilde{f}, \tilde{g})$, \mathcal{L} is linear (thus convex) and lower-semicontinuous.
- 1111 • For any (π_1, π_2) , \mathcal{L} is concave in $((f_\theta)_\theta, (g_\theta)_\theta, \tilde{f}, \tilde{g})$ because φ_i° is concave and thus \mathcal{L} is a sum of linear or concave functions.

1112
 1113
 1114
 1115
 1116
 1117 Furthermore, since we assumed e.g. that $0 \in \text{dom}(\varphi)$, then

$$1118 \sup_{((f_\theta)_\theta, (g_\theta)_\theta, \tilde{f}, \tilde{g})} \inf_{(\pi_1, \pi_2) \in \mathcal{M}_+(\mathbb{R}^d)^2} \mathcal{L} \leq \text{USOT}(\alpha, \beta) \leq \varphi_1(0)m(\alpha) + \varphi_2(0)m(\beta),$$

1119 because the marginals $(\pi_1, \pi_2) = (0, 0)$ are admissible and suboptimal. If we consider instead that $(m(\alpha), m(\beta)) \in \text{dom}(\varphi)$,
 1120 then we take the marginals $\pi_1 = \alpha/m(\alpha)$ and $\pi_2 = \beta/m(\beta)$, which yields an upper-bound by $m(\alpha)\varphi_1(\frac{1}{m(\alpha)}) +$
 1121 $m(\beta)\varphi_2(\frac{1}{m(\beta)})$. Then we consider an anchor dual point $b^* = ((f_\theta)_\theta, (g_\theta)_\theta, \tilde{f}, \tilde{g})$ to bound \mathcal{L} over a compact set. We take
 1122 $f_\theta = 0, g_\theta = 0$, which are always admissible since we take $C_1(x, y) \geq 0$. Then, since we assume there exists $p_i \leq 0$ in
 1123 $\text{dom}(\varphi_i^*)$, we take $\tilde{f} = p_1$ and $\tilde{g} = p_2$. For these potentials one has:

$$1124 \mathcal{L}((\pi_1, \pi_2), b^*) = \varphi_1^\circ(p_1)m(\alpha) - p_1m(\pi_1) + \varphi_2^\circ(p_2)m(\alpha) - p_2m(\pi_2).$$

1125
 1126
 1127
 1128
 1129 Note that the functional at this point only depends on the masses of the marginals (π_1, π_2) . Since $(p_1, p_2) \geq 0$ the
 1130 set of (π_1, π_2) such that $\mathcal{L}((\pi_1, \pi_2), b^*) \leq \varphi_1(0)m(\alpha) + \varphi_2(0)m(\beta)$ is non-empty (at least in a neighbourhood of
 1131 $(\pi_1, \pi_2) = (0, 0)$), and that $(m(\pi_1), m(\pi_2))$ are uniformly bounded by some constant $M > 0$. By the Banach-Alaoglu
 1132 theorem, such set of measures is compact for the weak* topology.

1133
 1134 Therefore, Proposition A.10 holds and we have strong duality, *i.e.*

$$1135 \text{USOT}(\alpha, \beta) = \sup_{\left\{ \begin{array}{l} f_\theta \oplus g_\theta \leq C_1 \\ (\tilde{f}, \tilde{g}) \in \mathcal{E}(\mathbb{R}^d) \end{array} \right\}} \inf_{(\pi_1, \pi_2) \in \mathcal{M}_+(\mathbb{R}^d)^2} \mathcal{L}((\pi_1, \pi_2), ((f_\theta)_\theta, (g_\theta)_\theta, \tilde{f}, \tilde{g})).$$

1136
 1137
 1138
 1139
 1140 To achieve the proof, note that taking the infimum in (π_1, π_2) (for fixed dual variables) reads

$$\begin{aligned}
 1141 \inf_{\pi_1, \pi_2 \geq 0} \int &\left(\int_{\mathbb{S}^{d-1}} f_\theta(\theta^*(x)) d\hat{\sigma}_K(\theta) \right) d\pi_1(x) - \int \tilde{f}(x) d\pi_1(x) \\
 1142 &+ \int \left(\int_{\mathbb{S}^{d-1}} g_\theta(\theta^*(y)) d\hat{\sigma}_K(\theta) \right) d\pi_2(y) - \int \tilde{g}(y) d\pi_2(y).
 \end{aligned}$$

1143
 1144
 1145
 1146
 1147
 1148 Note that we applied Fubini's theorem here, which holds here because all measures have compact support, thus all quantities
 1149 are finite. It allows to rephrase the minimization over $\pi_1, \pi_2 \geq 0$ as the following constraint

$$1150 \int_{\mathbb{S}^{d-1}} f_\theta(\theta^*(x)) d\hat{\sigma}_K(\theta) \geq \tilde{f}(x), \quad \int_{\mathbb{S}^{d-1}} g_\theta(\theta^*(y)) d\hat{\sigma}_K(\theta) \geq \tilde{g}(y),$$

1151
 1152
 1153
 1154

otherwise the infimum is $-\infty$. However, the function φ° is non-decreasing (see (Séjourné et al., 2019, Proposition 2)). Thus the maximization in (\tilde{f}, \tilde{g}) is optimal when the above inequality is actually an equality, i.e.

$$\int_{\mathbb{S}^{d-1}} f_\theta(\theta^*(x)) d\hat{\sigma}_K(\theta) = \tilde{f}(x), \quad \int_{\mathbb{S}^{d-1}} g_\theta(\theta^*(y)) d\hat{\sigma}_K(\theta) = \tilde{g}(y).$$

Plugging the above relation in the functional \mathcal{L} yields the desired result on the dual of USOT and ends the proof. \square

We mention a strong duality result which is very general and which we use in the proof of 3.7. This result is taken from (Liero et al., 2018, Theorem 2.4) which itself takes it from (Simons, 2006).

Proposition A.10. (Liero et al., 2018, Theorem 2.4) *Consider two sets A and B be nonempty convex sets of some vector spaces. Assume A is endowed with a Hausdorff topology. Let $L : A \times B \rightarrow \mathbb{R}$ be a function such that*

1. $a \mapsto L(a, b)$ is convex and lower-semicontinuous on A , for every $b \in B$
2. $b \mapsto L(a, b)$ is concave on B , for every $a \in A$.

If there exists $b_* \in B$ and $\kappa > \sup_{b \in B} \inf_{a \in A} L(a, b)$ such that the set $\{a \in A, L(a, b_*) < \kappa\}$ is compact in A , then

$$\inf_{a \in A} \sup_{b \in B} L(a, b) = \sup_{b \in B} \inf_{a \in A} L(a, b)$$

We also consider the following to swap the supremum in the integral which defines sliced-UOT (and in particular sliced-OT). In what follows we note sliced potentials as functions $f_\theta(z)$ with $(\theta, z) \in \mathbb{S}^{d-1} \times \mathbb{R}$, such that

$$\text{SUOT}(\alpha, \beta) = \int_{\mathbb{S}^{d-1}} \left[\sup_{f_\theta \oplus g_\theta \leq C_1} \int \varphi^\circ \circ f_\theta d(\theta_\#^* \alpha) + \int \varphi^\circ \circ g_\theta d(\theta_\#^* \beta) \right] d\hat{\sigma}_K(\theta).$$

Note that with the above definition, $z \mapsto f_\theta(z)$ is continuous for any θ , but $\theta \mapsto f_\theta(z)$ is only $\hat{\sigma}_K$ -measurable.

Lemma A.11. *Consider two sets X and Y , a measure σ such that $\sigma(X) < +\infty$. Assume Y is compact. Consider a function $\mathcal{F} : X \times Y \rightarrow \mathbb{R}$. Assume there exists a sequence (y_n) in Y such that $\mathcal{F}(\cdot, y_n) \rightarrow \sup_{y \in Y} \mathcal{F}(\cdot, y)$ uniformly. Then one has*

$$\sup_{y \in Y} \int_X \mathcal{F}(x, y) d\sigma(x) = \int_X \sup_{y \in Y} \mathcal{F}(x, y) d\sigma(x).$$

Proof. Define $\mathcal{G}(x) = \sup_{y \in Y} \mathcal{F}(x, y)$ and $\mathcal{H}(x, y) \triangleq \mathcal{G}(x) - \mathcal{F}(x, y)$. One has $\mathcal{H} \geq 0$ by definition, and the desired equality can be rewritten as

$$\begin{aligned} \sup_{y \in Y} \int_X \mathcal{F}(x, y) d\sigma(x) &= \int_X \sup_{y \in Y} \mathcal{F}(x, y) d\sigma(x) \\ &\Leftrightarrow \inf_{y \in Y} \int_X \mathcal{H}(x, y) d\sigma(x) = 0. \end{aligned}$$

Since the integral involving \mathcal{H} is non-negative, the infimum is zero if and only if we have a sequence (y_n) such that $\int_X \mathcal{H}(\cdot, y_n) d\sigma \rightarrow 0$. By assumption, one has $\mathcal{F}(\cdot, y_n) \rightarrow \sup_{y \in Y} \mathcal{F}(\cdot, y)$ uniformly, i.e. $\|\mathcal{H}(\cdot, y_n)\|_\infty \rightarrow 0$. This implies thanks to Holder's inequality that

$$0 \leq \int_X \mathcal{H}(\cdot, y_n) d\sigma \leq \sigma(X) \|\mathcal{H}(\cdot, y_n)\|_\infty$$

Thus by assumption one has $\int_X \mathcal{F}(\cdot, y_n) d\sigma \rightarrow \int_X \mathcal{G} d\sigma$, which indeed means that we have the desired permutation between supremum and integral. \square

1210 **Lemma A.12.** Let $p \in [1, +\infty)$ and assume that $C_1(x, y) = |x - y|^p$. Consider two positive measures (α, β) with compact
 1211 support. Assume that the measure $\hat{\sigma}_K$ is discrete, i.e. $\hat{\sigma}_K = \frac{1}{K} \sum_{i=1}^K \delta_{\theta_i}$ with $\theta_i \in \mathbb{S}^{d-1}$, $i = 1, \dots, n$. Then, one can swap
 1212 the integral over the sphere and the supremum in the dual formulation of SUOT, such that

$$1213 \text{SUOT}(\alpha, \beta) = \sup_{f_\theta \oplus g_\theta \leq C_1} \int_{\mathbb{S}^{d-1}} \left[\int \varphi^\circ \circ f_\theta d(\theta_\#^* \alpha) + \int \varphi^\circ \circ g_\theta d(\theta_\#^* \beta) \right] d\hat{\sigma}_K(\theta).$$

1214
 1215
 1216 In particular, this result is valid for SOT.

1217
 1218 *Proof.* The proof consists in applying Lemma A.11 for (X, Y) chosen as $X = \text{supp}(\hat{\sigma}_K) \subset \mathbb{S}^{d-1}$ and

$$1219 Y = \{ \forall \theta \in \text{supp}(\hat{\sigma}_K), f_\theta : \mathbb{R} \rightarrow \mathbb{R}, g_\theta : \mathbb{R} \rightarrow \mathbb{R}, f_\theta(x) + g_\theta(y) \leq C_1(x, y) \}.$$

1220
 1221 The functions in Y are dual potentials, and by definition are continuous for any θ . Let $\mathcal{F} : X \times Y \rightarrow \mathbb{R}$ be the functional
 1222 defined as

$$1223 \mathcal{F} : (\theta, (f_\theta)_\theta, (g_\theta)_\theta) \mapsto \int f_\theta d(\theta_\#^* \alpha) + \int g_\theta d(\theta_\#^* \beta).$$

1224
 1225 Since the measures (α, β) have compact support, then by Lemma A.13, the supremum is attained over a subset of dual
 1226 potentials of Y such that for any fixed $\theta \in X$, (f_θ, g_θ) are Lipschitz-continuous and bounded, thus uniformly equicontinuous
 1227 functions (with constants independent of θ). By the Ascoli-Arzelà theorem, the set of uniformly equicontinuous functions is
 1228 compact for the uniform convergence. Hence, for any $\theta \in X$, there exists a sequence of dual potentials $(f_{\theta,n}, g_{\theta,n})$ which
 1229 uniformly converges to optimal dual potentials (f_θ, g_θ) (up to extraction of subsequence). Besides, we have $\text{OT}(\theta_\#^* \alpha, \theta_\#^* \beta) =$
 1230 $\mathcal{F}(\theta, f_\theta, g_\theta)$ and $\mathcal{F}(\theta, (f_{\theta,n})_\theta, (g_{\theta,n})_\theta) \rightarrow \text{OT}(\theta_\#^* \alpha, \theta_\#^* \beta)$ as $n \rightarrow +\infty$. Denote $\mathcal{F}_n(\theta) \triangleq \mathcal{F}(\theta, (f_{\theta,n})_\theta, (g_{\theta,n})_\theta)$ and
 1231 $\text{OT}(\theta) \triangleq \text{OT}(\theta_\#^* \alpha, \theta_\#^* \beta)$. In order to apply Lemma A.11, we need to prove that the convergence of $(\mathcal{F}_n(\theta))_{n \in \mathbb{N}^*}$ to
 1232 $\text{OT}(\theta_\#^* \alpha, \theta_\#^* \beta)$ is uniform w.r.t. θ , i.e. $\sup_{\theta \in X} |\mathcal{F}_n(\theta) - \text{OT}(\theta)| \rightarrow 0$ as $n \rightarrow +\infty$.

1233
 1234 First, note that for any $\theta \in X$,

$$1235 |\mathcal{F}_n(\theta) - \text{OT}(\theta)| \leq m(\alpha) \|f_{\theta,n} - f_\theta\|_\infty + m(\beta) \|g_{\theta,n} - g_\theta\|_\infty.$$

1236
 1237 Since for a fixed $\theta \in X$, $(f_{\theta,n}, g_{\theta,n})_{n \in \mathbb{N}^*}$ uniformly converge to (f_θ, g_θ) , this means that

$$1238 \forall \theta \in X, \forall \varepsilon > 0, \exists N(\varepsilon, \theta), \forall n \geq N(\varepsilon, \theta), m(\alpha) \|f_{\theta,n} - f_\theta\|_\infty + m(\beta) \|g_{\theta,n} - g_\theta\|_\infty < \varepsilon.$$

1239
 1240 Since we assume that σ is supported on a discrete set, then the cardinal of X is finite and one can define $N(\varepsilon) \triangleq$
 1241 $\max_{\theta \in X} N(\varepsilon, \theta)$. This yields,

$$1242 \forall \varepsilon > 0, \exists N(\varepsilon), \forall n \geq N(\varepsilon), \sup_{\theta \in X} |\mathcal{F}_n(\theta) - \text{OT}(\theta)| < \varepsilon.$$

1243
 1244 which means that $\sup_{\theta \in X} |\mathcal{F}_n(\theta) - \text{OT}(\theta)| \rightarrow 0$, thus concludes the proof.

1245 □

1246
 1247 **Lemma A.13.** Let $p \in [1, +\infty)$ and $C_1(x, y) = |x - y|^p$. Consider two positive measures $(\alpha, \beta) \in \mathcal{M}_+(\mathbb{R}^d)$ whose
 1248 support is such that $C_d(x, y) = \|x - y\|^p \leq R$. Then for any $\theta \in \mathbb{S}^{d-1}$, one can restrict without loss of generality the
 1249 problem $\text{UOT}(\theta_\#^* \alpha, \theta_\#^* \beta)$ as a supremum over dual potentials satisfying $f_\theta(x) + g_\theta(y) \leq C_1(x, y)$, uniformly bounded by
 1250 M and uniformly L -Lipschitz, where M and L do not depend on θ .

1251
 1252 *Proof.* We adapt the proof of (Santambrogio, 2015, Proposition 1.11), and focus on showing that the uniform boundedness
 1253 and Lipschitz constant are independent of $\theta \in \mathbb{S}^{d-1}$ in this setting. Here we consider the translation-invariant formulation
 1254 of UOT from (Séjourné et al., 2022b), i.e. $\text{UOT}(\alpha, \beta) = \sup_{f \oplus g \leq C_d} \mathcal{H}(f, g)$, where $\mathcal{H}(f, g) = \sup_{\lambda \in \mathbb{R}} \mathcal{D}(f + \lambda, g - \lambda)$.
 1255 It is proved in (Séjourné et al., 2022b, Proposition 9) that the above problem has the same primal and is thus equivalent
 1256 to optimize \mathcal{D} . By definition one has $\mathcal{H}(f, g) = \mathcal{H}(f + \lambda, g - \lambda)$ for any $\lambda \in \mathbb{R}$, i.e. this formulation shares the same
 1257 invariance as Balanced OT. Thus we can reuse all arguments from (Santambrogio, 2015, Proposition 1.11), such that for
 1258
 1259
 1260
 1261
 1262
 1263
 1264

1265 UOT(α, β), one can use the constraint $f(x) + g(y) \leq C_d(x, y)$ and the assumption $C_d(x, y) \leq R$ to prove that without loss
 1266 of generality, one can restrict to potentials such that $f(x) \in [0, R]$ and $g(y) \in [-R, R]$. Furthermore if the cost satisfies in
 1267 \mathbb{R}^d

$$|C_d(x, y) - C_d(x', y')| \leq L(\|x - x'\| + \|y - y'\|),$$

1270 then one can also restrict w.l.o.g. to potentials which are L -Lipschitz. For the cost $C_d(x, y) = \|x - y\|^p$ with $p \geq 1$, this
 1271 holds with constant $L = pR^{p-1}$ because the support is bounded and the gradient of C_d is radially non-decreasing.

1273 Regarding OT($\theta_{\#}^* \alpha, \theta_{\#}^* \beta$), the bounds (M_{θ}, L_{θ}) could be refined by considering the dependence in $\theta \in \mathbb{S}^{d-1}$. However we
 1274 prove now these constants can be upper-bounded by a finite constant independent of θ . In this setting we consider the cost

$$C_1(\theta^*(x), \theta^*(y)) = |\langle \theta, x - y \rangle|^p \leq \|\theta\|^p \|x - y\|^p \leq \|x - y\|^p,$$

1277 by Cauchy-Schwarz inequality. Therefore, if (α, β) have supports such that $\|x - y\| \leq R$, then $(\theta_{\#}^* \alpha, \theta_{\#}^* \beta)$ also have
 1278 supports bounded by R in \mathbb{R} . Similarly note that the derivative of $h(x) = x^p$ is non-decreasing for $p \geq 1$. Hence the cost
 1279 $C_1(\theta^*(x), \theta^*(y))$ has a bounded derivative, which reads

$$p |\langle \theta, x - y \rangle|^{p-1} \leq p \|\theta\|^{p-1} \|x - y\|^{p-1} \leq p \|x - y\|^{p-1} \leq pR^{p-1}.$$

1283 Thus on the supports of $(\theta_{\#}^* \alpha, \theta_{\#}^* \beta)$ one can also bound the Lipschitz constant of the cost $C_1(x, y) = |x - y|^p$ by the same
 1284 constant L . □

1286 **Remark: Extending Theorem 3.7.** We conjecture that Theorem 3.7 also holds when σ is the uniform measures over \mathbb{S}^{d-1} ,
 1287 since the above holds for any $N \in \mathbb{N}^*$ and $\hat{\sigma}_N$ converges weakly* to σ . Proving this result would require that potentials
 1288 (f_{θ}, g_{θ}) are also regular (i.e., Lipschitz and bounded) w.r.t $\theta \in \mathbb{S}^{d-1}$. This regularity is proved in (Xi & Niles-Weed,
 1289 2022) assuming (α, β) have densities, but remains unknown for discrete measures. Since discretizing σ corresponds to
 1290 the computational approach, we assume it to be discrete, so that no additional assumption than boundedness on (α, β) is
 1291 required. For instance, such result remains valid for semi-discrete UOT computation.

1293 B. Additional details for Section 4

1295 B.1. Frank-Wolfe methodology for computing UOT

1296 **Background: FW for UOT.** Our approach to compute SUOT and USOT builds upon the construction of (Séjourné
 1297 et al., 2022b). It consists in applying a Frank-Wolfe (FW) procedure over the dual formulation of UOT. Such approach is
 1298 equivalent to solve a sequence of balanced OT problems between measures $(\tilde{\alpha}, \tilde{\beta})$ which are iterative renormalizations of
 1299 (α, β) . While the idea holds in wide generality, it is especially efficient in 1D where OT has low algorithmic complexity,
 1300 and we reuse it in our sliced setting.

1302 FW algorithm consists in optimizing a functional \mathcal{H} over a compact, convex set \mathcal{C} by optimizing its linearization $\nabla \mathcal{H}$.
 1303 Given a current iterate x^t of FW algorithm, one computes $r^{t+1} \in \arg \max_{r \in \mathcal{C}} \langle \nabla \mathcal{H}(x^t), r \rangle$, and performs a convex update
 1304 $x^{t+1} = (1 - \gamma_{t+1})x^t + \gamma_{t+1}r^{t+1}$. One typically chooses the learning rate $\gamma_t = \frac{2}{2+t}$. This yields the routine `FWStep` of
 1305 Section 4 which is detailed below.

1307 **Algorithm 3** – `FWStep`(f, g, r, s, γ)

1308 **Input:** $\alpha, \beta, f, g, \gamma$

1309 **Output:** Normalized measures (α, β) as in Equation (80)

1310 $f(x) \leftarrow (1 - \gamma)f(x) + \gamma r(x)$

1311 $g(y) \leftarrow (1 - \gamma)g(y) + \gamma s(y)$

1312 Return (f, g)

1314 In the setting of UOT, one would take $\mathcal{C} = \{f \oplus g \leq C_d\}$. However, this set is not compact as it contains $(\lambda, -\lambda)$
 1315 for any $\lambda \in \mathbb{R}$. Thus, (Séjourné et al., 2022b) propose to optimise a *translation-invariant* dual functional $\mathcal{H}(f, g) \triangleq$
 1316 $\sup_{\lambda \in \mathbb{R}} \mathcal{D}(f + \lambda, g - \lambda)$, with \mathcal{D} defined Equation (3). Similar to the balanced OT dual, one has $\mathcal{H}(f + \lambda, g - \lambda) = \mathcal{H}(f, g)$,
 1317 thus one can apply (Santambrogio, 2015, Proposition 1.11) to assume w.l.o.g. that e.g. $f(0) = 0$ and restrict to a compact
 1318 set of functions. We emphasize that FW algorithm is well-posed to optimize \mathcal{H} , but not \mathcal{D} .
 1319

1320 Note that once we have the dual variables (f, g) maximizing \mathcal{H} , we retrieve optimal dual variables maximizing \mathcal{D} as
 1321 $(f + \lambda^*(f, g), g - \lambda^*(f, g))$, where $\lambda^*(f, g) \triangleq \arg \max_{\lambda \in \mathbb{R}} \mathcal{D}(f + \lambda, g - \lambda)$. The KL setting where $D_{\varphi_1} = \rho_1 \text{KL}$ and
 1322 $D_{\varphi_2} = \rho_2 \text{KL}$ is especially convenient, because $\lambda^*(f, g)$ admits a closed form, which avoids iterative subroutines to compute
 1323 it. In that case, it reads

$$1324 \lambda^*(f, g) = \frac{\rho_1 \rho_2}{\rho_1 + \rho_2} \log \left(\frac{\int e^{-f(x)/\rho_1} d\alpha(x)}{\int e^{-g(y)/\rho_2} d\beta(y)} \right). \quad (77)$$

1325 We summarize the FW algorithm for UOT in the proposition below. We refer to (Séjourné et al., 2022b) for more details on
 1326 the algorithm and pseudo-code. We adapt this approach and result for SUOT and USOT.
 1327

1328 **Proposition B.1.** (Séjourné et al., 2022b) *Assume φ° is smooth. Given current iterates $(f^{(t)}, g^{(t)})$, the linear FW oracle
 1329 of UOT(α, β) is OT($\bar{\alpha}^{(t)}, \bar{\beta}^{(t)}$), where $\bar{\alpha}^{(t)} = \nabla \varphi^\circ(f^{(t)} + \lambda^*(f^{(t)}, g^{(t)}))\alpha$ and $\bar{\beta}^{(t)} = \nabla \varphi^\circ(g^{(t)} - \lambda^*(f^{(t)}, g^{(t)}))\beta$. In
 1330 particular, one has $m(\bar{\alpha}^{(t)}) = m(\bar{\beta}^{(t)})$, thus the balanced OT problem always has finite value. More precisely, the FW
 1331 update reads*

$$1332 (f^{(t+1)}, g^{(t+1)}) = (1 - \gamma^{(t+1)})(f^{(t)}, g^{(t)}) + \gamma^{(t+1)}(r^{(t+1)}, s^{(t+1)}), \quad (78)$$

$$1333 \text{ where } (r^{(t+1)}, s^{(t+1)}) \in \arg \max_{r \oplus s \leq C_d} \int r(x) d\bar{\alpha}^{(t)}(x) + \int s(y) d\bar{\beta}^{(t)}(y). \quad (79)$$

1334 Recall that the in KL setting one has $\varphi_i^\circ(x) = \rho_i(1 - e^{-x/\rho_i})$, thus $\nabla \varphi_i^\circ(x) = e^{-x/\rho_i}$. Thus in that case one normalizes
 1335 the measures as

$$1336 \bar{\alpha} = \exp \left(-\frac{f + \lambda^*(f, g)}{\rho_1} \right) \alpha, \quad \bar{\beta} = \exp \left(-\frac{g - \lambda^*(f, g)}{\rho_2} \right) \beta, \quad (80)$$

1337 where λ^* is defined in (77).

1338 This defines the `NORM` routine in Section 4, which we detail below.

1339 **Algorithm 4** – `NORM`($\alpha, \beta, f, g, \rho_1, \rho_2$)

1340 **Input:** $\alpha, \beta, f, g, \rho = (\rho_1, \rho_2)$

1341 **Output:** Normalized measures (α, β) as in eq. (80)

1342 Compute $\lambda^* = \lambda^*(f, g)$ as in eq. (77)

1343 $\bar{\alpha}(x) \leftarrow \exp \left(-\frac{f(x) + \lambda^*}{\rho_1} \right) \alpha(x)$

1344 $\bar{\beta}(y) \leftarrow \exp \left(-\frac{g(y) - \lambda^*}{\rho_2} \right) \beta(y)$

1345 Return (α, β)

1346 B.2. Frank-Wolfe methodology for computing SUOT

1347 **Proposition B.2.** *Given current iterates (f_θ, g_θ) , the linear Frank-Wolfe oracle of USOT(α, β) is
 1348 $\int_{\mathbb{S}^{d-1}} \text{OT}(\theta_\#^* \alpha^\theta, \theta_\#^* \beta^\theta) d\sigma(\theta)$, where*

$$1349 \alpha^\theta = \nabla \varphi^\circ \left(f_\theta + \lambda^*(f_\theta, g_\theta) \right) \alpha, \quad \beta^\theta = \nabla \varphi^\circ \left(g_\theta - \lambda^*(f_\theta, g_\theta) \right) \beta.$$

1350 As a consequence, given dual sliced potentials (r_θ, s_θ) solving OT($\theta_\#^* \alpha^\theta, \theta_\#^* \beta^\theta$), one can perform Frank-Wolfe updates (78)
 1351 on (f_θ, g_θ) .

1352 *Proof.* Our goal is to compute the first order variation of the SUOT functional. Given that $\text{SUOT}(\alpha, \beta) =$
 1353 $\int_{\mathbb{S}^{d-1}} \text{UOT}(\theta_\#^* \alpha, \theta_\#^* \beta) d\sigma(\theta)$, one can apply Proposition B.1 slice-wise. Since measures are assumed to have compact
 1354

1375 support, one can apply the dominated convergence theorem and differentiate under the integral sign. Furthermore, the
 1376 translation-invariant formulation in the setting of SUOT reads

$$1377 \text{SUOT}(\alpha, \beta) = \int_{\mathbb{S}^{d-1}} \sup_{f_\theta \oplus g_\theta \leq C_1} \left[\sup_{\lambda_\theta \in \mathbb{R}} \int \varphi^\circ(f_\theta(\cdot) + \lambda_\theta) d\theta_\#^* \alpha \right. \quad (81)$$

$$1380 \left. + \int \varphi^\circ(g_\theta(\cdot) - \lambda_\theta) d\theta_\#^* \beta \right], \quad (82)$$

1384 In the setting where φ° is smooth and strictly concave (such as $D_\varphi = \rho_{\text{KL}}$), there always exists a unique optimal λ_θ^* .
 1385 Furthermore, one can apply the envelope theorem such that the Fréchet differential w.r.t. to a perturbation (r_θ, s_θ) of (f_θ, g_θ)
 1386 reads

$$1387 \int_{\mathbb{S}^{d-1}} \left[\int r_\theta(\cdot) \times \nabla \varphi^\circ(f_\theta(\cdot) + \lambda_\theta^*(f_\theta, g_\theta)) d\theta_\#^* \alpha \right. \quad (83)$$

$$1390 \left. + \int s_\theta(\cdot) \times \nabla \varphi^\circ(g_\theta(\cdot) - \lambda_\theta^*(f_\theta, g_\theta)) d\theta_\#^* \beta \right] \quad (84)$$

1394 **Setting**

$$1395 \alpha_\theta = \nabla \varphi^\circ(f_\theta(\cdot) + \lambda^*(f_\theta, g_\theta)) \alpha, \quad \beta_\theta = \nabla \varphi^\circ(g_\theta(\cdot) - \lambda^*(f_\theta, g_\theta)) \beta,$$

1398 yields the desired result, *i.e.* the first order variation is

$$1400 \int_{\mathbb{S}^{d-1}} \left[\int r_\theta(\cdot) d(\theta_\#^* \alpha_\theta) + \int s_\theta(\cdot) d(\theta_\#^* \beta_\theta) \right]. \quad (85)$$

1404 \square

1406 B.3. Frank-Wolfe methodology for computing USOT

1407 To compute USOT, we leverage Theorem 3.7 and derive the linear Frank-Wolfe oracle based on its translation-invariant
 1408 formulation. We state the associated FW updates in the following proposition.

1410 **Proposition B.3.** *Given current iterates (f_θ, g_θ) , the linear Frank-Wolfe oracle of $\text{USOT}(\alpha, \beta)$ is $\text{SOT}(\bar{\alpha}, \bar{\beta})$, where*

$$1411 \bar{\alpha} = \nabla \varphi^\circ(f_{\text{avg}} + \lambda^*(f_{\text{avg}}, g_{\text{avg}})) \alpha, \quad \bar{\beta} = \nabla \varphi^\circ(g_{\text{avg}} - \lambda^*(f_{\text{avg}}, g_{\text{avg}})) \beta,$$

$$1412 f_{\text{avg}}(x) = \int_{\mathbb{S}^{d-1}} f_\theta(\theta^*(x)) d\hat{\sigma}_K(\theta), \quad g_{\text{avg}}(y) = \int_{\mathbb{S}^{d-1}} g_\theta(\theta^*(y)) d\hat{\sigma}_K(\theta)$$

1413 Thus given dual sliced potentials $(r_\theta(\cdot), s_\theta(\cdot))$ which solve $\text{SOT}(\bar{\alpha}, \bar{\beta})$, one can then perform Frank-Wolfe updates (78) on
 1414 (f_θ, g_θ) and thus $(f_{\text{avg}}, g_{\text{avg}})$.

1419 *Proof.* Our goal is to compute the first order variation of the USOT functional. First, we leverage Theorem 3.7 such that
 1420 USOT reads

$$1421 \text{USOT}(\alpha, \beta) = \sup_{f_\theta(\cdot) \oplus g_\theta(\cdot) \leq C_1} \int \varphi_1^\circ \left(\int_{\mathbb{S}^{d-1}} f_\theta(\theta^*(x)) d\hat{\sigma}_K(\theta) \right) d\alpha(x) \quad (86)$$

$$1422 + \int \varphi_2^\circ \left(\int_{\mathbb{S}^{d-1}} g_\theta(\theta^*(y)) d\hat{\sigma}_K(\theta) \right) d\beta(y) \quad (87)$$

$$1423 = \sup_{f_\theta(\cdot) \oplus g_\theta(\cdot) \leq C_1} \int \varphi_1^\circ(f_{\text{avg}}(x)) d\alpha(x) + \int \varphi_2^\circ(g_{\text{avg}}(y)) d\beta(y), \quad (88)$$

1429

1430 where

$$1431 \quad f_{avg}(x) = \int_{\mathbb{S}^{d-1}} f_{\theta}(\theta^*(x)) d\hat{\sigma}_K(\theta), \quad 1432 \quad g_{avg}(y) = \int_{\mathbb{S}^{d-1}} g_{\theta}(\theta^*(y)) d\hat{\sigma}_K(\theta). \quad 1433$$

1434 From this, we derive the translation-invariant formulation as follows.

$$1435 \quad \text{USOT}(\alpha, \beta) = \sup_{f_{\theta}(\cdot) \oplus g_{\theta}(\cdot) \leq C_1} \sup_{\lambda \in \mathbb{R}} \int \varphi_1^{\circ}(f_{avg}(x) + \lambda) d\alpha(x) \quad (89)$$

$$1436 \quad + \int \varphi_2^{\circ}(g_{avg}(y) - \lambda) d\beta(y), \quad (90)$$

1437 For smooth and strictly concave φ° , there exists a unique $\lambda^*(f_{avg}, g_{avg})$ attaining the supremum. Furthermore, one can
1438 apply the envelope theorem and differentiate under the integral sign (since the support is compact). Consider perturbations
1439 $(r_{\theta}(\cdot), s_{\theta}(\cdot))$ of $(f_{\theta}(\cdot), g_{\theta}(\cdot))$. Write

$$1440 \quad r_{avg}(x) = \int_{\mathbb{S}^{d-1}} r_{\theta}(\theta^*(x)) d\hat{\sigma}_K(\theta), \quad 1441 \quad s_{avg}(y) = \int_{\mathbb{S}^{d-1}} s_{\theta}(\theta^*(y)) d\hat{\sigma}_K(\theta).$$

1442 Given that $\varphi_1^{\circ}(f_{avg} + r_{avg}) = \varphi_1^{\circ}(f_{avg}) + r_{avg} \nabla \varphi_1^{\circ}(f_{avg}) + o(\|r_{avg}\|_{\infty})$, the first order variation reads

$$1443 \quad \int r_{avg}(x) \nabla \varphi_1^{\circ}(f_{avg}(x) + \lambda^*(f_{avg}, g_{avg})) d\alpha(x) \quad (91)$$

$$1444 \quad + \int s_{avg}(y) \nabla \varphi_2^{\circ}(g_{avg}(y) - \lambda^*(f_{avg}, g_{avg})) d\beta(y). \quad (92)$$

1445 Then we define

$$1446 \quad \bar{\alpha} = \nabla \varphi_1^{\circ}(f_{avg} + \lambda^*(f_{avg}, g_{avg})) \alpha, \quad 1447 \quad \bar{\beta} = \nabla \varphi_2^{\circ}(g_{avg} - \lambda^*(f_{avg}, g_{avg})) \beta,$$

1448 such that the first order variation reads

$$1449 \quad \int r_{avg}(x) d\bar{\alpha}(x) + \int s_{avg}(y) d\bar{\beta}(y). \quad (93)$$

1450 One can then explicit the definition of (r_{avg}, s_{avg}) , such that it reads

$$1451 \quad \int_{\mathbb{S}^{d-1}} \int r_{\theta}(\theta^*(x)) d\bar{\alpha}(x) + \int_{\mathbb{S}^{d-1}} \int s_{\theta}(\theta^*(y)) d\bar{\beta}(y) \quad (94)$$

$$1452 \quad = \int_{\mathbb{S}^{d-1}} \int r_{\theta} d\theta_{\#}^* \bar{\alpha}(x) + \int_{\mathbb{S}^{d-1}} \int s_{\theta} d\theta_{\#}^* \bar{\beta}(y). \quad (95)$$

1453 By optimizing the above over the constraint set $\{r_{\theta} \oplus s_{\theta} \leq C_1\}$, we identify the computation of $\text{SOT}(\bar{\alpha}, \bar{\beta})$, which concludes
1454 the proof. \square

1455 Since Proposition B.3 involves potentials averaged over σ , we thus need to define the `AvgPot` routine detailed below.

1456 **Algorithm 5** – `AvgPot`(f_{θ})

1457 **Input:** sliced potentials (f_{θ}) with $(\theta_k)_{k=1}^K$

1458 **Output:** Averaged potential f_{avg} as in Proposition B.3

$$1459 \quad \text{Average } f_{avg} = \frac{1}{K} \sum_{k=1}^K f_{\theta}$$

1484

B.4. Implementation of Sliced OT to return dual potentials

Recall from Section 4, Algorithms 1 and 2 and more precisely, Propositions B.2 and B.3, that FW linear oracle is a sliced OT program, *i.e.* a set of OT problems computed between univariate distributions of $\mathcal{M}_+(\mathbb{R})$. Therefore, a key building block of our algorithm is to compute the loss and dual variables of these univariate OT problems. We explain below how one can compute the sliced OT loss and dual potentials. The computation of the loss consists in implementing closed formulas of OT between univariate distributions, as detailed in (Santambrogio, 2015, Proposition 2.17). More precisely, when $C_1(x, y) = |x - y|^p$ and $(\mu, \nu) \in \mathcal{M}_+(\mathbb{R})$, then

$$\text{OT}(\mu, \nu) = \int_0^1 |F_\mu^{[-1]}(t) - F_\nu^{[-1]}(t)|^p dt, \quad (96)$$

where $F_\mu^{[-1]}$ denotes the inverse cumulative distribution function (ICDF) of μ .

Algorithm 6 – SlicedOTLoss($\alpha, \beta, \{\theta\}, p$)

Input: α, β , projections $\{\theta\}$, exponent p

Output: $\text{OT}(\theta_\#^* \alpha, \theta_\#^* \beta)$ as in eq. (96)

for $\theta \in \{\theta\}$ **do**

 Project support of $\theta_\#^* \alpha$ and $\theta_\#^* \beta$

 Sort weights of $(\theta_\#^* \alpha, \theta_\#^* \beta)$ and support $(\theta^*(x)), (\theta^*(y))$ s.t. support is non-decreasing

 Compute ICDF of $\theta_\#^* \alpha$ and $\theta_\#^* \beta$

 Compute $\text{OT}(\theta_\#^* \alpha, \theta_\#^* \beta)$ as in eq. (96) with exponent p

end for

To compute dual potentials using backpropagation, one computes the sliced OT losses (using Algorithm 6) then calls the backpropagation w.r.t to inputs (α, β) , because their gradients are optimal dual potentials (Santambrogio, 2015, Proposition 7.17). We describe this procedure in Algorithm 7.

Algorithm 7 – SlicedOTPotentialsBackprop($\alpha, \beta, \{\theta\}, p$)

Input: α, β , projections $\{\theta\}$, exponent p

Output: Dual potentials (f_θ, g_θ) solving $\text{OT}(\theta_\#^* \alpha, \theta_\#^* \beta)$

 Enable gradients w.r.t. $(\theta_\#^* \alpha, \theta_\#^* \beta)$

 Call $\text{SlicedOTLoss}(\alpha, \beta, \{\theta\}, p)$

 Sum (but do not average) losses $\mathcal{L} = \sum_\theta \text{OT}(\theta_\#^* \alpha, \theta_\#^* \beta)$.

 Backpropagate \mathcal{L} w.r.t. (α, β)

 Return (f_θ, g_θ) as gradients of \mathcal{L} w.r.t. (α, β) .

The implementation of the dual potentials using 1D closed forms relies on the north-west corner rule principle, which can be vectorized in PyTorch in order to be computed in parallel. The contribution of our implementation thus consists in making such algorithm GPU-compatible and allowing for a parallel computation for every slice simultaneously. We stress that this constitutes a non-trivial piece of code, and we refer the interested reader to the code in our supplementary material for more details on the implementation.

B.5. Output optimal sliced marginals

In all our algorithms, we focus on dual formulations of SUOT and USOT, which optimize the dual potentials. However, one might want the output variables of the primal formulation (See Definition 3.1). In particular, the marginals of optimal transport plans are interesting because they are interpreted as normalized versions of inputs (α, β) where geometric outliers have been removed. We detail where this interpretation comes from in the setting of UOT, and then give how it is adapted to SUOT and USOT. In particular, we justify that the `NORM` routine suffices to compute them.

Case of UOT. We focus on the $D_{\varphi_i} = \rho_i \text{KL}$. As per (Liero et al., 2018, Equation 4.21), we have at optimality that the optimal transport π^* plan solving $\text{UOT}(\alpha, \beta)$ as in Equation (2) has marginals (π_1^*, π_2^*) which read $\pi_1^* = e^{-f^*/\rho_1} \alpha$ and $\pi_2^* = e^{-g^*/\rho_2} \beta$, where (f^*, g^*) are the optimal dual potentials solving Equation (3). Since on $\text{supp}(\pi^*)$ one also has $f^*(x) + g^*(y) = C_d(x, y)$, if the transportation cost $C_d(x, y)$ is large (*i.e.* we are matching a geometric outlier), so are $f^*(x)$

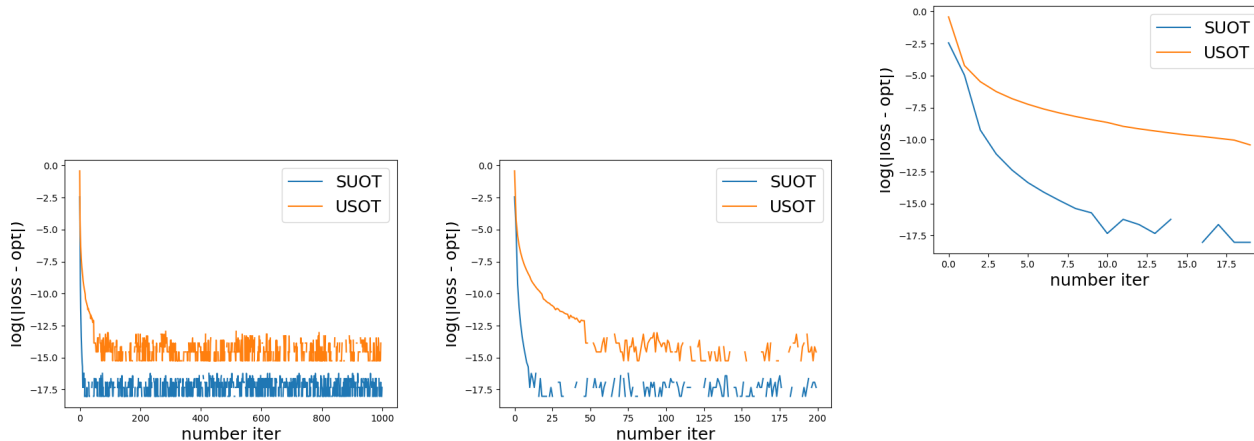


Figure 5: $|\text{SUOT}(\alpha, \beta) - \widehat{\text{SUOT}}_t|$ and $|\text{USOT}(\alpha, \beta) - \widehat{\text{USOT}}_t|$ against iteration t , where $\widehat{\text{SUOT}}_t, \widehat{\text{USOT}}_t$ are the estimated SUOT, USOT using t FW iterations. Plots are in log-scale. All figures are issued from the same run, but zoomed on a subset of first iterations: (left) 1000 iterations of FW, (middle) 200 iterations, (right) 20 iterations.

and $g^*(y)$, and eventually the weights $\pi_1^*(x)$ and $\pi_2^*(y)$ are small, hence the interpretation of the geometric normalization of the measures. Note that in that case, one obtain (π_1^*, π_2^*) by calling $\text{NORM}(\alpha, \beta, f^*, g^*, \rho_1, \rho_2)$.

Case of SUOT. Since $\text{SUOT}(\alpha, \beta)$ consists in integrating $\text{UOT}(\theta_{\#}^* \alpha, \theta_{\#}^* \beta)$ w.r.t. σ , it shares many similarities with UOT. For any θ , we consider π_{θ} and (f_{θ}, g_{θ}) solving the primal and dual formulation of $\text{UOT}(\theta_{\#}^* \alpha, \theta_{\#}^* \beta)$. The marginals of π_{θ} are thus given by $(e^{-f_{\theta}/\rho_1} \alpha, e^{-g_{\theta}/\rho_2} \beta)$. In particular, we retrieve the observation made in Figure 1 that the optimal marginals change for each θ . In that case we call for each θ the routine $\text{NORM}(\alpha, \beta, f_{\theta}, g_{\theta}, \rho_1, \rho_2)$.

Case of USOT. Recall that the optimal marginals (π_1, π_2) in $\text{USOT}(\alpha, \beta)$ do not depend on θ , contrary to $\text{SUOT}(\alpha, \beta)$. Leveraging the dual formulation of Theorem 3.7, and looking at the Lagrangian which is defined in the proof of Theorem 3.7 (see Appendix A.7), we have the optimality condition that $\pi_1 = e^{-f_{\text{avg}}/\rho_1} \alpha$ and $\pi_2 = e^{-g_{\text{avg}}/\rho_2} \beta$. Thus in that case, calling $\text{NORM}(\alpha, \beta, f_{\text{avg}}, g_{\text{avg}}, \rho_1, \rho_2)$ yields the desired marginals.

B.6. Convergence of Frank-Wolfe iterations: Empirical analysis

We display below an experiment on synthetic dataset to illustrate the convergence of Frank-Wolfe iterations. We also provide insights on the number of iterations that yields a reasonable approximation: a few iterations suffices in our practical settings, typically $F = 20$.

The results are displayed in Figure 5. We consider the empirical distributions (α, β) computed over respectively, $N = 400$ and $M = 500$ samples over the unit hypercube $[0, 1]^d$, $d = 10$. Moreover, β is slightly shifted by a vector of uniform coordinates $0.5 \times \mathbf{1}_d$. We choose $\rho = 1$ and report the estimation of $\text{SUOT}(\alpha, \beta)$ and $\text{USOT}(\alpha, \beta)$ through Frank-Wolfe iterations. We estimate the true values by running $F = 5000$ iterations, and display the difference between the estimated score and the 'true' values. Appendix B.6 shows that numerical precision is reached in a few tens of iterations. As learning tasks do not usually require an estimation of losses up to numerical precision, we think that it is hence reasonable to take $F \approx 20$ in numerical applications.

C. Additional details on Section 5

C.1. Document classification: Technical details and additional results

C.1.1. DATASETS

We sum up the statistics of the different datasets in Table 2.

Table 2: Dataset characteristics.

	BBCSport	Movies	Goodreads genre	Goodreads like
Doc	737	2000	1003	1003
Train	517	1500	752	752
Test	220	500	251	251
Classes	5	2	8	2
Mean words by doc	116 ± 54	182 ± 65	1491 ± 538	1491 ± 538
Median words by doc	104	175	1518	1518
Max words by doc	469	577	3499	3499

BBCSport. The BBCSport dataset contains articles between 2004 and 2005, and is composed of 5 classes. We average over the 5 same train/test split of (Kusner et al., 2015). The dataset can be found in <https://github.com/mkusner/wmd/tree/master>.

Movie Reviews. The movie reviews dataset is composed of 1000 positive and 1000 negative reviews. We take five different random 75/25 train/test split. The data can be found in <http://www.cs.cornell.edu/people/pabo/movie-review-data/>.

Goodreads. This dataset, proposed in (Maharjan et al., 2017), and which can be found at https://ritual.uh.edu/multi_task_book_success_2017/, is composed of 1003 books from 8 genres. A first possible classification task is to predict the genre. A second task is to predict the likability, which is a binary task where a book is said to have success if it has an average rating ≥ 3.5 on the website Goodreads (<https://www.goodreads.com>). The five train/test split are randomly drawn with 75/25 proportions.

C.1.2. TECHNICAL DETAILS

All documents are embedded with the Word2Vec model (Mikolov et al., 2013) in dimension $d = 300$. The embedding can be found in <https://drive.google.com/file/d/0B7XkCwpI5KDYN1NUTTL1SS21pQmM/view?resourcekey=0-wjGZdNAUop6WYkTtMip30g>.

In this experiment, we report the results averaged over 5 random train/test split. For discrepancies which are approximated using random projections, we additionally average the results over 3 different computations, and we report this standard deviation in Table 1. Furthermore, we always use 500 projections to approximate the sliced discrepancies. For Frank-Wolfe based methods, we use 10 iterations, which we found to be enough to have a good accuracy. We added an ablation of these two hyperparameters in Figure 7. We report the results obtained with the best ρ for USOT and SUOT computed among a grid $\rho \in \{10^{-4}, 5 \cdot 10^{-4}, 10^{-3}, 5 \cdot 10^{-3}, 10^{-2}, 10^{-1}, 1\}$. For USOT, the best ρ is consistently $5 \cdot 10^{-3}$ for the Movies and Goodreads datasets, and $5 \cdot 10^{-4}$ for the BBCSport dataset. For SUOT, the best ρ obtained was 0.01 for the BBCSport dataset, 1.0 for the movies dataset and 0.5 for the goodreads dataset. For UOT, we used $\rho = 1.0$ on the BBCSport dataset. For the movies dataset, the best ρ obtained on a subset was 50, but it took an unreasonable amount of time to run on the full dataset as the runtime increases with ρ (see (Chapel et al., 2021, Figure 3)). On the goodreads dataset, it took too much memory on the GPU. For Sinkhorn UOT, we used $\varepsilon = 0.001$ and $\rho = 0.1$ on the BBCSport and Goodreads datasets, and $\varepsilon = 0.01$ on the Movies dataset. For each method, the number of neighbors used for the k-NN method is obtained via cross-validation.

C.1.3. ADDITIONAL EXPERIMENTS

Runtime. We report in Figure 6 the runtime of computing the different discrepancies between each pair of documents. On the BBCSport dataset, the documents have in average 116 words, thus the main bottleneck is the projection step for sliced OT methods. Hence, we observe that OT runs slightly faster than SOT and the sliced unbalanced counterparts. Goodreads is a dataset with larger documents, with on average 1491 words by document. Therefore, as OT scales cubically with the number of samples, we observe here that all sliced methods run faster than OT, which confirms that sliced methods scale better w.r.t. the number of samples. In this setting, we were not able to compute UOT with the POT implementation in a reasonable time. Computations have been performed with a NVIDIA A100 GPU.

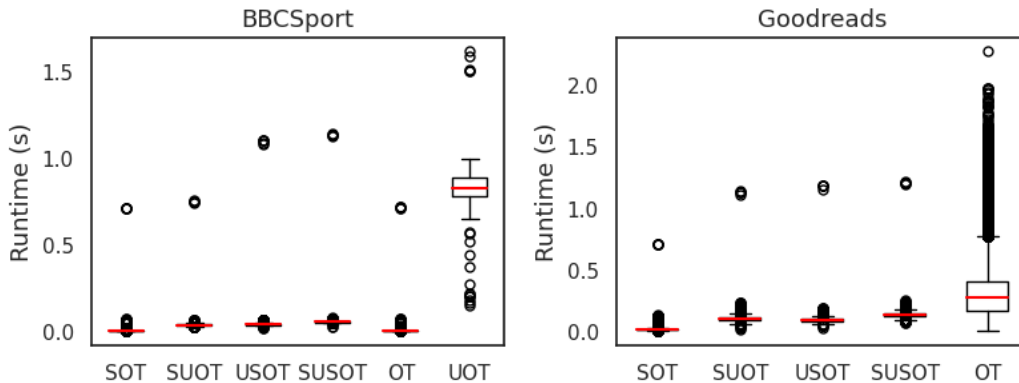


Figure 6: Runtime on the BBCSport dataset (*left*) and on the Goodreads dataset (*right*).

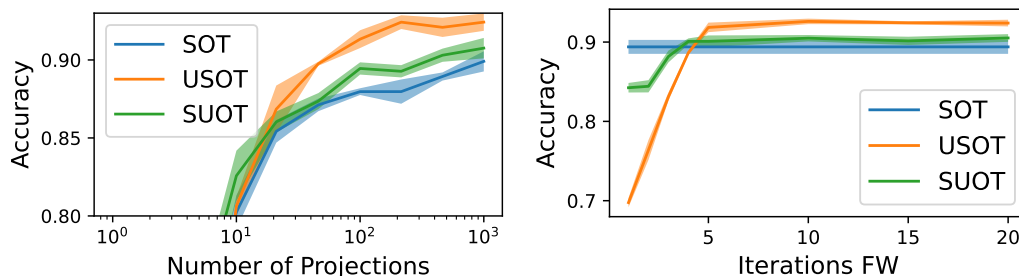


Figure 7: Ablation on BBCSport of the number of projections (*left*) and of the number of Frank-Wolfe iterations (*right*).

Ablations. We plot in Figure 7 accuracy as a function of the number of projections and the number of iterations of the Frank-Wolfe algorithm. We averaged the accuracy obtained with the same setting described in Appendix C.1.2, with varying number of projections $K \in \{4, 10, 21, 46, 100, 215, 464, 1000\}$ and number of FW iterations $F \in \{1, 2, 3, 4, 5, 10, 15, 20\}$. Regarding the hyperparameter ρ , we selected the one returning the best accuracy, *i.e.* $\rho = 5 \cdot 10^{-4}$ for USOT and $\rho = 10^{-2}$ for SUOT.

C.2. Unbalanced sliced Wasserstein barycenters

We define below the formulation of the USOT barycenter which was used in the experiments of Figure 4 to average predictions of geophysical data. We then detail how we computed it.

Definition C.1. Consider a set of measures $(\alpha_1, \dots, \alpha_B) \in \mathcal{M}_+(\mathbb{R}^d)^B$, and a set of non-negative coefficients $(\omega_1, \dots, \omega_B) \geq 0$ such that $\sum_{b=1}^B \omega_b = 1$. We define the barycenter problem (in the KL setting) as

$$\mathcal{B}((\alpha_b)_b, (\omega_b)_b) \triangleq \inf_{\beta \in \mathcal{P}(\mathbb{R}^d)} \sum_{b=1}^B \omega_b \text{USOT}(\alpha_b, \beta), \quad (97)$$

$$= \inf_{\beta \in \mathcal{P}(\mathbb{R}^d)} \sum_{b=1}^B \inf_{(\pi_{b,1}, \pi_{b,2})} \text{SOT}(\pi_{b,1}, \pi_{b,2}) + \rho_1 \text{KL}(\pi_{b,1} | \alpha_b) + \rho_2 \text{KL}(\pi_{b,2} | \beta), \quad (98)$$

where $\mathcal{P}(\mathbb{R}^d)$ denotes the set of probability measures.

To compute the barycenter, we aggregate several building blocks. First, since we consider that the barycenter $\beta \in \mathcal{P}(\mathbb{R}^d)$ is a probability, we perform mirror descent as in (Beck & Teboulle, 2003; Cuturi & Doucet, 2014b). More precisely, we use a Nesterov accelerated version of mirror descent. We also tried projected gradient descent, but it did not yield consistent outputs (due to convergence speed (Beck & Teboulle, 2003)). Second, we use a Stochastic-USOT version (see Section 4), *i.e.* we sample new projections at each iteration of the barycenter update (but not a each iteration of the FW subroutines in Algorithm 2). This procedure is described in Algorithm 8.

```

1705 Algorithm 8 – Barycenter $((\alpha_b)_b, (\omega_b)_b, \rho_1, \rho_2, lr)$ 
1706 Input: measures  $(\alpha_b)_b$ , weights  $(\omega_b)_b$ ,  $\rho_1, \rho_2$ , learning rate  $lr$ , FW iter  $F$ 
1707 Output: Optimal barycenter  $\beta$  of Equation (97)
1708    $t \leftarrow 1$ 
1709   Init  $(\beta, \tilde{\beta}, \hat{\beta})$  as uniform distribution over a grid
1710   while not converged do do
1711      $\gamma \leftarrow \frac{2}{(t+1)}$ ,
1712      $\beta \rightarrow (1 - \gamma)\hat{\beta} + \gamma\tilde{\beta}$ 
1713     Sample projections  $(\theta_k)_{k=1}^K$ 
1714     Compute  $\mathcal{B}((\alpha_b)_b, (\omega_b)_b)$  by calling USOT $(\alpha_b, \beta, F, (\theta_k)_{k=1}^K, \rho_1, \rho_2)$  in Algorithm 2 for each  $b$ 
1715     Compute  $g$  as the gradient of  $\mathcal{B}((\alpha_b)_b, (\omega_b)_b)$  w.r.t. variable  $\beta$ 
1716      $\tilde{\beta} \leftarrow \exp(-lr \times \gamma^{-1} \times g)\beta$ 
1717      $\hat{\beta} \leftarrow \tilde{\beta}/m(\tilde{\beta})$ 
1718      $\hat{\beta} \leftarrow (1 - \gamma)\hat{\beta} + \gamma\tilde{\beta}$ 
1719      $t \leftarrow t + 1$ 
1719   end while

```

```

1720
1721
1722 We illustrate this algorithm with several examples of interpolation in Figure 8. We propose to compute an interpolation
1723 between two measures located on a fixed grid of size  $200 \times 200$  with different values of  $\rho_i$  in  $D_{\varphi_i} = \rho_i \text{KL}$ . For illustration
1724 purposes, we construct the source distribution as a mixture of two Gaussians with a small and a larger mode, and the target
1725 distribution as a single Gaussian. Those distributions are normalized over the grid such that both total norms are equal to
1726 one (which is not required by our unbalanced sliced variants but grants more interpretability and possible comparisons with
1727 SOT). Figure 8a shows the result of the interpolation at three timestamps ( $t = 0.25, 0.5$  and  $0.75$ ) of a SOT interpolation
1728 (within this setting,  $\omega_1 = 1 - t$  and  $\omega_2 = t$ ). As expected, the two modes of the source distribution are transported over the
1729 target one. We verify in Figure 8b that for a large value of  $\rho_1 = \rho_2 = 100$ , the USOT interpolation behaves similarly as
1730 SOT, as expected from the theory. When  $\rho_1 = \rho_2 = 0.01$ , the smaller mode is not moved during the interpolation, whereas
1731 the larger one is stretched toward the target (Figure 8c). Finally, in Figure 8d, an asymmetric configuration of  $\rho_1 = 0.01$  and
1732  $\rho_2 = 100$  allows to get an interpolation when only the big mode of the source distribution is displaced toward the target. In
1733 all those cases, the mirror-descent algorithm 8 is run for 500 iterations. Even for a large grid of  $200 \times 200$ , those different
1734 results are obtained in a 2 – 3 minutes on a commodity GPU, while the OT or UOT barycenters are untractable with a
1735 limited computational budget.

```

```

1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759

```


1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814

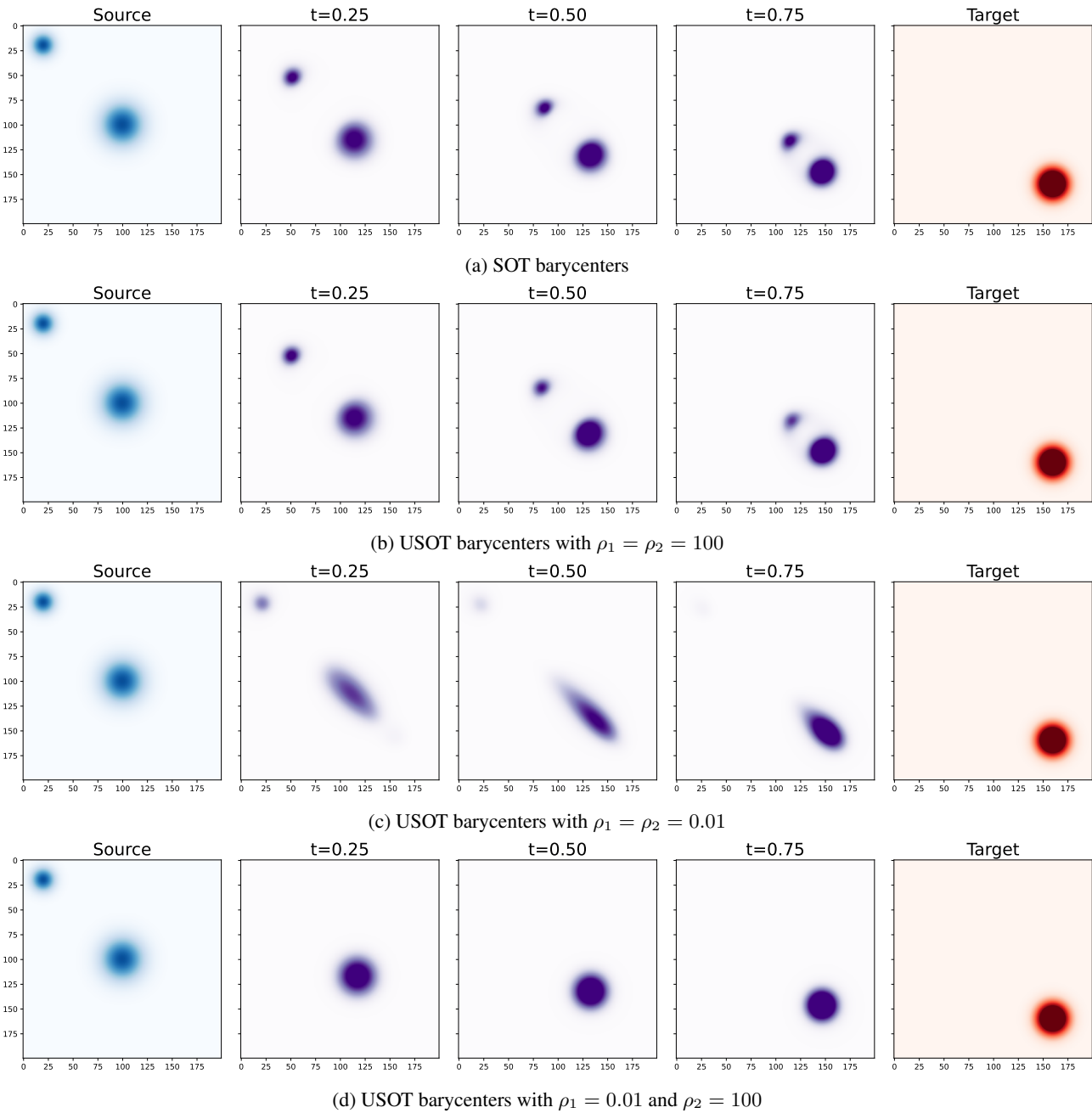


Figure 8: **Interpolation with USOT as a barycenter computation.** We compare different interpolations using SOT or USOT with different settings for the ρ values