

# PROJECTION OPTIMAL TRANSPORT ON TREE-ORDERED LINES

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Many variants of Optimal Transport (OT) have been developed to address its heavy computation. Among them, notably, Sliced Wasserstein (SW) is widely used for application domains by projecting the OT problem onto one-dimensional lines, and leveraging the closed-form expression of the univariate OT to reduce the computational burden. However, projecting measures onto low-dimensional spaces can lead to a loss of topological information. To mitigate this issue, in this work, we propose to replace one-dimensional lines with a more intricate structure, called *tree systems*. This structure is metrizable by a tree metric, which yields a closed-form expression for OT problems on tree systems. We provide an extensive theoretical analysis to formally define tree systems with their topological properties, introduce the concept of splitting maps, which operate as the projection mechanism onto these structures, then finally propose a novel variant of Radon transform for tree systems and verify its injectivity. This framework leads to an efficient metric between measures, termed Tree-Sliced Wasserstein distance on Systems of Lines (TSW-SL). By conducting a variety of experiments on gradient flows, image style transfer, and generative models, we illustrate that our proposed approach performs favorably compared to SW and its variants.

## 1 INTRODUCTION

Optimal transport (OT) (Villani, 2008; Peyré et al., 2019) is a naturally geometrical metric for comparing probability distributions. Intuitively, OT lifts the ground cost metric among supports of input measures into the metric between two input measures. OT has been applied in many research fields, including machine learning (Bunne et al., 2022; Takezawa et al., 2022; Fan et al., 2022; Hua et al., 2023; Nguyen & Ho, 2024), statistics (Mena & Niles-Weed, 2019; Weed & Berthet, 2019; Liu et al., 2022; Nguyen et al., 2022; Nietert et al., 2022; Wang et al., 2022; Pham et al., 2024), multimodal (Park et al., 2024; Luong et al., 2024), computer vision and graphics (Rabin et al., 2011; Solomon et al., 2015; Lavenant et al., 2018; Nguyen et al., 2021; Saleh et al., 2022).

However, OT has a supercubic computational complexity concerning the number of supports in input measures (Peyré et al., 2019). To address this issue, Sliced-Wasserstein (SW) (Rabin et al., 2011; Bonneel et al., 2015) exploits the closed-form expression of the one-dimensional OT to reduce its computational complexity. More concretely, SW projects supports of input measures onto a random line and leverage the fast computation of the OT on one-dimensional lines. SW is widely used in various applications, such as gradient flows (Bonet et al., 2021; Liutkus et al., 2019), clustering (Kolouri et al., 2018; Ho et al., 2017), domain adaptation (Courty et al., 2017), generative models (Deshpande et al., 2018; Wu et al., 2019; Nguyen & Ho, 2022), thanks to its computational efficiency. Due to relying on one-dimensional projection, SW limits its capacity to capture the topological structures of input measures, especially in high-dimensional domains.

**Related work.** Prior studies have aimed to enhance the Sliced Wasserstein (SW) distance (Nguyen et al., 2024a; 2020; Nguyen & Ho, 2024) or explore variants of SW (Bai et al., 2023; Kolouri et al., 2019; Quellmalz et al., 2023). These works primarily concentrate on improving existing components within the SW framework, including the sampling process (Nguyen et al., 2024a; 2020; Nadjahi et al., 2021), determining optimal lines for projection (Deshpande et al., 2019), and modifying the projection mechanism (Kolouri et al., 2019; Bonet et al., 2023). However, few studies have focused

on replacing one-dimensional lines, which play the role of integration domains, with more complex domains such as one-dimensional manifolds (Kolouri et al., 2019), or low-dimensional subspaces (Alvarez-Melis et al., 2018; Bonet et al., 2023; Paty & Cuturi, 2019; Niles-Weed & Rigollet, 2022; Lin et al., 2021; Huang et al., 2021; Muzellec & Cuturi, 2019). In this paper, we concentrate on the latter approach, aiming to discover novel geometrical domains that meet *two key criteria*: (i) pushing forward of high-dimensional measures onto these domains can be processed in a meaningful manner, and (ii) OT problems on these domains can be efficiently solved, ideally with a closed-form solution.

**Contribution.** In summary, our contributions are three-fold:

1. We introduce the concept of tree systems, which consist of copies of the real line equipped with additional structures, and study their topology. A key property of tree systems is that they form well-defined metric spaces, with metrics being tree metrics. This property is sufficient to guarantee that OT problems on tree systems admit closed-form solutions.
2. We define the space of integrable functions and probability measures on a tree system, and introduce a novel transform, called *Radon Transform on Systems of Lines*. This transform naturally transforms measures supported in high-dimensional space onto tree systems, and is a generalization of the original Radon transform. The injectivity of this variant holds, similar to other Radon transform variants in the literature.
3. We propose the Tree-Sliced Wasserstein distance on Systems of Lines (TSW-SL), and analyze its efficiency through the closed-form solution for the OT problem on tree systems, achieving a similar computational cost as the traditional SW.

**Organization.** The remainder of the paper is organized as follows. Section 2 provides necessary backgrounds of SW distance and Wasserstein distance on tree metric spaces. Section 3 provides a brief and intuitive introduction of tree systems and studies its properties, and Section 4 introduces the Radon Transform on System of Lines. The novel Tree-Sliced Wasserstein distance on Systems of Lines is proposed in Section 5. Finally, Section 6 contains empirical results for TSW-SL. Formal constructions, theoretical proofs of key results, and additional materials are presented in Appendix.

## 2 PRELIMINARIES

In this section, we review Sliced Wasserstein (SW) distance and Wasserstein distances on metric spaces with tree metrics (TW).

**Wasserstein Distance.** Let  $\Omega$  be a measurable space with a metric  $d$  on  $\Omega$ , and let  $\mu, \nu$  be two probability distributions on  $\Omega$ . Let  $\mathcal{P}(\mu, \nu)$  be the set of probability distributions  $\pi$  on the product space  $\Omega \times \Omega$  such that  $\pi(A \times \Omega) = \mu(A)$ ,  $\pi(\Omega \times B) = \nu(B)$  for all measurable sets  $A, B$ . For  $p \geq 1$ , the  $p$ -Wasserstein distance  $W_p$  between  $\mu$  and  $\nu$  (Villani, 2008) is defined as:

$$W_p(\mu, \nu) = \inf_{\pi \in \mathcal{P}(\mu, \nu)} \left( \int_{\Omega \times \Omega} d(x, y)^p d\pi(x, y) \right)^{\frac{1}{p}}. \quad (1)$$

**Sliced Wasserstein Distance.** For  $\mu, \nu \in \mathcal{P}(\mathbb{R}^d)$ , the Sliced  $p$ -Wasserstein distance (SW) (Bonneel et al., 2015) between  $\mu, \nu$  is defined by:

$$SW_p(\mu, \nu) := \left( \int_{\mathbb{S}^{d-1}} W_p^p(\mathcal{R}f_\mu(\cdot, \theta), \mathcal{R}f_\nu(\cdot, \theta)) d\sigma(\theta) \right)^{\frac{1}{p}}, \quad (2)$$

where  $\sigma = \mathcal{U}(\mathbb{S}^{d-1})$  is the uniform distribution on  $\mathbb{S}^{d-1}$ , operator  $\mathcal{R} : L^1(\mathbb{R}^d) \rightarrow L^1(\mathbb{R} \times \mathbb{S}^{d-1})$  is the Radon Transform (Helgason & Helgason, 2011)

$$\mathcal{R}f(t, \theta) = \int_{\mathbb{R}^d} f(x) \cdot \delta(t - \langle x, \theta \rangle) dx, \quad (3)$$

and  $f_\mu, f_\nu$  are the probability density functions of  $\mu, \nu$ , respectively. Max Sliced Wasserstein (MaxSW) distance is discussed in Appendix C.

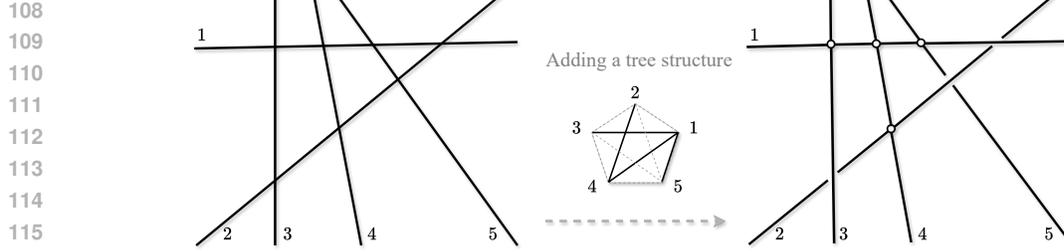


Figure 1: This illustration demonstrates the process of adding a tree structure to a system of lines. *Left*: An example of a system of 5 lines in  $\mathbb{R}^2$ , where the lines intersect, making the system connected. *Right*: Adding a tree structure to the connected system. In this example, only four pairs of lines are adjacent, shown by intersections, while the remaining pairs are disconnected, represented by gaps. This structure is derived by taking a spanning tree from a graph with five nodes (representing the five lines), with edges connecting nodes where lines intersect.

**Monte Carlo estimation for SW.** The Monte Carlo method is usually employed to approximate the intractable integral in Equation (2) as follows:

$$\widehat{\text{SW}}_p(\mu, \nu) = \left( \frac{1}{L} \sum_{l=1}^L \text{W}_p^p(\mathcal{R}f_\mu(\cdot, \theta_l), \mathcal{R}f_\nu(\cdot, \theta_l)) \right)^{\frac{1}{p}}, \quad (4)$$

where  $\theta_1, \dots, \theta_L$  are drawn independently from  $\mathcal{U}(\mathbb{S}^{d-1})$ . Using the closed-form expression of one-dimensional Wasserstein distance, when  $\mu$  and  $\nu$  are discrete measures that have supports of at most  $n$  supports, the computational complexity of  $\widehat{\text{SW}}_p$  is  $\mathcal{O}(Ln \log n + Ldn)$  (Peyré et al., 2019).

**Tree Wasserstein Distances.** Given a rooted tree  $(\mathcal{T}, r)$  ( $\mathcal{T}$  is a tree as a graph, with one certain node  $r$  called root) with non-negative edge lengths, and the ground metric  $d_{\mathcal{T}}$ , i.e. the length of the unique path between two nodes. For two distributions  $\mu, \nu$  supported on nodes of  $\mathcal{T}$ , the Wasserstein distance with ground cost  $d_{\mathcal{T}}$ , i.e., tree-Wasserstein (TW) (Le et al., 2019), admits a closed-form expression

$$\text{W}_{d_{\mathcal{T}}, 1}(\mu, \nu) = \sum_{e \in \mathcal{T}} w_e \cdot |\mu(\Gamma(v_e)) - \nu(\Gamma(v_e))|, \quad (5)$$

where  $v_e$  is the farther endpoint of edge  $e$  from  $r$ ,  $w_e$  is the length of  $e$ , and  $\Gamma(v_e)$  is the subtree of  $\mathcal{T}$  rooted at  $v_e$ , i.e. the subtree consists of all node  $x$  that the unique path from  $x$  to  $r$  contains  $v_e$ .

### 3 SYSTEM OF LINES WITH TREE STRUCTURES

This section provides an *intuitive and brief* introduction of systems of lines and their additional tree structures. These structures form metric spaces, called tree systems, which serve as a *generalization* of one-dimensional lines within the framework of the Sliced-Wasserstein distance. We then explore the topological properties and the construction of tree systems. The ideas are illustrated in Figures 1, 2, 3, and a *complete formal construction* with theoretical proofs are presented in Appendix A.

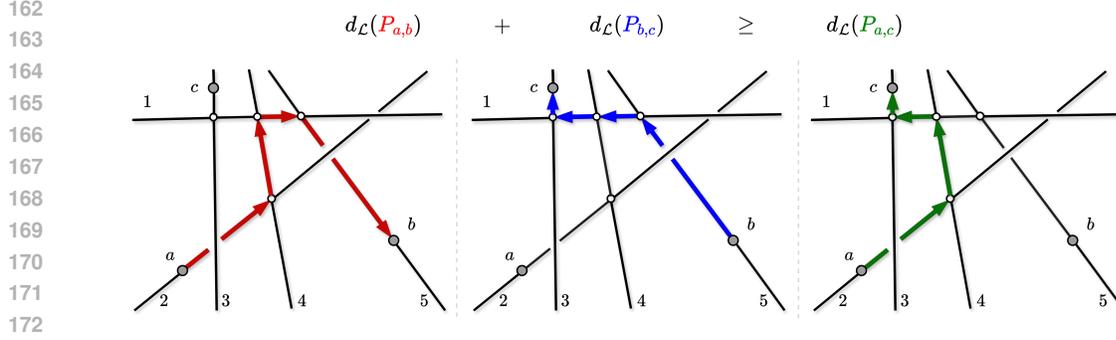
#### 3.1 SYSTEM OF LINES AND TREE SYSTEM

A line in  $\mathbb{R}^d$  can be fully described by specifying its direction and a point it passes through. Specifically, a line is determined by  $(x, \theta) \in \mathbb{R}^d \times \mathbb{S}^{d-1}$ , and is parameterized as  $x + t \cdot \theta$  for  $t \in \mathbb{R}$ .

**Definition 3.1** (Line and System of Lines in  $\mathbb{R}^d$ ). A *line* in  $\mathbb{R}^d$  is an element  $(x, \theta)$  of  $\mathbb{R}^d \times \mathbb{S}^{d-1}$ . For  $k \geq 1$ , a *system of  $k$  lines* in  $\mathbb{R}^d$  is a set of  $k$  lines in  $\mathbb{R}^d$ .

We denote a line in  $\mathbb{R}^d$  as  $l = (x_l, \theta_l)$ . Here,  $x_l$  and  $\theta_l$  are called *source* and *direction* of  $l$ , respectively. Denote  $(\mathbb{R}^d \times \mathbb{S}^{d-1})^k$  by  $\mathbb{L}_k^d$ , which is the *space of systems of  $k$  lines in  $\mathbb{R}^d$* , and an element of  $\mathbb{L}_k^d$  is usually denoted by  $\mathcal{L}$ . The *ground set* of a system of lines  $\mathcal{L}$  is defined by:

$$\bar{\mathcal{L}} := \{(x, l) \in \mathbb{R}^d \times \mathcal{L} : x = x_l + t_x \cdot \theta_l \text{ for some } t_x \in \mathbb{R}\}.$$



173 Figure 2: The same tree system  $\mathcal{L}$  shown in Figure 1, naturally has a topology derived from five  
174 copies of  $\mathbb{R}$ . Consider three points  $a, b, c$ . The red zigzag line presents the unique path from  $a$  to  $b$ .  
175 Here the distance between  $a, b$ , i.e.  $d_{\mathcal{L}}(a, b)$ , is the sum of four red line segments. Similar for paths  
176 between  $b$  and  $c$ ;  $a$  and  $c$ . This demonstrates that the triangle inequality is satisfied for  $d_{\mathcal{L}}$ .  
177

178 For each element  $\bar{\mathcal{L}}$ , we sometimes write  $(x, l)$  as  $(t_x, l)$ , where  $t_x \in \mathbb{R}$  presents the parameterization  
179 of  $x$  on  $l$  as  $x = x_l + t_x \cdot \theta_l$ . By a point of  $\mathcal{L}$ , we refer to a point of the ground set  $\bar{\mathcal{L}}$ . Now consider  
180 a system of distinct lines  $\mathcal{L}$  in  $\mathbb{R}^d$ .  $\mathcal{L}$  is said to be *connected* if its points form a connected set in  $\mathbb{R}^d$ .  
181 In this case,  $\mathcal{L}$  naturally has certain tree structures. Figure 1 gives an example of a system of lines  
182 with an added tree structure. A pair  $(\mathcal{L}, \mathcal{T})$  consists of a connected system of lines  $\mathcal{L}$  and its tree  
183 structure  $\mathcal{T}$  of  $\mathcal{L}$ , is called a *tree system*. We also denote it as  $\mathcal{L}$  for short.  
184

### 185 3.2 TOPOLOGICAL PROPERTIES OF TREE SYSTEMS

186 A tree system  $\mathcal{L}$  can be intuitively understood as a system of lines that are connected in certain  
187 ways. It naturally forms a topological space by *taking disjoint union copies of  $\mathbb{R}$*  and then *taking*  
188 *the quotient at intersections of these copies*. The disjoint union is straightforward, and the quotient  
189 follows the tree structure of  $\mathcal{L}$ . The topological space resulting from these actions is called the  
190 *topological space of a tree system  $\mathcal{L}$* , and is denoted by  $\Omega_{\mathcal{L}}$ . By its construction,  $\Omega_{\mathcal{L}}$  naturally  
191 carries a measure induced from the standard measure on each copy of  $\mathbb{R}$ . This measure is denoted  
192 by  $\mu_{\mathcal{L}}$ . Notice that, due to the tree structure, a *unique* path exists between any two points of  $\Omega_{\mathcal{L}}$ .  
193 This leads to an important result regarding the metrizable of  $\Omega_{\mathcal{L}}$ .

194 **Theorem 3.2** ( $\Omega_{\mathcal{L}}$  is metrizable by a tree metric). Consider  $d_{\mathcal{L}}: \Omega_{\mathcal{L}} \times \Omega_{\mathcal{L}} \rightarrow [0, \infty)$  defined by:

$$195 \quad d_{\mathcal{L}}(a, b) := \mu_{\mathcal{L}}(P_{a,b}), \quad \forall a, b \in \Omega_{\mathcal{L}}, \quad (6)$$

196 where  $P_{a,b}$  is the unique path between  $a$  and  $b$  in  $\Omega_{\mathcal{L}}$ . Then  $d_{\mathcal{L}}$  is a metric on  $\Omega_{\mathcal{L}}$ , which makes  
197  $(\Omega_{\mathcal{L}}, d_{\mathcal{L}})$  a metric space. Moreover,  $d_{\mathcal{L}}$  is a tree metric, and the topology on  $\Omega_{\mathcal{L}}$  induced by  $d_{\mathcal{L}}$  is  
198 identical to the topology of  $\Omega_{\mathcal{L}}$ .  
199

200 The proof is presented in Theorem A.11. Figure 2 illustrates an example of a unique path between  
201 two points on a tree system, providing an intuitive explanation of why  $d_{\mathcal{L}}$  is indeed a metric.  
202

### 203 3.3 CONSTRUCTION OF TREE SYSTEMS AND SAMPLING PROCESS

204 A tree system can be built inductively by sampling lines, ensuring that each new line intersects one  
205 of the previously sampled lines. We introduce a straightforward method to construct a tree system:  
206 start by sampling a line, and at each subsequent step, sample a new line that intersects the previously  
207 selected line. Specifically, the process is as follows:  
208

209 *Step 1.* Sampling  $x_1 \sim \mu_1$  for an  $\mu_1 \in \mathcal{P}(\mathbb{R}^d)$ , then  $\theta_1 \sim \nu_1$  for an  $\nu_1 \in \mathcal{P}(\mathbb{S}^{d-1})$ . The pair  
210  $(x_1, \theta_1)$  forms the first line;

211 *Step  $i$ .* At step  $i$ , sampling  $x_i = x_{i-1} + t_i \cdot \theta_{i-1}$  where  $t_i \sim \mu_i$  for an  $\mu_i \in \mathcal{P}(\mathbb{R})$ , then  $\theta_i \sim \nu_i$   
212 for an  $\nu_i \in \mathcal{P}(\mathbb{S}^{d-1})$ . The pair  $(x_i, \theta_i)$  forms the  $i^{\text{th}}$  line.  
213

214 The tree system produced by this construction has a *chain-like tree structure*, where the  $i^{\text{th}}$  line inter-  
215 sects the  $(i+1)^{\text{th}}$  line. A *general approach* for sampling tree systems is provided in Appendix A.4.  
In practice, we simply assume all the distributions  $\mu$ 's and  $\nu$ 's to be independent, and let:

1.  $\mu_1$  to be a distribution on a bounded subset of  $\mathbb{R}^d$ , for instance, the uniform distribution on the  $d$ -dimensional cube  $[-1, 1]^d$ , i.e.  $\mathcal{U}([-1, 1]^d)$ ;
2.  $\mu_i$  for  $i > 1$  to be a distribution on a bounded subset of  $\mathbb{R}$ , for instance, the uniform distribution on the interval  $[-1, 1]$ , i.e.  $\mathcal{U}([-1, 1])$ ;
3.  $\theta_n$  to be a distribution on  $\mathbb{S}^{d-1}$ , for instance, the uniform distribution  $\mathcal{U}(\mathbb{S}^{d-1})$ .

Using the distributions  $\mu$ 's and  $\nu$ 's, we get a distribution on the space of all tree systems that can be sampled by this way. We obtain a distribution over the space of all tree systems that can be sampled in this manner. The algorithm for sampling tree systems is summarized in Algorithm 1, and illustrated in Figure 3.

---

**Algorithm 1** Sampling (chain-like) tree systems.
 

---

**Input:** The number of lines in tree systems  $k$ .  
 Sampling  $x_1 \sim \mathcal{U}([-1, 1]^d)$  and  $\theta_1 \sim \mathcal{U}(\mathbb{S}^{d-1})$ .  
**for**  $i = 2$  to  $k$  **do**  
   Sample  $t_i \sim \mathcal{U}([-1, 1])$  and  $\theta_i \sim \mathcal{U}(\mathbb{S}^{d-1})$ .  
   Compute  $x_i = x_{i-1} + t_i \cdot \theta_{i-1}$ .  
**end for**  
**Return:**  $(x_1, \theta_1), (x_2, \theta_2), \dots, (x_k, \theta_k)$ .

---

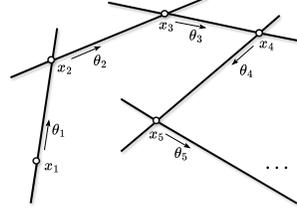


Figure 3: Illustration of Algorithm 1

#### 4 RADON TRANSFORM ON SYSTEMS OF LINES

In this section, we introduce the notions of the space of Lebesgue integrable functions and the Radon Transform for systems of lines. Let  $\mathcal{L} \in \mathbb{L}_k^d$  be a system of  $k$  lines. Denote  $L^1(\mathbb{R}^d)$  as the space of Lebesgue integrable functions on  $\mathbb{R}^d$  with norm  $\|\cdot\|_1$ , i.e.

$$L^1(\mathbb{R}^d) = \left\{ f: \mathbb{R}^d \rightarrow \mathbb{R} : \|f\|_1 = \int_{\mathbb{R}^d} |f(x)| dx < \infty \right\}. \quad (7)$$

Two functions  $f_1, f_2 \in L^1(\mathbb{R}^d)$  are considered to be identical if  $f_1(x) = f_2(x)$  almost everywhere on  $\mathbb{R}^d$ . As a counterpart, a *Lebesgue integrable function on  $\mathcal{L}$*  is a function  $f: \tilde{\mathcal{L}} \rightarrow \mathbb{R}$  such that:

$$\|f\|_{\mathcal{L}} := \sum_{l \in \mathcal{L}} \int_{\mathbb{R}} |f(t_x, l)| dt_x < \infty. \quad (8)$$

The *space of Lebesgue integrable functions on  $\mathcal{L}$*  is denoted by  $L^1(\mathcal{L})$ . Two functions  $f_1, f_2 \in L^1(\mathcal{L})$  are considered to be identical if  $f_1(x) = f_2(x)$  almost everywhere on  $\tilde{\mathcal{L}}$ . The space  $L^1(\mathcal{L})$  with norm  $\|\cdot\|_{\mathcal{L}}$  is a Banach space.

Recall that  $\mathcal{L}$  has  $k$  lines. Denote the  $(k-1)$ -dimensional standard simplex as  $\Delta_{k-1} = \{(a_l)_{l \in \mathcal{L}} : a_l \geq 0 \text{ and } \sum_{l \in \mathcal{L}} a_l = 1\} \subset \mathbb{R}^k$ . Denote  $\mathcal{C}(\mathbb{R}^d, \Delta_{k-1})$  as the space of continuous maps from  $\mathbb{R}^d$  to  $\Delta_{k-1}$ . A map in  $\mathcal{C}(\mathbb{R}^d, \Delta_{k-1})$  is referred to as a *splitting map*. Let  $\mathcal{L}$  be a system of  $k$  lines in  $\mathbb{L}_k^d$ ,  $\alpha$  be a splitting map in  $\mathcal{C}(\mathbb{R}^d, \Delta_{k-1})$ , we define an operator associated to  $\alpha$  that transforms a Lebesgue integrable functions on  $\mathbb{R}^d$  to a Lebesgue integrable functions on  $\mathcal{L}$ , analogous to the original Radon Transform. For  $f \in L^1(\mathbb{R}^d)$ , define:

$$\begin{aligned} \mathcal{R}_{\mathcal{L}}^{\alpha} f : \tilde{\mathcal{L}} &\longrightarrow \mathbb{R} \\ (x, l) &\longmapsto \int_{\mathbb{R}^d} f(y) \cdot \alpha(y)_l \cdot \delta(t_x - \langle y - x_l, \theta_l \rangle) dy, \end{aligned} \quad (9)$$

where  $\delta$  is the 1-dimensional Dirac delta function. For  $f \in L^1(\mathbb{R}^d)$ , we can show that  $\mathcal{R}_{\mathcal{L}}^{\alpha} f \in L^1(\mathcal{L})$ . Moreover, we have  $\|\mathcal{R}_{\mathcal{L}}^{\alpha} f\|_{\mathcal{L}} \leq \|f\|_1$ . In other words, the operator  $\mathcal{R}_{\mathcal{L}}^{\alpha}: L^1(\mathbb{R}^d) \rightarrow L^1(\mathcal{L})$  is well-defined, and is a linear operator. The proof for these properties is presented in Theorem B.2. We now propose a novel variant of Radon Transform for systems of lines.

**Definition 4.1** (Radon Transform on Systems of lines). For  $\alpha \in \mathcal{C}(\mathbb{R}^d, \Delta_{k-1})$ , the operator  $\mathcal{R}^{\alpha}$ :

$$\begin{aligned} \mathcal{R}^{\alpha} : L^1(\mathbb{R}^d) &\longrightarrow \prod_{\mathcal{L} \in \mathbb{L}_k^d} L^1(\mathcal{L}) \\ f &\longmapsto (\mathcal{R}_{\mathcal{L}}^{\alpha} f)_{\mathcal{L} \in \mathbb{L}_k^d}. \end{aligned}$$

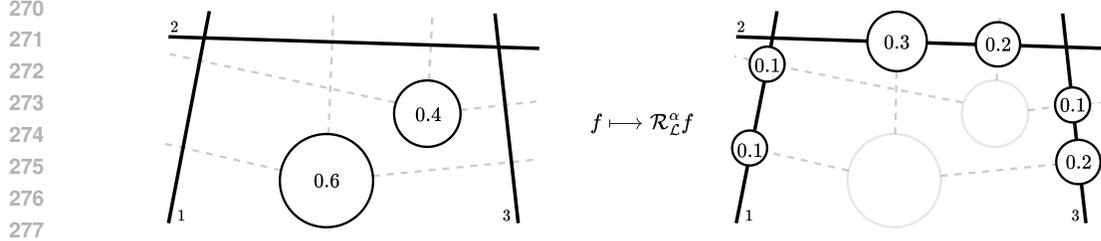


Figure 4: An illustration of Radon Transform on Systems of Lines. Given  $f \in L^1(\mathbb{R}^d)$  such that  $f(x) = 0.6$ ,  $f(y) = 0.4$ , and  $\mathcal{L}$  is a system of 3 lines. For a splitting map  $\alpha$  such that  $\alpha(x) = (1/6, 3/6, 2/6)$  and  $\alpha(y) = (1/4, 2/4, 1/4)$ ,  $f$  is transformed to  $\mathcal{R}_{\mathcal{L}}^{\alpha} f$ . By Equation (9), for instance, the value of  $\mathcal{R}_{\mathcal{L}}^{\alpha} f$  at the projection of  $x$  onto line (2) of  $\mathcal{L}$  is  $f(x) \cdot \alpha(x)_2 = 0.3$ .

is called the *Radon Transform on Systems of Lines*.

*Remark.* An illustration of splitting maps and the Radon Transform on Systems of Lines is presented in Figure 4. Intuitively, splitting map  $\alpha$  indicates *how the mass at a given point is distributed across all lines of a system of lines*. In the case  $k = 1$ , there is only one splitting map which is the constant function 1, and the Radon Transform for  $\mathbb{L}_1^d$  is identical to the traditional Radon Transform.

Many variants of the Radon transform require the transform to be injective. In the case of systems of lines, the injectivity also holds for  $\mathcal{R}^{\alpha}$ .

**Theorem 4.2.**  $\mathcal{R}^{\alpha}$  is injective for all splitting maps  $\alpha \in \mathcal{C}(\mathbb{R}^d, \Delta_{k-1})$ .

The proof of this theorem is presented in Theorem B.1. Denote  $\mathcal{P}(\mathbb{R}^d)$  as the space of all probability distribution on  $\mathbb{R}^d$ , and define a *probability distribution on  $\mathcal{L}$*  to be a function  $f \in L^1(\mathcal{L})$  such that  $f: \mathcal{L} \rightarrow [0, \infty)$  and  $\|f\|_{\mathcal{L}} = 1$ . The *space of probability distribution on  $\mathcal{L}$*  is denoted by  $\mathcal{P}(\mathcal{L})$ . Then  $\mathcal{R}_{\mathcal{L}}^{\alpha}$  transforms a distribution in  $\mathcal{P}(\mathbb{R}^d)$  to a distribution in  $\mathcal{P}(\mathcal{L})$ . In other words, the restricted operator  $\mathcal{R}_{\mathcal{L}}^{\alpha}: \mathcal{P}(\mathbb{R}^d) \rightarrow \mathcal{P}(\mathcal{L})$  is also well-defined.

## 5 TREE-SLICED WASSERSTEIN DISTANCE ON SYSTEMS OF LINES

In this section, we present a novel Tree-Sliced Wasserstein distance on Systems of Lines (TSW-SL). Consider  $\mathbb{T}$  as *the space of tree systems* consisting of  $k$  lines in  $\mathbb{R}^d$  that be sampled by Algorithm 1. By the remark at the end of Subsection 3.3, we have a *distribution  $\sigma$  on the space  $\mathbb{T}$* . General cases of  $\mathbb{T}$ , as in Appendix A.4, will be handled in a similar manner. For simplicity and convenience, we occasionally use the same notation to represent both a measure and its probability distribution function, provided the context makes the meaning clear.

### 5.1 TREE-SLICED WASSERSTEIN DISTANCE ON SYSTEMS OF LINES

Consider a splitting function  $\alpha$  in  $\mathcal{C}(\mathbb{R}^d, \Delta_{k-1})$ . Given two probability distributions  $\mu, \nu$  in  $\mathcal{P}(\mathbb{R}^d)$  and a tree system  $\mathcal{L} \in \mathbb{T}$ . By the Radon Transform  $\mathcal{R}_{\mathcal{L}}^{\alpha}$  in Definition 4.1,  $\mu$  and  $\nu$  are transformed to two probability distributions  $\mathcal{R}_{\mathcal{L}}^{\alpha} \mu$  and  $\mathcal{R}_{\mathcal{L}}^{\alpha} \nu$  in  $\mathcal{P}(\mathcal{L})$ . By Theorem 3.2,  $\mathcal{L}$  has a tree metric  $d_{\mathcal{L}}$ , we compute Wasserstein distance  $W_{d_{\mathcal{L}}, 1}(\mathcal{R}_{\mathcal{L}}^{\alpha} \mu, \mathcal{R}_{\mathcal{L}}^{\alpha} \nu)$  between  $\mathcal{R}_{\mathcal{L}}^{\alpha} \mu$  and  $\mathcal{R}_{\mathcal{L}}^{\alpha} \nu$  by Equation (5).

**Definition 5.1** (Tree-Sliced Wasserstein Distance on Systems of Lines). The *Tree-Sliced Wasserstein distance on Systems of Lines* between  $\mu, \nu$  in  $\mathcal{P}(\mathbb{R}^d)$  is defined by:

$$\text{TSW-SL}(\mu, \nu) := \int_{\mathbb{T}} W_{d_{\mathcal{L}}, 1}(\mathcal{R}_{\mathcal{L}}^{\alpha} \mu, \mathcal{R}_{\mathcal{L}}^{\alpha} \nu) d\sigma(\mathcal{L}). \quad (10)$$

*Remark.* Note that, the definition of TSW-SL depends on the space of sampled tree systems  $\mathbb{T}$ , the distribution  $\sigma$  on  $\mathbb{T}$ , and the splitting function  $\alpha$ . For simplifying the notation, we omit them.

TSW-SL is a metric on  $\mathcal{P}(\mathbb{R}^d)$ . The proof for the below theorem is provided in Appendix D.1.

**Theorem 5.2.** TSW-SL is a metric on  $\mathcal{P}(\mathbb{R}^d)$ .

*Remark.* If tree systems in  $\mathbb{T}$  consists consist only one line, i.e.  $k = 1$ , then in Definition 4.1, the splitting map  $\alpha$  is the constant map 1, and the Radon Transform  $\mathcal{R}^{\alpha}$  now becomes identical to the

original Radon Transform, as pushing forward measures onto lines depends only on their directions. Also, according to the sampling process described in Subsection 3.3,  $\sigma$  becomes the distribution of  $\theta_1$ , which is  $\mathcal{U}(\mathbb{S}^{d-1})$ . In this case, TSW-SL in Equation (10) is identical to SW in Equation (2). Furthermore, in Appendix C, we introduce *Max Tree-Sliced Wasserstein Distance on Systems of Lines* (MaxTSW-SL), an analog of MaxSW (Deshpande et al., 2019).

## 5.2 COMPUTING TSW-SL

We employ the Monte Carlo method to estimate the intractable integral in Equation (10) as follows:

$$\widehat{\text{TSW-SL}}(\mu, \nu) = \frac{1}{L} \sum_{i=1}^L W_{d_{\mathcal{L}_i, 1}}(\mathcal{R}_{\mathcal{L}_i}^\alpha \mu, \mathcal{R}_{\mathcal{L}_i}^\alpha \nu), \quad (11)$$

where  $\mathcal{L}_1, \dots, \mathcal{L}_L \stackrel{i.i.d.}{\sim} \sigma$  are referred to as projecting tree systems. We now discuss on how to compute  $W_{d_{\mathcal{L}, 1}}(\mathcal{R}_{\mathcal{L}}^\alpha \mu, \mathcal{R}_{\mathcal{L}}^\alpha \nu)$  for  $\mathcal{L} \in \mathbb{T}$ . In applications, consider  $\mu, \nu \in \mathcal{P}(\mathbb{R}^d)$  given as follows:

$$\mu(x) = \sum_{i=1}^n u_i \cdot \delta(x - a_i) \quad \text{and} \quad \nu(x) = \sum_{i=1}^m v_i \cdot \delta(x - b_i) \quad (12)$$

$\mathcal{R}_{\mathcal{L}}^\alpha$  projects  $\mu, \nu$  on  $\mathcal{L}$ , resulting in discrete measures  $\mathcal{R}_{\mathcal{L}}^\alpha \mu, \mathcal{R}_{\mathcal{L}}^\alpha \nu$  in  $\mathcal{P}(\mathcal{L})$ . In details, from definition of  $\mathcal{R}_{\mathcal{L}}^\alpha \mu$ , the support of  $\mathcal{R}_{\mathcal{L}}^\alpha \mu$  is the set of all projections of support of  $\mu$  onto lines of  $\mathcal{L}$ . Moreover, the value of  $\mathcal{R}_{\mathcal{L}}^\alpha \mu$  at projections of  $a_i$  onto  $l$  is equal to  $\alpha(a_i)_l \cdot u_i$ . Similar for  $\mathcal{R}_{\mathcal{L}}^\alpha \nu$ . From this detailed description of  $\mathcal{R}_{\mathcal{L}}^\alpha \mu, \mathcal{R}_{\mathcal{L}}^\alpha \nu$ , together with Equation (5), we derive a *closed-form expression* of  $W_{d_{\mathcal{L}, 1}}(\mathcal{R}_{\mathcal{L}}^\alpha \mu, \mathcal{R}_{\mathcal{L}}^\alpha \nu)$  as follows:

$$W_{d_{\mathcal{L}, 1}}(\mathcal{R}_{\mathcal{L}}^\alpha \mu, \mathcal{R}_{\mathcal{L}}^\alpha \nu) = \sum_{e \in \mathcal{T}} w_e \cdot \left| \mathcal{R}_{\mathcal{L}}^\alpha \mu(\Gamma(v_e)) - \mathcal{R}_{\mathcal{L}}^\alpha \nu(\Gamma(v_e)) \right|. \quad (13)$$

This expression enables an efficient and highly parallelizable implementation of TSW-SL, as it relies on fundamental operations like matrix multiplication and sorting.

*Remark.* Assume  $n \geq m$ , the time complexity for TSW-SL is  $\mathcal{O}(Lkn \log n + Lkdn)$  since it primarily involves projecting onto  $L \times k$  lines and sorting  $n$  projections on each line. This complexity is equivalent to that of SW when the number of projection directions is the same. Therefore, in our experiments, we ensure a fair comparison by evaluating the performance of TSW-SL against SW or its variants using *the same number of projection directions*.

We summarize this section with Algorithm 2 of computing TSW-SL.

---

### Algorithm 2 Tree Sliced Wasserstein distance on Systems of Lines.

---

**Input:**  $\mu$  and  $\nu$  in  $\mathcal{P}(\mathbb{R}^d)$ , the number of lines in each tree system  $k$ , the number of tree systems  $L$ , a splitting map  $\alpha: \mathbb{R}^d \rightarrow \Delta_{k-1}$ .  
**for**  $l = 1$  to  $L$  **do**  
  Sample tree system  $\mathcal{L}_l = ((x_1^{(l)}, \theta_1^{(l)}), \dots, (x_k^{(l)}, \theta_k^{(l)}))$ .  
  Project  $\mu$  and  $\nu$  onto  $\mathcal{L}_l$  to get  $\mathcal{R}_{\mathcal{L}_l}^\alpha \mu$  and  $\mathcal{R}_{\mathcal{L}_l}^\alpha \nu$ .  
  Compute  $W_{d_{\mathcal{L}_l, 1}}(\mathcal{R}_{\mathcal{L}_l}^\alpha \mu, \mathcal{R}_{\mathcal{L}_l}^\alpha \nu)$ .  
**end for**  
Compute  $\widehat{\text{TSW-SL}} = (1/L) \cdot \sum_{l=1}^L W_{d_{\mathcal{L}_l, 1}}(\mathcal{R}_{\mathcal{L}_l}^\alpha \mu, \mathcal{R}_{\mathcal{L}_l}^\alpha \nu)$ .  
**Return:**  $\widehat{\text{TSW-SL}}(\mu, \nu)$ .

---

## 6 EXPERIMENTAL RESULTS

In this section, we present empirical results demonstrating the advantages of our TSW-SL distance over traditional SW distance and its variants, and how MaxTSW-SL enhances the original MaxSW (Deshpande et al., 2019) through optimized tree construction. The splitting maps  $\alpha$  will be selected either as a trainable constant vector or a random vector, while the tree systems will be sampled such that the root is positioned near the mean of the target distribution, i.e. the data mean. It

Table 1: Average Wasserstein distance between source and target distributions of 10 runs on Swiss Roll and 25 Gaussians datasets. All methods use 100 projecting directions.

Methods	Swiss Roll						25 Gaussians					
	Iteration					Time/Iter(s)	Iteration					Time/Iter(s)
	500	1000	1500	2000	2500		500	1000	1500	2000	2500	
SW	5.73e-3	2.04e-3	1.23e-3	1.11e-3	1.05e-3	0.009	1.61e-1	9.52e-2	3.44e-2	2.56e-2	2.20e-2	0.006
MaxSW	2.47e-2	1.03e-2	6.10e-3	4.47e-3	3.45e-3	2.46	5.09e-1	2.36e-1	1.33e-1	9.70e-2	8.48e-2	2.38
SWG	3.84e-2	1.53e-2	1.02e-2	4.49e-3	3.57e-5	0.011	3.10e-1	1.17e-1	3.38e-2	3.58e-3	2.54e-4	0.009
LCVSW	7.28e-3	1.40e-3	1.38e-3	1.38e-3	1.36e-3	0.010	3.38e-1	6.64e-2	3.06e-2	3.06e-2	3.02e-2	0.009
TSW-SL	9.41e-3	<b>2.03e-7</b>	<b>9.63e-8</b>	<b>4.44e-8</b>	<b>3.65e-8</b>	0.014	3.49e-1	9.06e-2	2.96e-2	1.20e-2	<b>3.03e-7</b>	0.010
MaxTSW-SL	<b>2.75e-6</b>	<b>8.24e-7</b>	<b>5.14e-7</b>	<b>5.02e-7</b>	<b>5.00e-7</b>	2.53	<b>1.12e-1</b>	<b>8.28e-3</b>	<b>1.61e-6</b>	<b>7.32e-7</b>	<b>5.19e-7</b>	2.49

Table 2: Average Wasserstein distance between source and target distributions of 10 runs on high-dimensional datasets.

Dimension	Iteration 500		Iteration 1000		Iteration 1500		Iteration 2000		Iteration 2500		Time/Iter(s)	
	SW	TSW-SL	SW	TSW-SL	SW	TSW-SL	SW	TSW-SL	SW	TSW-SL	SW	TSW-SL
10	4.32e-3	<b>2.81e-3</b>	2.94e-3	<b>2.00e-3</b>	2.81e-3	<b>1.55e-3</b>	2.23e-3	<b>1.59e-3</b>	2.28e-3	<b>1.75e-3</b>	<b>0.010</b>	<b>0.015</b>
50	50.41	<b>39.26</b>	45.69	<b>21.91</b>	42.56	<b>11.91</b>	38.81	<b>4.08</b>	35.75	<b>1.72</b>	<b>0.014</b>	<b>0.018</b>
75	92.39	<b>79.71</b>	90.79	<b>67.99</b>	90.07	<b>53.92</b>	86.58	<b>44.91</b>	90.31	<b>31.61</b>	<b>0.015</b>	<b>0.018</b>
100	130.12	<b>117.66</b>	128.13	<b>103.23</b>	128.58	<b>93.41</b>	129.80	<b>80.46</b>	128.29	<b>75.28</b>	<b>0.018</b>	<b>0.019</b>
150	214.09	<b>203.30</b>	213.71	<b>190.62</b>	215.05	<b>186.77</b>	212.90	<b>183.52</b>	216.32	<b>182.63</b>	<b>0.020</b>	<b>0.022</b>
200	302.84	<b>289.83</b>	301.35	<b>283.34</b>	303.07	<b>276.94</b>	302.70	<b>279.24</b>	301.51	<b>279.08</b>	<b>0.020</b>	<b>0.021</b>

is worth noting that the paper presents a simple alternative by substituting lines in SW with tree systems, focusing mainly on comparing TSW-SL with the original SW, without expecting TSW-SL to outperform more recent SW variant. Further improvements to TSW-SL could be made by incorporating advanced techniques developed for SW, but we leave this for future research, choosing instead to focus on the fundamental aspects of TSW-SL.

## 6.1 GRADIENT FLOWS

First of all, we conduct experiments to compare the effectiveness of our methods with baselines in the gradient flow task. In this task, we aim to minimize  $\text{TSW-SL}(\mu, \nu)$ , where  $\nu$  is the target distribution and  $\mu$  represents the source distribution. The optimization process is carried out iteratively as  $\partial_t \mu_t = -\nabla \text{TSW-SL}(\mu_t, \nu)$  with  $\mu_0 = \mathcal{N}(0, 1)$ ,  $-\partial_t \mu_t$  represents the change in the source distribution over time and  $\nabla \text{TSW-SL}(\mu_t, \nu)$  is the gradient of TSW-SL with respect to  $\mu_t$ . We initialize with  $\mu_0 = \mathcal{N}(0, 1)$  and iteratively update  $\mu_t$  over 2500 iterations.

To compare the effectiveness of various distance metrics, we employ alternative distances as loss functions (SW (Bonnel et al., 2015), MaxSW (Deshpande et al., 2019), SWGG (Mahey et al., 2023) and LCVSW (Nguyen & Ho, 2023)) instead of TSW-SL. Over 2500 timesteps, we evaluate the Wasserstein distance between source and target distributions at iteration 500, 1000, 1500, 2000 and 2500. We use  $L = 100$  in SW variants and  $L = 25, k = 4$  in TSW-SL for a fair comparison. Detailed training settings are presented in Appendix E.1.

We first utilize both the Swiss Roll (a non-linear dataset) and 25 Gaussians (a multimodal dataset) as described in (Kolouri et al., 2019). In Table 1, we present the performance and runtime of various methods on these datasets, emphasizing the reduction of the Wasserstein distance over iterations. Notably, across both datasets, our TSW-SL method demonstrates superior performance by significantly reducing the Wasserstein distance. Moreover, our MaxTSW-SL method shows a significant decrease in the Wasserstein distance compared to MaxSW, highlighting its improved performance and effectiveness. Furthermore, we provide additional results from experiments of 10, 50, 75, 100, 150, and 200-dimensional Gaussian distributions, where target distribution supports were sampled from these high-dimensional spaces to showcase the empirical advantages of our TSW-SL in capturing topological properties. In this context, we compare the Tree Sliced Wasserstein distance on a System of Lines (TSW-SL) with Sliced Wasserstein distance (SW) to demonstrate TSW-SL’s effectiveness when distribution supports lie in high-dimensional spaces. The results presented in Table 2 highlight TSW-SL’s superior ability to preserve the original data’s topological properties compared to SW.

Table 3: Average FID and IS score of 3 runs on CelebA and STL-10 of SN-GAN.

	CelebA (64x64)			STL-10 (96x96)			CelebA (64x64)			STL-10 (96x96)		
	FID(↓)	FID(↓)	IS(↑)	FID(↓)	FID(↓)	IS(↑)	FID(↓)	FID(↓)	IS(↑)	FID(↓)	FID(↓)	IS(↑)
SW ( $L = 50$ )	9.97 ± 1.02	69.46 ± 0.21	9.08 ± 0.06	SW ( $L = 500$ )	9.62 ± 0.42	53.52 ± 0.61	10.56 ± 0.05					
TSW-SL ( $L = 10, k = 5$ )	9.63 ± 0.46	<b>61.15 ± 0.37</b>	<b>10.00 ± 0.03</b>	TSW-SL ( $L = 100, k = 5$ )	<b>8.90 ± 0.49</b>	<b>51.81 ± 1.02</b>	<b>10.74 ± 0.13</b>					
TSW-SL ( $L = 17, k = 3$ )	<b>8.98 ± 0.75</b>	65.91 ± 0.64	9.75 ± 0.10	TSW-SL ( $L = 167, k = 3$ )	<b>8.90 ± 0.38</b>	52.27 ± 0.96	10.62 ± 0.18					

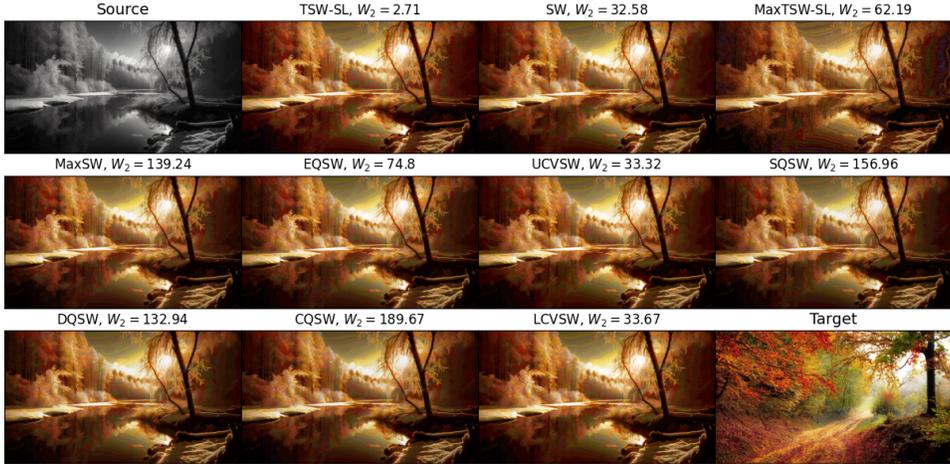


Figure 5: Style-transferred images from different models with 100 projecting directions.

## 6.2 COLOR TRANSFER

We continue by examining the performance of TSW-SL methods for transferring color between images to produce results that closely match the color distributions of the target images. Given a source image and a target image, we represent their respective color palettes as matrices  $X$  and  $Y$ , each with dimensions  $n \times 3$  (where  $n$  denotes the number of pixels). We traverse along the curve connecting  $P_X$  and  $P_Y$ , where  $P_X$  and  $P_Y$  denote the empirical distribution of the source and the target images respectively. More specifically, this curve (denoted as  $Z(t)$ ) starts from  $Z(0) = X$  and ends at  $Y$ . During optimization, we minimize the loss  $\mathcal{L}(Z(t), Y) = \text{Loss}(Z(t), Y)$  to make the color distribution of the obtained image close to that of the target image  $Y$ .

We evaluate the color-transferred images obtained by various loss  $\mathcal{L}$ , including SW (Bonneel et al., 2015), MaxSW (Deshpande et al., 2019), and SW variants proposed in (Nguyen et al., 2024a) to compare with our TSW-SL and MaxTSW-SL approaches. For consistency, we set  $L = 100$  for the SW variants and  $L = 25, k = 4$  for TSW-SL in our comparisons. We report the Wasserstein distances at the final time step along with the corresponding transferred images from various baselines Figure 5. TSW-SL produces images that most closely resemble the target, demonstrating a significant reduction compared to SW and its variants mentioned above with the same number of lines. In addition, MaxTSW-SL improves upon the original MaxSW, as highlighted by both qualitative and quantitative results.

## 6.3 GENERATIVE ADVERSARIAL NETWORK

We then explore the capabilities of our proposed TSW-SL framework within the context of generative adversarial networks (GANs). We employ the SNGAN architecture (Miyato et al., 2018). In detail, our approach is based on the methodology of the Sliced Wasserstein generator (Deshpande et al., 2018), with details provided in the Appendix E.3. Specifically, we conduct deep generative modeling experiments on the non-cropped CelebA dataset (Krizhevsky, 2009) with image size  $64 \times 64$ , and on the STL-10 dataset (Wang & Tan, 2016) with image size  $96 \times 96$ .

To demonstrate the empirical advantage of our method in enhancing generative adversarial networks, we employ two primary metrics: the Fréchet Inception Distance (FID) score (Heusel et al., 2017) and the Inception Score (IS) (Salimans et al., 2016). We omit to report the IS for the CelebA dataset

Table 4: Results for unconditional generation on CIFAR-10 of denoising diffusion models

Model	FID ↓	Time/Epoch(s) ↓	Time/Iter(s) ↓
DDGAN (Xiao et al. (2021))	3.64	136	0.45
SW-DD (Nguyen et al. (2024b))	2.90	140	0.47
DSW-DD (Nguyen et al. (2024b))	2.88	1059	3.53
EBSW-DD (Nguyen et al. (2024b))	2.87	145	0.48
RPSW-DD (Nguyen et al. (2024b))	2.82	159	0.53
IWRPSW-DD (Nguyen et al. (2024b))	2.70	152	0.51
TSW-SL-DD (Ours)	2.83	163	0.54

as it does not effectively capture the perceptual quality of face images (Heusel et al., 2017). Table 3 presents the results of SW and TSW-SL methodologies on the CelebA and STL-10 datasets, utilizing FID and IS as our metrics. We conduct experiments with two configurations of projecting directions: for 50 projecting directions, we use  $L = 50$  in SW compared to  $L = 10, k = 5$  and  $L = 17, k = 3$  in TSW-SL; for 500 projecting directions, we use  $L = 500$  in SW compared to  $L = 100, k = 5$  and  $L = 167, k = 3$  in TSW-SL. Our results reveal that TSW-SL significantly outperforms SW, demonstrating a considerable performance gap on both datasets in terms of IS and FID. We provide additional qualitative results in Appendix E.3.

#### 6.4 DENOISING DIFFUSION MODELS

Finally, we concentrate on denoising diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020), which are among the most complex generative frameworks for image generation. Diffusion models consist of a forward process that gradually adds Gaussian noise to data and a reverse process that learns to denoise the data. The forward process is defined as a Markov chain of  $T$  steps, where each step adds noise according to a predefined schedule. The reverse process, parameterized by  $\theta$ , aims to learn the denoising distribution. Traditionally, these models are trained using maximum likelihood by optimizing the evidence lower bound (ELBO). However, to accelerate generation, denoising diffusion GANs (Xiao et al., 2021) introduce an implicit denoising model and employ adversarial training. In our work, we build upon the framework in (Nguyen et al., 2024b) and replace the Augmented Generalized Mini-batch Energy distance with our novel TSW-SL distance as the kernel and conducting experiments on the CIFAR-10 dataset (Krizhevsky, 2009). For a detailed description of the model architecture and training loss, we refer readers to Appendix E.4.

Table 4 demonstrates that our TSW-SL loss function significantly enhances FID performance compared to conventional SW. It also outperforms RPSW and IWRPSW while yielding competitive results that are only marginally behind other state-of-the-art baselines in (Nguyen et al., 2024b). It is worth noting that our TSW-SL-DD maintains a competitive training time. This improvement underscores the efficacy of our approach in generating high-quality samples with improved fidelity.

## 7 CONCLUSION

This paper proposes a novel method called Tree-Sliced Wasserstein on Systems of Lines (TSW-SL), replacing the traditional one-dimensional lines in the Sliced Wasserstein (SW) framework with tree systems, providing a more geometrically meaningful space. This key innovation enables the proposed TSW-SL to capture more detailed structural information and geometric relationships within the data compared to SW while preserving computational efficiency. We rigorously develop the theoretical basis for our approach, verifying the essential properties of the Radon Transform and empirically demonstrating the benefits of TSW-SL across a range of application tasks. As this paper introduces a straightforward alternative by replacing one-dimensional lines in SW with tree systems, our primary comparison is between TSW-SL and the original SW, without anticipating that TSW-SL will surpass more recent SW variants. Future research on adapting recent advance techniques within the SW framework to TSW-SL remains an open area and is anticipated to lead to improved performance for Sliced Optimal Transport overall.

540 **Ethics Statement.** Given the nature of the work, we do not foresee any negative societal and ethical  
541 impacts of our work.

542 **Reproducibility Statement.** Source codes for our experiments are provided in the supplementary  
543 materials of the paper. The details of our experimental settings and computational infrastructure are  
544 given in Section 6 and the Appendix. All datasets that we used in the paper are published, and they  
545 are easy to access in the Internet.  
546

## 547 REFERENCES

548 David Alvarez-Melis, Tommi Jaakkola, and Stefanie Jegelka. Structured optimal transport. In  
549 *International Conference on Artificial Intelligence and Statistics*, pp. 1771–1780. PMLR, 2018.

550 Yikun Bai, Bernhard Schmitzer, Matthew Thorpe, and Soheil Kolouri. Sliced optimal partial trans-  
551 port. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,  
552 pp. 13681–13690, 2023.

553 Clément Bonet, Nicolas Courty, François Septier, and Lucas Drumetz. Efficient gradient flows in  
554 sliced-wasserstein space. *arXiv preprint arXiv:2110.10972*, 2021.

555 Clément Bonet, Laetitia Chapel, Lucas Drumetz, and Nicolas Courty. Hyperbolic sliced-wasserstein  
556 via geodesic and horospherical projections. In *Topological, Algebraic and Geometric Learning*  
557 *Workshops 2023*, pp. 334–370. PMLR, 2023.

558 Nicolas Bonneel, Julien Rabin, Gabriel Peyré, and Hanspeter Pfister. Sliced and radon wasserstein  
559 barycenters of measures. *Journal of Mathematical Imaging and Vision*, 51:22–45, 2015.

560 Charlotte Bunne, Laetitia Papaxanthos, Andreas Krause, and Marco Cuturi. Proximal optimal trans-  
561 port modeling of population dynamics. In *International Conference on Artificial Intelligence and*  
562 *Statistics*, pp. 6511–6528. PMLR, 2022.

563 Nicolas Courty, Rémi Flamary, Amaury Habrard, and Alain Rakotomamonjy. Joint distribution  
564 optimal transportation for domain adaptation. *Advances in neural information processing systems*,  
565 30, 2017.

566 Ishan Deshpande, Ziyu Zhang, and Alexander G Schwing. Generative modeling using the sliced  
567 wasserstein distance. In *Proceedings of the IEEE conference on computer vision and pattern*  
568 *recognition*, pp. 3483–3491, 2018.

569 Ishan Deshpande, Yuan-Ting Hu, Ruoyu Sun, Ayis Pyrros, Nasir Siddiqui, Sanmi Koyejo, Zhizhen  
570 Zhao, David Forsyth, and Alexander G Schwing. Max-sliced wasserstein distance and its use for  
571 gans. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,  
572 pp. 10648–10656, 2019.

573 Jiaojiao Fan, Isabel Haasler, Johan Karlsson, and Yongxin Chen. On the complexity of the optimal  
574 transport problem with graph-structured cost. In *International Conference on Artificial Intelli-*  
575 *gence and Statistics*, pp. 9147–9165. PMLR, 2022.

576 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,  
577 Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the*  
578 *ACM*, 63(11):139–144, 2020.

579 Allen Hatcher. *Algebraic topology*. 2005.

580 Sigurdur Helgason and Sigurdur Helgason. The radon transform on  $r$  n. *Integral Geometry and*  
581 *Radon Transforms*, pp. 1–62, 2011.

582 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter.  
583 Gans trained by a two time-scale update rule converge to a local nash equilibrium. 12 2017.

584 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*  
585 *neural information processing systems*, 33:6840–6851, 2020.

- 594 Nhat Ho, XuanLong Nguyen, Mikhail Yurochkin, Hung Hai Bui, Viet Huynh, and Dinh Phung.  
595 Multilevel clustering via wasserstein means. In *International conference on machine learning*,  
596 pp. 1501–1509. PMLR, 2017.
- 597 Xinru Hua, Truyen Nguyen, Tam Le, Jose Blanchet, and Viet Anh Nguyen. Dynamic flows on  
598 curved space generated by labeled data. In *Proceedings of the Thirty-Second International Joint*  
599 *Conference on Artificial Intelligence, IJCAI-23*, pp. 3803–3811, 2023.
- 601 Minhui Huang, Shiqian Ma, and Lifeng Lai. A riemannian block coordinate descent method for  
602 computing the projection robust wasserstein distance. In *International Conference on Machine*  
603 *Learning*, pp. 4446–4455. PMLR, 2021.
- 604 Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*,  
605 2014.
- 607 Soheil Kolouri, Gustavo K Rohde, and Heiko Hoffmann. Sliced wasserstein distance for learning  
608 gaussian mixture models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern*  
609 *Recognition*, pp. 3427–3436, 2018.
- 610 Soheil Kolouri, Kimia Nadjahi, Umut Simsekli, Roland Badeau, and Gustavo Rohde. Generalized  
611 sliced wasserstein distances. *Advances in neural information processing systems*, 32, 2019.
- 613 Alex Krizhevsky. Learning multiple layers of features from tiny images. 2009. URL <https://api.semanticscholar.org/CorpusID:18268744>.
- 614 Hugo Lavenant, Sebastian Claiici, Edward Chien, and Justin Solomon. Dynamical optimal transport  
615 on discrete surfaces. In *SIGGRAPH Asia 2018 Technical Papers*, pp. 250. ACM, 2018.
- 616 Tam Le, Makoto Yamada, Kenji Fukumizu, and Marco Cuturi. Tree-sliced variants of wasserstein  
617 distances. *Advances in neural information processing systems*, 32, 2019.
- 618 Tianyi Lin, Zeyu Zheng, Elynn Chen, Marco Cuturi, and Michael I Jordan. On projection robust  
619 optimal transport: Sample complexity and model misspecification. In *International Conference*  
620 *on Artificial Intelligence and Statistics*, pp. 262–270. PMLR, 2021.
- 621 Lang Liu, Soumik Pal, and Zaid Harchaoui. Entropy regularized optimal transport independence  
622 criterion. In *International Conference on Artificial Intelligence and Statistics*, pp. 11247–11279.  
623 PMLR, 2022.
- 624 Antoine Liutkus, Umut Simsekli, Szymon Majewski, Alain Durmus, and Fabian-Robert Stöter.  
625 Sliced-wasserstein flows: Nonparametric generative modeling via optimal transport and diffu-  
626 sions. In *International Conference on Machine Learning*, pp. 4104–4113. PMLR, 2019.
- 627 Manh Luong, Khai Nguyen, Nhat Ho, Reza Haf, Dinh Phung, and Lizhen Qu. Revisiting deep  
628 audio-text retrieval through the lens of transportation. *arXiv preprint arXiv:2405.10084*, 2024.
- 629 Guillaume Mahey, Laetitia Chapel, Gilles Gasso, Clément Bonet, and Nicolas Courty. Fast optimal  
630 transport through sliced generalized wasserstein geodesics. In *Thirty-seventh Conference on Neu-  
631 ral Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=n3XuYdvhNW>.
- 632 Gonzalo Mena and Jonathan Niles-Weed. Statistical bounds for entropic optimal transport: sample  
633 complexity and the central limit theorem. In *Advances in Neural Information Processing Systems*,  
634 pp. 4541–4551, 2019.
- 635 Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization  
636 for generative adversarial networks. In *International Conference on Learning Representations*,  
637 2018. URL <https://openreview.net/forum?id=BlQRgziT->.
- 638 James R Munkres. *Elements of algebraic topology*. CRC press, 2018.
- 639 Boris Muzellec and Marco Cuturi. Subspace detours: Building transport plans that are optimal on  
640 subspace projections. *Advances in Neural Information Processing Systems*, 32, 2019.

- 648 Kimia Nadjahi, Alain Durmus, Pierre E Jacob, Roland Badeau, and Umut Simsekli. Fast approx-  
649 imation of the sliced-wasserstein distance using concentration of random projections. *Advances*  
650 *in Neural Information Processing Systems*, 34:12411–12424, 2021.
- 651 Khai Nguyen and Nhat Ho. Amortized projection optimization for sliced wasserstein generative  
652 models. *Advances in Neural Information Processing Systems*, 35:36985–36998, 2022.
- 653 Khai Nguyen and Nhat Ho. Sliced wasserstein estimation with control variates. *arXiv preprint*  
654 *arXiv:2305.00402*, 2023.
- 655 Khai Nguyen and Nhat Ho. Energy-based sliced wasserstein distance. *Advances in Neural Informa-*  
656 *tion Processing Systems*, 36, 2024.
- 657 Khai Nguyen, Nhat Ho, Tung Pham, and Hung Bui. Distributional sliced-wasserstein and applica-  
658 tions to generative modeling. *arXiv preprint arXiv:2002.07367*, 2020.
- 659 Khai Nguyen, Nicola Bariletto, and Nhat Ho. Quasi-monte carlo for 3d sliced wasserstein. In  
660 *The Twelfth International Conference on Learning Representations*, 2024a. URL [https://](https://openreview.net/forum?id=Wd47f7HEXg)  
661 [openreview.net/forum?id=Wd47f7HEXg](https://openreview.net/forum?id=Wd47f7HEXg).
- 662 Khai Nguyen, Shujian Zhang, Tam Le, and Nhat Ho. Sliced wasserstein with random-path projecting  
663 directions. *arXiv preprint arXiv:2401.15889*, 2024b.
- 664 Tin D Nguyen, Brian L Trippe, and Tamara Broderick. Many processors, little time: MCMC for  
665 partitions via optimal transport couplings. In *International Conference on Artificial Intelligence*  
666 *and Statistics*, pp. 3483–3514. PMLR, 2022.
- 667 Trung Nguyen, Quang-Hieu Pham, Tam Le, Tung Pham, Nhat Ho, and Binh-Son Hua. Point-  
668 set distances for learning representations of 3d point clouds. In *Proceedings of the IEEE/CVF*  
669 *International Conference on Computer Vision (ICCV)*, pp. 10478–10487, 2021.
- 670 Sloan Nietert, Ziv Goldfeld, and Rachel Cummings. Outlier-robust optimal transport: Duality, struc-  
671 ture, and statistical analysis. In *International Conference on Artificial Intelligence and Statistics*,  
672 pp. 11691–11719. PMLR, 2022.
- 673 Jonathan Niles-Weed and Philippe Rigollet. Estimation of Wasserstein distances in the spiked trans-  
674 port model. *Bernoulli*, 28(4):2663–2688, 2022.
- 675 Jungin Park, Jiyoung Lee, and Kwanghoon Sohn. Bridging vision and language spaces with assign-  
676 ment prediction. *arXiv preprint arXiv:2404.09632*, 2024.
- 677 François-Pierre Paty and Marco Cuturi. Subspace robust wasserstein distances. In *International*  
678 *conference on machine learning*, pp. 5072–5081. PMLR, 2019.
- 679 Gabriel Peyré, Marco Cuturi, et al. Computational optimal transport: With applications to data  
680 science. *Foundations and Trends® in Machine Learning*, 11(5-6):355–607, 2019.
- 681 Thong Pham, Shohei Shimizu, Hideitsu Hino, and Tam Le. Scalable counterfactual distribution  
682 estimation in multivariate causal models. In *Conference on Causal Learning and Reasoning*  
683 *(CLear)*, 2024.
- 684 Michael Quellmalz, Robert Beinert, and Gabriele Steidl. Sliced optimal transport on the sphere.  
685 *Inverse Problems*, 39(10):105005, 2023.
- 686 Julien Rabin, Gabriel Peyré, Julie Delon, and Marc Berlot. Wasserstein barycenter and its applica-  
687 tion to texture mixing. In *International Conference on Scale Space and Variational Methods in*  
688 *Computer Vision*, pp. 435–446, 2011.
- 689 Joseph J Rotman. *An introduction to algebraic topology*, volume 119. Springer Science & Business  
690 Media, 2013.
- 691 Mahdi Saleh, Shun-Cheng Wu, Luca Cosmo, Nassir Navab, Benjamin Busam, and Federico  
692 Tombari. Bending graphs: Hierarchical shape matching using gated optimal transport. In *Pro-*  
693 *ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp.  
694 11757–11767, 2022.

- 702 Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen.  
703 Improved techniques for training gans. NIPS'16, pp. 2234–2242, Red Hook, NY, USA, 2016.  
704 Curran Associates Inc. ISBN 9781510838819.
- 705  
706 Tim Salimans, Han Zhang, Alec Radford, and Dimitris Metaxas. Improving gans using optimal  
707 transport. *arXiv preprint arXiv:1803.05573*, 2018.
- 708 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised  
709 learning using nonequilibrium thermodynamics. In *International conference on machine learn-*  
710 *ing*, pp. 2256–2265. PMLR, 2015.
- 711  
712 Justin Solomon, Fernando De Goes, Gabriel Peyré, Marco Cuturi, Adrian Butscher, Andy Nguyen,  
713 Tao Du, and Leonidas Guibas. Convolutional Wasserstein distances: Efficient optimal transporta-  
714 tion on geometric domains. *ACM Transactions on Graphics (TOG)*, 34(4):66, 2015.
- 715 Yuki Takezawa, Ryoma Sato, Zornitsa Kozareva, Sujith Ravi, and Makoto Yamada. Fixed sup-  
716 port tree-sliced Wasserstein barycenter. In *Proceedings of The 25th International Conference on*  
717 *Artificial Intelligence and Statistics*, volume 151, pp. 1120–1137. PMLR, 2022.
- 718  
719 C. Villani. *Optimal Transport: Old and New*, volume 338. Springer Science & Business Media,  
720 2008.
- 721 Dong Wang and Xiaoyang Tan. Unsupervised feature learning with c-svddnet. *Pattern Recognition*,  
722 60:473–485, 2016.
- 723  
724 Jie Wang, Rui Gao, and Yao Xie. Two-sample test with kernel projected Wasserstein distance. In  
725 *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume  
726 151, pp. 8022–8055. PMLR, 2022.
- 727  
728 Jonathan Weed and Quentin Berthet. Estimation of smooth densities in Wasserstein distance. In  
729 *Proceedings of the Thirty-Second Conference on Learning Theory*, volume 99, pp. 3118–3119,  
730 2019.
- 731 Jiqing Wu, Zhiwu Huang, Dinesh Acharya, Wen Li, Janine Thoma, Danda Pani Paudel, and Luc Van  
732 Gool. Sliced wasserstein generative models. In *Proceedings of the IEEE/CVF Conference on*  
733 *Computer Vision and Pattern Recognition*, pp. 3713–3722, 2019.
- 734  
735 Zhisheng Xiao, Karsten Kreis, and Arash Vahdat. Tackling the generative learning trilemma with  
736 denoising diffusion gans. *arXiv preprint arXiv:2112.07804*, 2021.
- 737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

756	NOTATION	
757		
758	$\mathbb{R}^d$	$d$ -dimensional Euclidean space
759	$\ \cdot\ _2$	Euclidean norm
760	$\langle \cdot, \cdot \rangle$	standard dot product
761	$\mathbb{S}^{d-1}$	$(d - 1)$ -dimensional hypersphere
762	$\theta$	unit vector
763	$\sqcup$	disjoint union
764	$L^1(X)$	space of Lebesgue integrable functions on $X$
765	$\mathcal{P}(X)$	space of probability distributions on $X$
766	$\mu, \nu$	measures
767	$\delta(\cdot)$	1-dimensional Dirac delta function
768	$\mathcal{U}(\mathbb{S}^{d-1})$	uniform distribution on $\mathbb{S}^{d-1}$
769	$\#$	pushforward (measure)
770	$\mathcal{C}(X, Y)$	space of continuous maps from $X$ to $Y$
771	$d(\cdot, \cdot)$	metric in metric space
772	$\mathbf{W}_p$	$p$ -Wasserstein distance
773	$\mathbf{SW}_p$	Sliced $p$ -Wasserstein distance
774	$\Gamma$	(rooted) subtree
775	$e$	edge in graph
776	$w_e$	weight of edge in graph
777	$l$	line, index of line
778	$\mathcal{L}$	system of lines, tree system
779	$\bar{\mathcal{L}}$	ground set of system of lines, tree system
780	$\Omega_{\mathcal{L}}$	topological space of system of lines
781	$\mathbb{L}_k^d$	space of systems of $k$ lines in $\mathbb{R}^d$
782	$\mathcal{T}$	tree structure in system of lines
783	$L$	number of tree systems
784	$k$	number of lines in a system of lines or a tree system
785	$\mathcal{R}$	original Radon Transform
786	$\mathcal{R}^\alpha$	Radon Transform on Systems of Lines
787	$\Delta_{k-1}$	$(k - 1)$ -dimensional standard simplex
788	$\alpha$	splitting map
789	$\mathbb{T}$	space of tree systems
790	$\sigma$	distribution on space of tree systems
791		
792		
793		
794		
795		
796		
797		
798		
799		
800		
801		
802		
803		
804		
805		
806		
807		
808		
809		

# Supplement for “Projection Optimal Transport on Tree-Ordered Lines”

## Table of Contents

---

810		
811		
812		
813		
814		
815	<b>A Tree System</b>	<b>16</b>
816		
817	A.1 System of Lines . . . . .	16
818	A.2 System of Lines with Tree Structures (Tree System) . . . . .	17
819	A.3 Topological Properties of Tree Systems . . . . .	17
820	A.4 Construction of Tree Systems . . . . .	20
821		
822		
823		
824	<b>B Radon Transform on Systems of Lines</b>	<b>21</b>
825	B.1 Space of Lebesgue integrable functions on a system of lines . . . . .	21
826	B.2 Probability distributions on a system of lines . . . . .	23
827		
828		
829	<b>C Max Tree-Sliced Wasserstein Distance on Systems of Lines.</b>	<b>23</b>
830		
831	<b>D Theoretical Proof for injectivity of TSW-SL</b>	<b>24</b>
832		
833	D.1 Proof of Theorem 5.2 . . . . .	24
834	D.2 Proof of Theorem C.2 . . . . .	25
835		
836	<b>E Experimental details</b>	<b>25</b>
837		
838	E.1 Gradient flows . . . . .	25
839	E.2 Color Transfer . . . . .	27
840	E.3 Generative adversarial network . . . . .	27
841	E.4 Denoising diffusion models . . . . .	30
842	E.5 Computational infrastructure . . . . .	31
843		
844		
845		
846	<b>F Broader Impact</b>	<b>32</b>
847		

## A TREE SYSTEM

In this section, we introduce the notion of a tree system, beginning with a collection of unstructured lines and progressively adding a tree structure to form a well-defined metric space with a tree metric. It is important to note that while some statements here differ slightly from those in the paper, the underlying ideas remain the same.

### A.1 SYSTEM OF LINES

We have a definition of lines by parameterization. Observe that, a line in  $\mathbb{R}^d$  is completely determined by a pair  $(x, \theta) \in \mathbb{R}^d \times \mathbb{S}^{d-1}$  via  $x + t \cdot \theta, t \in \mathbb{R}$ .

**Definition A.1** (Line and System of lines in  $\mathbb{R}^d$ ). A *line* in  $\mathbb{R}^d$  is an element  $(x, \theta)$  of  $\mathbb{R}^d \times \mathbb{S}^{d-1}$ , and the *image* of a line  $(x, \theta)$  is defined by:

$$\text{Im}(x, \theta) := \{x + t \cdot \theta : t \in \mathbb{R}\} \subset \mathbb{R}^d. \quad (14)$$

For  $k \geq 1$ , a *system of  $n$  lines* in  $\mathbb{R}^d$  is a sequence of  $k$  lines.

*Remark.* A line in  $\mathbb{R}^d$  is usually denoted, or indexed, by  $l = (x_l, \theta_l) \in \mathbb{R}^d \times \mathbb{S}^{d-1}$ . Here,  $x_l$  and  $\theta_l$  are called *source* and *direction* of  $l$ , respectively. Denote  $(\mathbb{R}^d \times \mathbb{S}^{d-1})^k$  by  $\mathbb{L}_k^d$ , which is the collection of systems of  $k$  lines in  $\mathbb{R}^d$ , and an element of  $\mathbb{L}_k^d$  is usually denoted by  $\mathcal{L}$ .

**Definition A.2** (Ground Set). The *ground set* of a system of lines  $\mathcal{L}$  is defined by:

$$\bar{\mathcal{L}} := \{(x, l) \in \mathbb{R}^d \times \mathcal{L} : x = x_l + t_x \cdot \theta_l \text{ for some } t_x \in \mathbb{R}\}.$$

For each element  $(x, l) \in \bar{\mathcal{L}}$ , we sometime write  $(x, l)$  as  $(t_x, l)$ , where  $t_x \in \mathbb{R}$ , which presents the parameterization of  $x$  on  $l$  by source  $x_l$  and direction  $\theta_l$ , as  $x = x_l + t_x \cdot \theta_l$ .

*Remark.* In other words, the ground set  $\bar{\mathcal{L}}$  is the disjoint union of images of lines in  $\mathcal{L}$ :

$$\bar{\mathcal{L}} = \bigsqcup_{l \in \mathcal{L}} \text{Im}(l).$$

This notation seems to be redundant, but will be helpful when we define functions on  $\bar{\mathcal{L}}$ .

## A.2 SYSTEM OF LINES WITH TREE STRUCTURES (TREE SYSTEM)

Consider a finite system of lines  $\mathcal{L}$  in  $\mathbb{R}^d$ . Assume that these lines are geometrically distinct, i.e. their images are distinct. Define the graph  $\mathcal{G}_{\mathcal{L}}$  associated with  $\mathcal{L}$ , where  $\mathcal{L}$  is the set of nodes in  $\mathcal{G}_{\mathcal{L}}$ , and two nodes are adjacent if the two corresponding lines intersect each other. Here, saying two lines in  $\mathbb{R}^d$  intersect means their images have exactly one point in common.

**Definition A.3** (Connected system of lines).  $\mathcal{L}$  is called *connected* if its associated graph  $\mathcal{G}_{\mathcal{L}}$  is connected.

*Remark.* Intuitively, each edge of  $\mathcal{G}_{\mathcal{L}}$  represents the intersection of its endpoints. If  $\mathcal{L}$  is connected, for every two points that each one lies on some lines in  $\mathcal{L}$ , one can travel to the other through lines in  $\mathcal{L}$ .

From now on, we will only consider the case  $\mathcal{L}$  is connected. Recall the notion of a spanning tree of a graph  $\mathcal{G}$ , which is a subgraph of  $\mathcal{G}$  that contains all nodes of  $\mathcal{G}$ , and also is a tree.

**Definition A.4** (Tree system of lines). Let  $\mathcal{L}$  be a connected system of lines. A spanning tree  $\mathcal{T}$  of  $\mathcal{G}_{\mathcal{L}}$  is called a *tree structure* of  $\mathcal{L}$ . A pair  $(\mathcal{L}, \mathcal{T})$  consists of a connected system of lines  $\mathcal{L}$  and a tree structure  $\mathcal{T}$  of  $\mathcal{L}$  is called a *tree system of lines*.

*Remark.* For short, we usually call a tree system of lines as a *tree system*. In a tree system  $(\mathcal{L}, \mathcal{T})$ , images of two lines of  $\mathcal{L}$  can intersect each other even when they are not adjacent in  $\mathcal{T}$ .

Let  $r$  be an arbitrary line of  $\mathcal{L}$ . Denote  $\mathcal{T}_r$  as the tree  $\mathcal{T}$  rooted at  $r$ , and denote the (rooted) tree system as  $(\mathcal{L}, \mathcal{T}_r)$  if we want to specify the root.

**Definition A.5** (Depth of lines in a tree system). Let  $(\mathcal{L}, \mathcal{T}_r)$  be a tree system. For each  $m \geq 0$ , a line  $l \in \mathcal{L}$  is called a *line of depth  $m$*  if the (unique) path from  $r$  to  $l$  in  $\mathcal{T}$  has length  $m$ . Denote  $\mathcal{L}_m$  as the *set of lines of depth  $m$* .

*Remark.* Note that  $\mathcal{L}_0 = \{r\}$ . Let  $T$  be the maximum length of paths in  $\mathcal{T}$  start from  $r$ , which is called the *depth of the line system*.  $\mathcal{L}$  has a partition as  $\mathcal{L} = \mathcal{L}_0 \sqcup \mathcal{L}_1 \sqcup \dots \sqcup \mathcal{L}_T$ .

For  $l \in \mathcal{L}$  that is not the root, denote  $\text{pr}(l) \in \mathcal{L}$  as the *parent of  $l$* , i.e. the (unique) node on the unique path from  $l$  to  $r$  that is adjacent to  $l$ . Note that, by definition,  $l$  and  $\text{pr}(l)$  intersect each other. We sometimes omit the root when the context is clear.

**Definition A.6** (Canonical tree system). A tree system  $(\mathcal{L}, \mathcal{T})$  is called a *canonical tree system* if for all  $l \in \mathcal{L}$  that is not the root, the intersection of  $l$  and  $\text{pr}(l)$  is the source  $x_l$  of  $l$ .

*Remark.* In other words, in a canonical tree system, a line that differs from the root will have its source lies on its parent. For the rest of the paper, a tree system  $(\mathcal{L}, \mathcal{T})$  will be considered to be a canonical tree system.

## A.3 TOPOLOGICAL PROPERTIES OF TREE SYSTEMS

We will introduce the notion of the topological space of a tree system. Let  $(\mathcal{L}, \mathcal{T})$  be a (canonical) tree system. Consider a graph where the nodes are elements of  $\bar{\mathcal{L}}$ ;  $(x, l)$  and  $(x', l')$  are adjacent if and only if one of the following conditions holds:

1.  $l = \text{pr}(l')$ ,  $x = x'$ , and  $x'$  is the source of  $l'$ .
2.  $l' = \text{pr}(l)$ ,  $x = x'$ , and  $x$  is the source of  $l$ .

Let  $\sim$  be the relation on  $\bar{\mathcal{L}}$  such that  $(x, l) \sim (x', l')$  if and only if  $(x, l)$  and  $(x', l')$  are connected in the above graph. By design,  $\sim$  is an equivalence relation on  $\bar{\mathcal{L}}$ . The set of all equivalence classes in  $\bar{\mathcal{L}}$  with respect to the equivalence relation  $\sim$  as  $\Omega_{\mathcal{L}} = \bar{\mathcal{L}} / \sim$ .

*Remark.* In other words, we identify the source of lines to the corresponding point on its parent.

We recall the notion of disjoint union topology and quotient topology in (Hatcher, 2005). For a line  $l$  in  $\mathbb{R}^d$ , the image  $\text{Im}(l)$  is a topological space, moreover, a metric space, that is homeomorphic and isometric to  $\mathbb{R}$  via the map  $t \mapsto x_l + t \cdot \theta_l$ . The metric on  $\text{Im}(l)$  is  $d_l(x, x') = |t_x - t_{x'}|$  for all  $x, x' \in \text{Im}(l)$ . For each  $l \in \mathcal{L}$ , consider the injection map:

$$f_l : \text{Im}(l) \longrightarrow \bigsqcup_{l \in \mathcal{L}} \text{Im}(l) = \bar{\mathcal{L}}$$

$$x \longmapsto (x, l).$$

$\bar{\mathcal{L}} = \bigsqcup_{l \in \mathcal{L}} \text{Im}(l)$  now becomes a topological space with the disjoint union topology, i.e. the finest topology on  $\bar{\mathcal{L}}$  such that the map  $f_l$  is continuous for all  $l \in \mathcal{L}$ . Also, consider the quotient map:

$$\pi : \bar{\mathcal{L}} \longrightarrow \Omega_{\mathcal{L}}$$

$$(x, l) \longmapsto [(x, l)].$$

$\Omega_{\mathcal{L}}$  now becomes a topological space with the quotient topology, i.e. the finest topology on  $\Omega_{\mathcal{L}}$  such that the map  $\pi$  is continuous.

**Definition A.7** (Topological space of a tree system). The topological space  $\Omega_{\mathcal{L}}$  is called the *topological space of a tree system*  $(\mathcal{L}, \mathcal{T})$ .

*Remark.* In other words,  $\Omega_{\mathcal{L}}$  is formed by gluing all images  $\text{Im}(l)$  along the relation  $\sim$ .

We show that the topological space  $\Omega_{\mathcal{L}}$  is metrizable.

**Definition A.8** (Paths in  $\Omega_{\mathcal{L}}$ ). For  $a$  and  $b$  in  $\Omega_{\mathcal{L}}$  with  $a \neq b$ , a *path from  $a$  to  $b$  in  $\Omega_{\mathcal{L}}$*  is a continuous injective map  $\gamma : [0, 1] \rightarrow \Omega_{\mathcal{L}}$  where  $\gamma(0) = a$  and  $\gamma(1) = b$ . By convention, for  $a$  in  $\Omega_{\mathcal{L}}$ , the path from  $a$  to  $a$  in  $\Omega_{\mathcal{L}}$  is the constant map  $\gamma : [0, 1] \rightarrow \Omega_{\mathcal{L}}$  such that  $\gamma(t) = a$  for all  $t \in [0, 1]$ . For a path  $\gamma$  from  $a$  to  $b$ , the *image of  $\gamma$*  is defined by:

$$\text{Im}(\gamma) := \gamma([0, 1]) \subset \Omega_{\mathcal{L}}. \quad (15)$$

**Theorem A.9** (Existence and uniqueness of path in  $\Omega_{\mathcal{L}}$ ). *For all  $a$  and  $b$  in  $\Omega_{\mathcal{L}}$ , there exist a path  $\gamma$  from  $a$  to  $b$  in  $\Omega_{\mathcal{L}}$ . Moreover,  $\gamma$  is unique up to a re-parameterization, i.e. if  $\gamma$  and  $\gamma'$  are two path  $\gamma$  from  $a$  to  $b$  in  $\Omega_{\mathcal{L}}$ , there exist a homeomorphism  $\varphi : [0, 1] \rightarrow [0, 1]$  such that  $\gamma = \gamma' \circ \varphi$ .*

*Proof.* All previous results we state in this proof can be found in (Munkres, 2018; Rotman, 2013; Hatcher, 2005). For two point  $a, b$  on the real line  $\mathbb{R}$ , all paths from  $a$  to  $b$  are homotopic to each other. In other words, all paths from  $a$  to  $b$  are homotopic to the canonical path:

$$\gamma_{a,b} : [0, 1] \longrightarrow \mathbb{R}$$

$$t \longmapsto (1-t) \cdot a + t \cdot b.$$

Now consider two point  $a, b$  on space  $\Omega_{\mathcal{L}}$ . Observe that  $\Omega_{\mathcal{L}}$  is path-connected by design and by the fact that  $\mathbb{R}$  is path-connected. Consider a curve from  $a$  to  $b$  on  $\Omega_{\mathcal{L}}$ , i.e. a continuous map  $f : [0, 1] \rightarrow \Omega_{\mathcal{L}}$ , and consider the set consists of sources of lines in  $\mathcal{L}$  that lie on the curve  $f$ , i.e. all the sources that belong to  $f([0, 1])$ . We choose the curve  $f$  that has the smallest set of sources. By the tree structure added to  $\mathcal{L}$ , all curves from  $a$  to  $b$  have the set of sources that contains the set of sources of  $f$ . We denote the sources belong to this set of  $f$  as  $s_1, \dots, s_{k-1}$ , and defined:

$$x_i = \inf f^{-1}(s_i) \text{ for all } 1 \leq i \leq k-1.$$

We reindex  $s_i$  such that:

$$x_1 \leq \dots \leq x_{k-1}$$

For convention, we define  $s_0 = a$  and  $s_k = b$ . By design, for  $i = 0, \dots, k-1$ , we have  $s_i$  and  $s_{i+1}$  line on the same line in  $\mathcal{L}$ . So by the result of paths on  $\mathbb{R}$ , there exist a path  $\gamma_i$  from  $s_i$  to  $s_{i+1}$  on  $\Omega_{\mathcal{L}}$ . Gluing  $\gamma_0, \gamma_1, \dots, \gamma_{k-1}$  to get a path  $\gamma$  from  $s_0 = a$  to  $s_k = b$  on  $\Omega_{\mathcal{L}}$  by:

$$\begin{aligned} \gamma : [0, 1] &\longrightarrow \Omega_{\mathcal{L}} \\ t &\longmapsto \gamma_i(k \cdot t - i) \text{ if } t \in \left[ \frac{i}{k}, \frac{i+1}{k} \right], i = 0, \dots, k-1. \end{aligned}$$

It is clear to check  $\gamma$  is a path from  $a$  to  $b$  on  $\Omega$ , and the uniqueness (up to re-parameterization) of  $\gamma$  comes from homotopy of paths in  $\mathbb{R}$ .  $\square$

*Remark.* The image of a path from  $a$  to  $b$  does not depend on the chosen path  $\gamma$  by the uniqueness property. Indeed, for a homeomorphism  $\varphi: [0, 1] \rightarrow [0, 1]$ , we have  $\gamma([0, 1]) = \gamma \circ \varphi([0, 1])$ . Denote the image of *any path* from  $a$  to  $b$  by  $P_{a,b}$ .

Let  $\mu$  be the standard Borel measure on  $\mathbb{R}$ , i.e.  $\mu((a, b]) = b - a$  for every half-open interval  $(a, b]$  in  $\mathbb{R}$ . For  $l \in \mathcal{L}$ , denote  $\mu_l$  as the pushforward of  $\mu$  by the map  $t \mapsto x_l + t \cdot \theta_l$ , which is a Borel measure on  $\text{Im}(l)$ . Denote the  $\sigma$ -algebra of Borel sets in  $\bar{\mathcal{L}}$  and  $\Omega_{\mathcal{L}}$  as  $\mathcal{B}(\bar{\mathcal{L}})$  and  $\mathcal{B}(\Omega_{\mathcal{L}})$ , respectively.

**Definition A.10** (Borel measure on  $\bar{\mathcal{L}}$  and  $\Omega_{\mathcal{L}}$ ). The map  $\mu_{\bar{\mathcal{L}}}: \mathcal{B}(\Omega_{\mathcal{L}}) \rightarrow [0, \infty)$  that is defined by:

$$\mu_{\bar{\mathcal{L}}}(B) := \sum_{l \in \mathcal{L}} \mu_l(f_l^{-1}(B)), \quad \forall B \in \mathcal{B}(\bar{\mathcal{L}}),$$

is called the *Borel measure on  $\bar{\mathcal{L}}$* . Define the *Borel measure on  $\Omega_{\mathcal{L}}$* , denoted by  $\mu_{\Omega_{\mathcal{L}}}$ , as the pushforward of  $\mu_{\bar{\mathcal{L}}}$  by the map  $\pi: \bar{\mathcal{L}} \rightarrow \Omega_{\mathcal{L}}$ .

It is straightforward to show that  $\mu_{\bar{\mathcal{L}}}$  is well-defined, and indeed a Borel measure of  $\bar{\mathcal{L}}$ . As a corollary,  $\mu_{\Omega_{\mathcal{L}}}$  is also a Borel measure of  $\Omega_{\mathcal{L}}$ .

*Remark.* By abuse of notation, we sometimes simply denote both of  $\mu_{\bar{\mathcal{L}}}$  and  $\mu_{\Omega_{\mathcal{L}}}$  as  $\mu_{\mathcal{L}}$ .

**Theorem A.11** ( $\Omega_{\mathcal{L}}$  is metrizable by a tree metric). Define the map  $d_{\Omega}: \Omega_{\mathcal{L}} \times \Omega_{\mathcal{L}} \rightarrow [0, \infty)$  by:

$$d_{\mathcal{L}}(a, b) := \mu_{\mathcal{L}}(P_{a,b}), \quad \forall a, b \in \Omega_{\mathcal{L}}. \quad (16)$$

Then  $d_{\mathcal{L}}$  is a metric on  $\Omega_{\mathcal{L}}$ , which makes  $(\Omega_{\mathcal{L}}, d_{\mathcal{L}})$  a metric space. Moreover,  $d_{\mathcal{L}}$  is a tree metric, and the topology on  $\Omega_{\mathcal{L}}$  induced by  $d_{\mathcal{L}}$  is identical to the topology of  $\Omega_{\mathcal{L}}$ .

*Proof.* It is straightforward to check that  $d_{\mathcal{L}}$  is positive definite and symmetry. We show the triangle inequality holds for  $d_{\mathcal{L}}$ . Let  $a, b, c$  be points of  $\Omega_{\mathcal{L}}$ . It is enough to show that  $P_{a,c}$  is a subset of  $P_{a,b} \cup P_{b,c}$ . Let  $\gamma_0, \gamma_1$  be paths on  $\Omega$  from  $a$  to  $b$  and from  $b$  to  $c$ , respectively. Consider the curve from  $a$  to  $c$  on  $\Omega$  defined by:

$$\begin{aligned} \gamma : [0, 1] &\longrightarrow \Omega_{\mathcal{L}} \\ t &\longmapsto \gamma_i(2 \cdot t - i) \text{ if } t \in \left[ \frac{i}{2}, \frac{i+1}{2} \right], i = 0, 1. \end{aligned}$$

It is clear that  $\gamma$  is a curve from  $a$  to  $c$ . We have  $\gamma([0, 1])$  is exactly the union of  $P_{a,b}$  and  $P_{b,c}$ . As in the proof of Theorem A.9, the set of sources of  $\gamma$  contains the set of sources lying on the path from  $a$  to  $c$ . So  $\gamma([0, 1])$  contains  $P_{a,c}$ .  $\square$

We have the below corollary says that: If we take finite points on  $\Omega_{\mathcal{L}}$ , together with the sources of lines, it induces a tree (as a graph) with nodes are these points; Moreover, we have a tree metric on this tree which is  $d_{\mathcal{L}}$ .

**Corollary A.12.** Let  $y_1, y_2, \dots, y_m$  be points on  $\Omega_{\mathcal{L}}$ . Consider the graph, where  $\{y_1, \dots, y_m\} \cup \{x_l : l \in \mathcal{L}\}$  is the node set, and two nodes are adjacent if the (unique) path between this two nodes on  $\Omega_{\mathcal{L}}$  does not contain any node, except them. Then this graph is a rooted tree at  $x_r$ , with an induced tree metric from  $d_{\mathcal{L}}$ .

1026 A.4 CONSTRUCTION OF TREE SYSTEMS  
1027

1028 We present a way to construct a tree system in  $\mathbb{R}^d$ . First, we have a way to describe the structure of  
1029 a rooted tree by a sequence of vectors.

1030 **Definition A.13** (Tree representation). Let  $T$  be a nonnegative integer, and  $n_1, \dots, n_T$  be  $T$  positive  
1031 integer. A sequence  $s = \{x_i\}_{i=0}^T$ , where  $x_i$  is a vector of  $n_i$  nonnegative numbers, is called a *tree*  
1032 *representation* if  $x_0 = [1]$ , and for all  $1 \leq i \leq T$ ,  $n_i$  is equal to the sum of all entries in vector  $x_{i-1}$ .

1033 **Example A.14.** For  $T = 5$  and  $\{n_i\}_{i=1}^5 = \{1, 3, 4, 2, 3\}$ , the sequence:  
1034

$$\begin{aligned} 1035 \quad s : x_0 &= [1] \\ 1036 \quad &\rightarrow x_1 = [3] \\ 1037 \quad &\rightarrow x_2 = [2, 1, 1] \\ 1038 \quad &\rightarrow x_3 = [1, 0, 2, 0] \\ 1039 \quad &\rightarrow x_4 = [1, 2] \\ 1040 \quad &\rightarrow x_5 = [0, 0, 1]. \end{aligned}$$

1041  
1042 is a tree representation.  
1043

1044 For a tree representation  $s = \{x_i\}_{i=0}^T$ , a *tree system of type  $s$*  is a tree system that is inductively  
1045 constructed step-by-step as follows:

1046  
1047 Step 0. Sample a point  $x_r \in \mathbb{R}^d$  and a direction  $\theta_r \in \mathbb{S}^{d-1}$ . Define  $r$  as the line that passes through  
1048  $x_r$  with direction  $\theta_r$ . We call  $r$  as the line of depth 0.

1049 Step  $i$ . On the  $j$ -th line of depth  $(i-1)$ , sample  $(x_i)_j$  points where  $(x_i)_j$  is the  $j$ -th entry of vector  
1050  $x_i$ . For each of these points, denoted as  $x_l$ , sample a direction  $\theta_l \in \mathbb{S}^{d-1}$ , and define  $l$  is  
1051 the line that passes through  $x_l$  with direction  $\theta_l$ . We call the set of all lines sampled at this  
1052 step as the set of lines of depth  $i$  and order them by the order that they are sampled.

1053 By this construction, we will get a system of lines  $\mathcal{L}$  in  $\mathbb{R}^d$ , together with a tree structure  $\mathcal{T}_r$ . The  
1054 pair  $(\mathcal{L}, \mathcal{T}_r)$  forms a tree system, which is canonical by design, and is said to be of type  $s$ . Denote  
1055  $\mathbb{T}_s$  as the *set of all tree systems of type  $s$* .

1056 **Remark.** A tree system in of type  $s$  has  $k = \sum_{i=0}^T \sum_{j=1}^{n_i} (x_i)_j$  lines. Observe that constructing a  
1057 tree system of type  $s$  only depends on sampling  $k$  points and  $k$  directions, so by some assumptions  
1058 on the probability distribution of these points and directions, we will have a probability distribution  
1059 on  $\mathbb{T}_s$ . Note that:

- 1061 1.  $x_r$  is sampled from a distribution on  $\mathbb{R}^d$ ;
- 1062 2. For all  $l \neq r$ ,  $x_l$  is sampled from a distribution on  $\mathbb{R}$ ;
- 1063 3. For all  $l$ ,  $\theta_l$  is sampled from a distribution on  $\mathbb{S}^{d-1}$ .

1064  
1065 We have some examples of tree presentations  $s$  and distribution on  $\mathbb{T}_s$ .  
1066

1067 **Example A.15** (Lines pass through origin). Consider the tree representation  $s$ :  
1068

$$1069 \quad s : [1], \tag{17}$$

1070 and the distributions on  $\mathbb{T}_s$  is determined by:  
1071

- 1072 1.  $x_r = 0 \in \mathbb{R}^d$ ;
- 1073 2.  $\theta_r \sim \mathcal{U}(\mathbb{S}^{d-1})$ .

1074  
1075 In this case,  $\mathbb{T}_s$  is identical to the set of lines that pass through the origin 0.  
1076

1077 **Example A.16** (Concurrent lines). Consider the tree representation  $s$ :  
1078

$$1079 \quad s : [1] \rightarrow [k-1], \tag{18}$$

and the distributions on  $\mathbb{T}_s$  is determined by:

- 1080 1.  $x_r \sim \mu_r$  for an  $\mu_r \in \mathcal{P}(\mathbb{R}^d)$ ;
- 1081 2. For all  $l \neq r$ ,  $x_l = x_r$ ;
- 1082 3. For all  $l$ ,  $\theta_l \sim \mathcal{U}(\mathbb{S}^{d-1})$ ;
- 1083 4.  $x_r$  and all  $\theta_l$ 's are pairwise independent.

1084 In this case,  $\mathbb{T}_s$  is identical to the set of all tuples of  $n$  concurrent lines.

1085 **Example A.17** (Series of lines). Consider the tree representation  $s$ :

$$1086 s : [1] \rightarrow [1] \rightarrow \dots \rightarrow [1], \quad (19)$$

1087 and the distributions on  $\mathbb{T}_s$  is determined by:

- 1088 1.  $x_r \sim \mu_r$  for an  $\mu_r \in \mathcal{P}(\mathbb{R}^d)$ ;
- 1089 2. For all  $l \neq r$ ,  $x_l \sim \mu_l$  for an  $\mu_l \in \mathcal{P}(\mathbb{R})$ ;
- 1090 3. For all  $l$ ,  $\theta_l \sim \mathcal{U}(\mathbb{S}^{d-1})$ ;
- 1091 4. All  $x_l$ 's and all  $\theta_l$ 's are pairwise independent.

1092 In this case, we call  $\mathbb{T}_s$  as the set of all series of  $k$  lines. This is the same as the sampling process in Subsection 3.3 and Algorithm 1.

1093 **Example A.18.** For an arbitrary tree representation  $s$ , the distributions on  $\mathbb{T}_s$  is determined by:

- 1094 1.  $x_r$  is sampled from the uniform distribution on a bounded subset of  $\mathbb{R}^d$ , for instance,  $\mu_r \sim \mathcal{U}([0, 1]^d)$ ;
- 1095 2. For  $l \neq r$ ,  $x_l$  will be sampled from the uniform distribution on a bounded subset of  $\mathbb{R}$ , for instance,  $\mu_l \sim \mathcal{U}([0, 1])$ ;
- 1096 3. For all  $l$ ,  $\theta_l$  will be sampled from the uniform distribution on  $\mathbb{S}^{d-1}$ , i.e  $\theta_l \sim \mathcal{U}(\mathbb{S}^{d-1})$ ;
- 1097 4. Together with assumptions on independence between all  $x_r$ 's and all  $\theta_l$ 's.

## 1100 B RADON TRANSFORM ON SYSTEMS OF LINES

1101 We introduce the notions of the space of Lebesgue integrable functions and the space of probability distributions on a system of lines. Let  $\mathcal{L}$  be a system of  $k$  lines.

### 1102 B.1 SPACE OF LEBESGUE INTEGRABLE FUNCTIONS ON A SYSTEM OF LINES

1103 Denote  $L^1(\mathbb{R}^d)$  as the space of Lebesgue integrable functions on  $\mathbb{R}^d$  with norm  $\|\cdot\|_1$ :

$$1104 L^1(\mathbb{R}^d) = \left\{ f: \mathbb{R}^d \rightarrow \mathbb{R} : \|f\|_1 = \int_{\mathbb{R}^d} |f(x)| dx < \infty \right\}. \quad (20)$$

1105 As usual, two functions  $f_1, f_2 \in L^1(\mathbb{R}^d)$  are considered to be identical if  $f_1(x) = f_2(x)$  almost everywhere on  $\mathbb{R}^d$ .

1106 **Definition B.1** (Lebesgue integrable function on a system of lines). A *Lebesgue integrable function* on  $\mathcal{L}$  is a function  $f: \bar{\mathcal{L}} \rightarrow \mathbb{R}$  such that:

$$1107 \|f\|_{\mathcal{L}} := \sum_{l \in \mathcal{L}} \int_{\mathbb{R}} |f(t_x, l)| dt_x < \infty. \quad (21)$$

1108 The space of Lebesgue integrable functions on  $\mathcal{L}$  is denoted by:

$$1109 L^1(\mathcal{L}) := \left\{ f: \bar{\mathcal{L}} \rightarrow \mathbb{R} : \|f\|_{\mathcal{L}} = \sum_{l \in \mathcal{L}} \int_{\mathbb{R}} |f(t_x, l)| dt_x < \infty \right\}. \quad (22)$$

1134 *Remark.* As an analog of integrable functions on  $\mathbb{R}^d$ , two functions  $f_1, f_2 \in L^1(\mathcal{L})$  are considered  
 1135 to be identical if  $f_1(x) = f_2(x)$  almost everywhere on  $\bar{\mathcal{L}}$ . The space  $L^1(\mathcal{L})$  with norm  $\|\cdot\|_{\mathcal{L}}$  is a  
 1136 Banach space.

1137 Recall that  $\mathcal{L}$  has  $k$  lines, we denote the  $(k-1)$ -dimensional standard simplex as  $\Delta_{k-1} =$   
 1138  $\{(a_l)_{l \in \mathcal{L}} : a_l \geq 0 \text{ and } \sum_{l \in \mathcal{L}} a_l = 1\} \subset \mathbb{R}^k$ . Let  $\mathcal{C}(\mathbb{R}^d, \Delta_{k-1})$  be the space of continuous function  
 1139 from  $\mathbb{R}^d$  to  $\Delta_{k-1}$ . A map in  $\mathcal{C}(\mathbb{R}^d, \Delta_{k-1})$  is called a *splitting map*. Let  $\mathcal{L}$  be a system of lines  
 1140 in  $\mathbb{L}_k^d$ ,  $\alpha$  be a map in  $\mathcal{C}(\mathbb{R}^d, \Delta_{k-1})$ , we define an operator associated to  $\alpha$  that transforms a Lebesgue  
 1141 integrable functions on  $\mathbb{R}^d$  to a Lebesgue integrable functions on  $\mathcal{L}$ . For  $f \in L^1(\mathbb{R}^d)$ , define:

$$1142 \mathcal{R}_{\mathcal{L}}^{\alpha} f : \bar{\mathcal{L}} \longrightarrow \mathbb{R}$$

$$1143 (x, l) \longmapsto \int_{\mathbb{R}^d} f(y) \cdot \alpha(y)_l \cdot \delta(t_x - \langle y - x_l, \theta_l \rangle) dy,$$

1144 where  $\delta$  is the 1-dimensional Dirac delta function.

1145 **Theorem B.2.** For  $f \in L^1(\mathbb{R}^d)$ , we have  $\mathcal{R}_{\mathcal{L}}^{\alpha} f \in L^1(\mathcal{L})$ . Moreover, we have  $\|\mathcal{R}_{\mathcal{L}}^{\alpha} f\|_{\mathcal{L}} \leq \|f\|_1$ . In  
 1146 other words, the operator:

$$1147 \mathcal{R}_{\mathcal{L}}^{\alpha} : L^1(\mathbb{R}^d) \rightarrow L^1(\mathcal{L}), \quad (23)$$

1148 is well-defined, and is a linear operator.

1149 *Proof.* Let  $f \in L^1(\mathbb{R}^d)$ . We show that  $\|\mathcal{R}_{\mathcal{L}}^{\alpha} f\|_{\mathcal{L}} \leq \|f\|_1$ . Indeed,

$$1150 \|\mathcal{R}_{\mathcal{L}}^{\alpha} f\|_{\mathcal{L}} = \sum_{l \in \mathcal{L}} \int_{\mathbb{R}} |\mathcal{R}_{\mathcal{L}}^{\alpha} f(t_x, l)| dt_x \quad (24)$$

$$1151 = \sum_{l \in \mathcal{L}} \int_{\mathbb{R}} \left| \int_{\mathbb{R}^d} f(y) \cdot \alpha(y)_l \cdot \delta(t_x - \langle y - x_l, \theta_l \rangle) dy \right| dt_x \quad (25)$$

$$1152 \leq \sum_{l \in \mathcal{L}} \int_{\mathbb{R}^d} \left( \int_{\mathbb{R}} |f(y)| \cdot \alpha(y)_l \cdot \delta(t_x - \langle y - x_l, \theta_l \rangle) \cdot dt_x \right) dy \quad (26)$$

$$1153 = \sum_{l \in \mathcal{L}} \int_{\mathbb{R}^d} |f(y)| \cdot \alpha(y)_l \cdot \left( \int_{\mathbb{R}} \delta(t_x - \langle y - x_l, \theta_l \rangle) dt_x \right) dy \quad (27)$$

$$1154 = \sum_{l \in \mathcal{L}} \int_{\mathbb{R}^d} |f(y)| \cdot \alpha(y)_l dy \quad (28)$$

$$1155 = \int_{\mathbb{R}^d} |f(y)| \cdot \sum_{l \in \mathcal{L}} \alpha(y)_l dy \quad (29)$$

$$1156 = \int_{\mathbb{R}^d} |f(y)| dy \quad (30)$$

$$1157 = \|f\|_1 < \infty. \quad (31)$$

1158 So the operator  $\mathcal{R}_{\mathcal{L}}^{\alpha}$  is well-defined, and is clearly a linear operator.  $\square$

1159 **Definition B.3** (Radon transform on system of lines). For  $\alpha \in \mathcal{C}(\mathbb{R}^d, \Delta_{k-1})$ , the operator  $\mathcal{R}^{\alpha}$ :

$$1160 \mathcal{R}^{\alpha} : L^1(\mathbb{R}^d) \longrightarrow \prod_{\mathcal{L} \in \mathbb{L}_k^d} L^1(\mathcal{L})$$

$$1161 f \longmapsto (\mathcal{R}_{\mathcal{L}}^{\alpha} f)_{\mathcal{L} \in \mathbb{L}_k^d}.$$

1162 is called the *Radon transform on a system of lines*.

1163 Many variants of Radon transform require the transforms to be injective. We show that the injectivity  
 1164 holds in the Radon transform on a system of lines.

1165 **Theorem B.4.**  $\mathcal{R}^{\alpha}$  is injective.

1188 *Proof.* Since  $\mathcal{R}^\alpha$  is linear, we show that if  $\mathcal{R}^\alpha f = 0$ , then  $f = 0$ . Let  $f \in L^1(\mathbb{R}^d)$  such that  
 1189  $\mathcal{R}^\alpha f = 0$ , which means  $\mathcal{R}_\mathcal{L}^\alpha f = 0$  for all  $\mathcal{L} \in \mathbb{L}_n^d$ . Fix a line index  $l$ , consider the operator:  
 1190

$$1191 \int_{\mathbb{R}^d} f(y) \cdot \alpha(y)_l \cdot \delta(t_x - \langle y - x_l, \theta_l \rangle) dy = 0, \forall t_x \in \mathbb{R}, (x_l, \theta_l) \in \mathbb{R}^d \times \mathbb{S}^{d-1}. \quad (32)$$

1192 Note that for index  $l$ ,  $f(y) \cdot \alpha(y)_l$  is a function of  $y$ . Let  $x_l$  be fixed and  $\theta_l$  varies in  $\mathbb{R}^d$ . By the  
 1193 injectivity of the usual Radon transform (Helgason & Helgason, 2011), we have  $f(x) \cdot \alpha(x)_l = 0$  for  
 1194 all  $x \in \mathbb{R}^d$ . This holds for all line index  $l$ , so  $f(x) = \sum_l f(x) \cdot \alpha(x)_l = 0$ . So  $\mathcal{R}^\alpha$  is injective.  $\square$   
 1195

1196 *Remark.* By the proof, we can show a stronger result as follows: Let  $A$  be a subset of  $\mathbb{L}_k^d$  such that  
 1197 for every index  $l$  and  $\theta \in \mathbb{S}^{d-1}$ , there exists  $\mathcal{L} \in A$  such that  $\theta_l = \theta$ , where  $\theta_l$  is the direction of line  
 1198 with index  $l$  in  $\mathcal{L}$ . Roughly speaking, the set of directions in  $\mathcal{L}$  is  $(\mathbb{S}^{d-1})^k$ .  
 1199

## 1201 B.2 PROBABILITY DISTRIBUTIONS ON A SYSTEM OF LINES

1202 Denote  $\mathcal{P}(\mathbb{R}^d)$  as the space of all probability distribution on  $\mathbb{R}^d$ :  
 1203

$$1204 \mathcal{P}(\mathbb{R}^d) = \{f: \mathbb{R}^d \rightarrow [0, \infty) : \|f\|_1 = 1\} \subset L^1(\mathbb{R}^d). \quad (32)$$

1205 **Definition B.5** (Probability distribution on a system of lines). Let  $\mathcal{L}$  be a system of lines. A *probabil-*  
 1206 *ity distribution on  $\mathcal{L}$*  is a function  $f \in L^1(\mathcal{L})$  such that  $f: \bar{\mathcal{L}} \rightarrow [0, \infty)$  and  $\|f\|_\mathcal{L} = 1$ . The *space*  
 1207 *of probability distribution on  $\mathcal{L}$*  is defined by:  
 1208

$$1209 \mathcal{P}(\mathcal{L}) := \{f: \bar{\mathcal{L}} \rightarrow [0, \infty) : \|f\|_\mathcal{L} = 1\} \subset L^1(\mathcal{L}). \quad (33)$$

1210 **Corollary B.6.** For  $f \in \mathcal{P}^1(\mathbb{R}^d)$ , we have  $\mathcal{R}_\mathcal{L}^\alpha f \in \mathcal{P}(\mathcal{L})$ . In other words, the restricted of Radon  
 1211 Transform:  
 1212

$$1213 \mathcal{R}_\mathcal{L}^\alpha: \mathcal{P}(\mathbb{R}^d) \rightarrow \mathcal{P}(\mathcal{L}), \quad (34)$$

1214 is well-defined.  
 1215

1216 *Proof.* Let  $f \in \mathcal{P}^1(\mathbb{R}^d)$ . It is clear that  $\mathcal{R}_\mathcal{L}^\alpha f: \bar{\mathcal{L}} \rightarrow [0, \infty)$ . We show that  $\|\mathcal{R}_\mathcal{L}^\alpha f\|_\mathcal{L} = 1$ . Indeed,  
 1217

$$1218 \|\mathcal{R}_\mathcal{L}^\alpha f\|_\mathcal{L} = \sum_{l \in \mathcal{L}} \int_{\mathbb{R}} \mathcal{R}_\mathcal{L}^\alpha f(t_x, l) dt_x \quad (35)$$

$$1219 = \sum_{l \in \mathcal{L}} \int_{\mathbb{R}} \left( \int_{\mathbb{R}^d} f(y) \cdot \alpha(y)_l \cdot \delta(t_x - \langle y - x_l, \theta_l \rangle) dy \right) dt_x \quad (36)$$

$$1220 = \int_{\mathbb{R}^d} f(y) dy = 1. \quad (37)$$

1221 So  $\mathcal{R}_\mathcal{L}^\alpha f \in \mathcal{P}(\mathcal{L})$ , and  $\mathcal{R}_\mathcal{L}^\alpha$  is well-defined.  $\square$   
 1222

## 1223 C MAX TREE-SLICED WASSERSTEIN DISTANCE ON SYSTEMS OF LINES.

1224 **Max Sliced Wasserstein distance.** Max Sliced Wasserstein (MaxSW) distance (Deshpande et al.,  
 1225 2019) uses only one maximal projecting direction instead of multiple projecting directions as SW.  
 1226

$$1227 \text{MaxSW}_p(\mu, \nu) := \max_{\theta \in \mathcal{U}(\mathbb{S}^{d-1})} \left[ \mathbb{W}_p(\mathcal{R}f_\mu(\cdot, \theta), \mathcal{R}f_\nu(\cdot, \theta)) \right], \quad (38)$$

1228 MaxSW requires an optimization procedure to find the projecting direction. It is a metric on space  
 1229 of probability distributions on  $\mathbb{R}^d$ .  
 1230

1231 We define the Max Tree-Sliced Wasserstein distance on System of Lines (MaxTSW-SL) as follows.  
 1232

1233 **Definition C.1** (Max Tree-Sliced Wasserstein Distance on Systems of Lines). The *Max Tree-Sliced*  
 1234 *Wasserstein Distance on Systems of Lines* between two probability distributions  $\mu, \nu$  in  $\mathcal{P}(\mathbb{R}^d)$  is  
 1235 defined by:  
 1236

$$1237 \text{MaxTSW-SL}(\mu, \nu) := \max_{\mathcal{L} \in \mathbb{T}} \left[ \mathbb{W}_{d_\mathcal{L}, 1}(\mathcal{R}_\mathcal{L}^\alpha \mu, \mathcal{R}_\mathcal{L}^\alpha \nu) \right], \quad (39)$$

1242 MaxTSW-SL is a metric on  $\mathcal{P}(\mathbb{R}^d)$ . The proof of the below theorem is in Appendix D.2.

1243 **Theorem C.2.** MaxTSW-SL distance is a metric on  $\mathcal{P}(\mathbb{R}^d)$ .

1244 We provide an algorithm to compute the MaxTSW-SL in Algorithm 3.

1245

1246 **Algorithm 3** Max Tree-Sliced Wasserstein distance on Systems of lines.

---

1247 **Input:** Probability measures  $\mu$  and  $\nu$ , the number of lines in tree system  $k$ , a splitting function  
1248  $\alpha: \mathbb{R}^d \rightarrow \Delta_{k-1}$ , learning rate  $\eta$ , max number of iterations  $T$ .  
1249 Initialize  $x_1 \in \mathbb{R}^d, t_2, \dots, t_k \in \mathbb{R}, \theta_1, \dots, \theta_k \in \mathbb{S}^{d-1}$ .  
1250 Compute  $\mathcal{L}$  corresponded to  $(x_1, t_2, \dots, t_k, \theta_1, \dots, \theta_k)$ .  
1251 **while**  $\mathcal{L}$  not converge or reach  $T$  **do**  
1252  $x_1 = x_1 + \eta \cdot \nabla_{x_1} \mathbb{W}_{d_{\mathcal{L}},1}(\mathcal{R}_{\mathcal{L}}^\alpha \mu, \mathcal{R}_{\mathcal{L}}^\alpha \nu)$ .  
1253 **for**  $i = 2$  to  $k$  **do**  
1254  $t_i = T_i + \eta \cdot \nabla_{t_i} \mathbb{W}_{d_{\mathcal{L}},1}(\mathcal{R}_{\mathcal{L}}^\alpha \mu, \mathcal{R}_{\mathcal{L}}^\alpha \nu)$ .  
1255 **end for**  
1256 **for**  $i = 1$  to  $k$  **do**  
1257  $\theta_i = \theta_i + \eta \cdot \nabla_{\theta_i} \mathbb{W}_{d_{\mathcal{L}},1}(\mathcal{R}_{\mathcal{L}}^\alpha \mu, \mathcal{R}_{\mathcal{L}}^\alpha \nu)$ .  
1258 Normalize  $\theta_i = \theta_i / \|\theta_i\|_2$ .  
1259 **end for**  
1260 **end while**  
1261 Compute  $\mathcal{L}$  corresponded to  $(x_1, t_2, \dots, t_k, \theta_1, \dots, \theta_k)$ .  
1262 Compute  $\mathbb{W}_{d_{\mathcal{L}},1}(\mathcal{R}_{\mathcal{L}}^\alpha \mu, \mathcal{R}_{\mathcal{L}}^\alpha \nu)$ .  
1263 **Return:**  $\mathcal{L}, \mathbb{W}_{d_{\mathcal{L}},1}(\mathcal{R}_{\mathcal{L}}^\alpha \mu, \mathcal{R}_{\mathcal{L}}^\alpha \nu)$ .

---

1264

## 1265 D THEORETICAL PROOF FOR INJECTIVITY OF TSW-SL

1266

1267 We will leave out “almost-surely-conditions” in the proofs, as they are straightforward to verify, and  
1268 including them would unnecessarily complicate the proofs.

1269

### 1270 D.1 PROOF OF THEOREM 5.2

1271

1272 *Proof.* Need to show that:

1273

$$1274 \text{TSW-SL}(\mu, \nu) := \int_{\mathbb{T}} \mathbb{W}_{d_{\mathcal{L}},1}(\mathcal{R}_{\mathcal{L}}^\alpha \mu, \mathcal{R}_{\mathcal{L}}^\alpha \nu) d\sigma(\mathcal{L}). \quad (40)$$

1275

1276 is a metric on  $\mathcal{P}(\mathbb{R}^d)$ .

1277

1278 **Positive definiteness.** For  $\mu, \nu \in \mathcal{P}(\mathbb{R}^n)$ , one has  $\text{TSW-SL}(\mu, \mu) = 0$  and  $\text{TSW-SL}(\mu, \nu) \geq 0$ . If  
1279  $\text{TSW-SL}(\mu, \nu) = 0$ , then  $\mathbb{W}_{d_{\mathcal{L}},1}(\mathcal{R}_{\mathcal{L}}^\alpha \mu, \mathcal{R}_{\mathcal{L}}^\alpha \nu) = 0$  for all  $\mathcal{L} \in \mathbb{T}$ . Since  $\mathbb{W}_{d_{\mathcal{L}},1}$  is a metric on  $\mathcal{P}(\mathcal{L})$ ,  
1280 we have  $\mathcal{R}_{\mathcal{L}}^\alpha \mu = \mathcal{R}_{\mathcal{L}}^\alpha \nu$  for all  $\mathcal{L} \in \mathbb{T}$ . Since  $\mathbb{T}$  is a subset of  $\mathbb{L}_k^d$  that satisfies the condition in the  
1281 remark at the end of the proof of Theorem B.4, we conclude that  $\mu = \nu$ .

1282

1283 **Symmetry.** For  $\mu, \nu \in \mathcal{P}(\mathbb{R}^n)$ , we have:

1284

$$1285 \text{TSW-SL}(\mu, \nu) = \int_{\mathbb{T}} \mathbb{W}_{d_{\mathcal{L}},1}(\mathcal{R}_{\mathcal{L}}^\alpha \mu, \mathcal{R}_{\mathcal{L}}^\alpha \nu) d\sigma(\mathcal{L}) \quad (41)$$

1286

$$1287 = \int_{\mathbb{T}} \mathbb{W}_{d_{\mathcal{L}},1}(\mathcal{R}_{\mathcal{L}}^\alpha \nu, \mathcal{R}_{\mathcal{L}}^\alpha \mu) d\sigma(\mathcal{L}) \quad (42)$$

1288

$$1289 = \text{TSW-SL}(\nu, \mu). \quad (43)$$

1290

1291 So  $\text{TSW-SL}(\mu, \nu) = \text{TSW-SL}(\nu, \mu)$ .

1296 **Triangle inequality.** For  $\mu_1, \mu_2, \mu_3 \in \mathcal{P}(\mathbb{R}^n)$ , we have:

$$1297 \text{TSW-SL}(\mu_1, \mu_2) + \text{TSW-SL}(\mu_2, \mu_3) \quad (44)$$

$$1299 = \int_{\mathbb{T}} \mathbf{W}_{d_{\mathcal{L}},1}(\mathcal{R}_{\mathcal{L}}^{\alpha}\mu_1, \mathcal{R}_{\mathcal{L}}^{\alpha}\mu_2) d\sigma(\mathcal{L}) + \int_{\mathbb{T}} \mathbf{W}_{d_{\mathcal{L}},1}(\mathcal{R}_{\mathcal{L}}^{\alpha}\mu_2, \mathcal{R}_{\mathcal{L}}^{\alpha}\mu_3) d\sigma(\mathcal{L}) \quad (45)$$

$$1301 = \int_{\mathbb{T}} \left( \mathbf{W}_{d_{\mathcal{L}},1}(\mathcal{R}_{\mathcal{L}}^{\alpha}\mu_1, \mathcal{R}_{\mathcal{L}}^{\alpha}\mu_2) + \mathbf{W}_{d_{\mathcal{L}},1}(\mathcal{R}_{\mathcal{L}}^{\alpha}\mu_2, \mathcal{R}_{\mathcal{L}}^{\alpha}\mu_3) \right) d\sigma(\mathcal{L}) \quad (46)$$

$$1304 \geq \int_{\mathbb{T}} \mathbf{W}_{d_{\mathcal{L}},1}(\mathcal{R}_{\mathcal{L}}^{\alpha}\mu_1, \mathcal{R}_{\mathcal{L}}^{\alpha}\mu_3) d\sigma(\mathcal{L}) \quad (47)$$

$$1306 = \text{TSW-SL}(\mu_1, \mu_3). \quad (48)$$

1307 The triangle inequality holds for TSW-SL. We conclude that TSW-SL is a metric on  $\mathcal{P}(\mathbb{R}^d)$ .  $\square$

## 1309 D.2 PROOF OF THEOREM C.2

1310 *Proof.* Need to show that:

$$1312 \text{MaxTSW-SL}(\mu, \nu) = \max_{\mathcal{L} \in \mathbb{T}} \left[ \mathbf{W}_{d_{\mathcal{L}},1}(\mathcal{R}_{\mathcal{L}}^{\alpha}\mu, \mathcal{R}_{\mathcal{L}}^{\alpha}\nu) \right] \quad (49)$$

1314 is a metric on  $\mathcal{P}(\mathbb{R}^d)$ .

1316 **Positive definiteness.** For  $\mu, \nu \in \mathcal{P}(\mathbb{R}^n)$ , one has  $\text{MaxTSW-SL}(\mu, \mu) = 0$  and  
 1317  $\text{MaxTSW-SL}(\mu, \nu) \geq 0$ . If  $\text{MaxTSW-SL}(\mu, \nu) = 0$ , then  $\mathbf{W}_{d_{\mathcal{L}},1}(\mathcal{R}_{\mathcal{L}}^{\alpha}\mu, \mathcal{R}_{\mathcal{L}}^{\alpha}\nu) = 0$  for all  $\mathcal{L} \in \mathbb{T}$ .  
 1318 Since  $\mathbf{W}_{d_{\mathcal{L}},p}$  is a metric, we have  $\mathcal{R}_{\mathcal{L}}^{\alpha}\mu = \mathcal{R}_{\mathcal{L}}^{\alpha}\nu$  for all  $\mathcal{L} \in \mathbb{T}$ . Since  $\mathbb{T}$  is a subset of  $\mathbb{L}_k^d$  that satisfies  
 1319 the condition in the remark at the end of the proof of Theorem B.4, we conclude that  $\mu = \nu$ .

1321 **Symmetry.** For  $\mu, \nu \in \mathcal{P}(\mathbb{R}^n)$ , we have:

$$1322 \text{MaxTSW-SL}(\mu, \nu) = \max_{\mathcal{L} \in \mathbb{T}} \left[ \mathbf{W}_{d_{\mathcal{L}},1}(\mathcal{R}_{\mathcal{L}}^{\alpha}\mu, \mathcal{R}_{\mathcal{L}}^{\alpha}\nu) \right] \quad (50)$$

$$1324 = \max_{\mathcal{L} \in \mathbb{T}} \left[ \mathbf{W}_{d_{\mathcal{L}},1}(\mathcal{R}_{\mathcal{L}}^{\alpha}\nu, \mathcal{R}_{\mathcal{L}}^{\alpha}\mu) \right] \quad (51)$$

$$1326 = \text{MaxTSW-SL}(\nu, \mu). \quad (52)$$

1327 So  $\text{MaxTSW-SL}(\mu, \nu) = \text{MaxTSW-SL}(\nu, \mu)$ .

1329 **Triangle inequality.** For  $\mu_1, \mu_2, \mu_3 \in \mathcal{P}(\mathbb{R}^n)$ , we have:

$$1331 \text{MaxTSW-SL}(\mu_1, \mu_2) + \text{TSW-SL}(\mu_2, \mu_3) \quad (53)$$

$$1332 = \max_{\mathcal{L} \in \mathbb{T}} \left[ \mathbf{W}_{d_{\mathcal{L}},1}(\mathcal{R}_{\mathcal{L}}^{\alpha}\mu_1, \mathcal{R}_{\mathcal{L}}^{\alpha}\mu_2) \right] + \max_{\mathcal{L}' \in \mathbb{T}} \left[ \mathbf{W}_{d_{\mathcal{L}'},1}(\mathcal{R}_{\mathcal{L}'}^{\alpha}\mu_2, \mathcal{R}_{\mathcal{L}'}^{\alpha}\mu_3) \right] \quad (54)$$

$$1334 \geq \max_{\mathcal{L} \in \mathbb{T}} \left[ \mathbf{W}_{d_{\mathcal{L}},1}(\mathcal{R}_{\mathcal{L}}^{\alpha}\mu_1, \mathcal{R}_{\mathcal{L}}^{\alpha}\mu_2) + \mathbf{W}_{d_{\mathcal{L}},1}(\mathcal{R}_{\mathcal{L}}^{\alpha}\mu_2, \mathcal{R}_{\mathcal{L}}^{\alpha}\mu_3) \right] \quad (55)$$

$$1336 \geq \max_{\mathcal{L} \in \mathbb{T}} \left[ \mathbf{W}_{d_{\mathcal{L}},1}(\mathcal{R}_{\mathcal{L}}^{\alpha}\mu_1, \mathcal{R}_{\mathcal{L}}^{\alpha}\mu_3) \right] \quad (56)$$

$$1338 = \text{MaxTSW-SL}(\mu_1, \mu_3). \quad (57)$$

1339 The triangle inequality holds for MaxTSW-SL. We conclude that MaxTSW-SL is a metric on  
 1340  $\mathcal{P}(\mathbb{R}^d)$ .  $\square$

## 1342 E EXPERIMENTAL DETAILS

### 1344 E.1 GRADIENT FLOWS

1346 Gradient flow is a concept in differential geometry and dynamical systems that describes the evolu-  
 1347 tion of a point or a curve under a given vector field. In the field of Sliced Wasserstein distance, this is  
 1348 a synthetic task that is used to evaluate the evolution of Wasserstein distance between 2 distributions  
 1349 (source and target distributions) while minimizing different distances (Mahey et al., 2023; Kolouri  
 et al., 2019) as a loss function.

1350  
 1351  
 1352  
 1353  
 1354  
 1355  
 1356  
 1357  
 1358  
 1359  
 1360  
 1361  
 1362  
 1363  
 1364  
 1365  
 1366  
 1367  
 1368  
 1369  
 1370  
 1371  
 1372  
 1373  
 1374  
 1375  
 1376  
 1377  
 1378  
 1379  
 1380  
 1381  
 1382  
 1383  
 1384  
 1385  
 1386  
 1387  
 1388  
 1389  
 1390  
 1391  
 1392  
 1393  
 1394  
 1395  
 1396  
 1397  
 1398  
 1399  
 1400  
 1401  
 1402  
 1403

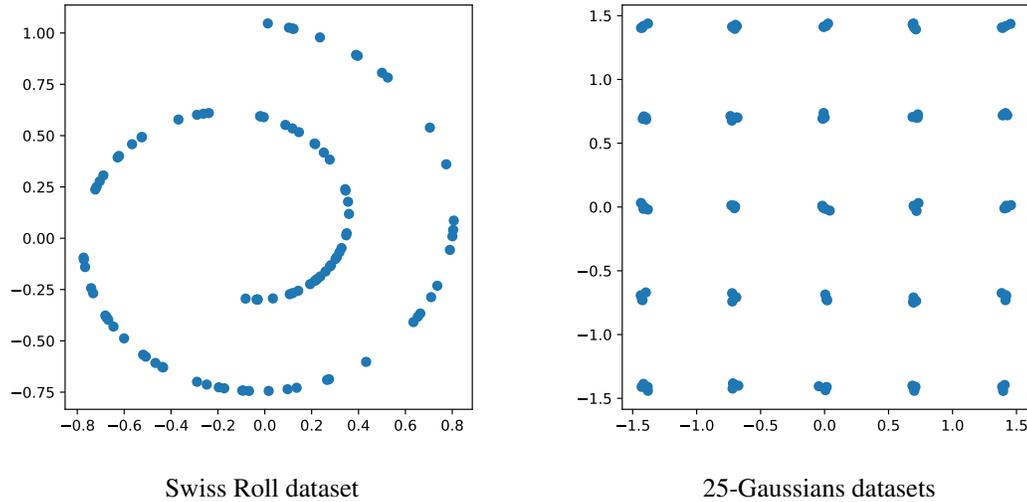


Figure 6: Swiss Roll and 25-Gaussians datasets for Gradient Flows task

**Datasets.** We use Swiss Roll, 25-Gaussians and high-dimensional Gaussian datasets as the target distribution as in (Kolouri et al., 2019). The details of these datasets can be described as follows.

- The **Swiss Roll dataset** is a popular synthetic dataset used in machine learning, particularly for visualizing and testing dimensionality reduction techniques. It is generated using the `make_swiss_roll` function of Pytorch, which creates a non-linear, three-dimensional dataset resembling a swiss roll or spiral shape. In the original version, it is a three-dimensional dataset with each dimension representing a coordinate in the 1 axis of the points. In order to simplify it, we follow (Kolouri et al., 2019) to just consider the two-dimensional Swiss Roll dataset by reducing the second coordinates and retaining only the first and the third coordinates. We set the number of samples to equal 100.
- The **25-Gaussians dataset** is obtained by first create a grid of 25 points spaced evenly in a  $5 \times 5$  arrangement. For each grid point, we generate a cluster by sampling points from a Gaussian distribution centered at that grid point, with a small standard deviation. All the points from the 25 clusters are combined, shuffled randomly, and scaled to form the final dataset.
- The **High-dimensional Gaussian datasets** are generated by initializing a mean vector,  $\mu_s$ , consisting of ones across all dimensions. Each element of this mean vector is scaled by a random value to introduce variability. The covariance matrix,  $\Sigma$ , is created as an identity matrix scaled by a constant, ensuring independence among dimensions. Points are then sampled from the multivariate normal distribution using these parameters, resulting in a dataset of  $N$  points in the specified high-dimensional space.

The Swiss Roll and 25-Gaussians datasets are presented in Figure 6.

**Hyperparameters.** For TSW-SL, we use  $L = 25$ ,  $k = 4$  in all experiments, while  $L = 100$  is set for SW and SW-variants, with 100 points generated per distribution across datasets. Following (Mahey et al., 2023), the global learning rate for all baselines is  $5 \times 10^{-3}$ . For our methods, we use  $5 \times 10^{-3}$  for the 25-Gaussians and Swiss Roll datasets, and  $5 \times 10^{-2}$  for the high-dimensional Gaussian datasets. We also follow (Mahey et al., 2023) in setting 100 iterations for both MaxSW and MaxTSW-SL, using a learning rate of  $1 \times 10^{-4}$  for both methods.

**Evaluation metrics.** We use the Wasserstein distance as a neutral metric to evaluate how close the model distribution  $\mu(t)$  is to the target distribution  $\nu$ . Over 2500 timesteps, we evaluate the

1404 Wasserstein distance between source and target distributions at iteration 500, 1000, 1500, 2000 and  
 1405 2500.

1406 We utilize the source code adapted from [Mahey et al. \(2023\)](#) for this task.

## 1408 E.2 COLOR TRANSFER

1409 This section extends our experiments to evaluate our methods against various baselines as discussed  
 1410 in [Nguyen et al. \(2024a\)](#). Similar to E.1, we set  $L = 100$  for all baselines and employ 25 trees and  
 1411 4 lines for our TSW-SL.

1412 **Settings.** Given a source image and a target image, we represent their respective color palettes as  
 1413 matrices  $X$  and  $Y$ , each with dimensions  $n \times 3$  (where  $n$  denotes the number of pixels).

1414 We follow [Nguyen et al. \(2024a\)](#) to first define the curve  $\dot{Z}(t) = -n\nabla_{Z(t)} [\text{SW}_2(P_{Z(t)}, P_Y)]$   
 1415 where  $P_X$  and  $P_Y$  are empirical distributions over  $X$  and  $Y$  in turn. Here, the curve starts from  
 1416  $Z(0) = X$  and ends at  $Y$ .

1417 We then reduce the number of colors in the images to 1000 using K-means clustering. After that,  
 1418 we iterate through the curve between the empirical distribution of colors in the source image  $P_X$   
 1419 and the empirical distribution of colors in the target image  $P_Y$  using the approximate Euler method.  
 1420 However, owing to the color palette values (RGB) lying within the set  $\{0, \dots, 255\}$ , an additional  
 1421 rounding step is necessary during the final Euler iterations. We increase the number of iterations  
 1422 to 2000 and utilize a step size of 1 as in ([Nguyen et al., 2024a](#)) for baselines and a step size of 16  
 1423 for our experiments. We use  $L = 100$  in SW variants and  $L = 25, k = 4$  in TSW-SL for a fair  
 1424 comparison.

1425 **Evaluation metrics.** We present the Wasserstein distances at the final time step alongside the  
 1426 corresponding transferred images to evaluate the performance of different methods. The results  
 1427 illustrated in Figure 7 demonstrate that our novel metrics substantially reduce the Wasserstein dis-  
 1428 tance of a large number of baselines. Our primary contribution is the development of a metric that  
 1429 effectively bridges SW and TSW, exhibiting superior performance over vanilla SW, MaxSW, and  
 1430 several enhanced variants of SW. This represents a significant breakthrough in the field of optimal  
 1431 transport and paves the way for further advancements.

1432 We utilize the source code adapted from [Nguyen et al. \(2021\)](#) for this task.

1433 **Additional Results.** We further provide in Figure 8 to show that our TSW-SL and MaxTSW-SL  
 1434 improve the performance of original SW and MaxSW both qualitatively and quantitatively.

## 1437 E.3 GENERATIVE ADVERSARIAL NETWORK

1438 **Architectures.** We denote  $\mu$  as our data distribution. Subsequently, we formulate the model distri-  
 1439 bution  $\nu_\phi$  as a resultant probability measure generated by applying a neural network  $G_\phi$  to transform  
 1440 a unit multivariate Gaussian ( $\epsilon$ ) into the data space. Additionally, we employ another neural network  
 1441  $T_\beta$  to map from the data space to a singular scalar value. More specifically,  $T_{\beta_1}$  represents the sub-  
 1442 set neural network of  $T_\beta$  that maps from the data space to a feature space, specifically the output of  
 1443 the last ResNet block, while  $T_{\beta_2}$  maps from the feature space (the image of  $T_{\beta_1}$ ) to a single scalar.  
 1444 Formally,  $T_\beta = T_{\beta_2} \circ T_{\beta_1}$ . We utilize the subsequent neural network architectures for  $G_\phi$  and  $T_\beta$ :

- 1445 • **CIFAR10:**

1446 -  $G_\phi : z \in \mathbb{R}^{128} (\sim \epsilon : \mathcal{N}(0, 1)) \rightarrow 4 \times 4 \times 256$  (Dense, Linear)  $\rightarrow$  ResBlock up 256  $\rightarrow$   
 1447 ResBlock up 256  $\rightarrow$  ResBlock up 256  $\rightarrow$  BN, ReLU,  $\rightarrow 3 \times 3$  conv, 3 Tanh .

1448 -  $T_{\beta_1} : x \in [-1, 1]^{32 \times 32 \times 3} \rightarrow$  ResBlock down 128  $\rightarrow$  ResBlock down 128  $\rightarrow$  ResBlock  
 1449 down 128  $\rightarrow$  ResBlock 128  $\rightarrow$  ResBlock 128.

1450 -  $T_{\beta_2} : x \in \mathbb{R}^{128 \times 8 \times 8} \rightarrow$  ReLU  $\rightarrow$  Global sum pooling (128)  $\rightarrow 1$  (Spectral  
 1451 normalization).

1452 -  $T_\beta(x) = T_{\beta_2}(T_{\beta_1}(x))$ .

- 1453 • **CelebA:**

1458  
 1459  
 1460  
 1461  
 1462  
 1463  
 1464  
 1465  
 1466  
 1467  
 1468  
 1469  
 1470  
 1471  
 1472  
 1473  
 1474  
 1475  
 1476  
 1477  
 1478  
 1479  
 1480  
 1481  
 1482  
 1483  
 1484  
 1485  
 1486  
 1487  
 1488  
 1489  
 1490  
 1491  
 1492  
 1493  
 1494  
 1495  
 1496  
 1497  
 1498  
 1499  
 1500  
 1501  
 1502  
 1503  
 1504  
 1505  
 1506  
 1507  
 1508  
 1509  
 1510  
 1511

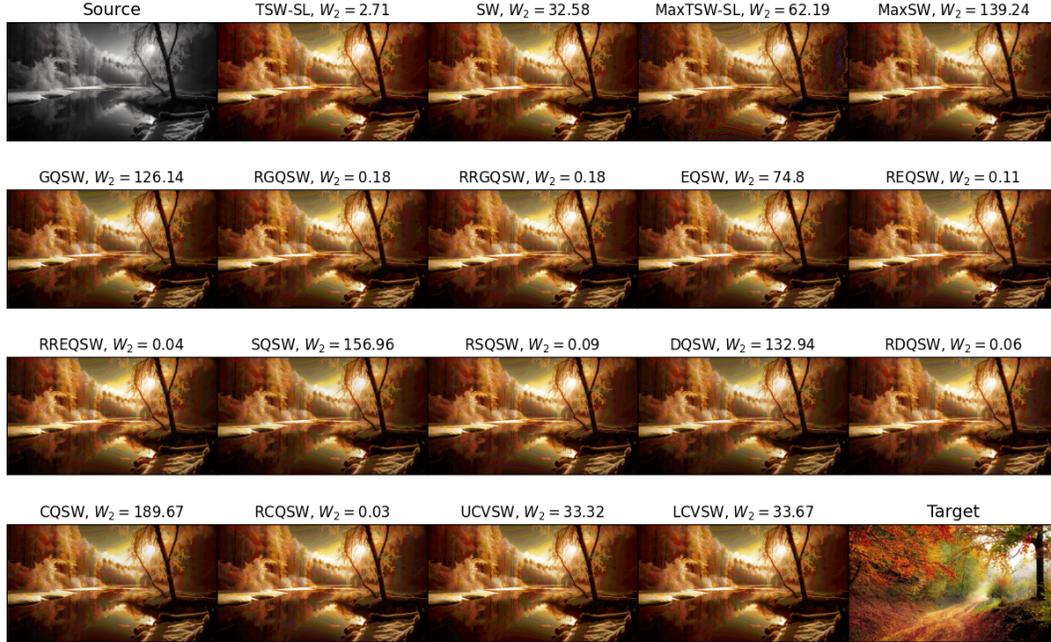


Figure 7: Style-transferred images from various baselines with 100 projecting directions.

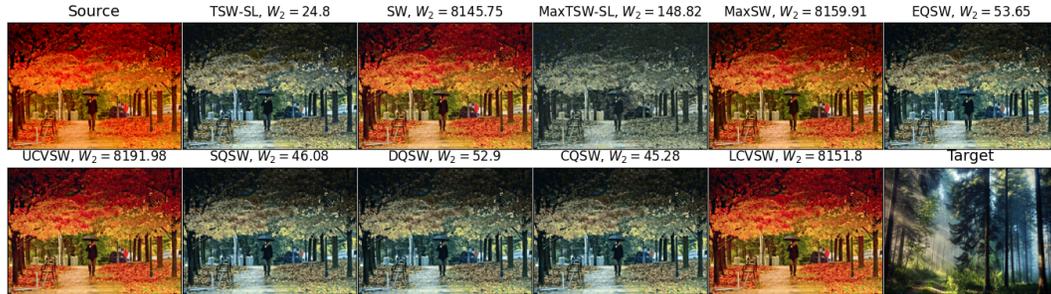


Figure 8: Additional style-transferred images from various baselines with 100 projecting directions.

-  $G_\phi : z \in \mathbb{R}^{128} (\sim \epsilon : \mathcal{N}(0, 1)) \rightarrow 4 \times 4 \times 256$  (Dense, Linear)  $\rightarrow$  ResBlock up 256  $\rightarrow$  ResBlock up 256  $\rightarrow$  ResBlock up 256  $\rightarrow$  ResBlock up 256  $\rightarrow$  BN, ReLU,  $\rightarrow$   $3 \times 3$  conv, 3 Tanh .

-  $T_{\beta_1} : x \in [-1, 1]^{32 \times 32 \times 3} \rightarrow$  ResBlock down 128  $\rightarrow$  ResBlock down 128  $\rightarrow$  ResBlock down 128  $\rightarrow$  ResBlock down 128  $\rightarrow$  ResBlock 128  $\rightarrow$  ResBlock 128.

-  $T_{\beta_2} : x \in \mathbb{R}^{128 \times 8 \times 8} \rightarrow$  ReLU  $\rightarrow$  Global sum pooling(128)  $\rightarrow$  1( Spectral normalization ).

-  $T_\beta(x) = T_{\beta_2}(T_{\beta_1}(x))$ .

• **STL-10:**

-  $G_\phi : z \in \mathbb{R}^{128} (\sim \epsilon : \mathcal{N}(0, 1)) \rightarrow 3 \times 3 \times 256$  (Dense, Linear)  $\rightarrow$  ResBlock up 256  $\rightarrow$  ResBlock up 256  $\rightarrow$  ResBlock up 256  $\rightarrow$  ResBlock up 256  $\rightarrow$  BN, ReLU,  $\rightarrow$   $3 \times 3$  conv, 3 Tanh .

-  $T_{\beta_1} : x \in [-1, 1]^{32 \times 32 \times 3} \rightarrow$  ResBlock down 128  $\rightarrow$  ResBlock down 128  $\rightarrow$  ResBlock down 128  $\rightarrow$  ResBlock down 128  $\rightarrow$  ResBlock 128  $\rightarrow$  ResBlock 128.

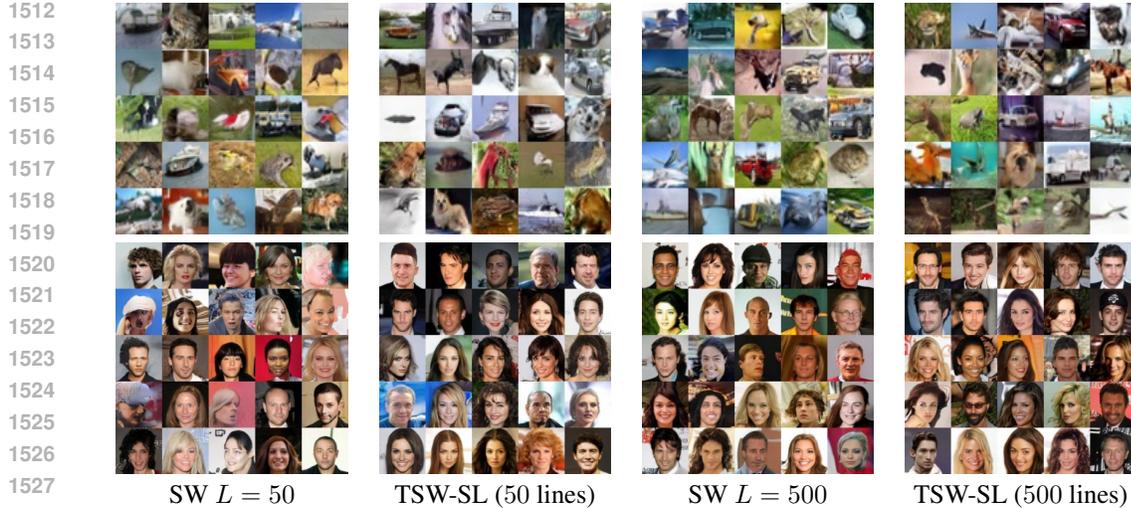


Figure 9: Randomly generated images of different methods on CIFAR10 and CelebA of SN-GAN

$-T_{\beta_2} : x \in \mathbb{R}^{128 \times 8 \times 8} \rightarrow \text{ReLU} \rightarrow \text{Global sum pooling}(128) \rightarrow 1$  (Spectral normalization).

$-T_{\beta}(x) = T_{\beta_2}(T_{\beta_1}(x))$ .

We use the following bi-optimization problem to train our neural networks:

$$\min_{\beta_1, \beta_2} (\mathbb{E}_{x \sim \mu} [\min(0, -1 + T_{\beta}(x))] + \mathbb{E}_{z \sim \epsilon} [\min(0, -1 - T_{\beta}(G_{\phi}(z)))]),$$

$$\min_{\phi} \mathbb{E}_{X \sim \mu^{\otimes m}, Z \sim \epsilon^{\otimes m}} \left[ \mathcal{S} \left( \tilde{T}_{\beta_1, \beta_2} \# P_X, \tilde{T}_{\beta_1, \beta_2} \# G_{\phi} \# P_Z \right) \right],$$

where the function  $\tilde{T}_{\beta_1, \beta_2} = [T_{\beta_1}(x), T_{\beta_2}(T_{\beta_1}(x))]$  which is the concatenation vector of  $T_{\beta_1}(x)$  and  $T_{\beta_2}(T_{\beta_1}(x))$ ,  $\mathcal{S}$  is an estimator of the sliced Wasserstein distance.

**Training setup.** In our experiments, we configured the number of training iterations to 100000 for CIFAR10, STL-10 and 50000 for CelebA. The generator  $G_{\phi}$  is updated every 5 iteration, while the feature function  $T_{\beta}$  undergoes an update each iteration. Across all datasets, we maintain a consistent mini-batch size of 128. We leverage the Adam optimizer (Kingma, 2014) with parameters  $(\beta_1, \beta_2) = (0, 0.9)$  for both  $G_{\phi}$  and  $T_{\beta}$  with the learning rate 0.0002. Furthermore, we use 50000 random samples generated from the generator to compute the FID and Inception scores. For FID score evaluation, the statistics of datasets are computed using all training samples.

**Results.** For qualitative analysis, Figure 9 displays a selection of randomly generated images produced by the trained models. It is evident that utilizing our TSW-SL as the generator loss significantly enhances the quality of the generated images. Additionally, increasing the number of projections further improves the visual quality of images across all estimators. This improvement in visual quality is corroborated by the FID and IS scores presented in Table 3.

We utilize the source code adapted from (Miyato et al., 2018) for this task.

**Additional results.** To fully show the empirical advantage of our methods, we conducted additional experiments on Adversarial Neural Networks on the CIFAR-10 dataset and STL-10 dataset. First of all, Table 6 presents the average FID and IS scores for different methods on the CIFAR-10 dataset. For 50 projecting directions, our TSW-SL method with 10 trees and 5 lines each ( $L = 10$ ,  $k = 5$ ) achieves the best performance, outperforming the standard SW method. Similarly, for 500 projecting directions, TSW-SL ( $L = 100$ ,  $k = 5$ ) shows superior results compared to SW. This demonstrates the consistent effectiveness of our approach across different numbers of projecting directions. Additionally, Table 5 illustrates the performance of generative models on the STL-10 ( $96 \times 96$ ) dataset with different numbers of trees and lines compared with SW and orthogonal-SW.

Table 5: Performance of different methods on STL-10 dataset on SN-GAN architecture

Methods	Total line	No. of lines per tree	No. of trees	FID	IS
SW	50	-	-	69.46	9.08
Orthogonal-SW	50	-	-	63.61	9.63
TSW-SL (ours)	51	3	17	65.93	9.75
TSW-SL (ours)	52	4	13	<u>62.91</u>	<u>9.95</u>
TSW-SL (ours)	50	5	10	<b>61.15</b>	<b>10.00</b>

In our experiments, we utilize the SN-GAN architecture (Miyato et al., 2018) for STL-10. For SW and Orthogonal-SW, we conduct experiments using 50 projecting directions. Our TSW-SL method is tested with three distinct configurations: 10 trees with 5 lines each, 13 trees with 4 lines each, and 17 trees with 3 lines each. All hyperparameters remain consistent with those used in our main paper. To evaluate the models, we generate 50000 random images.

Table 6: Average FID and IS score of 3 runs on CIFAR-10 of SN-GAN

	50 projecting directions		500 projecting directions		
	FID(↓)	IS(↑)	FID(↓)	IS(↑)	
SW ( $L = 50$ )	16.80 ± 0.45	7.97 ± 0.05	SW ( $L = 500$ )	14.23 ± 0.84	8.25 ± 0.05
TSW-SL ( $L = 10, k = 5$ )	<b>15.44 ± 0.42</b>	<b>8.14 ± 0.05</b>	TSW-SL ( $L = 100, k = 5$ )	<b>13.27 ± 0.23</b>	<b>8.30 ± 0.01</b>
TSW-SL ( $L = 17, k = 3$ )	15.9 ± 0.35	8.10 ± 0.04	TSW-SL ( $L = 167, k = 3$ )	14.18 ± 0.38	8.28 ± 0.07

#### E.4 DENOISING DIFFUSION MODELS

In this section, we provide details about denoising diffusion models, a class of generative models that have shown remarkable success in producing high-quality samples. We first describe the forward and reverse processes that form the foundation of these models. Then, we introduce the concept of denoising diffusion GANs, which aims to accelerate the generation process. Finally, we explain how our proposed TSW-SL distance can be integrated into this framework.

The process in diffusion models is typically divided into two main parts: the forward process and the reverse process.

The forward process is defined as:

$$q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1}), \quad q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t I)$$

where the variance schedule  $\beta_1, \dots, \beta_T$  can be constant or learned hyperparameters. The reverse process is defined as:

$$p_\theta(x_{0:T}) = p(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t), \quad p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t)),$$

where  $\mu_\theta(x_t, t)$  and  $\Sigma_\theta(x_t, t)$  are functions that provide the mean and covariance for the Gaussian and are defined using MLPs.

The model is trained by maximizing the variational lower bound on the negative log-likelihood:

$$\mathbb{E}_q[-\log p_\theta(x_0)] \leq \mathbb{E}_q \left[ -\log \frac{p_\theta(x_{0:T})}{q(x_{1:T}|x_0)} \right] = L,$$

While traditional models have successfully generated high-quality images without the need for adversarial training. However, their sampling process involves simulating a Markov chain for multiple

steps, which can be time-consuming. To accelerate the generation process by reducing the number of steps  $T$ , denoising diffusion GANs (Xiao et al., 2021) propose utilizing an implicit denoising model:

$$p_{\theta}(x_{t-1}|x_t) = \int p_{\theta}(x_{t-1}|x_t, \epsilon) G_{\theta}(x_t, \epsilon) d\epsilon, \quad \text{with } \epsilon \sim \mathcal{N}(0, I).$$

Subsequently, adversarial training is employed (Xiao et al., 2021) to optimize the model parameters

$$\min_{\phi} \sum_{t=1}^T \mathbb{E}_{q(x_t)} [D_{adv}(q(x_{t-1}|x_t) || p_{\phi}(x_{t-1}|x_t))],$$

where  $D_{adv}$  refers to either the GAN objective or the Jensen-Shannon divergence (Goodfellow et al., 2020). We follow the proposed Augmented Generalized Mini-batch Energy distances of Nguyen et al. (2024b) leverage our TSW-SL distance for  $D_{adv}$ .

More specifically, as described by Nguyen et al. (2024b), the adversarial loss is replaced by the augmented generalized Mini-batch Energy (AGME) distance. For two distributions  $\mu$  and  $\nu$ , with a mini-batch size  $n \geq 1$ , the AGME distance using a Sliced Wasserstein (SW) kernel is defined as:

$$AGME_b^2(\mu, \nu; g) = GME_b^2(\tilde{\mu}, \tilde{\nu}),$$

where  $\tilde{\mu} = f_{\#}\mu$  and  $\tilde{\nu} = f_{\#}\nu$ , with the mapping  $f(x) = (x, g(x))$  for a nonlinear function  $g : \mathbb{R}^d \rightarrow \mathbb{R}$ . The  $GME$  is the generalized Mini-batch Energy distance Salimans et al. (2018), given by:

$$GME_b^2(\mu, \nu) = 2\mathbb{E}[D(P_X, P_Y)] - \mathbb{E}[D(P_X, P'_X)] - \mathbb{E}[D(P_Y, P'_Y)],$$

where  $X, X' \stackrel{i.i.d.}{\sim} \mu^{\otimes m}$ ,  $Y, Y' \stackrel{i.i.d.}{\sim} \nu^{\otimes m}$ , and

$$P_X = \frac{1}{m} \sum_{i=1}^m \delta_{x_i}, \quad X = (x_1, \dots, x_m).$$

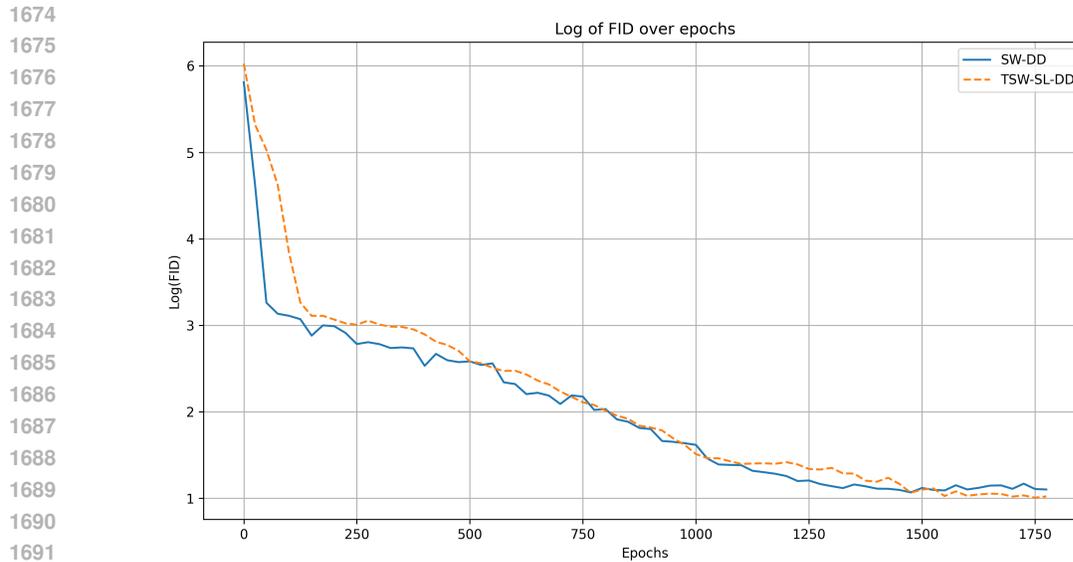
In the equation above,  $D$  denotes a distance function that can calculate the distance between two probability measures. To evaluate how well TSW-SL compares to other SW variants in capturing topological information, particularly when the supports lie in high-dimensional spaces, we replace  $D$  with both TSW and SW variants. We then train the generative model to assess which distance metric better quantifies the divergence between two probability distributions. A lower FID score indicates a more effective distance measure.

**Experimental setup.** For our experiments, we adopted the architecture and hyperparameters from Nguyen et al. (2024b), training our models for 1800 epochs. For TSW, we employed the following hyperparameters:  $L = 2500$ ,  $k = 4$ . For the vanilla SW and its variants, we adhered to the approach outlined in Nguyen et al. (2024b), using  $L = 10000$ . This consistent setup allowed us to effectively compare the performance of our proposed methods against existing approaches while maintaining experimental integrity.

**FID plot.** Figure 10 illustrates the FID scores of SW-DD and TSW-SL-DD across epochs. Due to the wide range of FID values, from over 400 in the initial epoch to less than 3.0 in the final epochs, we present the results on a logarithmic scale for improved visualization. The plot shows that TSW-SL-DD achieves a greater reduction in FID scores compared to SW-DD during the final 300 epochs.

## E.5 COMPUTATIONAL INFRASTRUCTURE

The experiments on gradient flow, color transfer and generative models using generative adversarial networks are conducted on a single NVIDIA A100 GPU. Training generative adversarial networks



1693 Figure 10: FID score over epochs between SW and TSW-SL

1694  
1695 on CIFAR10 requires 14 hours, while CelebA training takes 22 hours. Regarding gradient flows,  
1696 each dataset’s experiments take approximately 3.5 hours. For color transfer, the runtime is 15 min-  
1697 utes.

1698  
1699 The denoising diffusion experiments were conducted parallelly on 2 NVIDIA A100 GPUs and each  
1700 run takes us 81.5 hours.

## 1701 F BROADER IMPACT

1702  
1703  
1704 The novel Tree-Sliced Wasserstein distance on a System of Lines (TSW-SL) introduced in this  
1705 paper holds significant potential for societal advancement. By refining optimal transport method-  
1706 ologies, TSW-SL enhances their accuracy and versatility across diverse practical domains. This  
1707 approach, which synthesizes elements from both Sliced Wasserstein (SW) and Tree-Sliced Wasser-  
1708 stein (TSW), offers enhanced resilience and adaptability, particularly in dynamic scenarios. The  
1709 resulting improvements in gradient flows, color manipulation, and generative modeling yield more  
1710 potent computational tools. These advancements promise to catalyze progress across multiple sec-  
1711 tors. In healthcare, for instance, refined image processing could elevate the precision of medical  
1712 diagnostics. The creative industries stand to benefit from more sophisticated generative models, po-  
1713 tentially revolutionizing artistic expression. Moreover, TSW-SL’s proficiency in handling dynamic  
1714 environments opens new avenues for real-time analytics and decision-making in fields ranging from  
1715 finance to environmental monitoring. By expanding the applicability of advanced computational  
1716 techniques to a wider array of real-world challenges, TSW-SL contributes to technological innova-  
1717 tion and, consequently, to the enhancement of societal welfare.

1718  
1719  
1720  
1721  
1722  
1723  
1724  
1725  
1726  
1727