

## 644 A Proofs

### 645 A.1 Proof of Proposition 1

646 **Optimal  $\alpha$ .** Let  $\mu \in \mathcal{P}_2(S^1)$ ,  $\nu = \text{Unif}(S^1)$ . Since  $\nu$  is the uniform distribution on  $S^1$ , its cdf is  
 647 the identity on  $[0, 1]$  (where we identified  $S^1$  and  $[0, 1]$ ). We can extend the cdf  $F$  on the real line as  
 648 in [89] with the convention  $F(y + 1) = F(y) + 1$ . Therefore,  $F_\nu = \text{Id}$  on  $\mathbb{R}$ . Moreover, we know  
 649 that for all  $x \in S^1$ ,  $(F_\nu - \alpha)^{-1}(x) = F_\nu^{-1}(x + \alpha) = x + \alpha$  and

$$W_2^2(\mu, \nu) = \inf_{\alpha \in \mathbb{R}} \int_0^1 |F_\mu^{-1}(t) - (F_\nu - \alpha)^{-1}(t)|^2 dt. \quad (25)$$

650 For all  $\alpha \in \mathbb{R}$ , let  $f(\alpha) = \int_0^1 (F_\mu^{-1}(t) - (F_\nu - \alpha)^{-1}(t))^2 dt$ . Then, we have:

$$\begin{aligned} \forall \alpha \in \mathbb{R}, f(\alpha) &= \int_0^1 (F_\mu^{-1}(t) - t - \alpha)^2 dt \\ &= \int_0^1 (F_\mu^{-1}(t) - t)^2 dt + \alpha^2 - 2\alpha \int_0^1 (F_\mu^{-1}(t) - t) dt \\ &= \int_0^1 (F_\mu^{-1}(t) - t)^2 dt + \alpha^2 - 2\alpha \left( \int_0^1 x d\mu(x) - \frac{1}{2} \right), \end{aligned} \quad (26)$$

651 where we used that  $(F_\mu^{-1})_\# \text{Unif}([0, 1]) = \mu$ .

652 Hence,  $f'(\alpha) = 0 \iff \alpha = \int_0^1 x d\mu(x) - \frac{1}{2}$ .

653 **Closed-form for empirical distributions.** Let  $(x_i)_{i=1}^n \in [0, 1]^n$  such that  $x_1 < \dots < x_n$  and let  
 654  $\mu_n = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$  a discrete distribution.

655 To compute the closed-form of  $W_2$  between  $\mu_n$  and  $\nu = \text{Unif}(S^1)$ , we first have that the optimal  $\alpha$   
 656 is  $\alpha_n = \frac{1}{n} \sum_{i=1}^n x_i - \frac{1}{2}$ . Moreover, we also have:

$$\begin{aligned} W_2^2(\mu_n, \nu) &= \int_0^1 (F_{\mu_n}^{-1}(t) - (t + \hat{\alpha}_n))^2 dt \\ &= \int_0^1 F_{\mu_n}^{-1}(t)^2 dt - 2 \int_0^1 t F_{\mu_n}^{-1}(t) dt - 2\hat{\alpha}_n \int_0^1 F_{\mu_n}^{-1}(t) dt + \frac{1}{3} + \hat{\alpha}_n + \hat{\alpha}_n^2. \end{aligned} \quad (27)$$

657 Then, by noticing that  $F_{\mu_n}^{-1}(t) = x_i$  for all  $t \in [F(x_i), F(x_{i+1})]$ , we have

$$\int_0^1 t F_{\mu_n}^{-1}(t) dt = \sum_{i=1}^n \int_{\frac{i-1}{n}}^{\frac{i}{n}} t x_i dt = \frac{1}{2n^2} \sum_{i=1}^n x_i (2i - 1), \quad (28)$$

658

$$\int_0^1 F_{\mu_n}^{-1}(t)^2 dt = \frac{1}{n} \sum_{i=1}^n x_i^2, \quad \int_0^1 F_{\mu_n}^{-1}(t) dt = \frac{1}{n} \sum_{i=1}^n x_i, \quad (29)$$

659 and we also have:

$$\hat{\alpha}_n + \hat{\alpha}_n^2 = \frac{1}{n} \sum_{i=1}^n x_i - \frac{1}{2} + \left( \frac{1}{n} \sum_{i=1}^n x_i \right)^2 + \frac{1}{4} - \frac{1}{n} \sum_{i=1}^n x_i = \left( \frac{1}{n} \sum_{i=1}^n x_i \right)^2 - \frac{1}{4}. \quad (30)$$

660 Then, by plugging these results into (27), we obtain

$$\begin{aligned} W_2^2(\mu_n, \nu) &= \frac{1}{n} \sum_{i=1}^n x_i^2 - \frac{1}{n^2} \sum_{i=1}^n (2i - 1) x_i - 2 \left( \frac{1}{n} \sum_{i=1}^n x_i \right)^2 + \frac{1}{n} \sum_{i=1}^n x_i + \frac{1}{3} + \left( \frac{1}{n} \sum_{i=1}^n x_i \right)^2 - \frac{1}{4} \\ &= \frac{1}{n} \sum_{i=1}^n x_i^2 - \left( \frac{1}{n} \sum_{i=1}^n x_i \right)^2 + \frac{1}{n^2} \sum_{i=1}^n (n + 1 - 2i) x_i + \frac{1}{12}. \end{aligned} \quad (31)$$

## 661 A.2 Proof of Equation (17)

662 Let  $U \in \mathbb{V}_{d,2}$ . Then the great circle generated by  $U \in \mathbb{V}_{d,2}$  is defined as the intersection between  
 663  $\text{span}(UU^T)$  and  $S^{d-1}$ . And we have the following characterization:

$$\begin{aligned} x \in \text{span}(UU^T) \cap S^{d-1} &\iff \exists y \in \mathbb{R}^d, x = UU^T y \text{ and } \|x\|_2^2 = 1 \\ &\iff \exists y \in \mathbb{R}^d, x = UU^T y \text{ and } \|UU^T y\|_2^2 = y^T UU^T y = \|U^T y\|_2^2 = 1 \\ &\iff \exists z \in S^1, x = Uz. \end{aligned}$$

664 And we deduce that

$$\forall U \in \mathbb{V}_{d,2}, x \in S^{d-1}, P^U(x) = \underset{z \in S^1}{\operatorname{argmin}} d_{S^{d-1}}(x, Uz). \quad (32)$$

## 665 A.3 Proof of Lemma 1

666 Let  $U \in \mathbb{V}_{d,2}$  and  $x \in S^{d-1}$  such that  $U^T x \neq 0$ . Denote  $U = (u_1 \ u_2)$ , *i.e.* the 2-plane  $E$   
 667 is  $E = \text{span}(UU^T) = \text{span}(u_1, u_2)$  and  $(u_1, u_2)$  is an orthonormal basis of  $E$ . Then, for all  
 668  $x \in S^{d-1}$ , the projection on  $E$  is  $p^E(x) = \langle u_1, x \rangle u_1 + \langle u_2, x \rangle u_2 = UU^T x$ .

669 Now, let us compute the geodesic distance between  $x \in S^{d-1}$  and  $\frac{p^E(x)}{\|p^E(x)\|_2} \in E \cap S^{d-1}$ :

$$d_{S^{d-1}}\left(x, \frac{p^E(x)}{\|p^E(x)\|_2}\right) = \arccos\left(\left\langle x, \frac{p^E(x)}{\|p^E(x)\|_2} \right\rangle\right) = \arccos(\|p^E(x)\|_2), \quad (33)$$

670 using that  $x = p^E(x) + p^{E^\perp}(x)$ .

671 Let  $y \in E \cap S^{d-1}$  another point on the great circle. By the Cauchy-Schwarz inequality, we have

$$\langle x, y \rangle = \langle p^E(x), y \rangle \leq \|p^E(x)\|_2 \|y\|_2 = \|p^E(x)\|_2. \quad (34)$$

672 Therefore, using that  $\arccos$  is decreasing on  $(-1, 1)$ ,

$$d_{S^{d-1}}(x, y) = \arccos(\langle x, y \rangle) \geq \arccos(\|p^E(x)\|_2) = d_{S^{d-1}}\left(x, \frac{p^E(x)}{\|p^E(x)\|_2}\right). \quad (35)$$

673 Moreover, we have equality if and only if  $y = \lambda p^E(x)$ . And since  $y \in S^{d-1}$ ,  $|\lambda| = \frac{1}{\|p^E(x)\|_2}$ . Using  
 674 again that  $\arccos$  is decreasing, we deduce that the minimum is well attained in  $y = \frac{p^E(x)}{\|p^E(x)\|_2} =$   
 675  $\frac{UU^T x}{\|UU^T x\|_2}$ .

676 Finally, using that  $\|UU^T x\|_2 = x^T UU^T UU^T x = x^T UU^T x = \|U^T x\|_2$ , we deduce that

$$P^U(x) = \frac{U^T x}{\|U^T x\|_2}. \quad (36)$$

677 Finally, by noticing that the projection is unique if and only if  $U^T x = 0$ , and using [9, Proposition  
 678 4.2] which states that there is a unique projection for a.e.  $x$ , we deduce that  $\{x \in S^{d-1}, U^T x = 0\}$   
 679 is of measure null and hence, for a.e.  $x \in S^{d-1}$ , we have the result.

#### 680 A.4 Proof of Proposition 2

681 Let  $f \in L^1(S^{d-1})$ ,  $g \in C_0(S^1 \times \mathbb{V}_{d,2})$ , then by Fubini's theorem,

$$\begin{aligned}
\langle \tilde{R}f, g \rangle_{S^1 \times \mathbb{V}_{d,2}} &= \int_{V_{d,2}} \int_{S^1} \tilde{R}f(z, U) g(z, U) \, dz d\sigma(U) \\
&= \int_{V_{d,2}} \int_{S^1} \int_{S^{d-1}} f(x) \mathbb{1}_{\{z=P^U(x)\}} g(z, U) \, dx dz d\sigma(U) \\
&= \int_{S^{d-1}} f(x) \int_{V_{d,2}} \int_{S^1} g(z, U) \mathbb{1}_{\{z=P^U(x)\}} \, dz d\sigma(U) dx \\
&= \int_{S^{d-1}} f(x) \int_{V_{d,2}} g(P^U(x), U) \, d\sigma(U) dx \\
&= \int_{S^{d-1}} f(x) \tilde{R}^* g(x) \, dx \\
&= \langle f, \tilde{R}^* g \rangle_{S^{d-1}}.
\end{aligned} \tag{37}$$

#### 682 A.5 Proof of Proposition 3

683 Let  $g \in C_0(S^1 \times \mathbb{V}_{d,2})$ ,

$$\begin{aligned}
\int_{\mathbb{V}_{d,2}} \int_{S^1} g(z, U) (\tilde{R}\mu)^U(dz) \, d\sigma(U) &= \int_{S^1 \times \mathbb{V}_{d,2}} g(z, U) \, d(\tilde{R}\mu)(z, U) \\
&= \int_{S^{d-1}} \tilde{R}^* g(x) \, d\mu(x) \\
&= \int_{S^{d-1}} \int_{\mathbb{V}_{d,2}} g(P^U(x), U) \, d\sigma(U) d\mu(x) \\
&= \int_{\mathbb{V}_{d,2}} \int_{S^{d-1}} g(P^U(x), U) \, d\mu(x) d\sigma(U) \\
&= \int_{\mathbb{V}_{d,2}} \int_{S^1} g(z, U) \, d(P_{\#}^U \mu)(z) d\sigma(U).
\end{aligned} \tag{38}$$

684 Hence, for  $\sigma$ -almost every  $U \in \mathbb{V}_{d,2}$ ,  $(\tilde{R}\mu)^U = P_{\#}^U \mu$ .

#### 685 A.6 Study of the Spherical Radon transform $\tilde{R}$

686 In this Section, we first discuss the set of integration of the spherical Radon transform  $\tilde{R}$  (19). We  
687 further show that it is related to the hemispherical Radon transform and we derive its kernel.

688 **Set of integration.** While the classical Radon transform integrates over hyperplanes of  $\mathbb{R}^d$  and the  
689 generalized Radon transform integrates over hypersurfaces [60], the set of integration of the spherical  
690 Radon transform (19) is a half of a “big circle”, *i.e.* half of the intersection between a hyperplane and  
691  $S^{d-1}$  [96]. We illustrate this on  $S^2$  in Figure 7. On  $S^2$ , the intersection between a hyperplane and  $S^2$   
692 is a great circle.

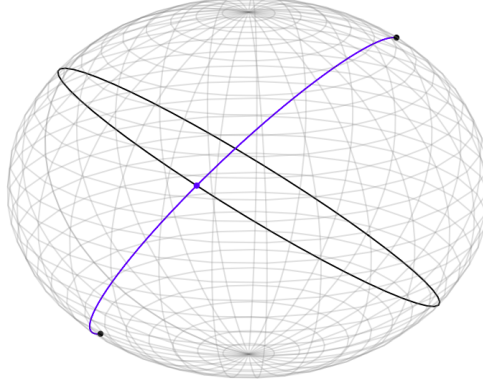


Figure 7: Set of integration of the spherical Radon transform (19). The great circle is in black and the set of integration in blue. The point  $Uz \in \text{span}(UU^T) \cap S^{d-1}$  is in blue.

693 **Proposition 6.** Let  $U \in \mathbb{V}_{d,2}$ ,  $z \in S^1$ . The set of integration of (19) is

$$\{x \in S^{d-1}, P^U(x) = z\} = \{x \in F \cap S^{d-1}, \langle x, Uz \rangle > 0\}, \quad (39)$$

694 where  $F = \text{span}(UU^T)^\perp \oplus \text{span}(Uz)$ .

695 *Proof.* Let  $U \in \mathbb{V}_{d,2}$ ,  $z \in S^1$ . Denote  $E = \text{span}(UU^T)$  the 2-plane generating the great circle,  
 696 and  $E^\perp$  its orthogonal complementary. Hence,  $E \oplus E^\perp = \mathbb{R}^d$  and  $\dim(E^\perp) = d - 2$ . Now, let  
 697  $F = E^\perp \oplus \text{span}(Uz)$ . Since  $Uz = UU^T Uz \in E$ , we have that  $\dim(F) = d - 1$ . Hence,  $F$  is a  
 698 hyperplane and  $F \cap S^{d-1}$  is a “big circle” [96], i.e. a  $(d - 2)$ -dimensional subsphere of  $S^{d-1}$ .

699 Now, for the first inclusion, let  $x \in \{x \in S^{d-1}, P^U(x) = z\}$ . First, we show that  $x \in F \cap S^{d-1}$ . By  
 700 Lemma 1 and hypothesis, we know that  $P^U(x) = \frac{U^T x}{\|U^T x\|_2} = z$ . By denoting by  $p^E$  the projection on  
 701  $E$ , we have:

$$p^E(x) = UU^T x = U(\|U^T x\|_2 z) = \|U^T x\|_2 Uz \in \text{span}(Uz). \quad (40)$$

702 Hence,  $x = p^E(x) + x_{E^\perp} = \|U^T x\|_2 Uz + x_{E^\perp} \in F$ . Moreover, as

$$\langle x, Uz \rangle = \|U^T x\|_2 \langle Uz, Uz \rangle = \|U^T x\|_2 > 0, \quad (41)$$

703 we deduce that  $x \in \{F \cap S^{d-1}, \langle x, Uz \rangle > 0\}$ .

704 For the other inclusion, let  $x \in \{F \cap S^{d-1}, \langle x, Uz \rangle > 0\}$ . Since  $x \in F$ , we have  $x = x_{E^\perp} + \lambda Uz$ ,  
 705  $\lambda \in \mathbb{R}$ . Hence, using Lemma 1,

$$P^U(x) = \frac{U^T x}{\|U^T x\|_2} = \frac{\lambda}{|\lambda|} \frac{z}{\|z\|_2} = \text{sign}(\lambda)z. \quad (42)$$

706 But, we also have  $\langle x, Uz \rangle = \lambda \|Uz\|_2^2 = \lambda > 0$ . Therefore,  $\text{sign}(\lambda) = 1$  and  $P^U(x) = z$ .

707 Finally, we conclude that  $\{x \in S^{d-1}, P^U(x) = z\} = \{x \in F \cap S^{d-1}, \langle x, Uz \rangle > 0\}$ .  $\square$

708 **Link with Hemispherical transform.** Since the intersection between a hyperplane and  $S^{d-1}$  is  
 709 isometric to  $S^{d-2}$  [56], we can relate  $\bar{R}$  to the hemispherical transform  $\mathcal{H}$  [96] on  $S^{d-2}$ . First, the  
 710 hemispherical transform of a function  $f \in L^1(S^{d-1})$  is defined as

$$\forall x \in S^{d-1}, \mathcal{H}f(x) = \int_{S^{d-1}} f(y) \mathbb{1}_{\{\langle x, y \rangle > 0\}} dy. \quad (43)$$

711 From Proposition 6, we can write the spherical Radon transform (19) as a hemispherical transform  
 712 on  $S^{d-2}$ .

713 **Proposition 7.** Let  $f \in L^1(S^{d-1})$ ,  $U \in \mathbb{V}_{d,2}$  and  $z \in S^1$ , then

$$\tilde{R}f(z, U) = \int_{S^{d-2}} \tilde{f}(x) \mathbb{1}_{\{\langle x, \tilde{U}z \rangle > 0\}} dx = \mathcal{H}\tilde{f}(\tilde{U}z), \quad (44)$$

714 where for all  $x \in S^{d-2}$ ,  $\tilde{f}(x) = f(O^T Jx)$  with  $O$  the rotation matrix such that for all  $x \in F$ ,  
 715  $Ox \in \text{span}(e_1, \dots, e_{d-1})$  where  $(e_1, \dots, e_d)$  denotes the canonical basis, and  $J = \begin{pmatrix} I_{d-1} \\ 0_{1,d-1} \end{pmatrix}$ , and  
 716  $\tilde{U} = J^T O U \in \mathbb{R}^{(d-1) \times 2}$ .

717 *Proof.* Let  $f \in L^1(S^{d-1})$ ,  $z \in S^1$ ,  $U \in \mathbb{V}_{d,2}$ , then by Proposition 6,

$$\tilde{R}f(z, U) = \int_{S^{d-1} \cap F} f(x) \mathbb{1}_{\{\langle x, U z \rangle > 0\}} dx. \quad (45)$$

718  $F$  is a hyperplane. Let  $O \in \mathbb{R}^{d \times d}$  be the rotation such that for all  $x \in F$ ,  $Ox \in \text{span}(e_1, \dots, e_{d-1}) =$   
 719  $\tilde{F}$  where  $(e_1, \dots, e_d)$  is the canonical basis. By applying the change of variable  $Ox = y$ , and since  
 720  $O^{-1} = O^T$ ,  $\det O = 1$ , we obtain

$$\tilde{R}f(z, U) = \int_{O(F \cap S^{d-1})} f(O^T y) \mathbb{1}_{\{\langle O^T y, U z \rangle > 0\}} dy = \int_{\tilde{F} \cap S^{d-1}} f(O^T y) \mathbb{1}_{\{\langle y, O U z \rangle > 0\}} dy. \quad (46)$$

721 Now, we have that  $O U \in \mathbb{V}_{d,2}$  since  $(O U)^T (O U) = I_2$ , and since  $U z \in F$ ,  $O U z \in \tilde{F}$ . For all  
 722  $y \in \tilde{F}$ , we have  $\langle y, e_d \rangle = y_d = 0$ . Let  $J = \begin{pmatrix} I_{d-1} \\ 0_{1,d-1} \end{pmatrix} \in \mathbb{R}^{d \times (d-1)}$ , then for all  $y \in \tilde{F} \cap S^{d-1}$ ,  
 723  $y = J \tilde{y}$  where  $\tilde{y} \in S^{d-2}$  is composed of the  $d-1$  first coordinates of  $y$ .  
 724 Let's define, for all  $\tilde{y} \in S^{d-2}$ ,  $\tilde{f}(\tilde{y}) = f(O^T J \tilde{y})$ ,  $\tilde{U} = J^T O U$ .

725 Then, since  $\tilde{F} \cap S^{d-1} \cong S^{d-2}$ , we can write:

$$\tilde{R}f(z, U) = \int_{S^{d-2}} \tilde{f}(\tilde{y}) \mathbb{1}_{\{\langle \tilde{y}, \tilde{U}z \rangle > 0\}} d\tilde{y} = \mathcal{H}\tilde{f}(\tilde{U}z). \quad (47)$$

726 □

727 **Kernel of  $\tilde{R}$ .** By exploiting the expression using the hemispherical transform in Proposition 7, we  
 728 can derive its kernel in Appendix A.7.

## 729 A.7 Proof of Proposition 4

730 First, we recall Lemma 2.3 of [94] on  $S^{d-2}$ .

731 **Lemma 2** (Lemma 2.3 [94]).  $\ker(\mathcal{H}) = \{\mu \in \mathcal{M}_{\text{even}}(S^{d-2}), \mu(S^{d-2}) = 0\}$  where  $\mathcal{M}_{\text{even}}$  is  
 732 the set of even measures, i.e. measures such that for all  $f \in C(S^{d-2})$ ,  $\langle \mu, f \rangle = \langle \mu, f^- \rangle$  where  
 733  $f^-(x) = f(-x)$  for all  $x \in S^{d-2}$ .

734 Let  $\mu \in \mathcal{M}_{ac}(S^{d-1})$ . First, we notice that the density of  $\tilde{R}\mu$  w.r.t.  $\lambda \otimes \sigma$  is, for all  $z \in S^1$ ,  $U \in \mathbb{V}_{d,2}$ ,  
 735

$$(\tilde{R}\mu)(z, U) = \int_{S^{d-1}} \mathbb{1}_{\{P^U(x)=z\}} d\mu(x) = \int_{F \cap S^{d-1}} \mathbb{1}_{\{\langle x, U z \rangle > 0\}} d\mu(x). \quad (48)$$

736 Indeed, using Proposition 2, and Proposition 6, we have for all  $g \in C_0(S^1 \times \mathbb{V}_{d,2})$ ,

$$\begin{aligned} \langle \tilde{R}\mu, g \rangle_{S^1 \times \mathbb{V}_{d,2}} &= \langle \mu, \tilde{R}^* g \rangle_{S^{d-1}} = \int_{S^{d-1}} R^* g(x) d\mu(x) \\ &= \int_{S^{d-1}} \int_{\mathbb{V}_{d,2}} \int_{S^1} g(z, U) \mathbb{1}_{\{z=P^U(x)\}} dz d\sigma(U) d\mu(x) \\ &= \int_{\mathbb{V}_{d,2} \times S^1} g(z, U) \int_{S^{d-1}} \mathbb{1}_{\{z=P^U(x)\}} d\mu(x) dz d\sigma(U) \\ &= \int_{\mathbb{V}_{d,2} \times S^1} g(z, U) \int_{F \cap S^{d-1}} \mathbb{1}_{\{\langle x, U z \rangle > 0\}} d\mu(x) dz d\sigma(U). \end{aligned} \quad (49)$$

737 Hence, using Proposition 7, we can write  $(\tilde{R}\mu)(z, U) = (\mathcal{H}\tilde{\mu})(\tilde{U}z)$  where  $\tilde{\mu} = J_{\#}^T O_{\#}\mu$ .

738 Now, let  $\mu \in \ker(\tilde{R})$ , then for all  $z \in S^1$ ,  $U \in \mathbb{V}_{d,2}$ ,  $\tilde{R}\mu(z, U) = \mathcal{H}\tilde{\mu}(\tilde{U}z) = 0$  and hence  
 739  $\tilde{\mu} \in \ker(\mathcal{H}) = \{\tilde{\mu} \in \mathcal{M}_{\text{even}}(S^{d-2}), \tilde{\mu}(S^{d-2}) = 0\}$ .

740 First, let's show that  $\mu \in \mathcal{M}_{\text{even}}(S^{d-1})$ . Let  $f \in C(S^{d-1})$  and  $U \in \mathbb{V}_{d,2}$ , then, by using the same  
 741 notation as in Propositions 6 and 7, we have

$$\begin{aligned}
 \langle \mu, f \rangle_{S^{d-1}} &= \int_{S^{d-1}} f(x) d\mu(x) = \int_{S^{d-1}} \int_{S^1} f(x) \mathbb{1}_{\{z=P^U(x)\}} dz d\mu(x) \\
 &= \int_{S^1} \int_{S^{d-1}} f(x) \mathbb{1}_{\{z=P^U(x)\}} d\mu(x) dz \\
 &= \int_{S^1} \int_{F \cap S^{d-1}} f(x) \mathbb{1}_{\{\langle x, Uz \rangle > 0\}} d\mu(x) dz \quad \text{by Prop. 6} \\
 &= \int_{S^1} \int_{S^{d-2}} \tilde{f}(y) \mathbb{1}_{\{\langle y, \tilde{U}z \rangle > 0\}} d\tilde{\mu}(y) dz \\
 &= \int_{S^1} \langle \mathcal{H}\tilde{\mu}, \tilde{f} \rangle_{S^{d-2}} dz \\
 &= \int_{S^1} \langle \tilde{\mu}, \mathcal{H}\tilde{f} \rangle_{S^{d-2}} dz \\
 &= \int_{S^1} \langle \tilde{\mu}, (\mathcal{H}\tilde{f})^- \rangle_{S^{d-2}} dz \quad \text{since } \tilde{\mu} \in \mathcal{M}_{\text{even}} \\
 &= \int_{S^{d-1}} f^-(x) d\mu(x) = \langle \mu, f^- \rangle_{S^{d-1}},
 \end{aligned} \tag{50}$$

742 using for the last line all the opposite transformations. Therefore,  $\mu \in \mathcal{M}_{\text{even}}(S^{d-1})$ .

743 Now, we need to find on which set the measure is null. We have

$$\begin{aligned}
 \forall z \in S^1, U \in \mathbb{V}_{d,2}, \tilde{\mu}(S^{d-2}) &= 0 \\
 \iff \forall z \in S^1, U \in \mathbb{V}_{d,2}, \mu(O^{-1}((J^T)^{-1}(S^{d-2}))) &= \mu(F \cap S^{d-1}) = 0.
 \end{aligned} \tag{51}$$

744 Hence, we deduce that

$$\begin{aligned}
 \ker(\tilde{R}) &= \{\mu \in \mathcal{M}_{\text{even}}(S^{d-1}), \forall U \in \mathbb{V}_{d,2}, \forall z \in S^1, F = \text{span}(UU^T)^\perp \cap \text{span}(Uz), \\
 &\quad \mu(F \cap S^{d-1}) = 0\}.
 \end{aligned} \tag{52}$$

745 Moreover, we have that  $\cup_{U,z} F_{U,z} \cap S^{d-1} = \{H \cap S^{d-1} \subset \mathbb{R}^d, \dim(H) = d-1\}$ .

746 Indeed, on the one hand, let  $H$  an hyperplane,  $x \in H \cap S^{d-1}$ ,  $U \in \mathbb{V}_{d,2}$ , and note  $z = P^U(x)$ . Then,  
 747  $x \in F \cap S^{d-1}$  by Proposition 6 and  $H \cap S^{d-1} \subset \cup_{U,z} F_{U,z}$ .

748 On the other hand, let  $U \in \mathbb{V}_{d,2}$ ,  $z \in S^1$ ,  $F$  is a hyperplane since  $\dim(F) = d-1$  and therefore  
 749  $F \cap S^{d-1} \subset \{H, \dim(H) = d-1\}$ .

750 Finally, we deduce that

$$\ker(\tilde{R}) = \{\mu \in \mathcal{M}_{\text{even}}(S^{d-1}), \forall H \in \mathcal{G}_{d,d-1}, \mu(H \cap S^{d-1}) = 0\}. \tag{53}$$

## 751 A.8 Proof of Proposition 5

752 Let  $p \geq 1$ . First, it is straightforward to see that for all  $\mu, \nu \in \mathcal{P}_p(S^{d-1})$ ,  $SSW_p(\mu, \nu) \geq 0$ ,  
 753  $SSW_p(\mu, \nu) = SSW_p(\nu, \mu)$ ,  $\mu = \nu \implies SSW_p(\mu, \nu) = 0$  and that we have the triangular

inequality since

$$\begin{aligned}
\forall \mu, \nu, \alpha \in \mathcal{P}_p(S^{d-1}), \quad SSW_p(\mu, \nu) &= \left( \int_{\mathbb{V}_{d,2}} W_p^p(P_{\#}^U \mu, P_{\#}^U \nu) d\sigma(U) \right)^{\frac{1}{p}} \\
&\leq \left( \int_{\mathbb{V}_{d,2}} (W_p(P_{\#}^U \mu, P_{\#}^U \alpha) + W_p(P_{\#}^U \alpha, P_{\#}^U \nu))^p d\sigma(U) \right)^{\frac{1}{p}} \\
&\leq \left( \int_{\mathbb{V}_{d,2}} W_p^p(P_{\#}^U \mu, P_{\#}^U \alpha) d\sigma(U) \right)^{\frac{1}{p}} \\
&\quad + \left( \int_{\mathbb{V}_{d,2}} W_p^p(P_{\#}^U \alpha, P_{\#}^U \nu) d\sigma(U) \right)^{\frac{1}{p}} \\
&= SSW_p(\mu, \alpha) + SSW_p(\alpha, \nu),
\end{aligned} \tag{54}$$

using the triangular inequality for  $W_p$  and the Minkowski inequality. Therefore, it is at least a pseudo-distance.

To be a distance, we also need  $SSW_p(\mu, \nu) = 0 \implies \mu = \nu$ . Suppose that  $SSW_p(\mu, \nu) = 0$ . Since, for all  $U \in \mathbb{V}_{d,2}$ ,  $W_p^p(P_{\#}^U \mu, P_{\#}^U \nu) \geq 0$ ,  $SSW_p(\mu, \nu) = 0$  implies that for  $\sigma$ -ae  $U \in \mathbb{V}_{d,2}$ ,  $W_p^p(P_{\#}^U \mu, P_{\#}^U \nu) = 0$  and hence  $P_{\#}^U \mu = P_{\#}^U \nu$  or  $(\tilde{R}\mu)^U = (\tilde{R}\nu)^U$  for  $\sigma$ -ae  $U \in \mathbb{V}_{d,2}$  since  $W_p$  is a distance on the circle. Therefore, it is a distance on the sets of injectivity of  $\tilde{R}$ .

## A.9 Convergence Properties

**Proposition 8.** Let  $(\mu_k), \mu \in \mathcal{P}_p(S^{d-1})$  such that  $\mu_k \xrightarrow[k \rightarrow \infty]{} \mu$ , then

$$SSW_p(\mu_k, \mu) \xrightarrow[k \rightarrow \infty]{} 0. \tag{55}$$

*Proof.* Since the Wasserstein distance metrizes the weak convergence (Corollary 6.11 [101]), we have  $P_{\#}^U \mu_k \xrightarrow[k \rightarrow \infty]{} P_{\#}^U \mu$  (by continuity)  $\iff W_p^p(P_{\#}^U \mu_k, P_{\#}^U \mu) \xrightarrow[k \rightarrow \infty]{} 0$  and hence by the dominated convergence theorem,  $SSW_p(\mu_k, \mu) \xrightarrow[k \rightarrow \infty]{} 0$ .  $\square$

## B Background on the Sphere

### B.1 Uniqueness of the Projection

Here, we discuss the uniqueness of the projection  $P^U$  for almost every  $x$ . For that, we recall some results of [9].

Let  $M$  be a closed subset of a complete finite-dimensional Riemannian manifold  $N$ . Let  $d$  be the Riemannian distance on  $N$ . Then, the distance from the set  $M$  is defined as

$$d_M(x) = \inf_{y \in M} d(x, y). \tag{56}$$

The infimum is a minimum since  $M$  is closed and  $N$  locally compact, but the minimum might not be unique. When it is unique, let's denote the point which attains the minimum as  $\pi(x)$ , i.e.  $d(x, \pi(x)) = d_M(x)$ .

**Proposition 9** (Proposition 4.2 in [9]). Let  $M$  be a closed set in a complete  $m$ -dimensional Riemannian manifold  $N$ . Then, for almost every  $x$ , there exists a unique point  $\pi(x) \in M$  that realizes the minimum of the distance from  $x$ .

From this Proposition, they further deduce that the measure  $\pi_{\#}\gamma$  is well defined on  $M$  with  $\gamma$  a locally absolutely continuous measure w.r.t. the Lebesgue measure.

In our setting, for all  $U \in \mathbb{V}_{d,2}$ , we want to project a measure  $\mu \in \mathcal{P}(S^{d-1})$  on the great circle  $\text{span}(UU^T) \cap S^{-1}$ . Hence, we have  $N = S^{d-1}$  which is a complete finite-dimensional Riemannian manifold and  $M = \text{span}(UU^T) \cap S^{d-1}$  a closed set in  $N$ . Therefore, we can apply Proposition 9 and the push-forward measures are well defined for absolutely continuous measures.

## 784 B.2 Optimization on the Sphere

785 Let  $F : S^{d-1} \rightarrow \mathbb{R}$  be some functional on the sphere. Then, we can perform a gradient descent on a  
 786 Riemannian manifold by following the geodesics, which are the counterpart of straight lines in  $\mathbb{R}^d$ .  
 787 Hence, the gradient descent algorithm [3, 14] reads as

$$\forall k \geq 0, x_{k+1} = \exp_{x_k}(-\gamma \text{grad} f(x)), \quad (57)$$

788 where for all  $x \in S^{d-1}$ ,  $\exp_x : T_x S^{d-1} \rightarrow S^{d-1}$  is a map from the tangent space  $T_x S^{d-1} = \{v \in$   
 789  $\mathbb{R}^d, \langle x, v \rangle = 0\}$  to  $S^{d-1}$  such that for all  $v \in T_x S^{d-1}$ ,  $\exp_x(v) = \gamma_v(1)$  with  $\gamma_v$  the unique geodesic  
 790 starting from  $x$  with speed  $v$ , i.e.  $\gamma(0) = x$  and  $\gamma'(0) = v$ .

791 For  $S^{d-1}$ , the exponential map is known and is

$$\forall x \in S^{d-1}, \forall v \in T_x S^{d-1}, \exp_x(v) = \cos(\|v\|_2)x + \sin(\|v\|_2)\frac{v}{\|v\|_2}. \quad (58)$$

792 Moreover, the Riemannian gradient on  $S^{d-1}$  is known as [3, Eq. 3.37]

$$\text{grad} f(x) = \text{Proj}_x(\nabla f(x)) = \nabla f(x) - \langle \nabla f(x), x \rangle x, \quad (59)$$

793  $\text{Proj}_x$  denoting the orthogonal projection on  $T_x S^{d-1}$ .

794 For more details, we refer to [3, 17].

## 795 B.3 Von Mises-Fisher Distribution

796 The von Mises-Fisher (vMF) distribution is a distribution on  $S^{d-1}$  characterized by a concentration  
 797 parameter  $\kappa > 0$  and a location parameter  $\mu \in S^{d-1}$  through the density

$$\forall \theta \in S^{d-1}, f_{\text{vMF}}(\theta; \mu, \kappa) = \frac{\kappa^{d/2-1}}{(2\pi)^{d/2} I_{d/2-1}(\kappa)} \exp(\kappa \mu^T \theta), \quad (60)$$

798 where  $I_\nu(\kappa) = \frac{1}{2\pi} \int_0^\pi \exp(\kappa \cos(\theta)) \cos(\nu\theta) d\theta$  is the modified Bessel function of the first kind.

799 Several algorithms allow to sample from it, see e.g. [100, 107] for algorithms using rejection sampling  
 800 or [62] without rejection sampling.

801 For  $d = 1$ , the vMF coincides with the von Mises (vM) distribution, which has for density

$$\forall \theta \in [-\pi, \pi[, f_{\text{vM}}(\theta; \mu, \kappa) = \frac{1}{I_0(\kappa)} \exp(\kappa \cos(\theta - \mu)), \quad (61)$$

802 with  $\mu \in [0, 2\pi[$  the mean direction and  $\kappa > 0$  its concentration parameter. We refer to [71, Section  
 803 3.5 and Chapter 9] for more details on these distributions.

804 In particular, for  $\kappa = 0$ , the vMF (resp. vM) distribution coincides with the uniform distribution on  
 805 the sphere (resp. the circle).

806 Jung [55] studied the law of the projection of a vMF on a great circle. In particular, they showed that,  
 807 while the vMF plays the role of the normal distributions for directional data, the projection actually  
 808 does not follow a von Mises distribution. More precisely, they showed the following theorem:

809 **Theorem 1** (Theorem 3.1 in [55]). *Let  $d \geq 3$ ,  $X \sim \text{vMF}(\mu, \kappa) \in S^{d-1}$ ,  $U \in \mathbb{V}_{d,2}$  and  $T = P^U(X)$   
 810 the projection on the great circle generated by  $U$ . Then, the density function of  $T$  is*

$$\forall t \in [-\pi, \pi[, f(t) = \int_0^1 f_R(r) f_{\text{vM}}(t; 0, \kappa \cos(\delta) r) dr, \quad (62)$$

811 where  $\delta$  is the deviation of the great circle (geodesic) from  $\mu$  and the mixing density is

$$\forall r \in ]0, 1[, f_R(r) = \frac{2}{I_\nu^*(\kappa)} I_0(\kappa \cos(\delta) r) r (1 - r^2)^{\nu-1} I_{\nu-1}^*(\kappa \sin(\delta) \sqrt{1 - r^2}), \quad (63)$$

812 with  $\nu = (d - 2)/2$  and  $I_\nu^*(z) = (\frac{z}{2})^{-\nu} I_\nu(z)$  for  $z > 0$ ,  $I_\nu^*(0) = 1/\Gamma(\nu + 1)$ .

813 Hence, as noticed by Jung [55], in the particular case  $\kappa = 0$ , i.e.  $X \sim \text{Unif}(S^{d-1})$ , then

$$f(t) = \int_0^1 f_R(r) f_{\text{vM}}(t; 0, 0) dr = f_{\text{vM}}(t; 0, 0) \int_0^1 f_R(r) dr = f_{\text{vM}}(t; 0, 0), \quad (64)$$

814 and hence  $T \sim \text{Unif}(S^1)$ .



## 815 B.4 Normalizing Flows on the Sphere

816 Normalizing flows [82] are invertible transformations. There has been a recent interest in defining  
817 such transformations on manifolds, and in particular on the sphere [23, 91, 92].

818 Here, we implemented the Exponential map normalizing flows introduced in [92]. The transformation  
819  $T$  is

$$\forall x \in S^{d-1}, z = T(x) = \exp_x(\text{Proj}_x(\nabla\phi(x))), \quad (65)$$

820 where  $\phi(x) = \sum_{i=1}^K \frac{\alpha_i}{\beta_i} e^{\beta_i(x^T \mu_i - 1)}$ ,  $\alpha_i \geq 0$ ,  $\sum_i \alpha_i \leq 1$ ,  $\mu_i \in S^{d-1}$  and  $\beta_i > 0$  for all  $i$ .  $(\alpha_i)_i$ ,  
821  $(\beta_i)_i$  and  $(\mu_i)_i$  are the learnable parameters.

822 The density of  $z$  can be obtained as

$$p_Z(z) = p_X(x) \det(E(x)^T J_T(x)^T J_T(x) E(x))^{-\frac{1}{2}}, \quad (66)$$

823 where  $J_f$  is the Jacobian in the embedded space and  $E(x)$  is the matrix whose columns form an  
824 orthonormal basis of  $T_x S^{d-1}$ .

825 The common way of training normalizing flows is to use either the reverse or forward KL divergence.  
826 Here, we use them with a different loss, namely SSW.

## 827 C Additional Experiments

### 828 C.1 Evolution of SSW between von Mises-Fisher distributions

829 The KL divergence between the von Mises-Fisher distribution and the uniform distribution has been  
830 derived analytically in [28, 110] as

$$\begin{aligned} \text{KL}(\text{vMF}(\mu, \kappa) || \text{vMF}(\cdot, 0)) &= \kappa \frac{I_{d/2}(\kappa)}{I_{d/2-1}(\kappa)} + \left(\frac{d}{2} - 1\right) \log \kappa - \frac{d}{2} \log(2\pi) - \log I_{d/2-1}(\kappa) \\ &\quad + \frac{d}{2} \log \pi + \log 2 - \log \Gamma\left(\frac{d}{2}\right). \end{aligned} \quad (67)$$

831 We plot on Figure 8 the evolution of KL and SSW *w.r.t.*  $\kappa$  for different dimensions. We observe a  
832 different trend. SSW seems to get lower with the dimension contrary to KL.

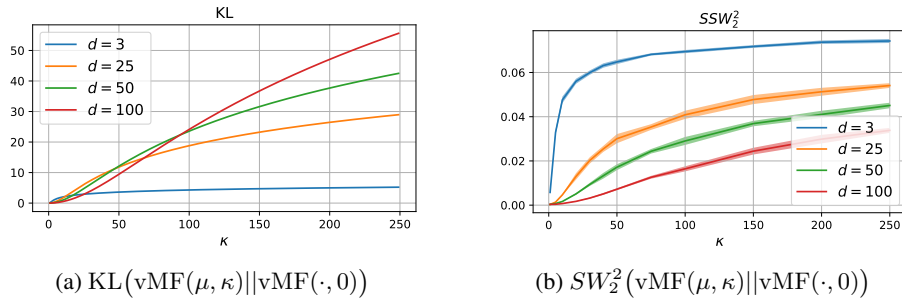


Figure 8: Evolution *w.r.t.*  $\kappa$  between  $\text{vMF}(\mu, \kappa)$  and  $\text{vMF}(\cdot, 0)$ . For SW, we used 100 projections (for memory reasons for  $d = 100$ ), and computed it for  $\kappa \in \{1, 5, 10, 20, 30, 40, 50, 75, 100, 150, 200, 250\}$ , 10 times by dimension and  $\kappa$ , and with 500 samples of both distributions.

833 As a sanity check, we compare on Figure 9 the evolution of SSW between vMF distributions  
834 where we fix  $\text{vMF}(\mu_0, 10)$  and we rotate the first vMF along a great circle. More precisely, we  
835 plot  $SW_2^2(\text{vMF}((1, 0, 0, \dots), 10), \text{vMF}((\cos(\theta), \sin(\theta), 0, \dots), 10))$  for  $\theta \in \{\frac{k\pi}{6}\}_{k \in \{0, \dots, 12\}}$ . As  
836 expected, we obtain a bell shape which is maximal when the second vMF distribution has for location

parameter  $-\mu_0$ . We observe a similar behavior between  $SSW_2$ ,  $SSW_1$  and  $SW_2$  with different scales.

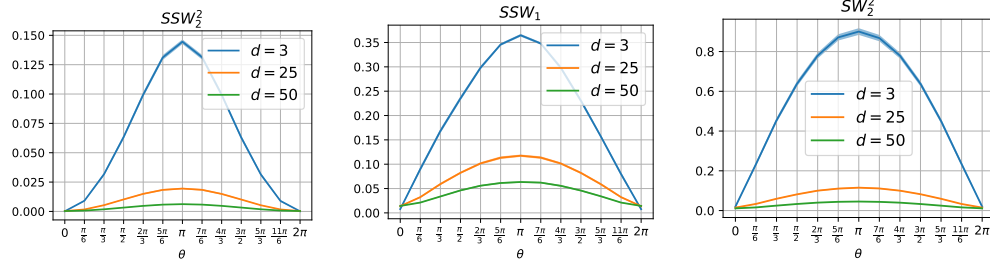


Figure 9: Evolution of  $SW$  between vMF samples in  $S^{d-1}$  (mean over 100 batch).

On Figure 10, we plot the evolution of SSW *w.r.t.* the number of projections for different dimensions. We observe that for around 100 projections, the variance seems to be low enough.

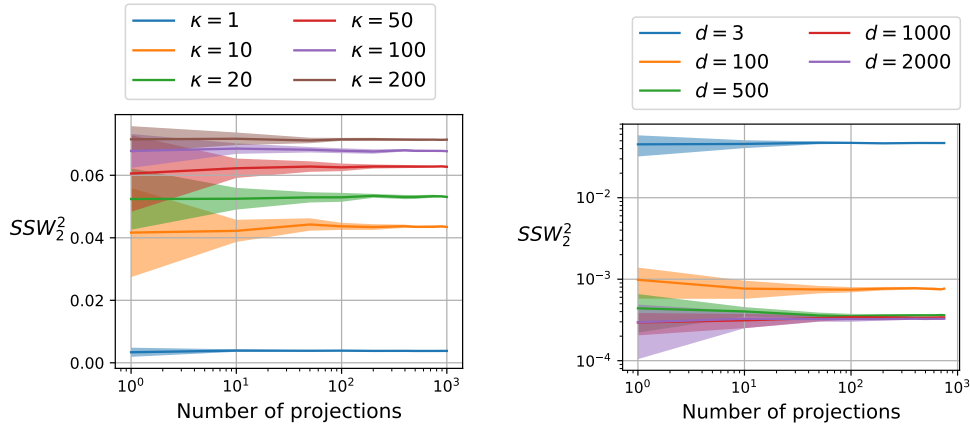


Figure 10: Influence of the number of projections. We compute  $SW_2^2(\text{vMF}(\mu, \kappa) || \text{vMF}(\cdot, 0))$  20 times, for  $n = 500$  samples in dimension  $d = 3$ .

Nadjahi et al. [76] proved that, contrary to the Wasserstein distance, the classical sliced-Wasserstein distance has a sample complexity independent of the dimension  $d$ . We show empirically on Figure 11 that we expect to have similar results for SSW by plotting SSW and the Wasserstein distance (with geodesic distance) between samples of the uniform distribution on the sphere *w.r.t.* the number of samples. We observe indeed that the convergence rate of SSW is independent of the dimension.

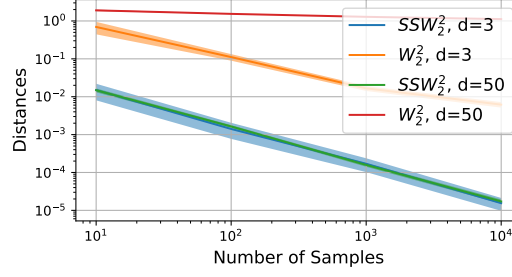


Figure 11: Spherical Sliced-Wasserstein and Wasserstein distance (with geodesic distance) between samples of the uniform distribution on the sphere. Results are averaged over 20 runs and the shaded are corresponds to the standard deviation.

## 846 C.2 Runtime Comparisons

847 We study here the evolution of the runtime *w.r.t.* different parameters. On Figure 12, we plot for  
 848 several dimensions the runtime to compute  $SSW_2$  *w.r.t.* the number of projections and the number of  
 849 samples. We observe the linearity *w.r.t.* the number of projections and the quasi-linearity *w.r.t.* the  
 850 number of samples.

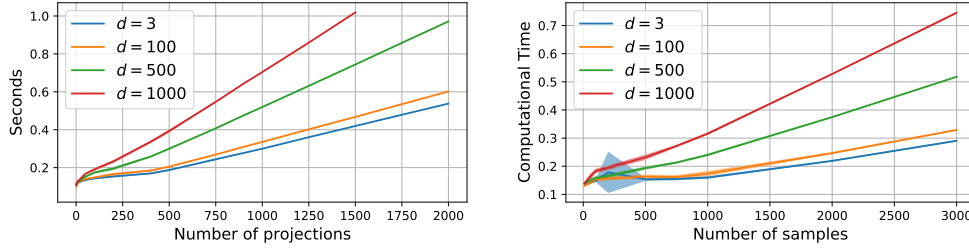


Figure 12: Computation time *w.r.t.* the number of projections or samples, taken for  $\kappa = 10$  and  $n = 500$  samples for the left figure, and  $\kappa = 10$  and 200 projections for the right figure, and for 20 times.

## 851 C.3 Gradient Flows

852 **Mixture of vMF distributions.** For the experiment in Section 5.1, we use as target distribution of  
 853 mixture of 6 vMF distributions from which we have access to samples. We refer to Appendix B.3 for  
 854 background on vMF distributions.

855 The 6 vMF distributions have weights  $1/6$ , concentration parameter  $\kappa = 10$  and location parameters  
 856  $\mu_1 = (1, 0, 0)$ ,  $\mu_2 = (0, 1, 0)$ ,  $\mu_3 = (0, 0, 1)$ ,  $\mu_4 = (-1, 0, 0)$ ,  $\mu_5 = (0, -1, 0)$  and  $\mu_6 = (0, 0, -1)$ .

857 We use two different approximation of the distribution. First, we approximate it using the empirical  
 858 distribution, *i.e.*  $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$  and we optimize over the particles  $(x_i)_{i=1}^n$ . To optimize over  
 859 particles, we can either use a projected gradient descent:

$$\begin{cases} x^{(k+1)} = x^{(k)} - \gamma \nabla_{x^{(k)}} SSW_2^2(\hat{\mu}_k, \nu) \\ x^{(k+1)} = \frac{x^{(k+1)}}{\|x^{(k+1)}\|_2}, \end{cases} \quad (68)$$

860 or a Riemannian gradient descent on the sphere [3] (see Appendix B.2 for more details). Note that  
 861 the projected gradient descent is a Riemannian gradient descent with retraction [17].

862 We can also use neural networks such as a multilayer perceptron (MLP). We used a MLP composed  
 863 of 5 layers of 100 units with leaky relu activation functions. The output of the MLP is normalized on  
 864 the sphere using a  $\ell^2$  normalization. We perform a gradient descent using Adam [57] as the optimizer

865 with a learning rate of  $10^{-4}$  for 2000 epochs. We approximate SSW with  $L = 1000$  projections and  
866 a batch size of 500. The base distribution is choose as the uniform distribution on the sphere.  
867 We report on Figure 13 a comparison of the 2 approximations where the density is estimated with a  
868 Gaussian kernel density estimator.

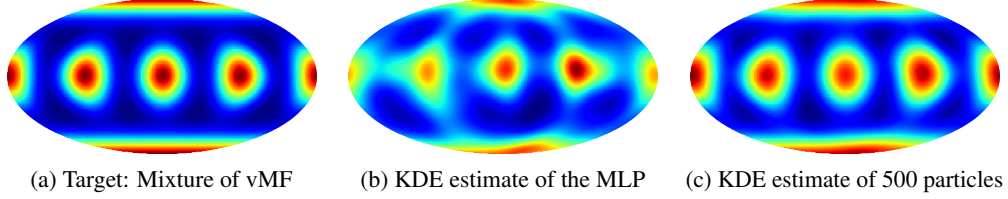


Figure 13: Minimization of SSW with respect to a mixture of vMF.

869 **vMF distribution.** A a simpler experiment, we choose a simple vMF distribution with  $\kappa = 10$ . We  
870 report on Figure 14 the evolution of the density approximated using a KDE, and on Figure 15 the  
871 evolution of particles.

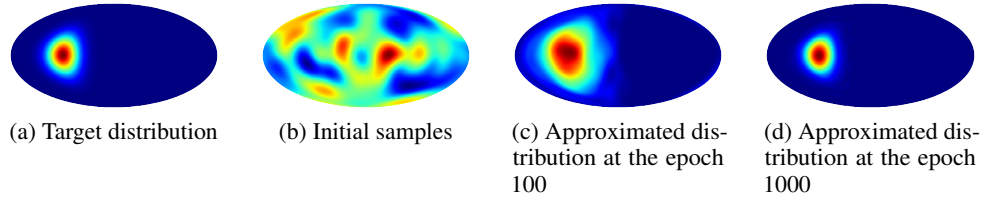


Figure 14: Gradient Flows on SW with a vMF target and Mollweide projections. The distributions are approximated using KDE.

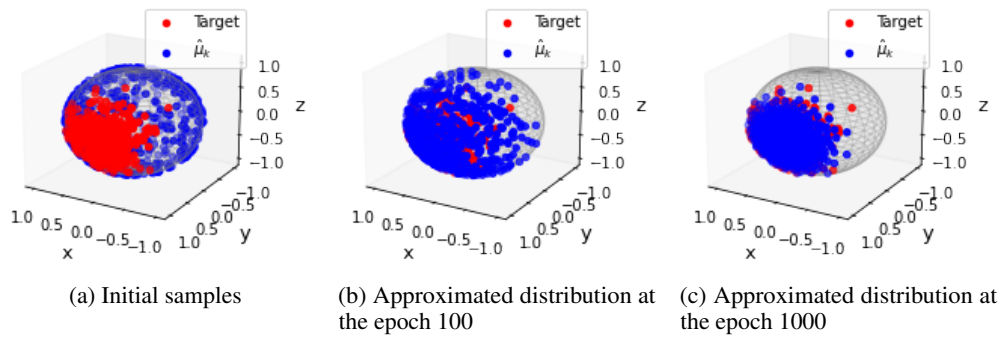


Figure 15: Gradient Flows on SW with a vMF target and Mollweide projections.

## 872 C.4 Sliced-Wasserstein Variational Inference

### 873 C.4.1 Variational Inference

874 In variational inference (VI) [12, 54], we have some observed data  $(x_i)_{i=1}^n$  and some latent data  
875  $(z_i)_{i=1}^n$ . The goal of variational inference is to approximate the posterior distribution  $p(\cdot|x)$  by some  
876 distribution  $q \in \mathcal{Q}$  where  $\mathcal{Q}$  is a family of probabilities. The usual way of doing that is to minimize

---

**Algorithm 2** SWVI [111]

---

**Input:**  $V$  a potential,  $K$  the number of iterations of SWVI,  $N$  the batch size,  $\ell$  the number of MCMC steps  
**Initialization:** Choose  $q_\theta$  a sampler  
**for**  $k = 1$  **to**  $K$  **do**  
  Sample  $(z_i^0)_{i=1}^N \sim q_\theta$   
  Run  $\ell$  MCMC steps starting from  $(z_i^0)_{i=1}^N$  to get  $(z_j^\ell)_{j=1}^N$   
  // Denote  $\hat{\mu}_0 = \frac{1}{N} \sum_{j=1}^N \delta_{z_j^0}$  and  $\hat{\mu}_\ell = \frac{1}{N} \sum_{j=1}^N \delta_{z_j^\ell}$   
  Compute  $J = SW_2^2(\hat{\mu}_0, \hat{\mu}_\ell)$   
  Backpropagate through  $J$  w.r.t.  $\theta$   
  Perform a gradient step  
**end for**

---

877 the Kullback-Leibler divergence among this family, *i.e.*

$$\min_{q \in \mathcal{Q}} \text{KL}(q||p(\cdot|x)) = \mathbb{E}_q[\log \left( \frac{q(Z)}{p(Z|x)} \right)]. \quad (69)$$

878 But the KL divergence suffers from some drawbacks, as it is only a divergence (*i.e.* it does not satisfy  
879 the triangular inequality, and it is non symmetric), but it also suffers from under estimating the target  
880 distribution (or over estimating it for the reverse KL).

881 Yi and Liu [111] propose to use an optimal transport distance instead, namely the SW distance  
882 which gives the sliced-Wasserstein variational inference method. Basically, given some unnormalized  
883 probability  $p(\cdot|x)$  that we want to approximate with some variational distribution  $q_\phi$ , we can first  
884 apply a MCMC algorithm and then learn  $q_\phi$  using a gradient descent on SW with the target being  
885 the empirical distributions of the samples given by the MCMC. But running long MCMC chain is  
886 time consuming and it might be difficult to diagnose burn-in period. Therefore, they propose to only  
887 run at each iteration some number of steps  $t$  of MCMC chain, and then learn by gradient descent the  
888 variational distribution. Therefore, the variational distribution is guided at each step by the MCMC  
889 samples toward the stationary distribution which is the target. This is called an amortized sampler  
890 (see Problem 1 in [103]). We sum up the procedure in Algorithm 2.

891 We propose here to substitute  $SW$  by  $SSW$  in order to perform SSWVI on the sphere. To do that,  
892 we first need a MCMC method on the sphere.

### 893 C.4.2 MCMC on the Sphere

894 Several MCMC methods on the sphere have been proposed. For example, Hamiltonian Monte-Carlo  
895 (HMC) methods were proposed in [18, 63, 68], and Riemannian Langevin algorithms were proposed  
896 in [65, 105].

897 In our experiments, we use the Geodesic Langevin algorithm (GLA) introduced by Wang et al.  
898 [105]. This algorithm is a natural generalization of the Unadjusted Langevin Algorithm (ULA) and it  
899 consists at simply following the geodesics of the regular ULA step, *i.e.*

$$\forall k > 0, x_{k+1} = \exp_{x_k}(\text{Proj}_{x_k}(-\gamma \nabla V(x_k) + \sqrt{2\gamma}Z)), Z \sim \mathcal{N}(0, I), \quad (70)$$

900 where for the sphere,

$$\forall x \in S^{d-1}, \forall v \in T_x S^{d-1}, \exp_x(v) = x \cos(\|v\|) + \frac{v}{\|v\|} \sin(\|v\|), \quad (71)$$

901  $\text{Proj}_x$  is the projection on the tangent space  $T_x S^{d-1} = \{v \in \mathbb{R}^d, \langle x, v \rangle = 0\}$  (which is the  
902 orthogonal space) and is defined as

$$\text{Proj}_x(v) = v - \langle x, v \rangle x. \quad (72)$$

903 For more details, we refer to [3].

We use GLA here for simplicity and as a proof of concept. But note that GLA, as ULA, is biased and therefore the distribution learned will not be the exact true stationary distribution. However, a Metropolis-Hastings step at each iteration could be used to enforce the reversibility *w.r.t.* the target distribution or we could use other MCMC with more appealing convergence properties (see *e.g.* [68]).

### C.4.3 Applications

**Target: Power spherical distribution.** First, as a simple example on  $S^2$ , we use the power spherical distribution introduced by De Cao and Aziz [29]. This distribution has the advantage over the vMF distribution to allow for the direct use of the reparameterization trick since it does not require rejection sampling. The pdf is obtained as,

$$\forall x \in S^{d-1}, p_X(x; \mu, \kappa) \propto (1 + \mu^T x)^\kappa \quad (73)$$

with  $\mu \in S^{d-1}$  and  $\kappa > 0$ . We can sample from drawing first  $Z \sim \text{Beta}(\frac{d-1}{2} + \kappa, \frac{d-1}{2})$ ,  $v \sim \text{Unif}(S^{d-2})$ , then constructing  $T = 2Z - 1$  and  $Y = [T, v^T \sqrt{1 - T^2}]^T$ . Finally, apply a Householder reflection about  $\mu$  to  $Y$ . All the operations are well differentiable and allow to apply the reparametrization trick. For the algorithm, see Algorithm 1 in [29]. Hence, in this case, if we denote  $g_\theta$  the map which takes samples from a uniform distribution on  $S^{d-2}$  and from a Beta distribution as input and outputs samples of power spherical distribution with parameters  $\theta = (\kappa, \mu)$ , we can use it as the sampler. We test the algorithm with a target being a power spherical distribution of parameter  $\mu = (0, 1, 0)$  and  $\kappa = 10$ , starting from  $\mu = (1, 1, 1)$  and  $\kappa = 0.1$ . Performing 2000 optimization steps with a gradient descent (Riemannian gradient descent on  $\mu$  to stay on the sphere), and 20 steps of the GLA algorithm, we are getting close enough to the true distribution as we can see on Figure 16.

For the hyperparameters, we used a step size of  $10^{-3}$  for GLA, 1000 projections to approximate SSW, a Riemannian gradient descent on the sphere [3] to learn the location parameter  $\mu$  with a learning rate of 2, and a learning of 200 for  $\kappa$ . We performed  $K = 2000$  steps and used  $N = 500$  particles.

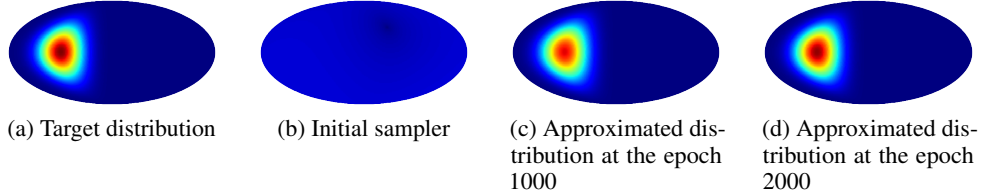


Figure 16: SWVI on Power Spherical Distributions with Mollweide projections.

**Target: mixture of vMFs.** In Section 5.1, we perform amortized variational inference with a mixture of vMF distributions as target. For this, we train exponential map normalizing flows (see [92] and Appendix B.4). Moreover, we use the same target as Rezende et al. [92], *i.e.* the target  $\nu$  has a density  $p(x) \propto \sum_{k=1}^4 e^{10x^T T_{s \rightarrow e}(\mu_k)}$  with  $\mu_1 = (0.7, 1.5)$ ,  $\mu_2 = (-1, 1)$ ,  $\mu_3 = (0.6, 0.5)$  and  $\mu_4 = (-0.7, 4)$ . These are spherical coordinates which are converted to euclidean using  $T_{s \rightarrow e}(\theta, \phi) = (\sin \phi \cos \theta, \sin \phi \sin \theta, \cos \phi)$ .

The exponential map normalizing flow is composed of  $N = 6$  blocks with  $K = 5$  components. We run the algorithm for 10000 iterations, with at each iteration 20 steps of GLA with  $\gamma = 10^{-1}$  as learning rate, and one step of backpropagation through SSW using the Adam [57] optimizer with a learning rate of  $10^{-3}$ .

We report on Figure 4 the Mollweide projection of the learned density. Since we learn to samples from a noise distribution, here the uniform distribution on the sphere, we do not have directly access to the density and we report a kernel density estimate with a Gaussian kernel using the implementation of Scipy [102].

940 We also report in Figure 5 the effective sample size (ESS) [33, 69] over the iterations. The ESS is  
 941 estimated by [92]

$$\text{ESS} = \frac{\text{Var}_{\text{Unif}}(e^{-\beta u(X)})}{\text{Var}_q\left(\frac{e^{-\beta u(X)}}{q_\eta(X)}\right)} \approx \frac{\left(\sum_{s=1}^S w_s\right)^2}{\sum_{s=1}^S w_s^2}, \quad (74)$$

942 where  $w_s = e^{-\beta u(x_s)/q_\eta(x_s)}$ . The ESS is reported as a percentage of the sample size. Higher ESS  
 943 indicates that the flow matches the target better [92].

## 944 C.5 Sliced-Wasserstein Autoencoder

945 We recall that in the WAE framework, we want to minimize

$$\mathcal{L}(f, g) = \int c(x, g(f(x))) d\mu(x) + \lambda D(f_{\#}\mu, p_Z), \quad (75)$$

946 where  $f$  is an encoder,  $g$  a decoder,  $p_Z$  a prior distribution,  $c$  some cost function and  $D$  is a divergence  
 947 in the latent space. Several  $D$  were proposed. For example, Tolstikhin et al. [99] proposed to use  
 948 the MMD, Kolouri et al. [59] used the SW distance, Patrini et al. [84] used the Sinkhorn divergence,  
 949 Kolouri et al. [60] used the generalized SW distance. Here, we use  $D = \text{SSW}_2^2$ .

950 **Architecture and procedure.** For the encoder  $f$  and the decoder  $g$ , we use the same architecture  
 951 as Kolouri et al. [59].

952 For both the encoder and the decoder architecture, we use fully convolutional architectures with 3x3  
 953 convolutional filters. More precisely, the architecture of the encoder is

$$\begin{aligned} x \in \mathbb{R}^{28 \times 28} &\rightarrow \text{Conv2d}_{16} \rightarrow \text{LeakyReLU}_{0.2} \\ &\rightarrow \text{Conv2d}_{16} \rightarrow \text{LeakyReLU}_{0.2} \rightarrow \text{AvgPool}_2 \\ &\rightarrow \text{Conv2d}_{32} \rightarrow \text{LeakyReLU}_{0.2} \\ &\rightarrow \text{Conv2d}_{32} \rightarrow \text{LeakyReLU}_{0.2} \rightarrow \text{AvgPool}_2 \\ &\rightarrow \text{Conv2d}_{64} \rightarrow \text{LeakyReLU}_{0.2} \\ &\rightarrow \text{Conv2d}_{64} \rightarrow \text{LeakyReLU}_{0.2} \rightarrow \text{AvgPool}_2 \\ &\rightarrow \text{Flatten} \rightarrow \text{FC}_{128} \rightarrow \text{ReLU} \\ &\rightarrow \text{FC}_{d_Z} \rightarrow \ell^2 \text{ normalization} \end{aligned}$$

954 where  $d_Z$  is the dimension of the latent space (either 11 for  $S^{10}$  or 3 for  $S^2$ ).

955 The architecture of the decoder is

$$\begin{aligned} z \in \mathbb{R}^{d_Z} &\rightarrow \text{FC}_{128} \rightarrow \text{FC}_{1024} \rightarrow \text{ReLU} \\ &\rightarrow \text{Reshape}(64 \times 4 \times 4) \rightarrow \text{Upsample}_2 \rightarrow \text{Conv}_{64} \rightarrow \text{LeakyReLU}_{0.2} \\ &\rightarrow \text{Conv}_{64} \rightarrow \text{LeakyReLU}_{0.2} \\ &\rightarrow \text{Upsample}_2 \rightarrow \text{Conv}_{64} \rightarrow \text{LeakyReLU}_{0.2} \\ &\rightarrow \text{Conv}_{32} \rightarrow \text{LeakyReLU}_{0.2} \\ &\rightarrow \text{Upsample}_2 \rightarrow \text{Conv}_{32} \rightarrow \text{LeakyReLU}_{0.2} \\ &\rightarrow \text{Conv}_1 \rightarrow \text{Sigmoid} \end{aligned}$$

956 To compare the different autoencoders, we used as the reconstruction loss the binary cross entropy,  
 957  $\lambda = 10$ , Adam [57] as optimizer with a learning rate of  $10^{-3}$  and Pytorch’s default momentum  
 958 parameters for 800 epochs with batch of size  $n = 500$ . Moreover, when using SW type of distance,  
 959 we approximated it with  $L = 1000$  projections.

960 We report in Table 1 the FID obtained using 10000 samples and we report the mean over 5 trainings.

961 For SSW, we used the formulation using the uniform distribution (12). To compute SW, we used the  
 962 POT library [39]. To compute the Sinkhorn divergence, we used the GeomLoss package [37].



963 **Additional experiments.** We report on Figure 17 samples obtained with SSW for a uniform prior  
 964 on  $S^{10}$ .

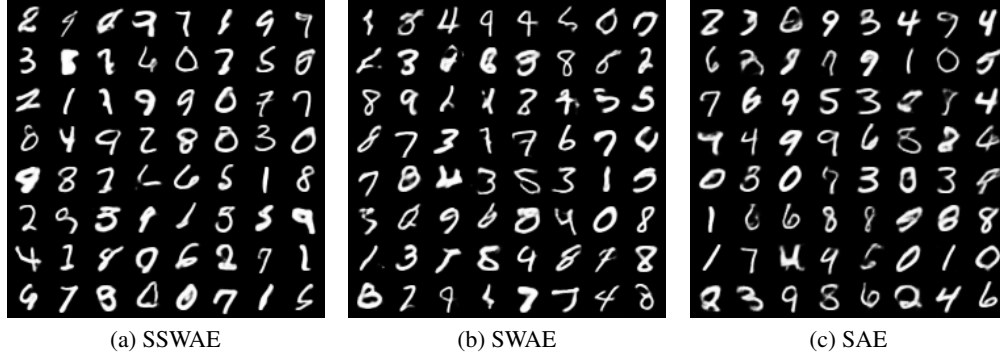


Figure 17: Samples generated with Sliced-Wasserstein Autoencoders with a uniform prior on  $S^{10}$ .

965 On Figure 18, we add the evolution over epochs of the Wasserstein distance between generated  
 966 images and samples from the test set.

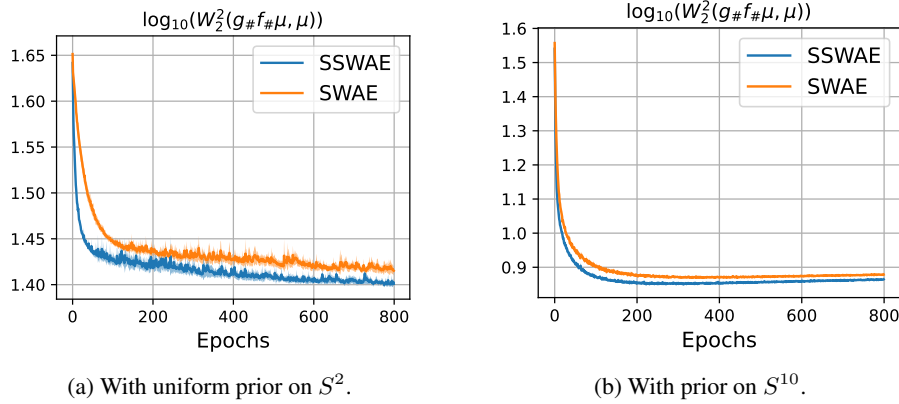


Figure 18: Comparison of the evolution of the Wasserstein distance over epochs between SWAE and SSWAE on MNIST (averaged over 5 trainings).

## 967 C.6 Self-supervised learning

968 We conduct experiments using SSW to  
 969 prevent collapsing representations in con-  
 970 trastive self-supervised learning (SSL)  
 971 models. Such contrastive losses on the hy-  
 972 persphere have exhibited great representa-  
 973 tive capacity [20, 21, 108] on unlabelled  
 974 datasets by learning robust image represen-  
 975 tations invariantly to augmentations. As  
 976 proposed in [104], the contrastive objec-  
 977 tive can be decomposed into an alignment  
 978 loss which forces positive representations  
 979 coming from the same image to be similar  
 980 and a uniformity loss which preserves maximal information of the feature distribution and hence  
 981 avoids collapsing representations. Without the uniformity loss, the representations tend to converge

Table 2: Linear evaluation on CIFAR10. The features are taken either on the encoder output or directly on the sphere  $S^2$ .

| Method                          | Encoder output | $S^2$ |
|---------------------------------|----------------|-------|
| Supervised                      | 82.26          | 81.43 |
| Chen et al. [21]                | 66.55          | 59.09 |
| Wang and Isola [104]            | 60.53          | 55.86 |
| SW-SSL, $\lambda = 1, L = 10$   | 62.65          | 57.77 |
| SW-SSL, $\lambda = 1, L = 3$    | 62.46          | 57.64 |
| SSW-SSL, $\lambda = 20, L = 10$ | 64.89          | 58.91 |
| SSW-SSL, $\lambda = 20, L = 3$  | 63.75          | 59.75 |



towards a constant representation which yields the best alignment loss possible but also contains no information about original images. Wang and Isola [104] propose to enforce uniformity by leveraging the Gaussian potential kernel which is bound to the uniform distribution on the sphere. This formulation is also related to the denominator of the contrastive loss as specified in Chen et al. [21]. We propose to replace the Gaussian kernel uniformity loss with SSW for which the complexity is more linear *w.r.t.* the number of batch samples. A simple choice of the alignment loss is to minimize the mean squared euclidean distance between pairs of different augmented versions of the same image. A self-supervised learning network is pre-trained using this alignment loss added with an uniformity term. Our overall self-supervised loss can be defined as:

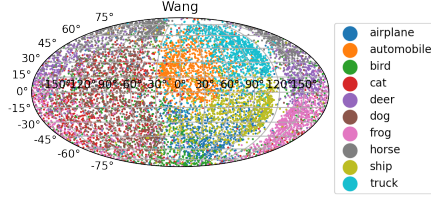
$$\mathcal{L}_{\text{SSW-SSL}} = \underbrace{\frac{1}{n} \sum_{i=1}^n \|z_i^A - z_i^B\|_2^2}_{\text{Alignment loss}} + \underbrace{\frac{\lambda}{2} (SSW_2^2(z^A, \nu) + SSW_2^2(z^B, \nu))}_{\text{Uniformity loss}}, \quad (76)$$

where  $z^A, z^B \in \mathbb{R}^{n \times d}$  are the representations from the network projected on the hypersphere of two augmented versions of the same images,  $\nu = \text{Unif}(S^{d-1})$  is the uniform distribution on the hypersphere and  $\lambda > 0$  is used to balance the two terms.

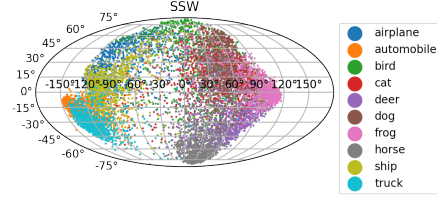
We pretrain a ResNet18 [47] model on the CIFAR10 [61] data with projections projected onto the sphere  $S^2$ . This feature dimension allow us to visualize the entire validation set of CIFAR10 and its distribution on the sphere. The visualization of the projections on  $S^2$  are visible on Figure 19. We then evaluate the performance of each contrastive objective by fitting a linear classifier on top of the output of the layer before the projection on the sphere on the training dataset as is common for SSL methods. For comparison, we also report the results when the features are taken directly on the sphere. As a baseline, we also train a predictive supervised encoder by training jointly the linear classifier and the image encoder in a supervised manner using cross entropy.

We use a ResNet18 [47] encoder which outputs 1024 features that are then projected onto the sphere  $S^2$  using a last fully connected layer followed by a  $\ell^2$  normalization. We pretrain the model for 200 epochs using minibatch stochastic gradient descent (SGD) with a momentum of 0.9, a weight decay of 0.001 and an initial learning rate of 0.05. We use a batch size of 512 samples. The images are augmented using a standard set of random augmentations for SSL: random crops, horizontal flipping, color jittering and gray scale transformation as done in Wang and Isola [104]. For the trade-off parameter  $\lambda$ , we  $\lambda = 20$  for SSW and  $\lambda = 1$  for SW.

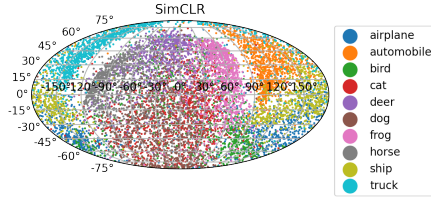
To evaluate the performance of representations, we use the common linear evaluation protocol where a linear classifier is fitted on top of the pre-trained representations and the best validation accuracy is reported. The linear classifiers are trained for 100 epochs using the Adam [57] optimizer with a learning rate of 0.001 with a decay of 0.2 at epoch 60 and 80. We compare our methods with two other contrastive objectives, Chen et al. [21] with the normalized temperature-scaled cross-entropy (NT-Xent) loss and Wang and Isola [104] which proposes to decompose the objective in two distinct terms  $\mathcal{L}_{\text{align}}$  and  $\mathcal{L}_{\text{uniform}}$ . We recall the respective uniformity loss of each method in Table 3. As one can see in Table 2, our method achieves here comparable performances to two state-of-the-art approaches, yet slightly under-performing compared to [21]. We suspect that a finer validation of the balancing parameter  $\lambda$  is needed. Especially since the representations on Figure 19b are not completely uniformly distributed around the sphere after pre-training compared to other contrastive methods. Nevertheless, these preliminary results show that SSW-SSL is a promising contrastive learning approach without explicit distances between negative samples, especially compared to SW on the sphere. To this end, further works should be devoted to finding a good balance between the alignment and uniformity objectives.



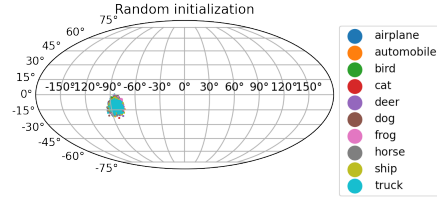
(a) Wang and Isola [104]



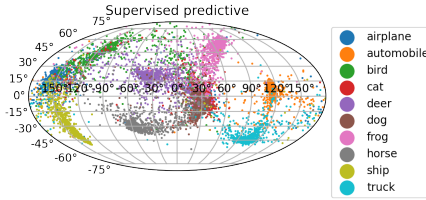
(b) SSW-SSL,  $\lambda = 6, L = 200$



(c) Chen et al. [21]



(d) Random initialization



(e) Supervised prediction

Figure 19: The CIFAR10 validation set on  $S^2$  after pre-training.

Table 3: Comparison of contrastive methods and their respective uniformity objective where  $z^A, z^B \in \mathbb{R}^{n \times d}$  are representations from two augmented versions of the same set of images and  $\nu = \text{Unif}(S^{d-1})$  is the uniform distribution on the hypersphere.

| Method               | $\mathcal{L}_{\text{uniform}}(z^A) + \mathcal{L}_{\text{uniform}}(z^B)$  | Complexity          |
|----------------------|--|---------------------|
| Chen et al. [21]     | $\frac{1}{2n} \sum_{i=1}^n \log \sum_{j \neq i} \exp(\frac{\langle \hat{z}_i, \hat{z}_j \rangle}{\tau}), \hat{z} = \text{cat}(z^A, z^B)$ | $O(n^2 d)$          |
| Wang and Isola [104] | $\sum_{z \in \{z^A, z^B\}} \log \frac{2}{n(n-1)} \sum_{i>j} \exp(-t \ z_i - z_j\ _2^2)$  | $O(n^2 d)$          |
| SSW-SSL (Ours)       | $\frac{1}{2} (SSW_2^2(z^A, \nu) + SSW_2^2(z^B, \nu))$  | $O(Ln(d + \log n))$ |