

1 **Reviewer #1 (R1):** *◆[...]*tree-sliced Wasserstein distance outperforms the original optimal transport distance[...] We
2 used the tree-sliced Wasserstein (TSW) and OT distances within an RBF kernel, $k = \exp(-td)$. The TSW kernel is
3 p.d. while the OT kernel k_{OT} is not. The indefiniteness of k_{OT} may affect its performances in some applications: k_{OT}
4 performs worse in TDA applications in §6.1, but works well for documents with word embedding applications in §6.2.
5 *◆[...]*include (in the supplement, at least) the full proof [of Proposition 1]? Agreed, we will follow your suggestion.
6 **Reviewer #2:** *◆[...]*Prop. 1 and 2 are not original[...] **◆R1:** proof of negative definiteness of the proposed OT[...] We
7 do state that the proof of Prop. 1 is not original in ℓ.39–41. Although Prop. 2 follows from Prop. 1, it follows the idea
8 underlying sliced W kernels (or Gaussian processes as you mention), and this result remains new to our knowledge.
9 *◆[...]*bound on the Wasserstein distance[...]interesting re-
10 sult[...]numeric analysis comparing the limit of the hypercube
11 tree-sliced metric to the Wasserstein metric and the sliced Wasser-
12 stein metric. **◆R1:** An upper bound on the Euclidean OT[...] The
13 hypercube tree-sliced metric is our suggestion to build practical
14 tree metrics for TSW when used on low-dimensional data spaces
15 (e.g. TDA in §6.1). We insist that we do not try to mimic the
16 Euclidean OT (W_2) or the sliced-Wasserstein (SW), but rather
17 propose a variant of OT distance. As stated in ℓ.173–179, and
18 ℓ.158–161 in §5, SW is a special case of TSW. Following your
19 point, we have carried out the following experiment: for a query
20 point q , let p be its nearest neighbor w.r.t. TSW. Figure 1 illus-
21 trates that p is very likely among the top 5 (on MPEG7), and
22 top 10 (on Orbit) near neighbors on the W_2 space (results are
23 averaged over 1000 runs of random split 90%/10% for training
24 and test). When the number of tree-slices in TSW increases,
25 the W_2 near-neighbor rank of p is improved. These empirical
26 results suggest that TSW may agree with some aspects of W_2 .
27 *◆[...]*experiment where the metric is used in a more conventional
28 OT problem such as color transfer or generative modeling. In the experiments, we used RBF kernels ($k = \exp(-td)$)
29 for a given metric d) with SVM which usually improves on k -NN results. We will add k -NN results. We are now
30 considering color transfer and barycenter applications. Gradients of TSW w.r.t. supports and weights of empirical
31 measures can be recovered pending some choices in how interpolations are defined. *◆a single multi paneled figure [...]*
32 *word embedding experiment[...]presented before the TDA* Many thanks for your suggestions. We will incorporate them.
33 **Reviewer #3:** *◆What is most troubling is that the paper seems to be completely unaware any literature of embedding*
34 *points into a distribution over metrics, and claims some standard and well-known techniques and novelties[...]* We
35 understand your point and will do everything to correct this misunderstanding. This was caused by a lack of care in
36 the presentation of §4. This was not the message we wanted to convey. We will rewrite this section following your
37 comments. As you have gathered from our algorithms, approximating an arbitrary metric using trees is *not* a key goal in
38 our submission, our goal is stated in ℓ.41–44 in §.1. Much like 1D projections do not offer interesting properties from a
39 distortion perspective but remain useful for SW, we do believe that trees with large distortion can still remain useful.
40 This is because metric approximations are used within another computation (Wasserstein) and therefore we do not gain
41 from overfitting too much our trees so that they match the true metric, as long as they provide guidance on the optimal
42 assignment. We will insist more on the importance of sampling tree metrics randomly, both for low-dimensional in §6.1
43 and high-dimensional §6.2 regimes. *◆Definite-negativity is mentioned and highlighted[...] explain why is it important*
44 *to you. Is this to ensure that the kernel is positive-definite?* Negative definiteness of a distance means essentially that
45 the space is flat and that positive definite kernels can be easily derived from them, following Berg et al.’s theorem.
46 This is why kernel methods kick in from §.6 (or Gaussian processes as per Reviewer #2’s suggestion). We will clarify
47 this motivation following your comment. *◆[...]*which kernel did you actually use in the experiments – $\exp(-TW)$ or
48 $\exp(-TSW)$?[...]kernel is called k_{TW} and not k_{TSW} ?[...]TSW is also negative-definite, simply because the average of
49 l_1 -metrics is an l_1 -metric – right?) In the experiments, we used the kernel $\exp(-td_{TSW})$ as stated in ℓ.225–227 for §6.1,
50 and ℓ.289–292 for §6.2. We will define k_{TSW} , and rename k_{TW} to k_{TSW} in the experiments following your suggestion.
51 Indeed, averaging of negative definite functions is trivially negative definite. Hence, d_{TSW} is negative definite. We will
52 clarify it in the updated version. *◆[...]*set the number of clusters in each level?[...]connection to fast Gauss transform
53 [...]more related than any other clustering method[...]use farthest-point clustering[...]downstream motivation and does
54 it effect the results? We fixed the same number of clusters κ when performing the farthest-point clustering (replace
55 step 9 in Algorithm 1) at different height levels. κ is typically chosen via cross-validation. Moreover, we also illustrate
56 the effect of κ in applications in Figure 5 (and Figures 5–7 in the supplement). In general, one can apply any favorite
57 clustering methods. We used the farthest-point clustering due to its fast computation, i.e. $O(n \log \kappa)$ as stated in
58 ℓ.148–149 to construct practical tree metrics for applications with high-dimensional data space, e.g. documents with
59 word embedding applications in §6.2.

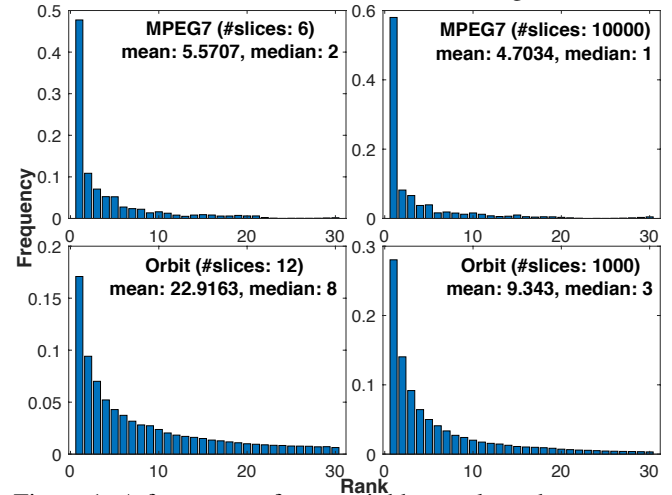


Figure 1: A frequency of near-neighbor rank on the W_2 space for the nearest neighbor w.r.t. TSW.