

1 We thank the reviewers for their insightful comments. We address their interrogations and comments below.

2 **R1.** - *Practical/concrete motivation for why/when online OT computation is needed.* Online OT estimation is useful for
3 training generative models. In this case, the samples are renewed at each iteration of a training algorithm, and may be
4 used to better evaluate the distance to minimize. Any applications that requires to estimate OT distances between large
5 point clouds can benefit from online OT estimation, that accelerates training. We will better motivate our work.

6 - *No experiments on high-dimensional settings, arguably those for which a streaming setting would be most compelling.*
7 This is indeed a weakness in our experiments, as we have measured performance up to $d = 10$. Real-world ML
8 applications would typically consider points in latent spaces, with typical dimension $d = 128$, with a certain "manifold"
9 structure that makes OT estimation possible. We may consider the output of a trained CNN on CIFAR10 images.

10 - *In L62 it is claimed that the memory complexity increases linearly on n_t . Should this be $O(n_t^2)$?* The memory
11 complexity is linear in n_t , as each potential is represented in memory by n_t points and weights. We will clarify.

12 - *I'm not entirely convinced that using "number of computations" in the x-axis makes sense for Figs 1, 3 etc.* We measure
13 the number of computations needed to obtain a *first estimate* of the OT potentials, which is roughlyly proportional to
14 wall-clock time (see answer to **R2**). It is of course higher for batch method than online method. Our intent in Fig. 3 is
15 to show that online Sinkhorn efficiently warms up OT computation. We will clarify.

16 - *I could not find a discussion or details on how the learning rate η_n is chosen in practice.* We give practical
17 recommendation regarding step-sizes and batch-sizes in Appendix B.3, and in particular Table 1. In experiments, we
18 found that setting $n(t) \propto (1 + 0.1t)^{1/2}$, and $\eta_t = 1$ work best, although the range of usable exponents is rather wide.
19 We will present Table 1 in the main text for clarification. See also App. C for details on hyper-parameters.

20 - *Further references.* We thank the reviewer for his insightful references on streaming method for EMD estimation, that
21 we will discuss in the related work section. In the batch or online setting, regularization permits a faster estimation of
22 OT distances, relying only on matrix-vector products. [39] fixes a spatial grid for the estimated barycenter, unrelated to
23 observed samples, while we define potentials based on observed samples. We will discuss this in details.

24 **R2.** - *How is convergence affected by [...] the distributions involved.* We have tried to give more insight on this aspect
25 in Appendix C. As predicted in the analysis, online Sinkhorn converges more slowly for lower ϵ (or equivalently, less
26 regular C , Fig. 5). For Gaussians distributions, online Sinkhorn outperforms batch Sinkhorn in all cases (Fig. 7).

27 - *Could the authors comment on actual runtime?* With proper GPU implementation of online Sinkhorn (using the
28 *pyKeops* library), the C-transform wall-clock time is indeed roughly in $O(n(t)^2)$. We have compared online Sinkhorn
29 to batch Sinkhorn in term of wall-clock time, and found similar curves as reported in the paper, using batch-sizes larger
30 than 1000. Batch Sinkhorn remains faster for small problems ($N < 10^4$), for which C can be precomputed and held in
31 GPU memory. We will add wall-clock time experiments to the appendix.

32 - *Confusion l.288-289.* f and g are fit until C is formed, and we then run batch Sinkhorn. We will clarify.

33 - *"It behaves like $\exp(1/\epsilon)$." Any intuition as to whether this can be improved?* This is a difficult open question.
34 Entropic regularization improves the sample complexity of Optimal Transport, going from a rate in $O(n^{-1/d})$ to a rate
35 in $O(1/\sqrt{n})$. This improvement is not free: as noted in [18], the constants before these rates explodes as $\exp(1/\epsilon)$ as
36 $\epsilon \rightarrow 0$. Both [18] and our work relies on the contractance modulus of the soft C-transform, hence the similar conclusion.

37 **R3.** - *The per-iteration cost of the classic Sinkhorn algorithm can be reduced [...].* We have been too elusive on this
38 aspect. Online Sinkhorn proves most useful in the case where the pairwise cost matrix must be computed on the fly due
39 to memory constraints (and serves as a sound warmup otherwise). We will recall and discuss this observation.

40 - *I find the Prop. 4 surprising.* In Prop. 4, we assume that $\iota > 0$, and therefore that the batch-size goes to infinity. This
41 is sufficient to ensure convergence, as the variance terms introduced by sampling are summable. For fixed batch-sizes,
42 convergence cannot be guaranteed, due to the fact that $\sum_t \frac{1}{t}$ is not summable. The proof of Prop. 4 established a
43 classical recursion between error terms, and requires $\iota > 0$ to conclude. We will discuss Prop. 4 more thoroughly.

44 **R4.** - *Experiments are performed in only simple cases.* This is a limitation of our work. The lack of gold-standard
45 for estimating continuous OT distances makes it hard to evaluate our method on hard settings, as we are forced to
46 approximate this gold-standard with very long runs of Sinkhorn algorithm. See also answer to **R1**.

47 - *Soft C-transform.* This term refers to Eq. (3). We will clarify.

48 - *Related work.* Bercu and Bigot's work is indeed relevant. It tackles the simpler problem of semi-discrete OT, that
49 rewrites as a expected risk-minimization problem. A single finite dimensional potential must be estimated, which
50 can be done through gradient descent. We will refer to Mena and Weed's refined sample complexities. The E-step of
51 Sinkhorn-EM could be implemented using online Sinkhorn, with potential gain from warm-starting.