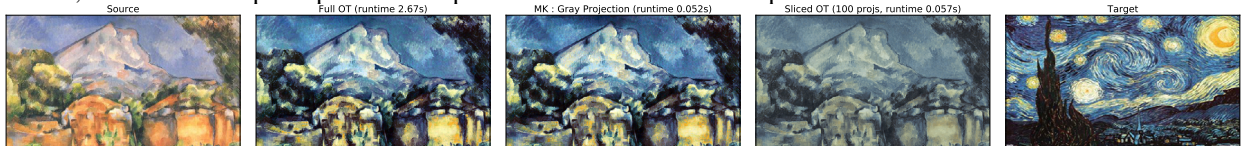


1 We thank reviewers for carefully reading our paper. We answer their questions
 2 below, but provide first two updates that are directly related to their remarks.
 3 **► Subspace Selection:** Alg. 1 from the paper was motivated by Prop. 6 (l.
 4 263). After benchmarking it carefully, we now believe it is not competitive
 5 with a projected gradient descent (PGD) on the basis vectors \mathbf{V} of E (see
 6 right). The projection of \mathbf{V} onto the set of unitary matrices is the unitary
 7 matrix in the polar decomposition of \mathbf{V} . The complexity per iteration is that
 8 of computing MK and the polar decomposition. We initialize $V = \text{Polar}(\mathbf{A}\mathbf{B})$
 9 because this is the optimal solution when \mathbf{A}, \mathbf{B} are co-diagonalizable. We
 10 tested this new algo. in the synthetic noisy setting (p.7), Fig.1 below. The
 11 PGD improves on the fixed direction (canonical basis) approach when $k < 4$, and remains competitive when $k \geq 4$.

Algorithm 1 MK Projected GD

Input: $\mathbf{A}, \mathbf{B} \in \text{PSD}, k \in \llbracket 1, d \rrbracket, \eta$
 $\mathbf{V} \leftarrow \text{Polar}(\mathbf{A}\mathbf{B})$
while not converged **do**
 $\mathcal{L} \leftarrow \text{MK}(\mathbf{V}^\top \mathbf{A} \mathbf{V}, \mathbf{V}^\top \mathbf{B} \mathbf{V}; k)$
 $\mathbf{V} \leftarrow \mathbf{V} - \eta \nabla_{\mathbf{V}} \mathcal{L}$
 $\mathbf{V} \leftarrow \text{Polar}(\mathbf{V})$
end while
Output: $E = \text{Span}\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$

12 **► Map visualization using color transfer:** All reviewers have pointed out that experiments in the paper did not illustrate
 13 the lifted transport maps/plans, but focused instead on distances. We experimented MK maps on color transfer, an
 14 illustrative task to visualize maps' properties. In the MK setting, we project images on the 1D space of grayscale images,
 15 relying on sorting-based algos for 1D-OT, before solving small 2D-OT problems on the corresponding *disintegrations*.
 16 We compare runtimes and visual results with vanilla OT and sliced OT below. MK results are visually very similar to
 full OT, with a $\sim \times 50$ speedup that is comparable to sliced OT. We will provide other illustrations.



17
 18 **Reviewer #1:** **► algo in terms of**
 19 **optimality, convergence, runtime,**
 20 **etc.** The runtime involves a complex-
 21 plexity per iteration equal to computing
 22 the polar decomp. and MK
 23 distance + gradient. Because the
 24 problem is non-convex we will
 25 stick to empirical evaluations and
 26 improve the presentation (p.7), as
 27 in Fig. 1 (right). **► applications do not seem to be terribly important [...]** more popular ones. Agreed. Color transfer
 28 was added as an illustrative example. We are now looking into applications to domain adaptation and biological datasets
 29 (Waddington-OT). **► experiments section [...]** a little confusing. We will add more context. The main purpose of the
 30 FID exp. (p.8) is to use data widely handled as samples from Gaussians. We show that even with a relatively small
 31 number of samples to estimate the covariance matrices, MK on the principal components has a stable behavior.

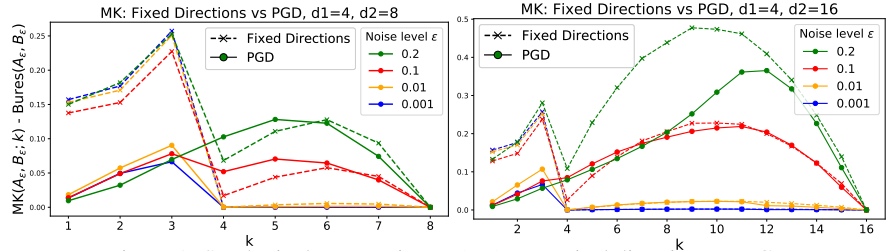


Figure 1: Synthetic data experiment (p.7): canonical directions vs PGD

32 **Reviewer #2:** **► 1:** E is indeed introduced later, l.62. We will fix this.
 33 **► 2: PCA with a (random) subset.** This counter-example is to show
 34 that the stability of MK is dependent on the chosen subspace. Permuting
 35 the principal directions is an adversarial setting used to showcase
 36 this. **► 3: what does 'underestimated' mean [...]** covariance matrices
 37 estimated [...] decent quality? In the setting of FID (p.8), $p = 2048$
 38 and we used $n = 2050$. Fig. 2 (right) shows the convergence of
 39 sample to full (on all 200K data points) covariance matrices in Bures
 40 and L2 distance (averaged over 20 sample matrices). At $n = 2050$
 41 the sample covariance matrices are close to having converged but
 42 not quite. However, the MK distance on the principal components
 43 is robust to the small amount of noise thus induced. We are glad to
 44 include this point in the discussion. **► 4: [...]** value of d_2 [...] role of d_2 in this context? As per the caption in the paper
 45 (Fig.4, p.7) $d_1 = 4$, top row is $d_2 = 8$ and bottom row $d_2 = 16$. We will make this more explicit. As d_2 increases, the
 46 MK distance for $d_1 \leq k \leq d_2$ increases as more noise is fitted by the transport map on the projection subspace.

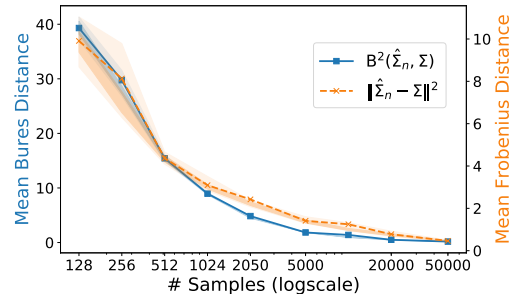


Figure 2: Mean distances from sample matrices to full covariance matrix (FID setting, p.8)

47 **Reviewer #3:** **► experimental verification of that chapter's suggestions, especially of Algo. 1?** Semantic mediation
 48 (p.8) is an example of using MK with prescribed directions (l.242-249), and FID experiments (p.8) of using principal
 49 components. We have added a verification of the new PGD algo in the experiment on noisy data (Fig. 1). **► Experiments**
 50 **with synthetic data seems informative, but semantic mediation etc are not convincing.** We added more semantic
 51 mediation examples. We are considering domain adaptation and biological datasets. **► Experiments on real data, and**
 52 **some more attention to selection of subspace E (experimentally).** Agreed. The PGD approach is a first step in that
 53 direction (Fig. 1). We will also try it first in color transfer, domain adaptation and in biology (Waddington-OT).