

1 **R1** We appreciate your very constructive comments. Good point. Some new experiments comparing COPT distance
2 to SOTA GW distance: With 6 classes of synthetic graphs, 18 queries on 60 graphs, with nearest-neighbors-based
3 classifier, COPT achieves acc. 83.5 ± 3.4 , and GW 62.2 ± 6.6 . This is repeated 20 times. Will expand and include these.
4 (The existing vectorize-and-retrieve experiments were designed to be highly practical on GPUs, leveraging the benefits
5 of COPT during the vectorization preprocessing.) • Will update paper to emphasize distance aspect as you suggest.
6 Algorithm#1 indeed already computes the COPT distance as the objective during optimization. • Indeed interesting
7 to vary optimizer. To sketch a 20-node graph down to 5, currently each iteration takes 3.7 ± 0.05 ms, after updating
8 to LBFGS, each iteration takes 128 ± 3 ms. Likely because LBFGS is a second order method, and may evaluate the
9 forward pass multiple times per iteration. Will try others. • There are in fact interesting relations to graph edit distance:
10 in COPT, the matrix P allows *continuous* (e.g. partial) flow between vertices, while edit distance involves matching the
11 vertices *one-to-one* (plus insertions & deletions). Edit distance gives the *same cost* to delete an arbitrary edge, while
12 COPT distance is higher if the edges deleted were important to the graph spectrum (i.e. edges connecting distinct
13 clusters). A continuous relaxation of edit distance would resemble [LFF+15] more than COPT for this reason. • Added
14 the fact the measure of variance is standard deviation, number of times 5.2 is repeated, will update broader impact &
15 add details to README. Thanks for these. [LFF+15: Lyzinski et al, Graph matching: Relax at your own risk. 2015.]

16 **R2** We thank you for your points. Indeed Batson-Spielman-Srivastava and related works from TCS are fascinating.
17 The primary reason for choosing the included literature was an alignment in goals. COPT was designed to reduce the
18 number of *vertices*, producing *dense* Laplacians; BSS and related works reduce the number of *edges*, producing *sparse*
19 Laplacians. One benefit of dense data is that they are well-suited for efficient processing on GPUs, such as in many
20 ML applications. • BSS focuses on proving theoretical guarantees on constructing graph sparsifiers, it uses a *discrete*
21 optimization method where one edge is added to the sparse graph per step (though weights are chosen continuously);
22 COPT uses a *continuous* optimization method where every variable can change at every time step. Though discrete
23 optimization can more readily produce provable guarantees, continuous optimization is often preferred in practice, e.g.
24 continuous relaxations of discrete problems such as LP relaxation. • Another effect of the proof-driven construction
25 is that the constructed sparsifier may not be useful for applications such as graph classification. E.g. BSS sparsifies
26 complete graphs to sparse expander graphs, which have intricate structures not present in original graph. Attempts to
27 classify the sparsifier could see this structure instead of properties of the original graph. COPT empirically reduces
28 complete graphs to complete or nearly-complete graphs. • We will add these. • Similarly, the approach (e.g. in
29 Jambulapati-Sidford) of defining *multiple* randomized sketches and then taking the median to obtain a good Laplacian
30 approximation may not be practical for applications. • Good call on motivating §3.2 with downstream applications,
31 *metric* approximations are good for applications where shortest paths between nodes are important, e.g. classifying road
32 networks, classifying graphs arising from physical objects where edge lengths carry geometric information. *Spectral*
33 approximations are useful for when the number of paths matter more, e.g. graph partitioning, where one aims to
34 minimize the number of edges cut (one way is to iteratively coarsen graphs and find cuts on smaller graphs). Similarly
35 COPT is applicable to graph clustering. • Will state nonconvexity earlier in the paper; comment on the theoretical fast
36 matrix multiplication methods vs practice (indeed the theoretical fastest method is not practical); add references as
37 suggested; comment on runtimes; and fix typos. • The only change to work with weighted graphs is to use the *weighted*
38 Laplacian instead of ordinary Laplacian. The proofs only use the fact that L_X is a symmetric PSD matrix. We brought
39 in the additional concept of weightings only when needed, but the potential for applications to weighted graphs is a
40 good point.

41 **R3** In Eq. (4), we define a loss function for a transport map between \mathbb{R}^X and \mathbb{R}^Y , as a generalization of the loss in GOT
42 from permutations to arbitrary matrices P . A priori this is just a formula, but in Lemma 3.1, we prove it is equal to a
43 distance in $\mathbb{R}^X \times \mathbb{R}^Y$, after embedding \mathbb{R}^X and \mathbb{R}^Y into $\mathbb{R}^{X \times Y}$, which we use to calculate it. In the special case when
44 P is a permutation, these embeddings have the same image, reducing the formula to optimal transport in just \mathbb{R}^X or
45 \mathbb{R}^Y as in GOT. • §3.2 compares COPT with GW, and discusses how COPT preserves paths-counting in relation to the
46 spectrum. • f is a free variable in the domain of integration, i.e. the space of functions on \mathbb{R} . Will specify this in the
47 paper. Thanks.

48 **R4** Thanks much for your comments. It is indeed useful to compare with the WL kernel. Some new experiments: On 60
49 synthetic graphs across 6 categories with 18 query graphs, using exact nearest-neighbor search with COPT as a distance:
50 acc. 83.7 ± 3.9 ; SVM with WL kernel: 80.5 ± 9.0 . This is repeated 20 times. • Choosing between the two depends on
51 the task and the dataset: For a large dataset of size N , we expect a WL kernel-based SVM to be more practical than
52 using linear scan with COPT distance. As an SVM can be trained during preprocessing, allowing *constant* time query
53 classification, whereas exact nearest-neighbor search requires time *linear* in N . • For small datasets, COPT distance
54 can be advantageous: with only three samples per class, using nearest-neighbor classification with COPT distance gives
55 accuracy 90.7 ± 5.2 , and the WL kernel-based SVM gives 77.8 ± 7.1 . • For tasks such as predicting *continuous* attributes,
56 an SVM with WL kernel can't readily be applied, but a nearest-neighbors classifier with COPT distance still can. • Will
57 expand on these and add to the paper.