1  We thank the reviewers for their time, their valuable and encouraging feedback, and their recommendations for
2  improvement. We remain confident that our work is of strong interest to the NeurIPS community and easily can
3  incorporate the suggested changes in a revision for the conference. Answers to specific comments appear below.

4  **Interpretability**   To address **R2**'s concerns about interpretability, we refer to Figure 1 in the paper, where we show an
5  example with two novels from the Gutenberg dataset. To interpret HOTT distance between a pair of documents, one
6  simply needs to look at the slice of the transport map corresponding to the dominant topics (typically 3-4 per document)
7  as we show in Figure 1. As a point of contrast, interpreting WMD would require investigating a transport map between
8  all unique words in a document pair (thousands of unique words for books).

9  **R1** requested clarifications about Figure 1. R1 correctly interpreted the percentages in the figure and pointed out that
10  topic titles are assigned by us, and not by any algorithm we use. We will improve the caption clarity accordingly.

11  **Theoretical contributions**   **R1** asked about the purpose of the RWMD-Hausdorff bound section. This section
12  strengthens the motivation for our work: WMD has been shown to be successful, however may be too slow in practice;
13  RWMD is a fast approximation empirically performing well in some cases, but we show that it may also be a very poor
14  choice of metric in practical scenarios. To give an extreme example, the documents [good, bad, bad, ..., bad] and [bad,
15  good, good, ..., good] are distance 0 from each other under RWMD. This motivates the study of alternative document
16  distance metrics utilizing word embeddings geometry, which we do by proposing HOTT. We will make this connection
17  more explicit in the paper.

18  **Technical details**   **R1** asked whether we can use *different cost metrics* to measure distance between topics. Our choice
19  of Wasserstein metric is motivated by the improvement seen in using the Earth Mover's (Word Mover's) Distance as a
20  document-to-document metric. Other distances between topics, i.e., Euclidean or cosine, do not exploit the geometry of
21  word embeddings. We conjecture (and can easily verify that) they will perform worse.

22  **R2** and **R3** had questions about the definition of *HOTT* on line 95. We will clarify this, since understanding this equation
23  is *fundamental* to the remainder of the paper. As we adopt a *hierarchical* approach, our documents are distributions
24  ($\vec{d}^i \in \Delta^{|T|}$) supported on topics ($\delta_{t_k}$, $k = 1, \ldots, K$ where *Dirac delta* ($\delta_x$) at $x$ is a probability distribution only
25  supported on the point $x$). Line 95 says that the distance between two documents is the Wasserstein distance between
26  their distributions over topics, where topics are also distributions, but over words, hence the ground metric is another
27  Wasserstein distance between topics represented as distributions over words. Thus the distance is *hierarchical*.

28  **Experiments and hyper-parameters**   To answer **R2**'s question on pruning for WMD, we refer to Figures 3 and 5
29  where WMD-T20 represents WMD truncated to the top 20 words: pruning heuristic *cannot* be efficiently applied to
30  WMD. While it helps with the run-time (see Table 1), it noticeably degrades the performance (see Figure 5).

31  **R3** asks how to interpret the *results* in Figure 4a. GloVe embeddings used in all of the experiments (including Figure 5)
32  are the *high quality* pre-trained 300d embeddings trained on 6 billion tokes, which can be downloaded online (please
33  see `main.py` file in the code). Figure 4 (a) quantifies sensitivity of different methods to *lower quality* word embeddings.
34  In particular, we trained 200d embeddings with word2vec algorithm using *only* documents of the Reuters dataset
35  (under 300k tokens). Figure 4 (a) shows that lower quality word embeddings significantly degrade performance of
36  the WMD-based methods. Our methods, on the contrary, maintain good performance because they are able to utilize
37  informative topic structure of the Reuters documents, which is independent of the word embeddings quality.

38  **R3** questions our choice of *20* words for WMD truncation. While the choice of 20 for WMD-T20 is somewhat arbitrary,
39  it is simply made to bring WMD complexity closer to HOTT and show that WMD *cannot* be made computationally
40  efficient using truncation without degrading its performance. WMD-T20 is already slower than HOTT (see Table 1) and
41  performs noticeably worse (see Figure 5); truncating it further will make the performance even worse, while truncating
42  less will quickly lead to impractical run-time, e.g., computing all pairwise WMD distances on the Gutenberg dataset
43  would take $\approx 178$ *days* on a single machine. We are happy to include a sensitivity analysis on the truncation of WMD
44  on one of the smaller datasets.

45  **R1** asked why we do not compare with TMD. Most importantly, the cubic complexity of the OT linear program remains
46  prohibitive for the number of topics used in TMD, i.e., from Table 1 in TMD paper it can be seen that number of
47  topics they use is only 3-4 times smaller than vocabulary size. We use 70 topics, i.e. over 100 fold vocabulary size
48  reduction on average across datasets. Quantitatively, Figure 3 of the TMD paper suggests that evaluating a kNN
49  classifier on the BBCSPORT (smallest dataset) takes 24h for WMD and 4h for TMD. First, the WMD implementation
50  we use takes 3-4min and, second, HOTT takes only about 40sec. We conclude that HOTT (and even a simply better
51  WMD implementation) is much faster than TMD. Discrepancy in the WMD speed may be due to authors of TMD not
52  fully utilizing sparsity of the documents when calling the linear program solver.