1 We are very grateful to the three reviewers for their careful reading of our paper, and for their insightful comments and
2 suggestions. Our responses are below with the reviewers' comments in italics.

3 • **Reviewer 1** speaks very positively about our paper tackling an *"important topic under a new perspective"*, and we
4 too anticipate that it will *"open up space and motivate new research"* in this area. In response to specific points:

5 – *"The paper needs more details on the experiments..."*Although we do give full details (of embedding models,
6 training data, embedding dimension, etc) in the Figure captions we appreciate that these might have gotten
7 somewhat lost; we will revise to clarify, and expand explanation in the main text.

8 – *"Give a clearer motivation for how embeddings were constrained in §4.1...the authors do not show why this*
9 *reduced set of solutions would be better than a random choice in the full set [nor give results for it, which] could*
10 *be presented in Table 1...Why not expand Table 1 to [include LSA, Word2Vec]".* As pointed out in Remark 3, the
11 constraints in §4.1 are related to the widely used Gram-Schmidt process, and are natural under the group-theoretic
12 formalism. We make no claim that this particular solution subset is optimal (or even at all good). We will amend
13 the text to make this clear, and also follow the excellent suggestion to add results corresponding to §4.1, as well as
14 Word2Vec, to Table 1.

15 – *"Do they optimize [$\Lambda$] on the same data D used later to get the scores?"* We did not use cross-validation for the
16 results in Table 1, because the purpose of these results is to understand how much $g$ can be increased without
17 changing $f$, as we say in the caption to indicate the "substantial scope for improving performance scores via an
18 appropriate choice of $\Lambda$." We will amend presentation to further emphasise this, and also expand our caution of
19 overfitting (line 250) to discuss specifically the need for cross-validation when producing predictive models.

20 – *"Fig 2 (a,b), the red lines do not match when $\alpha = 0$...".* Thank you for pointing this out (and conjecturing correctly
21 the reason!). We will amend.

22 – *"Fig 3, results for upper triangular...have better performance than diagonal ones...[Intuition?]...What are red*
23 *lines...original V results?"* Our intuition here is that the highest performing instances do (slightly) better for the
24 upper triangular case than the diagonal case since the extra degrees of freedom give greater scope for a solution
25 with large $g$; and that the left tail is long because with |N(0,1)| random diagonal elements there is a substantial
26 chance of small elements that essentially wipe out whole rows of $V$. We will add this discussion in the revision.
27 The red lines are indeed the results for the original $V^*$, which we will now state in the caption.

28 – *Various typos:* Thank you for pointing these out, all of which we will fix.

29 • We share **Reviewer 2**'s craving in answering the "meta question" of *why* embeddings work. Given the impressive
30 success of numerous embedding methods being used and compared, we investigate a related, albeit modest, meta
31 question 'for a fixed evaluation function, can we identify properties/issues of objective functions that drive the
32 apparent disparity in performances of embedding methods?'. To our knowledge we are the first to do this.

33 – *"Previous methodologies...implicity address [non-identifiability] already, as they place implicit constraints (like*
34 *requiring U and V to be...symmetric)...[hence] 'so what?'"* It is not clear what implicit constraints the reviewer is
35 referring to in general. But U and V are non-square so cannot be symmetric. Perhaps "requiring U and V to be
36 identical on account of X being symmetric" was intended, as in §4.1.1. But this strategy is limited to situations
37 with X symmetric. As we discuss in §4.1.1, the remedy suggested by the authors of GloVe is ad hoc (much less
38 natural than the alternative we suggest) and leads to embeddings that are not even optimal with respect to the
39 model's own objective!
40 Our view is that to even begin to address the reviewer's "meta question" the first step is to understand embedding
41 methods' objective functions and the properties of embeddings that result. We feel it is paramount to understand
42 the impact of the identifiability issue first, instead of adding to the growing list of embedding methods by throwing
43 out new objective functions.

44 – Suggestions: (i) *"looking at a carefully selected larger space, [e.g. Artetxe et al 2018]".* The linear transform in
45 their paper (in our notation $V = \Lambda^\alpha Q V^*$) is just identified with a one-dimensional subset of $GL(d)$, and since $g$ is
46 invariant to any $Q \in O(d)$ their approach is covered by our Proposition 1. (ii) *"modifying the training objective".*
47 Addressed previously. (iii) *"embedding methodologies that are theoretically more sound and performs better..."* It
48 is unclear what the reviewer means by 'theoretically more sound'. Perhaps it's related to the suggestion in (ii)?

49 • **Reviewer 3** makes some helpful suggestions about expanding §4.2, in particular to include *"NON-WORD-tasks (not*
50 *just word similarity...using [inner product])".* We emphasise it is not just similarity but analogy tasks too for which $g$
51 involves $V$ only via the inner product of its columns, though we acknowledge the lack of analogy results, which some
52 readers may find more interesting, in §4 and will include some in the revised version. One point of clarification: the
53 reviewer notes that we focus on the *"rather special linear case [of] LSA"*, though we do so only in §2 as a particular
54 example that aids exposition of the identifiability issue; the paper's key results in §3 and §4 are appropriate to any
55 model of the very general form (1), including GloVe and word2vec.