

1 We would first like to thank the reviewers for their thoughtful commentary and constructive feedback, particularly
2 given the constraints imposed by the COVID-19 pandemic. We are happy to see that the reviewers agree that studying
3 individual differences is an important problem for neuroimaging, and overall appear to lean towards acceptance.

4 The reviewers do have a number of comments and questions relating to baselines, the importance of characterizing
5 stimulus variation, and the number of available datasets for experiments. Many of these comments are actionable and
6 addressable by camera ready, as we discuss below.

7 **R2: Comparison to MN-SRM, tensor decompositions.** We appreciate R2's helpful reference to Shvartsman et
8 al. (2018), which appears to consider spatial variation more explicitly than the original SRM. We were not able to find
9 an open-source implementation of the MN-SRM and as a result were not able to perform a comparison in time for this
10 response. We will attempt to implement this method ourselves and include a comparison in the camera-ready. We will
11 also evaluate whether CP/Tucker tensor decomposition methods could serve as additional baselines, as suggested by R3.

12 **R3+R5: PCA baseline.** R3 notes that PCA is not a strong baseline. We agree and it is not intended as such; we will
13 emphasize this more clearly. We appreciate the suggestion from R3 and R5 to perform time-averaging before doing
14 PCA to improve it as a baseline. We did so and found that this did not result in qualitatively different embeddings from
15 the non-averaged analysis in appendix Section A.2. We will update the manuscript and the figures to reflect this change.

16 **R3: How important are the nonlinearities?** We have followed up and trained a version of NTFA in which we
17 replace dense networks with a single linear layer. On our datasets, this results in equivalent reconstruction performance,
18 at the cost of the inferred embeddings looking meaningless (the participant embeddings collapse into the Gaussian prior,
19 and stimulus embeddings lack the interpretable pattern seen in Figure 3).

20 **R5: Cross-validation in MVPA analysis.** R5 writes, "feature selection was performed within the cross-validation (as
21 it should be), but NTFA was performed outside." We will clarify that for MVPA, we treat NTFA as an unsupervised
22 feature extractor before any supervised training is performed. This means there is no supervised feature selection. We
23 believe that unsupervised training on all data (which was done as a computational shortcut) is likely not problematic.
24 To verify this, we will re-train NTFA independently for each fold and re-run our analysis.

25 **R2: Notation and dataset tables.** We appreciate these suggestions and will include tables that summarize notation, as
26 well as tables with summary statistics for each dataset.

27 **R2: SRM vs TFA.** We agree that SRM and TFA as methods have different intended cases (functional alignment and
28 connectivity respectively) and will explain this more clearly in the text. We will also cite Cai et al. (2020).

29 **R3+R5 Why examine stimulus variation?** In standard neuroimaging analysis, researchers designate categories of
30 stimuli prior to the experiment and model the BOLD signal with a Gaussian (or similar location-scale) distribution.
31 This approach assumes that variation among stimuli can be treated as noisy deviations from a single mode. However,
32 the assumption that experimenter-designated categories are optimal has been challenged. Cognitive and translational
33 neuroscience suggests that the same category of stimulus-induced state (e.g. fear) or phenotypic trait (e.g. depression),
34 may involve multiple distinct neural pathways. Rather than assume that experimenter-designated categories are correct,
35 NTFA enables researchers to test this assumption. If the stimuli group together into categories naturally, as we found
36 in some of our experiments, then the categorical assumption may be considered well-justified. Where significant
37 intra-category variance appears, NTFA enables researchers to then delve deeper into why stimuli show divergent effects
38 and how that impacts overall conclusions from the findings.

39 **R2+R3+R5: Limitations of NTFA and our present experiments.** The three reviewers point out that more interesting
40 evaluation results could be obtained by applying NTFA to further datasets, and that interpretation of participant
41 embeddings has been left to future work. We agree that NTFA is a first step in a longer research program. Participant
42 embeddings in the current model are difficult to interpret because they serve a dual role: they capture both the spatial
43 alignment of latent factors within voxel-space and the per-participant variations in BOLD response to stimuli. The main
44 factor that limits our ability to improve the current model is the availability of datasets that are suitable to the study
45 of individual differences. We are in the process of performing a more comprehensive version of the ThreadVids pilot
46 study, which we hope will allow us to improve upon our current results in future work. We will add discussion to this
47 effect to the manuscript.