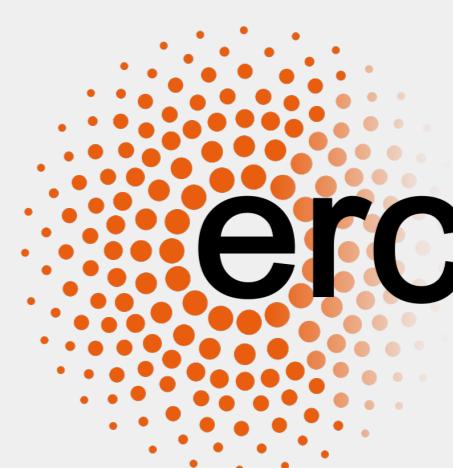
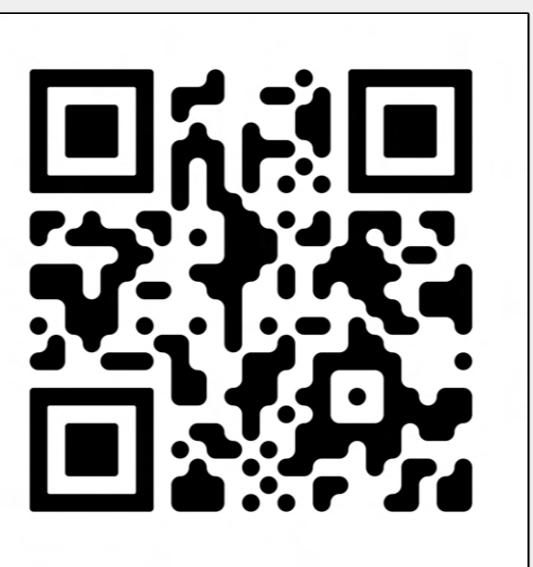


# The role of semantics in similarity judgements of scene stimuli

Katerina Marie Simkova<sup>a</sup>, Jasper van den Bosch<sup>a</sup>, Ian Charest<sup>a, b</sup>

<sup>a</sup>School of Psychology, University of Birmingham, <sup>b</sup>cerebrUM, Département de Psychologie, Université de Montréal



contact: k.m.simkova@bham.ac.uk

## Introduction

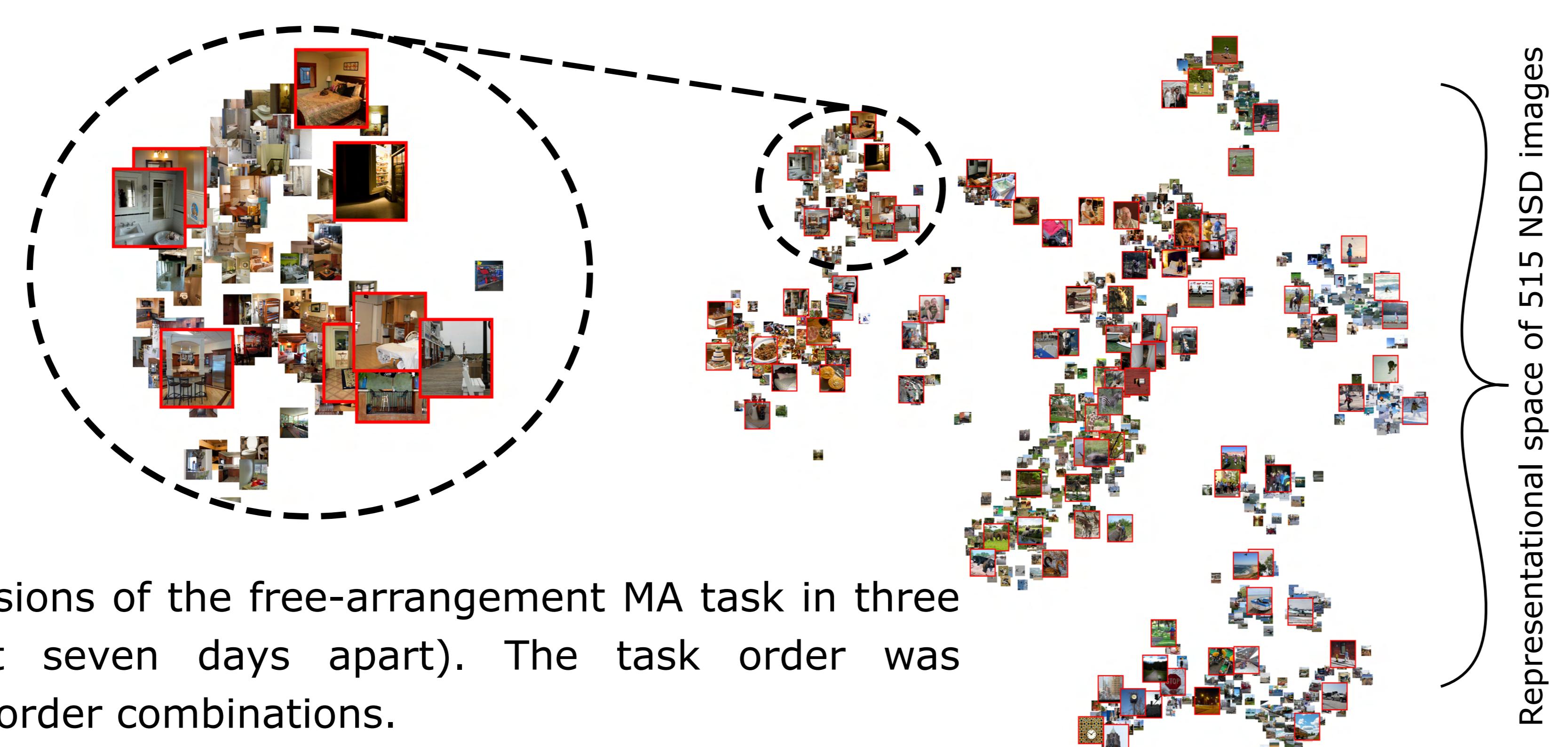
Classic accounts of human vision suggest functional segregation in two distinct cortical pathways: a “where” and a “what” pathway<sup>1, 2</sup>, a dorsal stream specialised in spatial information (“where”) and a ventral stream for category<sup>3–6</sup> or conceptual information<sup>7</sup>. The ventral stream could therefore be seen as a distributed system where overlapping feature maps<sup>8</sup> encode specific dimensions about particular objects<sup>9, 10</sup>. This view might stem from simple experimental paradigms involving single objects or simplified stimuli, but it is not clear whether this view accurately generalises to naturalistic viewing situations with complex scenes containing multiple objects and concepts.

Here we set out to investigate how human observers perceive semantic information that is being communicated via visual vs lexical stimuli using the multiple arrangements (MA) task: we isolated the semantic information by adapting the MA method into three different task formats - arrangements of images, hidden images, and captions - and explored how the similarity judgements relate to state-of-the-art deep neural networks and brain data.

## Methods

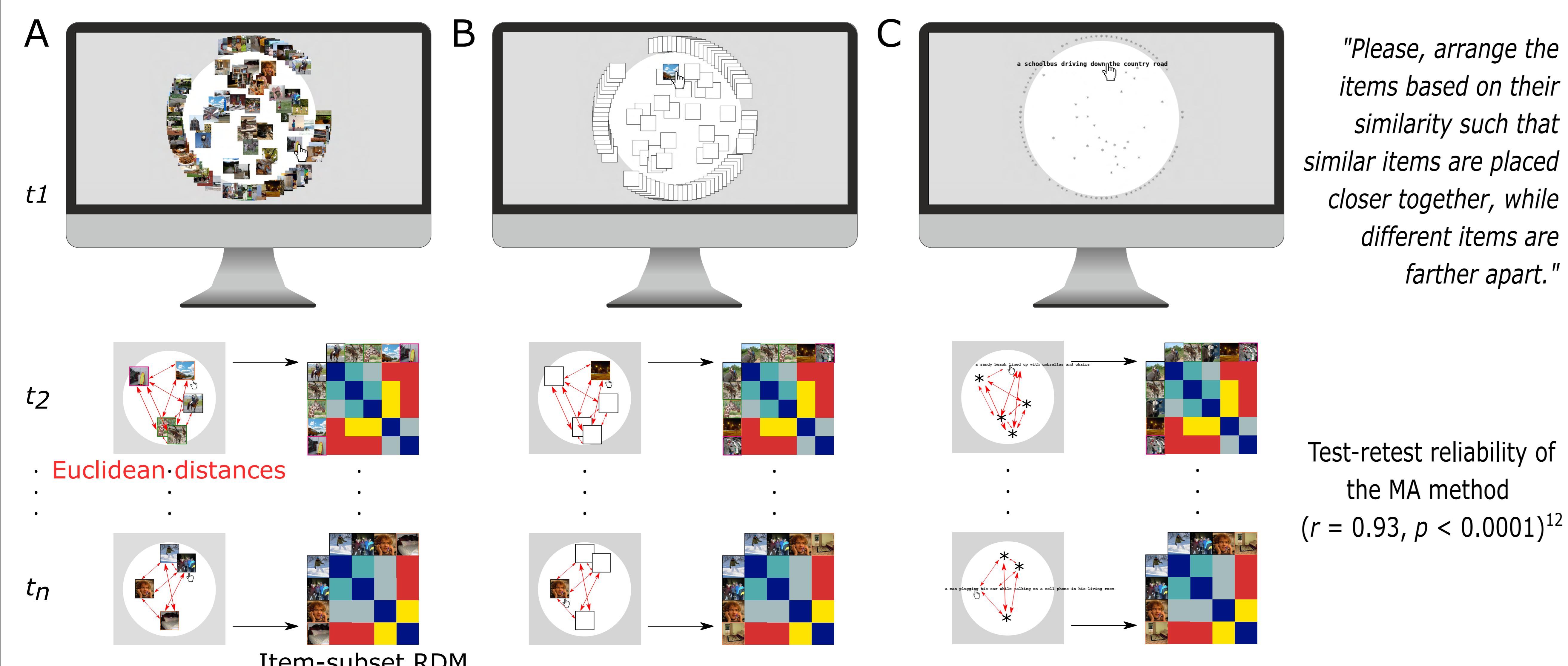
### Stimulus set

100 images (in red boxes) were sampled from the Natural Scenes Dataset (NSD)<sup>11</sup> so that they span the full representational space (as characterised from GUSE) of the NSD shared 515 images.



### Multiple arrangements task

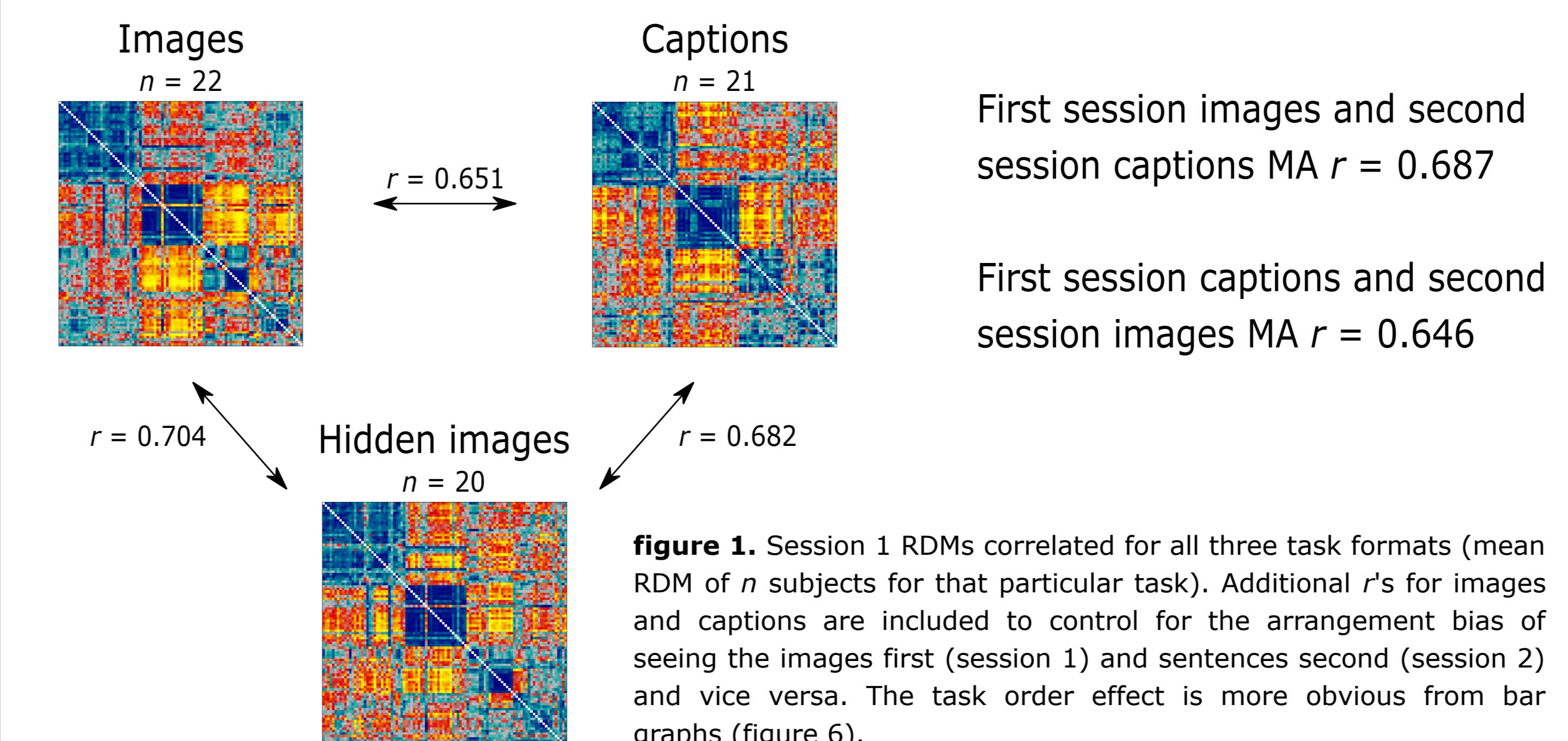
Participants ( $n = 63$ ) completed three versions of the free-arrangement MA task in three lab-based sessions (scheduled at least seven days apart). The task order was counterbalanced yielding six possible task order combinations.



The relative Euclidean distances between items were iteratively adjusted by the weighted-averaging approach: subsets of stimuli carrying the weakest weighted evidence were presented in subsequent n trials - reducing potential placement errors and yielding the final subject-level task-specific RDM. Trials terminated once 45 minutes elapsed.

## Results

### Task order effects

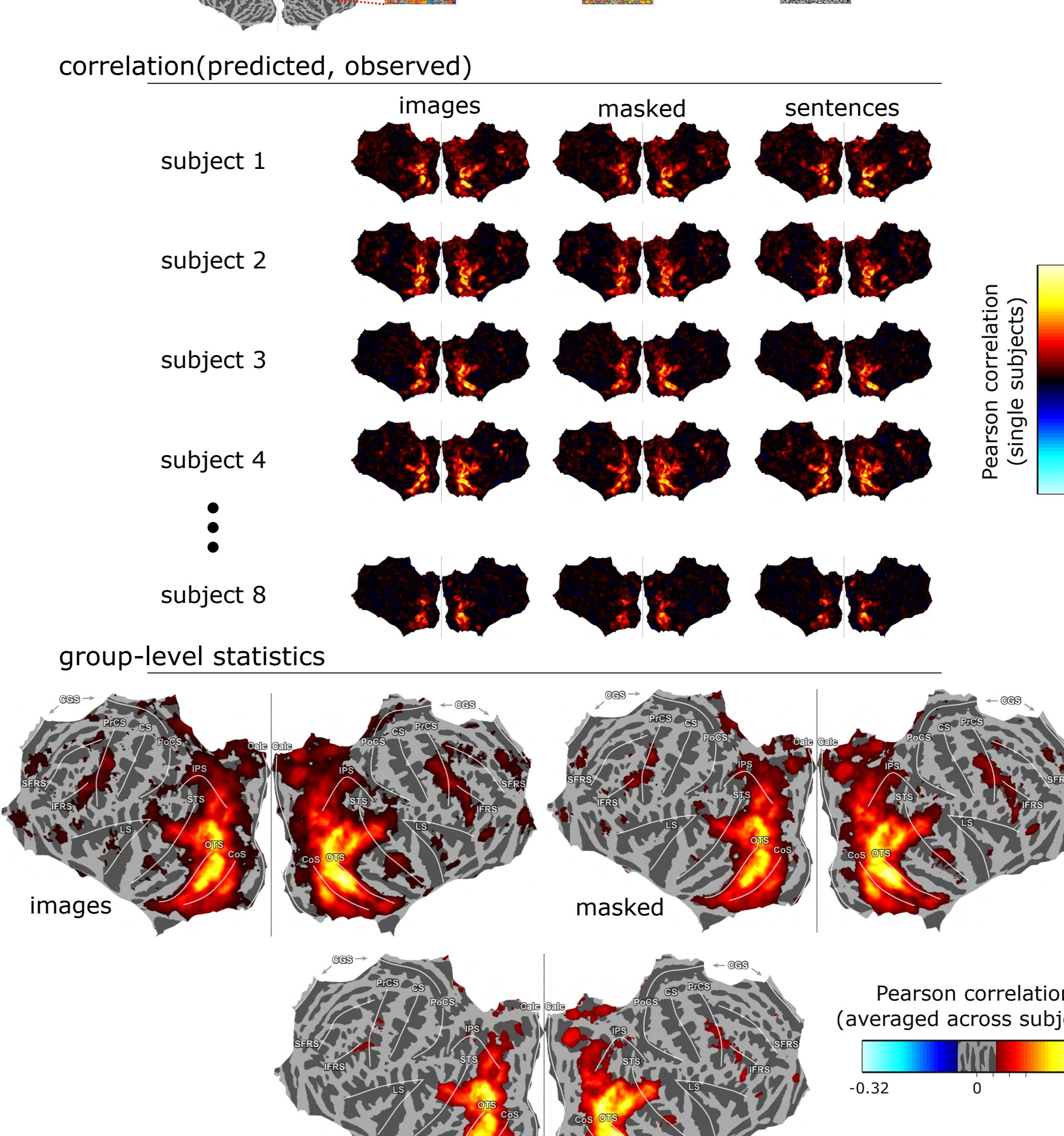


**figure 1.** Session 1 RDMs correlated for all three task formats (mean RDM of  $n$  subjects for that particular task). Additional  $r$ 's for images and captions are included to control for the arrangement bias of seeing the images first (session 1) and sentences second (session 2) and vice versa. The task order effect is more obvious from bar graphs (figure 6).

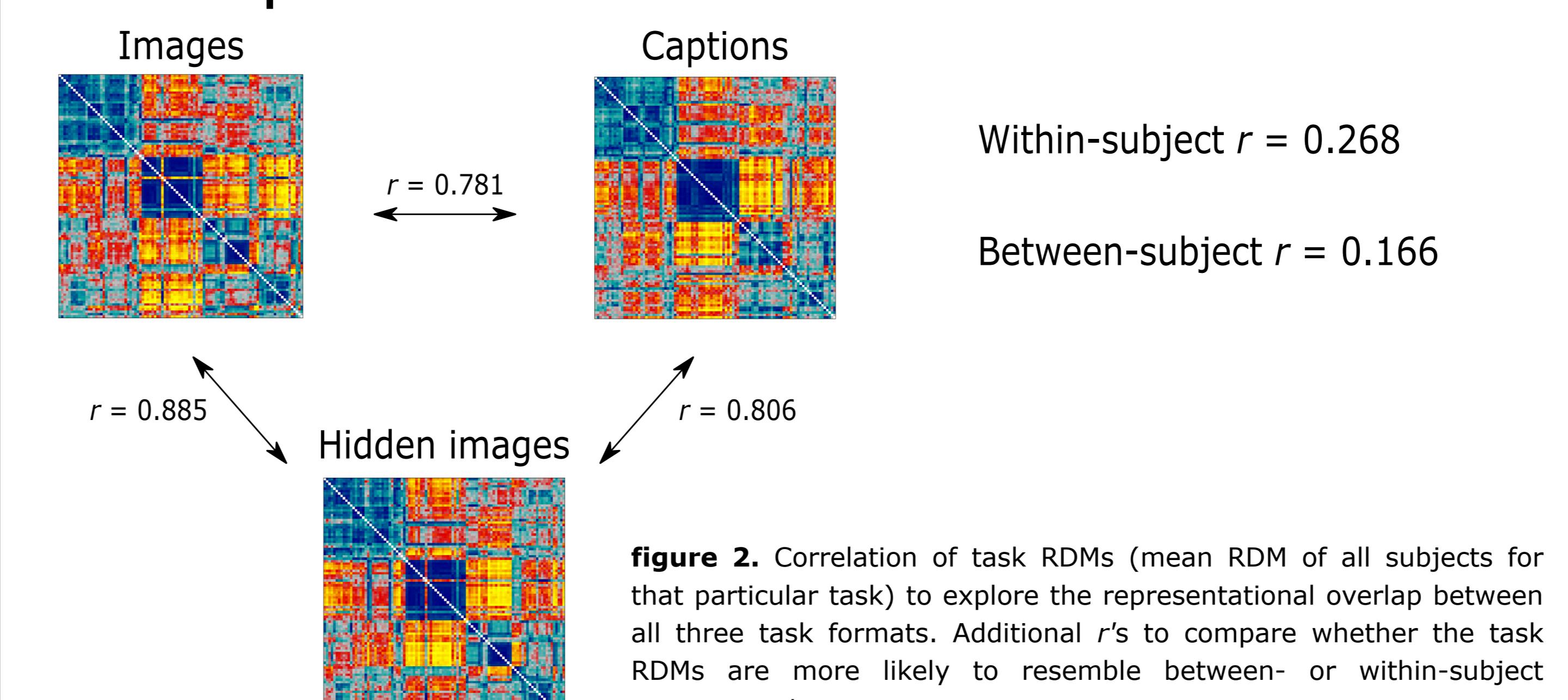
### fMRI searchlight RSA

fit NSD participants as a linear combination of behavioural subjects  
subjects 1 2 3 4 5 6 63  
MA task RDMs  
weights  $W_1, W_2, W_3, W_4, W_5, W_6, W_{63}$   
cross-validated non-negative least-squares

$$\text{searchlight RDM} = \text{pred. RDM} + \text{error}$$

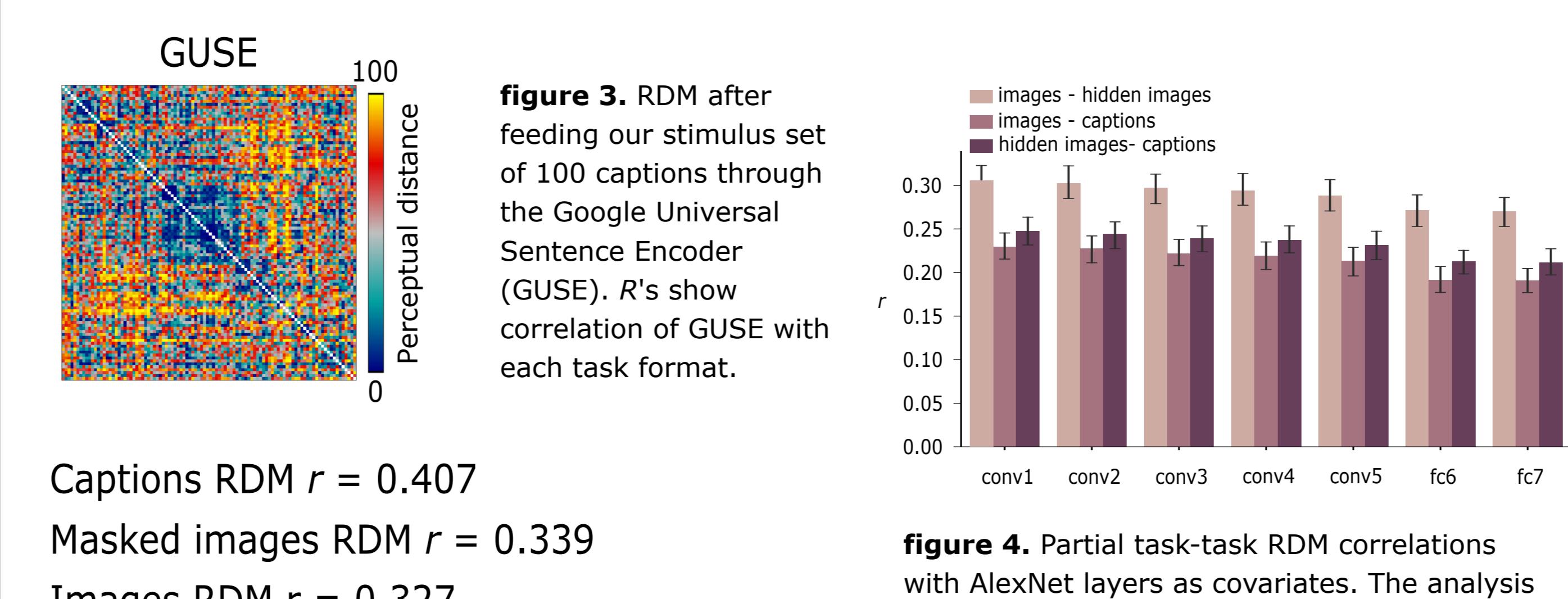


### Overlap between tasks

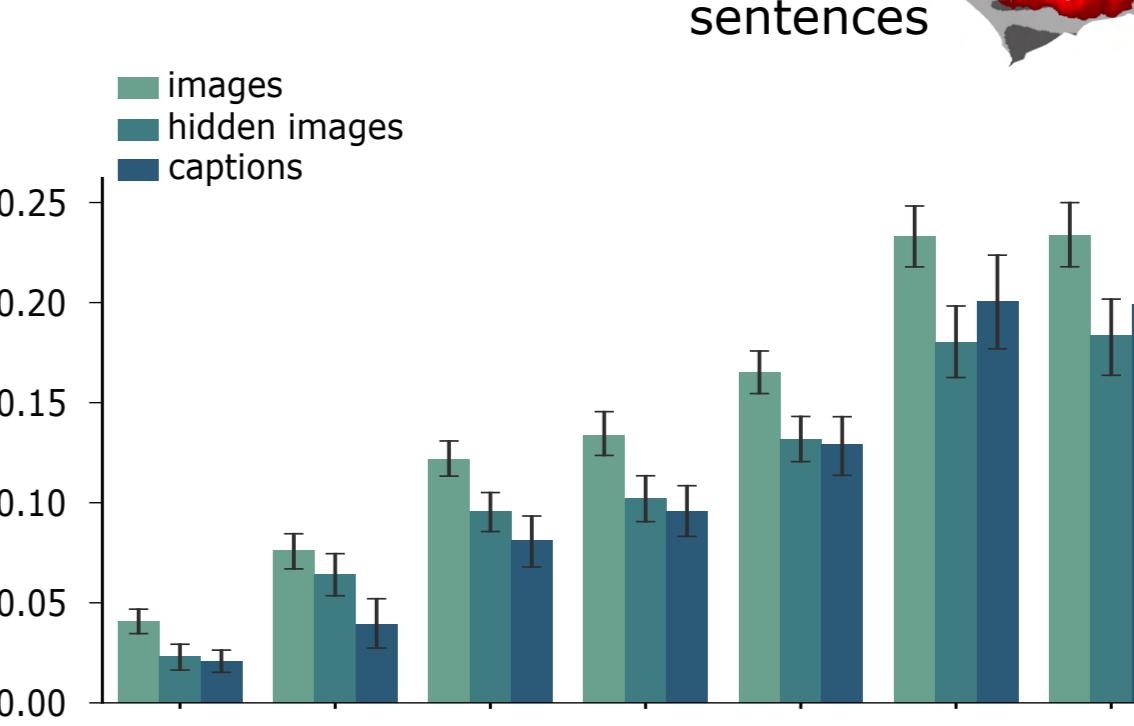


**figure 2.** Correlation of task RDMs (mean RDM of all subjects for that particular task) to explore the representational overlap between all three task formats. Additional  $r$ 's to compare whether the task RDMs are more likely to resemble between- or within-subject arrangements.

### Semantic features vs visual features

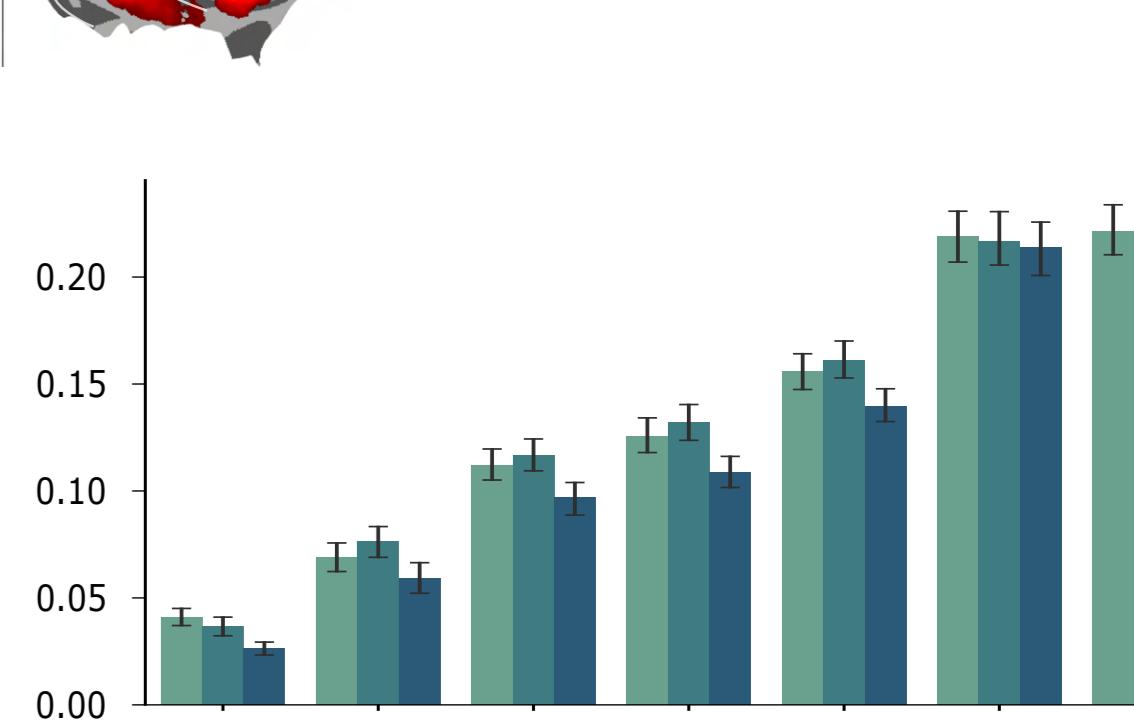


**figure 3.** RDM after feeding our stimulus set of 100 captions through the Google Universal Sentence Encoder (GUSE).  $r$ 's show correlation of GUSE with each task format.



**figure 4.**

Task RDM correlations with AlexNet layers based on first session only.



**figure 5.**

Task RDM correlations with AlexNet layers based on all sessions.

## Discussion

Here we investigated to what level human observers rely on purely visual versus semantic information when performing similarity judgements. Although the task order effects in the representational overlap is evident, analyses based on first session RDMs still reveal a significant level of correspondence between the tasks. The partial task-task RDM correlation with AlexNet layers show that the arrangements are not entirely dependent on visual features. This suggests that MA is an efficient method in capturing semantic level representations. Moreover, the fMRI searchlight analyses, using a linear combination of behavioural participants to predict NSD participants, revealed strikingly similar brain regions across the three MA tasks. Altogether this suggests that the representations in the visual ventral stream share a representational format with visual and lexical based similarity judgements.

## References

1. Sengpiel, L., & Hasley, J. V. (1991). “What” and “where” in the human brain. *Curr Opin Neurobiol*.
2. Goodale, M. A., & Milner, A. D. (1992). Separate visual pathways for perception and action. *Trends Neurosci*.
3. Kanwisher, N., et al. (1997). The fusiform Face Area: A module in human extrastriate cortex specialized for face perception. *J. Neurosci*.
4. Clarke, A., & Tyler, L. K. (2015). Understanding what we see: how we derive meaning from vision. *TICS*.
5. Hayby, J. V., et al. (2001). Distributed and overlapping representations of faces and objects in ventral temporal cortex. *Science*.
6. Tsoy, A., et al. (2010). The spatiotemporal neural dynamics underlying perceived similarity for real-world objects. *Neuroimage*.
7. Roberts, M. N., et al. (2020). Revealing the multidimensional mental representations of natural objects underlying human similarity judgments. *Nat Hum Behav*.
8. Allen, E. J., et al. (2022). A massive 7-t MRI dataset to bridge cognitive neuroscience and artificial intelligence. *Nature Neuroscience*.
9. Kriegeskorte, N., & Hurl, M. (2012). Inverse MDS: inferring dissimilarity structure from multiple item arrangements. *Front Psychol*.