Decoding topic vectors during memory encoding and retrieval

Jeremy R. Manning, David M. Blei, and Kenneth A. Norman

ing a collection of documents.1





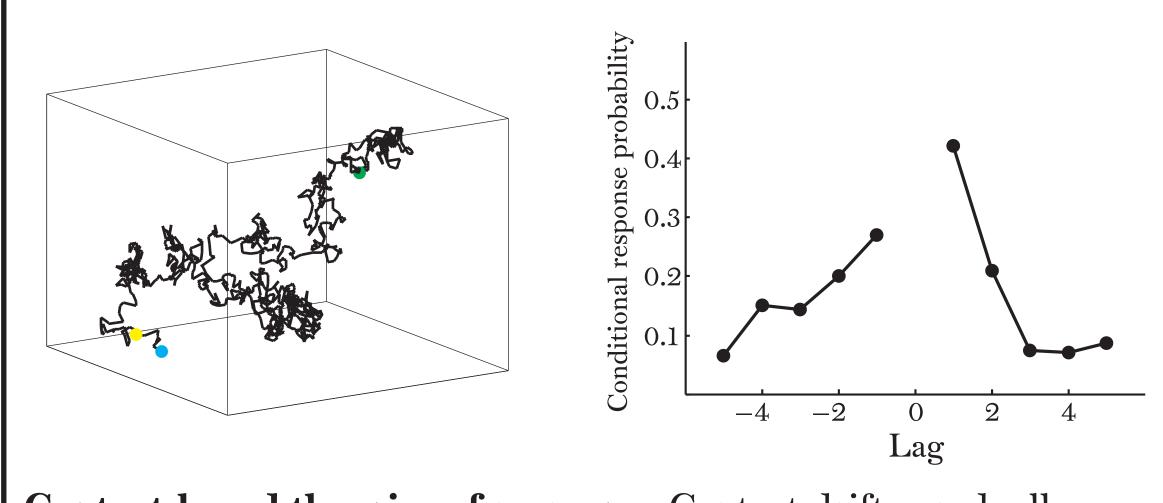
This research was supported by the NSF/NIH Collaborative Research in Computational Neuroscience Program, grant number NSF IIS-10009542

Princeton University

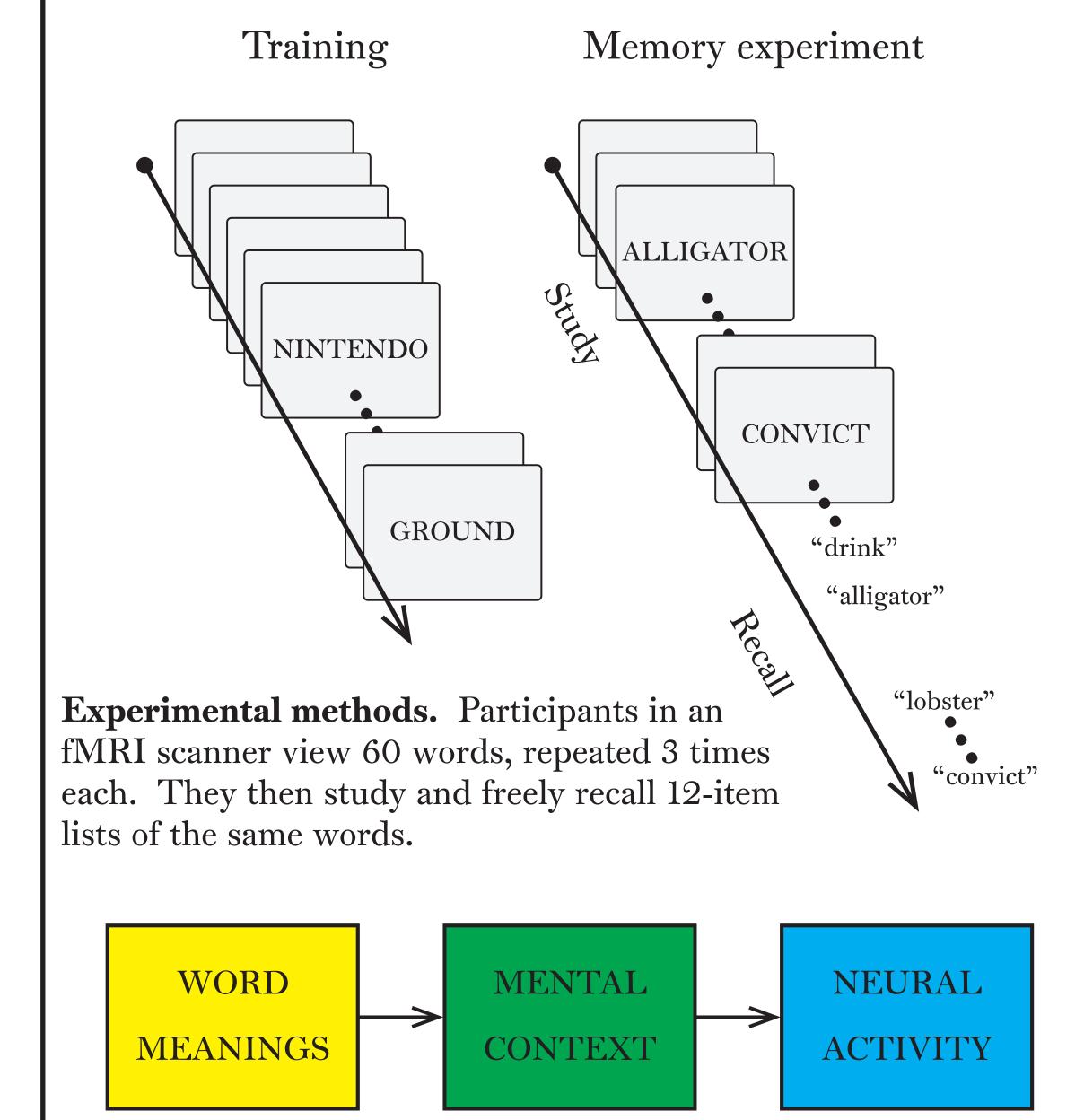
Overview & experimental methods

Context-based theories of memory posit that items on a studied list become associated with the mental contexts in which they are experienced.

We present a framework for tracking the neural correlates of individual items⁶ and the contexts in which they are experienced,⁵ during individual study and recall events.

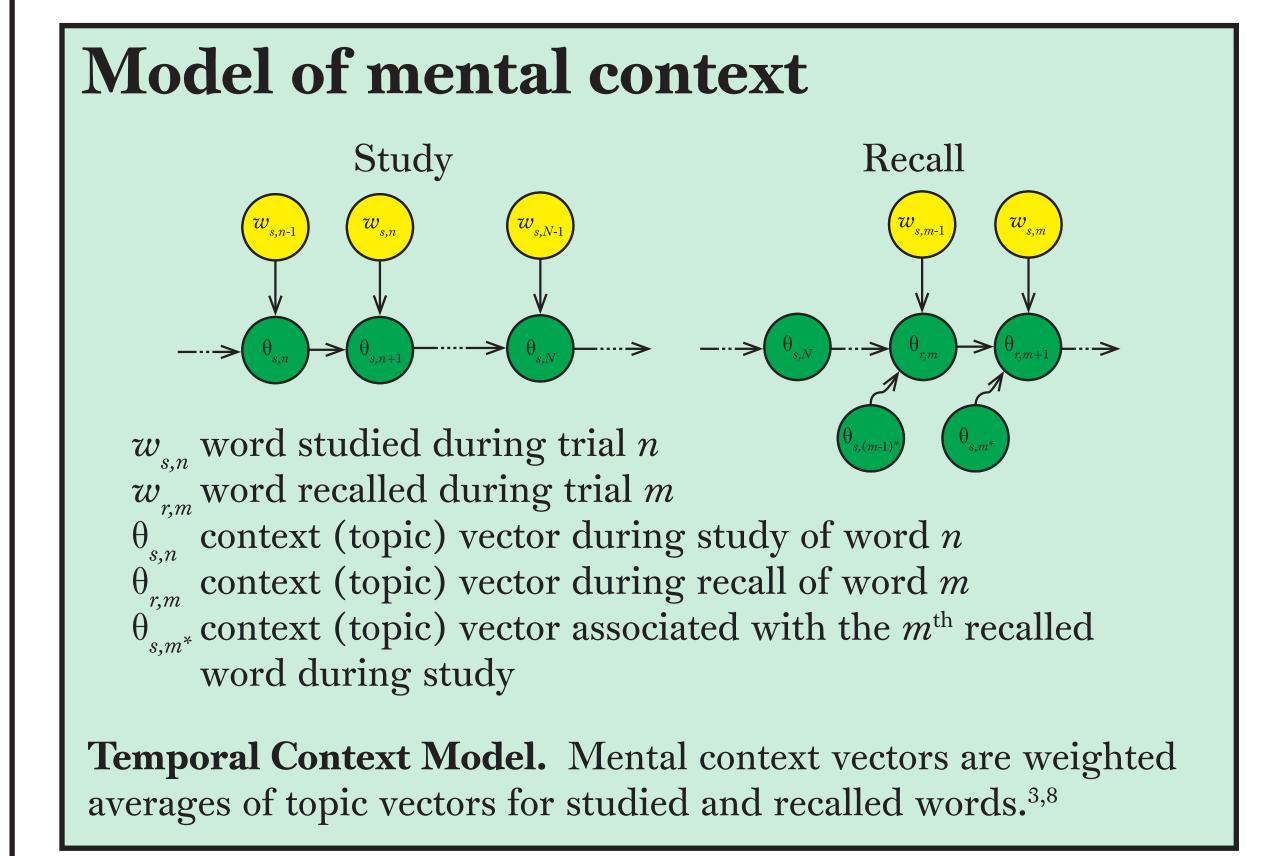


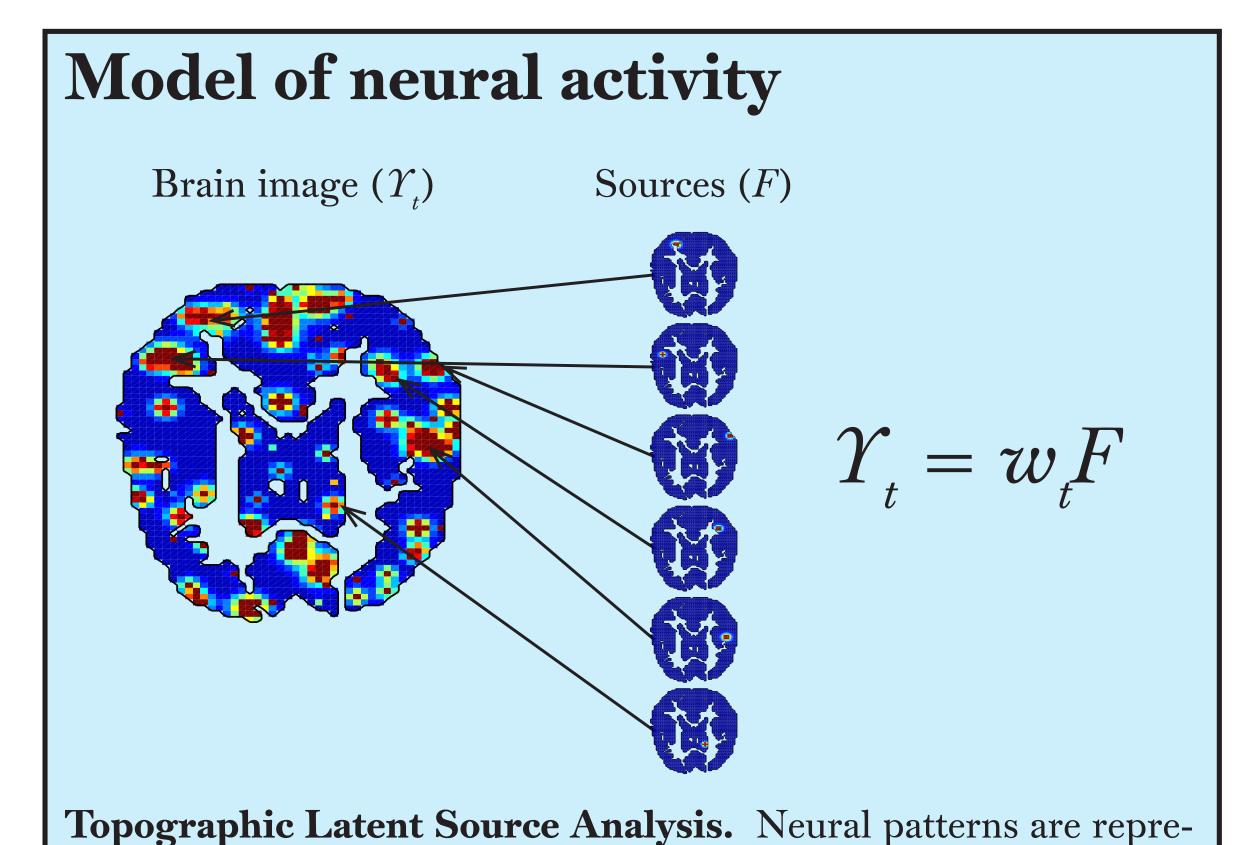
Context-based theories of memory. Context drifts gradually over time and becomes associated with each experienced event, giving rise to the contiguity effect in free recall.



Towards a unified model of corpora and cognition. Our approach attempts to infer the evolving state of mental context using text, behavioral, and neural data.

Model of word meanings Topic proportions and assignments Latent Dirichlet Allocation. Topic vectors are derived by analyz-





sented as linear combinations of spherical sources.^{2,4}

Decoding results

We collect brain images during each word presentation.

We use a topic model to obtain topic vectors for each presented word.

We infer the neural representation of each topic by computing a weighted average of the observed brain images. (Similar approach to 7.)

We interpret new neural patterns as topic vectors by computing the correlation between those patterns and the inferred representation of each topic.

Unfamiliar topics (biased towards less familiar words) are consistantly decoded more reliably than famliar topics.

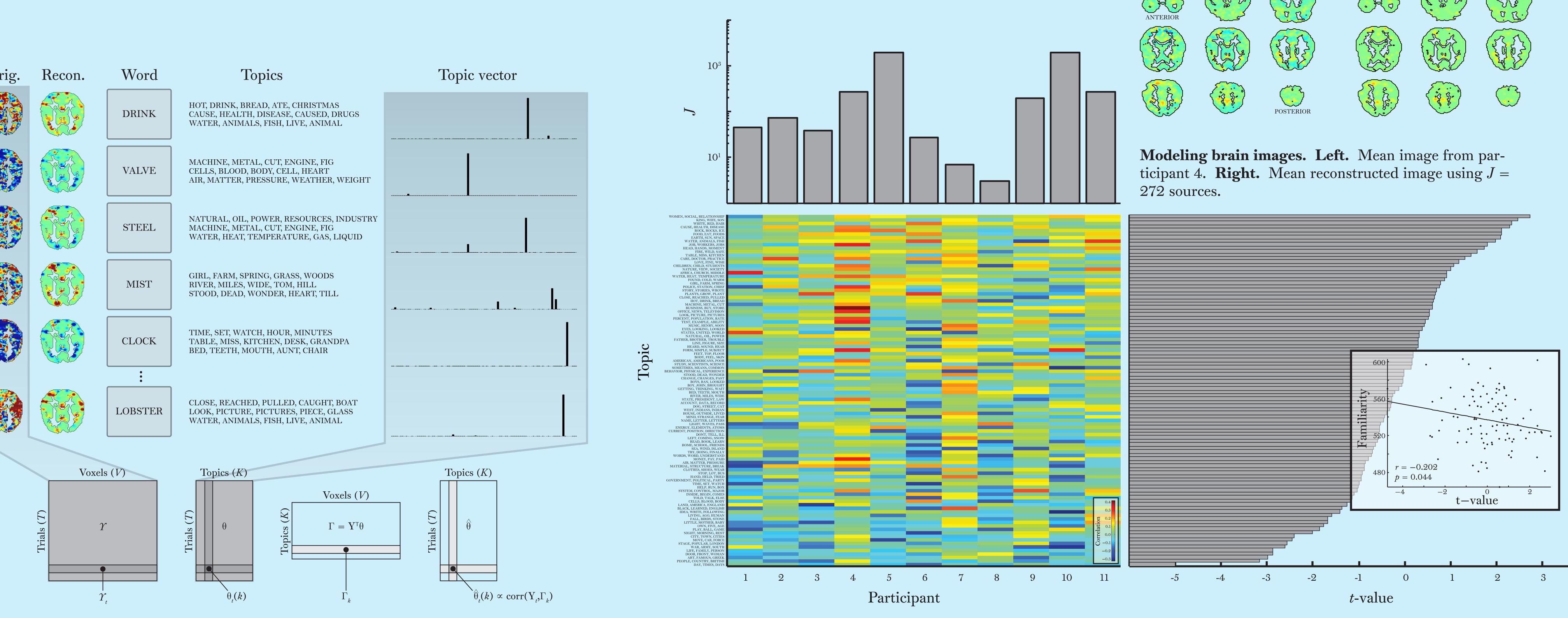
Brain images are modeled as linear combinations of spherical sources.

We use the observed brain images to infer the center locations and widths of each source.

We hold the center locations and widths fixed across images and fit per-image weights.

The vector of J weights provides a low dimensional representation of the original image.

The model can learn the neural representations of words or topics by incorporating weights that depend on which word or topic the participant is viewing. Reconstruction Original



Decoding topic vectors from brain images. Left. Decoding methods. We train and test the decoder using 10-fold cross validation. Right. Decoding results. The heatmap displays, for each participant and topic, the correlations between the decoded topic activations across trials (inferred using the brain images) and "true" topic vectors (inferred using a topic model fit to the TASA corpus). The bar plot on the right shows that some topics tend to be decoded well across participants; the inset shows that decodability varies with topic familiarity (MRC Psychological Database). The upper bar plot shows the number of sources that yielded the highest overall decoding accuracy for each participant.

Bibliography