

Mapping Between fMRI Responses and Natural Language Descriptions of Natural Stimuli

Kiran Vodrahalli*, Po-Hsuan Chen*, Yingyu Liang*, Christopher Baldassano*, Janice Chen♦, Esther Yong †, Christopher Honey♦, Peter J. Ramadge*, Kenneth A. Norman*, Sanjeev Arora*

Computation and the Brain F'18
October 3, 2018

* = Princeton, ♦ = Johns Hopkins, † = U. Toronto



Goal: **detect semantic meaning in fMRI signal.**

100 billion neurons in the brain

fMRI measures hemodynamic response at $\sim 10^5$ different 3mm x 3mm x 3mm voxels

Each voxel represents an average of the activity of the $\sim 10^6$ neurons it contains

Prior Work on Connecting a Semantic Space to fMRI Data

[Mitchell et al '08] predicts fMRI responses induced by **pictures of concrete nouns**.

[Naselaris et al '09] predicts fMRI responses induced by **images of scenes**.

[Pereira et al '11] uses the same dataset as Mitchell '08, but focuses on **generating words** related to the concrete nouns.

[Naselaris et al '11] tries to **reconstruct movie images** from fMRI signals measured while subjects watched movies.

[Wehbe et al '14] has subjects **read a chapter of Harry Potter** and predicts fMRI responses for held-out time points.

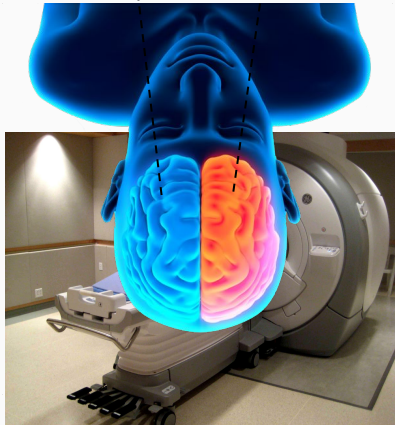
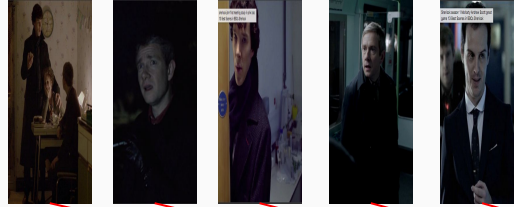
[Huth et al '16] reconstructs fMRI responses to **auditory stories**.

[Pereira et al '16] decodes fMRI responses to **word clouds and short sentences**.

Main Goal: Decode fMRI Response Semantics

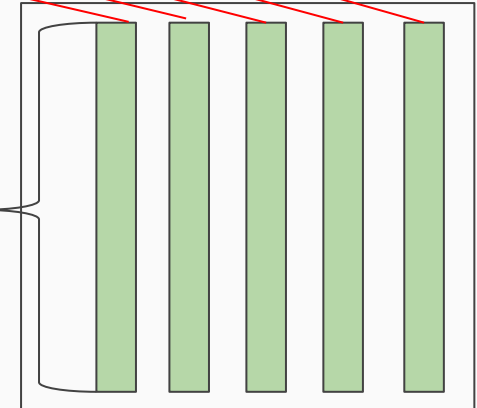


Movie scenes



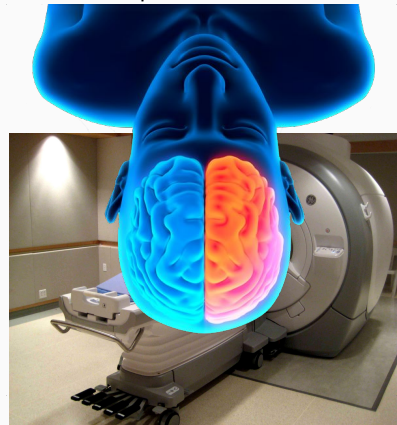
fMRI Machine

10^5
voxels



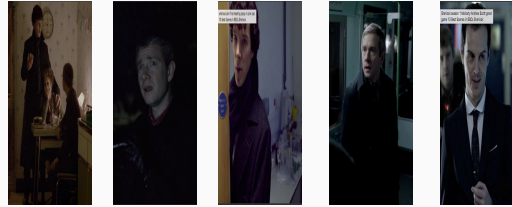
fMRI responses

Matching fMRI responses to annotations (Views: fMRI signal, text annotations)



fMRI Machine

Movie scenes



Annotations of movie scenes

Sherlock and John talk about the murder in an old room with Mrs. Hudson.

John is worried as Sherlock runs off.

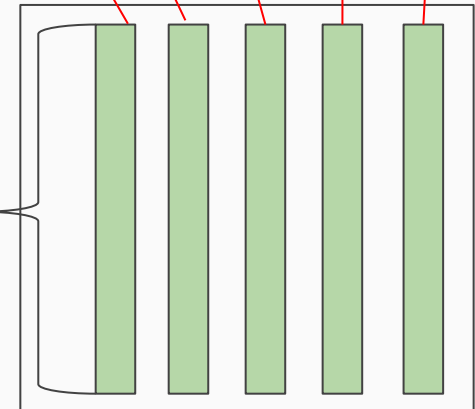
Sherlock enters the door to the chemistry lab, saying "John, I was here the whole time."

Once they get on the subway, John exclaims, "No you weren't!"

Moriarty arrives and says, "Hello Sherlock, John."

Each movie scene paired with text description from external party.

10^5
voxels



fMRI responses

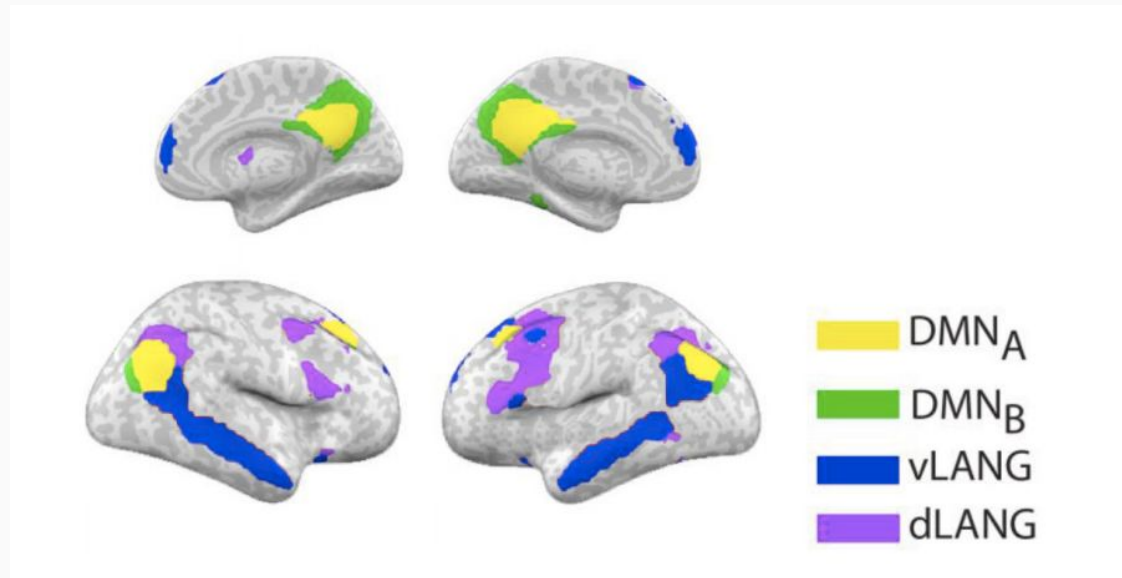
Representing Annotations: Word Embeddings

- To map between fMRI vectors and words, we would like to represent words in vector space
- Goals of embeddings: Preserve some notions of similarity and distance that apply to the words
- Idea: Assign to every word a 100-dim vector
- How?

- Idea: Train a predictive model on some task which “captures the meaning of natural language”.
 - Can be an unsupervised task, i.e. “language modeling”.
- Parameters of the model include a subset which are assigned uniquely to each word
 - Initialized randomly
- Training the model on external training set \rightarrow word vectors
- How to combine word vectors to get “annotation vectors”?

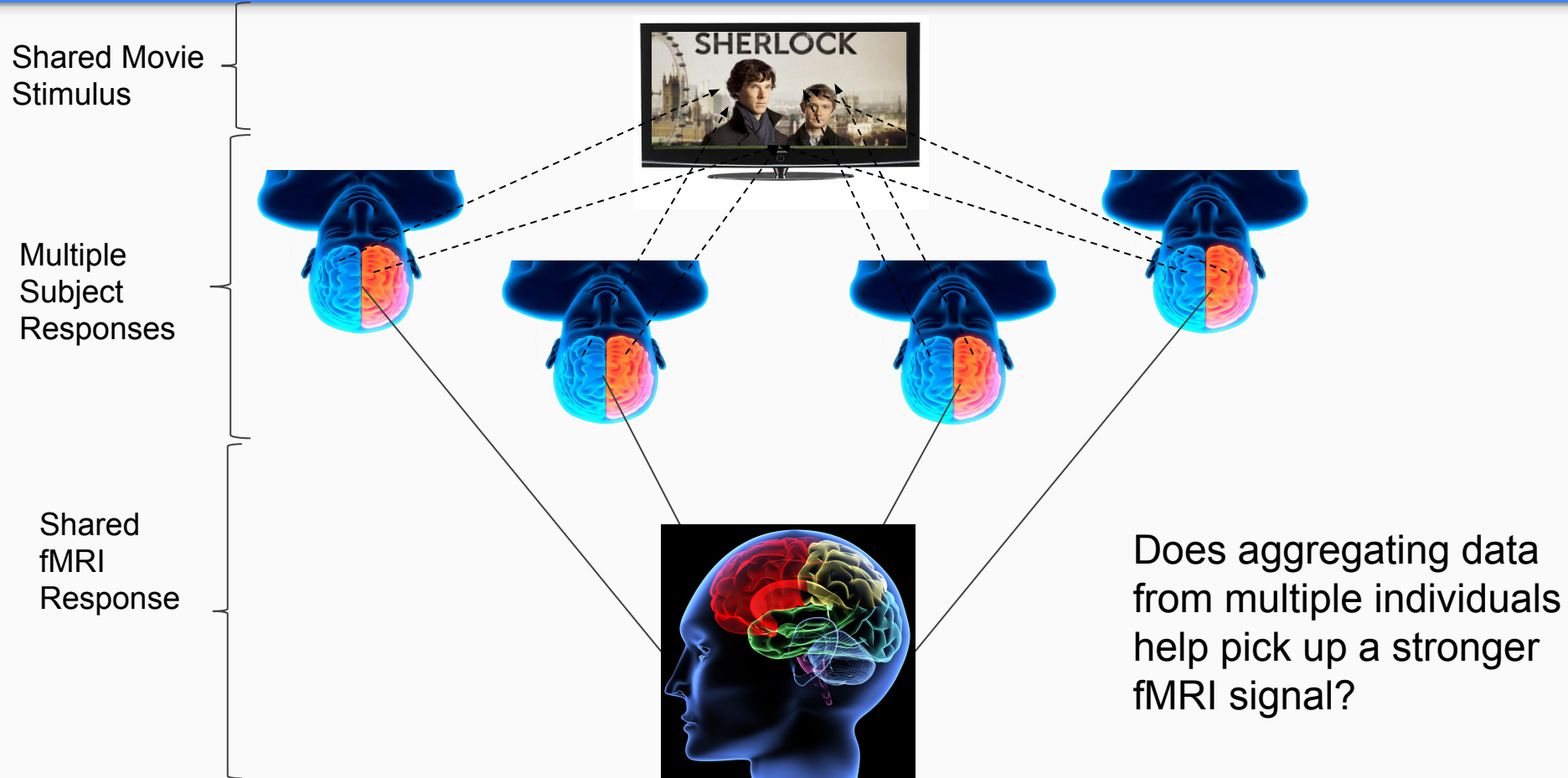
- The Shared Response Model (SRM, Chen et al. 2015) helps for decoding text!
- Weighted average word vectors → better semantic context vectors (ICLR 2017 paper, Arora et al)

Brain Regions (ROIs) Studied

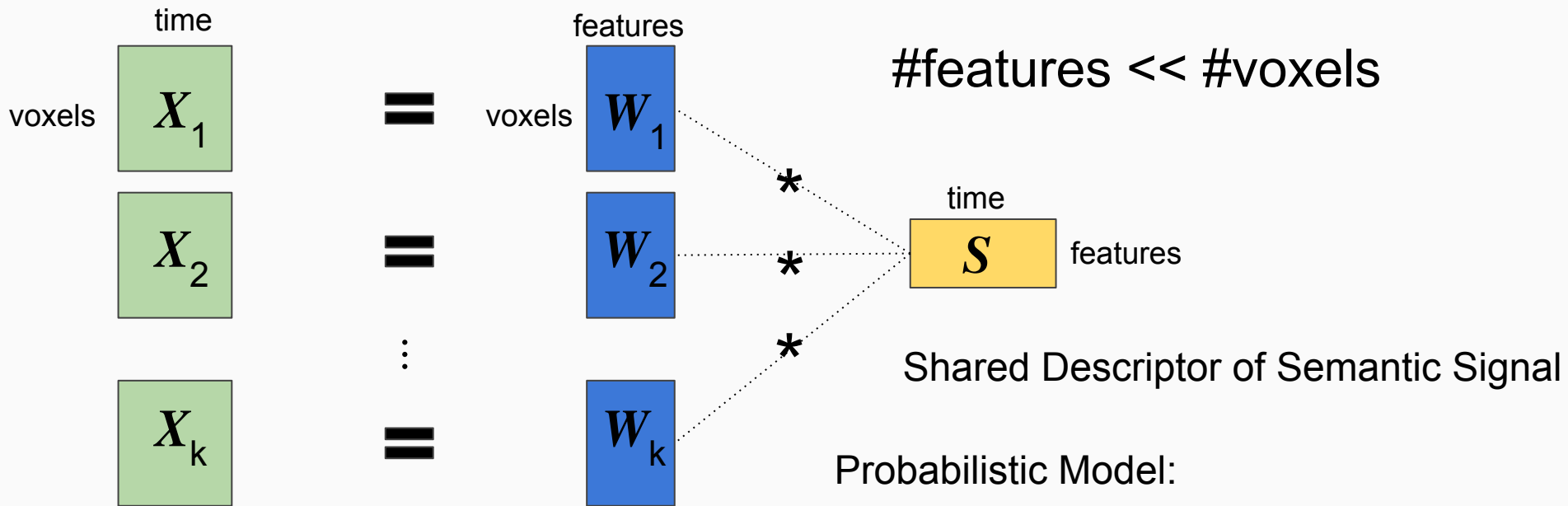


- Default Mode Network (DMN) standard area in literature
 - known to relate to narrative processing
 - DMN-A, -B (2000 voxels)
- Ventral/Dorsal Language (2000 voxels)
- Whole Brain (26000 voxels)
 - voxels with high inter-subject correlation
- Occipital Lobe (6000 voxels)

Leveraging Multiple Subject Views to Extract Better Semantics



Shared Response Model (SRM, [Chen, Chen, Yeshurun, Hasson, Haxby, Ramadge '15])



$$\operatorname{argmin}_{W^T W = I; S} \sum_{i=1}^k \|X_i - W_i S\|_F$$

$$s_t \sim \mathcal{N}(0, \Sigma_s)$$

$$x_{it} | s_t \sim \mathcal{N}(W_i s_t + \mu_i, \rho_i^2 I)$$

Embedding Annotations with Weighted Sums of Word Vectors



Fig. 3. Visualization of Semantic Annotation Vector Weightings: We display an example sentence from the Sherlock annotations, where we have colored important words red, and unimportant words blue. Brighter red means more important, and darker blue means less important.

Concatenating Previous Timepoints

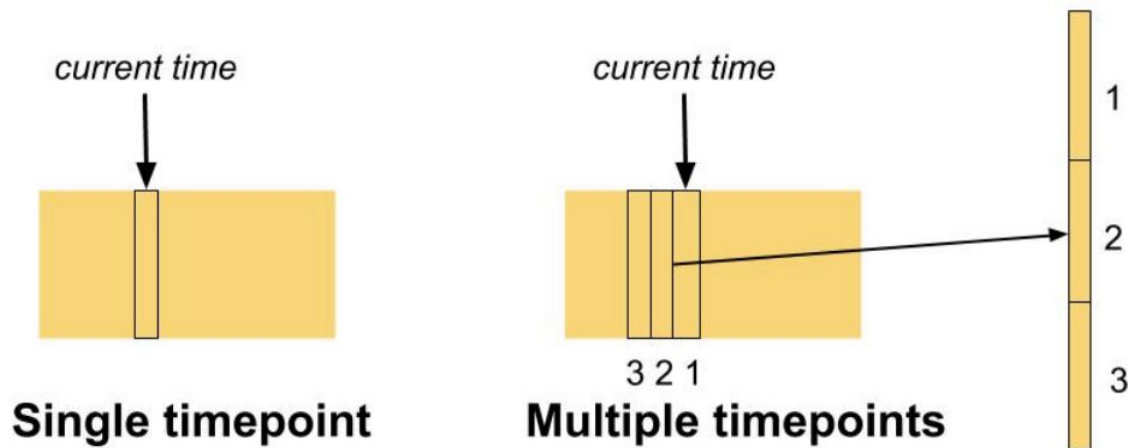


Fig. 4. Visualizing Concatenation: We visualize what the single timestep case looks like compared to a case where we use the previous two timesteps in our featurization as well. The latter case results in a more complicated model, since one of the dimensions of our linear map triples in size.

Basic Model:

$$WX = Y, \quad W \in \mathbb{R}^{m \times n}$$

X represents the fMRI data matrix (n x T)

Y represents the semantic annotation data matrix (m x T)

Learning the Map:

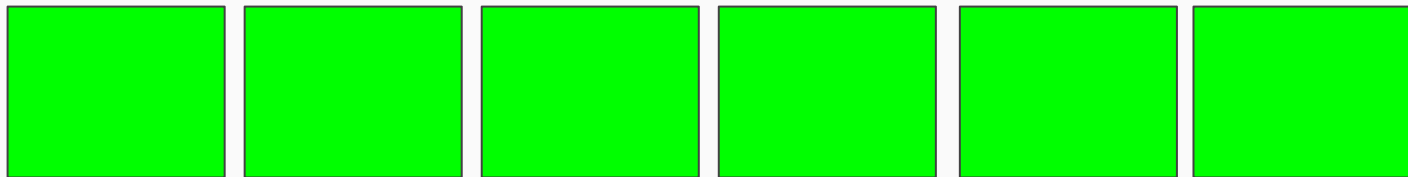
- Procrustes ($W^T W = I$)
 - Restrict map to “rotations” of the data.
 - Imposes strong constraint on map
- Ridge Regression (penalize l_2 norm $\|w\|_2$)
 - Classic linear regularization method
 - Restricts map weights to be uniformly small (not sparse)

Evaluation: Scene Classification/Ranking Experiments

25 test chunks from 1976 TRs

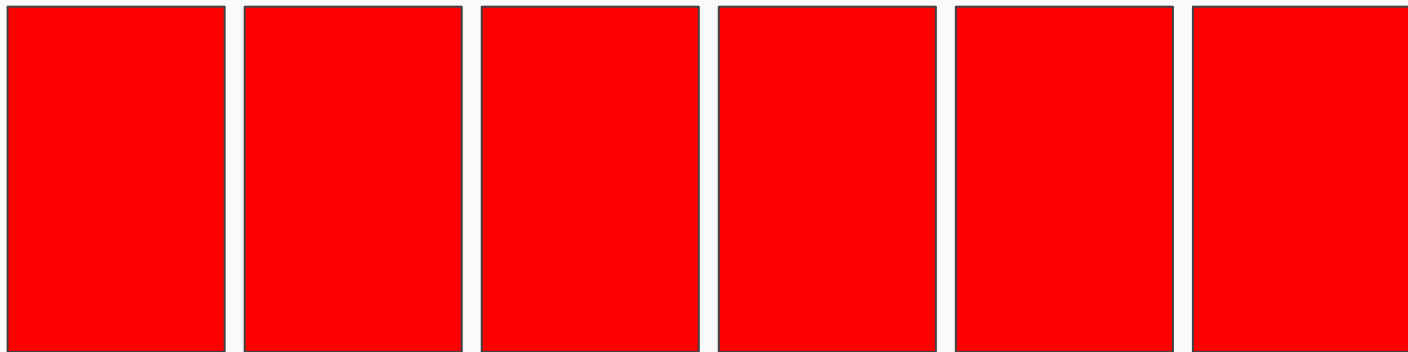
Shared fMRI
Space

20 dim



Semantic
Space

100 dim



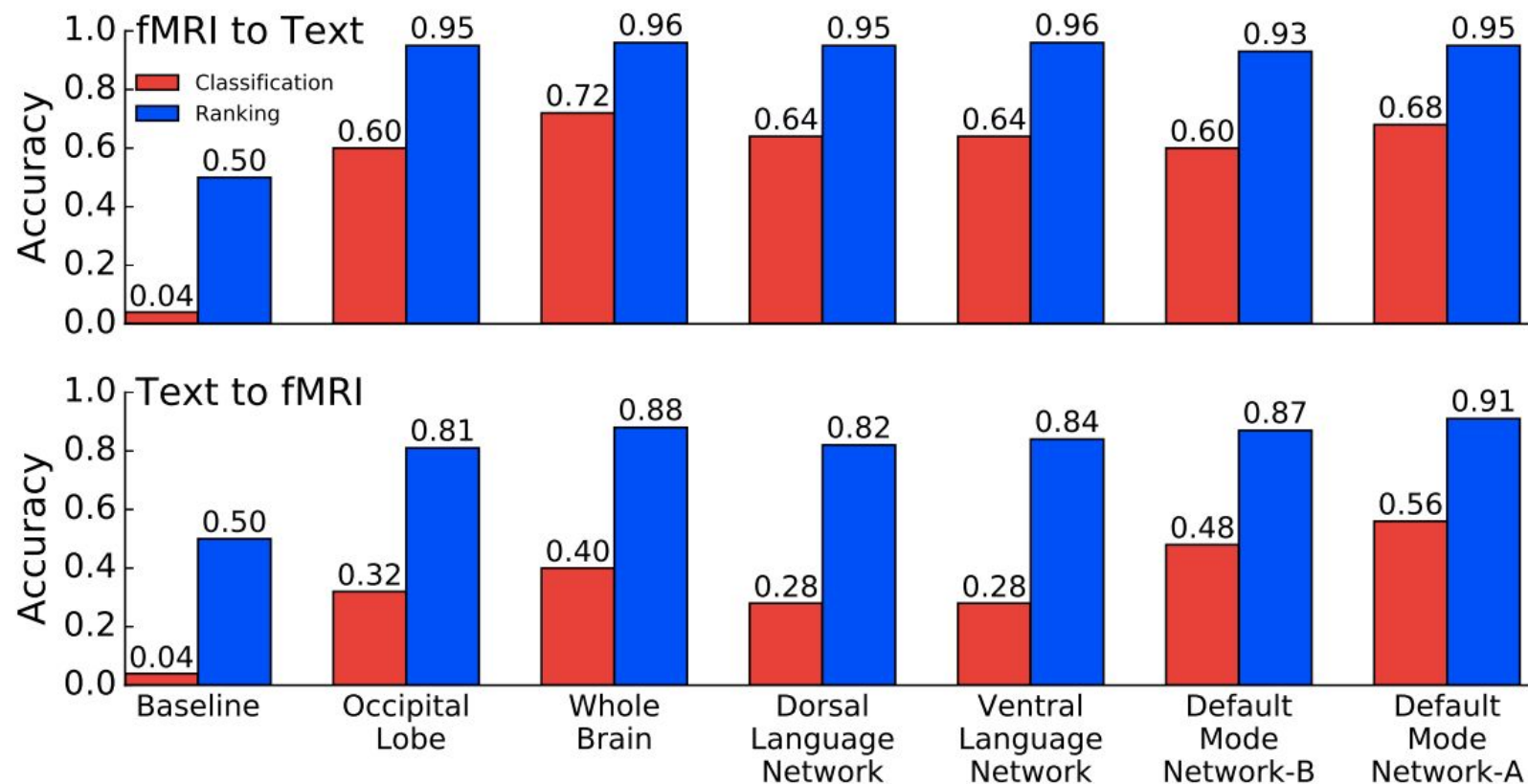
Results: Multiplicative Improvements with our Methods

Mapping Between fMRI Responses and Semantic Representations

fMRI \rightarrow Text	Maximum	Average
Previous Timesteps vs. None	5.3 \times	1.8 \times
Procrustes vs. Ridge	2.8 \times	1.3 \times
SRM/SRM-ICA vs. PCA	1.8 \times	1.3 \times
Weighted-SIF vs. Unweighted	1.6 \times	1.2 \times
Text \rightarrow fMRI	Maximum	Average
Previous Timesteps vs. None	2.5 \times	0.5 \times
Procrustes vs. Ridge	3.0 \times	0.8 \times
SRM/SRM-ICA vs. PCA	2.3 \times	1.2 \times
Weighted-SIF vs. Unweighted	1.8 \times	1.1 \times

Table 1. Table of Improvement Ratios for Various Algorithmic Parameters: In this table we give the maximum and average improvement ratios for a specific algorithmic technique over another, including usage of previous time steps, SRM/SRM-ICA versus PCA, SIF-weighted annotation embeddings versus unweighted annotation embeddings, and Procrustes versus ridge regression for both fMRI \rightarrow Text and Text \rightarrow fMRI. When we use previous timesteps, we consider the results for using 5 – 8 previous time steps. These numbers are all for the scene classification task. Note that the values from the maximum columns can be seen visually in Figures 6 and 7 respectively.

Results: Top-4% Classification and Average Rank



Results: Comparisons for fMRI → Text (4% Chance)

fMRI to Text (4% chance)

