# Text2Brain: Synthesis of Brain Activation Maps from Free-form Text Query

Gia H. Ngo[1*], Minh Nguyen[2*], Nancy F. Chen[3†], and Mert R. Sabuncu[1,4†]

[1] School of Electrical & Computer Engineering, Cornell University, USA
[2] Computer Science Department, University of California, Davis, USA
[3] Institute of Infocomm Research (I2R), A*STAR, Singapore
[4] Radiology, Weill Cornell Medicine, USA

**Abstract.** Most neuroimaging experiments are under-powered, limited by the number of subjects and cognitive processes that an individual study can investigate. Nonetheless, over decades of research, neuroscience has accumulated an extensive wealth of results. It remains a challenge to digest this growing knowledge base and obtain new insights since existing meta-analytic tools are limited to keyword queries. In this work, we propose Text2Brain, a neural network approach for coordinate-based meta-analysis of neuroimaging studies to synthesize brain activation maps from open-ended text queries. Combining a transformer-based text encoder and a 3D image generator, Text2Brain was trained on variable-length text snippets and their corresponding activation maps sampled from 13,000 published neuroimaging studies. We demonstrate that Text2Brain can synthesize anatomically-plausible neural activation patterns from free-form textual descriptions of cognitive concepts. Text2Brain is available at `https://braininterpreter.com` as a web-based tool for retrieving established priors and generating new hypotheses for neuroscience research.

**Keywords:** coordinate-based meta-analysis · transformers · information retrieval · image generation.

## 1 Introduction

Decades of neuroimaging research have yielded an impressive repertoire of findings and greatly enriched our understanding of the cognitive processes governing the mind. However, individual brain imaging experiments are often under-powered [1,2], constrained by the number of subjects and psychological processes that each experiment can probe [3]. To synthesize reliable trends across such experiments, researchers often perform meta-analysis on the coordinates of the most significant effect (such as 3D location of peak brain activation in response to a task). Most meta-analyses require the expert selection of relevant experiments (e.g. [4,5,6]). One key challenge with conducting meta-analysis on

---

*indicates equal contribution
†indicates equal contribution

neuroimaging experiments is the consolidation of synonymous terms. As neuroscientific research constantly evolves, different denominations might be used in different contexts or invented to refine existing ideas. For instance, "self-generated thought", one of the most highly studied functional domains of the human brain [7], can be referred to by different terms such as "task-unrelated thought" [8].

Manual selection of experiments for meta-analysis can be replaced by automated keyword search through data automatically scraped from the neuroimaging literature [9,10,11]. For example, Neurosynth [9] and more recently Neuroquery [10] both use automated keyword search to retrieve relevant studies to synthesize brain activation maps from text queries. However, Neurosynth and Neuroquery only allow for rigid queries formed out of predefined keywords and rely on superficial lexical similarity via co-occurrences of keywords for inference of longer or rarer queries. We propose an alternative approach named Text2Brain, which permits more flexible free-form text queries. Text2Brain also characterizes more fine-grained and implicit semantic similarity via vector representations from neural modeling in order to retrieve more relevant studies. Moreover, existing approaches estimate voxel-wise activations using either univariate statistical testing or regularized linear regression. In contrast, Text2Brain generates whole-brain activation maps using a 3D convolutional neural network (CNN) for more accurate construction of both coarse and fine details.

We compare Text2Brain's predictions with those from established baselines where we used article titles as free-form queries. Furthermore, we assess model predictions on independent test datasets, including reliable task contrasts and meta-analytic activation maps of well-studied cognitive domains predicted from their descriptions. Our analysis shows that Text2Brain generates activation maps that better match the target images than the baselines do. Given its flexibility in taking input queries, Text2Brain can be used as an educational aid as well as a tool for synthesizing prior maps for future research.

## 2    Materials and Methods

### 2.1    Overview

Figure 1 shows the overview of our approach. For each research article, full text and activation coordinates are extracted to create training samples (section 2.2). Text2Brain model consists of a transformer-based text encoder and a 3D CNN (section 2.3). The transformer uses attention to encode the input text into vector representation [12,13]. Thus, over many text-brain activation map pairs, the model automatically learns the association between activation at a spatial location with the most relevant words in the input text. Unlike classical keyword search that mainly exploits co-occurrence of keywords regardless of context, a transformer refines the vector representation depending on the specific phrasing of the text inputs (i.e. context) [14]. This allows Text2Brain to map synonymous text to a similar activation map. Instead of explicitly searching through articles, Text2Brain stores the articles' content in its parameters [15]

and outputs a relevant vector representation when presented with an input query. Thus, we use an augmented data sampling strategy to encourage the model to construct and store rich many-to-one mappings between textual description and activation maps (section 2.4).
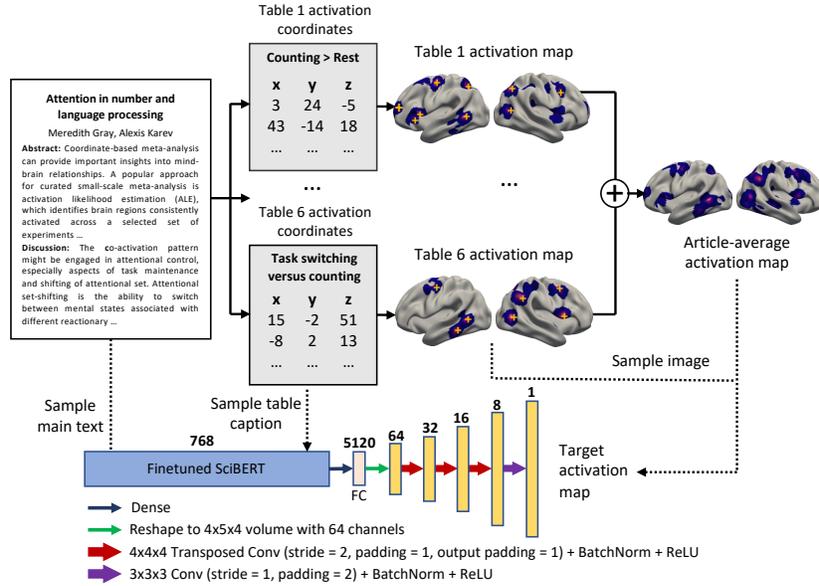


**Fig. 1.** Overview of data preprocessing, the Text2Brain model, and training procedure. All activation maps are 3D volumes, but projected to the surface for visualization.

## 2.2 Data Preprocessing

Coordinates of peak activation were scraped from tables of results reported in more than 13,000 neuroimaging articles and previously released in [10]. Each table has the corresponding article's PubMed ID, the table's ID as originally numbered in the article, and coordinates of peak activation converted to MNI152 coordinate system [16]. Following the preprocessing procedure of [10], a Gaussian sphere with full width at half maximum (FWHM) of 9mm is placed at each of the coordinates of peak activation. Thus, an activation map for each table in an article is generated from the set of activation foci associated with the table. An article-average activation map is also generated by averaging the activation maps of all the tables in the article. The articles' full text are scraped using their PubMedID via NCBI API [5] and Elsevier E-utilities API [6].

---

[5] https://www.ncbi.nlm.nih.gov/books/NBK25501/
[6] https://dev.elsevier.com/

### 2.3   Model

Figure 1 shows the basic schematic of Text2Brain, which consists of a text encoder based on SciBERT [17] and a 3D CNN as the image decoder. Output embedding of the text encoder is of dimension 768 and projected via a fully-connected layer, then reshaped to a 3D volume of dimension $4 \times 5 \times 4$ and 64 channels at each voxel. The image decoder consists of 3 transposed 3D convolutional layers with 32, 16, 8 channels respectively. The model was trained using mean-squared error for 2000 epochs, batch size of 24 with Adam [18]. The learning rate for the text encoder and image decoder are $10^{-5}$ and $3 \times 10^{-2}$, respectively. The model's source code is available at `https://github.com/sabunculab/text2brain`.

### 2.4   Training

During training, an activation map is sampled with equal probability from the set of table-specific activation maps and the article-average map. For each table-specific activation map, the first sentence of the table caption (as our data exploration suggested this to be the most useful description) is also extracted as the image's corresponding text. For each article-average activation map, one of four types of text is sampled with equal probability as the approximate description of the activation pattern, namely (1) the article's title (2) one of the article's keywords (2) abstract (3) a randomly chosen subset of sentences from the discussion section. This augmented sampling strategy encourages Text2Brain to generalize over input texts of different lengths. Furthermore, sampling multiple text snippets for an activation pattern encourages the model to automatically infer keywords present across queries and implicitly learn the association between different but synonymous words with an activation map. Supplemental Figure S2 shows an ablation study on the sampling strategy.

## 3   Experimental Setup

### 3.1   Predict activation maps from article titles

From the dataset of 13000 articles, 1000 articles are randomly sampled as the test set such that the keywords (defined by the articles' authors) are not included in the training and validation articles. Of the remaining articles, 1000 are randomly held out as a validation set for parameters tuning. For each article, the article-average activation map is predicted from its title using Text2Brain and the two baselines of Neurosynth and Neuroquery.

### 3.2   Predict activation maps from contrast descriptions

The Human Connectome Project (HCP) offers neuroimaging data from over 1200 subjects, including task fMRI (tfMRI) of 86 task contrasts from 7 domains [19]. While detailed descriptions of task contrasts are provided by HCP,

we instead use the more concise contrast descriptions provided by the Individual Brain Charting (IBC) project [20], which includes fMRI data from 12 subjects and 180 task contrasts, 43 of which are also studied in the HCP. The reason for using the IBC contrast descriptions is because they are more succinct and thus more favorable to the baselines. The target (ground-truth) activation maps are the group-average contrast maps provided by the HCP, as the large number of subjects provides more reliable estimates of the contrast maps. In our analyses, we use the agreement between the IBC and HCP maps as a measure of reliability. Note that despite using similar experimental protocols, there are subtle differences between the IBC and HCP experiments. For example, while the original HCP language task was conducted in English, the corresponding language task in the IBC project was conducted in French. Overall, Text2Brain and the two baselines were evaluated on the 43 HCP task contrasts.

### 3.3   Baselines

The first baseline, Neurosynth [9], collected all peak activation coordinates across neuroimaging articles that mention a given keyword and performed a statistical test at every voxel to determine a significant association. For longer query, we performed statistical test using activation coordinates reported in all articles that contain at least one of the keywords in the input text.

The second baseline, Neuroquery [10], builds upon Neurosynth by extending the vocabulary of keywords via manual selection from other sources. The keyword encoding is obtained after performing non-negative matrix factorization of the articles' abstract (as a bag of keywords) represented with term frequency - inverse document frequency (TF-IDF) features [21]. A ridge regression model was trained to learn the mapping from the text encoding to the activation at each voxel. The inference of a keyword is smoothed by a weighed summation with most related keywords (in the TF-IDF space). For longer queries, the predicted activation is the average of maps predicted from all keywords in the input.

### 3.4   Evaluation Metrics

To measure the similarity of predicted and target activation maps at different levels of detail, we compute Dice scores [22] at various thresholds. This evaluation procedure is similar to that used in [10] for a thresholded target map, but we apply the same thresholding to both the target and predicted map. For example, at a lower threshold (e.g., considering the 5% most activated voxels), the Dice score measures the correspondence of the fine-grained details between the target and predicted activation maps. At higher thresholds (e.g. 25% most activated voxels), this metric captures gross agreement of activation clusters. We also compute an approximated integration of Dice scores across all thresholds (from 5% up to 30%), i.e. the area under the Dice curve (AUC), as a summary measure. Supplemental Figure S1 shows the Dice curve for an example pair of target-predicted activation maps. We only consider up to 30% to be fair to the baselines,

as the portion of activated voxels predicted by Neuroquery only extends up to 30% of the gray matter mask.

## 4   Results

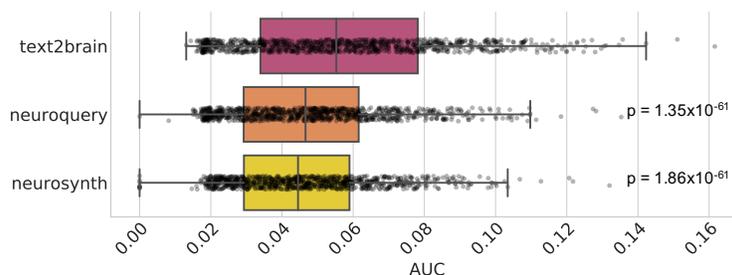### 4.1   Validate activation maps predicted from article title



**Fig. 2.** Evaluation of article-average activation maps predicted from their titles measured in area under the Dice curve (AUC) score. The p-values are computed from paired-sample t-tests between Text2Brain and each of the 2 baselines.

Figure 2 compares the quality of activation maps predicted from the titles of 1000 articles. Text2Brain model (mean AUC = 0.0576) outperforms Neuroquery (mean AUC = 0.0478) and Neurosynth (mean AUC = 0.0464). Paired-sample t-tests show that this performance gap is statistically very significant. The p-value for the comparison between Neuroquery and Neurosynth is $p = 0.015$. While Text2Brain can make a prediction for all samples, Neurosynth and Neuroquery fail to make prediction for some article titles, resulting in zero AUCs values.

### 4.2   Prediction of task contrast maps from description

Figure 3 shows the AUC scores for the prediction of the three models and the IBC average contrasts, against the HCP target maps. The 22 contrasts with the HCP-IBC's AUC score above the average, considered to be the reliable contrasts, are shown. Across all 43 HCP contrasts, Text2Brain (mean AUC = 0.082) performs better than the baselines, i.e. Neuroquery (mean AUC = 0.0755, $p = 0.08$), Neurosynth (mean AUC = 0.047, $p = 1.5 \times 10^{-5}$), where $p$-values are computed from the paired t-test between Text2Brain's and the baselines' prediction. As reference, IBC contrasts yield mean AUC = 0.094 ($p = 0.077$).

Figure 4 shows the prediction from three most reliable task contrasts (having the highest HCP-IBC AUC), thresholded at the top 20% most activated voxels. The three contrasts correspond to different HCP task groups, namely "WORKING MEMORY", "SOCIAL", and "MOTOR". Text2Brain's prediction
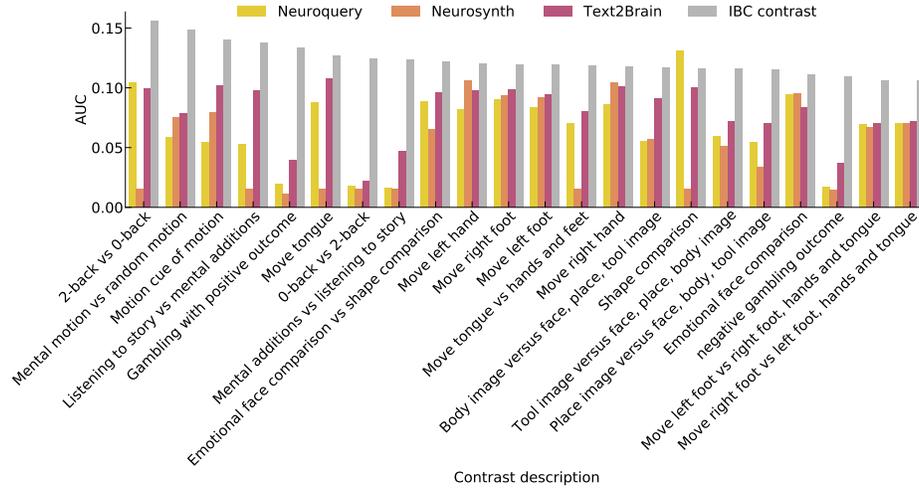
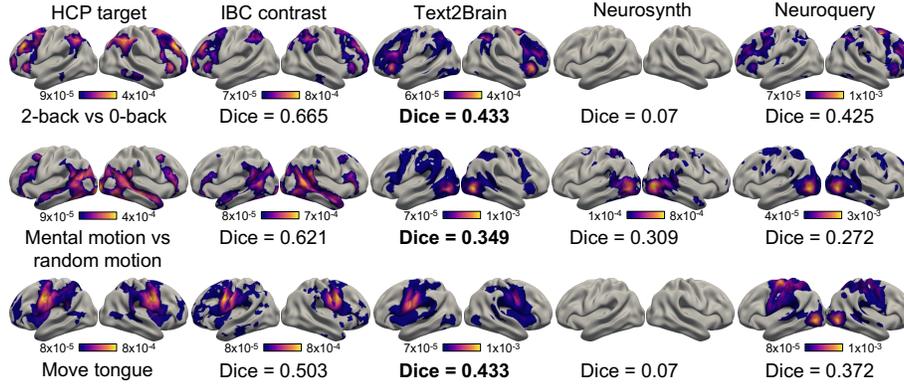**Fig. 3.** AUC of predicted HCP task activation maps from contrasts' description.



**Fig. 4.** Task activation maps predicted from contrasts' description. The Dice scores are computed between the binarized map of 20% most activated voxels in the predicted and target brain maps.

improves over the baselines for the three contrasts. Neurosynth was not able to generate activation maps for two of the contrast descriptions ("2-back vs 0-back" and "Move tongue"). On the other hand, for the "Move tongue" contrast, Neuroquery predicts activation in the primary cortex, but the peak is in the wrong location, shifted more toward the hand region of the homunculus. Additionally, there is a false positive prediction in the occipital cortex.

Finally, we are interested in examining the prediction for "Self-generated thought", which is one of the most commonly studied functional domains, due to its engagement in a wide range of cognitive processes that do not require external stimuli [8], and is associated with the default network [23]. The ground-truth map for self-generated thought, taken from [24], is estimated using activation likelihood estimation (ALE) [25,26,27], a well established tool for coordinate-based meta-analysis, applied on 1812 activation foci across 167 imaging studies over 7 tasks based on strict selection criteria [28,29,30]. Figure 5 shows the
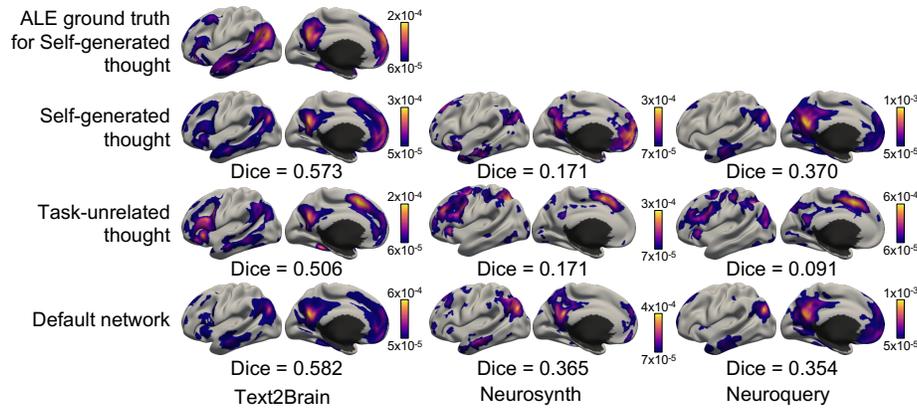


**Fig. 5.** Prediction of self-generated thought activation map using synonymous queries

prediction of self-generated thought activation map using three different query terms, thresholded at the top 20% most activated voxels. For the "self-generated thought" and "default network" queries, all approaches generate activation maps that are consistent with the ground-truth, which includes the precuneus, the medial prefrontal cortex, the temporo-parietal junction, and the temporal pole. Text2Brain's prediction best matches the ground-truth activation map compared to the baselines. Text2Brain can also replicate a similar activation pattern from the query "task-unrelated thought", evident by only a slight drop in the Dice score. However, Neuroquery and Neurosynth both produce activation maps that deviate from the typical default network's regions with increased activation in the prefrontal cortex, also evident by a large drop in the Dice scores.

## 5    Conclusion

In this work, we present a model named Text2Brain for generating activation maps from free-form text query. By finetuning a high-capacity SciBERT-based text encoder to predict coordinate-based meta-analytic maps, Text2Brain captures the rich relationship in the language representational space, allowing the model to generalize its prediction for synonymous queries. This is evident in the better performance of Text2Brain in predicting the self-generated thought activation map using different descriptions of the functional domain. Text2Brain's capability to implicitly learn relationships between terms and images will help the model stays relevant and useful even as neuroimaging literature continues to evolve with new information and rephrasing of existing concepts. We also show that Text2Brain accurately predicts most of the task contrasts included in the HCP dataset based on their description, validating its capability to make prediction for longer, arbitrary queries. Text2Brain also avoids the failure cases suffered by Neurosynth and Neuroquery in which they cannot predict if the input words are not defined in their vocabularies, even though the queries are relevant to neuroscience research such as the title of an article or a contrast description. In the future, we will work on the interpretability of the approach, such as to attribute regions of activation in the generated map to specific word in the input query, as well as to efficiently match activation maps and research text most relevant to the synthesized images.

## Acknowledgement

## References

1. Joshua Carp. The secret lives of experiments: methods reporting in the fMRI literature. Neuroimage, 63(1):289–300, 2012.
2. Katherine S Button, John PA Ioannidis, Claire Mokrysz, Brian A Nosek, Jonathan Flint, Emma SJ Robinson, and Marcus R Munafò. Power failure: why small sample size undermines the reliability of neuroscience. Nature reviews neuroscience, 14(5):365–376, 2013.
3. Jessica A Church, Steven E Petersen, and Bradley L Schlaggar. The "Task B problem" and other considerations in developmental functional neuroimaging. Human brain mapping, 31(6):852–862, 2010.
4. Sergi G Costafreda, Michael J Brammer, Anthony S David, and Cynthia HY Fu. Predictors of amygdala activation during the processing of emotional stimuli: a meta-analysis of 385 pet and fmri studies. Brain research reviews, 58(1):57–70, 2008.
5. Michael J Minzenberg, Angela R Laird, Sarah Thelen, Cameron S Carter, and David C Glahn. Meta-analysis of 41 functional neuroimaging studies of executive function in schizophrenia. Archives of general psychiatry, 66(8):811–822, 2009.

6. Alexander J Shackman, Tim V Salomons, Heleen A Slagter, Andrew S Fox, Jameel J Winter, and Richard J Davidson. The integration of negative affect, pain and cognitive control in the cingulate cortex. Nature Reviews Neuroscience, 12(3):154–167, 2011.
7. Jonathan Smallwood. Distinguishing how from why the mind wanders: a process–occurrence framework for self-generated mental activity. Psychological bulletin, 139(3):519, 2013.
8. Jessica R Andrews-Hanna, Jonathan Smallwood, and R Nathan Spreng. The default network and self-generated thought: component processes, dynamic control, and clinical relevance. Annals of the New York Academy of Sciences, 1316(1):29, 2014.
9. Tal Yarkoni, Russell A Poldrack, Thomas E Nichols, David C Van Essen, and Tor D Wager. Large-scale automated synthesis of human functional neuroimaging data. Nature methods, 8(8):665–670, 2011.
10. Jérome Dockès, Russell A Poldrack, Romain Primet, Hande Gözükan, Tal Yarkoni, Fabian Suchanek, Bertrand Thirion, and Gaël Varoquaux. NeuroQuery, comprehensive meta-analysis of human brain mapping. Elife, 9:e53385, 2020.
11. Timothy N Rubin, Oluwasanmi Koyejo, Krzysztof J Gorgolewski, Michael N Jones, Russell A Poldrack, and Tal Yarkoni. Decoding brain activity using a large-scale probabilistic functional-anatomical atlas of human cognition. PLoS computational biology, 13(10):e1005649, 2017.
12. Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. arXiv preprint arXiv:1706.03762, 2017.
13. Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805, 2018.
14. Ian Tenney, Patrick Xia, Berlin Chen, Alex Wang, Adam Poliak, R Thomas McCoy, Najoung Kim, Benjamin Van Durme, Samuel R Bowman, Dipanjan Das, et al. What do you learn from context? Probing for sentence structure in contextualized word representations. arXiv preprint arXiv:1905.06316, 2019.
15. Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. Language Models as Knowledge Bases? In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 2463–2473, 2019.
16. Jack L Lancaster, Diana Tordesillas-Gutiérrez, Michael Martinez, Felipe Salinas, Alan Evans, Karl Zilles, John C Mazziotta, and Peter T Fox. Bias between MNI and Talairach coordinates analyzed using the ICBM-152 brain template. Human brain mapping, 28(11):1194–1205, 2007.
17. Iz Beltagy, Kyle Lo, and Arman Cohan. SciBERT: A pretrained language model for scientific text. arXiv preprint arXiv:1903.10676, 2019.
18. Ilya Loshchilov and Frank Hutter. Decoupled Weight Decay Regularization. In Proceedings of ICLR, 2018.
19. Deanna M Barch, Gregory C Burgess, Michael P Harms, Steven E Petersen, Bradley L Schlaggar, Maurizio Corbetta, Matthew F Glasser, Sandra Curtiss, Sachin Dixit, Cindy Feldt, et al. Function in the human connectome: task-fMRI and individual differences in behavior. Neuroimage, 80:169–189, 2013.
20. Ana Luísa Pinho, Alexis Amadon, Baptiste Gauthier, Nicolas Clairis, André Knops, Sarah Genon, Elvis Dohmatob, Juan Jesús Torre, Chantal Ginisty, Séverine

Becuwe-Desmidt, et al. Individual Brain Charting dataset extension, second release of high-resolution fMRI data for cognitive mapping. Scientific Data, 7(1):1–16, 2020.

21. Gerard Salton and Christopher Buckley. Term-weighting approaches in automatic text retrieval. Information processing & management, 24(5):513–523, 1988.

22. Lee R Dice. Measures of the amount of ecologic association between species. Ecology, 26(3):297–302, 1945.

23. Randy L Buckner, Jessica R Andrews-Hanna, and Daniel L Schacter. The brain's default network: anatomy, function, and relevance to disease. 2008.

24. Gia H Ngo, Simon B Eickhoff, Minh Nguyen, Gunes Sevinc, Peter T Fox, R Nathan Spreng, and BT Thomas Yeo. Beyond consensus: embracing heterogeneity in curated neuroimaging meta-analysis. NeuroImage, 200:142–158, 2019.

25. Peter E Turkeltaub, Guinevere F Eden, Karen M Jones, and Thomas A Zeffiro. Meta-analysis of the functional neuroanatomy of single-word reading: method and validation. Neuroimage, 16(3):765–780, 2002.

26. Angela R Laird, P Mickle Fox, Cathy J Price, David C Glahn, Angela M Uecker, Jack L Lancaster, Peter E Turkeltaub, Peter Kochunov, and Peter T Fox. ALE meta-analysis: Controlling the false discovery rate and performing statistical contrasts. Human brain mapping, 25(1):155–164, 2005.

27. Simon B Eickhoff, Angela R Laird, Christian Grefkes, Ling E Wang, Karl Zilles, and Peter T Fox. Coordinate-based activation likelihood estimation meta-analysis of neuroimaging data: A random-effects approach based on empirical estimates of spatial uncertainty. Human brain mapping, 30(9):2907–2926, 2009.

28. R Nathan Spreng, Raymond A Mar, and Alice SN Kim. The common neural basis of autobiographical memory, prospection, navigation, theory of mind, and the default mode: a quantitative meta-analysis. Journal of cognitive neuroscience, 21(3):489–510, 2009.

29. Raymond A Mar. The neural bases of social cognition and story comprehension. Annual review of psychology, 62:103–134, 2011.

30. Gunes Sevinc and R Nathan Spreng. Contextual and perceptual brain processes underlying moral cognition: a quantitative meta-analysis of moral reasoning and moral emotions. PloS one, 9(2):e87427, 2014.

## Supplementary Materials for "Text2Brain: Synthesis of Brain Activation Maps from Free-form Text Query"

## A    Evaluation Metrics

Dice score [22] is used to measure the extent of overlap between a predicted activation map and the target activation map at a given threshold. At a given threshold of $x\%$, Dice score is computed as:

$$Dice(x) = \frac{2|Prediction(x) \cap Target(x)|}{|Prediction(x)| + |Target(x)|},\tag{1}$$

where $|Prediction(x)|$ denotes the number of top $x\%$ most activated voxels in the predicted activation map, $|Target(x)|$ denotes the number of top $x\%$ most activated voxels in the target map, and $|Prediction(x) \cap Target(x)|$ denotes the number of voxels that overlap between the predicted and target map at the given threshold.
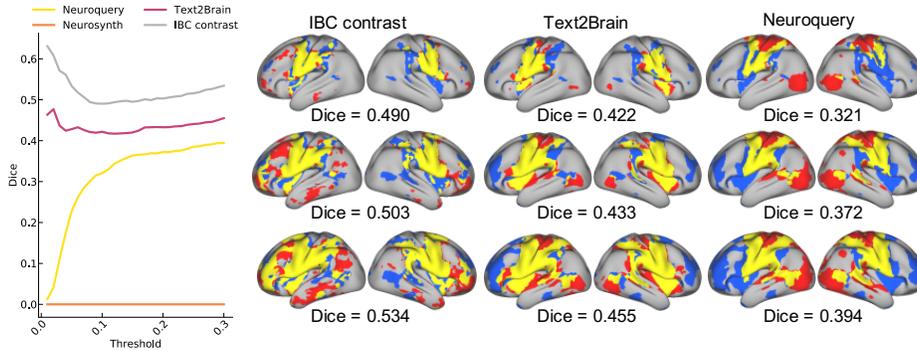


**Fig. S1.** Example Dice scores evaluated on the "Move tongue" contrast (in Fig. 4. The graph on the left shows the Dice scores computed between the target HCP activation map and Text2Brain's, Neurosynth's, Neuroquery's prediction, and the IBC contrast across thresholds ranging from 5% to 30%. Note that Neurosynth's Dice scores are all zeros as it fails to make an inference for the input text. The area under the Dice curve (AUC) was computed as the summary metrics of accuracy across all thresholds (e.g. Fig 3). The brain maps on the right are visualization of the extent of overlaps between predicted and target maps at 10%, 20% and 30% threshold of most activated voxels. Blue indicates activation in the target contrast, red is the predicted activation and yellow is the overlap.

# B    Ablation study of sampling strategy

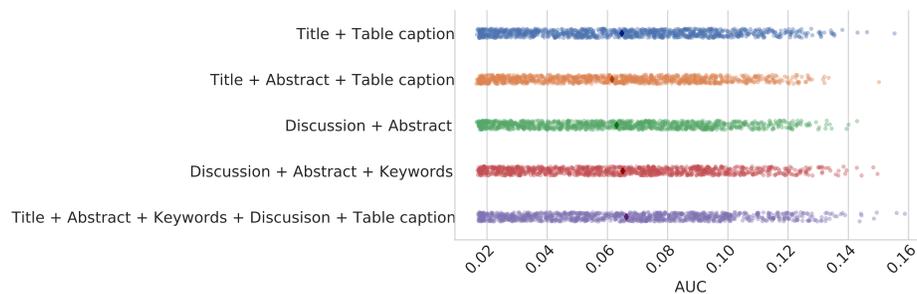| Text samples | Mean AUC |
|---|---|
| Title + Table caption | 0.0648 |
| Title + Abstract + Table caption | 0.0616 |
| Discussion + Abstract | 0.0631 |
| Discussion + Abstract + Keywords | 0.0651 |
| **Title + Abstract + Keywords + Discussion + Table caption** | 0.0663 |



**Fig. S2.** Performance of different sampling strategies in predicting article-average activation maps from the articles' titles in the validation set. All sampling strategies used the same model described in 2.3. The model parameters used for evaluation were chosen at the epoch with the best performance on the validation set.