

# GPT-2's activations predict the degree of semantic comprehension in the human brain

Charlotte Caucheteux<sup>1,2,\*</sup>, Alexandre Gramfort<sup>2</sup>, and Jean-Rémi King<sup>1,3</sup>

<sup>1</sup> Facebook AI Research, Paris, France; <sup>2</sup> Université Paris-Saclay, Inria, CEA, Palaiseau, France; <sup>3</sup> École normale supérieure, PSL University, CNRS, Paris, France

1 **Language transformers, like GPT-2, have demonstrated remarkable  
2 abilities to process text, and now constitute the backbone of  
3 deep translation, summarization and dialogue algorithms. However,  
4 whether these models encode information that relates to human com-  
5 prehension remains controversial. Here, we show that the represen-  
6 tations of GPT-2 not only map onto the brain responses to spoken  
7 stories, but also predict the extent to which subjects understand nar-  
8 ratives. To this end, we analyze 101 subjects recorded with func-  
9 tional Magnetic Resonance Imaging while listening to 70 min of short  
10 stories. We then fit a linear model to predict brain activity from  
11 GPT-2's activations, and correlate this mapping with subjects' com-  
12 prehension scores as assessed for each story. The results show that  
13 GPT-2's brain predictions significantly correlate with semantic com-  
14 prehension. These effects are bilaterally distributed in the language  
15 network and peak with a correlation of  $R=0.50$  in the angular gyrus.  
16 Overall, this study paves the way to model narrative comprehension  
17 in the brain through the lens of modern language algorithms.**

Neuroscience of language | Deep Neural Networks

1 In less than two years, language transformers like GPT-2  
2 have revolutionized the field of natural language processing  
3 (NLP). These deep learning architectures are typically trained  
4 on very large corpora to complete partially-masked texts, and  
5 provide a one-fit-all solution to translation, summarization,  
6 and question-answering tasks and algorithms (1).

7 Critically, their hidden representations have been shown to  
8 – at least partially – correspond to those of the brain: single-  
9 sample fMRI (2–4), MEG (2, 4), and intracranial responses to  
10 spoken and written texts (3, 5) can be significantly predicted  
11 from a linear combination of the hidden vectors generated  
12 by these deep networks. Furthermore, the quality of these  
13 predictions directly depends on the models' ability to complete  
14 text (3, 4).

15 In spite of these achievements, strong doubts subsist on  
16 whether language transformers actually encode meaningful  
17 constructs (6). When asked to complete "I had \$20 and gave  
18 \$10 away. Now, I thus have \$", GPT-2 predicts "20\*". Simi-  
19 lar trivial errors can be observed for geographical locations,  
20 temporal ordering, pronoun attribution and causal reasoning.  
21 These results have thus led some to argue that such "system  
22 has no idea what it is talking about" (7). Thus, how the rep-  
23 resentations of GPT-2 relate to a human-like understanding  
24 remains largely unknown.

25 Here, we propose to evaluate how the similarity between the  
26 brain and GPT-2 vary with semantic comprehension. Specific-  
27 ally, we first compare GPT-2's activations to the functional  
28 Magnetic Resonance Imaging of 101 subjects listening to  
29 70 min of seven short stories, and we quantify this similar-  
30 ity with a "brain score" ( $\mathcal{M}$ ) (8, 9). Second, we evaluate how

31 the brain scores systematically vary with semantic comprehen-  
32 sion, as individually assessed by a questionnaire at the end of  
33 each story.

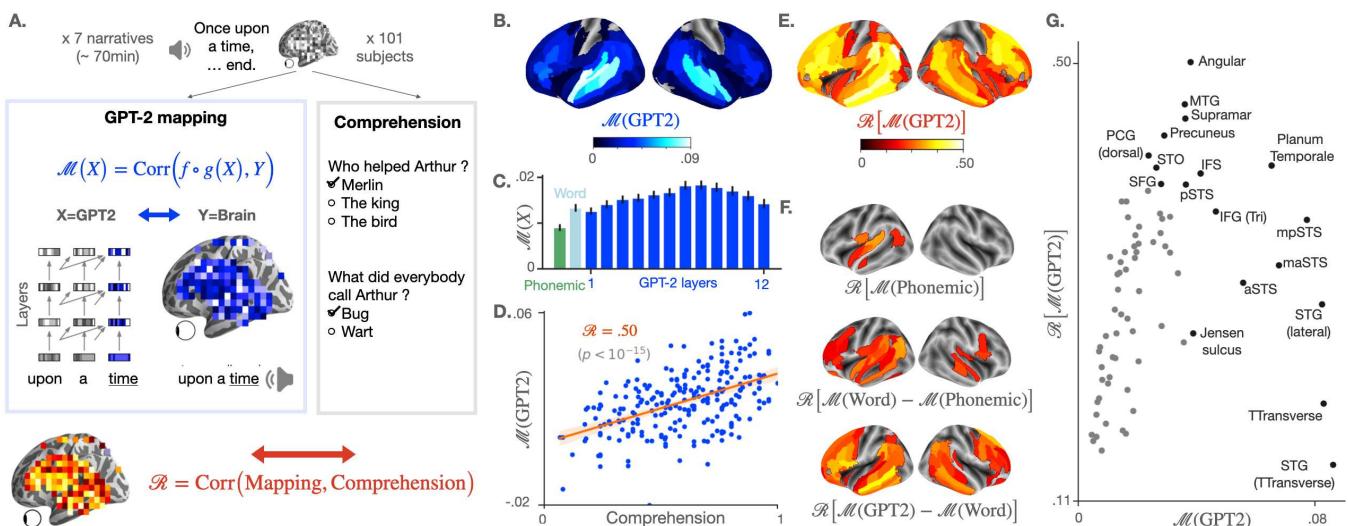
34 **GPT-2's activations linearly map onto fMRI responses to spoken nar-  
35 ratives.** To assess whether GPT-2 generates similar represen-  
36 tations to those of the brain, we first evaluate, for each voxel,  
37 subject and narrative independently, whether the fMRI re-  
38 sponds can be predicted from a linear combination of GPT-2's  
39 activations (Figure 1A). We summarize the precision of this  
40 mapping with a brain score  $\mathcal{M}$ : i.e. the correlation between  
41 the true fMRI responses and the fMRI responses linearly pre-  
42 dicted, with cross-validation, from GPT-2's responses to the  
43 same narratives (cf. Methods). To mitigate fMRI spatial  
44 resolution and the necessity to correct each observation by  
45 the number of statistical comparisons, we here report either 1)  
46 the average brain scores across voxels or 2) the average score  
47 within each region of interest ( $n = 314$ , following an automatic  
48 subdivision of Destrieux atlas (10), cf. SI.1). Consistent with  
49 previous findings (2, 4, 11, 12), these brain scores are signif-  
50 icant over a distributed and bilateral cortical network, and  
51 peak in middle- and superior-temporal gyri and sulci, as well  
52 as in the supra-marginal and the infero-frontal cortex (2, 4, 11)  
53 (Figure 1B).

54 By extracting GPT-2 activations from multiple layers (from  
55 layer one to layer twelve), we confirm that middle layers best  
56 map onto the brain (Figure 1C), as seen in previous studies  
57 (2, 4, 11). For clarity, the following analyses focus on the  
58 activations extracted from the *eighth* layer, i.e. GPT-2's most  
59 "brain-like" layer (Figure 1B).

60 **GPT-2's brain predictions correlate with semantic comprehension.**  
61 Does the linear mapping between GPT-2 and the brain reflect  
62 a fortunate correspondence (4)? Or, on the contrary, does  
63 it reflect similar representations of high-level semantics? To  
64 address this issue, we correlate these brain scores to the level of  
65 comprehension of the subjects, assessed for each subject-story  
66 pair. On average across all voxels, this correlation reaches  
67  $R = 0.50$  ( $p < 10^{-15}$ , Figure 1D, as assessed across subject-  
68 story pairs with the Pearson's test provided by SciPy). This  
69 correlation is significant across a wide variety of the bilateral  
70 temporal, parietal and prefrontal cortices typically linked to  
71 language processing (Figure 1E). Together, these results sug-  
72 gest that the shared representations between GPT-2 and the  
73 brain reliably vary with semantic comprehension.

74 **Low-level processing only partially accounts for the correlation be-  
75 tween comprehension and GPT-2's mapping** Low-level speech  
76 representations typically vary with attention (13, 14), and  
77 could thus, in turn, influence down-stream comprehension  
78 processes. Consequently, one can legitimately wonder whether

\*as assessed using Huggingface interface (<https://github.com/huggingface/transformers>) and  
GPT-2 pretrained model with temperature=0.



**Fig. 1. A.** 101 subjects listen to narratives (70 min of unique audio stimulus in total) while their brain signal is recorded using functional MRI. At the end of each story, a questionnaire is submitted to each subject to assess their understanding, and the answers are summarized into a comprehension score specific to each (narrative, subject) pair (grey box). In parallel (blue box on the left), we measure the mapping between the subject's brain activations and the activations of GPT-2, a deep network trained to predict a word given its past context, both elicited by the same narrative. To this end, a linear spatio-temporal model ( $f \circ g$ ) is fitted to predict the brain activity of one voxel  $Y$ , given GPT-2 activations  $X$  as input. The degree of mapping, called "brain score" is defined for each voxel as the Pearson correlation between predicted and actual brain activity on held-out data (blue equation, cf. Methods). Finally, we test the correlation between the comprehension scores of the subjects and their corresponding brain scores using Pearson's correlation (red equation). A positive correlation means that the representations shared across the brain and GPT-2 are key for the subjects to understand a narrative.

**B.** Brain scores (fMRI predictability) of the activations of the eighth layer of GPT-2. Scores are averaged across subjects, narratives, and voxels within brain regions (142 regions in each hemisphere, following a subdivision of Destrieux Atlas (10), cf. SI.1). Only significant regions are displayed, as assessed with a two-sided Wilcoxon test across (subject, narrative) pairs, testing whether the brain score is significantly different from zero (threshold: .05).

**C.** Brain scores, averaged across fMRI voxels, for different activation spaces: phonological features (word rate, phoneme rate, phonemes, tone and stress, in green), the non-contextualized word embedding of GPT-2 ("Word", light blue) and the activations of the contextualized layers of GPT-2 (from layer one to layer twelve, in blue). The error bars refer to the standard error of the mean across (subject, narrative) pairs ( $n=237$ ).

**D.** Comprehension and GPT-2 brain scores, averaged across voxels, for each (subject, narrative) pair. In red, Pearson's correlation between the two (denoted  $\mathcal{R}$ ), the corresponding regression line and the 95% confidence interval of the regression coefficient.

**E.** Correlations ( $\mathcal{R}$ ) between comprehension and brain scores over regions of interest. Brain scores are first averaged across voxels within brain regions (similar to B.), then correlated to the subjects' comprehension scores. Only significant correlations are displayed (threshold: .05).

**F.** Correlation scores ( $\mathcal{R}$ ) between comprehension and the subjects' brain mapping with phonological features  $\mathcal{M}(\text{Phonemic})$  (i), the share of the word-embedding mapping that is not accounted by phonological features  $\mathcal{R}[\mathcal{M}(\text{Word}) - \mathcal{M}(\text{Phonemic})]$  (ii) and the share of the GPT-2 eighth layer's mapping not accounted by the word-embedding  $\mathcal{R}[\mathcal{M}(\text{GPT2}) - \mathcal{M}(\text{Word})]$  (iii).

**G.** Relationship between the average GPT-2-to-brain mapping (eighth layer) per region of interest (similar to B.), and the corresponding correlation with comprehension ( $\mathcal{R}$ , similar to D.). Only regions of the left hemisphere, significant in both B. and E. are displayed. In black, the top ten regions in terms of brain and correlation scores (cf. SI.1 for the acronyms). Significance in D, E and F is assessed with Pearson's p-value provided by SciPy<sup>†</sup>. In B, E and F, p-values are corrected for multiple comparison using a False Discovery Rate (Benjamini/Hochberg) over the  $2 \times 142$  regions of interest.

79 the correlation between comprehension and GPT-2's brain  
 80 mapping is simply driven by variations in low-level auditory  
 81 processing. To address this issue, we evaluate the predictabil-  
 82 ity of fMRI given low-level phonological features: the word  
 83 rate, phoneme rate, phonemes, stress and tone of the narrative  
 84 (cf. Methods). The corresponding brain scores correlate with  
 85 the subjects' understanding ( $\mathcal{R} = 0.17, p < 10^{-2}$ ) but less so  
 86 than the brain scores of GPT-2 ( $\Delta\mathcal{R} = 0.32$ ). These low-level  
 87 correlations with comprehension peak in the left superior tem-  
 88 poral cortex (Figure 1F). Overall, this result suggests that the  
 89 link between comprehension and GPT-2's brain mapping may  
 90 be partially explained by – but not reduced to – the variations  
 91 of low-level auditory processing.

92 **The reliability of high-level representations best predict compre-  
 93 hension** Is the correlation between comprehension and GPT-2's  
 94 mapping driven by a *lexical* process and/or by an ability to  
 95 meaningfully *combine* words? To tackle this issue, we compare  
 96 the correlations obtained from GPT-2's word embedding (i.e.  
 97 layer 0) to those obtained from GPT-2's eighth layer, i.e. a  
 98 contextual embedding. On average across voxels, the corre-  
 99 lation with comprehension is 0.12 lower with GPT-2's word  
 100 embedding than with its contextual embedding. An analogous  
 101 analysis, comparing word embedding to phonological features

is displayed in 1F. Strictly lexical effects (word-embedding  
 102 versus phonological) peak in the superior-temporal lobe and  
 103 in pars triangularis. By contrast, higher-level effects (GPT-2  
 104 eighth layer versus word-embedding) peak in the superior-  
 105 frontal, posterior superior-temporal gyrus, in the precuneus  
 106 and in both the triangular and opercular parts of the inferior  
 107 frontal gyrus – a network typically associated with high-level  
 108 language comprehension (4, 15–19).

**Comprehension effects are mainly driven by individuals' variability**  
 110 The variability in comprehension scores could result from  
 111 exogeneous factors (e.g. some stories may be harder to com-  
 112 prehend than others for GPT-2) and/or from endogeneous  
 113 factors (e.g. some subjects may better understand specific  
 114 texts because of their prior knowledge). To address this issue,  
 115 we fit a linear mixed model to predict comprehension scores  
 116 given brain scores, specifying the narrative as a random effect  
 117 (cf. SI.1). The fixed effect of brain score (shared across nar-  
 118 ratives) is highly significant:  $\beta = 0.04, p < 10^{-29}$ , cf. SI.1).  
 119 However, the random effect (slope specific to each single nar-  
 120 rative) is not ( $\beta < 10^{-2}, p > 0.11$ ). We also replicate the  
 121 main analysis (Figure 1D) within each single narrative: the  
 122 correlation with comprehension reaches 0.76 for the 'sherlock'  
 123 story and is above 0.40 for every story (cf. SI.1). Overall,  
 124

125 these analyses confirm that the link between GPT-2 and semantic  
126 comprehension is mainly driven by subjects' individual  
127 differences in their ability to make sense of the narratives.

128 **Discussion** Our analyses reveal a positive correlation between  
129 semantic comprehension and the degree to which GPT-2 maps  
130 onto brain responses to spoken narratives.

131 These results strengthen and complete prior work on the  
132 brain bases of semantic comprehension. In particular, previous  
133 studies have used inter-subject brain correlation to reveal the  
134 brain regions associated with understanding (17). For example,  
135 Lerner et al. recorded subjects' fMRI while they listened  
136 to normal texts or texts scrambled at the word, sentence or  
137 paragraph level, in order to parametrically manipulate their  
138 level of comprehension (15). The corresponding fMRI signals  
139 correlated across subjects in the primary and secondary auditory  
140 areas even when the input was scrambled below the lexical  
141 level. By contrast, fMRI signals also became correlated in the  
142 bilateral infero-frontal and temporo-parietal cortex when the  
143 scrambling was either not performed, or performed at the level of  
144 sentences and paragraphs. Our results are consistent with  
145 this hierarchical organization, and thus make an important  
146 step towards the development of a cerebral model of narrative  
147 comprehension.

148 The relationship between GPT-2's representations and human  
149 comprehension remains to be qualified. First, although  
150 highly significant, our brain scores are relatively low (2, 9, 17).  
151 This phenomenon likely results from a mixture of different  
152 elements: i) we ran our analyses across *all* voxels to avoid  
153 selection biases, which automatically reduces the average effect  
154 sizes and ii) we report the results without correcting for  
155 a noise ceiling (cf. SI.1), as our pilot analyses suggest that  
156 such noise-ceiling can greatly vary depending on how it is  
157 implemented (i.e. fit from mean across subjects, from all or on  
158 voxels etc). Second, the correlation between semantic comprehension  
159 and GPT-2's mapping is robust ( $p < 10^{-15}$ ) but far from perfect ( $R = 0.50$ ). Such correlation thus indicates that  
160 the modeling of brain responses with GPT-2 does not *fully*  
161 account for the variation in comprehension. While this result  
162 is expected (7), our study provides a promising framework to  
163 evaluate the extent to which deep language models represent  
164 and understand texts like we do.

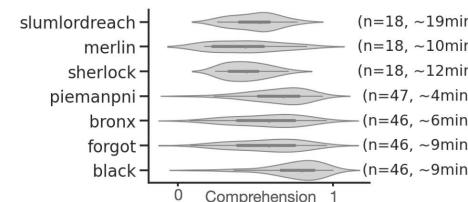
166 Finally, our results suggest that the neural bases of comprehension  
167 relate to the *high-level* representations of deep language models.  
168 While the mapping of phonological features and word embeddings do correlate with comprehension,  
169 GPT-2's contextual embeddings provides brain maps that  
170 more reliably predict comprehension (Figure 1F). The superiority  
171 of contextual-embedding in predicting comprehension suggests that i) GPT-2 encodes features supporting  
172 comprehension and ii) our finding are not solely driven by low- or  
173 mid-level processing (13, 14). These elements remain solely  
174 based on correlations, however. The factors that *causally* influence  
175 comprehension, ranging from prior knowledge, attention and  
176 language complexity should be explicitly manipulated in  
177 future work.

178 Overall, the present study strengthens and clarifies the similarity  
179 between the brain and deep language models, repeatedly  
180 observed in the past three years (2–4, 11, 20). Together, these  
181 findings reinforce the relevance of deep language models in  
182 unraveling the neural bases of narrative comprehension.

## Materials and Methods

186 Our analyses rely on the "Narratives" dataset (21), composed of  
187 the brain signals, recorded using fMRI, of 345 subjects listening to  
188 27 narratives.

189 **Narratives and comprehension score** Among the 27 stories of the  
190 dataset, we selected the seven stories for which subjects were asked  
191 to answer a comprehension questionnaire at the end, and for which  
192 the answers varied across subjects (more than ten different com-  
193 prehension scores across subjects), resulting in 70 minutes of audio  
194 stimuli in total, from four to 19 minutes per story (Figure 2). Ques-  
195 tionnaires were either multiple-choice, fill-in-the-blank, or open  
196 questions (answered with free text) rated by humans (21). Here,  
197 we used the comprehension score computed in the original dataset  
198 which was either a proportion of correct answers or the sum of the  
199 human ratings, scaled between 0 and 1 (21). It summarizes the  
200 comprehension of one subject for one narrative (specific to each  
201 (narrative, subject) pair).



202 **Fig. 2.** For each  
203 of the seven nar-  
204 ratives: number of  
205 subjects ( $n$ ), dis-  
206 tribution of compre-  
207 hension scores across  
208 subjects and length  
209 of the narrative.

210 **Brain activations** The brain activations of the 101 subject who  
211 listened to the selected narratives were recorded using fMRI, as de-  
212 scribed in (21). As suggested in the original paper, pairs of (subject,  
213 narrative) were excluded because of noisy recordings, resulting in  
214 237 pairs in total.

215 **GPT-2 activations** GPT-2 (1) is a high-performing neural language  
216 model trained to predict a word given its previous context (it does  
217 not have access to succeeding words), given millions of examples  
218 (e.g. Wikipedia texts). It consists of multiple Transformer modules  
219 (twelve, each of them called "layer") stacked on a non-contextual  
220 word embedding (a look-up table that outputs a single vector per  
221 vocabulary word) (1). Each layer  $l$  can be seen as a nonlinear  
222 system that takes a sequence of  $w$  words as input, and outputs  
223 a contextual vector of dimension  $(w, d)$ , called the "activations"  
224 of layer  $l$  ( $d = 768$ ). Intermediate layers were shown to better  
225 encode syntactic and semantic information than input and output  
226 layers (22), and to better map onto brain activity (2, 4). Here, we  
227 show that the *eighth* layer of GPT-2 best predicts brain activity  
228 1C. We thus select the eighth layer of GPT-2 for our analyses.  
229 Our conclusions remain unchanged with other intermediate-to-deep  
230 layers of GPT-2 (from 6<sup>th</sup> to 12<sup>th</sup> layers).

231 In practice, the narratives' transcripts were formatted (replacing  
232 special punctuation marks such as "-" and duplicated marks "?" by  
233 dots), tokenized using GPT-2 tokenizer and input to the GPT-2  
234 pretrained model provided by Huggingface <sup>‡</sup>. The representation of  
235 each token is computed separately using a context window a 1024.  
236 For instance, to compute the representation of the third token of  
237 the story, we input GPT-2 with the third, second and first token,  
238 and then extract the activations corresponding to the third token.  
239 To compute the representation of a token  $w_k$  at the end of the  
240 story, GPT-2 is input with this token combined with the 1,023  
241 preceding tokens. Then, we extract the activations corresponding  
242 to  $w_k$ . The procedure results in a vector of activations of size  $(w, d)$   
243 with  $w$  the number of tokens in the story and  $d$  the dimensionality  
244 of the model. There are fewer fMRI scans than words. Thus,  
245 the activation vectors between successive fMRI measurements are  
246 summed to obtain one vector of size  $d$  per measurement. To match  
247 the fMRI measurements and the GPT-2 vectors over time, we used  
248 the speech-to-text correspondences provided in the fMRI dataset  
249 (21).

250 <sup>‡</sup><https://github.com/huggingface/transformers>

243 **Linear mapping between GPT-2 and the brain** For each (subject,  
244 narrative) pair, we measure the mapping between i) the fMRI  
245 activations elicited by the narrative and ii) the activations of GPT-2  
246 (layer nine) elicited by the same narrative. To this end, a linear  
247 spatiotemporal model is fitted on a train set to predict the fMRI  
248 scans given the GPT-2 activations as input. Then, the mapping is  
249 evaluated by computing the Pearson correlation between predicted  
250 and actual fMRI scans on a held out set  $I$ :

251 
$$\mathcal{M}^{(s,w)} : I \mapsto \mathcal{L} \left( f \circ g(X^{(w)})_{i \in I}, (Y_i^{(s,w)})_{i \in I} \right) \quad [1]$$

252 With  $f \circ g$  the fitted estimator ( $g$ : temporal and  $f$ : spatial  
253 mappings),  $\mathcal{L}$  Pearson's correlation,  $X^{(w)}$  the activations of GPT-2  
254 and  $Y^{(s,w)}$  the fMRI scans of subjects  $s$ , both elicited by the  
255 narrative  $w$ .

256 In practice,  $f$  is a  $\ell_2$ -penalized linear regression. We follow scikit-  
257 learn implementation<sup>§</sup> with ten possible regularization parameters  
258 log-spaced between  $10^{-1}$  and  $10^8$ , one optimal parameter per voxel  
259 and leave-one-out cross-validation.  $g$  is a finite impulse response  
260 (FIR) model with 5 delays, where each delay sums the activations  
261 of GPT-2 input with the words presented between two TRs. For  
262 each (subject, narrative) pair, we split the corresponding fMRI time  
263 series into five contiguous chunks using scikit-learn cross-validation.  
264 The procedure is repeated across the five train (80% of the fMRI  
265 scans) and disjoint test folds (20% of the fMRI scans). Pearson  
266 correlations are averaged across folds to obtain a single score per  
267 (subject, narrative) pair. This score, denoted  $\mathcal{M}(X)$  in Figure 1A,  
268 measures the mapping between the activations space  $X$  and the  
269 brain of one subject, elicited by one narrative.

270 **Phonological features** To account for low-level speech processing,  
271 we computed the alignment (Eq. (1)) between the fMRI brain recordings  
272  $Y$  and phonological features  $X$ : the word rate (of dimension  
273  $d = 1$ , the number of words per fMRI scan), the phoneme rate  
274 (of dimension  $d = 1$ , the number of phonemes per fMRI scan) and the concat-  
275 enation of phonemes, stresses and tones of the words in the stimuli  
276 (categorical feature,  $d = 117$ ). The latter features are provided in  
277 the original Narratives database (21), and computed using Gentle<sup>¶</sup>  
278 forced-alignment algorithm.

279 **Significance** Significance was either assessed by using either (i) a  
280 second-level Wilcoxon test (two-sided) across subject-narrative pairs,  
281 testing whether the mapping (one value per pair) was significantly  
282 different from zero (Figure 1B), or (ii) by using the first-level Pearson  
283 p-value provided by SciPy<sup>||</sup> (Figure 1D-G). In Figure 1B, E, F, p-  
284 values were corrected for multiple comparison ( $2 \times 142$  ROIs) using  
285 False Discovery Rate (Benjamini/Hochberg)\*\*.

## References

287 1. Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language Models are Unsupervised Multitask Learners. page 24, 2018.  
288 2. Mariya Toneva and Leila Wehbe. Interpreting and improving natural-language processing  
289 (in machines) with natural-language-processing (in the brain). *arXiv:1905.11833 [cs, q-bio]*, November 2019. arXiv: 1905.11833.  
290 3. Martin Schrimpf, Idan Blank, Greta Tuckute, Carina Kauf, Eghbal A. Hosseini, Nancy Kan-  
291 wisher, Joshua Tenenbaum, and Evelina Fedorenko. Artificial Neural Networks Accurately  
292 Predict Language Processing in the Brain. *bioRxiv*, page 2020.06.26.174482, June 2020. .  
293 Publisher: Cold Spring Harbor Laboratory Section: New Results.  
294 4. Charlotte Caucheteux and Jean-Rémi King. Language processing in brains and deep neural  
295 networks: computational convergence and its limits. *bioRxiv*, page 2020.07.03.186288, July  
296 2020. . Publisher: Cold Spring Harbor Laboratory Section: New Results.  
297 5. Ariel Goldstein, Zaid Zada, Elav Buchnik, Mariano Schain, Amy Price, Bobbi Aubrey,  
298 Samuel A. Nastase, Amir Feder, Dotan Emanuel, Alon Cohen, Aren Jansen, Harshvardhan  
299 Gazula, Gina Choe, Aditi Rao, Catherine Kim, Colton Castro, Fonda Lora, Adeem Flinker,  
300 Sasha Devore, Werner Doyle, Patricia Dugan, Daniel Friedman, Avinatan Hassidim, Michael  
301 Brenner, Yossi Matias, Ken A. Norman, Orrin Devinsky, and Uri Hasson. Thinking ahead:  
302 prediction in context as a keystone of language in humans and machines. *bioRxiv*, page  
303 2020.12.02.403477, January 2021. . Publisher: Cold Spring Harbor Laboratory Section:  
304 New Results.  
305 6. Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela.  
306 Adversarial nli: A new benchmark for natural language understanding. *arXiv preprint  
arXiv:1910.14599*, 2019.  
307 7. Gary Marcus. Gpt-2 and the nature of intelligence. *The Gradient*, 2020.  
308 8. D. L. K. Yamins, H. Hong, C. F. Cadieu, E. A. Solomon, D. Seibert, and J. J. DiCarlo.  
309 Performance-optimized hierarchical models predict neural responses in higher visual cortex.  
310 *Proceedings of the National Academy of Sciences*, 111(23):8619–8624, June 2014. ISSN  
311 0027-8424, 1091-6490. .  
312 9. Alexander G. Huth, Wendy A. de Heer, Thomas L. Griffiths, Frédéric E. Theunissen, and  
313 Jack L. Gallant. Natural speech reveals the semantic maps that tile human cerebral cortex.  
314 *Nature*, 532(7600):453–458, April 2016. ISSN 0028-0836, 1476-4687. .  
315 10. Christoph Destrieux, Bruce Fischl, Anders Dale, and Eric Halgren. Automatic parcellation  
316 of human cortical gyri and sulci using standard anatomical nomenclature. *NeuroImage*, 53  
317 (1):1–15, October 2010. ISSN 1053-8119. .  
318 11. Shailee Jain and Alexander G Huth. Incorporating Context into Language Encoding Models  
319 for fMRI. preprint, Neuroscience, May 2018.  
320 12. Martin Schrimpf, Jonas Kubilius, Ha Hong, Najib J. Majaj, Rishi Rajalingham, Elias B.  
321 Issa, Kohitij Kar, Pouya Bashivan, Jonathan Prescott-Roy, Franziska Geiger, Kailyn Schmidt,  
322 Daniel L. K. Yamins, and James J. DiCarlo. Brain-Score: Which Artificial Neural Network for  
323 Object Recognition is most Brain-Like? preprint, Neuroscience, September 2018.  
324 13. Nima Mesgarani and Edward F. Chang. Selective cortical representation of attended speaker  
325 in multi-talker speech perception. *Nature*, 485(7397):233–236, May 2012. ISSN 1476-4687.  
326 . Bandiera\_abtest: a Cg\_type: Nature Research Journals Number: 7397 Primary\_atype: Re-  
327 search Publisher: Nature Publishing Group Subject\_term: Auditory system;Neuronal physiol-  
328 ogy;Perception Subject\_term\_id: auditory-system;neuronal-physiology;perception.  
329 14. Laurent Cohen, Philippe Salondy, Christophe Pallier, and Stanislas Dehaene. How does  
330 inattention affect written and spoken language processing? *Cortex*, 138:212–227, 2021.  
331 15. Y. Lerner, C. J. Honey, L. J. Silbert, and U. Hasson. Topographic Mapping of a Hierarchy  
332 of Temporal Receptive Windows Using a Narrated Story. *Journal of Neuroscience*, 31(8):  
333 2906–2915, February 2011. ISSN 0270-6474, 1529-2401. .  
334 16. C. Pallier, A.-D. Devauchelle, and S. Dehaene. Cortical representation of the constituent  
335 structure of sentences. *Proceedings of the National Academy of Sciences*, 108(6):2522–  
336 2527, February 2011. ISSN 0027-8424, 1091-6490. .  
337 17. Evelina Fedorenko, Terri Scott, Peter Brunner, William Coon, Brianna Pritchett, Gerwin  
338 Schalk, and Nancy Kanwisher. Neural correlate of the construction of sentence meaning.  
339 *Proceedings of the National Academy of Sciences of the United States of America*, 113,  
340 September 2016. .  
341 18. Angela D. Friederici. The Brain Basis of Language Processing: From Structure to Function.  
342 *Physiological Reviews*, 91(4):1357–1392, October 2011. ISSN 0031-9333, 1522-1210. .  
343 19. Gregory Hickok and David Poeppel. The cortical organization of speech processing. *Nature  
Reviews Neuroscience*, 8(5):393–402, May 2007. ISSN 1471-0048. . Number: 5 Publisher:  
344 Nature Publishing Group.  
345 20. Jon Gauthier and Anna Ivanova. Does the brain represent words? An evaluation of brain  
346 decoding studies of language understanding. *arXiv:1806.00591 [cs]*, June 2018. arXiv:  
347 1806.00591.  
348 21. Samuel A. Nastase, Yun-Fei Liu, Hanna Hillman, Asieh Zadbood, Liat Hasenfratz, Neggin  
349 Keshavarzian, Janice Chen, Christopher J. Honey, Yaara Yeshenko, Mor Regev, Mai Nguyen,  
350 Claire H. C. Chang, Christopher Baldassano, Olga Lositsky, Erez Simony, Michael A. Chow,  
351 Yuan Chang Leong, Paula P. Brooks, Emily Micicche, Gina Choe, Ariel Goldstein, Tamara  
352 Vanderwal, Yaroslav O. Halchenko, Kenneth A. Norman, and Uri Hasson. Narratives: fMRI  
353 data for evaluating models of naturalistic language comprehension. preprint, Neuroscience,  
354 December 2020.  
355 22. Ganesh Jawahar, Benoît Sagot, and Djamel Seddah. What Does BERT Learn about the  
356 Structure of Language? In *Proceedings of the 57th Annual Meeting of the Association for  
357 Computational Linguistics*, pages 3651–3657, Florence, Italy, 2019. Association for Computational  
358 Linguistics. .  
359 23. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
360 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
361 24. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
362 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
363 25. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
364 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
365 26. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
366 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
367 27. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
368 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
369 28. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
370 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
371 29. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
372 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
373 30. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
374 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
375 31. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
376 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
377 32. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
378 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
379 33. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
380 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
381 34. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
382 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
383 35. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
384 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
385 36. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
386 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
387 37. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
388 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
389 38. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
390 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
391 39. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
392 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
393 40. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
394 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
395 41. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
396 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
397 42. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
398 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
399 43. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
400 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
401 44. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
402 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
403 45. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
404 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
405 46. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
406 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
407 47. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
408 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
409 48. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
410 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
411 49. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
412 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
413 50. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
414 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
415 51. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
416 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
417 52. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
418 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
419 53. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
420 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
421 54. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
422 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
423 55. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
424 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
425 56. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
426 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
427 57. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
428 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
429 58. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
430 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
431 59. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
432 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
433 60. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
434 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
435 61. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
436 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
437 62. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
438 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
439 63. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
440 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
441 64. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
442 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
443 65. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
444 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
445 66. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
446 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
447 67. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
448 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
449 68. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
450 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
451 69. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
452 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
453 70. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
454 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
455 71. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
456 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
457 72. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
458 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
459 73. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
460 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
461 74. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
462 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
463 75. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
464 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
465 76. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
466 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
467 77. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
468 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
469 78. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
470 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
471 79. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
472 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
473 80. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
474 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
475 81. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
476 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
477 82. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
478 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
479 83. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
480 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
481 84. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
482 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
483 85. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
484 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
485 86. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
486 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
487 87. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
488 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
489 88. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
490 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
491 89. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
492 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
493 90. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
494 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
495 91. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
496 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
497 92. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
498 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
499 93. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
500 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
501 94. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
502 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
503 95. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
504 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
505 96. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
506 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
507 97. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
508 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
509 98. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
510 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
511 99. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
512 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
513 100. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
514 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
515 101. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
516 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
517 102. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
518 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
519 103. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
520 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
521 104. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
522 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
523 105. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
524 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
525 106. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
526 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
527 107. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
528 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
529 108. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
530 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
531 109. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
532 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .  
533 110. G. Caucheteux, J. R. King, and U. Hasson. Language processing in the brain: A review of  
534 the fMRI literature. *bioRxiv*, page 2020.07.03.186288, July 2020. .

## 363 Supporting Information (SI)

364 **Brain parcellation.** In Figure 1B, E, and F, we used a subdivision  
 365 of the parcellation from Destrieux Atlas (10). Regions  
 366 with more than 400 vertices were split into smaller regions (so  
 367 that each regions contains less than 400 vertices). The original  
 368 parcellation consists of 75 regions per hemisphere. Our custom  
 369 parcellation consists in 142 regions per hemisphere.

370 In Figure 1G, we use the original parcellation for simplicity,  
 371 and the following acronyms:

Acronym	Definition
STG / STS	Superior temporal gyrus / sulcus
aSTS	Anterior STS
maSTS	Mid-anterior STS
mpSTS	Mid-posterior STS
pSTS	Posterior STS
Angular / Supramar	Angular / Supramarginal inferior parietal gyrus
MTG / MTS	Medial temporal gyrus / sulcus
SFG / SFS	Superior frontal gyrus / sulcus
IFG / IFS	Inferior frontal gyrus / sulcus
Tri / Op	Pars triangularis / opercularis (IFG)
TTransverse	Temporal transverse sulcus
PCG	Posterior cingulate gyrus
STO	Temporo-occipital lateral sulcus

372 **Mixed-effect model.** Not all subjects listened to the same stories.  
 373 To check that the  $\mathcal{R}$  scores (correlation between comprehension  
 374 and brain mapping) were not driven by the narratives  
 375 and questionnaires' variability, a linear mixed-effect model was  
 376 fit to predict the comprehension of a subject given its brain  
 377 mapping scores, specifying the narrative as a random effect.  
 378 More precisely, if  $w_i \in \mathbb{R}$  corresponds to the mapping scores  
 379 of the  $i^{th}$  subject that listened to the story  $w$ , and  $C_{w_i} \in \mathbb{R}$   
 380 refers to the comprehension scores, we estimate the fixed effect  
 381 parameters  $\tilde{\beta} \in \mathbb{R}$  and  $\tilde{\eta} \in \mathbb{R}$  (shared across narratives), and  
 382 the random effect parameter  $\beta_w \in \mathbb{R}$  and  $\eta_w \in \mathbb{R}$  (specific to  
 383 the narrative  $w$ ) such that:

$$C_{w_i} = (\tilde{\beta} + \beta_w) \times w_i + (\tilde{\eta} + \eta_w) + \epsilon_{w_i}$$

384 with  $\epsilon_{w_i}$  a vector of i.i.d normal errors with mean 0 and variance  $\sigma^2$ . In practice, we use the statsmodels<sup>††</sup> implementation  
 385 of linear mixed-effect models. Significance of the coefficients  
 386 were assessed with a t-test, as implemented in statsmodels.  
 387

388 **Replication across single narratives.** To further support that  
 389 the  $\mathcal{R}$  were not driven by the narratives' variability, we repli-  
 390 cate the analysis of Figure 1D within single narratives. In  
 391 Figure 3, we show that correlation scores between brain scores  
 392 and comprehension scores are positive for each of the seven  
 393 narratives.  
 394

395 **Noise Ceiling Estimates.** fMRI recordings are inherently noisy.  
 396 Thus, we estimate an upper bound of the best brain score that  
 397 can be obtained given the level of noise in the Narrative dataset.  
 398 To this end, for each (subject, narrative) pair, we linearly  
 399 map the fMRI recordings, not with the GPT-2 activations,  
 400 but with the average fMRI recordings of the other subjects  
 401 who listened to that narrative. More precisely, we use the

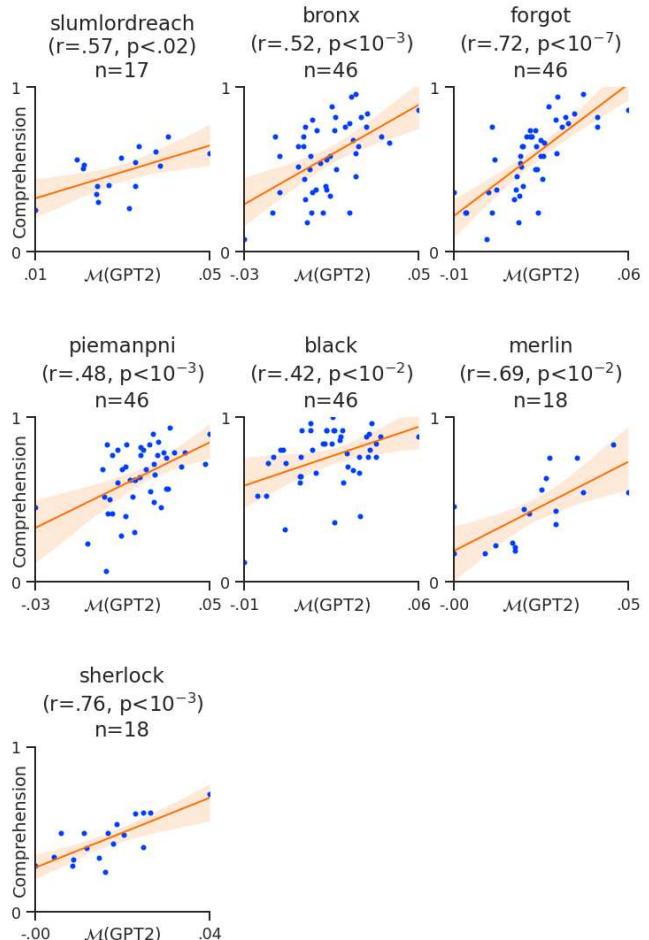
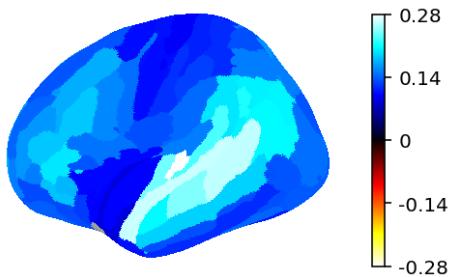


Fig. 3. Replication within single narratives. Same as Figure 1D for each single narrative.

402 exact same setting as in Eq. (1), but we predict  $Y^{(s)}$ , not  
 403 from  $g(X)$  (GPT-2's features after temporal alignment, of size  
 404  $n_{\text{times}} \times n_{\text{dim}}$ ), but from the mean of the other subject's brains  
 405  $\bar{Y} = \frac{1}{|S|} \sum_{s' \neq s} Y^{(s')}$  (of size  $n_{\text{times}} \times n_{\text{voxels}}$ ). This score is  
 406 called the noise ceiling for the (subject, narrative) pair. The  
 407 noise ceilings for each brain region are displayed in Figure 4,  
 408 and correspond to upper bounds of the brain scores displayed  
 409 in Figure 1B.

<sup>††</sup><https://www.statsmodels.org/>



**Fig. 4. Noise ceiling estimates.** Noise ceilings averaged across subjects, narratives and voxels within each region of interest. They are upper bounds of the brain scores in Figure 1B.