
What Are Large Language Models Mapping to in the Brain? A Case Against Over-Reliance on Brain Scores

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Given the remarkable capabilities of large language models (LLMs), there has
2 been a growing interest in evaluating their similarity to the human brain. One
3 approach towards quantifying this similarity is by measuring how well a model
4 predicts neural signals, also called "brain score". Internal representations from
5 LLMs achieve state-of-the-art brain scores, leading to speculation that they share
6 computational principles with human language processing. This inference is only
7 valid if the subset of neural activity predicted by LLMs reflects core elements
8 of language processing. Here, we question this assumption by analyzing three
9 neural datasets used in an impactful study on LLM-to-brain mappings, with a
10 particular focus on an fMRI dataset where participants read short passages. We
11 first find that when using shuffled train-test splits, as done in previous studies
12 with these datasets, a trivial feature that encodes temporal autocorrelation not only
13 outperforms LLMs but also accounts for the majority of neural variance that LLMs
14 explain. We therefore caution against shuffled train-test splits, and use contiguous
15 test splits moving forward. Second, we explain the surprising result that untrained
16 LLMs have higher-than-expected brain scores by showing they do not account
17 for additional neural variance beyond two simple features: sentence length and
18 sentence position. This undermines evidence used to claim that the transformer
19 architecture biases computations to be more brain-like. Third, we find that brain
20 scores of trained LLMs on this dataset can largely be explained by sentence
21 position, sentence length, and static word vectors; a small, additional amount is
22 explained by sense-specific word embeddings and contextual representations of
23 sentence structure. We conclude that over-reliance on brain scores can lead to
24 over-interpretations of similarity between LLMs and brains, and emphasize the
25 importance of deconstructing what LLMs are mapping to in neural signals.

26 1 Introduction

27 Recent developments in large language models (LLMs) have led many to wonder whether LLMs
28 process language like humans do. Whereas LLMs acquire many abstract linguistic generalizations, it
29 remains unclear to what extent their internal machinery bears resemblance to the human brain [1]. A
30 number of studies have attempted to answer this question through the framework of neural encoding
31 [2–4]. Within this framework, an LLM’s internal representations of some linguistic stimuli are used
32 to predict brain activity during comprehension of the same stimuli. Results have been uniformly
33 positive, showing that LLM representations are highly effective at predicting neural signals [5, 6].

34 In one impactful study, authors evaluated the brain scores of 43 models on three neural datasets [2].
35 They found that GPT2-XL [7] achieved the highest brain score and, in one neural dataset, accounted
36 for 100% of the "explainable" neural variance (i.e., taking into account the noise inherent in the data)

37 [8]. This result was interpreted as evidence that the brain may be optimizing for the same objective
38 as GPT2, namely, next-word prediction. Surprisingly, the authors further found that untrained (i.e.
39 randomly initialized) LLMs predict neural activity well, leading to speculations that the transformer
40 architecture biases computations to be more brain-like. The finding that untrained LLMs predict
41 neural signals significantly above chance has been replicated in other studies [9, 4, 10].

42 More generally, many studies have compared models to brain activity and concluded that high
43 prediction performance reveals correspondence between some interesting aspect of the model and
44 biological linguistic processing [4, 11–14]. One issue with this approach is that it assumes that the
45 subset of neural activity predicted by a model reflects core processes of the human language system
46 [15]. However, this assumption is not necessarily true. For example, a recent paper found that, when
47 participants listen to stories, the fMRI signal includes an initial ramping, positional artifact [16].
48 It is likely that LLMs which contain absolute positional embeddings would be able to predict this
49 ramping signal, whereas a simpler model such as a static word embedding (e.g. GloVe, [17]) would
50 not, leading to exaggerated differences between LLMs and GloVe due to reasons of little theoretical
51 interest. This issue relates to a more general trend in machine learning research: a complex algorithm
52 solves a task, but it is later discovered that the key innovation was a very simple component of the
53 algorithm [18]. Analogous to Weinberger [18], without attempting to rigorously deconstruct the
54 mapping between LLMs and brains, it is possible to draw erroneous conclusions about the brain’s
55 mechanisms for processing language.

56 We analyze the same three neural datasets used in [2]. These include the Pereira fMRI dataset, where
57 participants read short passages [8]; the Fedorenko electrocorticography (ECoG) dataset, where
58 participants read isolated sentences [19]; and the Blank fMRI dataset, where participants listened to
59 short stories [20]. As in Schrimpf et al. [2], we focus our analyses on the Pereira dataset. In order to
60 deconstruct the mapping between LLMs and the brain, we follow Reddy and Wehbe [21] and de Heer
61 et al. [22] by building a set of predictors that describe simple features of the linguistic input, and
62 gradually add features that increase in complexity. Our goal is to find the simplest set of features
63 which account for the greatest portion of the mapping between LLMs and brains.

64 2 Methods

65 2.1 Experimental data

66 For all three neural datasets, we used the same version as used by [2]. For additional details, refer to
67 A.1.

68 **Pereira (fMRI):** The Pereira dataset is composed of two experiments. Experiment 1 (EXP1) consists
69 of 96 passages each containing 4 sentences, with $n = 9$ participants. Experiment 2 (EXP2) consists of
70 72 passages each consisting of 3 or 4 sentences, with $n = 6$ participants. Passages in each experiment
71 were evenly divided into 24 semantic categories which were not related across experiments (4
72 passages per category in EXP1, and 3 passages per category in EXP2). A single fMRI scan (TR)
73 was taken after visual presentation of each sentence. Unless otherwise noted, we focus our results
74 on voxels from within the "language network" in the main paper. EXP1 was a 384×92450 matrix
75 (number of sentences \times number of voxels) and EXP2 was a 243×60100 matrix. All analyses were
76 conducted separately for each experiment.

77 **Fedorenko (ECoG):** Participants ($n = 5$) read 52 sentences of length 8 words. A total of 97
78 language-responsive electrodes were used across 5 participants: 47, 8, 9, 15, and 18, for participants
79 1 through 5, respectively. Neural activity was temporally averaged across the full presentation of each
80 word after extracting high gamma, and the entire dataset was a 416×97 matrix.

81 **Blank (fMRI):** The dataset consisted of 5 participants listening to 8 stories from the publicly
82 available Natural Stories Corpus [23]. An fMRI scan was taken every 2 seconds, resulting in a total
83 of 1317 TRs across the 8 stories. fMRI BOLD signals were averaged across voxels within each
84 functional region of interest (fROI). There were 60 fROIs across all 5 participants, resulting in a
85 1317×60 matrix.

86 2.2 Language models

87 We focus our analyses on GPT2-XL [7], as it was shown to be the best-performing model on the
88 Pereira dataset [10, 24, 2]. GPT2 is an auto-regressive transformer model, meaning that it can
89 only attend to current and past inputs, trained on next token prediction. The XL variant has $\sim 1.5\text{B}$
90 parameters and 48 layers. We replicate some of our key findings on Pereira with RoBERTa-Large[25]
91 (A.6). RoBERTa is a transformer model with bidirectional attention trained on masked token
92 prediction, meaning that it can attend to past and future tokens. The large variant contains 335M
93 parameters and 24 layers. Both GPT2 and RoBERTa use learned absolute positional embeddings,
94 such that a unique vector corresponding to each token position is added to the input static embeddings.

95 2.3 LLM feature pooling

96 **Pereira:** Each sentence was fed into an LLM, with previous sentences from the same passage also fed
97 as input. Since each fMRI scan was taken at the end of the sentence, we converted LLM token-level
98 embeddings to sentence-level embeddings by summing across all tokens within a sentence (sum
99 pooling). We used the sum pooling method because it is consistent with other neural encoding studies
100 [26, 27], and it performed better than taking the representation at the last token which was done in
101 [2] A.5.

102 **Fedorenko:** The current and previous tokens from within the same sentence were fed into the LLM
103 as context. We converted LLM token-level embeddings to word embeddings, since each word has a
104 neural response, by summing across tokens in multi-token words, and leaving single token words
105 unmodified.

106 **Blank:** For each story, we fed the current and all preceding tokens up to a maximum context size of
107 512 tokens. As in Schrimpf et al. [2], for each TR, we took the representation of the word that was
108 closest to being 4 seconds before the TR. For multi-token words, we took the representation of the
109 last token of that word.

110 2.4 Banded ridge regression

111 We used ridge regression (linear regression with an L2 penalty) to predict activations for each
112 voxel/electrode/fROI independently. We did not use "vanilla" ridge regression because it applies a
113 single L2 penalty for all weights, whereas our analyses use multiple sets of distinct features. In such
114 a case, a single penalty causes the regression will be biased against small feature spaces. Moreover,
115 different L2 penalties are likely optimal for each feature space. To remedy this, we employed banded
116 ridge regression which effectively allows a different L2 penalty to be applied to each feature space
117 [28] (for further details, refer to A.2).

118 2.5 Out of sample R^2 metric

119 We define the brain score of a model as the out-of-sample R^2 metric (R_{oos}^2) [29]. R_{oos}^2 quantifies
120 how much better a set of features performs at predicting held-out data compared to a model which
121 simply predicts the mean of the training data (i.e. a regression with only an intercept term). To be
122 precise, given mean squared error (MSE) values from a model using features M and MSE values
123 from an intercept only regression (I), then:

$$R_{oos}^2 = 1 - \frac{MSE_M}{MSE_I}. \quad (1)$$

124 A positive (negative) value indicates that M was more (less) helpful than predicting the mean of
125 training data. We elected to use R_{oos}^2 over the standard R^2 because of this clear interpretation
126 and because it is a less biased estimate of test set performance [29]. We use R_{oos}^2 over Pearson's
127 correlation coefficient (r) because R_{oos}^2 can be interpreted as the fraction of variance explained,
128 which lends more straightforwardly to estimating how much variance one feature space explains
129 over others. Whenever averaging across voxels, we set R_{oos}^2 values to be non-negative to prevent
130 differences in performance on noisy voxels/electrodes/fROIs from significantly impacting the results.
131 We refer to R_{oos}^2 as R^2 throughout the rest of the paper for brevity, and use the notation R_M^2 to refer
132 to the performance of features M .

133 2.6 Selection of best layer

134 We evaluate the R^2 for each LLM layer, and select the layer that performs best across vox-
135 els/electrodes/fROIs. Due to the stochastic nature of untrained LLMs, we selected the best layer for
136 10 random seeds and computed the average R^2 across seeds. When reporting the best layer, we refer
137 to layer 0 as the input static layer, and layer 1 as the first intermediate layer.

138 2.7 Train, validation, and test folds:

139 For each dataset, we construct contiguous train-test splits by ensuring neural data from the same
140 passage/sentence/story is not included in both train and test data. Due to low sample sizes, we
141 employed a nested cross-validation procedure for each dataset (A.3). When computing R^2 across
142 inner or outer folds, we pooled predictions across folds and computed a single R^2 as recommended by
143 Hawinkel et al. [29]. The optimal parameters for banded regression were selected based on validation
144 data.

145 We created shuffled train-test splits, as done in [2], of the same size as the contiguous train-test splits.
146 Unless explicitly noted, all results are performed using contiguous train-test splits.

147 2.8 Correcting for decreases in test-set performance due to addition of feature spaces

148 It is possible for a "full" encoding model to perform worse than a "sub-model" (which consists
149 of only a subset of the predictors) because we are evaluating performance on a held-out test set
150 [22]. To address this problem, in some analyses we select the best performing sub-model for each
151 voxel/electrode/fROI which includes a given feature of interest. For instance, to examine how much
152 feature space C adds onto features spaces A and B , we select the best sub-model which includes C
153 and denote it as $A + B + C^*$. More precisely, the R^2 of $A + B + C^*$ is:

$$R_{A+B+C^*}^2 = \max(R_C^2, R_{A+C}^2, R_{B+C}^2, R_{A+B+C}^2). \quad (2)$$

154 2.9 Orthogonal Auto-correlated Sequences Model (OASM)

155 To model temporal auto-correlation in neural activity, we construct a feature matrix for each dataset
156 by (i) forming an n -dimensional identity matrix, where n is the total number of time points in the
157 dataset (per voxel / electrode / TR), and (ii) applying a Gaussian filter within "chunks" along the
158 diagonal that correspond to temporally contiguous time points (i.e., within each passage in Pereira,
159 each sentence in Fedorenko, and each story in Blank). This generates an auto-correlated sequence for
160 each passage/sentence/story that is orthogonal to that of each other passage/sentence/story (A.7).

161 3 Pereira dataset

162 3.1 Shuffled train-test splits are severely affected by temporal auto-correlation

163 Prior LLM encoding studies using this dataset [24, 2, 10, 30, 11] used shuffled train-test splits. Here,
164 we demonstrate that this approach compromises the evaluation of the neural predictivity of LLMs.
165 First, we replicated the pattern of neural predictivity across GPT2-XL's layers reported in [2] and [24]
166 when using shuffled splits. Using this procedure, early and late layers perform best and intermediate
167 layers perform worst. Strikingly, when using the alternative approach of contiguous train-test splits,
168 the opposite pattern is observed: intermediate layers perform best. Across layers, neural predictivity
169 using the shuffled method is highly anti-correlated with neural predictivity using the contiguous
170 method ($r = -.929$ in EXP1, $r = -.764$ in EXP2) (Fig. 1a).

171 Next, we hypothesized that much of what LLMs might be mapping to when using shuffled splits
172 could be accounted for by OASM, a model which only represents within passage auto-correlation
173 and between passage orthogonality. OASM out-performed GPT2-XL on both EXP1 and EXP2
174 (Fig. 1b, blue and red bars), revealing that a completely non-linguistic feature space can achieve
175 absurdly high brain scores in the context of shuffled splits. This strongly challenges the assumption
176 of multiple previous studies [2, 11, 10] that performance on this benchmark is an indication of a
177 model's brain-likeness, .

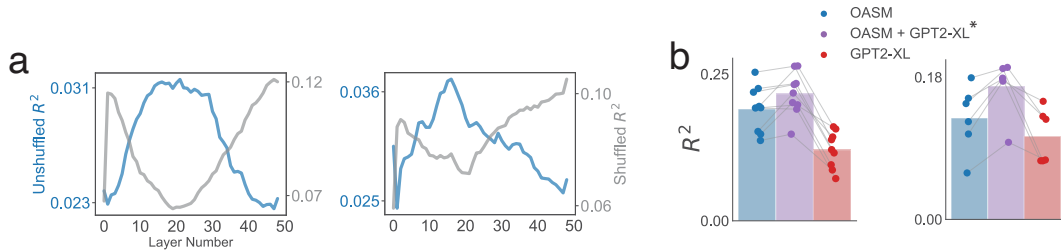


Figure 1: Comparing different approaches for creating train-test splits in the Pereira dataset. Within each panel, EXP1 results are on the left and EXP2 results are on the right (same formatting in Figure 2,3) (a) R^2 values across layers for GPT2-XL on shuffled train-test splits (gray) and contiguous (unshuffled) splits (blue). (b) Each dot shows the mean R^2 value across voxels within a participant, with bars indicating mean R^2 across participants.

178 Moreover, we find that the unique neural variance that GPT2-XL explains over OASM is very small
 179 relative to what OASM explains alone. To calculate this, we combine OASM with GPT2-XL and
 180 observe how much neural variance they explain together. To prevent OASM from ever weakening
 181 the reported performance of GPT2-XL for any voxel, we correct the R^2 value for each voxel with
 182 the OASM+GPT2-XL model to be at least as high as with GPT2-XL alone (denoted OASM+GPT2-
 183 XL*) (2.8). Even with these corrections, we find that $R^2_{OASM+GPT2-XL^*}$ was 13.6% higher than
 184 R^2_{OASM} in EXP1, and 31.5% higher than R^2_{OASM} in EXP2 (Fig. 1b) (% differences after averaging
 185 R^2 across participants). To be clear, this means that any linguistically-driven neural variance that
 186 GPT2-XL uniquely explains over OASM is far smaller (13.6% on EXP1 and 31.5% on EXP2) than
 187 what is predicted solely by OASM, a model with no linguistic features that completely lacks the
 188 ability to generalize to fully held out passages. Thus, it appears that the largest determinant of
 189 model predictivity on this dataset when using shuffled train-test splits is whether a model contains
 190 autocorrelated sequences within passages that are orthogonal between passages.

191 3.2 Untrained LLM neural predictivity is fully accounted for by sentence length and position

192 We next sought to deconstruct what explains the neural predictivity of untrained GPT2-XL (GPT2-
 193 XLU) in the Pereira dataset. We hypothesized that $R^2_{GPT2-XLU}$ could be explained by two simple
 194 features: sentence length (SL) and sentence position within the passage (SP). Sentence length is
 195 captured by GPT2-XLU because the GELU nonlinearity in the first layer’s MLP transforms normally
 196 distributed inputs with zero mean into outputs with a non-zero mean. This introduces a non-zero
 197 mean component to each token’s representation in the residual stream. When these representations
 198 are sum-pooled, this non-zero mean component accumulates in a way that reflects the sentence length,
 199 making the length decodable in the intermediate layers (see A.9 for a formal proof). Sentence position
 200 is encoded within GPT2-XLU due to absolute positional embeddings which, although untrained, still
 201 result in sentences at the same position having similar representations when tokens are sum-pooled.
 202 We represent sentence position as a 4-dimensional one-hot vector, where each element corresponds
 203 to a given position within a passage, and sentence length as the number of words in a passage.

204 To obtain representations from GPT2-XLU, we selected the best-performing layer for each of the 10
 205 untrained seeds. For EXP1 the best performing layer was layer 0 for 6 seeds, layer 1 for 3 seeds (first
 206 intermediate layer), and layer 19 for one seed. For EXP2 the best layer was layer 1 for 5 seeds, layer
 207 2 for 4 seeds, and layer 5 for 1 seed.

208 We fit a regression using all subsets of the following feature spaces, SL, SP, GPT2-XLU, resulting in
 209 7 models. For both experiments, R^2_{SP+SL} was descriptively higher than all other models, including
 210 the best-performing model with GPT2-XLU (SP+SL+GPT2-XLU) (Fig. 2a). Sentence position was
 211 particularly important in EXP1, and sentence length was particularly important in EXP2. This may
 212 explain why the static layer often outperformed intermediate layer representations in EXP1 despite
 213 encoding sentence length more poorly. Overall, these results suggest that, when averaging across
 214 voxels within the language network in this dataset, GPT2-XLU does not improve neural encoding
 215 performance over sentence length and position.

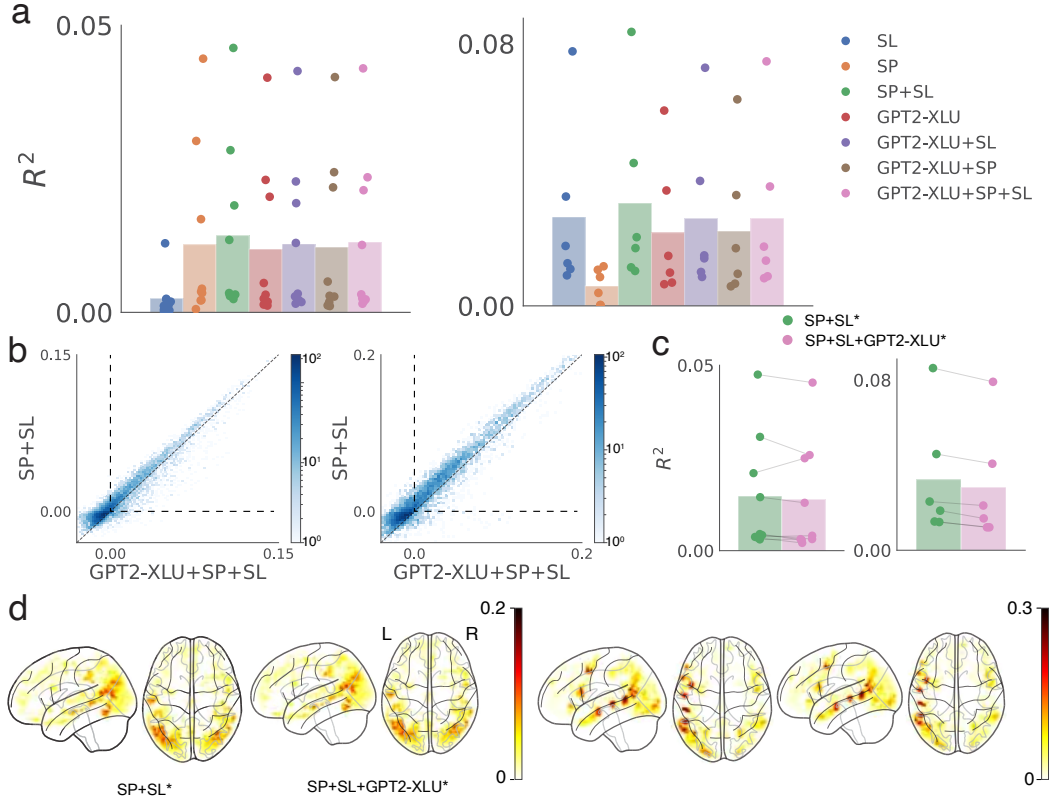


Figure 2: For all panels, EXP1 results are on the left and EXP2 results are on the right. **(a)** Brain score (R^2) for different combinations of features. Each dot represents R^2 values averaged across voxels in a single participant, with bars showing mean across participants. **(b)** 2D histogram of R^2 values for the best model without GPT2-XLU (SP+SL), and the best model with GPT2-XLU (GPT2-XLU+SP+SL). The dotted lines show $y = x$, $y = 0$, and $x = 0$. Values below $y = 0$ or left of $x = 0$ were clipped when averaging, but are shown here to visualize the full distribution. **(c)** Same as **(a)**, but after voxel-wise correction; lines connect data-points from the same participant. **(d)** Glass brain plots showing R^2 values of SP+SL (left) and GPT2-XLU+SP+SL (right) after voxel-wise correction. Conventions are the same as Figure 1.

216 Although GPT2-XLU did not enhance encoding performance when averaging across voxels, there
 217 may be a subset of voxels where GPT2-XLU does explain significant additional neural vari-
 218 ance. To examine this possibility, we plotted a 2D histogram of voxel-wise R^2_{SP+SL} values vs.
 219 $R^2_{SP+SL+GPT2-XLU}$ values in the language network (Fig. 2b). Values were clustered around the
 220 identity line, and there was no cluster of voxels where $R^2_{SP+SL+GPT2-XLU}$ appeared significantly
 221 higher. Next, for each voxel, we performed a one-sided paired t -test between the squared error
 222 values obtained over sentences (EXP1: $N = 384$, EXP2: $N = 243$) between SP+SL+GPT-XLU
 223 and SP+SL. Across all functional networks, only 1.26% (EXP1) and 1.42% (EXP2) of voxels were
 224 significantly ($\alpha = 0.05$) better explained by the GPT2-XLU model before false discovery rate
 225 (FDR) correction; these numbers dropped to 0.001% (EXP1) and 0.078% (EXP2) after performing
 226 FDR correction within each participant and network [31]. None of the significant voxels after FDR
 227 correction were inside the language network. Taken together, these results suggest GPT2-XLU does
 228 not enhance neural prediction performance over sentence length and position even at the voxel level.

229 To control for voxels where the neural encoding performance of GPT2-XLU is weakened by the
 230 addition of SP+SL, we compared SP+SL* and SP+SL+GPT2-XLU*. When averaging across voxels,
 231 $R^2_{SP+SL}^*$ still exceeded $R^2_{GPT2-XLU+SP+SL}^*$ (Fig. 2c). Furthermore, the values for $R^2_{SP+SL}^*$
 232 and $R^2_{GPT2-XLU+SP+SL}^*$ across brain areas were highly similar in both experiments (Fig. 2d).
 233 Only 1.00% (EXP1) and 1.18% (EXP2) of voxels were significantly better explained by the addition
 234 of GPT2-XLU before FDR correction; 0% (EXP1) and 0.05% (EXP2) of voxels were better explained

Table 1: Mean R^2 values (across participants) for each model. For models composed of multiple features, the best sub-model is used which includes the last feature.

Features	EXP1	EXP2
GPT2-XL	0.032	0.036
SP+SL	0.013	0.031
SP+SL+WORD	0.024	0.039
SP+SL+WORD+SENSE	0.026	0.040
SP+SL+WORD+SENSE+SYNT	0.027	0.043
SP+SL+WORD+SENSE+SYNT+GPT2-XL	0.032	0.045

235 after FDR correction (once again, no significant voxels were inside the language network). Thus, our
 236 results hold even when controlling for decreases in performance due to the addition of feature spaces.

237 3.3 Sentence length, sentence position, and static word embeddings account for the majority 238 of trained LLM encoding performance

239 We next turned to explaining the neural predictivity of the trained GPT2-XL. In addition to sentence
 240 position and sentence length, we added static word embeddings (WORD). Together, these features
 241 defined a baseline model which does not account for any form of linguistic processing of words
 242 in context. We next included three more complex features which involved contextual processing.
 243 First, we added sense-specific word embeddings from RoBERTa-Large using the LMMS package
 244 [32]. Sense embeddings contain distinct representations for different senses of the same word (e.g.,
 245 mouse: *computer device*, and mouse: *rodent*). LMMS generates sense embeddings by averaging over
 246 contextual embeddings corresponding to the same sense of a word (see A.10 for further details).

247 Whereas sense embeddings help disambiguate many content words, they do not disambiguate
 248 pronouns, i.e., do not encode the entities that they refer to. Therefore, our sense embeddings were
 249 generated for a version of the Pereira text where pronouns were dereferenced (i.e., replaced by
 250 the words that they referred to). To maintain consistency with these sense embeddings, our static
 251 word embeddings were created (1) by taking a frequency-weighted average of sense embeddings
 252 for the same word, where frequency values were obtained from WordNet [33]; and (2) based on the
 253 dereferenced Pereira texts. Importantly, this means the impact of pronoun dereferencing and word
 254 and sense embeddings are not decoupled in this study. Finally, we created an abstract representation
 255 of the syntax of each sentence (SYNT), using an approach highly similar to that of Caucheteux
 256 et al. [34]: we collected sentences that are syntactically equivalent but semantically dissimilar to the
 257 original sentence, and averaged their representations from the best layer of GPT2-XL (A.11). We
 258 selected the best layer based on averaged R^2 across language voxels on test data (EXP1: layer 21,
 259 EXP2: layer 16).

260 We fit a regression to the fMRI data using all subsets of the feature spaces SL+SP, WORD, SENSE,
 261 SYNT, GPT2-XL, resulting in 64 models. In this list, features are ranked from least to most complex.
 262 For each feature, we took the model that exhibited the best performance in the language network
 263 which included that feature but did not include features more complex than it. For instance, values
 264 reported for $R^2_{SL+SP+WORD+SENSE}$ were taken from the best model which included SENSE,
 265 excluding models which included SYNT and GPT2-XL. By doing so, we were able to examine
 266 the impact of adding more complex features in explaining $R^2_{GPT2-XL}$ while still accounting for
 267 decreases in test performance due to adding redundant features. We note that since this procedure is
 268 not performed at the voxel-level, we do not add a * to the R^2 notation.

269 Table 1 displays the performance of each model, including GPT2-XL on its own (Fig. 2a, 2b). The
 270 baseline SP+SL+WORD model, which does not account for any form of contextual processing,
 271 performs 75% as well as GPT2-XL in EXP1, and outperforms GPT2-XL in EXP2. When adding
 272 contextual features, namely SENSE and SYNT, our model performs 84.4% as well as GPT2-XL and
 273 the full model in EXP1, and better than GPT2-XL and 95.5% as well as the full model in EXP2,
 274 indicating that SENSE and SYNT play a modest role in accounting for GPT2-XL brain scores beyond
 275 simple features in this dataset.

276 Similar to previous sections, we perform voxel-wise correction by selecting the best sub-model with
 277 GPT2-XL and the best sub-model without GPT2-XL for each voxel. We focus only on sentence

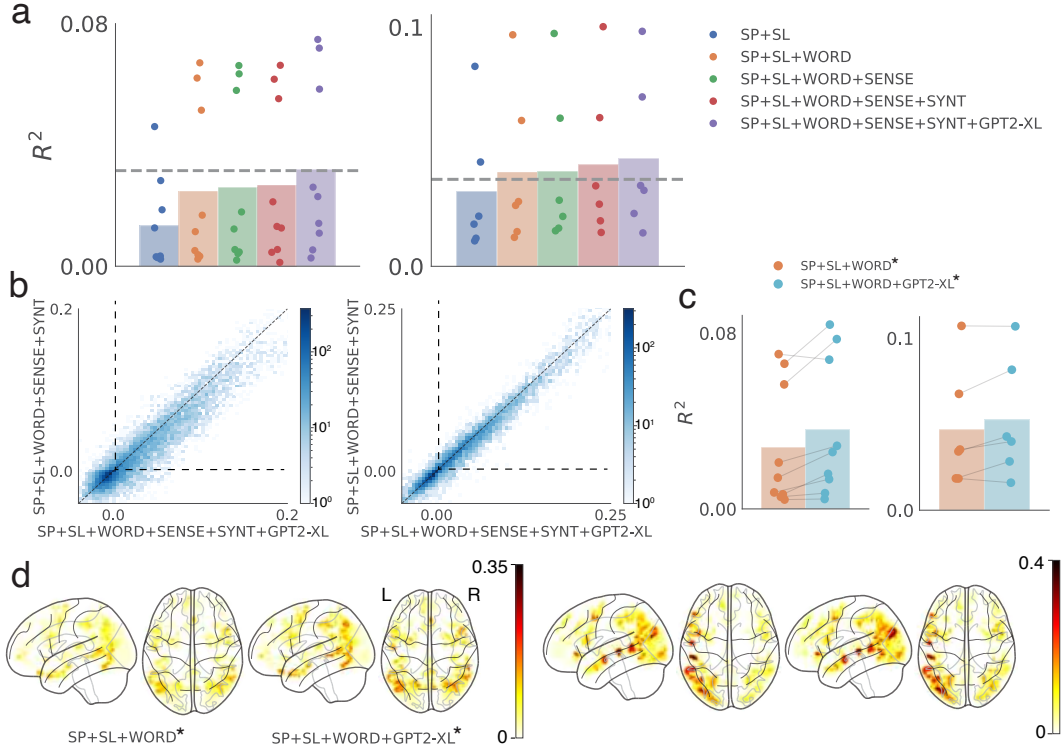


Figure 3: For all panels, EXP1 results are on the left and EXP2 results are on the right. (a) For each model, we display the sub-model which includes the added feature. Dots represent participants and bars are mean across participants. Grey dashed line is the performance of GPT2-XL alone. (b) 2d histogram comparing full model and full model with GPT2-XL. (c) Same as (a) but after voxel-wise correction for SP+SL+WORD and SP+SL+WORD+GPT2-XL. (d) Glass brain plots showing R^2 values of SP+SL+WORD (left) and SP+SL+WORD+GPT2-XLU (right) after voxel-wise correction.

278 position, sentence length, and static word embeddings because sense and syntax had modest con-
 279 tributions beyond these features. $R^2_{SP+SL+WORD*}$ was 0.028 in EXP1 and 0.048 in EXP2, and
 280 $R^2_{SP+SL+WORD+GPT2-XL*}$ was 0.036 in EXP1 and 0.056 in EXP2 (mean across participants)
 281 (Fig. 3c). This indicates that even after controlling for a reduction in GPT2-XL performance from
 282 the addition of simple features, GPT2-XL only explains an additional 28.57% (EXP1) and 16.7%
 283 (EXP2) neural variance over a model composed of features that are all non-contextual.

284 **4 Fedorenko dataset**

285 **4.1 Shuffled train-test splits also impact ECoG datasets, but less than with fMRI**

286 We first evaluated the impact of shuffled train-test splits on the Fedorenko dataset. Unlike in Pereira,
 287 the across-layer performance is well correlated between shuffled and contiguous splits ($r = 0.622$)
 288 (Fig. 4a). The OASM model performs 93.1% as well as GPT2-XL when averaging R^2 values across
 289 participants (Fig. 4b). $R^2_{OASM+GPT2-XL*}$ was 45.3% better than OASM, meaning that the unique
 290 contribution of GPT2-XL is less than half the total contribution of a simple, auto-correlated model.
 291 Therefore, shuffled train-test splits also impact results on Fedorenko, albeit less than Pereira. This
 292 may be due to lower autocorrelation of ECoG compared to fMRI. We use contiguous splits for the
 293 remainder of the Fedorenko analyses.

294 **4.2 Word position explains all of untrained, and most of trained, GPT2-XL brain score**

295 As noted in [35], there was a strong positional signal in the ECoG dataset during comprehension of
 296 sentences that is likely related to the construction of sentence meaning. We therefore hypothesized

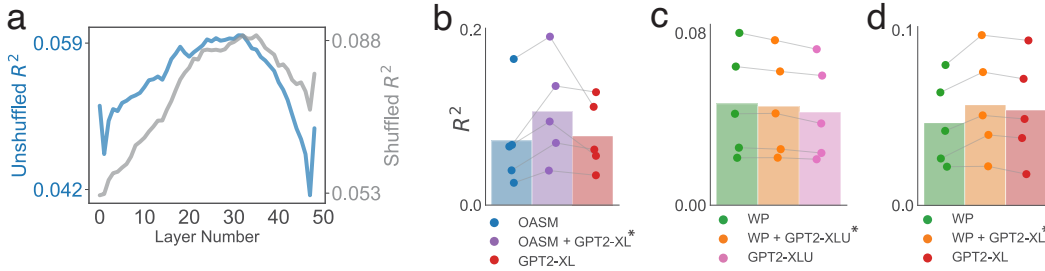


Figure 4: **(a)** Across-layer R^2 , averaged across electrodes in the Fedorenko dataset, for GPT2-XL with and without shuffled splits. **(b)** Each dot is a participant, lines connect data-points from the same participant. Bars display mean across participants. **(c)** and **(d)** Same guidelines as **(b)**.

297 that a feature space that accounted for word position (WP) would do well relative to untrained and
 298 trained GPT2-XL. We generated a simple feature space that encodes word position, such that words
 299 in nearby positions were given similar representations (A.12). When performing a one-sided paired
 300 t -test between the squared error predictions of WP+GPT2-XLU* and WP, three electrodes were
 301 significantly better explained by the addition of GPT2-XLU before FDR correction, and none were
 302 better explained after FDR correction within each participant. Moreover, WP performs 86.7% as well
 303 as GPT2-XL, and 82.1% as well as WP+GPT2-XL*. Our results therefore suggest that the mapping
 304 between GPT2-XL and neural activity on the Fedorenko dataset is largely driven by positional signals.
 305

306 5 Blank dataset is predicted at near chance levels

307 Lastly, we address the Blank dataset. We find that OASM achieves an R^2 that is 103.6 times
 308 larger than that of GPT2-XL when using shuffled splits A.13, demonstrating that such splits are
 309 massively contaminated by temporal autocorrelation. We next turn to using contiguous splits, and test
 310 whether GPT2-XL performs better than an intercept only model by applying a one-sided paired t -test
 311 between the squared error values obtained from GPT2-XL and the intercept only model ($N = 1317$
 312 TRs). GPT2-XL predicts 1 fROI significantly better than an intercept only model, and 0 fROIs are
 313 significantly better after FDR correction. Our results therefore suggest that GPT2-XL performs at
 314 near chance levels on the version of the Blank dataset used by [2, 10, 11].

315 6 Limitations and Conclusions

316 Our study has three main limitations. First, our method of examining how much neural variance
 317 an LLM predicts over simple features scales poorly when the number of features is large. Second,
 318 although we attempted to correct for cases where adding features decreases test set performance and
 319 employed banded regression, fitting regressions with large feature spaces on noisy neural data with
 320 low sample sizes can lead to poor estimations of the neural variance explained. Finally, we did not
 321 analyze datasets with large amounts of neural data per participant, for instance [36], in which the gap
 322 between the neural predictivity of simple and complex features might be much larger.

323 In summary, we find that on the Pereira dataset, shuffled splits are heavily impacted by temporal
 324 autocorrelation, untrained GPT2-XL brain score is explained by sentence length and position, and
 325 trained GPT2-XL brain score is largely explained by non-contextual features. We find that the
 326 majority of GPT2-XL brain score on the Fedorenko dataset is accounted for by word position, and
 327 on the Blank dataset GPT2-XL predicts neural activity at near chance levels. These results suggest
 328 that (i) brain scores on these datasets should be interpreted with caution; and (ii) more generally,
 329 analyses using brain scores should be accompanied by a systematic deconstruction of neural encoding
 330 performance, and an evaluation against simple and theoretically uninteresting features. Only after
 331 such deconstruction can we be somewhat confident that the neural predictivity of LLMs reflects core
 332 aspects of human linguistic processing.

333 References

- 334 [1] Kyle Mahowald, Anna A Ivanova, Idan A Blank, Nancy Kanwisher, Joshua B Tenenbaum, and
335 Evelina Fedorenko. Dissociating language and thought in large language models. *Trends Cogn.*
336 *Sci.*, March 2024.
- 337 [2] Martin Schrimpf, Idan Asher Blank, Greta Tuckute, Carina Kauf, Eghbal A Hosseini, Nancy
338 Kanwisher, Joshua B Tenenbaum, and Evelina Fedorenko. The neural architecture of language:
339 Integrative modeling converges on predictive processing. *Proc. Natl. Acad. Sci. U. S. A.*, 118
340 (45), November 2021.
- 341 [3] Mariya Toneva, Tom M Mitchell, and Leila Wehbe. Combining computational controls with
342 natural text reveals aspects of meaning composition. *Nat Comput Sci*, 2(11):745–757, November
343 2022.
- 344 [4] Charlotte Caucheteux and Jean-Rémi King. Brains and algorithms partially converge in natural
345 language processing. *Commun Biol*, 5(1):134, February 2022.
- 346 [5] Shailee Jain and Alexander Huth. Incorporating context into language encoding models for
347 fMRI. *Adv. Neural Inf. Process. Syst.*, 31, 2018.
- 348 [6] Mariya Toneva and Leila Wehbe. Interpreting and improving natural-language processing (in
349 machines) with natural language-processing (in the brain). *Adv. Neural Inf. Process. Syst.*, pages
350 14928–14938, May 2019.
- 351 [7] Alec Radford, Jeff Wu, R Child, D Luan, Dario Amodei, and I Sutskever. Language models are
352 unsupervised multitask learners. 2019.
- 353 [8] Francisco Pereira, Bin Lou, Brianna Pritchett, Samuel Ritter, Samuel J Gershman, Nancy
354 Kanwisher, Matthew Botvinick, and Evelina Fedorenko. Toward a universal decoder of linguistic
355 meaning from brain activation. *Nat. Commun.*, 9(1):963, March 2018.
- 356 [9] Alexandre Pasquiou, Yair Lakretz, John Hale, Bertrand Thirion, and Christophe Pallier. Neural
357 language models are not born equal to fit brain data, but training helps. July 2022.
- 358 [10] Eghbal A Hosseini, Martin Schrimpf, Yian Zhang, Samuel Bowman, Noga Zaslavsky, and
359 Evelina Fedorenko. Artificial neural network language models predict human brain responses
360 to language even after a developmentally realistic amount of training. *Neurobiol Lang (Camb)*,
361 5(1):43–63, April 2024.
- 362 [11] Khai Loong Aw, Syrielle Montariol, Badr AlKhamissi, Martin Schrimpf, and Antoine Bosselut.
363 Instruction-tuned LLMs with world knowledge are more aligned to the human brain, 2024.
364 URL <https://openreview.net/forum?id=DZ6B5u4vfe>.
- 365 [12] Charlotte Caucheteux, Alexandre Gramfort, and Jean-Rémi King. Evidence of a predictive
366 coding hierarchy in the human brain listening to speech. *Nat Hum Behav*, 7(3):430–441, March
367 2023.
- 368 [13] Ariel Goldstein, Eric Ham, Mariano Schain, Samuel Nastase, Zaid Zada, Avigail Dabush,
369 Bobbi Aubrey, Harshvardhan Gazula, Amir Feder, Werner K Doyle, Sasha Devore, Patricia
370 Dugan, Daniel Friedman, Roi Reichart, Michael Brenner, Avinatan Hassidim, Orrin Devinsky,
371 Adeen Flinker, Omer Levy, and Uri Hasson. The temporal structure of language processing in
372 the human brain corresponds to the layered hierarchy of deep language models, 2024. URL
373 <https://openreview.net/forum?id=950bXevgHx>.
- 374 [14] Refael Tikochinski, Ariel Goldstein, Yoav Meiri, Uri Hasson, and Roi Reichart. Incremental ac-
375 cumulation of linguistic context in artificial and biological neural networks. *bioRxiv*, 2024. doi:
376 10.1101/2024.01.15.575798. URL [https://www.biorxiv.org/content/early/2024/
377 01/17/2024.01.15.575798](https://www.biorxiv.org/content/early/2024/01/17/2024.01.15.575798).
- 378 [15] Jeffrey S Bowers, Gaurav Malhotra, Federico Adolfi, Marin Dujmović, Milton L Montero,
379 Valerio Biscione, Guillermo Puebla, John H Hummel, and Rachel F Heaton. On the importance
380 of severely testing deep learning models of cognition. *Cogn. Syst. Res.*, 82:101158, December
381 2023.

- 382 [16] Richard Antonello, Aditya R Vaidya, and Alexander G Huth. Scaling laws for language
383 encoding models in fMRI. *Adv. Neural Inf. Process. Syst.*, abs/2305.11863, May 2023.
- 384 [17] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for
385 word representation. In *Proceedings of the 2014 conference on empirical methods in natural*
386 *language processing (EMNLP)*, pages 1532–1543, 2014.
- 387 [18] Kilian Weinberger. On the importance of deconstruction in machine learning research.
388 ML-Retrospectives @ NeurIPS 2020, 2020. URL [https://slideslive.com/38938218/](https://slideslive.com/38938218/the-importance-of-deconstruction)
389 [the-importance-of-deconstruction](https://slideslive.com/38938218/the-importance-of-deconstruction).
- 390 [19] Evelina Fedorenko, Michael K Behr, and Nancy Kanwisher. Functional specificity for high-level
391 linguistic processing in the human brain. *Proc. Natl. Acad. Sci. U. S. A.*, 108(39):16428–16433,
392 September 2011.
- 393 [20] Idan Blank, Nancy Kanwisher, and Evelina Fedorenko. A functional dissociation between
394 language and multiple-demand systems revealed in patterns of BOLD signal fluctuations. *J.*
395 *Neurophysiol.*, 112(5):1105–1118, September 2014.
- 396 [21] Aniketh Janardhan Reddy and Leila Wehbe. Can fMRI reveal the representation of syntactic
397 structure in the brain? *Adv. Neural Inf. Process. Syst.*, 34:9843–9856, December 2021.
- 398 [22] Wendy A de Heer, Alexander G Huth, Thomas L Griffiths, Jack L Gallant, and Frédéric E
399 Theunissen. The hierarchical cortical organization of human speech processing. *J. Neurosci.*,
400 37(27):6539–6557, July 2017.
- 401 [23] Richard Futrell, Edward Gibson, Harry J Tily, Idan Blank, Anastasia Vishnevetsky, Steven
402 Piantadosi, and Evelina Fedorenko. The natural stories corpus. In Nicoletta Calzolari, Khalid
403 Choukri, Christopher Cieri, Thierry Declerck, Sara Goggi, Koiti Hasida, Hitoshi Isahara, Bente
404 Maegaard, Joseph Mariani, H el ene Mazo, Asuncion Moreno, Jan Odijk, Stelios Piperidis,
405 and Takenobu Tokunaga, editors, *Proceedings of the Eleventh International Conference on*
406 *Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan, May 2018. European
407 Language Resources Association (ELRA).
- 408 [24] Carina Kauf, Greta Tuckute, Roger Levy, Jacob Andreas, and Evelina Fedorenko. Lexical-
409 Semantic content, not syntactic structure, is the main contributor to ANN-Brain similarity of
410 fMRI responses in the language network. *Neurobiol Lang (Camb)*, 5(1):7–42, April 2024.
- 411 [25] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy,
412 Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. RoBERTa: A robustly optimized BERT
413 pretraining approach. July 2019.
- 414 [26] Alexander G Huth, Wendy A de Heer, Thomas L Griffiths, Frédéric E Theunissen, and Jack L
415 Gallant. Natural speech reveals the semantic maps that tile human cerebral cortex. *Nature*, 532
416 (7600):453–458, April 2016.
- 417 [27] Shailee Jain, Vy A Vo, Shivangi Mahto, Amanda LeBel, Javier Turek, and Alexander G Huth.
418 Interpretable multi-timescale models for predicting fMRI responses to continuous natural speech.
419 *Adv. Neural Inf. Process. Syst.*, 33, October 2020.
- 420 [28] Tom Dupr e la Tour, Michael Eickenberg, Anwar O Nunez-Elizalde, and Jack L Gallant. Feature-
421 space selection with banded ridge regression. *Neuroimage*, 264:119728, December 2022.
- 422 [29] Stijn Hawinkel, Willem Waegeman, and Steven Maere. Out-of-Sample r^2 : Estimation and
423 inference. *Am. Stat.*, pages 1–11.
- 424 [30] Subba Reddy Oota, Jashn Arora, Veeral Agarwal, Mounika Marreddy, Manish Gupta, and Bapi
425 Surampudi. Neural language taskonomy: Which NLP tasks are the most predictive of fMRI
426 brain activity? In Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz,
427 editors, *Proceedings of the 2022 Conference of the North American Chapter of the Association*
428 *for Computational Linguistics: Human Language Technologies*, pages 3220–3237, Seattle,
429 United States, July 2022. Association for Computational Linguistics.

- 430 [31] Yoav Benjamini and Yosef Hochberg. Controlling the false discovery rate: A practical and
 431 powerful approach to multiple testing. *J. R. Stat. Soc. Series B Stat. Methodol.*, 57(1):289–300,
 432 1995.
- 433 [32] Daniel Loureiro, Alípio Mário Jorge, and Jose Camacho-Collados. LMMS reloaded:
 434 Transformer-based sense embeddings for disambiguation and beyond. *Artif. Intell.*, 305:103661,
 435 April 2022.
- 436 [33] George A Miller. WordNet: a lexical database for english. *Commun. ACM*, 38(11):39–41,
 437 November 1995.
- 438 [34] Charlotte Caucheteux, Alexandre Gramfort, and Jean-Remi King. Disentangling syntax and
 439 semantics in the brain with deep networks. March 2021.
- 440 [35] Evelina Fedorenko, Terri L Scott, Peter Brunner, William G Coon, Brianna Pritchett, Gerwin
 441 Schalk, and Nancy Kanwisher. Neural correlate of the construction of sentence meaning. *Proc.*
 442 *Natl. Acad. Sci. U. S. A.*, 113(41):E6256–E6262, October 2016.
- 443 [36] Amanda LeBel, Lauren Wagner, Shailee Jain, Aneesh Adhikari-Desai, Bhavin Gupta, Allyson
 444 Morgenthal, Jerry Tang, Lixiang Xu, and Alexander G Huth. A natural language fMRI dataset
 445 for voxelwise encoding models. *Sci Data*, 10(1):555, August 2023.
- 446 [37] Zachary Mineroff, Idan Asher Blank, Kyle Mahowald, and Evelina Fedorenko. A robust
 447 dissociation among the language, multiple demand, and default mode networks: Evidence from
 448 inter-region correlations in effect size. *Neuropsychologia*, 119:501–511, October 2018.
- 449 [38] Jonathan D Power, Alexander L Cohen, Steven M Nelson, Gagan S Wig, Kelly Anne Barnes,
 450 Jessica A Church, Alecia C Vogel, Timothy O Laumann, Fran M Miezin, Bradley L Schlaggar,
 451 and Steven E Petersen. Functional network organization of the human brain. *Neuron*, 72(4):
 452 665–678, November 2011.
- 453 [39] Thomas Lumley, Paula Diehr, Scott Emerson, and Lu Chen. The importance of the normality
 454 assumption in large public health data sets. *Annu. Rev. Public Health*, 23:151–169, 2002.
- 455 [40] Simon Musall, Matthew T Kaufman, Ashley L Juavinett, Steven Gluf, and Anne K Churchland.
 456 Single-trial neural dynamics are dominated by richly varied movements. *Nat. Neurosci.*, 22(10):
 457 1677–1686, October 2019.
- 458 [41] Matthew Honnibal and Ines Montani. spaCy 2: Natural language understanding with Bloom
 459 embeddings, convolutional neural networks and incremental parsing. To appear, 2017.

460 A Appendix

461 A.1 Experimental data

462 **Pereira:** For both experiments, each sentence was visually presented for 4 s with 4 s between
 463 sentences and an additional 4 s between passages. A single fMRI scan was taken in the interval
 464 between each sentence. Because fMRI data is noisy, each experiment was repeated three times and
 465 fMRI data was averaged across the repetitions. A single fMRI scanning session consisted of 8 runs,
 466 where each run contained 12 passages in EXP1 and 9 passages in EXP2. Participants performed
 467 a total of 3 scanning sessions. The division of passages into runs and the order of the runs was
 468 randomized for each participant and scanning session.

469 **Fedorenko:** Participants read sentence on word at a time, and each word was visually displayed for
 470 450 or 700 ms. For each electrode, high gamma signal was extracted using gaussian filter banks at
 471 center frequencies ranging from 73 – 144 Hz, the envelope of the high gamma signal was computed
 472 through a hilbert-transform, and the envelope was z-scored within each electrode. For each participant,
 473 language-selective electrodes were selected where the z-scored envelope of the gamma activity was
 474 significantly higher during the sentences than a condition where participants read nonword lists.
 475 Z-scored high gamma activity from these language-selective electrodes were used in subsequent
 476 analyses.

477 **Blank:** Text was split into 2 s segments corresponding to each TR, with words that were on the
478 boundary being assigned to the later TR. Due to the delay in the hemodynamic response function
479 (HRF), neural activity was predicted using stimuli from 2 TRs (4 s) previous.

480 **Functional localization:** For Pereira and Blank, the language network was defined by the following
481 procedure [19]. First, voxels were identified in each participant which showed stronger responses
482 to sentences compared to lists of non-words (sentences > non-word lists contrast). These voxels
483 were then constrained by data-driven language activation maps formed by applying the same contrast
484 to many other participants. Finally, the top 10% of the voxels were selected which showed the
485 greatest sentences > non-word lists difference. For Pereira, we perform some analyses using four
486 other networks: multiple demand (MD), default mode network (DMN), auditory, and visual network.
487 The multiple demand (MD) and default mode network (DMN) networks were defined using the same
488 procedure, except that the contrast involved a spatial working memory task, where a hard > easy
489 condition contrast was used for MD and a fixation > hard contrast was used for DMN [37]. Auditory
490 and visual networks were defined using resting state connectivity [38].

491 **A.2 Banded ridge regression**

492 We used a random search method to optimize the banded regression hyperparameters [28]. Banded
493 regression has two hyperparameters, γ , which is a vector of shape number of feature spaces that
494 determines how much each feature space is scaled, and α , which is the L2 penalty applied across
495 feature spaces. Values for γ are drawn from a Dirichlet distribution and hence sum to 1. Down-scaling
496 a certain feature space relative to others is functionally equivalent to assigning a separate L2 penalty
497 for each feature space. This is because when a feature space is down-scaled, the L2 magnitude of
498 the weights must increase for it to have a meaningful contribution to the predictions, which equates
499 to increasing the L2 penalty for that feature space. The optimal γ and α combination was found
500 for each voxel/electrode/fROI by performing a random search over γ values, storing the α value
501 that performed best for that γ on validation data, and then selecting the best performing γ and α
502 combination.

503 Before starting the random search, we tried all combinations of γ values that removed feature spaces
504 (i.e. down-scaled at least one feature space to 0) to ensure the regression had an opportunity to
505 remove features which hurt performance. In theory, this should obviate the need for the procedure
506 implemented in 2.8. This is because the banded regression procedure can remove feature spaces
507 based on validation data, meaning if a model performs worse than a sub-model the banded procedure
508 has the opportunity to set the γ value corresponding to the additional feature spaces to 0. However,
509 because neural data is noisy and there is often little data per subject, performance on validation data is
510 not always indicative of performance on test-data. Therefore it is possible for the banded regression
511 procedure to include a feature space (since it helps on validation data), and for this feature space to
512 ultimately hurt test set performance, necessitating the correction procedure detailed in 2.8.

513 We ran banded ridge regression for a maximum of 1000 random search iterations with early stopping
514 if the mean R^2 did not improve by more than 10^{-4} after 50 iterations. We treated feature spaces
515 with many dimensions as one features because preliminary results showed this performed better.
516 Specifically, we always treated the following feature spaces as one feature space: static word
517 embeddings, sense-specific word embeddings, syntactic representations, and GPT2-XL and Roberta-
518 Large representations. All other features were treated as their own feature space.

519 We z-score all features across samples before training regressions, as is standard when using ridge
520 regression in neural encoding studies.

521 **A.3 Additional details on train, validation, and test folds**

522 **Pereira:** During each outer fold, a single passage from each of the 24 semantic categories from one
523 experiment was selected, and half of these passages were designated as the test set. This equated to 8
524 test folds for experiment 1 (4 passages per semantic category) and 6 test folds for experiment 2 (3
525 passages per semantic category). During each inner fold, we again selected one passage from each
526 semantic category, and half of these passages were designated as validation (leading to 7 inner folds
527 for experiment 1, and 5 inner folds for experiment 2).

528 **Fedorenko:** For each outer fold, we selected 4 sentences as the test fold, resulting in 13 outer folds.
 529 For each inner fold, we once again select 4 sentences as the validation set, resulting in 12 inner folds
 530 per outer fold.

531 **Blank:** For each outer fold, we selected a single story as the test fold, resulting in 8 outer folds. For
 532 each inner fold, each of the remaining stories served in turn as the validation set, resulting in 7 inner
 533 folds.

534 A.4 Justification of statistical tests

535 We performed a *t*-test between squared error values from two models to determine if one model
 536 performs better than another. While squared error values are not always normally distributed, our
 537 sample sizes were large (the minimum sample size was 243) and so we still opted to use a *t*-test over
 538 a non-parametric alternative [39]. One issue with a *t*-test is that relies on the assumption that samples
 539 are not correlated, which is not true for time-series data. However, we note that correlated samples
 540 leads one to underestimate the standard error of the mean and exaggerate differences between two
 541 models. Since we only perform one-sided *t*-tests to examine whether adding GPT2-XL representations
 542 improves performance, the net impact of this on our results is to overestimate how much GPT2-XL
 543 contributes over simple features.

544 A.5 Across layer R^2 values in the Pereira dataset

545 Across layer performances in the Pereira dataset for GPT2-XLU and GPT2-XL when using the sum
 546 pooling method (Fig. 5a,b) and the last token method (Fig. 5c,d). Performance in language network
 547 is higher across the board than performance in DMN, MD, and visual networks. We do not show
 548 auditory network results because participants read passages in Pereira and hence auditory brain scores
 549 are near 0. Furthermore, performance is lower with the last token method in every case except in
 EXP1 trained results where the last token method performs slightly better.

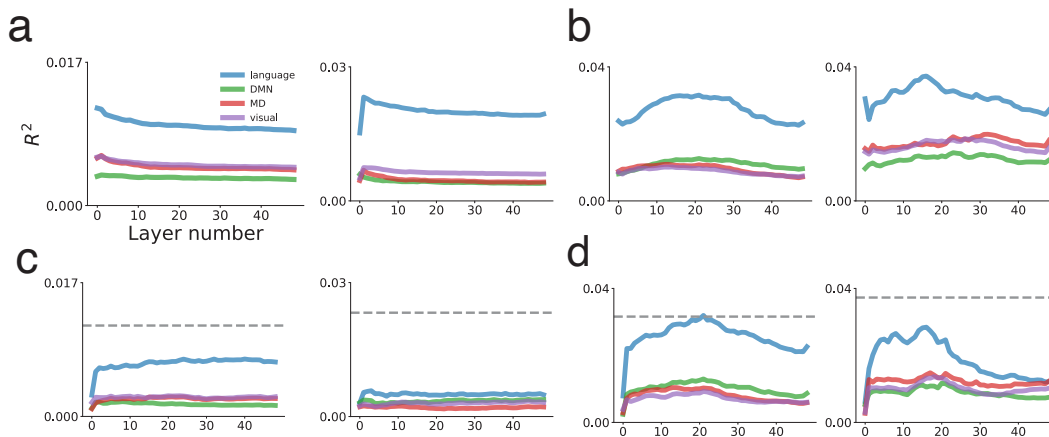


Figure 5: **a)** Across layer performances in Pereira dataset for GPT2-XLU for each functional network when using the sum-pooling method. EXP1 is on the left, and EXP2 is on the right. **b)** Same as **a)** but for GPT2-XL, also using the sum-pooling method. **c)** Same as **a)** but when using the last token method. Dotted grey line shows performance of best layer of GPT2-XLU in language network when sum pooling. **d)** Same as **b)** but when using the last token method. Dotted grey line shows performance of best layer of GPT2-XL in language network when sum pooling.

550

551 A.6 RoBERTa-Large shows similar results as GPT2-XL

552 To examine whether our results depending on the choice of LLM, we replicated all of our Pereira
 553 trained analyses with RoBERTa-Large (ROB). The overall trend in results was the same as
 554 with GPT2-XL (Fig. 6). Namely, SP+SL+WORD performed 76.8% as well as the full model
 555 (SP+SL+WORD+SENSE+SYNT+ROB) and 80.0% as well as ROB alone in EXP1, and in EXP2 it

556 performed 88.0% as well as the full model and better than ROB. Furthermore, SENSE and SYNT
 557 bridge the gap to the full model by a small amount. In sum, our main conclusion that a large amount
 558 of trained LLM brain score in the Pereira dataset is accounted for by non-contextual features also
 559 applies to RoBERTa-Large.

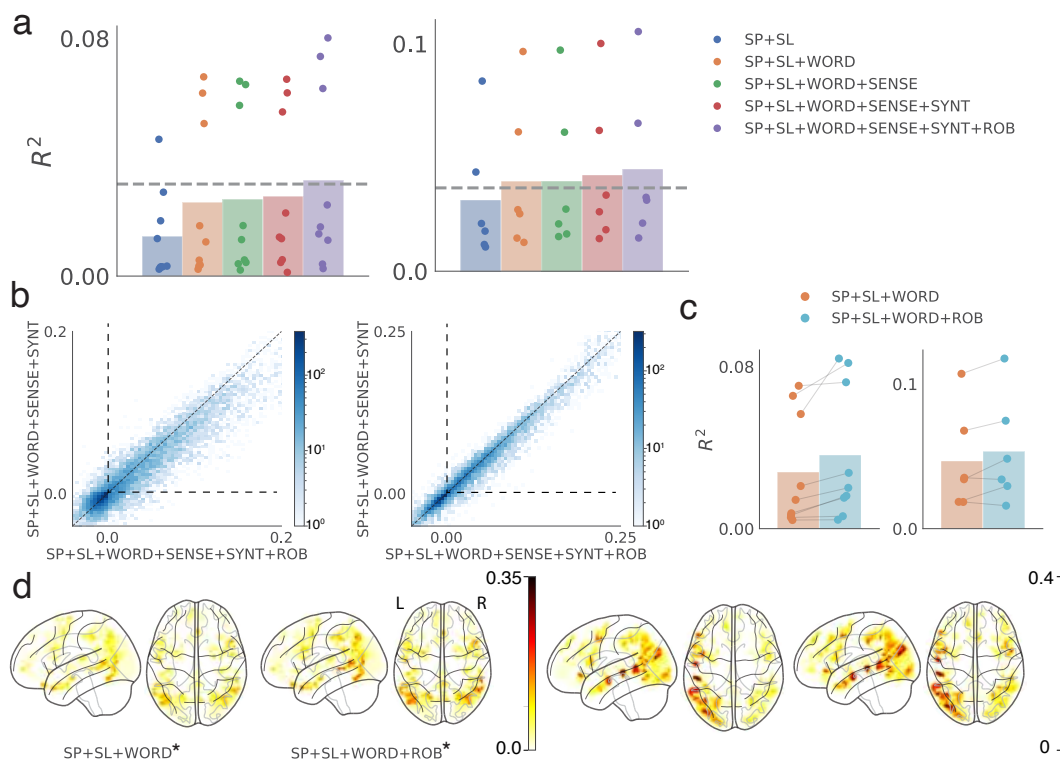


Figure 6: All panels are the same as Figure 3, except GPT2-XL is replaced with RoBERTa-Large (ROB).

560 A.7 Orthogonal autocorrelated sequences model (OASM) hyperparameters

561 The width of the Gaussian filter used for within-block smoothing was $\sigma = 2.2$ in Pereira, $\sigma = 1.8$ in
 562 Fedorenko, and $\sigma = 1.5$ in Blank. Gaussian widths were determined by sweeping σ across 50 evenly
 563 spaced values between 0.1 and 5.0 and choosing the best-performing σ for each dataset.

564 A.8 Shuffled train test splits confound task-relevant and task-irrelevant neural activity

565 OASM is a model which clearly lacks any linguistic representations that would allow it generalize to
 566 fully held-out passages. However, this is not to say that OASM is not correlated with linguistic
 567 features. For instance, sentences in a given passage are more semantically related with each other
 568 than with sentences in other passages. Nonetheless, using shuffled train-test splits almost certainly
 569 exaggerates the variance explained by a model which, on the basis of semantic similarity, arrives
 570 at a similar representational structure as OASM. This is because task-irrelevant neural responses
 571 make up a large fraction of neural activity [40], and shuffled train-test splits allow a model with
 572 OASM-like representational structure to predict not just the task-relevant neural responses driven
 573 by the participant reading the passage, but also any task-irrelevant neural activity that was present
 574 throughout the reading of the passage. Hence, we strongly urge researchers to avoid shuffled train
 575 test splits when evaluating the neural predictivity of language models, and we surmise that previous
 576 studies using shuffled train-test splits to compare neural predictivity between models might have
 577 come to erroneous conclusions.

578 **A.9 Linear decodability of sentence length**

579 Here, we show that the MLP block adds a linearly decodable component with non-zero mean to the
580 residual stream in the GPT2 architecture.

581 **Proof :**

582 We denote the i 'th input to the MLP block in the first layer of GPT2-XL as x_i . The output of the
583 MLP block is defined as follows:

$$MLP(x_i) = x_i + W_d(GELU(W_u(LayerNorm(x_i))))$$

584 We assume that the elements of x_i are normally distributed. For a given x_i , it then follows that the
585 distribution of elements in $LayerNorm(x_i)$ is normal with $\mu = 0$ and $\sigma = 1$ (assuming the standard
586 $LayerNorm$ initialization).

587 Because W_u is initialized from a zero-mean normal distribution, $W_u(LayerNorm(x_i))$ also has
588 zero-mean.

589 Note that $GELU$ is a function for which $\mathbb{E}[Y] > 0$ for Y normally distributed with mean 0. Hence, the
590 mean value across elements following the $GELU$ is non-zero. Let us denote this mean value across
591 all elements of $GELU(W_u(LayerNorm(x)))$ and across all tokens x as m . Then, for an MLP
592 with up-projected dimension d_u , we can take the dot product of $GELU(W_u(LayerNorm(x_i)))$
593 and $\frac{1}{d_u m} \times \hat{k}$, where \hat{k} is a d_u -dimensional vector of 1s. The resulting value will have mean 1.

However, we cannot decode this value directly from the MLP in practice; first, this vector is down-
projected back to the residual stream by W_d . Nonetheless, we can still closely approximate it,
assuming it is approximately orthogonal to x_i , by using the pseudo-inverse of W_d . More specifically,
we can extract a scalar with mean 1 as follows:

$$\sqrt{\frac{d_u}{d_d}} \times \frac{1}{d_u m} \times \hat{k} W_d^\dagger MLP(x_i)$$

594 where d_d is the down-projected dimension. Because this extracted scalar value is distributed with
595 mean 1 across token representations x_i , assuming independence of token representations within a
596 sentence, the sum of the extracted scalar value across the tokens of a sentence is distributed with
597 mean equaling the number of tokens in the sentence.

598 **A.10 LMMS**

599 LMMS generates a sense embedding for each word by averaging across contextual embeddings (in
600 our case from RoBERTa-Large) of that sense derived from a sense-annotated corpus. For words in
601 WordNet where labeled senses don't exist, LMMS sets their sense embeddings equal to the average
602 of sense embeddings with the same sense (or same hypernym/lexname if that approach fails). Finally,
603 the sense embeddings are averaged together with the gloss embeddings for that sense of the word
604 generated using the same LLM. For additional details refer to Loureiro et al. [32].

605 **A.11 Contextual syntactic representations**

606 Syntactic embeddings are derived by substituting content words (nouns, verbs, adjectives, and
607 adverbs) in the original sentences with words from the Generics KB corpus, matching their part-of-
608 speech and dependency tag via the SpaCy transformer-based tagger [41]. For each sentence in the
609 Pereira dataset, we generate 170 new sentences, ensuring the subtree token indices from each token
610 match those of the original sentence. The top 100 sentences, selected based on summed surprisal
611 with GPT2-XL, are retained. Each sentence's syntactic embedding is then computed by summing
612 token representations within each sentence and then averaging across the 100 sentences.

613 **A.12 Word position feature in Fedorenko dataset**

614 The primary finding in the paper which first collected the Fedorenko dataset [35] was a ramping of
615 neural activity across the words of sentences, where each sentence was 8 words long. Hence, we
616 concatenate a linearly ramping 1-dimensional positional signal to an 8-dimensional 1-hot positional

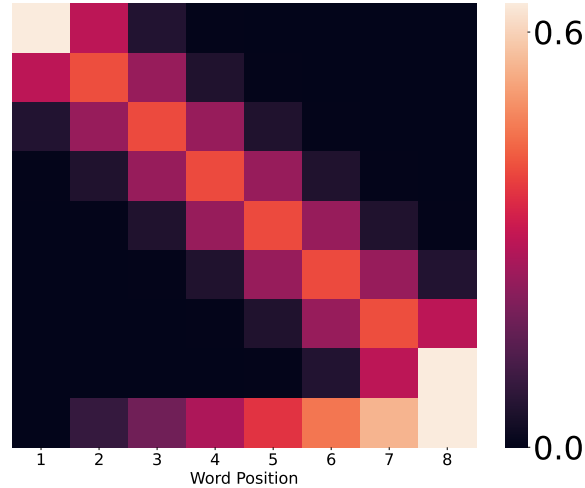


Figure 7: Word Position feature for a single sentence in the Fedorenko dataset.

617 signal. Because we expect positional signals to be more similar between adjacent words than more
 618 distant words, we apply a Gaussian filter ($\sigma = 1$) to the 8-dimensional positional signal. The resulting
 619 feature space, which we refer to as "word position" in the main text, is shown for a single sentence in
 620 the above figure.

621 **A.13 OASM and GPT2 Model Comparison on Blank Dataset**

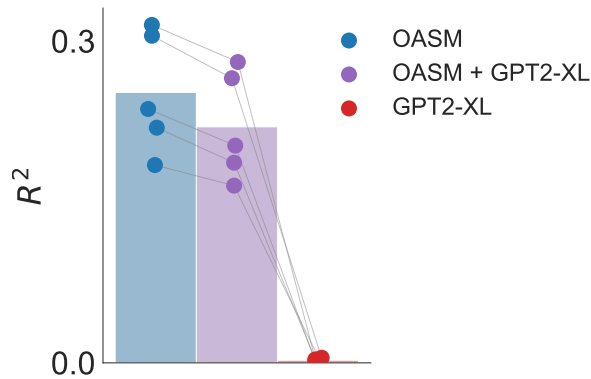


Figure 8: OASM far outperforms GPT2-XL on the Blank dataset, and GPT2-XL does not appear to explain any variance beyond that explained by OASM.

622 We find that OASM achieves 103.6 times higher neural predictivity than GPT2-XL on the Blank
 623 dataset when using shuffled train-test splits. There could be several reasons for this. First, it might
 624 be that the method for pooling representations from GPT2-XL used here 2.3 and in [2, 10, 11]
 625 did not yield useful enough representations for GPT2-XL to map effectively to the brain data. An
 626 additional likely culprit is that, of the three datasets we study here, Blank has the greatest potential for
 627 autocorrelation in temporally adjacent samples. This is because, while the Pereira dataset typically
 628 has a TR every 8 seconds, the Blank dataset has a TR every 2 seconds. We note that our results here
 629 are not completely surprising; given that [2, 10] observed untrained GPT2 models perform far better
 630 than trained models on this dataset, it did not seem likely that GPT2-XL would map onto neural
 631 representations of linguistic features here.

632 **A.14 Computational Resources**

633 All analyses were done between 2 machines: One with 2 RTX 3090 GPUs, and another with 1
634 RTX 4090 GPU. The most computationally demanding parts of our analyses were fitting the banded
635 ridge regressions used to generate Figure 3, collecting untrained model results across 10 seeds, and
636 generating syntactic representations, which each took around 3 hours to complete.

637 **A.15 Dataset Licenses**

638 The Blank dataset was originally released as part of the Natural Stories Corpus, which is provided
639 under the CC BY-NC-SA license [23]. The Pereira dataset is released under the Creative Commons
640 License [8]. The version of the Fedorenko dataset used here is provided under the MIT license. All
641 datasets used are the same versions as in [2] and can be downloaded using the neural-nlp repository:
642 <https://github.com/mschrimpf/neural-nlp/tree/master>. All datasets were collected with
643 IRB approval at their respective institutions.

644 **NeurIPS Paper Checklist**

645 **1. Claims**

646 Question: Do the main claims made in the abstract and introduction accurately reflect the
647 paper’s contributions and scope?

648 Answer: [\[Yes\]](#)

649 Justification: We support each of the three claims made in the abstract regarding shuffled
650 train-test splits, untrained LLM brain scores, and trained LLM brain scores in the Results
651 section. These results support the claim that it is important to deconstruct the mapping
652 between LLMs and the brain.

653 Guidelines:

- 654 • The answer NA means that the abstract and introduction do not include the claims
655 made in the paper.
- 656 • The abstract and/or introduction should clearly state the claims made, including the
657 contributions made in the paper and important assumptions and limitations. A No or
658 NA answer to this question will not be perceived well by the reviewers.
- 659 • The claims made should match theoretical and experimental results, and reflect how
660 much the results can be expected to generalize to other settings.
- 661 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
662 are not attained by the paper.

663 **2. Limitations**

664 Question: [\[Yes\]](#)

665 Justification: We discuss the three main limitations in the paper in the section titled "Limita-
666 tions and Conclusions", and additionally include limitations throughout the Appendix (e.g.
667 Justification of statistical tests).

668 Guidelines:

- 669 • The answer NA means that the paper has no limitation while the answer No means that
670 the paper has limitations, but those are not discussed in the paper.
- 671 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 672 • The paper should point out any strong assumptions and how robust the results are to
673 violations of these assumptions (e.g., independence assumptions, noiseless settings,
674 model well-specification, asymptotic approximations only holding locally). The authors
675 should reflect on how these assumptions might be violated in practice and what the
676 implications would be.
- 677 • The authors should reflect on the scope of the claims made, e.g., if the approach was
678 only tested on a few datasets or with a few runs. In general, empirical results often
679 depend on implicit assumptions, which should be articulated.
- 680 • The authors should reflect on the factors that influence the performance of the approach.
681 For example, a facial recognition algorithm may perform poorly when image resolution
682 is low or images are taken in low lighting. Or a speech-to-text system might not be
683 used reliably to provide closed captions for online lectures because it fails to handle
684 technical jargon.
- 685 • The authors should discuss the computational efficiency of the proposed algorithms
686 and how they scale with dataset size.
- 687 • If applicable, the authors should discuss possible limitations of their approach to
688 address problems of privacy and fairness.
- 689 • While the authors might fear that complete honesty about limitations might be used by
690 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
691 limitations that aren’t acknowledged in the paper. The authors should use their best
692 judgment and recognize that individual actions in favor of transparency play an impor-
693 tant role in developing norms that preserve the integrity of the community. Reviewers
694 will be specifically instructed to not penalize honesty concerning limitations.

695 **3. Theory Assumptions and Proofs**

696 Question: For each theoretical result, does the paper provide the full set of assumptions and
697 a complete (and correct) proof?

698 Answer: [Yes]

699 Justification: Our only theoretical result is that the MLP layer introduces a non-zero mean
700 component in the residual stream. We provide both a rough sketch in the main paper as well
701 as a formal proof.

702 Guidelines:

- 703 • The answer NA means that the paper does not include theoretical results.
- 704 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
705 referenced.
- 706 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 707 • The proofs can either appear in the main paper or the supplemental material, but if
708 they appear in the supplemental material, the authors are encouraged to provide a short
709 proof sketch to provide intuition.
- 710 • Inversely, any informal proof provided in the core of the paper should be complemented
711 by formal proofs provided in appendix or supplemental material.
- 712 • Theorems and Lemmas that the proof relies upon should be properly referenced.

713 4. Experimental Result Reproducibility

714 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
715 perimental results of the paper to the extent that it affects the main claims and/or conclusions
716 of the paper (regardless of whether the code and data are provided or not)?

717 Answer: [Yes]

718 Justification: We include all details regarding the following: banded regression proce-
719 dure, construction of feature spaces, train, validation, and test splits, and selection of
720 voxels/electrodes/fROIs in neural data. These are all the elements needed to reproduce our
721 results, with the exception of slight variability due to stochasticity in untrained LLM seeds
722 and the random search process in banded regression.

723 Guidelines:

- 724 • The answer NA means that the paper does not include experiments.
- 725 • If the paper includes experiments, a No answer to this question will not be perceived
726 well by the reviewers: Making the paper reproducible is important, regardless of
727 whether the code and data are provided or not.
- 728 • If the contribution is a dataset and/or model, the authors should describe the steps taken
729 to make their results reproducible or verifiable.
- 730 • Depending on the contribution, reproducibility can be accomplished in various ways.
731 For example, if the contribution is a novel architecture, describing the architecture fully
732 might suffice, or if the contribution is a specific model and empirical evaluation, it may
733 be necessary to either make it possible for others to replicate the model with the same
734 dataset, or provide access to the model. In general, releasing code and data is often
735 one good way to accomplish this, but reproducibility can also be provided via detailed
736 instructions for how to replicate the results, access to a hosted model (e.g., in the case
737 of a large language model), releasing of a model checkpoint, or other means that are
738 appropriate to the research performed.
- 739 • While NeurIPS does not require releasing code, the conference does require all submis-
740 sions to provide some reasonable avenue for reproducibility, which may depend on the
741 nature of the contribution. For example
 - 742 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
743 to reproduce that algorithm.
 - 744 (b) If the contribution is primarily a new model architecture, the paper should describe
745 the architecture clearly and fully.
 - 746 (c) If the contribution is a new model (e.g., a large language model), then there should
747 either be a way to access this model for reproducing the results or a way to reproduce
748 the model (e.g., with an open-source dataset or instructions for how to construct
749 the dataset).

750 (d) We recognize that reproducibility may be tricky in some cases, in which case
751 authors are welcome to describe the particular way they provide for reproducibility.
752 In the case of closed-source models, it may be that access to the model is limited in
753 some way (e.g., to registered users), but it should be possible for other researchers
754 to have some path to reproducing or verifying the results.

755 5. Open access to data and code

756 Question: Does the paper provide open access to the data and code, with sufficient instruc-
757 tions to faithfully reproduce the main experimental results, as described in supplemental
758 material?

759 Answer: [Yes]

760 Justification: We will release all our code on Github, and all neural datasets are openly
761 available for use. We also provide anonymized code.

762 Guidelines:

- 763 • The answer NA means that paper does not include experiments requiring code.
- 764 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
765 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 766 • While we encourage the release of code and data, we understand that this might not be
767 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
768 including code, unless this is central to the contribution (e.g., for a new open-source
769 benchmark).
- 770 • The instructions should contain the exact command and environment needed to run to
771 reproduce the results. See the NeurIPS code and data submission guidelines ([https://
772 nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 773 • The authors should provide instructions on data access and preparation, including how
774 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 775 • The authors should provide scripts to reproduce all experimental results for the new
776 proposed method and baselines. If only a subset of experiments are reproducible, they
777 should state which ones are omitted from the script and why.
- 778 • At submission time, to preserve anonymity, the authors should release anonymized
779 versions (if applicable).
- 780 • Providing as much information as possible in supplemental material (appended to the
781 paper) is recommended, but including URLs to data and code is permitted.

782 6. Experimental Setting/Details

783 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
784 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
785 results?

786 Answer: [Yes]

787 Justification: We dedicate sections towards explaining the data splits in the main paper, and
788 the necessary details to run the banded ridge regression in the main paper and Appendix.

789 Guidelines:

- 790 • The answer NA means that the paper does not include experiments.
- 791 • The experimental setting should be presented in the core of the paper to a level of detail
792 that is necessary to appreciate the results and make sense of them.
- 793 • The full details can be provided either with the code, in appendix, or as supplemental
794 material.

795 7. Experiment Statistical Significance

796 Question: Does the paper report error bars suitably and correctly defined or other appropriate
797 information about the statistical significance of the experiments?

798 Answer: [Yes]

799 Justification: We perform a paired t-test and justify its use in the Appendix. For all plots
800 which show the average across participants we show individual dots for each participant,
801 and for this reason we do not include standard deviation values for the values in Table 1.

802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We provide a section in the appendix describing the GPUs and CPUs used for our analyses, and we describe how long each experiment took to run.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

Answer: [Yes]

Justification: We did not conduct any direct interactions with human participants, none of the data-related concerns apply for us, and we do not see any direct societal impacts from our work. We make our methods clear to the best of our ability and provide anonymized code.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

854 Question: Does the paper discuss both potential positive societal impacts and negative
855 societal impacts of the work performed?

856 Answer: [NA]

857 Justification: We do not develop any novel technology that can be used for good or bad, but
858 rather show that some high-profile previous results have been over-interpreted. While our
859 results are relevant for the cognitive neuroscience community, we do not see a direct path to
860 any larger societal impacts.

861 Guidelines:

- 862 • The answer NA means that there is no societal impact of the work performed.
- 863 • If the authors answer NA or No, they should explain why their work has no societal
864 impact or why the paper does not address societal impact.
- 865 • Examples of negative societal impacts include potential malicious or unintended uses
866 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
867 (e.g., deployment of technologies that could make decisions that unfairly impact specific
868 groups), privacy considerations, and security considerations.
- 869 • The conference expects that many papers will be foundational research and not tied
870 to particular applications, let alone deployments. However, if there is a direct path to
871 any negative applications, the authors should point it out. For example, it is legitimate
872 to point out that an improvement in the quality of generative models could be used to
873 generate deepfakes for disinformation. On the other hand, it is not needed to point out
874 that a generic algorithm for optimizing neural networks could enable people to train
875 models that generate Deepfakes faster.
- 876 • The authors should consider possible harms that could arise when the technology is
877 being used as intended and functioning correctly, harms that could arise when the
878 technology is being used as intended but gives incorrect results, and harms following
879 from (intentional or unintentional) misuse of the technology.
- 880 • If there are negative societal impacts, the authors could also discuss possible mitigation
881 strategies (e.g., gated release of models, providing defenses in addition to attacks,
882 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
883 feedback over time, improving the efficiency and accessibility of ML).

884 11. Safeguards

885 Question: Does the paper describe safeguards that have been put in place for responsible
886 release of data or models that have a high risk for misuse (e.g., pretrained language models,
887 image generators, or scraped datasets)?

888 Answer: [NA]

889 Justification: We release no new models or datasets, and do not see any potential for our
890 results being misused in unsafe ways.

891 Guidelines:

- 892 • The answer NA means that the paper poses no such risks.
- 893 • Released models that have a high risk for misuse or dual-use should be released with
894 necessary safeguards to allow for controlled use of the model, for example by requiring
895 that users adhere to usage guidelines or restrictions to access the model or implementing
896 safety filters.
- 897 • Datasets that have been scraped from the Internet could pose safety risks. The authors
898 should describe how they avoided releasing unsafe images.
- 899 • We recognize that providing effective safeguards is challenging, and many papers do
900 not require this, but we encourage authors to take this into account and make a best
901 faith effort.

902 12. Licenses for existing assets

903 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
904 the paper, properly credited and are the license and terms of use explicitly mentioned and
905 properly respected?

906 Answer: [Yes]

907 Justification: We cite the papers in which all datasets used were first published. We provide
908 the licenses for the Blank and Pereira datasets in the supplement (we could not find a license
909 for the Fedorenko dataset). We also specify the version of the datasets used and provide a
910 link.

911 Guidelines:

- 912 • The answer NA means that the paper does not use existing assets.
- 913 • The authors should cite the original paper that produced the code package or dataset.
- 914 • The authors should state which version of the asset is used and, if possible, include a
915 URL.
- 916 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 917 • For scraped data from a particular source (e.g., website), the copyright and terms of
918 service of that source should be provided.
- 919 • If assets are released, the license, copyright information, and terms of use in the
920 package should be provided. For popular datasets, `paperswithcode.com/datasets`
921 has curated licenses for some datasets. Their licensing guide can help determine the
922 license of a dataset.
- 923 • For existing datasets that are re-packaged, both the original license and the license of
924 the derived asset (if it has changed) should be provided.
- 925 • If this information is not available online, the authors are encouraged to reach out to
926 the asset’s creators.

927 13. New Assets

928 Question: Are new assets introduced in the paper well documented and is the documentation
929 provided alongside the assets?

930 Answer: [NA]

931 Justification: We do not release any new assets with this paper.

932 Guidelines:

- 933 • The answer NA means that the paper does not release new assets.
- 934 • Researchers should communicate the details of the dataset/code/model as part of their
935 submissions via structured templates. This includes details about training, license,
936 limitations, etc.
- 937 • The paper should discuss whether and how consent was obtained from people whose
938 asset is used.
- 939 • At submission time, remember to anonymize your assets (if applicable). You can either
940 create an anonymized URL or include an anonymized zip file.

941 14. Crowdsourcing and Research with Human Subjects

942 Question: For crowdsourcing experiments and research with human subjects, does the paper
943 include the full text of instructions given to participants and screenshots, if applicable, as
944 well as details about compensation (if any)?

945 Answer: [Yes]

946 Justification: We use open source datasets where neural data is obtained from consenting
947 human adults. Information regarding research protocols is detailed in the references for
948 these datasets.

949 Guidelines:

- 950 • The answer NA means that the paper does not involve crowdsourcing nor research with
951 human subjects.
- 952 • Including this information in the supplemental material is fine, but if the main contribu-
953 tion of the paper involves human subjects, then as much detail as possible should be
954 included in the main paper.
- 955 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
956 or other labor should be paid at least the minimum wage in the country of the data
957 collector.

958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [\[Yes\]](#)

Justification: All datasets used here were collected with IRB approval at their respective institutions, and this is stated in the appendix. We do not collect any data of our own from human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.