

Research Articles: Behavioral/Cognitive

An integrated neural decoder of linguistic and experiential meaning

<https://doi.org/10.1523/JNEUROSCI.2575-18.2019>

Cite as: J. Neurosci 2019; 10.1523/JNEUROSCI.2575-18.2019

Received: 4 October 2018

Revised: 26 August 2019

Accepted: 31 August 2019

This Early Release article has been peer-reviewed and accepted, but has not been through the composition and copyediting processes. The final version may differ slightly in style or formatting and will contain links to any extended data.

Alerts: Sign up at www.jneurosci.org/alerts to receive customized email alerts when the fully formatted version of this article is published.

1

2 An integrated neural decoder of linguistic and experiential meaning

3

4 Abbreviated title: A neural decoder of linguistic/experiential meaning

5 Andrew James Anderson¹, Jeffrey R. Binder², Leonardo Fernandino², Colin J. Humphries²,
6 Lisa L. Conant², Rajeev D. S. Raizada³, Feng Lin^{4,5}, Edmund C. Lalor^{1,6,7}

7 ¹Department of Neuroscience, University of Rochester, Rochester, NY 14642, USA

8 ²Department of Neurology, Medical College of Wisconsin, Milwaukee, WI, 53226, USA

9 ³Department of Brain and Cognitive Sciences, University of Rochester, Rochester, NY, 14627, USA

10 ⁴School of Nursing, University of Rochester, Rochester, NY 14642, USA

11 ⁵Psychiatry, University of Rochester, Rochester, NY 14642, USA

12 ⁶Department of Biomedical Engineering, University of Rochester, Rochester, NY, 14627, USA

13 ⁷School of Engineering, Trinity Centre for Bioengineering, and Trinity College Institute of Neuroscience, Trinity
14 College Dublin, Dublin, Ireland

15 Corresponding author with complete address, including an email address and postal code: Andrew James
16 Anderson, Department of Neuroscience, University of Rochester, Rochester, NY 14627,
17 aander41@ur.rochester.edu

18

19 Number of pages: 48

20 Number of figures: 11

21 Number of tables: 0

22 Number of multimedia and 3D models: 0

23 Number of words for Abstract: 237, Significance Statement 120 words, Introduction: 675, and
24 Discussion: 1502

25

26

27

28 **Conflict of Interest:** None

29 **Acknowledgements**

30 This work was supported in part by a Schmitt Program on Integrative Neuroscience (University
31 of Rochester Medical Center) award and the Intelligence Advanced Research Projects Activity
32 (IARPA) via the Air Force Research Laboratory under grant FA8650-14-C-7357 and NSF

33 CAREER award 1652127. We thank 3 reviewers for their insightful comments and time and
34 Douwe Kiela and Katrin Erk for models used in supporting analyses.

35

36 **Abstract 237/250 words**

37 The brain is thought to combine linguistic knowledge of words and non-linguistic knowledge of their
38 referents to encode sentence meaning. However, functional neuroimaging studies aiming at decoding
39 language meaning from neural activity have mostly relied on distributional models of word semantics,
40 which are based on patterns of word co-occurrence in text corpora. Here, we present initial evidence
41 that modeling non-linguistic “experiential” knowledge contributes to decoding neural representations of
42 sentence meaning. We model attributes of peoples’ sensory, motor, social, emotional and cognitive
43 experiences with words using behavioral ratings. We demonstrate that functional Magnetic Resonance
44 Imaging (fMRI) activation elicited in sentence reading is more accurately decoded when this
45 experiential attribute model is integrated with a text-based model than when either model is applied in
46 isolation (participants were 5 males and 9 females). Our decoding approach exploits a representation-
47 similarity-based framework which benefits from being parameter free, whilst performing at accuracy
48 levels comparable to those from parameter fitting approaches such as ridge-regression. We find that
49 the text-based model contributes particularly to the decoding of sentences containing linguistically
50 oriented “abstract” words and reveal tentative evidence that the experiential model improves decoding
51 of more concrete sentences. Finally, we introduce a cross-participant decoding method to estimate an
52 upper-bound on model-based decoding accuracy. We demonstrate that a substantial fraction of neural
53 signal remains unexplained, and leverage this gap to pinpoint characteristics of weakly decoded
54 sentences and hence identify model weaknesses to guide future model development.

55

56 **Significance statement 120 words**

57 Language gives humans the unique ability to communicate about historical events, theoretical concepts
58 and fiction. Although words are learnt through language and defined by their relations to other words in
59 dictionaries, our understanding of word meaning presumably draws heavily on our non-linguistic
60 sensory, motor, interoceptive and emotional experiences with words and their referents. Behavioral
61 experiments lend support to the intuition that word meaning integrates aspects of linguistic and non-

62 linguistic “experiential” knowledge. However behavioral measures do not provide a window on how
63 meaning is represented in the brain and tend to necessitate artificial experimental paradigms. We
64 present a model-based approach that reveals early evidence that experiential and linguistically
65 acquired knowledge can be detected in brain activity elicited in reading natural sentences.

66

67 **Introduction 675/650 words**

68 Humans’ knowledge of historical events, theoretical concepts and fiction is acquired through language.
69 Whilst a considerable body of linguistic information is stored in text and media repositories, the
70 meaning of language is a biological construct, instantiated in our brains during comprehension. Recent
71 advances in neuroimaging technology, big data and computational modeling, have led to an ability to
72 decode brain activity associated with linguistic meaning. The dominant decoding approach achieves
73 this using only *text-based* information. Word meaning is modeled as a vector of values reflecting how
74 often each word co-occurred with other words across a huge body of text (Landauer and Dumais 1997,
75 Burgess and Lund 1997, Mikolov et al. 2013, Pennington et al. 2014). Despite never having
76 experienced walking or eating, the model “learns” that walk and eat mean different things because they
77 appear in different textual contexts, but that walking relates to legs/shoes and that eating relates to
78 hunger/food because these words frequently co-occur (and end up with similar semantic vectors). This
79 approach supports the construction of conceptual knowledge hierarchies, e.g., a dragonfly is an insect
80 is an animal (Fu et al. 2014), and enables some level of analogical reasoning, e.g., Einstein is to
81 scientist as Picasso is to painter (Mikolov et al. 2013). By registering model and brain activity for
82 corresponding words, and mapping between the two, brain activity for new words, sentences and
83 narratives can be predicted and decoded (Mitchell et al. 2008; Pereira et al. 2013, Wehbe et al. 2014,
84 Huth et al. 2016, de Heer et al. 2017, Pereira et al. 2018).

85

86 Behavioral experiments suggest that in addition to linguistic experience, word meaning is shaped by
87 non-linguistic perceptual, motor and interoceptive experiences (Paivio, 1971, Barsalou 1999, Barsalou
88 et al. 2008, Vigliocco et al. 2009, Andrews et al. 2009, Riordan and Jones, 2010, Louwerse and

89 Jeuniaux 2010, Kousta et al. 2011, Dove 2014, Vigliocco et al. 2014, Andrews et al. 2014, Zwaan 2014,
90 Louwerse, 2018). Indirect evidence that sentence comprehension induces perceptual/motor simulations
91 related to sentence content has been amassed in reaction time studies (Stanfield and Zwaan 2001,
92 Zwaan et al. 2002, Glenberg and Kaschak 2002, Kaschak et al. 2005, Kaschak et al. 2006, Connell
93 2007, Glenberg et al. 2008, Winter and Bergen 2012, Zwaan and Pecher 2012, Speed and Vigliocco
94 2014). However, there is little direct *neural* evidence concerning how linguistic and non-linguistic
95 “experiential” sources of knowledge are encoded in brain activity (though see Anderson et al. 2015 and
96 Wang et al. 2018). The ability to estimate linguistic and non-linguistic contributions to neural
97 representations of meaning is necessary to fully characterize human language and related clinical
98 conditions (Patterson et al. 2007; Fernandino et al. 2013; Lambon Ralph et al. 2017, Anderson and Lin,
99 2019, Bruffaerts et al. 2019). Critically, neural measures can also help provide a window on natural
100 language comprehension, removing the necessity for artificial behavioral response tasks (which may
101 perturb linguistic systems), and bespoke stimulus materials (that may not reflect natural language). See
102 Hasson et al. (2018) and Hamilton and Huth (2018) for related discussion.

103
104 We reveal initial evidence that non-linguistic experiential knowledge can be detected in brain activity
105 elicited in sentence reading by combining an *experiential attribute* model (Binder et al. 2016) with a
106 state-of-the-art text-based semantic model (Pennington et al. 2014) to enhance decoding of a large
107 fMRI dataset. The experiential model is based on peoples’ ratings of their
108 sensory/motor/affective/cognitive experiences with words and their referents (building on Cree and
109 McRae 2003, Vinson et al. 2003, Lynott and Connell, 2013). Whilst experiential models have provided a
110 basis for neural decoding (Chang et al. 2010, Fernandino et al. 2015, Fernandino et al. 2016, Anderson
111 et al. 2016a, Yang et al. 2017, Wang et al. 2017, Anderson et al. 2018) it has never been clear whether
112 and how text-based and experiential approaches complement one other. We newly demonstrate that
113 text-based and experiential models differentially contribute to decoding sentences that do/don’t contain
114 linguistically oriented “abstract” words. Finally, we introduce a cross-participant decoding method to

115 estimate the room for improvement in model-based decoding and pinpoint model weaknesses for future
116 development.

117

118 **Methods**

119 **Overview**

120 We reanalyzed an fMRI data set scanned as 14 people read 240 sentences describing everyday
121 situations (Anderson et al. 2016a and summarized below). Sentences were 3 to 9 words long and
122 formed from 242 different content words. 10 participants saw the set of sentences repeated 12 times in
123 total (randomly shuffled each time), and the remaining 4 participants who attended half the number of
124 visits saw the sentences 6 times. Following standard fMRI preprocessing steps, each sentence
125 presentation was represented as a single fMRI volume (there were 12 replicate volumes per sentence
126 for 10 participants and 6 replicates for the remaining 4 participants. Analyses were focused on a
127 “semantic network” of 22 anatomical regions-of-interest (ROIs) that had been detected in previous
128 analysis of the same data (Anderson et al. 2018) testing for regional sensitivity to experiential semantic
129 features associated with words with different grammatical roles. ROIs included left temporal, inferior
130 parietal, inferior/superior frontal cortex as well as some right hemispheric homologs (illustrated later in
131 **Figure 6**). These regions have well established associations with semantic processing (e.g. Binder et
132 al., 2009, Binder and Desai, 2011) and broadly adhere to the “language network” identified by
133 Fedorenko and Thompson-Schill (2014). fMRI activation across the network of brain regions was then
134 decoded using a text-based model, an experiential model, and the two models integrated as detailed
135 below. We refer to the integrated text/experiential approach as “multimodal” to reflect the combination
136 of linguistic information with behavioral ratings. Though to dispel any confusion the experiential model
137 serves as a proxy for knowledge acquired through multiple modalities of experience, just each modality
138 is estimated through the same rating procedure.

139

140 **Materials**

141 All sentences were pre-selected as experimental materials for the Knowledge Representation in Neural
142 Systems (KRNS) project (Glasgow et al. 2016, www.iarpa.gov/index.php/research-programs/krns),
143 sponsored by the Intelligence Advanced Research Projects Activity (IARPA). The stimuli consisted of
144 240 written sentences containing 3-9 words and 2-5 (mean \pm sd = 3.33 \pm .76) content words, formed from
145 different combinations of 141 nouns, 62 verbs, and 39 adjectives (242 words). The sentences are listed
146 in full in the Supplementary Materials of Anderson et al. (2016) and Anderson et al. (2018). Sentences
147 were in active voice and consisted of a noun phrase followed by a verb phrase in past tense, with no
148 relative clauses. Most sentences (200/240) contained an action verb and involved interactions between
149 humans, animals and objects, or described situations involving different entities, events, locations, and
150 affective connotations. The remaining 40 sentences contained only a linking verb ("was"). The entire list
151 is in Table S1. Each word occurs a mean \pm sd [range] of 3.3 \pm 1.7 [1 7] times throughout the entire set of
152 sentences and co-occurs with 8.1 \pm 4.3 [1 19] other unique words. The same two words rarely co-occur
153 in more than one sentence, and 213/242 words never co-occur more than once with any other single
154 word. Forty-two sentences contained instances of words not found in any of the other 239 sentences,
155 and 3 of these sentences contained 2 unique words. There were thus 45 words that occurred in only
156 one sentence, of which 29 were nouns, 7 were verbs and 9 were adjectives.

157

158 **Participants**

159 Participants were 14 healthy, native speakers of English (5 males, 9 females; mean age = 32.5, range
160 21-55) with no history of neurological or psychiatric disorders. All were right-handed according to the
161 Edinburgh Handedness Inventory (Oldfield 1971). Participants received monetary compensation and
162 gave informed consent in conformity with the protocol approved by the Medical College of Wisconsin
163 Institutional Review Board.

164

165 **Procedure**

166 Participants took part in either 4 or 8 scanning visits. The mean interval between sessions was 3.5 days
167 (SD = 3.14). The range of the intervals between first and last visits was 15 - 43 days. In each visit, the

168 entire list of sentences was presented 1.5 times, resulting in 12 presentations of each sentence over
169 the 8 visits in 10 participants, and 6 presentations over 4 visits in 4 participants. Each visit consisted of
170 12 scanning runs, each run containing 30 trials (one sentence per trial) and lasting approximately 6
171 minutes. The presentation order of each set of 240 sentences was randomly shuffled.

172

173 The stimuli were back-projected on a screen in white Courier font on a black background. Participants
174 viewed the screen while in the scanner through a mirror attached to the head coil. Sentences were
175 presented word-by-word using a rapid serial visual presentation paradigm. Nouns, verbs, adjectives,
176 and prepositions were presented for 400 ms each, followed by a 200-ms inter-stimulus interval (ISI).
177 Articles ("the") were presented for 150 ms followed by a 50-ms ISI. Mean sentence duration was 2.8 s.
178 Words subtended an average horizontal visual angle of approximately 2.5°. A jittered inter-trial interval,
179 ranging from 400 ms to 6000 ms (mean = 3200 ms), was used to facilitate deconvolution of the BOLD
180 signal. Participants were instructed to read the sentences and think about their overall meaning. They
181 were told that some sentences would be followed by a probe word, and that in those trials they should
182 respond whether the probe word was semantically related to the overall meaning of the sentence by
183 pressing one of two response keys (10% of trials contained a probe). Participants' mean accuracy was
184 86% correct, with a minimum accuracy of 81%. Participants were given practice with the task outside
185 the scanner with a different set of sentences. Response hand was counterbalanced across scanning
186 visits.

187

188 **MRI parameters and preprocessing**

189 MRI data were acquired with a whole-body 3T GE 750 scanner at the Center for Imaging Research of
190 the Medical College of Wisconsin using a GE 32-channel head coil. Functional T2*-weighted
191 echoplanar images (EPI) were collected with TR = 2000 ms, TE = 24 ms, flip angle = 77°, 41 axial
192 slices, FOV = 192 mm, in-plane matrix = 64 x 64, slice thickness = 3 mm, resulting in 3 x 3 x 3 mm
193 voxels. T1-weighted anatomical images were obtained using a 3D spoiled gradient-echo sequence with
194 voxel dimensions of 1 x 1 x 1 mm. fMRI data were pre-processed using AFNI (Cox, 1996). EPI volumes

195 were corrected for slice acquisition time and head motion. Functional volumes were aligned to the T1-
196 weighted anatomical volume, transformed into a standardized space (Talairach and Tournoux 1988),
197 and smoothed with a 6mm FWHM Gaussian kernel. The data were analyzed using a general linear
198 model with a duration-modulated HRF, and the model included one regressor for each sentence.
199 Neural activity was modeled as a gamma function convolved with a square wave with the same
200 duration as the presentation of the sentence, as implemented in AFNI's 3dDeconvolve with the option
201 dmBLOCK. Duration was coded separately for each individual sentence. Finally a single sentence-level
202 fMRI representation was created for each unique sentence by taking the voxelwise mean of all
203 replicates of the sentence.

204

205 **Experimental Design and Statistical Analysis**

206 **Semantic models**

207 Analyses used the following two semantic models of word meaning. As a proxy for the representational
208 structure that can be acquired the distributional statistics of words in language we used “GloVe”
209 (Pennington et al. 2014). GloVe is a freely downloadable text-based semantic model that represents
210 individual words as 300 dimensional floating point vectors derived by factorizing a word co-occurrence
211 matrix (vocabulary size is 2.2million words and co-occurrences were measured across 840 billion
212 tokens from Common Crawl <https://commoncrawl.org>). GloVe in particular was used because it yielded
213 state-of-the-art performance decoding fMRI activation associated with sentences in Pereira et al.'s
214 (2018) “universal neural decoder of linguistic meaning”, although we found there to be relatively minor
215 differences between using GloVe and other models such as “word2vec” (Mikolov et al. 2013) – see also
216 the final **Results** section.

217

218 As a proxy for the representational structure that can be acquired from direct experience with the world
219 we used an experiential attribute model (Binder et al. 2016). This model represents words in terms of
220 human ratings of their degree of association with different attributes of experience (e.g., “On a scale of
221 0 to 6, to what degree do you think of a banana as having a characteristic or defining color?”). Ratings

222 were collected via Amazon Mechanical Turk for a total of 65 attributes spanning sensory, motor,
223 affective, spatial, temporal, causal, social and abstract cognitive experiences. Ratings for each attribute
224 were averaged across workers to derive a single 65 dimensional vector for each word. As such this
225 model broadly aligns with “embodiment” theories that posit representations of word meaning reflect a
226 summarization of the brain states involved in experiencing that word, partially reenacted across
227 sensory/motor/affective/cognitive subsystems (e.g. Barsalou et al. 2008, Glenberg 2010, Pulvermüller
228 2013, Binder et al. 2016). This same experiential model has previously been used as the basis for
229 predicting and decoding the same fMRI dataset as the current study (Anderson et al. 2016a; Anderson
230 et al. 2018).

231

232 Some overlap in the semantic information content of text-based and experiential models is very much
233 expected (see also Riordan and Jones, 2010). Because text describes worldly experiences we expect it
234 to partially capture the structure of experiential knowledge. On the flipside, the experiential attribute
235 model seeks to comprehensively model experiential knowledge, and of course language contributes to
236 our experience. Beyond this, the experiential model was itself built through a linguistically guided rating
237 procedure. However, systematic differences between the models are also expected. This is in part
238 because a lot of experiential information goes unstated in natural verbal communication. For instance,
239 borrowing an example from Bruni and Baroni (2014), it is rarely useful to communicate the color of
240 bananas because it is obvious to all those with experience of bananas. Likewise, it would be unusual to
241 specify that dropping things involves movement. Consequently, whilst an analysis of natural text may
242 indicate that a banana is an edible berry, it may not capture the dominance of color as a perceptual
243 attribute. Therefore, despite being derived via language, attribute ratings can potentially anchor to
244 experiential neural systems and access information that would not otherwise have been reported or
245 experienced in natural verbal communication. Conversely, the experiential attribute model as it stands
246 may be less well suited to capturing the meaning of so-called abstract words which tend to be more
247 amenable to verbal description in terms of their relationships with other words (e.g. “fiction” is an
248 imaginary story), rather than through physical example (e.g. being presented with a cat).

249

250 We consider the extent to which experiential information is available in text (and language) to be an
251 empirical question. For instance, representational similarity analysis (Kriegeskorte et al. 2008) can be
252 applied to compare the information structure of the text-based to the experiential models. This yields a
253 statistically significant correlation coefficient of $p=0.2$, $p<10^{-6}$ (reflecting Spearman correlation between
254 the below diagonal triangles of inter-word Pearson correlation matrices derived using the text-based
255 and experiential models, $n=29161$ correlation coefficients per triangle as computed from 242 words). A
256 related regression analysis conducted by Utsumi (2018) demonstrated that text-based models were
257 weakly capable of predicting spatial, temporal and affective attributes of the current experiential model.
258 However, demonstrating that text-based models do not contain aspects of experiential information does
259 not also entail that the missing information is relevant for explaining semantic brain activity.
260 Consequently, we conduct our forthcoming analyses by testing for brain activity that is explainable
261 using the experiential model but not using the text-based model and vice versa (see also Anderson et
262 al. 2015, Popov et al. 2018). This approach takes the *assumption* that the text-based model accurately
263 captures all semantic information that can be extracted from text alone. Whilst we acknowledge that
264 this is a strong assumption that is not likely to have been strictly met here, we believe that current
265 models are sufficiently advanced to begin to segregate experiential from linguistic aspects of semantic
266 representation in brain activity (see also the **Discussion**).

267

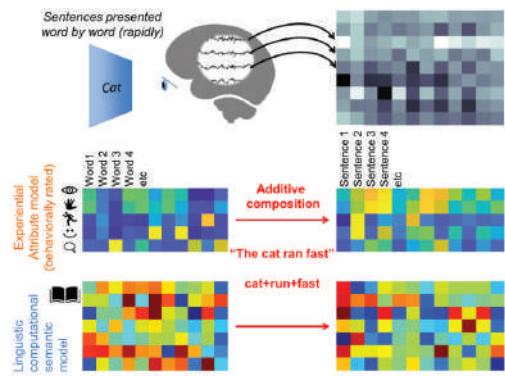
268 **Modeling sentences by summing word-level semantic vectors**

269 To turn word-level semantic vectors into representations of sentences we identified all constituent
270 content words in each sentence, and then pointwise summed together these semantic vectors (**Figure**
271 **1**). Although such additive composition is obviously an oversimplification that neglects the effects of
272 word order, syntax and morphology, it has endured as a practically successful technique in both
273 computational linguistics (Mitchell and Lapata 2010, Kiela et al. 2014), and fMRI analyses (Anderson et
274 al. 2016a; Yang et al. 2016; Wang et al. 2017; Pereira et al. 2018; Anderson et al. 2018). Indeed,

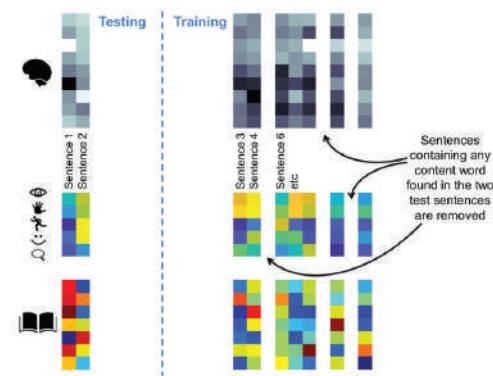
275 attempts to incorporate other linguistic factors such as syntax into models have yet to make appreciable
 276 difference to neural decoding performance (Pereira et al. 2018; Anderson et al. 2018).

277

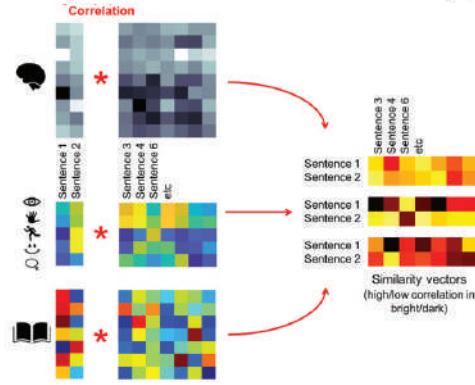
1. Neural, experiential and text-based sentences



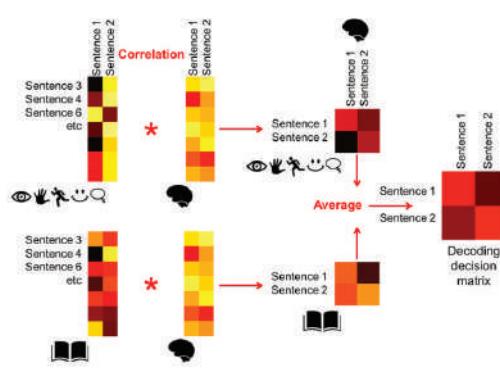
2. Cross validation training/test split



3. Re-representation in common similarity space



4. Multimodal model-based decoding



278
 279
 280

281 **Figure 1. Representational similarity-based decoding algorithm** set up to support multiple model-
 282 based decoding. Multimodal model combination takes place in stage 4 by averaging 2×2 decoding
 283 decision matrices generated by the different models. An alternative approach would have been to
 284 pointwise average together the two similarity vectors for the experiential model with those of the text-
 285 based model in stage 3. This was disfavored to avoid having to introduce an extra normalization step to
 286 deal with correlation coefficients arising from the different models being on different scales (correlation
 287 coefficient magnitudes tend to diminish as the number of features correlated becomes large, and here
 288 the experiential and text-based models widely differ in the number of features – 65 and 300
 289 respectively). This problem is naturally dealt with in stage 4, because the 2×2 decision matrices are
 290 based on correlations between similarity vectors that are all matched in their dimensions.

291

292

293

294

295 **Representational similarity-based neural decoding set up**

296 To decode fMRI activation, we applied the representational similarity (see also Kriegeskorte et al. 2008)
297 decoding framework introduced in Anderson et al. (2016b) and Anderson et al. (2017) and further
298 extended here to support simultaneous multi-model / multi-ROI / multi-participant decoding of
299 sentences. This was a considered departure from the more commonplace strategy of using multiple
300 regression to map between model and brain (Mitchell et al. 2008, Chang et al. 2010, Sudre et al. 2012,
301 Pereira et al. 2013, Wehbe et al. 2014, Fernandino et al. 2015, Fernandino et al. 2016, Huth et al.
302 2016, Anderson et al. 2016a, Yang et al. 2017, Wang et al. 2017, Pereira et al. 2018, Anderson et al.
303 2018). The main reason for choosing the similarity-based approach over (ridge) regression here was for
304 simplicity: to avoid repeating the analyses multiple times over with different regularization penalties and
305 the need to introduce a decision over which penalty to use. For the current analyses, this would
306 complicate the process of integrating information across models, ROIs and individuals (because in
307 each case there would be multiple results associated with the different penalties and multiple decisions
308 to be made). The current focus on the similarity-based approach should not be misconstrued as a claim
309 that similarity-based methods are superior to regression, and we report on some strengths and
310 weaknesses of the two approaches in a supporting analysis using ridge regression (see Results, and
311 **Figure 10**). Other comparative analyses are in Anderson et al. (2016b) and Bulat et al. (2017).

312

313 fMRI data was decoded according to a commonly used leave-2-item-out cross validation procedure
314 (Mitchell et al. 2008, Chang et al. 2010, Sudre et al. 2012, Pereira et al. 2013, Wehbe et al. 2014,
315 Anderson et al. 2016a,b, Yang et al. 2017, Wang et al. 2017, Pereira et al. 2018, Anderson et al. 2018).
316 At each cross-validation iteration the 240 sentences were split into a test set of 2 sentences, and a
317 training set of 238 sentences. Then both fMRI and model data for any of the 238 training sentences that
318 contained any content word within the 2 test sentences was deleted from the training data (**Figure 1**)

319 **stage 2).** This was to enable testing of how well the approach generalized to decoding novel fMRI
320 sentences using sentence vectors built from an entirely novel set of semantic vectors. The mean+/-SD
321 number of sentences in the training set for each iteration was 218+/-5, containing a mean+/-SD of
322 232+/-2 words). Model and fMRI data corresponding to the training set was featurewise/voxelwise z-
323 scored. Model and fMRI test sentence data were also normalized by subtracting and dividing by the
324 feature/voxelwise mean and SD of the training data.

325

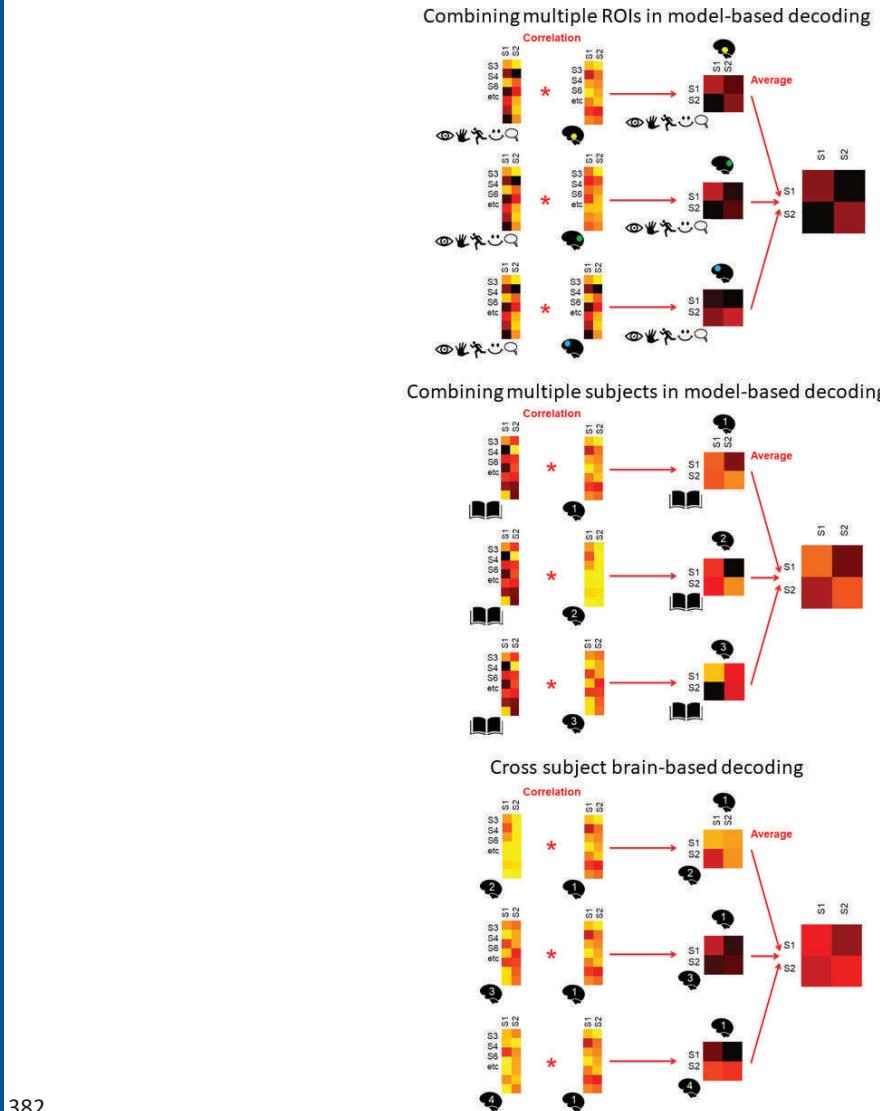
326 Because not all fMRI voxels contain informative signal, we estimated which ones were likely to be
327 informative using a commonly used strategy (e.g. Mitchell et al. 2008, Chang et al. 2010, Pereira et al.
328 2013, Anderson et al. 2015, Anderson et al. 2016b, Anderson et al. 2017, Yang et al. 2017, Wang et al.
329 2017). For each participant, and separately for each ROI, we took each of the 12 (or 6) fMRI runs
330 through the entire set of sentences, selected only the 218 (or so) training sentences from this and then
331 voxelwise correlated each unique pair of runs together. For the ten participants with 12 runs, this left 66
332 pairwise correlation coefficients per voxel, and for the four participants with 6 runs this left 15 pairwise
333 correlation coefficients per voxel. A single score was assigned to each voxel by taking the mean of
334 these (66 or 15) correlation coefficients. The 50 voxels with the largest mean value per ROI were
335 selected for analysis. This choice of 50 voxels was ultimately arbitrary though guided by previous work
336 (e.g. Anderson et al. 2013; Anderson et al. 2015, Anderson et al. 2017). In subsequent post hoc
337 analyses (reported later) the interpretation of results was found to be unchanged when using 100
338 voxels, however in both cases voxel selection systematically improved decoding accuracies over no
339 voxel selection at all.

340

341 Decoding of the fMRI data was accomplished by independently re-representing both model and fMRI
342 test data in a common representational similarity space, and then matching model to fMRI sentences in
343 this space. For each fMRI/model dataset in turn, we correlated the two test sentence vectors with each
344 one of the 218 (or so) training sentence vectors. This enabled us to newly represent every single test
345 sentence (whether model or fMRI) as a similarity vector of 218 (or so) correlation coefficients (**Figure 1**

346 **stage 3).** We transformed all correlation coefficients in all similarity vectors using Fisher's r-to-z
347 (arctanh) as is a customary treatment when comparing correlation values (though this had only a
348 marginal effect on the current results). Then decoding was achieved by cross correlating the two model
349 similarity-vectors with the two fMRI similarity-vectors (**Figure 1 stage 4**). This resulting 2 * 2 matrix of
350 correlation coefficients was r-to-z transformed (arctanh) and this constituted the "decoding decision".
351 This was evaluated as correct (and scored as a 1 as opposed to zero) only if the sum of z-transformed
352 correlations between the correctly matched model vs fMRI test sentence pair (the matrix diagonal in
353 **Figure 1 stage 4**), exceeded the sum for the incongruent pair (anti-diagonal in **Figure 1 stage 4**).
354 Decoding was repeated for all possible unique pairs of the 240 sentences (28 680 cross validation
355 iterations in total) and the mean score used as the final metric of decoding accuracy. When operating at
356 random (e.g. if the fMRI data contained no semantic signal, or the semantic vectors did not reflect
357 meaning) a mean decoding accuracy of 0.5 is expected. Permutation testing was applied to test
358 whether decoding accuracies were significantly better than chance (by randomly shuffling semantic
359 model sentences and repeating the entire cross-validation analysis 1000 times over, as described in
360 Anderson et al. 2017). Typically, final decoding accuracies above 0.56 were significant at the p=0.05
361 level (95% of metrics derived in analyzing randomly shuffled sentences were less than or equal to
362 0.56).
363
364 For each participant 2 * 2 decoding decisions were computed in parallel for all 22 ROIs and each of the
365 2 models. To generate a single decoding decision corresponding to the entire semantic network of all
366 22 ROIs, we applied an "ensemble averaging" strategy and pointwise averaged together decision
367 matrices across the 22 ROIs (**Figure 2 top row**). We use this network-level decoding estimate as the
368 basis for our main analyses but also report results for individual ROIs.
369
370 In a similar vein, we used ensemble averaging as the basis for testing whether models have
371 complementary strengths in decoding. If the models have complementary strengths, then integrating
372 their decisions together as an ensemble will counteract individual model's weaknesses and boost

373 overall decoding accuracy. Multimodal model integration was achieved by pointwise averaging decision
374 matrices associated with the different models as illustrated in **Figure 1 stage 4**. In all ensemble
375 averaging tests (whether integrating ROIs and/or models) decoding decision matrices were scored as
376 correct precisely as before by testing whether correlations on the diagonal were greater than the
377 antidiagonal. To produce a final metric, scores were then averaged across all 28680 cross validation
378 iterations. Importantly, this ensemble averaging strategy is *not* guaranteed to produce equivalent or
379 better final accuracies than the strongest model of the pair (which would limit its applicability for testing
380 for a multimodal decoding advantage). Specifically, if one model is sufficiently noisy, then the final
381 multimodal decoding accuracy will be lower than the strongest model.



382

383 **Figure 2. Representational similarity-based algorithm set ups for ensemble decoding.** Top,
 384 model-based decoding of multiple brain regions in the same participant (see also results in **Figures**
 385 **4,5,6,8,9,11**). Middle, model-based decoding of multiple participants (see also results in **Figure 5**).
 386 Bottom, cross-subject decoding (see also results in **Figure 8**).

387

388 The entire cross-validation procedure described above was repeated for each ROI (22) within each
 389 participant (14) using both semantic models (2). Ensemble averaging of model decision matrices was

390 used to derive multimodal decoding accuracies for each participant and ROI. Ensemble averaging of all
391 22 ROIs' decision matrices was used to generate a single network-level decoding accuracy for each
392 participant and model combination (text, experiential and multimodal). Differences in accuracy between
393 different models were evaluated using t-tests, and p-values associated with multiple ROIs corrected for
394 multiple comparisons according to False Discovery Rate (FDR) (Benjamini and Yekutieli, 2001).

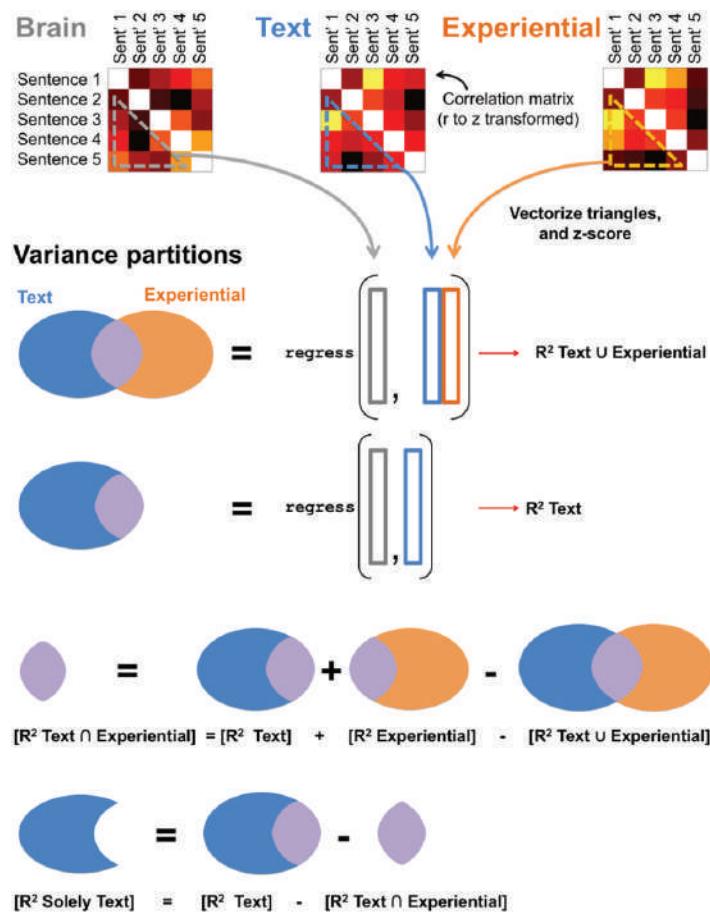
395

396 **Post hoc analyses partitioning the variance in neural similarity structure explained by each
397 model**

398 ROIs for which decoding accuracy was boosted through model integration were taken to post hoc
399 analyses. We further estimated the unique contribution made by each model to explaining variance in
400 the neural sentence similarity structure and what could be explained equally by both models. This
401 analysis is inspired by de Heer et al. (2017) who partitioned the variance in heard speech fMRI data
402 that was unique to acoustic, articulatory and semantic voxelwise encoding models and shared across
403 them. Besides using different models to de Heer et al. (2017) the forthcoming analysis differs in that it is
404 conducted in representational similarity space. The analysis is illustrated in **Figure 3**.

405 The analysis was conducted on the entire similarity space defined by all 240 sentences as is a standard
406 approach in Representational Similarity Analysis (RSA) (Kriegeskorte, 2008). Sentence similarity
407 matrices were computed for each participant, for each ROI, by inter-correlating the neural
408 representations of the 240 sentences. In each case this yielded a 240 by 240 correlation matrix.
409 Pearson correlation was used and the correlation coefficients were subsequently r to z transformed
410 (arctanh). Prior to this, voxel selection was conducted to estimate the 50 informative voxels per ROI.
411 Voxel selection used the same correlation-based approach as detailed above for the previous decoding
412 analysis. However, in the case at hand voxel selection was computed only once on all 240 sentences
413 together (because the current RSA did not employ cross validation). To generate a single correlation
414 matrix capturing the 22 ROI ensemble, the 22 r to z transformed correlation matrices corresponding to
415 each ROI were pointwise averaged.

416 Correlation matrices for the text-based and experiential models were computed in the same fashion by
 417 inter-correlating the 240 sentences. Again, Pearson correlation was used and all coefficients were r to z
 418 transformed. Next, the unique information contained within every single correlation matrix was
 419 extracted by segmenting the below diagonal matrix triangle. Each triangle was then vectorized to create
 420 a 28680 element similarity vector. Similarity vectors were subsequently normalized by z-scoring to
 421 support the forthcoming regression analysis.



423 **Figure 3 Partitioning the variance in neural similarity structure that is solely accounted for
 424 individual models and shared between them.**

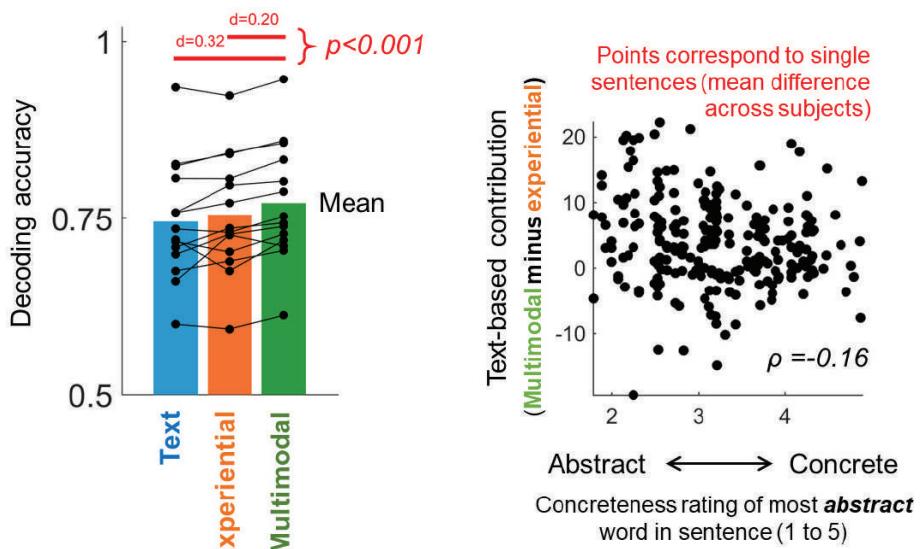
425

426 We applied the set theoretic approach of de Heer et al. (2017) to estimate the variance in neural
427 similarity structure that was explained by the *union* of both models, the *shared* variance that is equally
428 explained by either model, and the variance *solely* accountable to one model (see **Figure 3**). The union
429 [R^2 Text \cup Experiential] was first estimated using a multiple regression in which the similarity vectors for
430 both models were used as predictors and the neural similarity vector was the response variable. The
431 variance associated with but not necessarily exclusive to the individual models ($[R^2$ Text] and [R^2
432 Experiential]) was estimated in 2 separate regression analyses. In each a single similarity vector
433 (associated with one model) was the predictor. To estimate the shared variance [R^2 Text \cap Experiential]
434 we subtracted away the variance explained by the model union [R^2 Text \cup Experiential] from the sum of
435 the variance explained by the Text and Experiential models ($[R^2$ Text] + [R^2 Experiential]). Then, to
436 estimate the variance solely accountable to the text and experiential models ($[R^2$ Solely Text] and [R^2
437 Solely Experiential]) we subtracted the shared variance away from the variance explained by the
438 models (e.g. [R^2 Text] - R^2 Text \cup Experiential]). To produce positive correlation coefficients from these
439 measures, we took the square root of R^2 values as undertaken by de Heer et al. (2017). However,
440 because variance/positive correlation measures do not facilitate testing whether individual models
441 made significantly greater than zero contribution to explaining neural data (because they are always
442 greater than or equal to zero) we also undertook a partial correlation analysis (see also Anderson et al.
443 2015, Wang et al. 2018). Here we computed the correlation between neural similarity structure and one
444 model, whilst controlling for the other model. To test the generality of partial correlation coefficients
445 across participants, we compared them to zero using one sample t-tests (n=14).

446 **Code accessibility**

447 MATLAB similarity-based decoding code is available on request from the corresponding author.

449



450
451
452 **Figure 4. Integrating text-based and experiential models produces stronger decoding.** Individual-
453 level accuracies arising from decoding the 22 ROI ensemble (see **Figure 2** top row). The contribution of
454 the text-based model to multimodal decoding was particularly pronounced for sentences containing
455 abstract words (right). Effect sizes (d) were estimated according to Dunlap et al. (1996) as $d=t^*(2*(1-r)/n)$, where t is the t-statistic arising from the corresponding paired t-test, r is Pearson correlation, and n
456 is the number of participants (14).
457

458

459

460

461 Results

462 Integrating text-based and experiential semantic models produces stronger decoding than 463 either alone

464 To test for evidence that both linguistic and non-linguistic experiential aspects of meaning were present
465 in neural activation, we tested whether combining the two models together improved decoding
466 accuracies above using either model in isolation. Decoding accuracies for both models and their
467 combination are illustrated in **Figure 4**. T-tests revealed that whilst there were no differences in
468 decoding accuracy between the text-based and experiential models, when the models were combined
469 accuracies were significantly greater than for either model in isolation. (Difference between multimodal
470 and text, $t=6.9$, $p<0.0001$; difference between multimodal and experiential, $t=4.8$, $p=0.0004$). This key

471 result provides direct neural evidence that linguistic and experiential semantic information were present
472 in brain activation elicited in sentence reading. Otherwise, decoding accuracies were significant for all
473 participants and all models.

474

475 **Text-based model enhances discrimination of sentences containing abstract words**

476 We next examined the nature of the contribution made by the text-based and experiential model in
477 more detail. An obvious area to anticipate divergence between models is for more linguistically oriented
478 “abstract” words. These words do not directly correspond to “concrete” entities in the world that can be
479 directly sensed. As such an abstract word’s meaning is language dependent and most amenable to
480 description in terms of other words which could contribute different parts of that meaning (Brysbaert et
481 al. 2014). Examples of words that are relatively abstract within the current dataset include: ‘negotiate’,
482 ‘agreement’, ‘wealthy’, ‘famous’ and ‘clever’. We hypothesized that the text-based model (which is built
483 from word co-occurrence statistics) would make a particular contribution to multimodal decoding of
484 sentences containing abstract words.

485

486 To test for an abstractness advantage associated with integrating the text-based model in decoding, we
487 looked up concreteness ratings (Brysbaert et al. 2014) for each of the 242 content words in the
488 sentence set. Each of the 240 sentences was then scored according to (1) the concreteness of the
489 least concrete (most abstract) word in the sentence; (2) the concreteness of the most concrete word in
490 the sentence; (3) the mean concreteness of all words in the sentence. Because the experimental
491 sentences also varied across many other factors which through coincidence might be confounded with
492 concreteness measures, we attempted to take these into account. Specifically, we additionally scored
493 each sentence according to the least frequent, most frequent and mean frequency of content words in
494 the sentence (derived from log2 transformed SUBTLEX-US counts of Brysbaert and New, 2009), the
495 number of words in each sentence, and the minimum, maximum and mean word length (number of
496 characters). After this each sentence was represented with 10 measures (including the 3 concreteness
497 measures).

498

499 We then identified how well each individual sentence was decoded using the different models. For each
500 participant ($n=14$), this yielded a vector of 240 sentence decoding scores for the text-based model, the
501 experiential model and the multimodal combination. The maximum score attainable (and maximum
502 value in the vector) was 239, which could have been achieved if the corresponding sentence was
503 successfully discriminated from all 239 other sentences during the entire leave-2-out cross validation
504 analysis.

505

506 To estimate the independent contribution made by the text-based model to multimodal decoding, we
507 pointwise subtracted the experiential model-based decoding scores for individual sentences away from
508 corresponding multimodal model scores (and repeated for each participant). Positive scores arising
509 from this subtraction indicate sentences that were better discriminated by integrating the text-based
510 model, which in turn we hypothesized relate to a measure of sentence abstractness. To test this, we
511 correlated (Spearman) each participant's vector of "boost" values with the three different sentence
512 concreteness measures (leaving $14 * 3$ correlation coefficients). We repeated these correlations for the
513 other seven sentence measures (re: word frequencies and lengths of constituent words and the number
514 of words in the sentence).

515

516 To test for the generality of positive correlations across participants, participants' correlation coefficients
517 were r -to- z transformed (arctanh) and then compared to zero using a one sample t-test ($n=14$). T-tests
518 were repeated on correlation coefficients associated with each of the 10 sentence measures. The
519 resultant 10 p-values were FDR corrected. Only one of the ten t-tests yielded a statistically significant
520 result. This was the test based on correlations between the concreteness rating of the most abstract
521 word in the sentence and the decoding boost ($t=-3.7$, $p=0.04$, FDR corrected). The mean correlation
522 coefficient across participants was -0.08. This provided evidence that the decoding advantage brought
523 by integrating the text-based model was related to better discrimination of sentences containing
524 abstract words. This relationship is illustrated in **Figure 4** (note that the figure illustrates the mean

525 accuracy boost per sentence across all 14 participants with associated correlation coefficient of -0.16,
526 whilst the above t-test was based on 14 individual-level correlation coefficients, with a mean value of -
527 0.08).

528

529 In an attempt to gain an intuition of whether there was a common theme to the sentences containing
530 abstract words that had received a decoding boost we listed them. However, we were unable to
531 confidently pinpoint any systematic pattern. The ten sentences that received the greatest decoding
532 boost are as follows with the most abstract word in each sentence in italics. The abstract word's
533 concreteness rating (C) and the mean decoding boost (B) associated with the sentence is in brackets
534 after the sentence. "The family was *happy*." (C=2.6, B=22); "The team *celebrated*." (C=2.9, B=21); "The
535 *patient* survived." (C=2.5, B=21); "The family *survived* the powerful hurricane." (C=2.6, B=20); "The pilot
536 was *friendly*." (C=2.3, B=20); "The flood was *dangerous*." (C=2.1, B=20); "The man *lost* the ticket to
537 soccer." (C=2.3, B=20); "The jury watched the *witness*." (C=4.1, B=19); "The council read the
538 *agreement*." (C=2.2, B=18); "The artist *shouted* in the hotel." (C=4.2 B=18). On face value a
539 commonality would appear to be that many of the sentences have affective connotations. However,
540 such affective connotations were not exclusive to the boosted sentences. Indeed, three of the four
541 sentences that were most disadvantaged by integrating the text-based model were also valenced.
542 These four sentences were "The team *lost* the football in the forest." (C=2.3 B=-19), "The teacher broke
543 the *small* camera." (C=3.2, B=-15), "The *aggressive* team took the baseball." (C=2.5, B=-12.4), "The
544 actor *gave* the football to the team." (C=2.8, B=-12.5). We therefore presume simply that the text-based
545 model helped explain neural signal reflecting the additional linguistic information exposed in accessing
546 abstract words and integrating their meaning into sentences. However, in the future it will also be
547 valuable to consider stimuli that are more controlled in their content and less dominated by concrete
548 sentences than the current dataset.

549

550 We next ran an analogous analysis attempting to understand the contribution of the experiential model
551 which we anticipated might be associated with concreteness. Sentence-wise decoding score vectors

552 (length 240, maximum value 239) for the text-based model were subtracted from the multimodal model
553 to generate an “accuracy boost” vector for each participant. Correlations between boost vectors and all
554 ten sentence measures were computed. Whilst a significant positive correlation was detected with the
555 concreteness of the most abstract word in the sentence – i.e. the contribution of the experiential model
556 was especially associated with concrete sentences without any abstract words (average correlation
557 across participants=-0.045, p=0.02, disappointingly this result did not survive FDR correction for
558 multiple comparisons and should be treated tentatively for the time being.

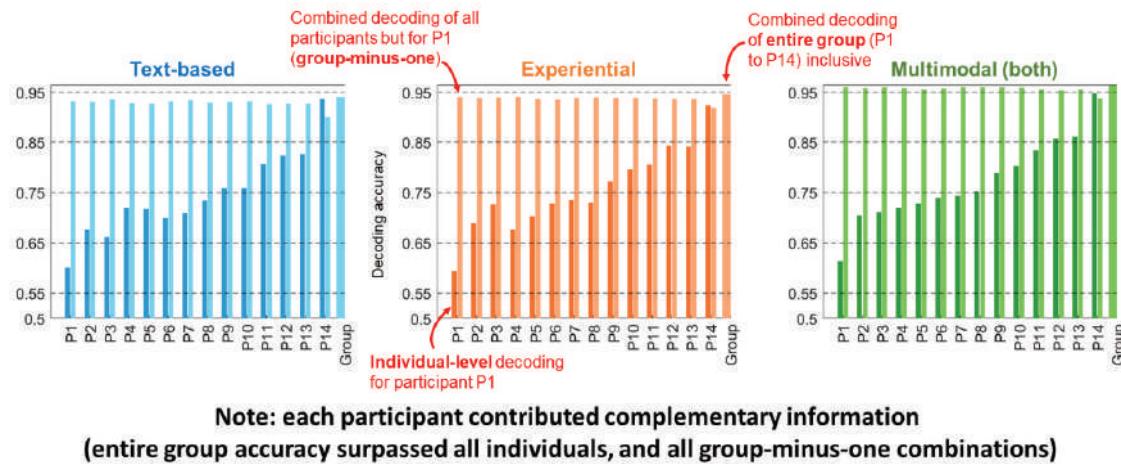
559

560 **Model-based neural decoding at group-level leverages cross-participant regularities**

561 A fundamental difference between either of the semantic models and individuals' fMRI data was that
562 models were generic representations of concepts built from group-level information (either from texts
563 written by many authors, or ratings given by many people), whereas fMRI activation captured
564 snapshots of an individual's interpretation of a sentence at the particular time of scanning, and in the
565 broader context of their own personal experiences. We therefore reasoned that combining individuals'
566 fMRI data together as a group should expose regularities in semantic representations across
567 individuals and lead to a stronger pattern match to the group-level models (**Figure 2 middle**). Aside
568 from this the group-level combination should also iron out noise (whether this arises due to
569 technological reasons or participants attention levels and/or compliance with the task). Either way, this
570 should lead towards a less noisy comparison between model and fMRI data, albeit at the expense that
571 the results may not generalize to individuals (though this has already been tested in **Figure 4**).

572

573



574
 575
 576
577 Figure 5. Decoding neural data at group-level exploits cross participant regularities. Model-
 578 based decoding accuracies at group-level (**Figure 2 middle**) and for each individual corresponding to
 579 all 22 ROIs decoded as an ensemble. Individual participants' decoding accuracies (dark) are plotted
 580 beside group-minus-one (light) decoding accuracies derived using all other participants. "Group" is all
 581 14 participants combined.

582
 583
 584 Combining fMRI activation across participants is in general complicated by both anatomical and
 585 functional differences between individuals. Whilst sophisticated "hyper-alignment" methods for
 586 combining group fMRI data exist (Haxby et al. 2011, Guntupalli et al. 2016), in the current case it is also
 587 possible to combine individuals together as a group, by averaging together individual's decoding
 588 decision matrices (see **Figure 2 middle**) in precisely the same way as we have combined models and
 589 ROIs.

590
 591 Group-level decoding accuracies, illustrated in **Figure 5**, were unanimously greater than average
 592 individual-level decoding accuracies. For example, where the mean individual-level decoding accuracy
 593 was 0.77 the group-level result was 0.97 (as a side note, although this may seem surprisingly high
 594 accuracy, this score is consistent with previous word-level decoding studies (Anderson et al. 2017,
 595 Anderson et al. 2016b) and also reflects the high sensitivity of the leave-2-out test). To confirm the
 596 statistical significance of this effect we used on sample t-tests to test the individual-level decoding

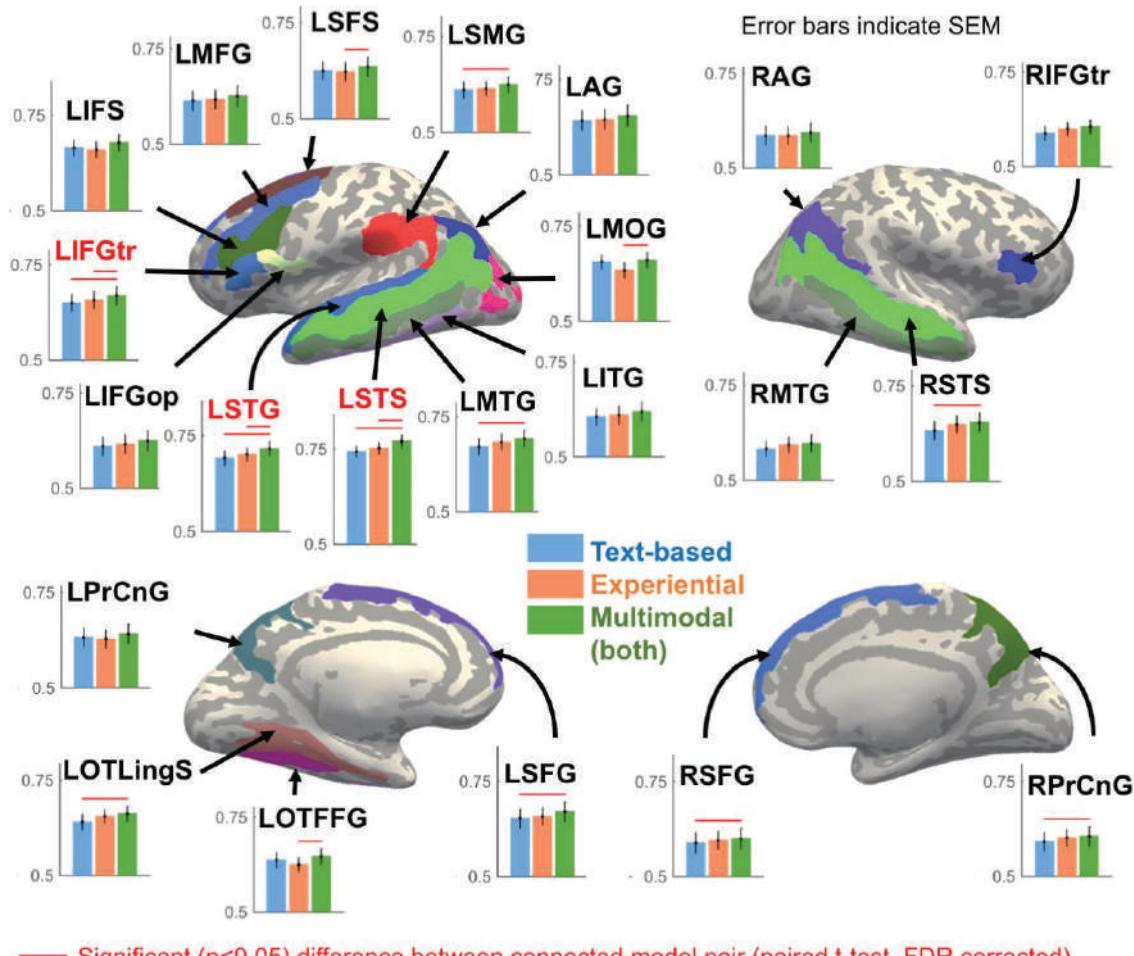
597 results against the single group-level score. The outcome was significant ($p<0.05$) in every test
598 following FDR correction (for all ROIs and for both models and their multimodal combination).

599

600 In conducting group-level analyses it is possible that individual participants play a dominant role in
601 results (an individual may have just happened to elicit semantic representations that match the models
602 well, or have been particularly attentive to the task). To examine the influence of individual participants
603 we re-ran the group-level analysis holding out each individual from the group in turn. Group minus
604 participant decoding accuracies are plotted next to the held out participant in **Figure 5**. It became clear
605 that one participant (P14) indeed played a dominant role relative to the other participants (the decoding
606 accuracy for P14 was slightly greater than group-level decoding based on P1 to P13). However, when
607 P14 was excluded from the group, decoding accuracies remained high (0.94) and this accuracy was
608 substantially greater than individual results for P1 to P13. Therefore, although P14 played a dominant
609 role, the group-level advantage persisted when this participant was excluded. Interestingly the decoding
610 accuracy for the entire group (all 14 participants) always exceeded every single individual-level
611 decoding accuracy, and every group-minus-one decoding accuracy, which indicates that every single
612 participant, including the poorest decoded (P1) beneficially contributed to the group-level decision.

613

614



— Significant ($p < 0.05$) difference between connected model pair (paired t-test, FDR corrected)

615
616
617
618 **Figure 6. Multimodal model integration improves decoding of superior temporal and inferior**
619 **frontal regions.** Mean \pm SEM decoding accuracies across 14 participants derived using the text-based
620 and experiential model independently, and then when combined together (i.e. multimodal decoding, see
621 **Figure 1 stage 4**).

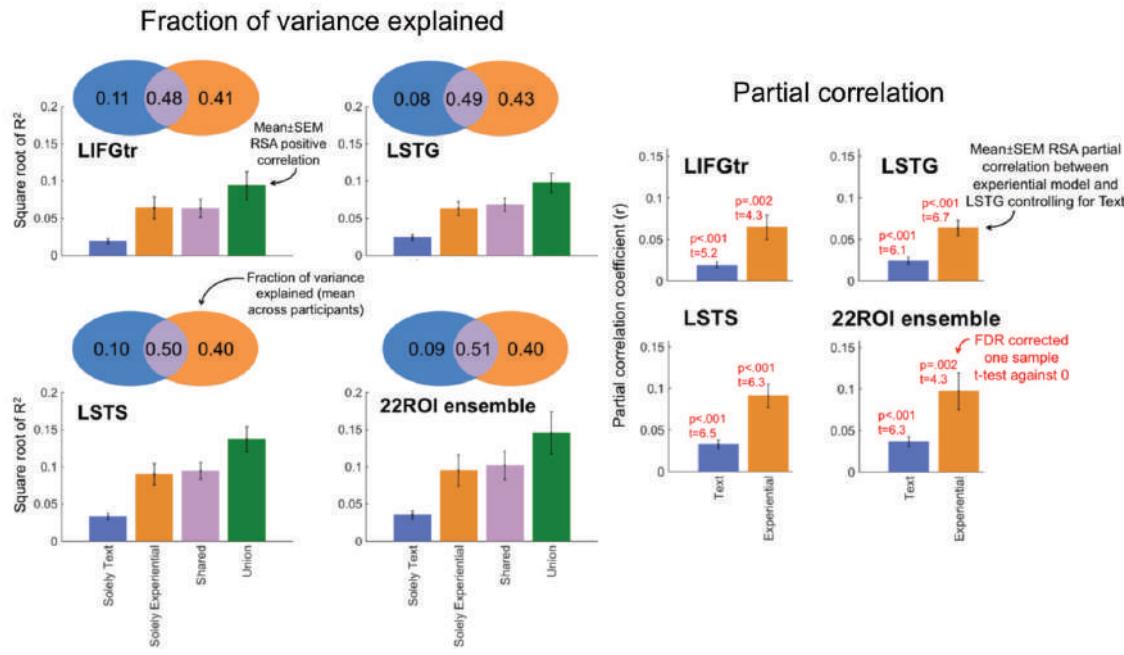
622
623
624
625 **Model integration improves decoding of superior temporal and inferior frontal brain regions**
626 Thus far analyses have been based on decoding the semantic network of all 22 ROIs as an ensemble.
627 Mean \pm SEM decoding accuracies across participants for individual ROIs are in **Figure 6**. As displayed
628 on **Figure 6** decoding accuracies arising from the multimodal approach were significantly greater than

629 those for either model in isolation in the left superior temporal sulcus (LSTS), left superior temporal
630 gyrus (LSTG) and the triangular part of the left inferior frontal gyrus (LIFGtr), and not reduced for any
631 ROI. LSTS and LSTG yielded the highest and second highest decoding accuracies across all ROIs. All
632 ROIs were on average decoded better than chance by both models. However, no statistically significant
633 differences in overall decoding accuracy were detected between models for any ROI. Otherwise
634 decoding accuracies varied across ROIs in a similar fashion to that observed in previous results
635 (Anderson et al. 2016a, Anderson et al. 2018) and were strongest in LSTS with a mean accuracy of
636 0.75 (permutation tests revealed all individuals were decoded at accuracies significantly above chance
637 level ($p=0.05$). The lowest average decoding accuracy was for the right angular gyrus (RAG),
638 mean=0.59, with 7/14 participants returning results that were significantly above chance ($p=0.05$).
639

640 **Post hoc tests partitioning the variance in the neural similarity structure explained by each
641 model**

642 The fraction of variance in the overall representational similarity structure (derived from all 240
643 sentences) of LSTS, LSTG, LIFGtr and the 22 ROI ensemble that was uniquely explained by each
644 model or commonly explained by both models is in the Venn diagrams of **Figure 7**. In each of the four
645 tests, approximately 50% of the captured variance in neural similarity structure could be explained
646 equally by either model (within participant percentage averaged across participants). Around 40% of
647 the remaining variance was associated with the experiential model and the other 10% associated with
648 the text-based model. Also displayed in **Figure 7 (right)** are partial correlation coefficients, reflecting
649 the correlation between the neural sentence-level similarity structure, and the text-based model, whilst
650 controlling for the experiential model (and vice versa). For all ROIs, partial correlation coefficients for
651 both the text-based and experiential models were found to be significantly greater than zero (across
652 participants, test statistics are in **Figure 7**). This provides further evidence that both models
653 independently contributed to explaining the neural data.

654
655



656

657 **Figure 7. Partitioning the contribution made by the text-based and experiential models to**
 658 **explaining neural similarity structure across the entire set of 240 sentences.** The Venn diagrams
 659 (left) illustrate the mean (across participants) fraction of variance that is solely accounted for by the
 660 individual models and shared between them in the RSA analyses (see **Figure 3**). The bar plots (left)
 661 display the associated mean \pm SEM positive correlations (square root of R^2). Deserving of additional
 662 explanation, in LIFGtr the mean experiential coefficient is marginally greater than the shared coefficient
 663 whereas the mean fraction of variance explained by the experiential model is less than the shared
 664 component. This occurred because the experiential model tended to uniquely explain more variance in
 665 participants with large Union R^2 values (relative to shared variance) and vice versa. The averages of
 666 raw coefficients (in the bar plots) reflect this trend, but the Venn diagrams do not, because the trend
 667 was removed by computing fractions within each participant, prior to averaging across participants.
 668 Mean \pm SEM partial correlation coefficients for the two models in the same RSA analyses are displayed
 669 in the four bar plots to the right (and tested against zero).

670

671 In both of the post hoc RSA analyses the contribution of the experiential model is on face value
 672 stronger and the text-based model weaker than would have been anticipated from the previous
 673 decoding analyses in which accuracies arising from each individual model were fairly well balanced

674 (Figure 6). We are currently unsure of the precise reason for this. The decoding analysis differs from
675 the current RSA in two key respects. First it repeatedly tests different sub-samples of the
676 representational similarity matrix that correspond to particular sentence pairs. Specifically, each cross-
677 validation iteration is based on a test of $238 * 2$ correlation coefficients rather than the global set of all
678 28680 coefficients as tested by the RSA. Second the decoding analysis incorporates a decision
679 function to best match up model sentences with fMRI sentences. Thus, it seems that one, other or both
680 of these differences render the decoding analysis less sensitive to emphasizing the contribution of the
681 experiential model. We leave detailed investigation of these differences over to future work. However,
682 in the meantime the current post hoc analyses underlines the value of the experiential model in
683 explaining the neural data.

684

685 **Cross-participant neural decoding estimates an upper bound on decoding accuracy**

686 Having demonstrated the benefits of combining models, we questioned how much room there is left
687 over for improvement in decoding. We asserted that in the general case the best decoder of an
688 individual's fMRI activation will be other people's fMRI activation (at least in the absence of a
689 personalized semantic model). To this end we decoded each individual's fMRI data using all other
690 individuals in turn and then combined the collective decoding decisions together as an estimate for the
691 upper bound on decoding accuracy. In advance, such cross-participant decoding is liable to be
692 advantaged over the semantic models by additionally decoding non-semantic information – e.g.
693 activation reflecting orthographic or syntactic processing. Nevertheless, this information is still useful to
694 identify ROIs for which there is room for improvement in decoding.

695

696 As illustrated in **Figure 2 bottom**, we first used each participant to decode each other participant
697 (repeated for each ROI, and the combination of 22 ROIs). For each individual, this left 13 decision
698 matrices (2^*2) at each cross-validation iteration (for each of the 13 other participants). These 13
699 matrices were pointwise averaged, scored as previously by comparing the sums on the matrix diagonal
700 and antidiagonal, and then scores were averaged across all cross-validation iterations to give a final

701 cross-participant decoding accuracy. We considered this final accuracy as an estimate on the upper
702 bound decoding achievable for that individual (in absence of a personalized semantic model).

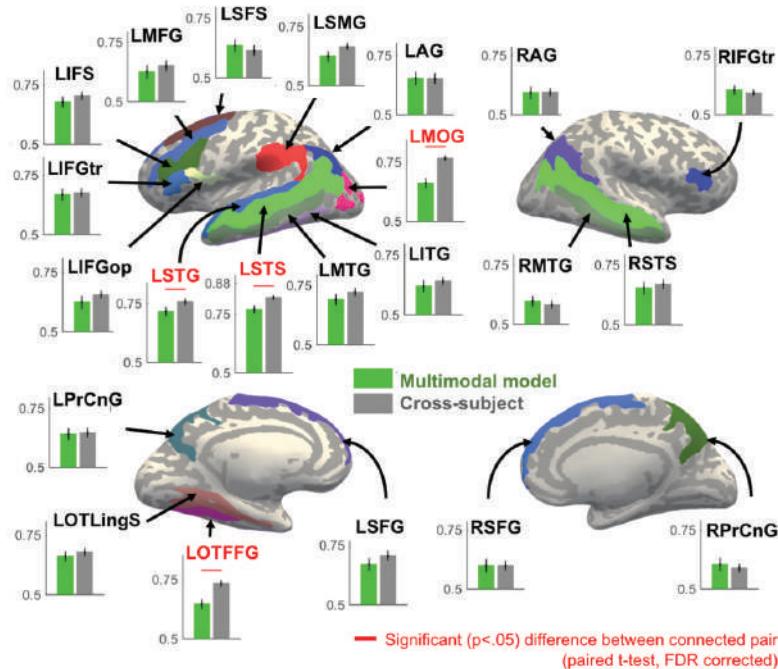
703

704 Mean+/-SEM cross participant decoding accuracies for each ROI are compared to the multimodal
705 decoding accuracies in **Figure 8**. Only multimodal decoding accuracies (i.e. the best decoder so far)
706 are displayed as a comparison to avoid visual clutter (the isolated model results are in **Figure 6**). FDR
707 corrected paired t-tests between cross participant and multimodal model-based accuracies revealed
708 significantly stronger decoding for the cross-participant approach in LSTS, LSTG, left occipitotemporal
709 fusiform gyrus (LOTFFG) and left mid occipital gyrus (LMOG) and for the set of 22 ROIs decoded as an
710 ensemble (all $p < 0.05$). This indicated that there was sentence-related signal present in the fMRI data
711 within these regions that had not been decoded by the models. This was particularly the case for
712 LMOG and LOTFFG ($d = 1.67$ and 1.24 respectively) however improvements for LSTS and LSTG were
713 notable ($d = 0.81$ and 0.56 respectively).

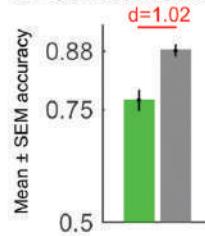
714

715 Also noteworthy was the gap separating cross participant decoding accuracies for the 22ROI ensemble
716 and individual ROIs. For instance, the difference between cross participant decoding of the highest
717 accuracy ROI (LSTS) and the 22ROI ensemble was significant and sizeable ($t = 7.7$, $p < 0.0001$, $d = 1.47$).
718 This provides evidence that complementary sentence-related information (which could be semantic or
719 orthographic or syntactic) was distributed across the different ROIs. In contrast for the semantic model-
720 based analysis most of the decodable information appears to be in LSTS (see also **Figures 9-11**).

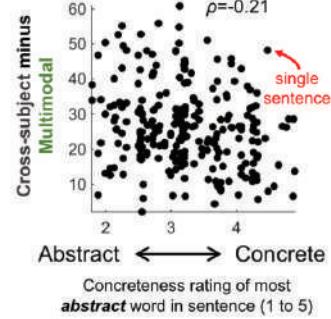
Note: Cross-subject (brain activation-based) decoding is liable to decode semantic and non-semantic (e.g. orthographic/syntactic) elements of neural activation



22 ROIs decoded as an ensemble



Cross-subject approach particularly improves on decoding of sentences containing **abstract** words



721

722

Figure 8. Estimating the room for improvement: how cross-participant decoding improves on the multimodal model-based approach. Mean \pm SEM cross participant (brain-based) decoding accuracies (**Figure 2 bottom**) across all 14 participants beside comparative results for the multimodal model (also shown in **Figure 4**). Detailed results arising from decoding using the combination of all 22 ROIs (**Figure 2 top**) are to the right. The scatter plots indicate characteristics of sentences for which cross-participant (brain-based) decoding was advantaged over the multimodal model-based approach. The effect size (d) was estimated as described in **Figure 4**.

730

731

Sentences decoded weakly by the models tended to contain abstract words

To attempt to get a handle on what semantic information could have been left undecoded by the models, and in so doing identify model weaknesses we finally tested whether particular types of sentences were better decoded by the cross-participant approach. We concentrated analyses on decoding all 22 ROIs as an ensemble. Then for each participant (14) we extracted sentence-wise decoding scores (240) for both the cross-participant and multimodal approaches. The vector of 240

738 multimodal model-based decoding scores was pointwise subtracted from the cross-participant decoding
739 scores to generate an “accuracy boost vector” indicating which sentences were better decoded by
740 cross participant decoding. For each participant we correlated (Spearman) this accuracy boost vector
741 with min/mean/max concreteness of constituent content words in sentences, the min/mean/max (log2
742 transformed) word frequency per sentence, the minimum/mean/maximum word length per sentence
743 and the number of words per sentence. The resulting correlation coefficients were Fisher’s r to z
744 transformed (arctanh). For each set of 10 tests, transformed coefficients for all 14 participants were
745 compared to zero using a t-test. P-values were then FDR corrected. Four tests yielded significant
746 results. These were the concreteness rating of the most abstract word in the sentence (mean $r=-0.08$,
747 $t=-4.1$, $p=0.01$ FDR corrected, see also **Figure 8**), mean word frequency (mean $r=-0.06$, $t=3.6$, $p=0.03$
748 FDR corrected), minimum word length (mean $r=-0.05$, $t=-3.5$, $p=0.03$ FDR corrected) and the number
749 of words in the sentence (mean $r=0.07$, $t=4.5$, $p=0.01$, FDR corrected). To counteract possible effects
750 of spurious intercorrelations between these four sentence measures, a second round of partial
751 correlation analyses was run. Each participant’s accuracy boost vectors were partially correlated
752 (Spearman) with vectors associated with each of the four measures in turn whilst controlling for the
753 other three. Of the four measures, only correlations with the concreteness rating of the most abstract
754 word in the sentence were found to be significantly lower than zero (mean partial $r=-0.06$, $t=-3.6$,
755 $p=.003$). This provided evidence that cross participant decoding gained a particular advantage over the
756 multimodal model approach for sentences containing abstract words. The relationship between the
757 decoding advantage and the number of words per sentence, constituent word frequencies and word
758 lengths is unclear due to their intercorrelation.

759

760 **Secondary supporting analyses**

761 The following two sections present the outcome of a series of secondary supporting analyses that are
762 possibly best suited to the dedicated reader.

763

764 **Supporting analysis: multimodal decoding advantage is preserved in analyses of different**
765 **cortical networks**

766 Our main analysis was focused on a semantic network of 22 regions that had been identified in
767 Anderson et al. (2018) which differently used a two-stage regression-based analysis and the
768 experiential model alone. Reasons for this were to maintain continuity, so as to enable current results to
769 be referenced back to the results of Anderson et al. (2018), and also for simplicity. However, because
770 the network of 22 regions was derived using the experiential model only it is possible that the current
771 analysis could have been biased towards the experiential model.

772 To confirm that results were not specific to the 22ROI network, we repeated our initial decoding
773 analysis using different network configurations. We first derived decoding accuracies for all 150 ROIs in
774 the Destrieux atlas using the text-based and experiential models independently. At this stage the
775 analysis was conducted without voxel selection to cut down on computational overheads. We then
776 identified a set of “high accuracy” brain regions that were decoded significantly at the p=0.01 level
777 (FDR corrected) by either model. We then reperformed the decoding analysis, this time with voxel
778 selection (50 voxels per ROI). We firstly decoded the *union* of high accuracy ROIs associated with
779 either or both of the models. Secondly, we decoded the *intersection* of high accuracy ROIs associated
780 with both models.

781 The “intersection” network contained 12 ROIs. Of these, the following 9/12 ROIs also appeared in the
782 22 ROI ensemble: LSTS, LSTG, LMTG, LIFGtr, LIFS, LAG, LPrChG, LOTLingS, RSTS. Destrieux atlas
783 names for the 3 ROIs that were not in the 22 were: ctx_lh_G_cingul-Post-dorsal, ctx_lh_S_oc-temp_lat,
784 ctx_lh_S_subparietal. The “union” network contained 15 ROIs. These included the 12 listed above as
785 well as LOTFFG and LMOG which were in the original 22 ROI ensemble and ctx_lh_S_precentral-inf-
786 part which was not. The 4 new ROIs that had not been in the original 22 were all relatively low scoring
787 (ranked 11th highest or below out of the union of 15 ROIs). Mean±SEM decoding accuracies for the
788 text-based, experiential and multimodal for the 4 new ROIs were respectively: ctx_lh_G_cingul-Post-
789 dorsal 0.63±0.02, 63±0.02, 64±0.02; ctx_lh_S_oc-temp_lat 0.63±0.02, 63±0.01, 65±0.02;

790 ctx_lh_S_subparietal: 0.62 ± 0.02 , 63 ± 0.02 , 64 ± 0.02 and ctx_lh_S_precentral-inf-part 0.62 ± 0.02 ,
791 61 ± 0.02 , 63 ± 0.02 . T-tests revealed no significant differences in decoding accuracy between the text-
792 based and experiential models for these 4 ROIs.

793 Decoding accuracies arising from the intersection and union networks are displayed in **Figure 9**. For
794 comparison we display results of the original 22 ROI ensemble analysis conducted on either 50 or 100
795 voxels selected per region. Additionally, we display results for LSTS (the highest scoring ROI) and also
796 results when the analysis was undertaken on the entire cortex without any voxel selection. Also
797 displayed are correlations relating the contribution of the text-based model to multimodal decoding to
798 ratings of sentence abstractness echoing **Figure 3**.

799 **Figure 9** reveals that the decoding advantage brought through model integration is preserved across all
800 tests irrespective of the network configuration. Also preserved is the “abstractness” advantage brought
801 by the text-based model to decoding sentences containing abstract words. It is visually apparent that
802 that decoding accuracy was modulated by the network configuration tested. Because it is not a focus of
803 the current article we leave in depth treatment of these differences and how to select the optimal
804 network over to future work. Suffice it to say that activity in LSTS appears to have been the linchpin of
805 network-based decoding and there was no dominant network configuration that yielded significantly
806 greater accuracy than all others.

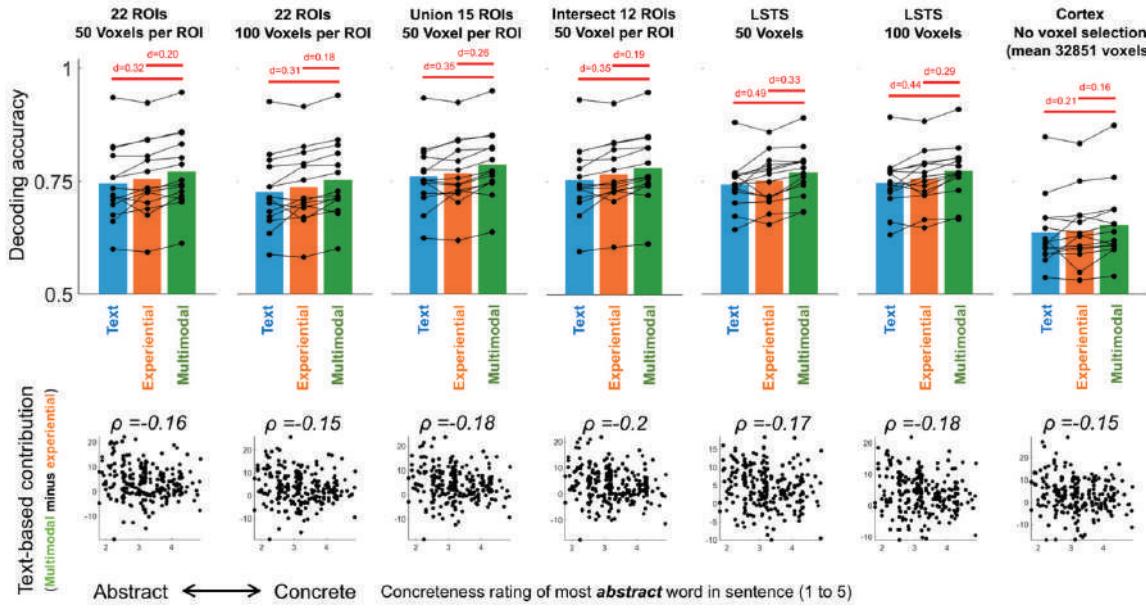


Figure 9. Decoding accuracies arising from decoding different networks of ROIs using different model combinations. This Figure is companion to **Figure 4** which describes how effect sizes (d) were estimated.

Supporting analysis: comparison and integration of ridge regression with similarity-based decoding

Prompted by a concern that the current RSA approach might inherently bias results in favor of one or other model we used ridge regression (Hoerl and Kennard, 1970) to reanalyze multimodal ROIs: LSTS, LSTG, LIFGtr and the 22ROI ensemble. Ridge regression was selected because it has been popular in recent fMRI studies using text-based models (e.g. Huth et al. 2016; de Heer et al. 2017; Pereira et al. 2018). Importantly, this reanalysis provides a quantitative estimate of the decoding advantages brought by using experiential attributes and similarity-based methods comparative to a state-of-the-art text-based/ridge regression approach.

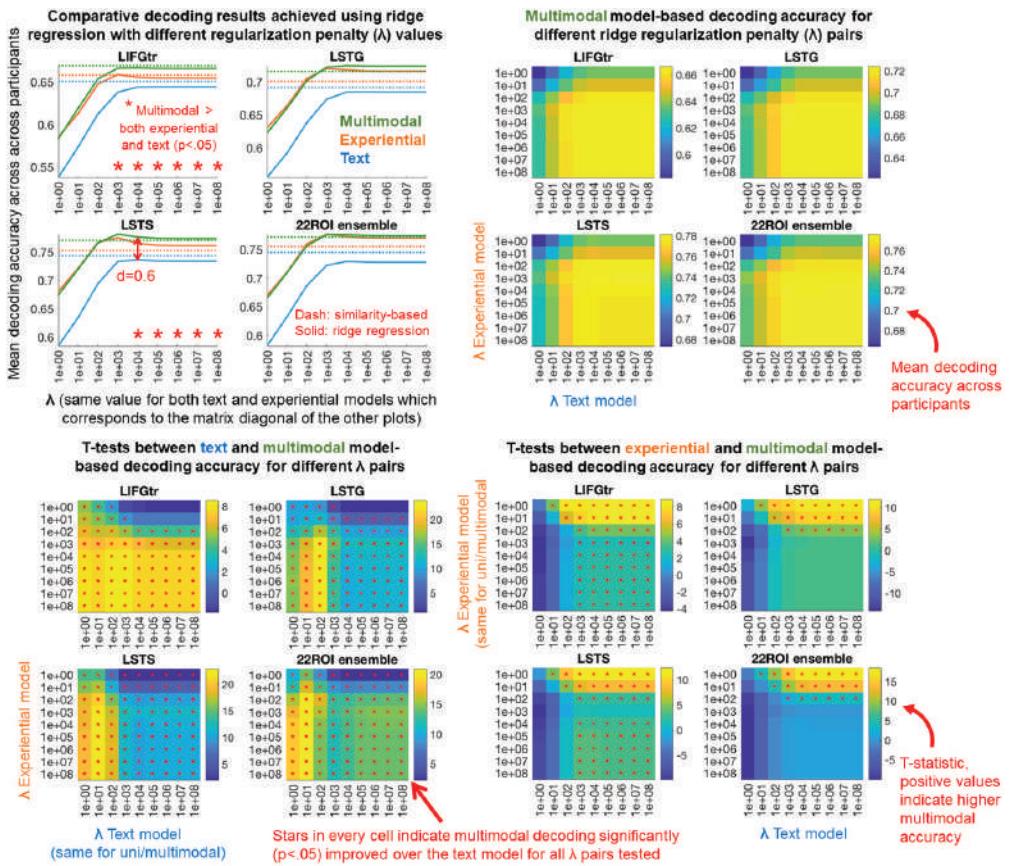
For each participant and each ROI we reran precisely the same leave-2-sentence-out cross validation procedure with the same training/testing data splits and same voxel selection (50 voxels per ROI) as

824 our main similarity-based analysis. At each of the 28 680 cross validation iterations, we fit a separate
825 ridge regression for the text-based model and then for the experiential model to predict activation in
826 each individual voxel (i.e. forward encoding). We repeated regression fitting using each of the following
827 9 regularization penalties (λ) [1 10 100 10^3 10^4 10^5 10^6 10^7 10^8]. As before, at each cross validation
828 iteration, we computed a 2*2 decoding decision matrix for both the text-based and experiential models.
829 This was repeated for each λ by correlating predicted and actual fMRI activation patterns and r-to-z
830 transforming correlation coefficients (9 decoding decision matrices per model). To create a multimodal
831 2*2 decoding decision matrix we pointwise averaged together decision matrices arising from each of
832 the two models precisely as described in **Figure 1 stage 4**. We repeated this for every combination of λ
833 pairs leaving $9*9=81$ multimodal decision matrices per iteration per ROI. We created a 22 ROI
834 ensemble decoding matrix for each model and λ combination by pointwise averaging each of the $9 + 9$
835 + 81 decision matrices across ROIs. Note that for the multimodal approach there were 81^{22} possible
836 ways that λ could have been combined across ROIs and therefore the pointwise averaging approach
837 we have taken could have missed out on the ideal combination. Alternative approaches could have
838 included selecting an optimal λ for each ROI using nested cross-validation, or building a single decision
839 matrix from all voxels. Because regression is not our prime focus we leave detailed investigation of this
840 to future work. As for all our other analyses, at each iteration decoding decision matrices were
841 evaluated as correct (1) if the sum of coefficients on the diagonal exceeded the sum on the
842 antidiagonal, otherwise they were incorrect (0). A final decoding accuracy score was assigned as the
843 mean correctness across trials for each model and λ combination (yielding $9 + 9 + 81$ accuracies per
844 ROI per participant).

845
846 Decoding accuracies were in general qualitatively similar to those observed for the similarity-based
847 analysis and are illustrated in **Figure 10**. Decoding accuracies for both models tended to reach a
848 maxima and flatten off at λ values greater than $\lambda=10^3$ or 10^4 . For “flattened” λ values, in tests on LSTS
849 and LIFGtr the multimodal approach was significantly more accurate than either of the text-based and
850 experiential models alone. For LSTG and the 22ROI ensemble, whilst the multimodal approach yielded

851 significantly greater decoding accuracies than the text-based model (**Figure 10 bottom left**),
852 accuracies were not significantly greater than the experiential model (at least for “flattened” λ values
853 **Figure 10 bottom right**). This weak multimodal improvement reflects both the comparatively weak
854 decoding accuracies achieved for the text-based model using ridge regression, and the relatively strong
855 accuracies for the experiential model (see caption of **Figure 10** for statistical test results). Ridge
856 regression’s weak performance with the text-based model could reflect a combination of difficulties
857 surrounding: scaling up to the 300 text-based features (in comparison to 65 experiential attributes, with
858 240 sentences); and/or the distributional properties of the text-based data; and/or our current selection
859 of λ missing out on the optimal value. Ridge regression’s stronger performance with the experiential
860 model evidences that parameter fitting can lead to decoding improvements with the current data. There
861 was little to separate multimodal decoding accuracies for the two decoders, although for a particular λ
862 value (10^3) ridge regression yielded stronger performance in LSTS alone.

863



864 **Figure 10. Comparative text-based, experiential and multimodal decoding accuracies acquired**
865 **using ridge regression** *Multimodal advantages* for a particular selection of λ values are in the top left.
866 Multimodal decoding accuracies for all λ configurations are in the top right. Results of paired t-tests
867 comparing multimodal decoding accuracies to text-based decoding accuracies and experiential
868 decoding accuracies for each λ configuration are in the bottom left and right plots. All tests were one-
869 tailed, in anticipation of the multimodal advantage observed in our initial analyses. The illustrated effect
870 size (d) provides a conservative estimate of the benefit of integrating experiential features into a
871 conventional text-based ridge regression approach. d was computed as described in the caption of
872 **Figure 4. Differences between ridge regression and similarity-based decoding:** The top left plot
873 illustrates both similarity-based (dashed lines) and regression-based (solid lines) results. Decoding
874 accuracies using ridge regression with the text model and the top performing λ (always $\lambda=10^4$) were
875 unanimously significantly lower than for the similarity-based approach (LSTS: $t=5.8$, $p=6*10^{-5}$, $df=13$;
876 LSTG: $t=5.3$, $p=1.4*10^{-4}$, $df=13$; LIFGtr: $t=5.2$, $p=1.8*10^{-4}$, $df=13$; 22ROI: $t=8$, $p=2*10^{-6}$, $df=13$; all 2-
877 tailed paired t-tests, $df=13$). Conversely, in 75% of tests using the experiential model with the top

879 scoring λ (always $\lambda=10^3$), ridge regression yielded significantly stronger decoding accuracies than the
880 similarity-based analysis (LSTS: $t=5.6$, $p=8.5*10^{-5}$; LSTG: $t=4$, $p=0.001$; LIFGtr: $t=0.1$, $p=0.9$; 22ROI:
881 $t=7$, $p=9.8*10^{-6}$; all 2-tailed paired t-tests, $df=13$). For multimodal decoding, ridge regression yielded
882 stronger decoding in LSTS with the top scoring $\lambda=10^3$ ($t=3.3$, $p=0.006$) but not other λ values, and
883 otherwise there were no significant differences for the other ROIs (LSTG: $t=1.9$, $p=0.07$; LIFGtr: $t=-.88$,
884 $p=0.4$; 22ROI: $t=1.4$, $p=0.2$; all 2-tailed paired t-tests, $df=13$). All p-values in plots and captions are not
885 corrected for multiple comparisons.

886

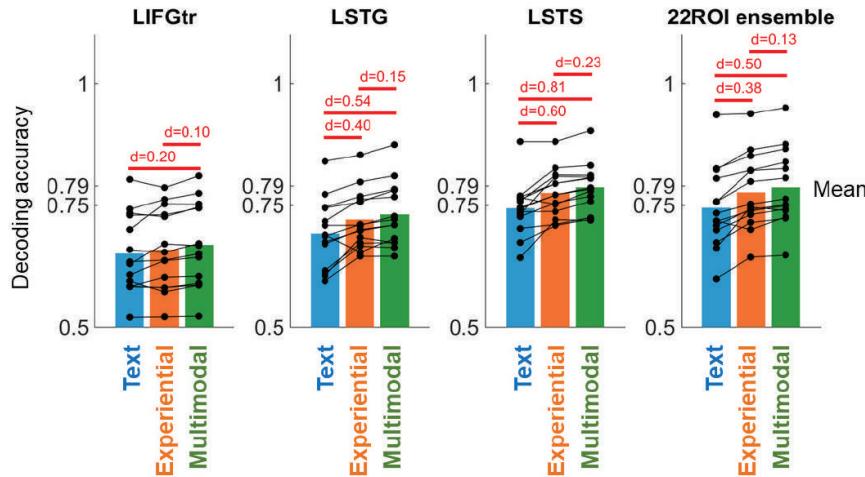
887 The comparatively weak decoding accuracies obtained with ridge regression for the text-based model
888 begged the question of whether the ridge regression multimodal decoder was failing to capitalize on
889 neural information that was decoded successfully by the similarity-based approach. Relatedly, whether
890 decoding of LSTG and the 22ROI ensemble would be advantaged by a “best of both worlds”
891 multimodal approach that jointly leverages similarity-based decoding and ridge regression. To answer
892 this question, we reran the cross-validation analysis using the similarity-based approach with the text-
893 based model in parallel with ridge regression on the experiential model (using top scoring $\lambda=10^3$). At
894 each cross validation iteration we integrated the respective decoding decisions made by the different
895 models/decoders by pointwise summing respective $2*2$ decoding decision matrices (**as in Figure 1**
896 **stage 4**). The multimodal joint similarity/regression approach indeed yielded significantly greater
897 decoding accuracies than both the regression-based experiential decoder and the text similarity-based
898 decoder in LSTG and the 22ROI ensemble as well as LSTS and LIFGtr (see **Figure 11**, statistical test
899 results are in the caption).

900

901 In sum this section has provided further evidence that experiential semantic features explain variance in
902 sentence-level fMRI data that cannot be accounted for by state-of-the-art text-based regression
903 approaches, and further support for the claim that multimodal approaches provide the most accurate
904 models of fMRI to date.

905

Text-based decoding using similarity approach
Experiential attribute decoding using ridge regression
Multimodal decoding using a fusion of the above



906
907 **Figure 11. Multimodal decoder integrating the best text-based decoder (similarity) with the best**
908 **experiential decoder (ridge regression, $\lambda=10^3$)** Effect sizes (d) are displayed in cases of statistically
909 significant differences (paired t-test, all $p \leq .01$, FDR corrected). d was computed as described in the
910 caption of **Figure 4**. Paired t-test results were: for contrasts between multimodal decoding and
911 experiential regression-based decoding: LSTS: $t=4.8$, $p=0.002$; LSTG: $t=3.6$, $p=0.01$; LIFGtr: $t=4.5$,
912 $p=0.003$; 22ROI: $t=5$, $p=0.002$; for contrasts between multimodal decoding and text similarity-based
913 decoding: LSTS: $t=8.3$, $p<10^{-4}$; LSTG: $t=8.0$, $p<10^{-4}$; LIFGtr: $t=4.0$, $p=0.006$; 22ROI: $t=5$, $p<10^{-4}$; for
914 contrasts between experiential regression-based decoding and text similarity-based decoding: LSTS:
915 $t=4.2$, $p<0.005$; LSTG: $t=3.8$, $p<0.009$; LIFGtr: $t=1.39$, $p=0.58$; 22ROI: $t=4.8$, $p=0.002$. All p-values FDR
916 corrected.
917

918 **Supporting analysis: persistence of multimodal advantage using different text-based models**

919 The claim that the experiential model enhanced decoding by capturing non-linguistic experiential
920 knowledge rests on the assumption that the current text-based model captured all of the experiential
921 structure that is possible to obtain from word use statistics (see **Methods**). As both the text-based and

922 experiential modeling approaches are in an ongoing state of development it cannot be concluded that
923 the current results will be the same for all future models and/or neural data sets. Additionally, it is
924 possible that idiosyncrasies surrounding how the experiential model was constructed and its statistical
925 properties (e.g. representational sparseness) could have contributed to the decoding advantage it
926 conferred. Vice versa for the text-based model. Consequently, we consider that the current results
927 provide “early evidence” that linguistic and non-linguistically acquired knowledge is represented in fMRI
928 activation elicited in sentence comprehension. However, our core finding that an integrated experiential
929 and text-based decoding approach yields significantly higher accuracy than either model alone has held
930 true for all text-based models we have tested thus far, which have been built using different algorithms
931 from different text corpora. We are not aware of another experiential model suitable for the current
932 analysis.

933

934 In preliminary investigations we had tested word co-occurrence models (e.g. Roller et al. 2014),
935 word2vec (Mikolov et al. 2013, Baroni et al. 2014), that yielded similar decoding accuracy levels (e.g.
936 LSTS: mean ≤ 0.75) and a similar core result that integrating experiential and text-based models
937 yielded significantly greater decoding accuracies (e.g. LSTS: Multimodal > cooccurrence $t=10.1$, $p<10^{-5}$,
938 $d=0.78$; Multimodal>word2vec $t=6.9$, $p<10^{-5}$, $d=.44$, both 2-tailed paired t-tests). We focused on GloVe
939 (Pennington et al. 2014) principally because it was the basis for Pereira et al.’s (2018) fMRI sentence
940 decoding study.

941

942 In the interim a number of new computational text-based approaches have emerged (e.g. Conneau et
943 al. 2017, Peters et al. 2018, Subramanian et al. 2018, Devlin et al. 2019). These typically leverage deep
944 artificial neural networks to derive sentence representations that reflect word order and within-sentence
945 contexts. A thorough comparison of deep network text-based approaches and the experiential model is
946 beyond the scope of the current work, and perhaps would be best undertaken using an experiential
947 model that also accommodates word order and context effects (which could be achieved by rating
948 entire sentences, words in context, or entering experiential word vectors as deep network input).

949 Nevertheless, to date we have tested one deep model, InferSent (Conneau et al. 2017), which is
950 notable in having recently yielded state-of-the-art decoding of Pereira et al.'s (2018) sentence-level
951 fMRI dataset (Sun et al. 2019). We hope to present a full treatment of results in future work, however to
952 foreshadow those, the current core finding still holds: Integrating the current sentence-level experiential
953 model with InferSent yields significantly greater decoding accuracy than either model in isolation (e.g.
954 LSTS: Multimodal > InferSent $t=7.2$, $p<10^{-5}$, $d=0.34$, 2-tailed paired t-test).

955

956 **Discussion 1502/1500 words**

957 This study has revealed early evidence that modeling both linguistic knowledge of word usage and
958 experiential knowledge of words' referents enhances decoding of brain activation patterns associated
959 with sentence meaning. This suggests that non-linguistic experiential knowledge is represented in
960 sentence-level fMRI activation. Importantly, because this result is based on direct measures of brain
961 activation elicited during the comprehension of natural sentences, it is an advance on previous
962 behavioral evidence that has been indirectly inferred from experimental responses (e.g. Paivio, 1971,
963 Stanfield and Zwaan, 2001 Zwaan et al. 2002; Andrews et al. 2009, Louwerse and Jeuniaux, 2010,
964 Kousta et al. 2011). More generally this is evidence that it is now possible to use brain data to
965 quantitatively estimate the contribution that linguistic and non-linguistically acquired knowledge make to
966 representing the meaning of natural language. This is especially relevant to theories that conceptual
967 representations are acquired through and partially embodied within experiential neural systems
968 (Barsalou et al. 2008, Glenberg, 2010, Pulvermüller 2013, Binder et al. 2016). Relatedly, it suggests we
969 can begin to estimate how close "ungrounded" semantic models (e.g. text-based) can get to
970 representing human conceptual knowledge (see Harnad, 1990 for discussion of the "symbol grounding
971 problem"). With respect to questions of grounding and embodiment we should be clear that the current
972 analyses provide no guarantee that brain activation that was selectively decoded by the experiential
973 model was actually represented within primary perceptual/modal processing systems, or was critical to
974 comprehension rather than epiphenomenal (Mahon and Caramazza 2009, Mahon 2015).

975

976 At a more practical level, the current study advances on previous state-of-the-art text-based neural
977 encoders/decoders (Huth et al. 2016, Huth et al. 2016b, Pereira et al. 2018) because multimodal
978 integration boosts decoding performance. This provides further evidence that combining multiple
979 modalities of information in semantic models leads to more human-like representations of meaning
980 (Andrews et al. 2009, Bruni et al. 2014, Anderson et al. 2015).

981

982 We have also demonstrated that the text-based model contributed particularly to decoding sentences
983 containing abstract words. Although this was hypothesized (Anderson et al. 2016a) and text-based
984 models have previously helped explain abstract concept fMRI (Anderson et al. 2017, Wang et al. 2018,
985 Pereira et al. 2018), it was not a foregone conclusion. This was because abstract concepts are thought
986 to be grounded relatively strongly upon affective experiences (Kousta et al. 2011, Vigliocco et al. 2014)
987 and contemporary text-based models generate relatively weak predictions of affective experiential
988 attributes (Utsumi, 2018). As it turned out many sentences that the text-based model helped explain
989 contained words with affective connotations (e.g. “happy”, “celebrated”, “survived”). It seems likely that
990 the advantage conferred in these cases (and others) was down to the extra linguistic/contextual
991 information in the text-based model. Otherwise, Utsumi (2018) found text-based models to be
992 disadvantaged in predicting spatial/temporal attributes. Testing how spatial/temporal attributes
993 contribute to semantic representations in the brain may thus provide an interesting avenue for future
994 investigation. Whilst we here detected tentative evidence that the experiential model contributed to
995 decoding concrete sentences (without any abstract words), this result did not survive correction for
996 multiple comparisons.

997

998 The demonstration that sentence-level neural activation is best decoded using a multimodal approach
999 is not without foreshadow. Anderson et al. (2015) found that a visually grounded semantic model
1.000 (derived from natural images) and a text-based model differentially correlated with fMRI activation in
1.001 brain regions with known visual/linguistic processing roles. However, because participants were tasked
1.002 to read concrete nouns and actively contemplate their semantic properties (Just et al. 2010) it is not

.003 clear whether the results reflect active visual imagery as opposed to more passive language
.004 comprehension (see also Willems et al. 2010). In other work, Abnar et al. (2018) used a joint text-
.005 based/experiential approach to better predict fMRI elicited by drawings of nouns alongside their names,
.006 and Wang et al. (2018) revealed partial correlations between fMRI elicited by abstract Chinese words, a
.007 Chinese text-based model and a model built from 12 behavioral ratings that interestingly included
.008 valence, space and time. In both cases it is not clear how results would extend to decoding read
.009 sentences (in English).

.010

.011 The current study extended a representational similarity-based decoding method (Anderson et al.
.012 2016b, Anderson et al. 2017) to the neural decoding of sentences using parallelized combinations of
.013 multiple models, brain regions and participants. Combination of multiple participants' neural data was
.014 achieved by "ensemble averaging" of decoding decisions. This sets the similarity-decoding method
.015 apart from "hyper alignment" methods (Haxby et al. 2011, Guntupalli et al. 2016) that represent neural
.016 responses using a common representational space. A disadvantage of integrating decoding decisions
.017 is that this does not generate predictions of individual voxel's activity (unlike regression-based
.018 *encoding*). Similarity-based approaches can be configured to estimate voxel activity by applying the
.019 correlation coefficients comprising similarity vectors (e.g. **Figure 1 stage 3**) as weights in a weighted
.020 average of corresponding brain activation patterns (as described by Anderson et al. 2016b). However,
.021 we leave a comparative investigation of this over to future work. We did not run the current analysis
.022 using similarity-based encoding in part to avoid the additional data normalization step that would have
.023 been required to combine data (see **Figure 1** caption). For the *decoding* case at hand the similarity-
.024 based approach performs competitively with ridge regression (better for the text-based model and
.025 worse for the experiential model, **Figure 10**) whilst cutting out overheads associated with repeating the
.026 analyses with different regularization penalties, and picking the appropriate one.

.027

.028 Cross participant neural decoding was introduced as a method to estimate an upper bound on
.029 decoding accuracy achievable with (group-level) models. This followed the reasoning that on average

.030 the most accurate neural decoder will be based on neural data. Indeed, for LSTS, LSTG, LOTFFG, and
.031 LMOG decoding accuracy was significantly greater for the cross-participant approach, and there were
.032 no ROIs for which accuracy was significantly worse. Practically speaking this result was however not
.033 guaranteed. Had the fMRI data been too noisy and the group been too small, the model-based
.034 approach could have yielded stronger decoding. It is also important to recall that the upper bound
.035 estimate provided by cross participant decoding does not apply to decoding semantics per se but to
.036 decoding the entire linguistic processing stream from stimulus perception to semantic interpretation.
.037 Also, the cross-participant approach does not apply to decoding person specific aspects of semantic
.038 representation, so there may well be decodable neural signal left over that could only be revealed by
.039 personal information.

.040

.041 The upper bound decoding estimate was used to identify that weakly decoded sentences tended to
.042 contain abstract words, which suggests that the neural data contains undecoded aspects of abstract
.043 conceptual representations. This presents a challenge for future modeling to improve on models of
.044 abstract concepts. Given limitations in current understanding of abstract knowledge representation
.045 there may be an interesting opportunity to move forward here in a different way, by incorporating
.046 features of brain activation into artificial semantic models and in so doing provide a new way for
.047 neuroscience to feed back to AI (see also Fyshe et al. 2014 and Hassabis et al. 2017).

.048

.049 One limitation of the current study is the assumption that the additional neural activation decoded by
.050 the experiential model reflects semantic information that cannot be extracted from natural language
.051 data. This is not strictly guaranteed and it is possible that future text-based approaches will account for
.052 the signal decoded by the experiential model. A limitation of the current experiential approach is the
.053 assumption that experiential knowledge can be comprehensively estimated through introspective
.054 ratings of the relationship between concepts and putative neural systems. Indeed, we have revealed
.055 evidence that the text-based language model captures information the experiential model did not,
.056 however there may be other semantic features that cannot be verbally described and/or introspectively

.057 accessed, in which case models that are truly grounded in modal information (e.g. Bruni and Baroni,
.058 2014; Anderson et al. 2015) may come to the fore. Ultimately the answers to the above questions will
.059 be borne out through future work that incorporates different modalities of information into semantic
.060 models (e.g. Andrews et al. 2009, Bruni et al. 2014, Kiela and Clark 2017), and compares this to brain
.061 data (e.g. Anderson et al. 2013, Anderson et al. 2015, Anderson et al. 2017, Bulat et al. 2017). The
.062 current study has contributed methods that we hope will assist in this enterprise.

.063

.064 In conclusion, the current study has provided initial evidence that linguistic and non-linguistic
.065 experiential knowledge can be detected in sentence-level brain activation by extending a similarity-
.066 based framework to exploit respective models in fMRI decoding. It has also presented a cross
.067 participant decoding method which has demonstrated that a substantial amount of neural signal
.068 remains unexplained. This decoding gap is likely to be filled by modeling advances that take word
.069 order, syntax, morphology and polysemy into account in semantic composition and begin to
.070 accommodate pragmatic inferences and theory of mind. For the future, in all of these endeavors we
.071 contend that model-based approaches that integrate information across multiple modalities of
.072 experience will be necessary for the fullest interpretation of neural activation patterns associated with
.073 meaning.

.074

.075 **References**

.076

- .077 Abnar S, Ahmed R, Mijnheer M, Zuidema W. (2018). Experiential, Distributional and Dependency-
.078 based Word Embeddings have Complementary Roles in Decoding Brain Activity. Proceedings of the
.079 8th Workshop on Cognitive Modeling and Computational Linguistics (CMCL 2018), Salt Lake City,
.080 Utah, USA: Association for Computational Linguistics. pp 57–66.
.081
.082 Andrews M, Vigliocco G, Vinson D. (2009). Integrating experiential and distributional data to learn
.083 semantic representations. *Psychol. Rev.* 116 (3), 463–498.
.084 Andrews M, Frank S, Vigliocco G. (2014). Reconciling embodied and distributional accounts of
.085 meaning in language. *Top Cogn Sci.* 6:359–370.
.086 Anderson AJ, Bruni E, Bordignon U, Poesio M, Baroni M. (2013). Of words, eyes and brains:
.087 Correlating image-based distributional semantic models with neural representations of concepts.
.088 Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP 2013);
.089 Seattle, WA: Association for Computational Linguistics. pp. 1960–1970.

- .090 Anderson AJ, Bruni E, Lopopolo A, Poesio M, Baroni M. (2015). Reading visually embodied meaning
.091 from the brain: visually grounded computational models decode visual-object mental imagery induced
.092 by written text. *NeuroImage*. 120:309–322.
- .093 Anderson AJ, Zinzser BD, Raizada, RDS. (2016b). Representational similarity encoding for fMRI:
.094 pattern-based synthesis to predict brain activity using stimulus-model-similarities. *NeuroImage*. 128:44–
.095 53.
- .096 Anderson AJ, Binder JR, Fernandino L, Humphries CJ, Conant LL, Aguilar M, Wang X, Doko D,
.097 Raizada, RDS. (2016). Predicting neural activity patterns associated with sentences using a
.098 neurobiologically motivated model of semantic representation. *Cerebral Cortex*. doi:
.099 10.1093/cercor/bhw240.
- .100 Anderson AJ, Kiela D, Clark S, Poesio M. (2017). Visually Grounded and Textual Semantic Models
.101 Differentially Decode Brain Activity Associated with Concrete and Abstract Nouns. *Transactions of the
.102 Association for Computational Linguistics*. 5, 17–30.
- .103 Anderson AJ, Lalor EC, Lin F, Binder JR, Fernandino L, Raizada RDS, Grimm S, Wang X. (2018).
.104 Multiple regions of a cortical network commonly encode the meaning of words in multiple grammatical
.105 positions of read sentences. *Cerebral Cortex*. doi: 10.1093/bphy110
- .106 Anderson AJ, Lin F. (2019). How pattern information analyses of semantic brain activity elicited in
.107 language comprehension could contribute to the early identification of Alzheimer's Disease.
.108 *NeuroImage: Clinical* 22, 101788.
- .109 Barsalou LW. (1999). Perceptual symbol systems. *Behavioral and Brain Sciences*, 22(4), 637–660.
- .110 Barsalou LW, Santos A, Simmons WK, Wilson CD. (2008). Language and simulation in conceptual
.111 processing. In: De Vega M, Glenberg AM, Graesser AC, editor. *Symbols, embodiment, and meaning*.
.112 Oxford: Oxford University Press. pp. 245–283.
- .113 Binder JR, Desai RH, Graves WW, Conant LL. (2009). Where is the semantic system? A critical review
.114 and meta-analysis of 120 functional neuroimaging studies. *Cereb Cortex*. 19:2767–2796.
- .115 Binder JR, Desai RH. (2011). The neurobiology of semantic memory. *Trends Cogn Sci*. 15(11):527–
.116 536.
- .117 Binder JR, Conant LL, Humphries CJ, Fernandino L, Simons S, Aguilar M, Desai R. (2016). Toward a
.118 brain-based componential semantic representation. *Cognitive neuropsychology* 33 (3-4), 130-174.
- .119 Bulat L, Clark S, Shutova E. (2017). Speaking, Seeing, Understanding: Correlating Semantic Models
.120 with Conceptual Representation in the Brain. *Proc. Conf. Emp. Meth. Nat. Lang. Proc.* 1532–1543
.121 (Association for Computational Linguistics, Doha, Qatar, 2014).
- .122 Bruffaerts R, De Deyne S, Meersmans K, Liuzzi AG, Storms G Vandenberghe R. 2019. Redefining the
.123 resolution of semantic knowledge in the brain: advances made by the introduction of models of
.124 semantics in neuroimaging. *Neuroscience and Biobehavioral Reviews*.
.125 <https://doi.org/10.1016/j.neubiorev.2019.05.015>
- .126 Bruni E, Tran N, Baroni M. (2014). Multimodal distributional semantics. *J. Artif. Intell. Res.* 49, 1–47.
- .127 Brysbaert M, New B. (2009). Moving beyond Kucera and Francis: A critical evaluation of current word
.128 frequency norms and the introduction of a new and improved word frequency measure for American
.129 English. *Behavior Research Methods*, 41: 977–990.
- .130
- .131
- .132
- .133
- .134
- .135
- .136
- .137
- .138
- .139
- .140

- .141
.142 Brysbaert M, Warriner AB, Kuperman V. (2014). Concreteness ratings for 40 thousand generally known
.143 English word lemmas. *Behav. Res. Methods* 46, 904–911.
- .144
.145 Benjamini Y, Yekutieli D (2001). The control of the false discovery rate in multiple testing under
.146 dependency. *Annals of Statistics*. 29 (4): 1165–1188.
- .147
.148 Caramazza A, Hillis A. (1991). Lexical organization of nouns and verbs in the brain. *Nature*.
.149 349(6312):788–90.
- .150
.151 Conneau A, Kiela D, Schwenk H, Barrault L, Bordes A. (2017). Supervised learning of universal
.152 sentence representations from natural language inference data. In *Proceedings of the 2017 Conference*
.153 on Empirical Methods in Natural Language Processing, pages 670–680, Copenhagen, Denmark.
.154 Association for Computational Linguistics.
- .155
.156 Cox RW. (1996). AFNI: software for analysis and visualization of functional magnetic resonance
.157 neuroimages. *Comput Biomed Res*. 29:162–173.
- .158
.159 Cree GS, McRae K. (2003). Analyzing the factors underlying the structure and computation of the
.160 meaning of chipmunk, cherry, chisel, cheese, and cello (and many other such concrete nouns). *J Exp*
.161 *Psychol Gen*. 132(2):163–201.
- .162
.163 Chang KM, Mitchell TM, Just MA. (2010). Quantitative modeling of the neural representations of
.164 objects: How semantic feature norms can account for fMRI activation. *Neuroimage: Special Issue on*
.165 *Multivariate Decoding and Brain Reading*. 56:716–727.
- .166
.167 Connell L. (2007). Representing object colour in language comprehension. *Cognition* 102: 476–485.
- .168
.169 Desai R, Binder JR, Conant LL, Seidenberg MS. (2009). Activation of sensory-motor and visual areas
.170 by sentence comprehension. *Cerebral Cortex*. 20: 468-478.
- .171
.172 Devlin J, Chang MW, Lee K, Toutanova K. (2019.) BERT: Pre-training of Deep Bidirectional
.173 Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North*
.174 *American Chapter of the Association for Computational Linguistics: Human Language Technologies*,
.175 *Volume 1 (Long and Short Papers)* 2019 Jun (pp. 4171-4186).
- .176
.177 Dove G. (2014). Thinking in Words: Language as an Embodied Medium of Thought. *Topics in Cognitive*
.178 *Science* 6, 371–389.
- .179
.180 Dunlap WP, Cortina JM, Vaslow JB, Burke MJ. 1996. Meta-analysis of experiments with matched
.181 groups or repeated measures designs. *Psychological methods*. 1(2):170.
- .182
.183 Fedorenko E, Thompson-Schill SL. (2014). Reworking the language network. *Trends in cognitive*
.184 *sciences* 18 (3), 120-126.
- .185
.186 Fedorenko E, Scott TL, Brunner P, Coon WG, Pritchett B, Schalk G, Kanwisher N. (2016). Neural
.187 correlate of the construction of sentence meaning. *Proceedings of the National Academy of Sciences*
.188 113 (41), E6256-E6262.
- .189
.190 Fernandino L, Conant LL, Binder JR, Blindauer K, Hiner B, Spangler K, Desai RH. (2013) Where is the
.191 action? Action sentence processing in Parkinson's disease. *Neuropsychologia*. 51(8):1510-7.
- .192

- .193 Fernandino L, Humphries CJ, Seidenberg MS, Gross WL, Conant LL, Binder JR. (2015). Predicting
.194 brain activation patterns associated with individual lexical concepts based on five sensory-motor
.195 attributes. *Neuropsychologia*. doi:10.1016/j.neuropsychologia.2015.04.009.
- .196
- .197 Fernandino L, Humphries CJ, Conant LL, Seidenberg MS, Binder JR. (2016). Heteromodal cortical
.198 areas encode sensory-motor features of word meaning. *Journal of Neuroscience* 36(38): 9763-9769.
- .199
- .200 Fu R, Guo J, Qin B, Che W, Wang H, Liu T. (2014). Learning semantic hierarchies via word
.201 embeddings. In *Proceedings of the 52nd Annual Meeting of the Association for Computational
.202 Linguistics*. Baltimore, Maryland: Association for Computational Linguistics. pp 1199–1209.
- .203
- .204 Fyshe A, Talukdar PP, Murphy B, Mitchell TM. (2014). Interpretable Semantic Vectors from a Joint
.205 Model of Brain- and Text-Based Meaning. *Proceedings of the Meeting of the Association for
.206 Computational Linguistics*. Baltimore, Maryland: Association for Computational Linguistics. pp 489–499.
- .207
- .208 Glasgow K, Roos M, Haufler A, Chevillet M, Wolmetz, M. (2016). Evaluating semantic models with
.209 word-sentence relatedness. arXiv:1603.07253.
- .210
- .211 Glenberg AM, Kaschak MP. (2002). Grounding language in action. *Psychonomic Bulletin & Review*, 9,
.212 558–565.
- .213
- .214 Glenberg AM, Sato M, Cattaneo L, Riggio L, Palumbo D, Buccino G. (2008). Processing abstract
.215 language modulates motor system activity. *Quarterly Journal of Experimental Psychology*, 61, 905–
.216 919.
- .217
- .218 Glenberg A. (2010). Embodiment as a unifying perspective for psychology. *Wiley Interdisciplinary
.219 Reviews: Cognitive Science*, 1(4), 586–596.
- .220
- .221 Guntupalli JS, Hanke M, Halchenko YO, Connolly AC, Ramadge PJ, Haxby JV. (2016). A model of
.222 representational spaces in human cortex. *Cereb. Cortex*, 26, 2919-2934.
- .223
- .224 Hamilton LS, Huth AG. (2018). The revolution will not be controlled: Natural stimuli in speech
.225 neuroscience. *Language, Cognition and Neuroscience*. 21:1-0.
- .226
- .227 Handjaras G, Leo A, Cecchetti L, Papale P, Lenci A, Marotta G, Pietrini P, Ricciardi E. (2017). Modality-
.228 independent encoding of individual concepts in the left parietal cortex. *Neuropsychologia*. 105. 39-49.
- .229
- .230 Handjaras G, Ricciardi E, Leo A, Lenci A, Cecchetti L, Cosottini M, Marotta G, Pietrini P. (2016). How
.231 concepts are encoded in the human brain: a modality independent, category-based cortical
.232 organization of semantic knowledge. *NeuroImage*, 135, 232-242.
- .233
- .234 Harnad, S. (1990). The symbol grounding problem. *Physica D: Nonlinear Phenomena*, 42(1-3), 335–
.235 346.
- .236
- .237 Hassabis D, Kumaran D, Summerfield C, Botvinick M. (2017). Neuroscience-inspired artificial
.238 intelligence. *Neuron* 95 (2), 245-258.
- .239
- .240 Hasson U, Egidi G, Marelli M, Willems RM. (2018). Grounding the neurobiology of language in first
.241 principles: The necessity of non-language-centric explanations for language comprehension. *Cognition*.
.242 180:135-57.
- .243
- .244 Hauk O, Johnsrude I, Pulvermüller F. (2004). Somatotopic representation of action words in human
.245 motor and premotorcortex. *Neuron*. 41(2):301–7.

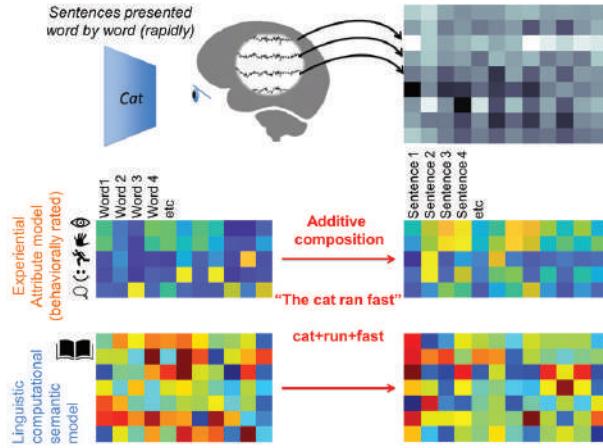
- .246
.247 Haxby JV, Guntupalli JS, Connolly AC, Halchenko YO, Conroy BR, Gobbini MI, Hanke M, Ramadge
.248 PJ. (2011). A common, high-dimensional model of the representational space in human ventral
.249 temporal cortex. *Neuron*, 72:404-416.
- .250
.251 de Heer WA, Huth AG, Griffiths TL, Gallant JL, Theunissen FE. (2017). The hierarchical cortical
.252 organization of human speech processing. 2017. *J. Neurosci.*, 3267-3216.
- .253
.254 Hoerl AE, Kennard RW. (1970). Ridge regression: Biased estimation for nonorthogonal problems.
.255 *Technometrics*, 12(1), 55-67.
- .256
.257 Huth AG, de Heer WA, Griffiths TL, Theunissen FE, Gallant JL. (2016). Natural speech reveals the
.258 semantic maps that tile human cerebral cortex. *Nature* 532, 453-458.
- .259
.260 Huth, AG, Lee T, Nishimoto S, Bilenko, NY, Vu AT, Gallant JL. (2016b). Decoding the semantic content
.261 of natural movies from human brain activity. *Frontiers in systems neuroscience*, 10, 81.
- .262
.263 Just MA, Cherkassky VL, Aryal S, Mitchell TM. (2010). A neurosemantic theory of concrete
.264 noun representation based on the underlying brain codes. *PLoS ONE* 5 (1), e8622
- .265
.266 Kaschak MP, Madden CP, Therriault DJ, Yaxley RH, Aveyard ME, Blanchard AA, Zwaan RA. (2005).
.267 Perception of motion affects language processing. *Cognition*, 94, B79–B89.
- .268
.269 Kaschak MP, Zwaan RA, Aveyard M, Yaxley RH. (2006). Perception of auditory motion affects
.270 language processing. *Cognitive Science*, 30(4), 733–744.
- .271
.272 Kiela D, Clark S. (2014). A systematic study of semantic vector space model parameters. *Proceedings*
.273 of the 2ndWorkshop on Continuous Vector Space Models and their Compositionality (CVSC) at EACL,
.274 pp. 21–30.
- .275
.276 Kiela D, Clark S. (2017). Learning Neural Audio Embeddings for Grounded Semantics in Auditory
.277 Perception. *Journal of Artificial Intelligence Research*. doi.org/10.1613/jair.5665.
- .278
.279 Kriegeskorte N, Mur M, Bandettini P. (2008). Representational similarity analysis—Connecting the
.280 branches of systems neuroscience. *Frontiers in Systems Neuroscience*, 2, 4.
- .281
.282 Kousta ST, Vigliocco G, Vinson DP, Andrews M, Del Campo E. (2011). The representation of abstract
.283 words: why emotion matters. *Journal of Experimental Psychology: General* 140 (1), 14.
- .284
.285 Landauer T, Dumais S. (1997). A solution to Plato's problem: The latent semantic analysis
.286 theory of acquisition, induction, and representation of knowledge. *Psychol. Rev.* 104(2),
.287 211–240.
- .288
.289 Lambon Ralph MA, Jefferies E, Patterson K, Rogers TT. (2017) The neural and computational bases of
.290 semantic cognition. *Nature Reviews Neuroscience*. 2017 Jan;18(1):42.
- .291
.292 Louwerse MM, Jeuniaux P. (2010). The linguistic and embodied nature of conceptual processing.
.293 *Cognition*. 114:96–104.
- .294
.295 Louwerse MM. (2018). Knowing the Meaning of a Word by the Linguistic and Perceptual Company It
.296 Keeps. *Topics in Cognitive Science*. doi:10.1111/tops.12349

- .297 Lund K, Burgess C. (1996). Producing high-dimensional semantic spaces from lexical cooccurrence.
.298 Behav. Res. Methods Instrum. Comput. 28, 203–208.
- .299
- .300 Lynott D, Connell L. (2013). Modality exclusivity norms for 400 nouns: The relationship between
.301 perceptual experience and surface word form. Behavior research methods 45 (2), 516-526.
- .302
- .303 Mahon BZ, Caramazza A. (2008). A critical look at the embodied cognition hypothesis and a new
.304 proposal for grounding conceptual content. Journal of physiology-Paris 102 (1-3), 59-70.
- .305
- .306 Mahon BZ. (2015). What is embodied about cognition? Language, cognition and neuroscience 30 (4),
.307 420-429.
- .308
- .309 Mikolov T, Chen K, Corrado G, Dean, J. (2013). Efficient estimation of word representations in vector
.310 space. arXiv preprint arXiv:1301.3781.
- .311
- .312 Mitchell J, Lapata M. (2010). Composition in distributional models of semantics. Cogn. Sci. 34, 1388–
.313 1429.19.
- .314
- .315 Mitchell TM, Shinkareva SV, Carlson A, Chang K-M, Malave VL, Mason RA, Just MA. (2008).
.316 Predicting human brain activity associated with the meaning of nouns. Science. 320:1191–1195.
- .317
- .318 Oldfield RC. (1971). The assessment and analysis of handedness: The Edinburgh Inventory.
.319 Neuropsychologia. (9):97-113.
- .320
- .321 Paivio A. (1971). Imagery and verbal processes. Holt, Rinehart, and Winston, New York.
- .322
- .323 Patterson K, Nestor PJ, Rogers TT. (2007) Where do you know what you know? The representation of
.324 semantic knowledge in the human brain. Nature Reviews Neuroscience. 8(12):976.
- .325
- .326 Pennington J, Socher R, Manning CD. (2014). GloVe: Global Vectors for Word Representation. In Proc.
.327 Conf. Emp. Meth. Nat. Lang. Proc. Doha, Qatar: Association for Computational Linguistics. pp 1532–
.328 1543.
- .329
- .330 Pereira F, Botvinick M, Detre G. (2013). Using Wikipedia to learn semantic feature representations of
.331 concrete concepts in neuroimaging experiments. Artif Intell. 194:240–252.
- .332
- .333 Pereira F, Lou B, Pritchett B, Ritter S, Gershman S, Kanwisher N, Botvinick M, Fedorenko E. 2018.
.334 Toward a universal decoder of linguistic meaning from brain activation. Nature Communications 9 (963)
- .335
- .336 Peters ME, Neumann M, Iyyer M, Gardner M, Clark C, Lee K, Zettlemoyer L. (2018). Deep
.337 contextualized word representations. In Proceedings of NAACL-HLT 2018 (pp. 2227-2237).
- .338
- .339 Popov V, Ostarek M, Tenison C. (2018). Practices and Pitfalls in Inferring Neural Representations.
.340 NeuroImage, 174, 340-351. doi:10.1016/j.neuroimage.2018.03.041.
- .341
- .342 Pulvermüller F. (2013). How neurons make meaning: brain mechanisms for embodied and abstract-
.343 symbolic semantics. Trends Cogn. Sci. 17 (9), 458–470.
- .344
- .345 Riordan B, Jones MN. (2010). Redundancy in perceptual and linguistic experience: Comparing feature
.346 based and distributional models of semantic information. Topics in Cognitive Science, 3, 303–345.
- .347

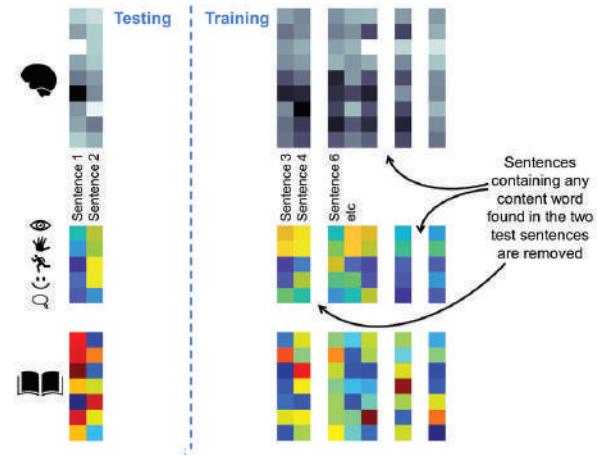
- .348 Roller S, Erk K, Boleda G. Inclusive yet selective: Supervised distributional hypernymy detection. In
.349 Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics:
.350 Technical Papers 2014 Aug (pp. 1025-1036).
- .351
- .352 Speed LJ, Vigliocco G. (2014). Eye movements reveal the dynamic simulation of speed in language.
.353 Cognitive science 38 (2), 367-382.
- .354
- .355 Stanfield RA, Zwaan RA. (2001). The effect of implied orientation derived from verbal context on picture
.356 recognition. Psychological Science, 12, 153–156.
- .357 Sudre G, Pomerleau D, Palatucci M, Wehbe L, Fyshe A, Salminen R, Mitchell T. (2012). Tracking neural
.358 coding of perceptual and semantic features of concrete nouns. Neuroimage. 62:451– 463.
- .359
- .360 Subramanian S, Trischler A, Bengio Y, Pal, C. J. (2018). Learning general purpose distributed sentence
.361 representations via large scale multi-task learning. In Proceedings of the 2018 Inter- national
.362 Conference on Learning Representations (ICLR).
- .363 Sun J, Wang S, Zhang J, Zong C. (2019). Towards sentence-level brain decoding with distributed
.364 representations. In Proceedings of the AAAI Conference on Artificial Intelligence 2019 Jul 17 (Vol. 33,
.365 pp. 7047-7054).
- .366 Talairach J, Tournoux P. (1988). Co-planar stereotaxic atlas of the human brain. 3-Dimensional
.367 proportional system: an approach to cerebral imaging. New York: Thieme.
- .368 Turney P, Pantel P. (2010). From frequency to meaning: Vector space models of semantics. J. Artif.
.369 Intell. Res. 37, 141–188.
- .370
- .371 Utsumi A. (2018). A Neurobiologically Motivated Analysis of Distributional Semantic Models.
.372 arXiv:1802.01830.
- .373
- .374 Vigliocco G, Meteyard L, Andrews M, Kousta S. (2009). Toward a theory of semantic representation.
.375 Language and Cognition 1 (2), 219-247.
- .376
- .377 Vigliocco G, Vinson DP, Druks J, Barber H, Cappa SF. (2011). Nouns and verbs in the brain: a review
.378 of behavioural, electrophysiological, neuropsychological and imaging studies. Neurosci Biobehav Rev.
.379 35(3), 407–26.
- .380
- .381 Vigliocco G, Kousta ST, Della Rosa PA, Vinson DP, Tettamanti M, Devlin JT, Cappa SF. (2014). The
.382 neural representation of abstract words: the role of emotion. Cerebral Cortex 24 (7):1767-1777.
- .383
- .384 Vinson DP, Vigliocco G, Cappa S, Siri S. (2003). The breakdown of semantic knowledge: insights from
.385 a statistical model of meaning representation. Brain Lang. 86(3):347–365.
- .386
- .387 Wang J, Cherkassky VL, Just MA. (2017). Predicting the brain activation pattern associated with the
.388 propositional content of a sentence: Modeling neural representations of events and states. Human
.389 brain mapping 38 (10), 4865-4881.
- .390
- .391 Wang X, Wu W, Ling Z, Xu Y, Fang Y, Wang X, Binder JR, Men W, Gao JH, Bi Y. 2018. Organizational
.392 principles of abstract words in the human brain. Cerebral Cortex. 28(12):4305-18.
- .393
- .394 Wehbe L, Murphy B, Talukdar P, Fyshe A, Ramdas A, Mitchell T. (2014). Simultaneously uncovering
.395 the patterns of brain regions involved in different story reading subprocesses. PloS one. 9(11), e11257.
- .396

- .397 Willems RM, Toni I, Hagoort P, Casasanto D. (2010). Neural dissociations between action verb
.398 understanding and motor imagery. *Journal of Cognitive Neuroscience*, 22(10), 2387–2400.
.399
- .400 Winter B, Bergen B. (2012). Language comprehenders represent object distance both visually and
.401 auditorily. *Language and Cognition*, 4(1), 1–16.
.402
- .403 Yang Y, Wang J, Bailer C, Cherkassky V, Just MA. (2017). Commonality of neural representations of
.404 sentences across languages: Predicting brain activation during Portuguese sentence comprehension
.405 using an English-based model of brain function. *NeuroImage* (146):658-666.
.406
- .407 Zwaan RA, Stanfield RA, Yaxley RH. (2002). Language comprehenders mentally represent
.408 the shapes of objects. *Psychol. Sci.* 13 (2), 168–171.
.409
- .410 Zwaan RA, Pecher D. (2012). Revisiting mental simulation in language comprehension: Six replication
.411 attempts. *PLoS ONE*, 7, 12.
.412
- .413 Zwaan RA. (2014). Embodiment and language comprehension: Reframing the discussion. *Trends in
Cognitive Sciences*, 18(5), 229–234.

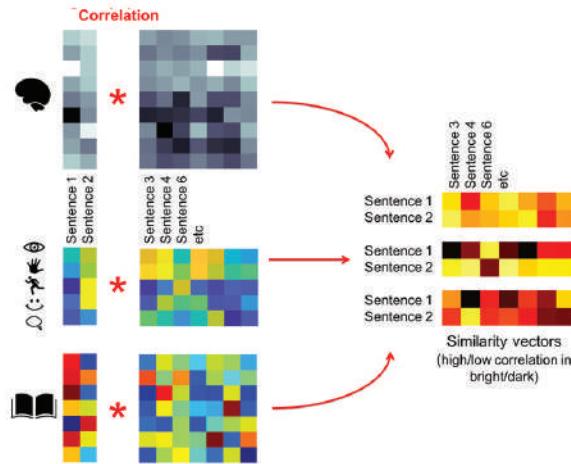
1. Neural, experiential and text-based sentences



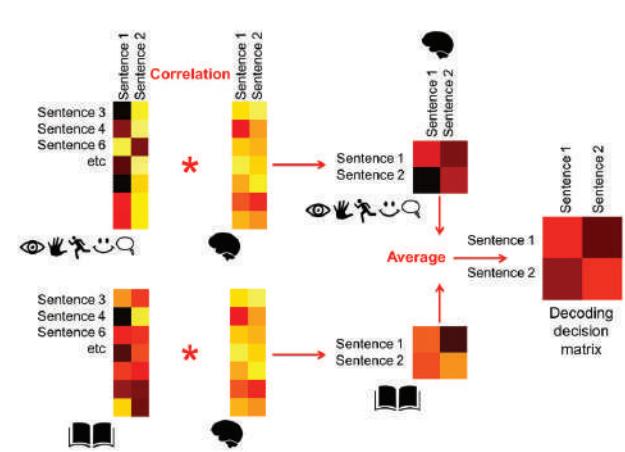
2. Cross validation training/test split



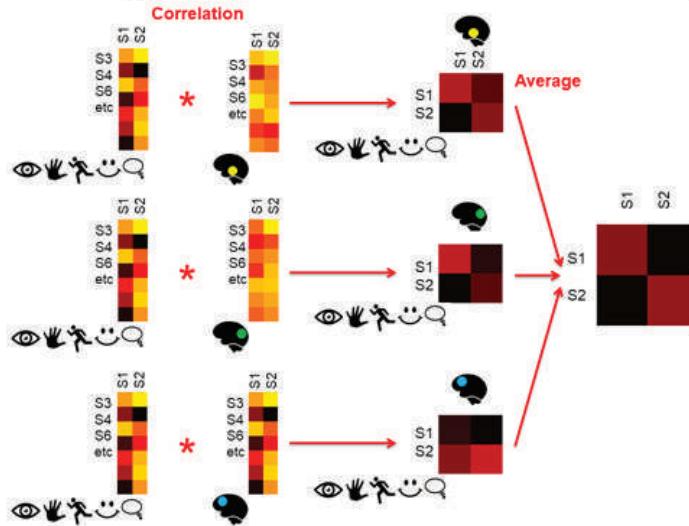
3. Re-representation in common similarity space



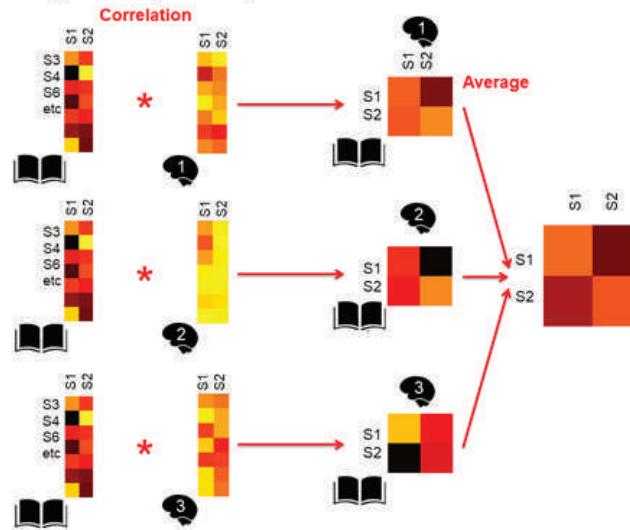
4. Multimodal model-based decoding



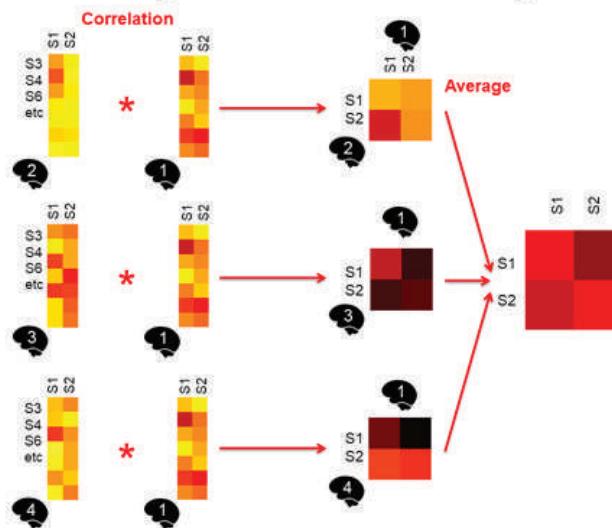
Combining multiple ROIs in model-based decoding

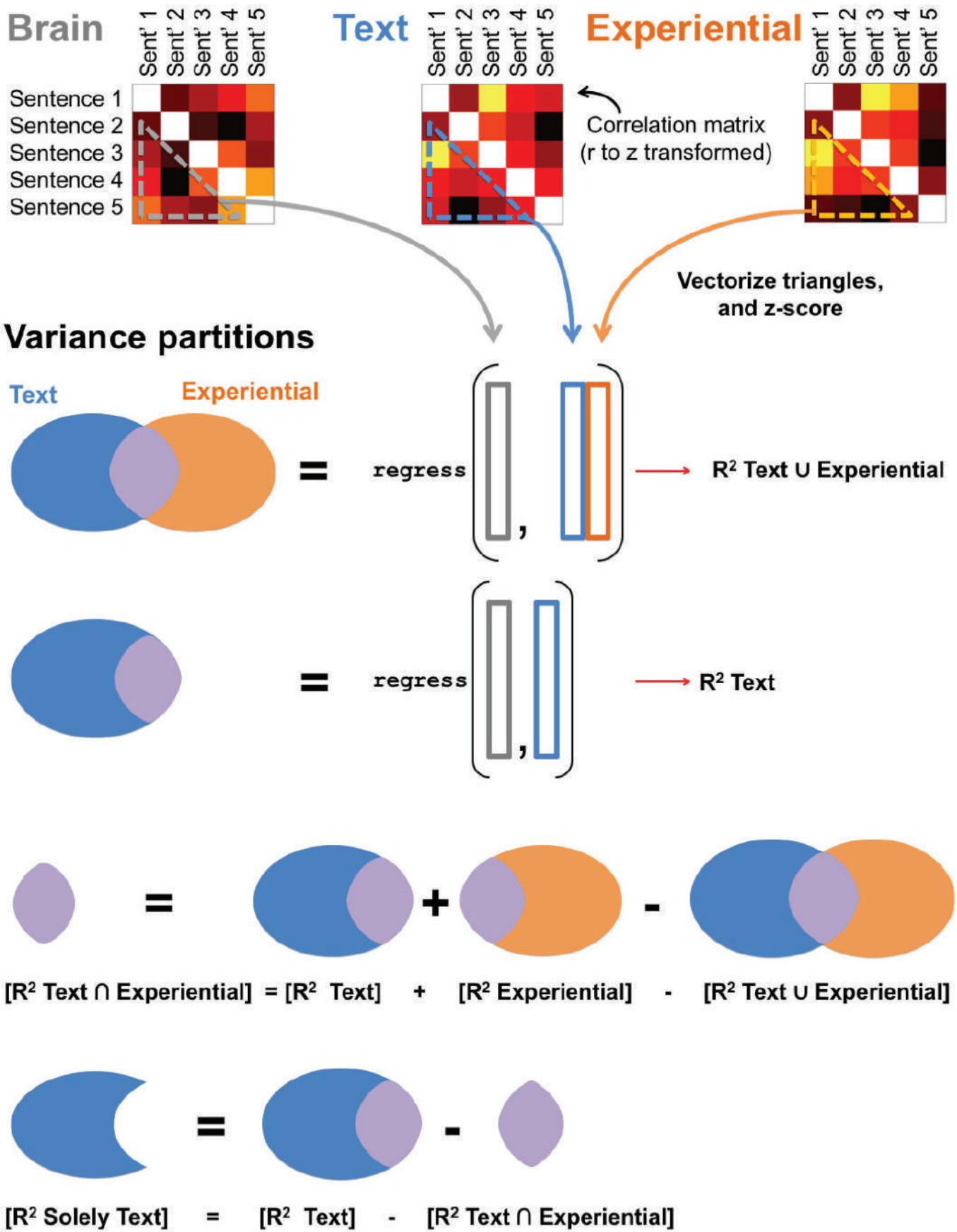


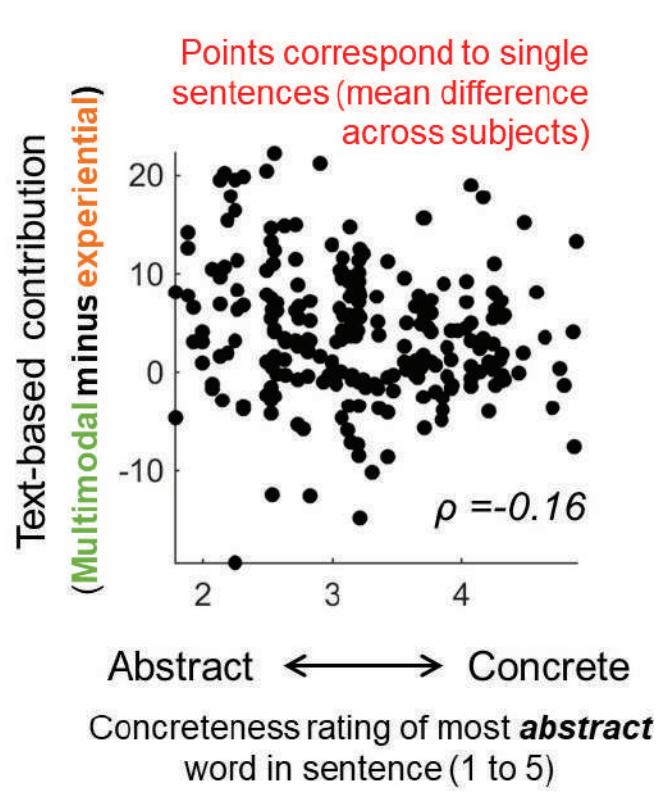
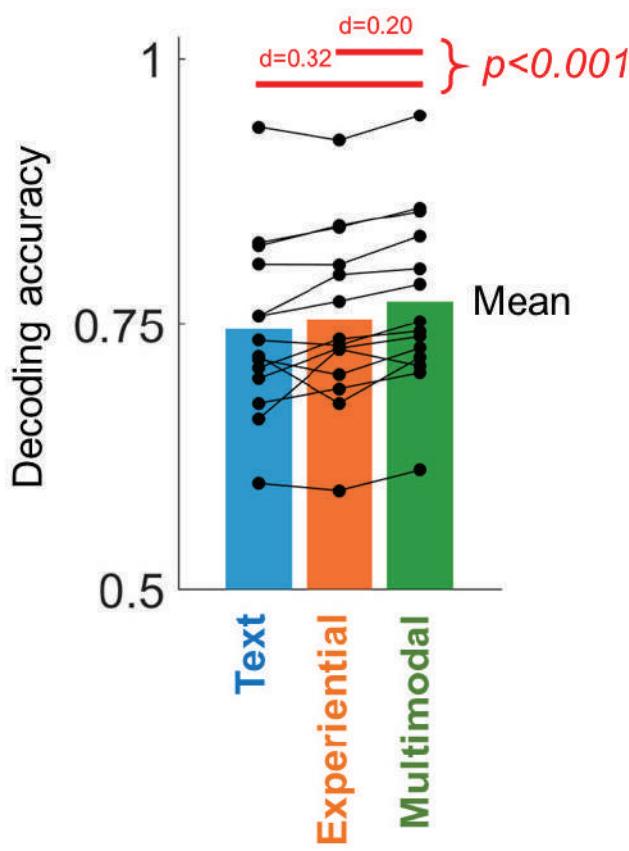
Combining multiple subjects in model-based decoding

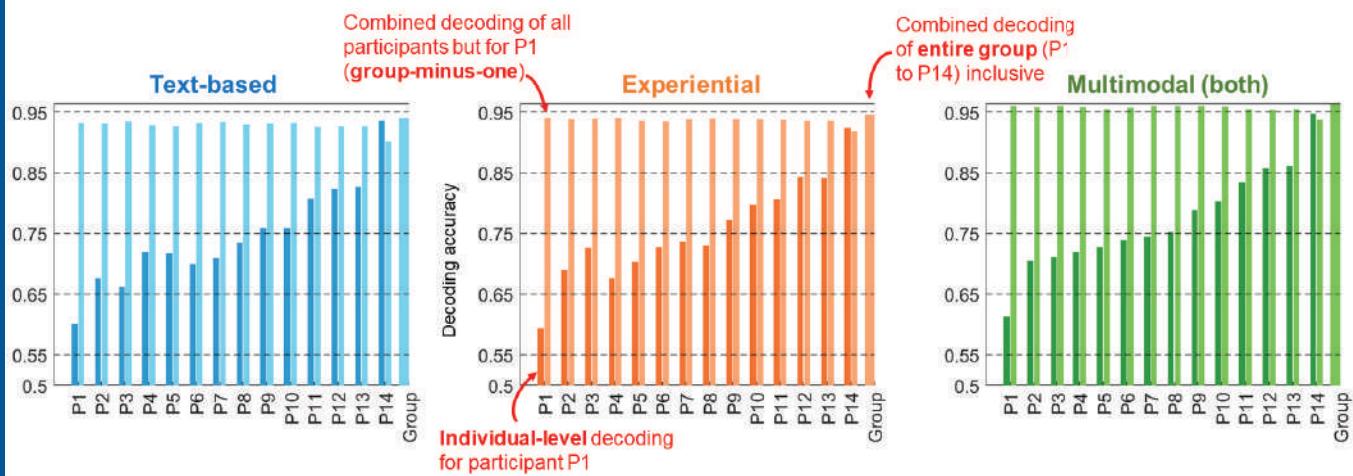


Cross subject brain-based decoding

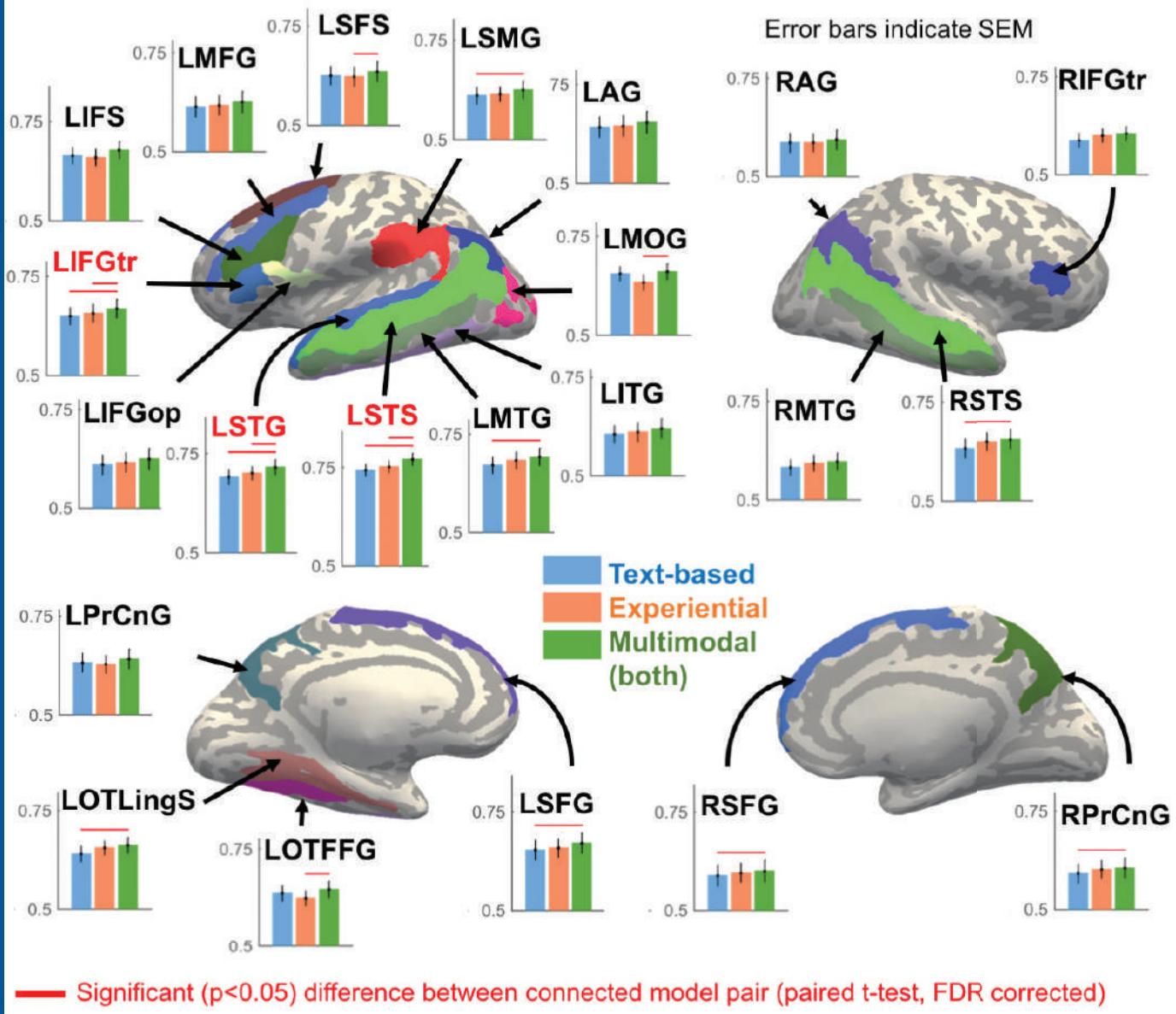


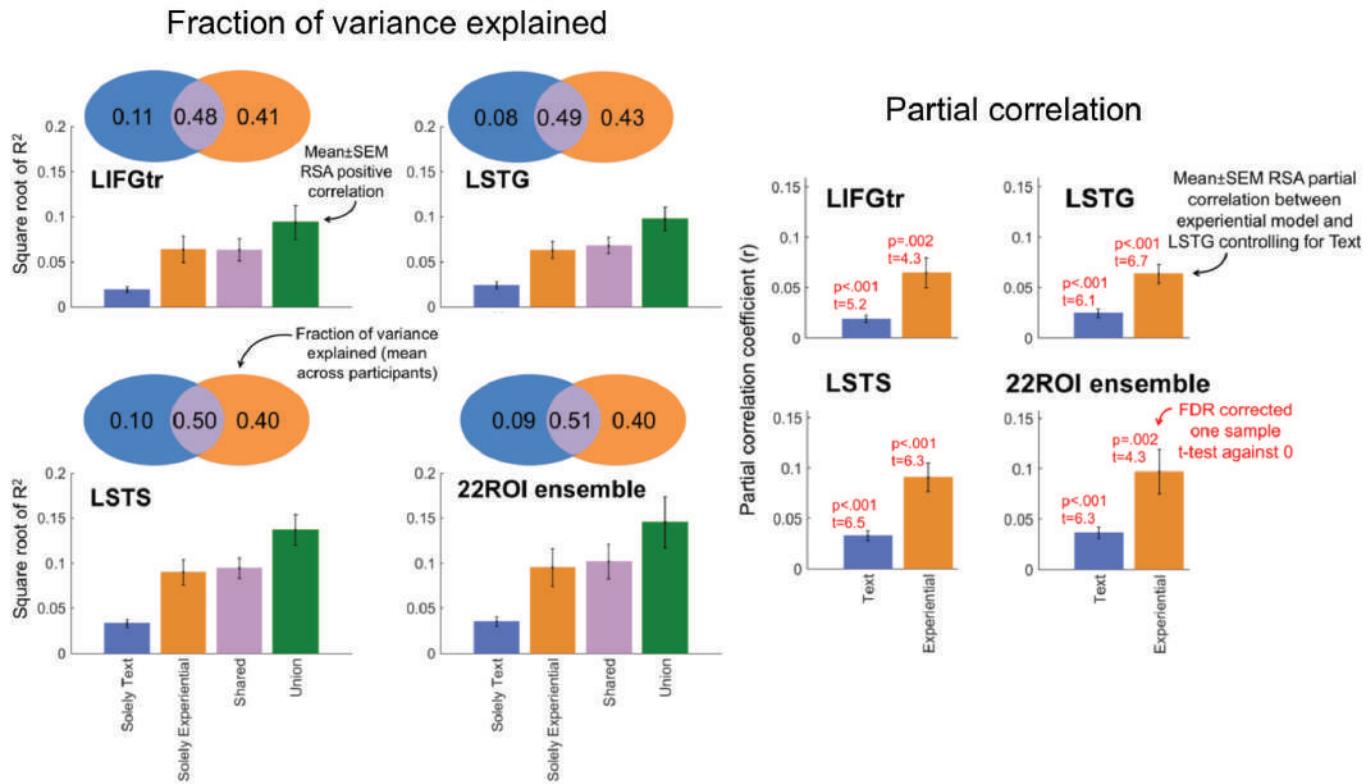


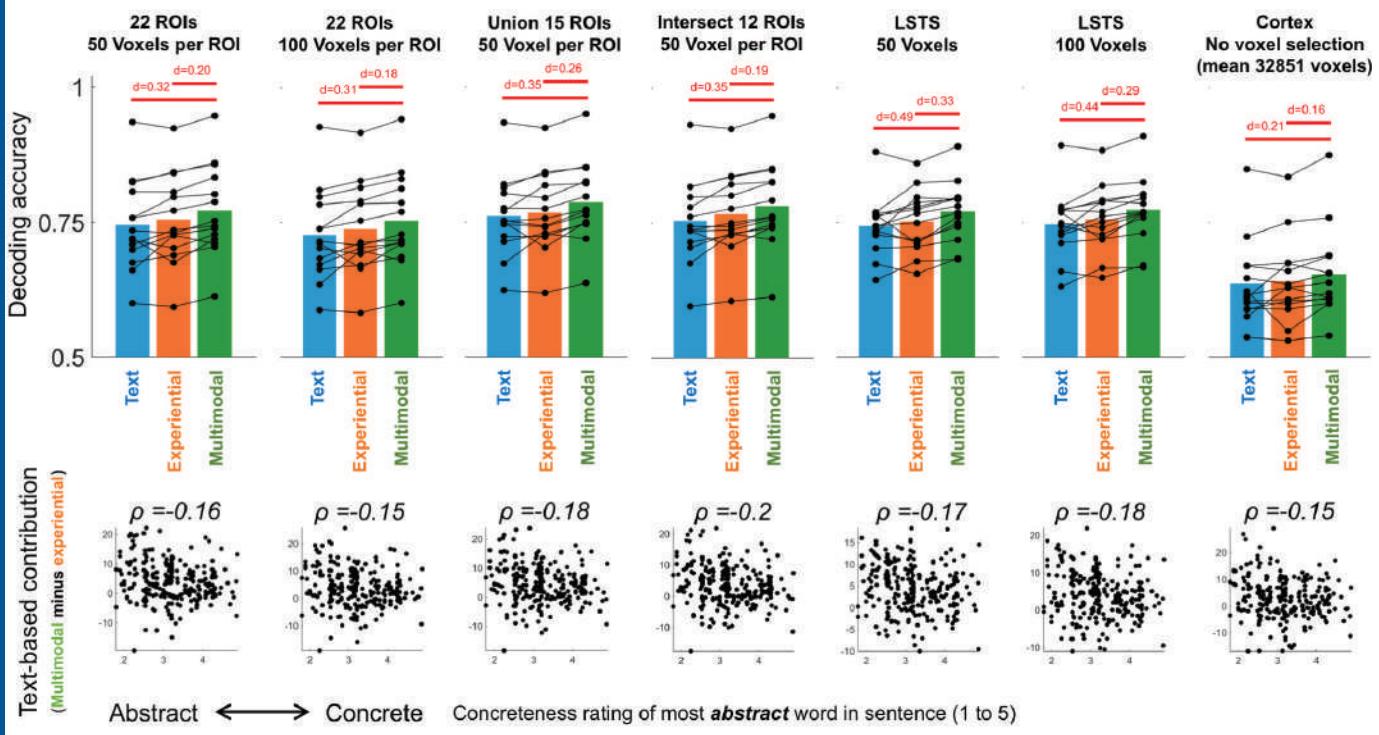




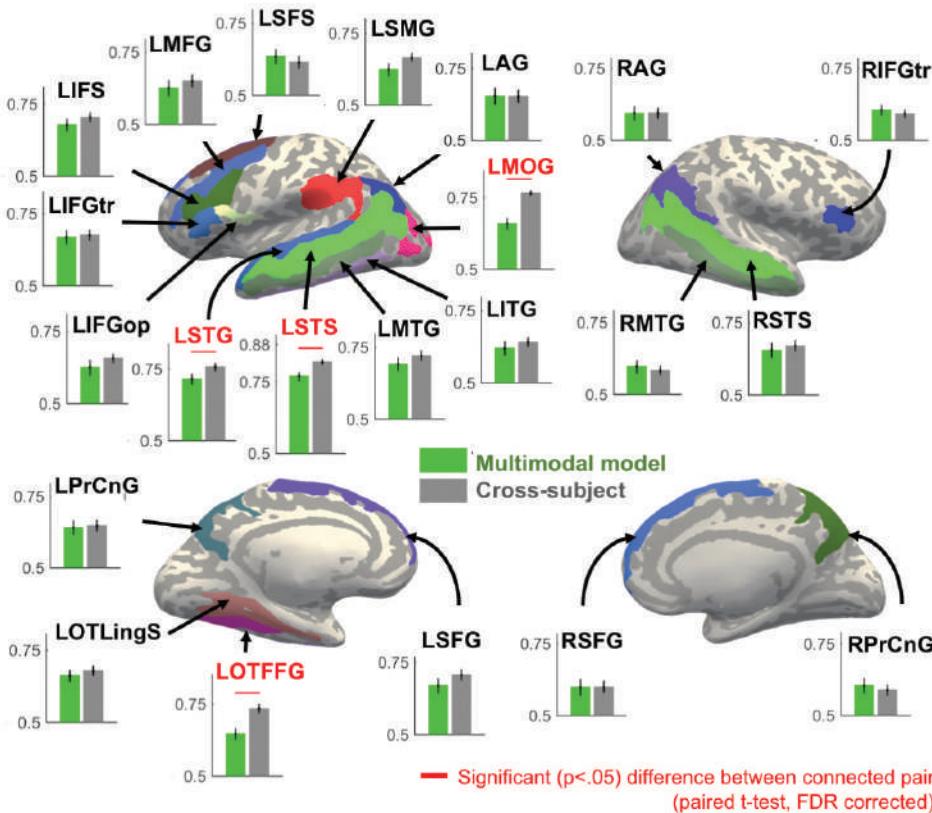
**Note: each participant contributed complementary information
(entire group accuracy surpassed all individuals, and all group-minus-one combinations)**



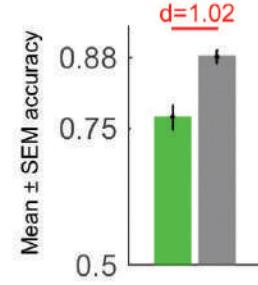




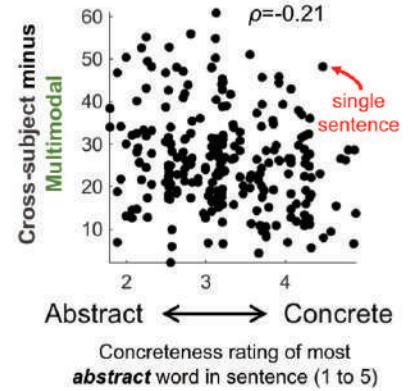
Note: Cross-subject (brain activation-based) decoding is liable to decode semantic and non-semantic (e.g. orthographic/syntactic) elements of neural activation



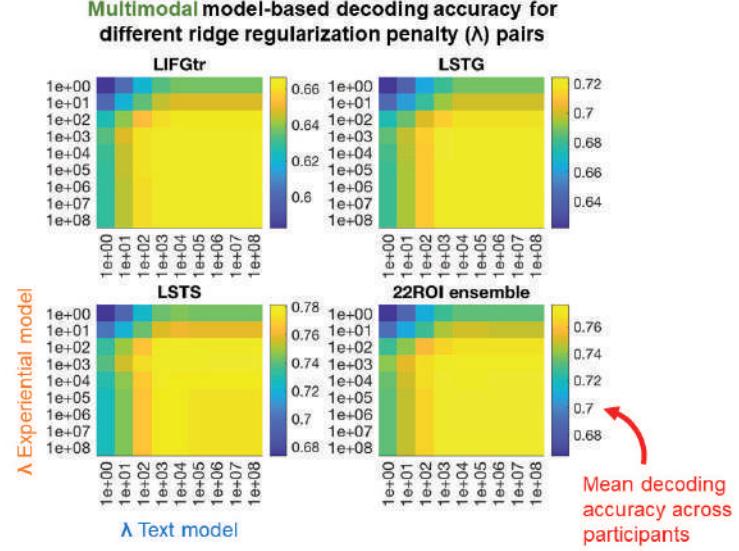
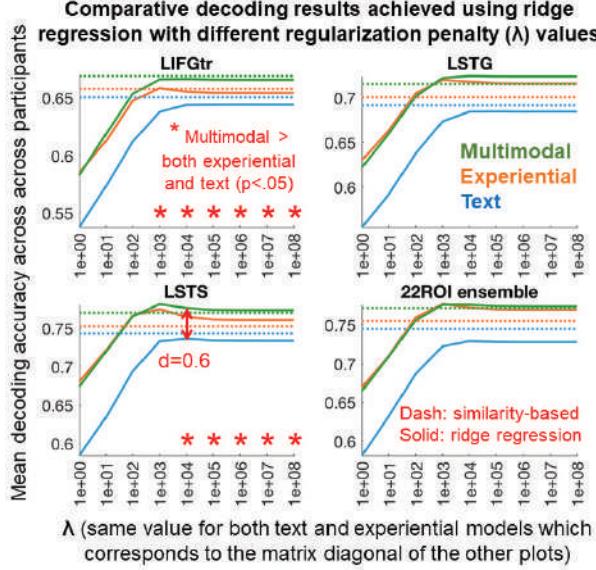
22 ROIs decoded as an ensemble



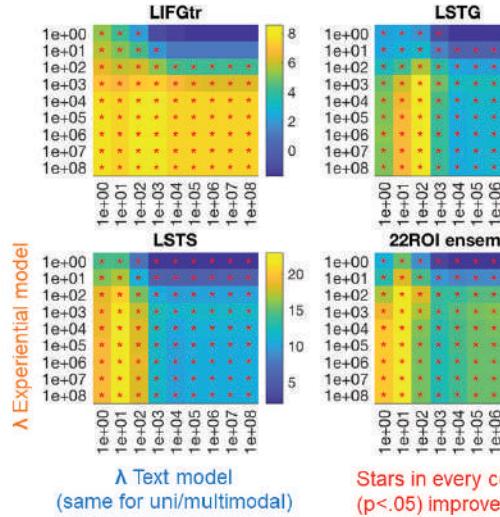
Cross-subject approach particularly improves on decoding of sentences containing **abstract** words



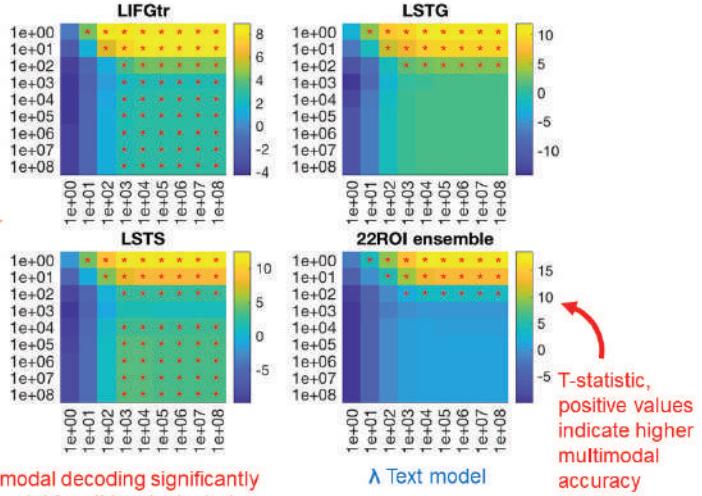
Concreteness rating of most **abstract** word in sentence (1 to 5)



T-tests between text and multimodal model-based decoding accuracy for different λ pairs



T-tests between experiential and multimodal model-based decoding accuracy for different λ pairs



Text-based decoding using similarity approach
Experiential attribute decoding using ridge regression
Multimodal decoding using a fusion of the above

