# Can computers understand word meanings like the human brain does? Comparable semantic representation in neural and computer systems

Linmin Zhang[1, 2, *], Lingting Wang[2, 3], Jinbiao Yang[4], Peng Qian[5], Xuefei Wang[6], Xipeng Qiu[6], Zheng Zhang[1, 7], and Xing Tian[1, 2, 3, *]

[1]*NYU Shanghai*
[2]*NYU-ECNU Institute of Brain & Cognitive Science at NYU Shanghai*
[3]*East China Normal University*
[4]*Max Planck Institute for Psycholinguistics*
[5]*MIT*
[6]*Fudan University*
[7]*AWS Shanghai AI Lab*
[*]*Corresponding authors: linmin.zhang@nyu.edu (L. Z.), xing.tian@nyu.edu (X. T.).*
*Mailing address: NYU Shanghai, 1555 Century Avenue, 200122, Shanghai, China.*

## Abstract

Semantic representation, a crucial window into human cognition, has been studied independently in neuroscience and computer science. A deep understanding of neural computations in the human brain and the revolution to a strong artificial intelligence appeal for a necessity of joint force in the language domain. We investigated the representational formats of comparable lexical semantic features between these two complex systems with fine temporal resolution neural recordings. We found semantic representations generated from computational models significantly correlated with EEG responses at an early stage of a typical semantic processing time window in a two-word semantic priming paradigm. Moreover, three selected computational models differentially predicted EEG responses along the dynamics of word processing in the human brain. Our study provided a finer-grained understanding of the neural dynamics underlying semantic processing and developed an objective biomarker for assessing human-like computation in computational models. Our novel framework trailblazed a promising way to bridge across disciplines in the investigation of higher-order cognitive functions in human and artificial intelligence.

# 1 Introduction

Humans intuitively know that the meaning of the word *moon* is more related to *stars* than to *apples*. Establishing semantic similarity among concepts is a rudimentary adaptive trait for generalization. As an initial step for simulating human intelligence, computational models need to establish semantic relationship among words as well. To leap towards real artificial intelligence, we need to bridge representational formats independently developed from two complex systems – our brain and the computer.

Bridging the representational formats between computers and human brain has recently obtained promising breakthroughs. For example, in vision, the representations in visual hierarchy have been mapped onto distinct layers in deep neural networks (Khaligh-Razavi and Kriegeskorte 2014, Yamins et al. 2014). However, the important branch of artificial intelligence – natural language processing (NLP) – has yet to make substantial connections to higher-level cognitive function of language. The lack of fine-grained neurolinguistic processing models and granular neural recording methods constrains the progress in the language domain (Poeppel 2012). In this project, we proposed a novel approach to join forces across computer science and cognitive neuroscience. By searching for the correlations between neural activity recorded by electroencephalography (EEG) and semantic similarity learned by deep learning models of NLP, our work pioneered in bridging the gap in two ways. Specifically, (a) semantic information encoded in computational models unveiled the neural dynamics of semantic processing; (b) neural data quantified a biomarker for objectively assessing human-like semantic similarity in NLP models.

**Semantics in computer science and cognitive neuroscience**

Within computer science, semantic representation is the cornerstone of complex tasks such as information retrieval, question answering, machine translation, document clustering, etc. Earlier approaches were typically confined to algorithms that require the use of expert-knowledge-based corpus like WordNet (e.g., Resnik 1995, 1999, Lin 1998). Recent development in deep learning NLP models creates embedding representations based on the idea that lexical semantic information is reflected by word distribution (Harris 1954, Firth 1957, Miller 1986). Specifically, embedding models learn semantic representation from words' distribution in their context in a large corpus. Distributional information of words is compressed into dense, lower-dimensional vectors. The
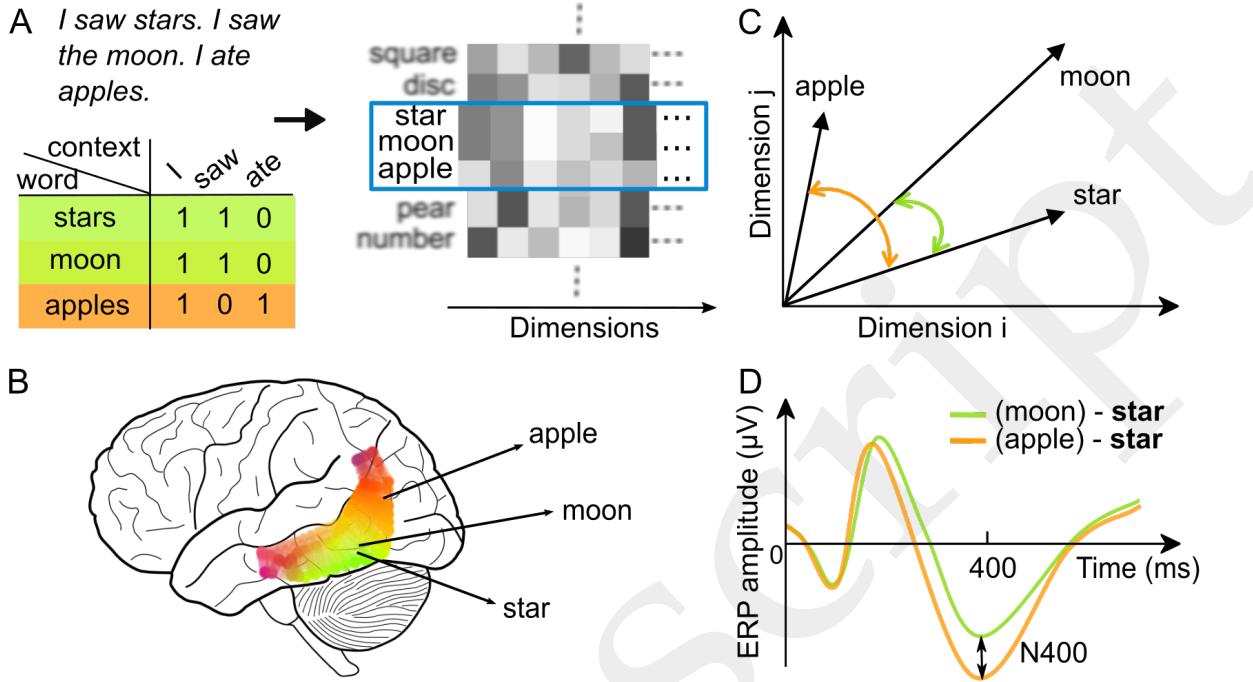
Figure 1: Schematic diagram of semantic representations in the human brain and word embedding models. **A)** A schematic diagram showing how the frequency of words in context yields embedding representations in computational models. Semantically similar words share higher distributional similarity, as illustrated by the counts of neighboring words in the sample mini corpus. Computational models learn semantic representation from words' distribution and generate embedding representations. **B)** A schematic diagram showing the semantic space in the human brain à la Huth et al. (2016). Semantically more similar concepts are represented with more cortical overlaps, indicating shared features. **C)** A schematic diagram showing how the angle between high-dimensional vectors represents semantic similarity in computational models. The angle between two high-dimensional vectors (only two dimensions are used for demonstration) represents semantic similarity. The smaller an angle is (i.e., a higher cosine value), the higher semantic similarity (e.g., the angle between *star* and *moon* is smaller than the one between *star* and *apple*, because *star* and *moon* share more features, as shown in Fig. 1A). **D)** A schematic diagram showing how the amplitudes of neural responses represent semantic similarity in the human brain (e.g., N400, see Kutas and Federmeier 2011). The smaller amplitudes in neural responses to a word are observed when it shares more semantic features with its preceding word (e.g., *star* shares more features with *moon* than with *apple*, as shown in Fig. 1B).

50 similarity between two words can be represented by the cosine value of the angle

51 between the vectors (see Turney and Pantel 2010, and see Fig. 1 for an illustration).

52     Within cognitive science, ample empirical evidence has shown that the similarity of

53 semantic representation has a profound impact on human behavior (Neely 1976, 1977,

54 Lorch Jr 1982, Balota 1983, Anderson 1983, Roelofs 1992, Kiefer 2002). For example, in
55 lexical decision, widely observed priming effects consist in that humans react to a word
56 (e.g., *star*) faster when it is preceded by a semantically related word (e.g., in the pair *moon*
57 – *star*) than by a semantically unrelated word (e.g., in the pair *apple – star*). The
58 presentation of the first word (i.e., the prime) activates a node in a semantic network,
59 which automatically spreads to neighboring nodes, facilitating the processing of the
60 second word (i.e., the target) if it is semantically related (Collins and Loftus 1975)

61 These behavioral findings on semantic similarity were further supported by
62 neuroimaging studies. For example, the voxel-wise modelling neuroimaging study has
63 yielded a semantic map in the human brain, on which concepts sharing more semantic
64 features are mapped to closer brain areas (Huth et al. 2016). In electrophysiological
65 studies, N400 effects – less neural activity around 400 ms after the onset of a more
66 semantically expected word – were observed in both contextual and priming settings
67 (e.g., Bentin et al. 1985, Kutas and Hillyard 1989, Holcomb 1993, Brown and Hagoort
68 1993, Federmeier and Kutas 1999, Deacon et al. 2000, Kiefer 2002) (see Fig. 1).

69 So far semantic representations have been investigated independently in computer
70 science and cognitive neuroscience. It remains unclear to what extent representations
71 yielded from computer models resemble to the ground truth of human representation.

**Bridging semantic representations in the human brain and computer models**

73 Recently, computational models have started advancing our understanding of language
74 processing in human brain(Brennan 2016). The bridging of representations between
75 neural activity and computational models has been preliminarily investigated in
76 sentential context using N400 effects (Ettinger et al. 2016, Broderick et al. 2018). However,
77 neural activity recorded during the comprehension of sentential stimuli and continuous
78 speech was driven by both compositional processing (e.g., the composition between *lamb*
79 and *stew*, yielding *lamb stew*, see e.g., Bemis and Pylkkänen 2011, Zhang and Pylkkänen
80 2015, Pylkkänen 2019) and semantic processing (e.g., similarity-based spreading from
81 *lamb* to *stew*), making neural data hardly comparable with pure semantic representations
82 yielded from word embedding models.

83 Therefore, our study focused on the representation of lexical semantics in the human
84 brain and computer models. We adopted a canonical semantic priming design that
85 elicited the measures of semantic similarity in the brain, directly comparable to semantic
86 representations yielded by computational models without confounding factors from

4

87 compositional processing. We predicted that the two measures from the brain and

88 computers would correlate in a rather narrow time window within classical N400

89 component, presumably at the beginning of the processing purely related to semantic

90 representation without contamination from compositional processing.

91 Moreover, we selected three representative word embedding models, differing in the

92 way of learning semantic representation. The CBOW (Continuous Bag-of-words) model

93 (Mikolov et al. 2013) solely uses local context – a number of words immediately

94 preceding and following a word. The other two models are based on CBOW. The GloVe

95 (Global Vectors) model (Pennington et al. 2014) combines both local context and global

96 corpus statistics for learning word representation. The CWE (Character-enhanced Word

97 Embedding) model (Chen et al. 2015) captures both word-external local contextual

98 information and word-internal character information. We predicted that both GloVe and

99 CWE would correlate with brain responses better than CBOW. The better correlation

100 would occur at different times because of particular features of the models – CWE at an

101 earlier perceptual stage due to its inclusion of character-level information, whereas

102 GloVe at a later stage reflecting semantic processing.

103 By assessing the representational formats with a well-controlled experiment and

104 millisecond-level neural recordings, we provided a framework directly bridging

105 semantic representations between the human brain and computers. Our aim was

106 twofold: (a) information encoded in NLP models contributed to a finer-grained

107 understanding of the neural dynamics underlying semantic processing; (b) neural data

108 contributed an objective assessment for human-like language processing in NLP models.

## 2 Methods

### 2.1 Participants

111 A group of 30 healthy right-handed native Chinese speakers participated in the study.

112 All had normal or corrected-to-normal vision. Five participants were excluded from data

113 analyses: three due to excessive noise during recording, and two for being outliers in

114 terms of accuracy in the behavioral task (more than 3 standard deviations below the

115 average). Thus, 25 participants were included in EEG data analyses (14 females; average

116 age = 22.6 years, $SD$ = 2.8 years). All data were collected at the EEG lab at the

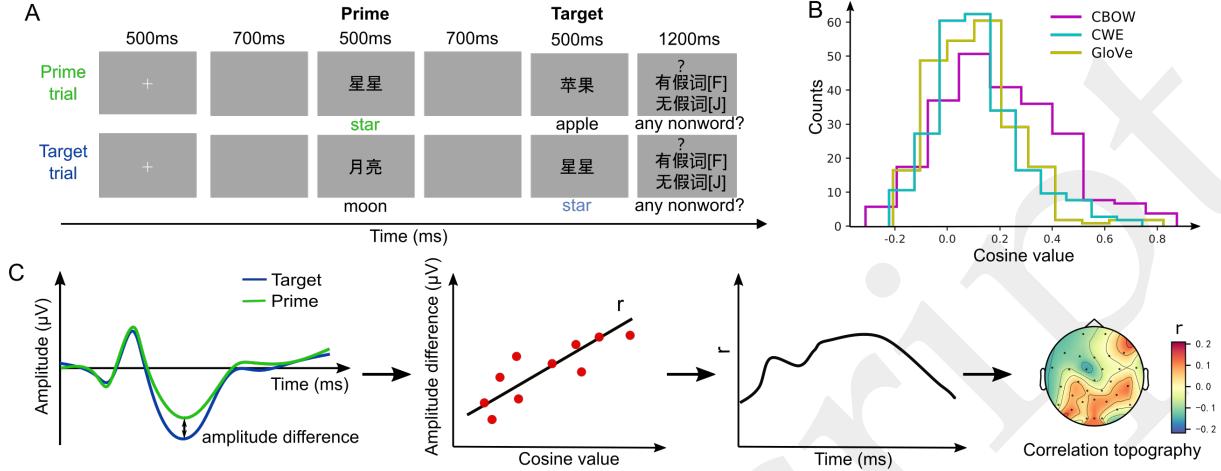117 NYU-ECNU Institute of Brain and Cognitive Science at NYU Shanghai (Shanghai,

Figure 2: Experimental procedure and single-trial correlation analysis. **A)** The trial structure of the experiment. Sample trials are illustrated for the two-word priming paradigm. In each trial, a prime word was followed by a target word. Each word (here 星星 (star)) was used once at the prime position (in the prime trial) and once at the target position (it the target trial). English translations below the screens are for demonstration only, but not included in the expreiment. **B)** Stimuli statistics of semantic similarity generated by the three computational models. **C)** The flowchart of single-trial correlational analysis, (i) Computing the amplitude differences between single-trial EEG responses to the same word at its target vs. prime presentation (target minus prime); (ii) For the 240 word pairs, calculating the correlation between cosine values generated from computational models and amplitude differences from step (i) at each time point in each sensor; (iii) The obtained correlation coefficients form a waveform across time for each sensor; (iv) The distribution of correlation coefficients from all sensors is plotted in a topography at each latency.

China). This study was approved by the local ethical committee at NYU Shanghai. Written consents were obtained from each participant.

## 2.2   Experimental design and stimuli

Our EEG experiment adopted a canonical two-word priming paradigm, with stimuli visually presented to the participants. We used 240 pairs of two-character Chinese nouns as critical stimuli. We randomly selected nouns to form 'prime-target' pairs. Among these 'prime-target' word pairs, some pairs (e.g., 月亮 (moon) – 星星 (star)) are intuitively of a higher semantic similarity than others (e.g., 苹果 (apple) – 月亮 (moon)). This random selection procedure yielded a distribution of semantic similarity (between prime and target) shown in Fig. 2B (see the entire stimuli list at https://ray306.github.io/brain_NLP/).

6

129  To construct 240 critical trials, we used 240 distinct nouns. Each noun appeared at the
130  prime position once and at the target position once (see Fig. 2A). For each noun (e.g., 月
131  亮 (moon)), the EEG responses elicited at the prime position (e.g., in the trial 月亮
132  (moon) – 星星 (star)) represent semantic retrieval of its out-of-context meaning.
133  Whereas, the EEG responses elicited at the target position (e.g., in the trial 苹果 (apple) –
134  月亮 (moon)) include the influence of the preceding word. Thus, the difference between
135  these two EEG responses to the same word at different positions is priming effects,
136  reflecting semantic similarity without the contamination from semantic retrieval.
137  Therefore, we extracted the neural measure directly comparable to the semantic
138  similarity computed from NLP models. Moreover, we extended the previous
139  condition-level computation of ERP differences to the trial-level and provided a
140  trial-level measurement of semantic priming effects.

141  We added 120 additional pairs of stimuli as fillers, in which either the prime or the
142  target was a two-character non-word (e.g., 害天, 粽七). Thus, a total of 360 trials were
143  included in this experiment. Participants were instructed to perform a lexical decision
144  task, judging whether a trial contained a non-word. The purpose was to keep
145  participants alert, encouraging them to process the stimuli at least to the lexical
146  semantics level.

147  The trial structure is illustrated in Fig. 2A. Each trial started with a fixation lasting for
148  500 ms. After a 700 ms blank screen, the prime was presented for 500 ms. After another
149  700 ms blank screen, the target was also presented for 500 ms, followed by a question
150  mark '?' and a prompt for the lexical decision task. The stimuli were in a white 40-point
151  Songti font on a gray background. The 360 trials were divided into 6 blocks, each
152  containing 60 trials. The critical trials and fillers were pseudo-randomized and
153  quasi-evenly distributed in each block. The blocks were also pseudo-randomized.
154  Between blocks, participants could take a short rest. The experimental presentation was
155  programmed with a Python package – Expy (https://github.com/ray306/expy), an
156  in-house software for presenting and controlling psychological experiments, available at
157  http://slang.science.

## 2.3  Procedure of data collection

159  EEG recordings took place in an electrically-shielded and sound-proof room. EEG data
160  were continuously recorded via a 32-channel ActiChamp system (Brain Products).

7

161 Electrodes were held in place on the scalp by an elastic cap (ActiCap) in a 10-20
162 configuration as shown in Fig. 3A. Two more electrodes were placed below the left eye
163 and at the outer canthus of the right eye to monitor vertical and horizontal eye
164 movements (electro-oculogram, EOG). Impedance was kept less than 10 kΩ for all
165 electrodes. The EEG signal was recorded in single DC mode, digitized at a sampling rate
166 of 1000 Hz and online referenced to the vertex (Cz), with the use of the software
167 BrainVision PyCoder. The recording session lasted approximately 30 minutes.

## 2.4 Data pre-processing

169 Only the 240 critical trials were included in EEG analysis. EEG data were processed and
170 analyzed with EasyEEG toolbox (Yang et al. 2018,
171 https://github.com/ray306/EasyEEG). Raw EEG data were bandpass filtered between
172 0.1 and 30 Hz and epoched from 200 ms before to 800 ms after the onset of a word.
173 Epochs were baseline corrected with the 200 ms interval before word onset. We removed
174 those epochs affected by large vertical or horizontal eye movements, based on data
175 recorded from the two electrodes monitoring EOG. We further visually inspected the
176 epochs and removed those with large artifacts. The data were re-referenced to the
177 average reference.

## 2.5 Data analyses

### 2.5.1 Behavioral data

180 We checked the accuracy and reaction times for all 360 trials. Reaction times were
181 measured from the onset of prompt for each trial and for each participant. We ran a
182 two-tailed $t$-test on the data of accuracy and reaction times between critical trials and
183 fillers, to verify whether participants paid attention to the stimuli.

### 2.5.2 EEG data

185 The analysis of EEG data constituted two parts. The first part aimed to examine the
186 validity of the data by checking the ERP components in reading as well as N400 priming
187 effects with the use of data averaged across trials (see Section 2.5.2.1). The second part
188 was at the trial level, aiming to test (a) whether EEG responses can be predicted by a
189 computational model within the typical time window for N400 priming effects (see

8

190  Section 2.5.2.2) and (b) among CBOW, GloVe, and CWE, which computational model
191  was the best predictor at which time point (see Section 2.5.2.3).

192  **2.5.2.1 ERP analysis**

193  Trials were averaged for prime and target respectively. We plotted the ERP
194  waveforms in a representative channel (Cz) for ERP to compare our data with N400
195  effects reported in literature. To summarize and visualize the distributed energy
196  fluctuation, we plotted the dynamics of Global Field Power (GFP, see Lehmann and
197  Skrandies 1980), calculated as a geometric mean of electric potentials across all sensors.
198  To reveal and visualize ERP components during word processing, we used an automatic
199  segregation method (Topography-based Temporal-analysis Toolbox, TTT) to detect
200  component boundaries and plotted boundaries along with average ERP responses of
201  each channel and GFP (Wang et al. 2019). To visualize the dynamics of activation
202  patterns, we plotted the topographies across time for ERP responses to prime and target
203  as well as the differences between the two (i.e., target minus prime).

204  **2.5.2.2 EEG data analysis at trial level (a): testing whether EEG responses can be**
205  **predicted by a computational model**

206  All the three selected word embedding models (i.e., CBOW, GloVe, and CWE) were
207  trained on Chinese Wikipedia. These models calculated cosine similarities for the 240
208  word pairs used as critical stimuli, and we correlated the model-generated cosine
209  similarities with single-trial EEG responses, according to the following procedure (see
210  Fig. 2C):

211  First, for each word, we subtracted the EEG responses to its presentation at the prime
212  position from those responses at the target position. This EEG difference for each word
213  represents priming effects with no contamination of semantic retrieval.

214  Second, we calculated the correlation co-efficient $r$ between ERP differences
215  (computed from 240 critical trials by Step 1) and model-generated semantic similarities
216  (cosine values). This calculation of correlation was performed at each time point in each
217  channel.

218  Third, the correlations of all time points at a channel yielded a temporal progression
219  of correlations at this channel.

220  Fourth, based on the previous three steps, we calculated the temporal progression of
221  correlations for all channels and obtained a series of topographies of correlations along

9

222 the time course.

223 We obtained a null distribution of $r$ values by shuffling the pairing among the 240
224 EEG response differences and the 240 cosine similarities for 1000 times. Empirical $r$
225 values were checked against this null distribution to determine the statistic significance
226 (at the level of $p < 0.05$) at each time point.

227 **2.5.2.3 EEG data analysis at trial level (b): testing which computational model was the**
228 **best predictor at which time point**

229 When testing which word embedding model (among CBOW, GloVe, and CWE)
230 was the best predictor at which time point, we conducted permutation tests on
231 correlation $r$ values averaged across channels to estimate the overall predictability of
232 each model. We did the same permutation tests on correlation $r$ values for each channel
233 to examine the spatial distribution of the predictability of each model.

234 Specifically, from the correlation between EEG responses in each of the 32 channels at
235 each of the 800 milliseconds and cosine similarities computed from each of the three
236 computational models, we obtained a $32 \times 800$-dimensional matrix of $r$ values for each
237 model.

238 To estimate the overall predictability of each model, we averaged the absolute $r$
239 values across channels, yielding a line of temporal progression of $r$ for each word
240 embedding model. At each time point, we randomly shuffled the pairing between EEG
241 responses and cosine values generated by the three models for 1000 times. The shuffling
242 yielded a null distribution of $r$ differences between any two models. Empirical $r$
243 differences were checked against this null distribution at each time point.

244 We did the same permutation tests for each channel to further compare the
245 predictability of models and investigate the site of effects.

# 3   Results

## 3.1   Behavioral data

248 The mean accuracy of lexical decision task was 94.6% (SD = 2.4%). The two-tailed $t$-test
249 revealed significant differences between critical trials and fillers (mean accuracy and SD
250 for critical trials: 96% (2.5%); mean accuracy and SD for fillers: 91% (4.6%); $t$ (24) = 4.99; $p$
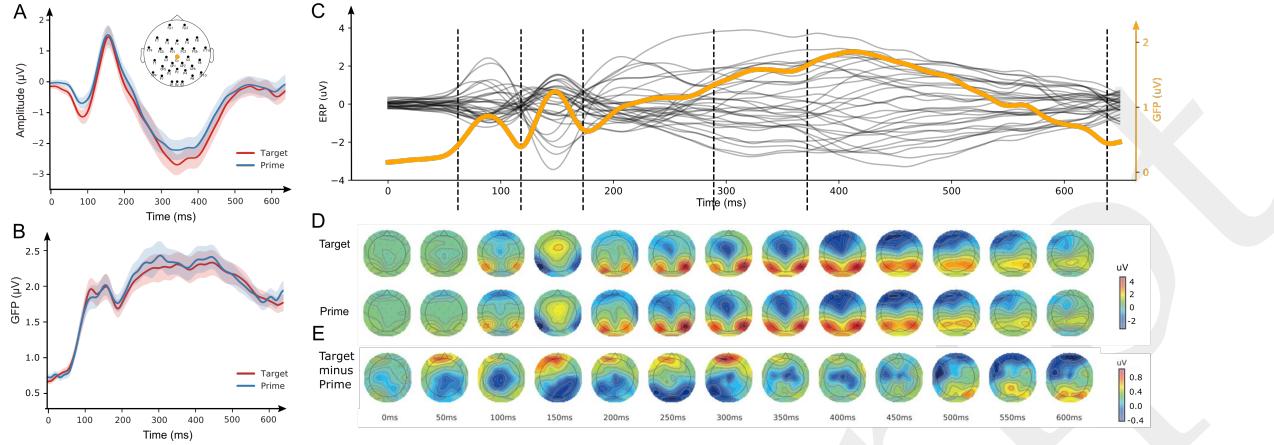251 $< 0.001$).

10

Figure 3: Event-related waveform and topographic responses consistent with perceptual and semantic processes in language comprehension. **A)** The waveform responses in a representative channel (Cz). Typical N400 profile was observed in both prime and target. The montage of sensor locations is inserted with the selected channel Cz highlighted. **B)** The dynamics of GFP. The aggregated neural activity across all sensors represented in GFP shows the similar dynamics that has clear perceptual and semantic activation. **C)** The temporal components revealed in the grand averaged ERP responses across targets and primes. Each black line represents ERP responses in each channel. The orange line represents the GFP across all sensors. The vertical dashed lines label the temporal boundaries between ERP components revealed by an automatic segregation method. **D)** The temporal progression of topographies. The topographies for target and prime were represented in the upper and lower rows respectively. Similar topographic patterns and temporal progressions were observed in both target and prime. **E)** The temporal progression of topographic differences. Differences resulted from subtracting prime from target revealed classic N400 topographic patterns from 250 to 600 ms.

The mean reaction time was 289 ms (SD = 102 ms). The two-tailed $t$-test also revealed significant differences between critical trials and fillers (mean reaction time and SD for critical trials: 296 ms (103 ms); mean reaction time and SD for fillers 274 ms (101 ms); $t(24) = 4.475$; $p < 0.001$).

Behavioral data indicated that participants reacted differently towards critical trials and fillers, suggesting that they fully processed lexical semantic information.

## 3.2  EEG data

### 3.2.1  Results from ERP analysis

ERP responses were obtained after averaging trials for prime and target respectively (Fig. 3). The waveform ERP responses at a representative channel, Cz, clearly indicate the

11

262 evolution of ERP components associated with reading a word (Fig. 3A). Responses to

263 both target and prime showed early visual responses N1 and P2 as well as

264 semantics-related N400 effects, consistent with well-established literature (Kutas and

265 Federmeier 2011). Similar evolution of ERP components was also observed in the

266 dynamics of GFP which included activity of all sensors (Fig. 3B), demonstrating the

267 reliability of elicited data without the potential pitfalls of subjective bias. The boundaries

268 of ERP components were detected based on an automatic segregation method (Wang

269 et al. 2019) and plotted in Fig. 3C. The component after visual processing was further

270 segregated into three sub-components.

271 Topographic responses to prime and target demonstrate consistent evolution of

272 response patterns (Fig. 3D), suggesting common cognitive functions unfolding over time

273 during the reading of these words at prime and target positions. Topographic differences

274 between target and prime showed magnitude differences in sensors over frontal and

275 temporo-parietal regions around 300 ms (Fig. 3E), consistent with the pattern of typical

276 N400 priming effects (see Kutas and Federmeier 2011)

277 Our ERP responses were temporally and spatially consistent with well-established

278 N400 priming effects, demonstrating the reliability and validity of neural measures on

279 semantic similarity.

### 3.2.2 Results from trial-level analysis (a): single-trial EEG responses can be predicted by a computational model

282 We selected GloVe as a representative NLP model. The generated measure of semantic

283 similarity was correlated with single-trial EEG response differences between prime and

284 target (Fig. 4). The correlation was significant at 300 ms after word onset at channel Oz: $r$

285 = 0.173 ($p$ = 0.007) (Fig. 4A). The dynamics of $r$ was obtained in the same channel (Fig.

286 4B). A non-parametric statistics revealed that the GloVe-generated semantic similarity

287 values significantly correlated with EEG response differences between 226 and 306 ms.

288 The spatial distribution of $r$ value was further investigated, by computing the

289 correlations in all sensors (Fig. 4C). The heamap shows that correlations in about half of

290 the sensors were significant between 200 to 300 ms, consistent with the results in Fig. 4B.

291 The distribution of significant correlations in this time window was scrutinized by

292 delineating the evolution of topographies. Most robust correlations were found at

293 sensors over the left frontal and occipital regions, consistent with the typical pattern of

294 N400 effects. The observed semantic processing in a narrow and early time window was
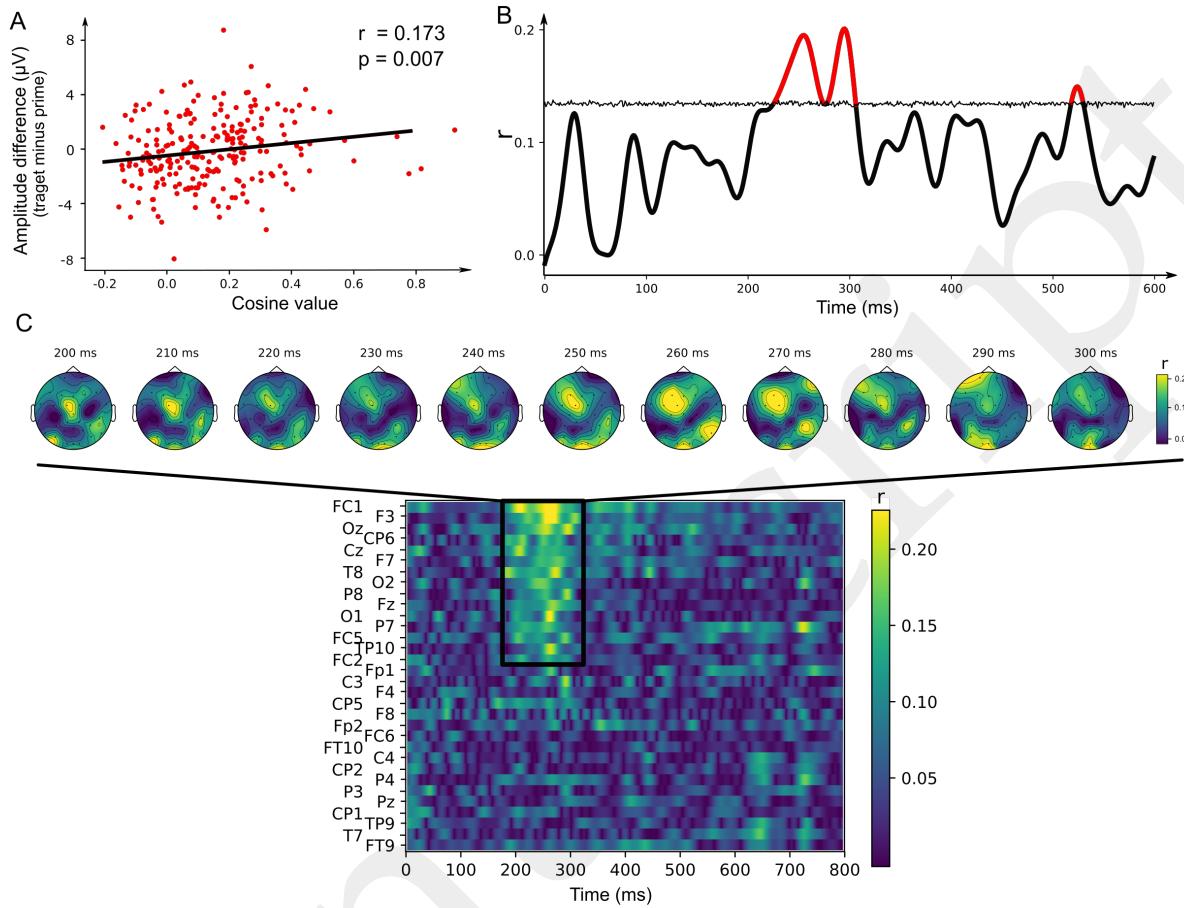
12

Figure 4: Correlations between EEG responses and a word embedding models reveals the dynamics of semantic processing. **A)** Significant correlation was observed between EEG responses in channel Oz at the latency of 300 ms and cosine values computed by the model GloVe. **B)** The temporal progression of correlations (channel Oz). Significant correlations were observed between 226 and 274 ms, between 279 and 306 ms, and between 518 and 529 ms (in red). The significance was determined by the threshold (horizontal line) obtained in a non-parametric permutation test at each time point (alpha level at 0.05). **C)** The spatio-temporal characteristics of correlations. The heatmap of correlations across time and channels revealed significance between 200 and 300 ms in about half of the sensors. The progression of topographies in the time window of significance is zoomed in above. Significant correlations were concentrated in the sensors above the left frontal and tempo-parietal regions.

295  consistent with the findings of semantic dynamics in ERP responses after removing

296  temporal variance among trials (Wang et al. 2019). Taken together, these results

297  demonstrated that NLP models can predict EEG responses, suggesting the common

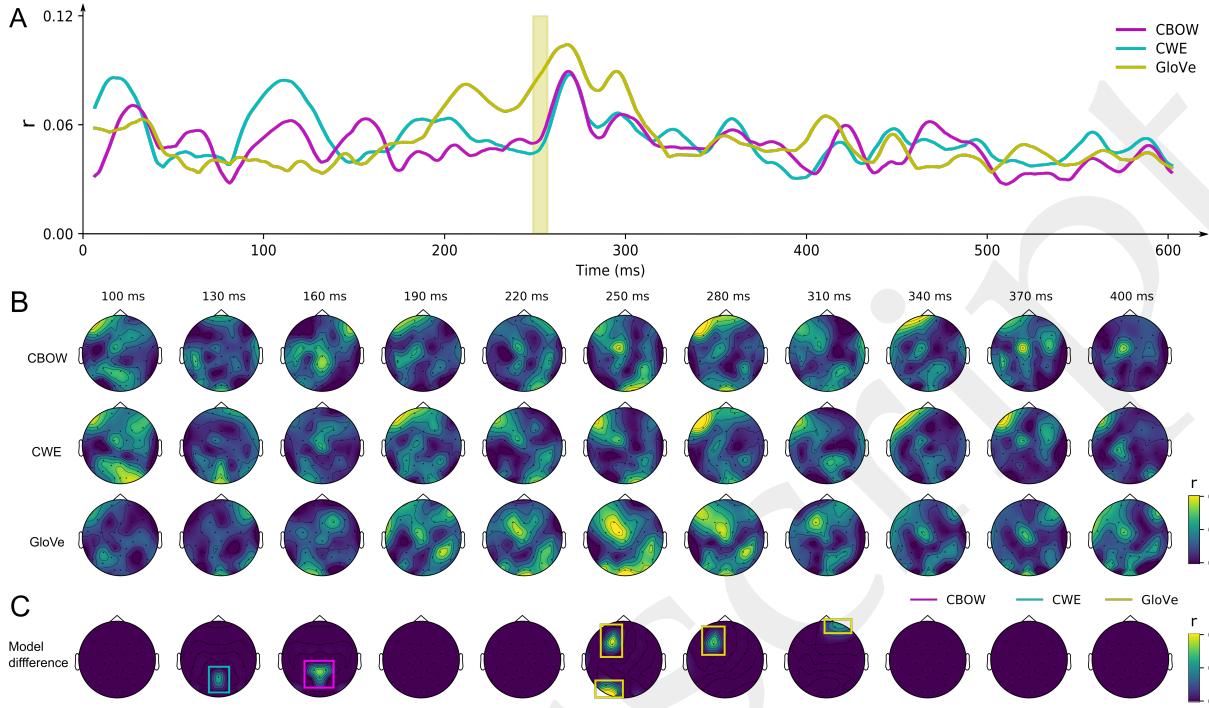298  semantic representations between two complex systems.

13

Figure 5: Three word embedding models distinctively correlate with EEG responses. **A)** The temporal progression of averaged correlations across sensors for each computational model. The correlation for GloVe was significantly better than the other two models between 244 and 251 ms, as highlighted in the shaded window. The significance was determined by non-parametric permutation tests. **B)** The temporal progression of correlation topographies for each computational model. Similar patterns were observed among all models. **C)** The tempo-spatial characteristics of correlation differences among the three computational models. Pairwise non-parametric permutation tests in each sensor revealed distinct predictability at different latencies for each model.

### 3.2.3 Results from trial-level analysis (b): The three NLP models distinctively correlated with EEG responses

We compared the predictability of three selected NLP models (CBOW, GloVe, and CWE) with permutation tests along the temporal progression. Averaged $r$ values across channels in any two of the three models were subject to pairwise comparisons. The results revealed three time windows (lasting for at least 10 ms) within which one model was a significantly better predictor than another one at each time point: (a) CWE predicted significantly better than GloVe between 94 and 122 ms; (b) GloVe predicted significantly better than CWE between 244 and 256 ms; (c) GloVe predicted significantly better than CBOW between 202 and 251 ms. GloVe was a significantly better predictor than the other two models between 244 and 251 ms (yellow shaded area in Fig. 5A).

14

The topographies of $r$ values for all three models were plotted in Fig. 5B, demonstrating that the correlation patterns were spatially consistent among the three models. In particular, high correlations were observed in sensors over the left frontal and occipital regions around 250 to 300 ms, similar as the observation in Fig. 4. The similar spatio-temporal configuration was obtained in permutation tests at channel level, which further revealed that GloVe was the best predictor at sensors over the left frontal and occipital regions around 250 to 300 ms (Fig. 5C). These consistent results in temporal and spatial domains provide strong evidence for the dynamics of semantic processing.

Moreover, CWE was the best predictor around 130 ms in posterior channels. Consistent spatio-temporal configurations for this earlier effect were also observed across all the three models (Fig. 5B). CBOW was the best predictor around 160 ms in posterior channels. Taken together, these results show that the three word embedding models distinctively correlated with ERP differences at distinct latencies.

# 4   Discussion

In this study, we investigated whether and how the lexical semantic representation that independently established in the human brain and computational models share similar formats. We found that semantic similarities computed by word embedding models correlated with EEG semantic priming responses in an early and narrow time window of N400 component. Moreover, distinct word embedding models that include different weighting of orthographic and semantic information correlated with neural responses at perceptual and semantic processing stages. Our study provided strong evidence suggesting that the dynamic processing of lexical semantics can be characterized by word embedding models based on the commonality of semantic representation between two complex systems.

With a better controlled two-word semantic priming paradigm and non-invasive electrophysiological recordings, we provided an analytical approach to collaboratively investigate the semantic representations in two independent complex systems. Computational models can yield quantitative hypothesis to investigate neural processing, and neuroscience data can back-feed to computer models towards creating a stronger artificial intelligence that better emulate neural processes and human behavior. The current study provided a novel framework on how cognitive neuroscience and computer science can be bridged in a bi-directional investigation of the computational

15

342 mechanisms in language research.

343 Computer science can help investigating neuroscientific theories. Granular aspects of
344 linguistic information, such as lexical semantics, can be captured by computational
345 models precisely, without contamination from other factors. Such dedicated and
346 quantitative linking hypothesis between computers and brain provides lens to scrutinize
347 neural computations. The millisecond-by-millisecond single-trial correlational analysis
348 in the current study strikingly narrowed down the time window associated with
349 well-established N400 component that commonly lasts from 250 to 600 ms after a word
350 onset. The observation of significant correlation in a narrow and early time window
351 remarkably reflected the processing of lexical semantics per se. These results can resolve
352 a long lasting debate regarding to one of the most investigated linguistic processing
353 components, N400 – whether it is integration (e.g., Hagoort et al. 2004) or semantic
354 retrieval (Kutas and Federmeier 2011). Our results based on semantic representation
355 extracted independently from computational models suggest that the commonly
356 observed long duration of N400 presumably contains several sub-processes, and
357 semantics-related processing starts at the beginning.

358 Neuroscience can facilitate the journey to strong artificial intelligence. The current
359 study advances in this direction from three aspects. First, neural measures can provide a
360 biomarker for objectively assessing android performance of computational models. The
361 correlations between two complex systems vary as a function of model selections (Fig.
362 5A). The model GloVe correlated with neural data significantly better at around 250 ms
363 than the other two models, suggesting that the implementation of global context yielded
364 more human-like semantic representation. Second, the characteristics of neural
365 dynamics can dissect computational models to probe their features. Distinct models
366 showed better correlations at different latencies (Fig. 5C), suggesting CWE that
367 correlated best at around 130 ms weighted more on lexical-orthographic features,
368 whereas GloVe weighted more on lexical semantics.

369 Third, this study trailblazes a database that will integrate research communities that
370 vary across disciplines, cultures, and societies (https://ray306.github.io/brain_NLP/).
371 The database can help computer scientists to evaluate how human-like their models are
372 and to assess in which aspects the human-like features are. Moreover, the obtained
373 millisecond-level, continuous neural data can help improve model performance and
374 generalize across tasks by optimally integrating the best aspects of models based on
375 dynamic featural processing. Our database (currently only in Mandarin Chinese and

16

English) is expected to expand to many other languages and dialects. We welcome the whole research community to contribute. This joint force will broaden the horizon and provide a unique opportunity to generalize computational models for language processing.

Relating AI models and cognitive neuroscience has brought fruitful findings in other domains of cognition. For example, in vision, the state-of-the-art works by Kriegeskorte's and DiCarlo's groups (Kriegeskorte and Kievit 2013, Khaligh-Razavi and Kriegeskorte 2014, Yamins et al. 2014) have established a mapping between features in different layers of deep neural network model and neural representation in the hierarchical processing in the brain. Our current study was an attempt to create such mapping in the domain of language. Unlike research in vision that can obtain from animal models using invasive methods, linking NLP models and language processing in human brain is constrained by the limits of neuroimaging methods. We carefully chose semantic features and a functional paradigm that can establish direct mapping between computational models and human brain in the linguistic domain. This endeavor opened a brand new door towards a full understanding of computational mechanisms of language processing in both complex systems.

# 5   Conclusion

By investigating the representational formats of comparable lexical semantic features between complex systems with fine temporal resolution neural recordings, we provided a novel framework directly bridging neuroscience and computer science in the domain of language. This framework brought a finer-grained understanding of the neural dynamics underlying semantic processing and developed an objective biomarker for assessing human-like computation in NLP models. Our study suggested a promising way to join forces across disciplines in the investigation of higher-order cognitive functions in human and artificial intelligence.

17

# Acknowledgements

# Conflict of interest

The authors declare no competing financial interests.

# Supplementary materials

Supplementary materials of this study are available at
https://ray306.github.io/brain_NLP/.

18

# References

414 Anderson, John R. 1983. A spreading activation theory of memory. *Journal of verbal*
415 *learning and verbal behavior* 22:261–295.

416 Balota, David A. 1983. Automatic semantic activation and episodic memory encoding.
417 *Journal of verbal learning and verbal behavior* 22:88–104.

418 Bemis, Douglas K, and Liina Pylkkänen. 2011. Simple composition: A
419 magnetoencephalography investigation into the comprehension of minimal linguistic
420 phrases. *Journal of Neuroscience* 31:2801–2814.

421 Bentin, Shlomo, Gregory McCarthy, and Charles C. Wood. 1985. Event-related
422 potentials, lexical decision and semantic priming. *Electroencephalography and clinical*
423 *Neurophysiology* 60:343–355.

424 Brennan, Jonathan. 2016. Naturalistic sentence comprehension in the brain. *Language and*
425 *Linguistics Compass* 10:299–313.

426 Broderick, Michael P, Andrew J Anderson, Giovanni M Di Liberto, Michael J Crosse, and
427 Edmund C Lalor. 2018. Electrophysiological correlates of semantic dissimilarity reflect
428 the comprehension of natural, narrative speech. *Current Biology* 28:803–809.

429 Brown, Colin, and Peter Hagoort. 1993. The processing nature of the N400: Evidence
430 from masked priming. *Journal of cognitive neuroscience* 5:34–44.

431 Chen, Xinxiong, Lei Xu, Zhiyuan Liu, Maosong Sun, and Huanbo Luan. 2015. Joint
432 learning of character and word embeddings. In *Twenty-Fourth International Joint*
433 *Conference on Artificial Intelligence*.

434 Collins, Allan M., and Elizabeth F. Loftus. 1975. A spreading-activation theory of
435 semantic processing. *Psychological review* 82:407.

436 Deacon, Diana, Sean Hewitt, Chien-Ming Yang, and Masanouri Nagata. 2000.
437 Event-related potential indices of semantic priming using masked and unmasked
438 words: evidence that the n400 does not reflect a post-lexical process. *Cognitive Brain*
439 *Research* 9:137–146.

19

Ettinger, Allyson, Naomi Feldman, Philip Resnik, and Colin Phillips. 2016. Modeling n400 amplitude using vector space models of word representation. In *CogSci*.

Federmeier, Kara D, and Marta Kutas. 1999. A rose by any other name: Long-term memory structure and sentence processing. *Journal of memory and Language* 41:469–495.

Firth, John R. 1957. A synopsis of linguistic theory, 1930-1955, *Studies in linguistic analysis*.

Hagoort, Peter, Lea Hald, Marcel Bastiaansen, and Karl Magnus Petersson. 2004. Integration of word meaning and world knowledge in language comprehension. *science* 304:438–441.

Harris, Zellig S. 1954. Distributional structure. *Word* 10:146–162.

Holcomb, Phillip J. 1993. Semantic priming and stimulus degradation: Implications for the role of the n400 in language processing. *Psychophysiology* 30:47–61.

Huth, Alexander G., Wendy A. de Heer, Thomas L. Griffiths, Frédéric E. Theunissen, and Jack L. Gallant. 2016. Natural speech reveals the semantic maps that tile human cerebral cortex. *Nature* 532:453.

Khaligh-Razavi, Seyed-Mahdi, and Nikolaus Kriegeskorte. 2014. Deep supervised, but not unsupervised, models may explain IT cortical representation. *PLoS computational biology* 10:e1003915.

Kiefer, Markus. 2002. The n400 is modulated by unconsciously perceived masked words: Further evidence for an automatic spreading activation account of n400 priming effects. *Cognitive Brain Research* 13:27–39.

Kriegeskorte, Nikolaus, and Rogier A. Kievit. 2013. Representational geometry: Integrating cognition, computation, and the brain. *Trends in cognitive sciences* 17:401–412.

Kutas, Marta, and Kara D. Federmeier. 2011. Thirty years and counting: finding meaning in the N400 component of the event-related brain potential (ERP). *Annual review of psychology* 62:621–647.

Kutas, Marta, and Steven A. Hillyard. 1989. An electrophysiological probe of incidental semantic association. *Journal of Cognitive Neuroscience* 1:38–49.

469 Lehmann, Dietrich, and Wolfgang Skrandies. 1980. Reference-free identification of
470     components of checkerboard-evoked multichannel potential fields.
471     *Electroencephalography and Clinical Neurophysiology* 48:609–621.

472 Lin, Dekang. 1998. An information-theoretic definition of similarity. In *ICML*,
473     volume 98, 296–304. Citeseer.

474 Lorch Jr, Robert F. 1982. Priming and search processes in semantic memory: A test of
475     three models of spreading activation. *Journal of verbal learning and verbal behavior*
476     21:468–492.

477 Mikolov, Tomas, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of
478     word representations in vector space, arxiv preprint arxiv:1301.3781.

479 Miller, George A. 1986. Dictionaries in the mind. *Language and cognitive processes*
480     1:171–185.

481 Neely, James H. 1976. Semantic priming and retrieval from lexical memory: Evidence for
482     facilitatory and inhibitory processes. *Memory & Cognition* 4:648–654.

483 Neely, James H. 1977. Semantic priming and retrieval from lexical memory: Roles of
484     inhibitionless spreading activation and limited-capacity attention. *Journal of*
485     *experimental psychology: general* 106:226.

486 Pennington, Jeffrey, Richard Socher, and Christopher Manning. 2014. GloVe: Global
487     vectors for word representation. In *Proceedings of the 2014 conference on empirical*
488     *methods in natural language processing (EMNLP)*, 1532–1543.

489 Poeppel, David. 2012. The maps problem and the mapping problem: two challenges for
490     a cognitive neuroscience of speech and language. *Cognitive neuropsychology* 29:34–55.

491 Pylkkänen, Liina. 2019. The neural basis of combinatory syntax and semantics. *Science*
492     366:62–66.

493 Resnik, Philip. 1995. Using information content to evaluate semantic similarity in a
494     taxonomy.

495 Resnik, Philip. 1999. Semantic similarity in a taxonomy: An information-based measure
496     and its application to problems of ambiguity in natural language. *Journal of Artificial*
497     *Intelligence Research* 11:95–130.

21

498 Roelofs, Ardi. 1992. A spreading-activation theory of lemma retrieval in speaking.
499     *Cognition* 42:107–142.

500 Turney, Peter D., and Patrick Pantel. 2010. From frequency to meaning: Vector space
501     models of semantics. *Journal of Artificial Intelligence Research* 37:141–188.

502 Wang, Xuefei, Hao Zhu, and Xing Tian. 2019. Revealing the temporal dynamics in
503     non-invasive electrophysiological recordings with topography-based analyses. biorxiv
504     preprint. https://doi.org/10.1101/779546.

505 Yamins, Daniel, Ha Hong, Charles Cadieu, Ethan Solomon, Darren Seibert, and James
506     DiCarlo. 2014. Performance-optimized hierarchical models predict neural responses in
507     higher visual cortex. *Proceedings of the National Academy of Sciences* 111:8619–8624.

508 Yang, Jinbiao, Hao Zhu, and Xing Tian. 2018. Group-level multivariate analysis in
509     EasyEEG toolbox: Examining the temporal dynamics using topographic responses.
510     *Frontiers in Neuroscience* 12:468.

511 Zhang, Linmin, and Liina Pylkkänen. 2015. The interplay of composition and concept
512     specificity in the left anterior temporal lobe: An meg study. *NeuroImage* 111:228–240.

22

# Legends

**Figure 1.** Schematic diagram of semantic representations in the human brain and word embedding models. **A)** A schematic diagram showing how the frequency of words in context yields embedding representations in computational models. Semantically similar words share higher distributional similarity, as illustrated by the counts of neighboring words in the sample mini corpus. Computational models learn semantic representation from words' distribution and generate embedding representations. **B)** A schematic diagram showing the semantic space in the human brain à la Huth et al. (2016). Semantically more similar concepts are represented with more cortical overlaps, indicating shared features. **C)** A schematic diagram showing how the angle between high-dimensional vectors represents semantic similarity in computational models. The angle between two high-dimensional vectors (only two dimensions are used for demonstration) represents semantic similarity. The smaller an angle is (i.e., a higher cosine value), the higher semantic similarity (e.g., the angle between *star* and *moon* is smaller than the one between *star* and *apple*, because *star* and *moon* share more features, as shown in Fig. 1A). **D)** A schematic diagram showing how the amplitudes of neural responses represent semantic similarity in the human brain (e.g., N400, see Kutas and Federmeier 2011). The smaller amplitudes in neural responses to a word are observed when it shares more semantic features with its preceding word (e.g., *star* shares more features with *moon* than with *apple*, as shown in Fig. 1B).

**Figure 2.** Experimental procedure and single-trial correlation analysis. **A)** The trial structure of the experiment. Sample trials are illustrated for the two-word priming paradigm. In each trial, a prime word was followed by a target word. Each word (here 星星 (star)) was used once at the prime position (in the prime trial) and once at the target position (it the target trial). English translations below the screens are for demonstration only, but not included in the expreiment. **B)** Stimuli statistics of semantic similarity generated by the three computational models. **C)** The flowchart of single-trial correlational analysis, (i) Computing the amplitude differences between single-trial EEG responses to the same word at its target vs. prime presentation (target minus prime); (ii) For the 240 word pairs, calculating the correlation between cosine values generated from computational models and amplitude differences from step (i) at each time point in each sensor; (iii) The obtained correlation coefficients form a waveform across time for each

23

545  sensor; (iv) The distribution of correlation coefficients from all sensors is plotted in a

546  topography at each latency.

547  **Figure 3.** Event-related waveform and topographic responses consistent with

548  perceptual and semantic processes in language comprehension. **A)** The waveform

549  responses in a representative channel (Cz). Typical N400 profile was observed in both

550  prime and target. The montage of sensor locations is inserted with the selected channel

551  Cz highlighted. **B)** The dynamics of GFP. The aggregated neural activity across all

552  sensors represented in GFP shows the similar dynamics that has clear perceptual and

553  semantic activation. **C)** The temporal components revealed in the grand averaged ERP

554  responses across targets and primes. Each black line represents ERP responses in each

555  channel. The orange line represents the GFP across all sensors. The vertical dashed lines

556  label the temporal boundaries between ERP components revealed by an automatic

557  segregation method. **D)** The temporal progression of topographies. The topographies

558  for target and prime were represented in the upper and lower rows respectively. Similar

559  topographic patterns and temporal progressions were observed in both target and

560  prime. **E)** The temporal progression of topographic differences. Differences resulted

561  from subtracting prime from target revealed classic N400 topographic patterns from 250
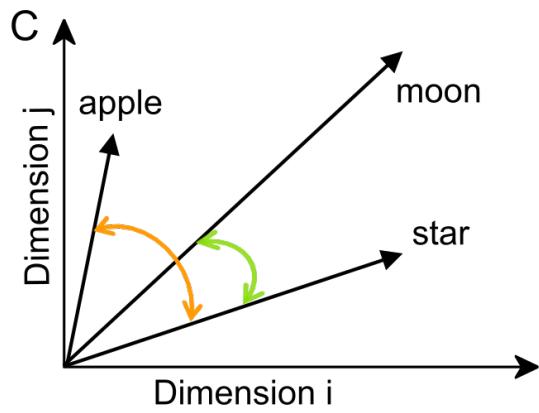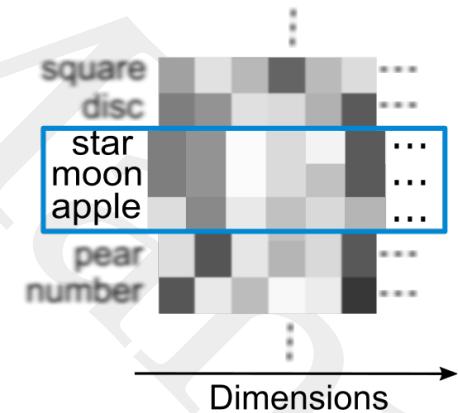
562  to 600 ms.

563  **Figure 4.** Correlations between EEG responses and a word embedding model reveals

564  the dynamics of semantic processing. **A)** Significant correlation was observed between

565  EEG responses in channel Oz at the latency of 300 ms and cosine values computed by the

566  model GloVe. **B)** The temporal progression of correlations (channel Oz). Significant

567  correlations were observed between 226 and 274 ms, between 279 and 306 ms, and

568  between 518 and 529 ms (in red). The significance was determined by the threshold

569  (horizontal line) obtained in a non-parametric permutation test at each time point (alpha

570  level at 0.05). **C)** The spatio-temporal characteristics of correlations. The heatmap of

571  correlations across time and channels revealed significance between 200 and 300 ms in

572  about half of the sensors. The progression of topographies in the time window of

573  significance is zoomed in above. Significant correlations were concentrated in the

574  sensors above the left frontal and tempo-parietal regions.

24

**Figure 5.** Three word embedding models distinctively correlate with EEG responses. **A)** The temporal progression of averaged correlations across sensors for each computational model. The correlation for GloVe was significantly better than the other two models between 244 and 251 ms, as highlighted in the shaded window. The significance was determined by non-parametric permutation tests. **B)** The temporal progression of correlation topographies for each computational model. Similar patterns were observed among all models. **C)** The tempo-spatial characteristics of correlation differences among the three computational models. Pairwise non-parametric permutation tests in each sensor revealed distinct predictability at different latencies for each model.

**Figure 1.** Schematic diagram of semantic representations in the human brain and word embedding models
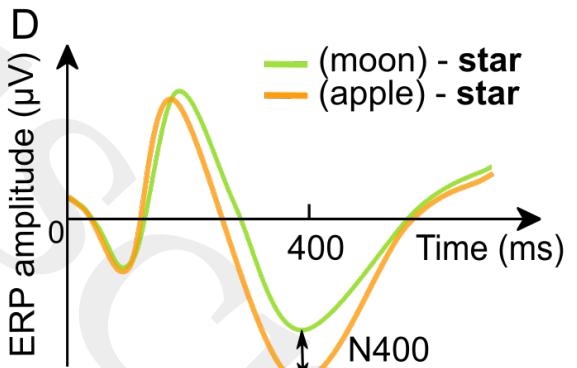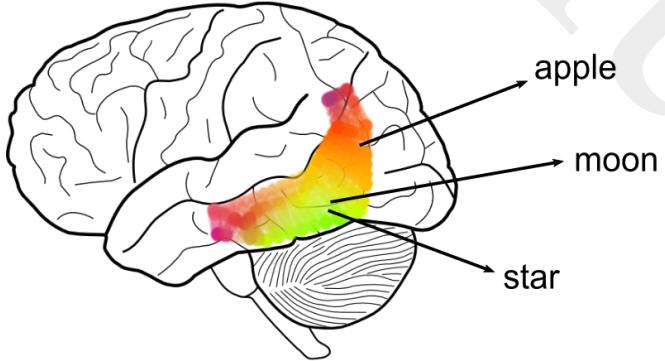
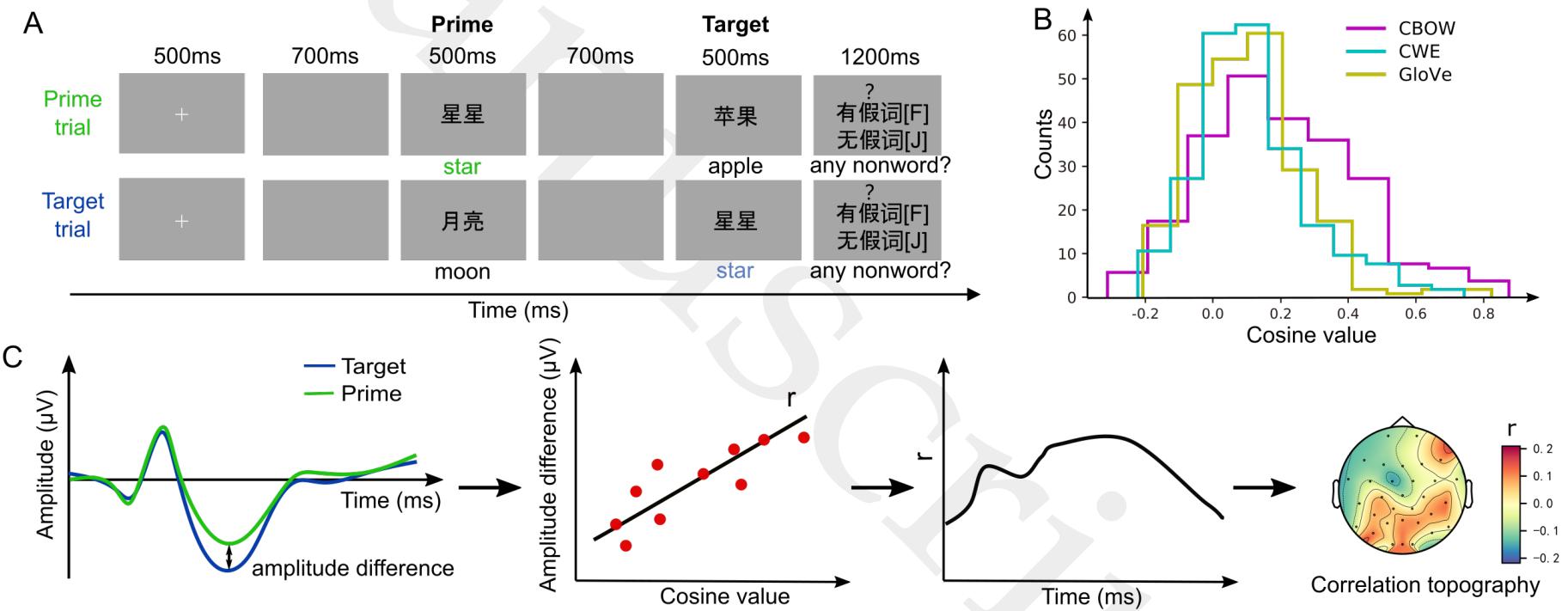**Figure 2.** Experimental procedure and single-trial correlation analysis

**Figure 3.** Event-related waveform and topographic responses consistent with perceptual and semantic processes in language comprehension
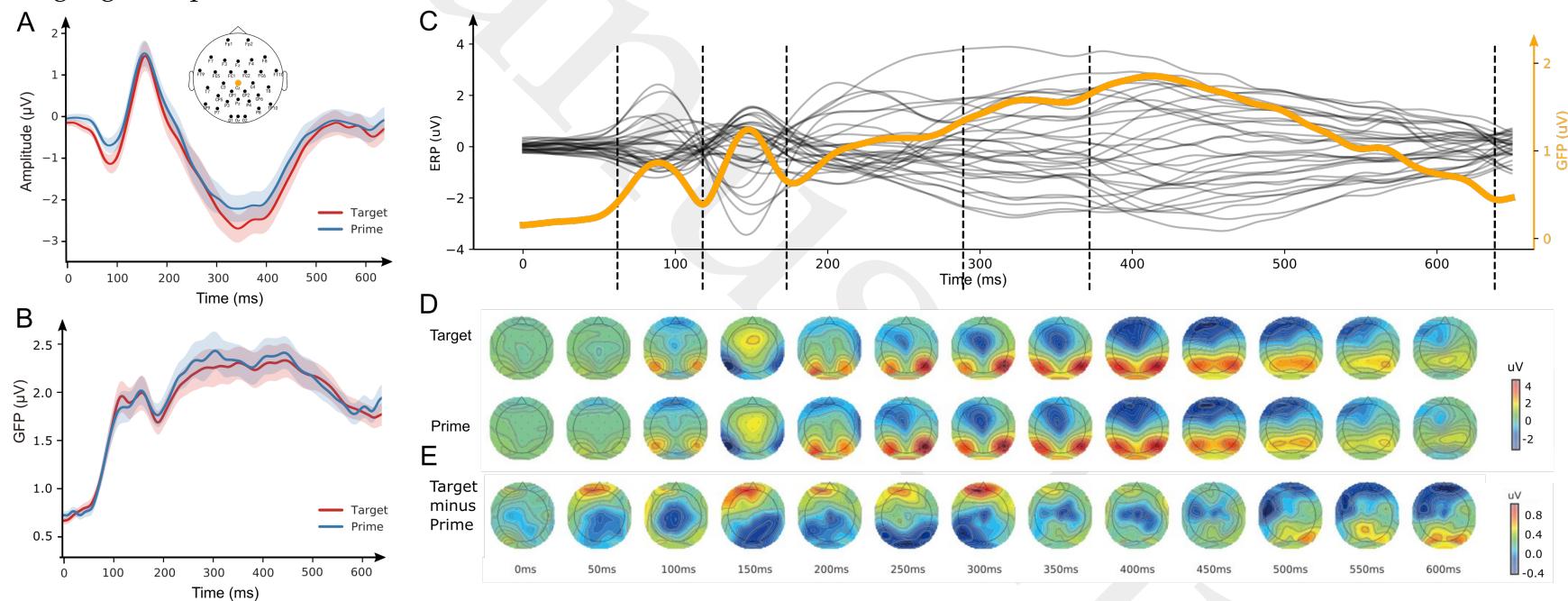
**Figure 4.** Correlations between EEG responses and a word embedding models reveals the dynamics of semantic processing
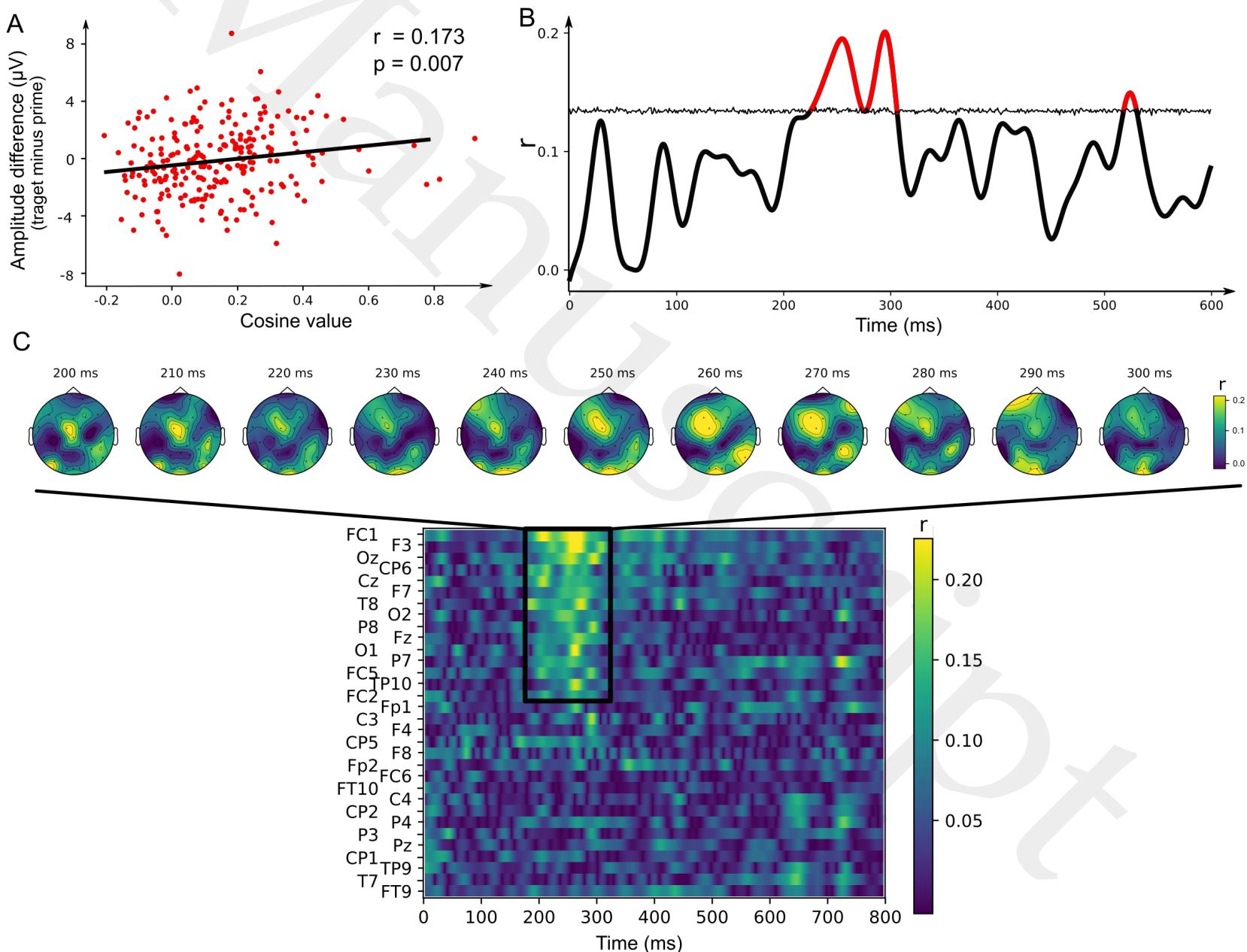
**Figure 5.** Three word embedding models distinctively correlate with EEG responses