# Invariant neural representation of parts of speech in the human brain

Pranav Misra[1,5], Yen-Cheng Shih[2,3], Hsiang-Yu Yu[2,3], Daniel Weisholtz[4], Akshay Sharma[5], Demitre Serletis[5], Juan Bulacio[5], William Bingaman[5], Joseph R Madsen[6], Sceillig Stone[6], Gabriel Kreiman[6,7*]


[1]Harvard University, Cambridge, MA, USA

[2]Department of Neurology, Taipei Veterans General Hospital, Taipei, Taiwan

[3]School of Medicine, National Yang Ming Chiao Tung University College of Medicine, Taipei, Taiwan

[4]Brigham and Women's Hospital, Harvard Medical School, Boston, MA, USA

[5]Cleveland Clinic, Cleveland, OH, USA

[6]Boston Children's Hospital, Harvard Medical School, Boston, MA, USA

[7]Center for Brains, Minds and Machines, Cambridge, MA, USA

*To whom correspondence should be addressed: Gabriel.kreiman@tch.harvard.edu

Number of figures:     5
Number of supplementary figures: 11
Number of supplementary tables: 12


**Abstract**

Elucidating the internal representation of language in the brain has major implications for cognitive science, brain disorders, and artificial intelligence. A pillar of linguistic studies is the notion that words have defined functions, often referred to as parts of speech. Here we recorded invasive neurophysiological responses from 1,801 electrodes in 20 patients with epilepsy while they were presented with two-word phrases consisting of an adjective and a noun. We observed neural signals that distinguished between these two parts of speech. The selective signals were circumscribed within a small region in the left lateral orbitofrontal cortex. The representation of parts of speech showed invariance across visual and auditory presentation modalities, robustness to word properties like length, order, frequency, and semantics, and even generalized across different languages. Furthermore, we extended these ideas by evaluating how parts of speech are processed within full sentences.  Recording activity from additional 1,593 electrodes in 17

34    participants, we found neural signals that separate nouns from verbs in sentences. This selective,

35    invariant, and localized representation of parts of speech provides a foundation to understand

36    how the brain orchestrates more complex aspects of language.

37

38

## Introduction

Language plays a central role in almost all of our daily activities and is at the heart of how we interact with others[1-3]. Early neurological studies and subsequent work using electrical stimulation demonstrated that there exist specific brain regions that play essential roles in language understanding and production[4-9]. Despite the critical importance of language, progress towards elucidating the neural circuits underlying its representation has remained elusive, in part due to the difficulties in investigating animal models, and in part due to the challenges associated with examining the neurophysiological responses in the human brain.

Several neurophysiological experiments have begun to investigate neural signals associated with presentation of individual words or short phrases[10-20]. There has been work examining the orthographic features of real versus pseudowords[21-23], phonetic features of word comprehension[11,21,24,25] and production[26-30], retrieval of semantic information for audio-visual naming[11], and semantic encoding[31]. These studies have shed light on the early processes associated with detecting, comprehending and producing words. Beyond individual words, at the heart of linguistic structures is the notion that words serve specific functions within a sentence, including articles, nouns, adjectives, and verbs. These parts of speech (POS) are widely shared across languages, are combined according to defined grammatical rules, and play critical roles in natural language processing algorithms[1,12,18,19,32-38]. Furthermore, recent work has suggested that POS may be implicitly learned and represented in modern large language models[39,40]. Many studies in patients with brain lesions have shown deficits in the retrieval of individual nouns versus verbs[20,22,41-45]. However, previous studies could not resolve explicit neural circuits associated with POS processing due to insufficient spatial or temporal resolution[20].

What would a representation for parts of speech like nouns and adjectives in the brain look like? Consider the adjective "green" and the noun "apple", combined to create the simple phrase "green apple." Fundamental constraints for such a representation should include the basic invariances underlying the cognitive understanding of this phrase. The basic desiderata for the representation of parts of speech in language includes invariance to: (i) presentation modality (e.g., auditory versus visual), (ii) specific noun or adjective (e.g., green or red), (iii) position within a phrase (e.g., "green apple" versus "apple green"), (iv) specific language in bilingual speakers (e.g., "green apple" in English versus "manzana verde" in Spanish), (v) other word properties like their written length, number of syllables, and phoneme composition.

73

74    Here we set out to investigate the representation of parts of speech in the human brain by
75    recording intracranial field potential responses with high spatiotemporal resolution and high
76    signal-to-noise ratio from large dataset encompassing 3,394 electrodes implanted in 37
77    participants with pharmacologically-resistant epilepsy. We describe neural signals, especially in
78    the left lateral orbitofrontal cortex, that selectively distinguish between nouns and adjectives.
79    These part-of-speech selective signals are robust when words are matched for orthography (e.g.,
80    word length), acoustic features (e.g., number of syllables), word sequence (e.g., noun or adjective
81    at first or second position within a phrase), and frequency of occurrence. Interestingly, the
82    representation of nouns versus adjectives generalizes across audio and visual modalities, across
83    different semantic categories within each part of speech, and across different languages.
84    Furthermore, we extend the work to whole sentences where we show neural signals that
85    distinguish between nouns and verbs.

86

87    **Results**

88

89    We conducted two experiments, one with minimal two-word phrases and a second one
90    with full sentences composed of four words. We start by describing the experiment with two-word
91    phrases and extend the work to full sentences in the last section. We recorded intracranial field
92    potentials from 1,801 electrodes (840 in gray matter, 961 in white matter) implanted in 20
93    participants via stereoelectroencephalography. Participants heard (auditory modality) or read
94    (visual modality) two words that were sequentially presented and were asked to indicate whether
95    the words were the same or not (**Figure 1a**, **Methods**). Participants performed the task correctly
96    on 93.6±7.7% of the trials (here and throughout, mean±std, unless stated otherwise). All electrode
97    locations are shown in **Figure 1b-g** (see also **Tables S1-S2** and **Methods**). We use a bipolar
98    reference, and we focus on the intracranial field potential signals filtered in the high gamma
99    frequency band, referred to as neural responses throughout and reported in the plots as gamma
100   power (65-150 Hz, **Methods**).

101

102   **Neural signals reflect visual, auditory and multimodal inputs**

103

104   We observed 565 electrodes (31.4% of the total) that responded to auditory stimuli (**Figure**
105   **S1a-c, g-i**) and 532 electrodes (29.5% of the total) that responded to visual stimuli (**Figure S1d-**
106   **f, g-i**). The overall proportions and dynamics of visual and auditory responsive signals are

107      consistent with previous work[27,46]. Of these electrodes, there were 293 electrodes that responded

108      to *both* auditory and visual stimuli (**Figure S1g-i**). These 293 electrodes represent 16.3% of the

109      total, 51.9% of the auditory responsive electrodes, and 55.0% of the visually responsive

110      electrodes. This number of audiovisual electrodes is unlikely to arise by chance from the number

111      of auditory and visual electrodes ($p<10^{-4}$, permutation test, $n=10^6$ iterations). Of these 293

112      electrodes, 147 (50.2%) were in the left hemisphere and 146 (49.8%) were in the right

113      hemisphere. Of the 41 the regions in the Desikan-Killiani Atlas where we had sampling (34 defined

114      regions and 7 extra regions representing deep gray matter structures, **Methods**, **Figure 1**, **Tables**

115      **S1-S2**), 13 regions had a significantly higher number of multimodal electrodes than expected from

116      the number of audio or visual electrodes ($p<0.01$, permutation test, $n=10^6$ iterations). These

117      regions are indicated in bold in **Table S2**. **Figure 1h-j** shows the responses of an example

118      audiovisual responsive electrode located in the left rostral middle-frontal gyrus (**Figure 1k**). This

119      electrode showed strong evoked responses evident in the trial-averaged responses (**Figure 1h**),

120      and even in individual trials for both auditory stimuli (**Figure 1i**) and visual stimuli (**Figure 1j**).

121

122      To compare the response dynamics of auditory and visual responses, we calculated the

123      time at which the neural signals reached half of the max amplitude (half-maximum time, arrows

124      in **Figure 1h**, **Methods**) and the average area under the curve (AUC) for neural responses such

125      as those in **Figure 1h**. **Figure S2a** shows the half-maximum time for auditory-only electrodes

126      (left), visual-only electrodes (middle), and audiovisual electrodes on audio trials (right light-gray

127      half) or visual trials (right black half). There was no significant difference between the half-

128      maximum time for auditory-only electrodes ($329\pm187$ ms) and visual only electrodes ($336\pm174$

129      ms) ($p>0.05$, ranksum test). Similarly, there was no significant difference between the half-

130      maximum time for the audio and visual responses of audiovisual electrodes ($379\pm193$ ms versus

131      $341\pm174$ ms, $p>0.05$, ranksum test). However, there was a small but significant difference

132      between the half-maximum time for audio only electrodes and auditory responses of audiovisual

133      electrodes ($p<0.01$, ranksum test).

134

135      As expected, for the audio-only electrodes, the response AUC to auditory stimuli ($108\pm100$

136      $\mu V^2/Hz\text{-}ms$) was larger than to visual stimuli ($44\pm16$ $\mu V^2/Hz\text{-}ms$) ($p<10^{-4}$, ranksum test, **Figure**

137      **S2b**). Similarly, for the visual-only electrodes, the response AUC to auditory stimuli ($40\pm23$

138      $\mu V^2/Hz\text{-}ms$) was smaller than to visual stimuli ($53\pm43$ $\mu V^2/Hz\text{-}ms$) ($p<10^{-4}$, ranksum test, **Figure**

139    **S2c**). For the audiovisual electrodes, the response AUC to auditory stimuli ($71\pm72$ $\mu V^2/\text{Hz-ms}$)

140    was slightly larger than to visual stimuli ($54\pm39$ $\mu V^2/\text{Hz-ms}$) (p<0.01, ranksum test, **Figure S2d**).

141

142    **Multimodal neural signals distinguish different parts of speech**

143

144        We evaluated whether the neural signals differentiated between nouns and adjectives.

145    Nouns and adjectives were matched for their number of syllables and word length to control for

146    potential confounds not specific to parts of speech (**Table S3**, **Methods**). **Figure 2** shows the

147    responses of an example electrode located in the orbital H-shaped sulcus within the left lateral

148    orbitofrontal cortex (**Figure 2i** depicts the electrode location). The orbital H-shaped sulcus lies

149    above the bone of the eye socket where a butterfly-like gyrus can be seen, formed along H-

150    shaped recessions of the sulcus. The neural responses are aligned to the word onset (vertical

151    dashed line) for auditory presentation (**Figure 2a, b**) or visual presentation (**Figure 2c, d**), for the

152    first (**Figure 2a, c**), or second (**Figure 2b, d**) word in each trial. This electrode showed multimodal

153    responses triggered by both auditory and visual stimuli. The responses to nouns (blue) were

154    stronger than adjectives (red) across all four conditions, including both word 1 and word 2, and

155    both for visual and auditory stimuli. The differences between nouns and adjectives can be readily

156    appreciated even in individual trials (**Figures 2e-h**). These differences became significant at

157    approximately 430 ms after word onset for visual presentation and about 610 ms for auditory

158    presentation.

159

160        In all, there were 89 electrodes, 97 electrodes, and 48 electrodes that showed a difference

161    between nouns and adjectives for auditory stimuli only, visual stimuli only, or both modalities,

162    respectively. The 48 electrodes cannot be ascribed to randomly sampling from the total of audio

163    and visual electrodes ($p<10^{-4}$, permutation test, $n=10^6$ iterations).

164

165        **Neural selectivity for nouns versus adjectives was robust to word properties,**

166    **phrase grammar, usage frequency, and word subcategory**

167

168        Even though nouns and adjectives were matched in their average number of syllables and

169    word length, we asked whether these variables could still contribute to the neural responses

170    differentiating nouns and adjectives. Additionally, each trial could be grammatically correct (e.g.,

171    "green apple"), or incorrect (e.g., "apple green") (**Methods**); therefore, we asked whether

172    grammar could contribute to the neural differences between nouns and adjectives. To address

173    these questions, we built a generalized linear model (GLM) for each electrode to predict its

174    response AUC between 200 ms and 800 ms after word onset using four predictors: nouns versus

175    adjectives, grammatically correct or not, word length (vision) or number of syllables (audition)

176    (**Methods**). The predictor coefficients in the GLM model for the example electrode in **Figure 2a-**

177    **d** show that only the nouns versus adjectives label significantly explained the neural responses

178    for both auditory and visual presentation (**Figure 2j**). A total of 14 electrodes showed nouns

179    versus adjectives as the *only* statistically significant predictor in the GLM analysis; 13/14 (93%)

180    of these electrodes distinguished nouns versus adjectives for both auditory and visual inputs,

181    such as the example electrode in **Figure 2a-j**.

182

183         The locations of electrodes that robustly distinguished nouns and adjectives (orange in

184    **Figure 2k**) revealed a cluster enriched in the left lateral orbitofrontal cortex (LOF). Within the left

185    LOF, 8 out of the 8 (100%) electrodes were in the posterior part of the orbital H-shaped sulcus.

186    We recorded from a total of 113 electrodes in the lateral orbitofrontal region, 38 electrodes in the

187    left hemisphere and 75 electrodes in the right hemisphere (**Figure 1b-g, Table S1**). Of the 38 left

188    hemisphere electrodes, 21% distinguished nouns from adjectives during both audio and visual

189    presentation. In stark contrast, only 1.3% of the 75 electrodes in the right hemisphere

190    distinguished nouns from adjectives in both audio and vision (these hemispheric differences were

191    statistically significant: $p<10^{-4}$, permutation test, $n=10^6$ iterations). **Table S4** shows the distribution

192    of electrodes distinguishing part of speech between the left and right hemispheres for all brain

193    regions and **Table S5** shows the distribution of electrodes separating nouns versus adjectives in

194    different participants.

195

196         We had initially assumed that distinguishing parts of speech constitutes a core component

197    of language and would therefore be reflected *exclusively in both* visual and auditory modalities.

198    Indeed, 13/14 (93%) of electrodes differentiating nouns from adjectives in the GLM did so in both

199    modalities. In addition to these 13 electrodes there was a small number of electrodes (2 auditory

200    only and 1 visual only) that showed differences between nouns and adjectives in one modality

201    but not the other. Unlike the electrodes in **Figure 2k**, for the 2 auditory-only electrodes, the

202    number of syllables also significantly contributed towards explaining the neural responses. **Figure**

203    **S3** shows the responses of an example electrode located in the right insula that showed a

204    difference between nouns and adjectives during auditory presentation but *not* during visual

205    presentation. Conversely, **Figure S4** shows the responses of an example electrode located in the

206    left lateral orbitofrontal cortex that showed a clear difference between nouns and adjectives during

207 visual presentation but *not* during auditory presentation. **Figure S4 k,l** shows the locations of
208 auditory only (white circles) and visual only electrodes (black circle) in the left and the right
209 hemispheres, respectively.

210

211       Nouns and adjectives differ in their usage frequency. We asked whether the differences
212 in the neural responses to nouns versus adjectives depended on usage frequency. To address
213 this question, we randomly subsampled the trials to match the distribution of Google Ngram
214 frequency (**Methods**). **Figure S5a** shows matched noun and adjective distributions for the
215 example electrode shown in **Figure 2a-k.** This electrode showed differential responses between
216 parts of speech for auditory (**Figure S5b,c**) and visual (**Figure S5d,e**) stimuli during word1
217 (**Figure S5b,d**) and word2 (**Figure S5c,e**), even after nouns and adjectives were matched for
218 their frequency of occurrence. Of the 13 audiovisual electrodes where nouns versus adjectives
219 was the only significant predictor in the GLM analysis, 6 electrodes (43%, 4 in the left-LOF, and
220 2 in left superior temporal gyrus) robustly distinguished nouns and adjectives matched for their
221 frequency of occurrence, like the example electrode in **Figures 2** and **S5** whereas the other
222 electrodes maintained their selectivity in most but not all conditions.

223

224       Within our stimulus set, there were two subcategories of nouns, animals and food, and
225 there were two subcategories of adjectives, concrete and abstract (**Table S3**). We asked whether
226 the electrodes that showed differential responses generalized across different word
227 subcategories. The example electrode in **Figure 2a-j** did not show differences between the two
228 noun or adjective subcategories for either auditory stimuli (**Figure S6a, b, f, g**), visual stimuli
229 (**Figure S6c, d**, **h**, **i**), word 1 (**Figure S6a**, **c**, **f**, **h**), or word 2 (**Figure S6b**, **d**, **g**, **i**). Of the 13
230 audiovisual electrodes where nouns versus adjectives was the only significant predictor in the
231 GLM analysis, 8 electrodes (62%) showed generalization across different noun or adjective
232 subcategories. The remaining 6 electrodes (38%) showed a significant difference between the
233 two noun subcategories or between the two adjective subcategories (**Table S5**). **Figure S7** shows
234 one of the exceptions, i.e., an electrode in the left LOF which showed a significant response only
235 for food nouns. This selectivity was particularly pronounced for the visual stimuli (**Figure S7c**, **d**,
236 **h**, **i**), but was also apparent for auditory stimuli (**Figure S7a**, **b**, **f**, **g**), and was evident both for
237 word 1 and word 2.

238

239       In sum, differences in selective responses to nouns versus adjectives were particularly
240 prominent and clustered in the left lateral orbitofrontal cortex, persisted across different word

241 lengths, whether the word was used in a grammatically correct phrase or not, after equalizing
242 word occurrence frequency, and generalized across different noun or adjective subcategories.
243
244 **Neural signals enhanced for nouns versus adjectives were anatomically segregated**
245
246     Of those electrodes uniquely selective for part of speech, 77% showed responses that
247 were significantly stronger for nouns compared to adjectives ($\beta_{NvsA} > 0$) as illustrated by the
248 example in **Figure 2a-j**. The remaining 23% showed responses that were stronger for adjectives
249 compared to nouns ($\beta_{NvsA} < 0$) as illustrated by the example in **Figure S8 a-i** (**Table S5**). For
250 auditory stimuli, the difference in the onset time between nouns and adjectives was larger for
251 noun-preferring electrodes (550 ± 107 ms) than adjective-preferring electrodes (312 ± 94 ms,
252 ranksum test, p<0.05). For visual stimuli, the difference in the onset time between nouns and
253 adjectives was not different between noun-preferring electrodes (425 ± 107 ms) and adjective-
254 preferring electrodes (437 ± 134 ms, ranksum test, p>0.05). There was a significant correlation
255 between auditory and visual difference onset times for noun-preferring electrodes (Pearson $R^2$ =
256 0.80, p<0.01) but not for adjective-preferring electrodes (Pearson $R^2$ = -0.70, p>0.05).
257
258     When we displayed the electrode locations on the brain, we observed an anatomical
259 separation between these two groups of responses (**Figure 2l,m,** x-axis: lateral to medial, y-axis:
260 anterior to posterior, z-axis: ventral to dorsal). We compared noun- versus adjective- preferring
261 electrodes along 3 axes of Montreal Neurological Institute 305 Coordinates (MNI305, units
262 abbreviated as m.u.)[47]. Along the lateral to medial axis (x-axis in **Figure 2l,m,** zero being more
263 medial), noun-preferring electrodes had a mean of 25.3±6.2 m.u. and adjective-preferring
264 electrodes had a mean of 47.3±7.7 m.u. (p<0.01, ranksum test). Along the ventral-dorsal axis (z-
265 axis in **Figure 2l**), noun electrodes had a mean of -12.17±5.3 m.u. and adjective electrodes had
266 a mean of -3.7±1.7 m.u. (p<0.05, ranksum test). Along the posterior-anterior axis (y-axis in **Figure
267 2m**), noun electrodes had a mean of 21.4±18.9 m.u. and adjective electrodes had a mean of -
268 2.7±25.8 m.u. (p<0.05, ranksum test). **Table S6** summarizes the locations of noun- vs adjective-
269 preferring electrodes across brain regions. A permutation test combining all brain regions for
270 these electrodes showed that that electrodes in the LOF tended to show stronger responses to
271 nouns (~90% $\beta_{NvsA} > 0$, p<$10^{-4}$, permutation test, n=$10^6$ iterations, **Methods**).
272
273 **A population of electrodes in the lateral orbitofrontal cortex can distinguish nouns from
274 adjectives in individual trials and generalizes across words and modalities**

275

To assess whether information about part of speech was available in individual trials, we used a machine learning pseudopopulation approach by combining electrodes within anatomically defined brain regions in the Desikan-Killiany Atlas[48]. We binned the response in 100 ms time bins and used the top-N principal components that explained more than 70% of the variance in the training data for all the electrodes. We trained an SVM classifier with a linear kernel to distinguish between nouns and adjectives and tested the classifier on held-out data (**Methods**). **Figure 3** shows decoding accuracy for the left (**Figure 3a,d,g**) and the right (**Figure 3b,e,h**) LOF as a function of time from word onset. When trained using data from both word1 and word2 with combined auditory and visual features, there was a statistically significant decoding performance starting approximately at ~300 ms after word onset and reaching a peak of 63.6±1.1% at ~500 ms after word onset in the left LOF (**Figure 3a**). Statistical significance was assessed by comparing with a control where noun and adjective labels were randomly shuffled (**Methods**). Even though there were almost twice as many electrodes in the right LOF compared to the left LOF (**Table S2, Figure 1b-g**), decoding performance was much higher for the left LOF compared to the right LOF (compare **Figure 3a** versus **Figure 3b**). The differences between the left and right LOF persisted after randomly subsampling to equalize the number of electrodes across hemispheres for all regions (**Figure S9a,b**).

In **Figure 3a,b**, word 1 and word 2 are combined. Decoding performance in the left LOF was also high when separately considering word 1 (**Figure S10a-c**) and word 2 (**Figure S10d-f**). Furthermore, the machine learning classifier was able to generalize across words, as evidenced by the decoding performance when training on the responses to word 1 and testing on the responses to word 2 (**Figure 3d,e**), and vice versa (**Figure 3g,h**). Similarly, auditory and visual trials are combined in **Figure3a,b**. Decoding performance in the left LOF was also high when separately considering auditory trials (**Figure S10g-i**) and visual trials (**Figure S10j-l**). Furthermore, the machine learning classifier was able to generalize across modalities as evidenced by the decoding performance when training on auditory trials and testing on vision trials (**Figure S10m-o**) and vice versa (**Figure S10p-r**).

We extended the analyses in **Figure 3a,b,d,e,g,h** to all other regions in the Desikan-Killiany atlas. In addition to the left LOF, the left superior temporal cortex and the left fusiform cortex also showed statistically significant decoding performance (**Figure 3c**). However, in contrast to the results for the left LOF, the decoding results for other regions were less robust

309    (**Figure S9c**) and did not generalize across words (**Figure 3f,i**) or across modalities (**Figure**
310    **S10o,r**).

311

312    **Multimodal neural signals distinguishing different parts of speech are conserved across**
313    **languages**

314

315    One of the participants was fluent in two languages, English and Spanish. Therefore, this
316    patient provided an opportunity to ask whether the neural signals discriminating between different
317    parts of speech were language-specific or showed invariance across languages. All the words
318    were translated into Spanish by a native Spanish speaker and the task was repeated in both
319    languages. **Figure 4a-h** shows the responses of an example electrode located in the left LOF
320    (**Figure 4k**). This electrode showed a stronger response to nouns compared to adjectives for
321    auditory stimuli (**Figure 4a**, **b**, **e**, **f**), for visual stimuli (**Figure 4c**, **d**, **g**, **h**), for Word 1 (**Figure 4a**,
322    **c**, **e**, **g**), and for Word 2 (**Figure 4b**, **d**, **f**, **h**). Interestingly, the separation between nouns and
323    adjectives was evident both when the words were presented in English (**Figure 4a-d**) and when
324    the words were presented in Spanish (**Figure 4e-h**). The GLM analysis showed that nouns versus
325    adjectives was the only significant predictor in English trials (**Figure 4i**), and Spanish trials (**Figure**
326    **4j**). All in all, there were three electrodes in this participant that showed a multimodal response
327    selective for part of speech. All three of these electrodes were in the left orbital H-shaped sulcus
328    within the LOF (**Figure 4k,** green).

329

330    In addition to this bilingual participant, the task was run in monolingual participants who
331    spoke English (n=16 participants) and monolingual participants who spoke Taiwanese (n=3
332    participants, **Table S1**). In **Figure 4k**, we show all electrodes from the left LOF that showed part-
333    of-speech encoding from different participants (**Table S7**). We also indicate the language in which
334    this difference was observed whether it be English (pink), Taiwanese (brown) or bilingual
335    English/Spanish (green). All participants in **Figure 4k** were right-handed. Electrodes separating
336    parts of speech from monolingual participants were also clustered in the same region. Thus, the
337    left LOF distinguished between parts of speech for both auditory and visual presentations of
338    stimuli across participants speaking different languages.

339

340    **Multimodal neural signals distinguished nouns and verbs in full sentences**

341

342    The experiment presented thus far concerned the responses to nouns and adjectives
343    within minimal phrases. We extended these results in two ways: (1) by evaluating whether there
344    are multimodal signals that distinguish between nouns and verbs; (2) by evaluating the neural
345    signals to words embedded within full sentences. We recorded intracranial field potentials from
346    1,563 electrodes (844 in gray matter, 719 in white matter) implanted in 17 patients via
347    stereoelectroencephalography. Participants heard (auditory modality) or read (visual modality)
348    four-word sentences that were sequentially presented (**Figure 5a**, **Methods**). To assess
349    comprehension, participants were asked to indicate whether the sentence adequately described
350    an image that followed the last word after a 1,000 ms interval. Participants performed the task
351    correctly on 85.7±14.3% of the trials. We considered two types of sentences, semantic (e.g., "the
352    girls ate cakes") or non-semantic (e.g., "the cakes ate girls"). All electrode locations are shown in
353    **Figure S11 (**see also **Table S11, Methods).**
354
355    Following the procedures described in the analyses of neural responses to nouns versus
356    adjectives, we evaluated whether neural signals differentiated between nouns and verbs. **Figure**
357    **5** shows the responses of an example electrode located in the pars triangularis (**Figure 5f** denotes
358    the electrode location). The neural responses are aligned to word onset for auditory presentation
359    (**Figure 5b**) or visual presentation (**Figure 5c**). The responses to nouns (blue) were stronger than
360    verbs (black) for auditory and visual stimuli. The differences between nouns and verbs can be
361    readily appreciated even in individual trials (**Figures 5d,e**). These differences became significant
362    at approximately 140 ms after word onset for auditory presentation and about 320 ms for visual
363    presentation. In all, there were 121 electrodes that showed selective responses distinguishing
364    nouns from verbs both for auditory and visual presentation.
365
366    Even though we tested these electrodes for nouns versus verbs differences, it is possible
367    that auditory features (like number of syllables) or orthographic features (like word length) could
368    contribute to the neural responses. Further, each sentence could either be semantic (S, e.g., "the
369    girls ate cakes") or not (NS, e.g., "the cakes ate girls). To evaluate whether word features and
370    semantic features contributed to the neural signals underlying parts of speech, we built a GLM for
371    each electrode to predicts its response AUC between 200 ms and 800 ms after word onset using
372    four predictors: nouns versus verbs, semantic or not, number of syllables, and word length
373    (**Methods**). The predictor coefficients in the GLM model for the example electrode in **Figure 5b-**
374    **e** show that only the nouns versus verbs label significantly explained the neural responses
375    (**Figure 5g**). A total of 41 audiovisual electrodes showed nouns versus verbs as the *only*

376      statistically significant predictor in the GLM analysis, such as the example electrode **Figure 5b-**

377      **g**. The locations of these electrodes are shown in **Figure 5h,i.** The electrode locations reveal two

378      clusters enriched in the left pars triangularis and precentral regions. The difference in the number

379      of significant electrodes between the right and left hemispheres was statistically significant: $p<10^{-4}$

380      , permutation test, $n=10^6$ iterations, see **Table S12**).

381

382          Many electrodes (63%) showed responses that were significantly stronger for nouns

383      compared to adjectives ($\beta_{NvsV} > 0$), as illustrated by the example in **Figure 5b-e**. We observed an

384      anatomical separation between these two groups of responses (**Figure 5j,k,** x-axis: lateral to

385      medial, y-axis: anterior to posterior, z-axis: ventral to dorsal). We compared noun- versus verb-

386      preferring electrodes along 3 axes of Montreal Neurological Institute 305 Coordinates (MNI305,

387      units abbreviated as m.u.)[47]. Along the anterior-posterior axis (y-axis in **Figure 5j,k**), noun

388      electrodes had a mean of -11.8±21.8 m.u. and verb electrodes had a mean of 10.1±26.4 m.u.

389      (p<0.01, ranksum test). Along the ventral-dorsal axis (z-axis in **Figure 5j**), noun electrodes had a

390      mean of -5.1±29.5 m.u. and verb electrodes had a mean of 15.9±27.0 m.u. (p<0.05, ranksum

391      test). Along the lateral to medial axis (x-axis in **Figure 5k,** zero being more medial), noun-

392      preferring electrodes had a mean of 38.3±16 m.u. and adjective-preferring electrodes had a mean

393      of 40.3±11.6 m.u. (not significant, p>0.05, ranksum test).

394

395      **Discussion**

396

397          We described neurophysiological signals that selectively discriminate between two

398      parts of speech, nouns and adjectives (**Figure 2**). This selectivity was robust to

399      orthographic variables such as word length, phonetic features such as number of

400      syllables, and word occurrence statistics (**Figure 2**). This selectivity for part of speech

401      generalized across sensory modalities (**Figures 2**, **3**, **4**), word positions, grammatical

402      correctness and motor outputs (**Figures 2**, **3**, **4**), and semantic groups of nouns and

403      adjectives (**Figure S6**). These neurophysiological signals enable discrimination between

404      parts of speech even in single trials (**Figures 2, 3**). Electrodes that uniquely distinguished

405      nouns from adjectives were clustered within a small, circumscribed region of the lateral

406      orbitofrontal cortex, lateralized to the left hemisphere (**Figure 2**, **4**). Neural discrimination

407      of nouns from adjectives was apparent in the LOF in English-speaking and Taiwanese-

408    speaking participants (**Figure 4**). Interestingly, in a bilingual participant, the same

409    electrodes within the left LOF distinguished nouns and adjectives in both English and in

410    Spanish (**Figure 4**). Extending the study of minimal phrases, we conducted an additional

411    experiment where we showed neural signals that distinguished nouns from verbs within

412    full sentences (**Figures 5**).

413

414        In English and other languages, some words can be used both as a noun or as an

415    adjective (e.g., long *race* versus *race* horse). In most instances, one usage is more

416    frequent than the other. In particular, the nouns and adjectives in this study are highly

417    overrepresented in their labeled part of speech (**Table S8**). Similarly, some words can be

418    used both as a noun or as a verb (e.g., "long *race*" versus "*race* you to the top"); all the

419    nouns in this study are highly overrepresented in their usage as nouns or verbs (**Table

420    S8**). Thus, the words used in this study had a prototypical interpretation as either, noun,

421    adjective, or verb. The distinction between POS includes their grammatical roles but also

422    their associated semantic connotations (e.g., nouns typically refer to things and adjectives

423    to the attributes of those things).

424

425        In languages like English, nouns and adjectives follow a specific grammatical order

426    (i.e., adjectives precede nouns). Other languages reverse this order. In Spanish,

427    adjectives typically follow nouns, though the English order can also be used. It is thus

428    interesting to observe that many electrodes demonstrated strong selectivity for nouns

429    versus adjectives, irrespective of their position within the two-word phrases. Furthermore,

430    in the bilingual participant, the neural responses separated nouns and adjectives in both

431    languages despite the fact that the grammatical order is typically reversed between

432    English and Spanish. It is conceivable that the strong part-of-speech selectivity

433    independent of grammar shown here could be linked to the two-word phrase structures.

434    Another possibility is that the representation of nouns versus adjectives is invariant to

435    grammatical usage rules. The results in **Figure 5** demonstrates a selective representation

436    of parts of speech that extends to full sentences, invariant to changes in semantics.

437

438    Non-invasive scalp electroencephalography and magnetoencephalography
439    signals have revealed correlates of language processing with a wide range of onset times
440    from approximately 100 ms all the way to well over 600 ms (for a review, see[49]). The
441    earliest onset signals commencing between 100 and 300 ms after stimulus onset,
442    sometimes referred to as early left anterior negativity, have been associated with
443    grammatical violations, but previous studies have not documented any invariance in the
444    representation of parts of speech and there is disagreement about whether these early
445    signals are even associated with language[49]. Our work reports an invariant distinction
446    between nouns and adjectives in the LOF commencing at approximately 400 ms after
447    stimulus onset, which is consistent with part-of-speech being represented well after the
448    onset of modality-specific purely visual and auditory signals.

449

450    A remarkable hallmark of language is its universality. We can interpret the word
451    *cat* when uttering the word, writing it, listening to it, reading it, and even when examining
452    a photograph of a cat. It is therefore tempting to speculate that there may be an invariant
453    representation of language concepts in the brain. Several studies have examined putative
454    correlates of language processing using only unimodal signals (e.g.,[12-16,21,27,32,37,38]).
455    While we observed electrodes that distinguished between parts of speech only in the
456    auditory stimuli or only in the visual stimuli, the responses of those electrodes could be
457    partly explained by other variables including number of syllables, word frequency, or
458    grammar. Using strict criteria and after controlling for confounding variables, most
459    electrodes that distinguished nouns from adjectives showed selectivity during both
460    auditory and visual presentation. Future work should evaluate whether the same
461    electrodes also distinguish parts of speech when participants utter words, write them, or
462    when examining images. An intriguing study described neurons in the human medial
463    temporal lobe that respond selectively to images and their corresponding text and sound
464    descriptions[50,51]. However, these medial temporal lobe neurons do not seem to
465    distinguish between different parts of speech and their responses seem to be connected
466    with the formation of memories rather than the internal representation of language[52].
467    Indeed, there exist strong anatomical and functional connections between the medial
468    temporal lobe and frontal regions that could link language and memory formation[53].

469

The lateral orbitofrontal (LOF) cortex constitutes a large expanse of neocortex within the frontal lobe, spanning Brodmann areas (BA) 10, 11, 12 (called BA47 in humans due to cytoarchitectural differences from monkeys) and 13[54-56]. Neurobiological tracings from rats, mice, and macaques have identified LOF as a nexus of many inputs[56] conveying olfactory, gustatory, visual, auditory, somatosensory, and visceral-sensory information. The LOF has been associated with many cognitive functions, including multisensory integration, working memory, long-term memory consolidation, reward processing, social interactions, memory, decision making, and emotion processing[53,55,57-61]. This heterogeneity might be partly ascribed to investigations probing different cognitive tasks, as in the case of the proverbial blind men sampling different parts of an elephant. Given the prominent role of language in cognition, it is conceivable that previous studies that describe other roles of the LOF did not probe its possible associations with language. However, it is even more likely that descriptors like LOF that refer to such large brain areas would inevitably fail to uncover specific functionality. The current results point to a rather well circumscribed location within LOF, the posterior part of the H-shaped sulcus in the left hemisphere. In humans, this location overlaps with BA13-lateral and BA 47-medial and has been shown to have a strong convergence of auditory and visual inputs[55,56,62]. Interestingly, work on Primary Progressive Aphasia, and frontotemporal lesions implicate the orbitofrontal cortex in word and sentence comprehension deficits[8,62-64] (see also[31]). In these studies, the orbitofrontal cortex, dorsal premotor cortex, temporoparietal junction (canonical Wernicke's area), and pars opercularis were associated with sentence comprehension and grammatical production aphasias (evaluated with complex grammatical output requiring planning and motor production). Word comprehension and naming deficits were assessed using binary perceptual choice tasks, implicating the orbitofrontal cortex and the anterior temporal lobe (ATL). Consistent with extensive work documenting the lateralization of language functions, the results presented here also show a strong predominance of the left hemisphere in the representation of part of speech, despite the fact that there were more electrodes sampling signals from the right hemisphere.

499

500    Several limitations in the current work are worth noting. First, all the results
501    reported here are derived from patients with epilepsy. The invasive study of epilepsy
502    patients constitutes the predominant way to access neurophysiological signals from the
503    human brain[65,66]. Neurophysiological studies in other patient populations (e.g., paraplegic
504    patients, Parkinson's patients, brain tumor patients), typically target specific regions that
505    are not known to be associated with language processing.  Caution should be exercised
506    in the interpretation of results from patient populations. To the best of our knowledge, all
507    patients used language fluently and had no language impediments, but one should be
508    aware of the possibility that epilepsy could potentially impact the representation of
509    language. Second, the electrode locations are strictly dictated by clinical criteria. Our
510    sampling of brain activity is extensive but not exhaustive (**Figure 1**, **Tables S1-S2**). It is
511    quite possible that other areas not examined here may also reveal neural correlates of
512    parts of speech and that the regions we found interact with other relevant brain areas. A
513    critical goal of cortical resections in epilepsy patients is to cure seizures without interfering
514    with cognitive function. As such, given the strong lateralization and ubiquitous role for
515    language in cognition, it is extremely important to precisely understand the neural
516    structures that support language in these patients and the current results could help guide
517    surgical approaches for epilepsy. Third, the current work focuses on three parts of
518    speech. Nouns, adjectives, and verbs do *not* constitute an exhaustive list of parts of
519    speech and future work should examine the representation of pronouns[67], adverbs,
520    prepositions, and conjunctions. Finally, our work provides a *correlate* of the
521    representation of POS, but future work should evaluate whether any such signals are
522    causally required for online language interpretation.
523
524    These results provide initial glimpses into highly localized structures that represent
525    a fundamental component of language that has been extensively studied by linguists for
526    decades, the functional role of different words within a sentence. The representation of
527    nouns versus adjectives in the human brain is invariant to the presentation modality, word
528    properties, grammar, and semantics. Furthermore, the representation even generalizes
529    across different languages. These observations open the doors to begin to elucidate the

530 neural representation of more complex language concepts and to bridge the extensive
531 work in language and linguistics to their underlying neural representations.

532 **Methods**
533 **Preregistration**
534 This study was preregistered on the Open Science Framework (OSF) website. The
535 preregistration DOI is: https://doi.org/10.17605/OSF.IO/8TU2G.

536

537 **Data availability**
538 All data and code will be made publicly available through the following link:
539 https://klab.tch.harvard.edu/resources/Misraetal_POS.html
540 The pseudocode can be found within the Readme.docx file in the above link.

541

542 **Participants**
543 We recorded data from 20 participants (9 male, 9-60 years old, 2 left-handed, 2
544 ambidextrous, **Table S1**) with pharmacologically resistant epilepsy for the minimal phrase
545 task (**Fig. 1a**) and 17 participants (7 male, 13-50 years old, 3 left-handed, Table S9) for
546 the full sentence task (**Fig. 5a**). All experiments were conducted while participants stayed
547 at Children's Hospital Boston (CHB), Brigham and Women's Hospital (BWH), Taipei
548 Veterans General Hospital (TVGH), or Cleveland Clinic (CC). All studies were approved
549 by each hospital's institutional review boards and were carried out with the participants'
550 informed consent.

551

552 **Recordings and Electrode Locations**
553 Participants were implanted with intracranial electrodes via stereo
554 electroencephalography (sEEG) (Ad-Tech, Racine, WI, USA). Neurophysiological data
555 were recorded using XLTEK (Oakville, ON, Canada), Bio-Logic (Knoxville, TN, USA),
556 Nihon Kohden (Tokyo, Japan), and Natus (Pleasanton, CA). The sampling rate was 2048
557 Hz at BCH and TVGH, 1024 Hz or 512 Hz at BWH, and 1000 Hz at CC. All data were
558 referenced in a bipolar montage. There were no seizure events in any of the sessions.
559 Electrode locations were decided based on clinical criteria for each participant. Electrodes
560 in the epileptogenic foci, as well as pathological areas, were removed from analyses. The

561 total number of electrodes after bipolar referencing and removing electrodes with no
562 signal, line noise or recording artifacts was 1,801 for the minimal phrase task and 1,593
563 for the full sentence task[68].

564

565 Following implantation, electrodes were localized by co-registration of pre-operative T1
566 MRI and post-operative CT scans using the iELVis software[47]. We used FreeSurfer to
567 segment MRI images, upon which post implant CT was rigidly registered[69]. Electrodes
568 were marked in the CT aligned to pre-operative MRI using the Bioimage Suite[70]. The
569 Desikan-Killiany (DK) atlas was used to assign the electrodes locations. **Figure 1b-g** and
570 **Table S2** show the locations of all the electrodes.

571

572 **Experiment Design**
573 All visual stimuli were displayed on a 15.4 inch 2,880 × 1,800 pixel LCD screen using the
574 Psychtoolbox in MATLAB (Natick, MA) and a MacBook Pro laptop (Cupertino, CA). The
575 stimuli were positioned at eye level at about 80 cm from the participant and each word
576 subtended approximately 3 degrees of visual angle. Sounds were played from the
577 speakers of a MacBook Pro 15.4 at 80% loudness using the Psychtoolbox in MATLAB[71].
578 We used the USB-1208FS-Plus device from Measurement Computing Corporation
579 (Norton, Massachusetts) to send trigger pulses that enabled us to align stimuli onsets and
580 behavioral responses to neural recordings.

581

582 *Minimal Phrase Task:*
583 A schematic of the task is shown in **Figure 1**. Participants were presented two words,
584 875 ms presentation time, with a 400 ms blank screen between them. At the end of each
585 trial, participants were asked to indicate via a button press whether the two words were
586 same or different. Word presentation was either visual or auditory. On average, we
587 presented 1500 ± 710 trials (**Table S1** shows the number of trials per participant).

588

589 There were three types of trials: Noun followed by Adjective (42% of trials, e.g., "apple
590 green"), Adjective followed by Noun (42% of trials, e.g., "green apple"), Repeated Noun
591 (8% of trials, e.g., "apple apple"), and Repeated Adjective (8% of trials, e.g., "green

592  green"). The order of trials (stimulus presentation modality and noun/adjective structure)
593  was randomly interleaved. Each word combination was presented in a randomized
594  manner 5 times in the audio modality and 5 times in the visual modality. The nouns
595  belonged to two categories, animals (e.g., "cat") and food (e.g., "apple"). The adjectives
596  belonged to two categories, concrete adjectives (e.g., "big") and abstract adjectives (e.g.,
597  "good"). A list of all the nouns and adjectives is included in **Table S3**. We selected only
598  high frequency English words that were more frequent than $10^{-6}$ in Google Ngram and
599  were shorter than 7 letters and had no more than 1 or 2 syllables. We used the max
600  frequency of a word between 2006 and 2019. Finally, we created a balanced selection of
601  nouns and adjectives such that noun and adjectives were indistinguishable from each
602  other using word length or number of syllables ($p > 0.05$ ranksum test). We conducted the
603  experiment in 3 languages, English (16 monolingual and 1 bilingual participants), Spanish
604  (1 bilingual participants) and Taiwanese (3 monolingual participants). Two bilingual
605  international scholars whose native language was Spanish (MAG), and Taiwanese (YLK)
606  translated the words in the task. For non-English languages, we also kept nouns and
607  adjectives indistinguishable based on word-length and number of syllables.

608

609  Participants had to indicate whether the two words in a trial were the same or not. The
610  motor responses were the same for nouns or adjectives. The motor responses were also
611  the same for noun followed by adjective or adjective followed by noun trials. Thus, the
612  motor responses were orthogonal to parts of speech and grammar and differences
613  between nouns and adjectives cannot be attributed to motor signals.

614

615  *Sentence Task:*
616  A schematic of the task is shown in **Figure 5a**. Participants were presented with four-
617  word sentences. There was a 600 ms fixation, followed by four words presented
618  sequentially for 875 ms each. After the last word there was a 1,000 ms delay with a gray
619  screen and then an image was presented.

620

621  There were two types of trials: semantical (50% of trials, e.g., "the girls ate cakes"), and
622  non-semantic (50% of trials, e.g., "the cakes ate girls"). The non-semantic sentences were

623    formed by swapping the nouns of the correct sentences, without changing the

624    grammatical correctness of the sentence (see **Table S10,** for example sentences). The

625    order of trials (stimulus presentation modality and semantical/non-semantic structure)

626    was randomly interleaved. Participants were instructed to indicate via a button-press

627    whether the sentence described the image (green button) or not (red button), ignoring

628    notions of singular or plural. An accurately described image was *only* possible for

629    semantic sentences We conducted the experiment in 3 languages, English (13

630    monolingual), Taiwanese (3 monolingual) and Hindi (1 monolingual participant).

631

632    **Data Analyses**

633    *Preprocessing*

634    For the minimal phrase task, a total of 2,428 electrode contacts were implanted, 627 of

635    which were excluded from analysis due to bipolar referencing, presence of line noise or

636    recording artifacts[68]. Similarly, for the sentence task, a total of 2,507 electrode contacts

637    were implanted, 914 of which were excluded from the analysis. We removed 60 Hz line

638    noise and its harmonics using a fifth-order Butterworth filter. We focus on the high-gamma

639    band of the intracranial field potential signals obtained by bandpass filtering raw data of

640    each electrode in the 65–150 Hz range (fifth-order Butterworth filter). The high gamma

641    band (65-150 Hz) power was computed using the Chronux toolbox[72]. We used a time-

642    bandwidth product of 3 and 4 leading tapers, a moving window size of 200 ms, and a step

643    size of 5 ms. For every trial, we computed the normalized high gamma activity by

644    subtracting the mean activity from -150 to 50 ms from the onset of the first fixation and

645    then dividing by the standard deviation. This normalized response is reported as "gamma

646    power" on the y-axis when showing electrode responses.

647

648    *Responsive Electrodes*

649    We evaluated whether an electrode was responsive to visual or auditory stimuli by

650    comparing the 100 to 400 ms post stimulus onset to the -400 to -100 ms before stimulus

651    onset (e.g., **Figure S1**). The responsiveness threshold was set using Cohen's d prime

652    coefficient and based on the number of trials for a statistical power of 80% and $p<0.01$

653 (one-tailed z-test). We also computed the time at which the neural signals reach half of
654 the maximum amplitude.
655
656 *Part-of-speech selectivity*
657 We compared the neural responses to nouns versus adjectives. Periods of significant
658 selective activation were tested using a one-tailed t-test with p<0.05 at each time point to
659 differentiate between nouns and adjectives and were corrected for multiple comparisons
660 with a Benjamini-Hochberg false detection rate (FDR) corrected threshold of q<0.05,
661 separately for auditory and visual trials. After fixing the FDR with q<0.05, an electrode
662 was considered to be selective for part of speech if there was a significant difference
663 between nouns and adjectives for a minimum contiguous window of 65 ms.
664
665 *General Linear Model (GLM)*
666 *Minimal Phrase Task:*
667 We created a GLM to tease out the experiment variables that significantly contribute to
668 explaining the responses of a given electrode.

669 $$AUC \ = \ \beta_0 \ + \ \beta_{NvsA} NvsA \ + \ \beta_{GvsUG} GvsUG \ + \ \beta_{NSyllables} NSyllables$$
670 $$+ \ \beta_{WordLength} WordLength \quad (1)$$

671 where AUC is the area under the response curve (e.g., **Figure 2a**) from 200 ms to 800
672 ms after word onset, $\beta_0$ is a constant additive term, NvsA is 1 for Nouns and -1 for
673 Adjectives, GvsUG is 1 for Grammatical trials and -1 for Ungrammatical trials,
674 NumberOfSyllables is 1 or 2 (and 0 for visual trials), or WordLength goes from 3 to 7 (and
675 0 for auditory trials) as the task predictors. We fit this GLM model for each electrode
676 separately using the MATLAB function fitglm and report the corresponding β coefficients
677 (e.g., **Figure 2j**). We assessed whether each coefficient was significantly different from
678 zero when compared to β coefficients generated from shuffled labels (p<0.01, corrected
679 for multiple comparisons).
680
681 *Sentence Task:*
682 The GLM followed the one in the minimal phrase task. The GLM in this case was:

683
$$AUC = \beta_0 + \beta_{NvsV}NvsV + \beta_{SvsNS}SvsNS + \beta_{NSyllables}NSyllables$$

684
$$+ \beta_{WordLength}WordLength \quad (2)$$

685 where AUC is the area under the response curve (e.g., **Figure 5b**) from 200 ms to 800

686 ms after the word onset, $\beta_0$ is a constant additive term, NvsV is 1 for Nouns and -1 for

687 Verbs, SvsNS is 1 for semantic trials and -1 for non-semantic trials, NumberOfSyllables

688 is 1, 2 or 3 (and 0 for visual trials), and WordLength goes from 3 to 10 (and 0 for auditory

689 trials).

690

691 *Anatomical comparisons*

692 To assess the degree of anatomical specificity in the neural responses, we compared the

693 percentage of significant electrodes in each brain region to the null distribution expected

694 given the number of electrodes in each area using a permutation test ($p < 0.01$, $10^6$

695 iterations). A similar approach was followed to compare the same region between the left

696 and right hemispheres.

697

698 *Decoding Analysis*

699 We performed a machine learning decoding analysis[73] to decode parts of speech in

700 individual words combining all the electrodes in each brain region as defined by the

701 Desikan-Killiany atlas[48] (**Figure S9**). The top-N principal components of all electrodes

702 that explained more than 70% of the variance in the training data for the area under curve

703 of non-overlapping 100 ms time-windows of the signal following word onset were used

704 for decoding. The signal for decoding comprised of features from different frequency

705 bands (beta:12-30 Hz, low gamma:30-65 Hz, and high gamma power: 65-150Hz). The

706 analysis was repeated for 30 random splits of the data with 80% of the data used for

707 training a Support Vector Machine with a linear kernel. Significant decoding performance

708 was found by comparing performance from the original data at each time-window with a

709 null distribution obtained by shuffling labels ($p < 0.01$, ranksum test). Regions with

710 statistically significant decoding performance were found by comparing the average of

711 the maximum decoding performance across time for 30 random iterations of the original

712 data with that of the null distribution, separately for both hemispheres ($p < 0.01$, ranksum

713    test corrected for multiple comparisons) (**Figure 3c,f,i, Figure S9c, Figure S10c,f,i,l,o,r**).

714    We also applied a threshold such that for a given region R

715
$$\left[\mu_R - 3 * \sigma_R\right]_{original\ data} > \left[\mu_R + 3 * \sigma_R\right]_{null\ data}$$

716    where μ and σ represent the average and standard deviation in region R. For the

717    significant regions, the average max-performance between the left and right hemispheres

718    was compared to find if decoding performance was lateralized (p<0.01, ranksum test,

719    corrected for multiple comparisons).

720

725

**References**

1    Chomsky, N. *The minimalist program*.  (MIT Press, 1995).
2    Scott, S. K. From speech and talkers to the social world: The neural processing of human spoken language. *Science* **366**, 58-62, doi:10.1126/science.aax0288 (2019).
3    Pylkkanen, L. The neural basis of combinatory syntax and semantics. *Science* **366**, 62-66, doi:10.1126/science.aax0050 (2019).
4    Heilman, K. M. & Valenstein, E. *Clinical Neuropsychology*.  (Oxford University Press, 1993).
5    Ojemann, G., Ojemann, J., Lettich, E. & Berger, M. Cortical language localization in left, dominant hemisphere. An electrical stimulation mapping investigation in 117 patients. *J Neurosurg* **71**, 316-326, doi:10.3171/jns.1989.71.3.0316 (1989).
6    Petrides, M. *Neuroanatomy of language regions in the human brain*.  (Elsevier, 2014).
7    Mahon, B. Z. & Caramazza, A. Concepts and categories: a cognitive neuropsychological perspective. *Annual review of psychology* **60**, 27-51, doi:10.1146/annurev.psych.60.110707.163532 (2009).
8    Mesulam, M. M., Thompson, C. K., Weintraub, S. & Rogalski, E. J. The Wernicke conundrum and the anatomy of language comprehension in primary progressive aphasia. *Brain* **138**, 2423-2437, doi:10.1093/brain/awv154 (2015).
9    Warrington, E. & Shallice, T. Category specific semantic impairments. *Brain* **107**, 829-854 (1984).
10    Nourski, K. V. *et al.* Gamma Activation and Alpha Suppression within Human Auditory Cortex during a Speech Classification Task. *J Neurosci* **42**, 5034-5046, doi:10.1523/JNEUROSCI.2187-21.2022 (2022).
11    Forseth, K. J. *et al.* A lexical semantic hub for heteromodal naming in middle fusiform gyrus. *Brain* **141**, 2112-2126, doi:10.1093/brain/awy120 (2018).
12    Murphy, E. *et al.* Minimal Phrase Composition Revealed by Intracranial Recordings. *J Neurosci* **42**, 3216-3227, doi:10.1523/JNEUROSCI.1575-21.2022 (2022).
13    Keshishian, M. *et al.* Joint, distributed and hierarchically organized encoding of linguistic features in the human auditory cortex. *Nat Hum Behav* **7**, 740-753, doi:10.1038/s41562-023-01520-0 (2023).
14    Sinai, A. *et al.* Electrocorticographic high gamma activity versus electrical cortical stimulation mapping of naming. *Brain* **128**, 1556-1570, doi:10.1093/brain/awh491 (2005).
15    Ding, N., Melloni, L., Zhang, H., Tian, X. & Poeppel, D. Cortical tracking of hierarchical linguistic structures in connected speech. *Nat Neurosci* **19**, 158-164, doi:10.1038/nn.4186 (2016).
16    Cometa, A. *et al.* Event-related causality in stereo-EEG discriminates syntactic processing of noun phrases and verb phrases. *J Neural Eng* **20**, doi:10.1088/1741-2552/accaa8 (2023).
17    Artoni, F. *et al.* High gamma response tracks different syntactic structures in homophonous phrases. *Sci Rep* **10**, 7537, doi:10.1038/s41598-020-64375-9 (2020).
18    Bemis, D. & Pylkkänen, L. Basic linguistic composition recruits the left anterior temporal lobe and left angular gyrus during both listening and reading. . *Cereb Cortex* **23**, 1859-1873 (2013).
19    Bemis, D. K. & Pylkkanen, L. Simple composition: a magnetoencephalography investigation into the comprehension of minimal linguistic phrases. *J Neurosci* **31**, 2801-2814, doi:10.1523/JNEUROSCI.5003-10.2011 (2011).

772  20  Crepaldi, D., Berlingeri, M., Paulesu, E. & Luzzatti, C. A place for nouns and a place for
773      verbs? A critical review of neurocognitive data on grammatical-class effects. *Brain and*
774      *Language* **116**, 33-49, doi:10.1016/j.bandl.2010.09.005 (2011).
775  21  Woolnough, O. *et al.* Spatiotemporal dynamics of orthographic and lexical processing in
776      the ventral visual pathway. *Nat Hum Behav* **5**, 389-398, doi:10.1038/s41562-020-00982-w
777      (2021).
778  22  Vigliocco, G., Vinson, D. P., Druks, J., Barber, H. & Cappa, S. F. Nouns and verbs in the
779      brain: a review of behavioural, electrophysiological, neuropsychological and imaging
780      studies. *Neurosci Biobehav Rev* **35**, 407-426, doi:10.1016/j.neubiorev.2010.04.007 (2011).
781  23  Castellucci, G. A., Kovach, C. K., Howard, M. A., 3rd, Greenlee, J. D. W. & Long, M. A.
782      A speech planning network for interactive language use. *Nature* **602**, 117-122,
783      doi:10.1038/s41586-021-04270-z (2022).
784  24  Yi, H. G. *et al.* Learning nonnative speech sounds changes local encoding in the adult
785      human cortex. *Proc Natl Acad Sci U S A* **118**, doi:10.1073/pnas.2101777118 (2021).
786  25  Gwilliams, L., King, J. R., Marantz, A. & Poeppel, D. Neural dynamics of phoneme
787      sequences reveal position-invariant code for content and order. *Nat Commun* **13**, 6606,
788      doi:10.1038/s41467-022-34326-1 (2022).
789  26  Forseth, K., Pitkow, X., Fischer-Baum, S. & Tandon, N. What the brain does as we speak.
790      *bioRxiv* **bioRxiv 2021.02.05.429841**, doi:https://doi.org/10.1101/2021.02.05.429841
791      (2021).
792  27  Forseth, K. J., Hickok, G., Rollo, P. S. & Tandon, N. Language prediction mechanisms in
793      human auditory cortex. *Nat Commun* **11**, 5240, doi:10.1038/s41467-020-19010-6 (2020).
794  28  Bhaya-Grossman, I. & Chang, E. F. Speech Computations of the Human Superior
795      Temporal Gyrus. *Annual review of psychology* **73**, 79-102, doi:10.1146/annurev-psych-
796      022321-035256 (2022).
797  29  Hamilton, L. S., Oganian, Y., Hall, J. & Chang, E. F. Parallel and distributed encoding of
798      speech across human auditory cortex. *Cell* **184**, 4626-4639 e4613,
799      doi:10.1016/j.cell.2021.07.019 (2021).
800  30  Khanna, A. R. *et al.* Single-neuronal elements of speech production in humans. *Nature*
801      **626**, 603-610, doi:10.1038/s41586-023-06982-w (2024).
802  31  Jamali, M. *et al.* Semantic encoding during language comprehension at single-cell
803      resolution. *Nature*, doi:10.1038/s41586-024-07643-2 (2024).
804  32  Chomsky, N., Gallego, A. & Ott, D. Generative grammar and the faculty of language:
805      insights, questions, and challenges. *Catalan Journal of Linguisttics*, 226-261 (2019).
806  33  OpenAI. GPT-4 Technical Report. *arXiv* **2303.08774** (2023).
807  34  Hagoort, P. The neurobiology of language beyond single-word processing. *Science* **366**,
808      55-58, doi:10.1126/science.aax0289 (2019).
809  35  Hagoort, P. & Indefrey, P. The neurobiology of language beyond single words. *Annu Rev*
810      *Neurosci* **37**, 347-362, doi:10.1146/annurev-neuro-071013-013847 (2014).
811  36  Calinescu, L., Ramchand, G. & Baggio, G. How (not) to look for meaning composition in
812      the brain: A reassessment of current experimental paradigms. *Frontiers in Language*
813      *Sciences* **2**, 1096110 (2023).
814  37  Goldstein, A. *et al.* Shared computational principles for language processing in humans
815      and deep language models. *Nat Neurosci* **25**, 369-380, doi:10.1038/s41593-022-01026-4
816      (2022).

817    38    Cai, J., Hadjinicolau, A., Paulik, A., Williams, Z. & Cash, S. Natural language processing
818        models reveal neural dynamics of human communication. *Biorxiv* **2023.03.10.531.095**
819        (2023).

820    39    Tenney, I., Dipanjan, D. & Pavlick, E. BERT rediscovers the classical NLP pipeline. *arXiv*
821        **1905.05950** (2019).

822    40    Elazar, Y., Ravfogel, S., Jacovi, A. & Goldberg, Y. Amnesic probing: behavioral
823        explanation with amnesic counterfactuals. *arXiv* **2006.00995** (2021).

824    41    Rapp, B. & Caramazza, A. Selective difficulties with spoken nouns and written verbs: A
825        single case study. *Journal of neurolinguistics* **15**, 373-402 (2002).

826    42    Caramazza, A. & Hillis, A. E. Lexical organization of nouns and verbs in the brain. *Nature*
827        **349**, 788-790, doi:10.1038/349788a0 (1991).

828    43    Woolnough, O., Forseth, K. J., Rollo, P. S., Roccaforte, Z. J. & Tandon, N. Event-related
829        phase synchronization propagates rapidly across human ventral visual cortex. *Neuroimage*
830        **256**, 119262, doi:10.1016/j.neuroimage.2022.119262 (2022).

831    44    Aflalo, T. *et al.* A shared neural substrate for action verbs and observed actions in human
832        posterior parietal cortex. *Sci Adv* **6**, doi:10.1126/sciadv.abb3984 (2020).

833    45    Damasio, A. R. & Tranel, D. Nouns and verbs are retrieved with differently distributed
834        neural systems. *Proc Natl Acad Sci U S A* **90**, 4957-4960, doi:10.1073/pnas.90.11.4957
835        (1993).

836    46    Bansal, A. *et al.* Neural Dynamics Underlying Target Detection in the Human Brain.
837        *Journal of Neuroscience* **34**, 3042-3055, doi: 10.1523/JNEUROSCI.3781-13.2014 (2014).

838    47    Groppe, D. M. *et al.* iELVis: An open source MATLAB toolbox for localizing and
839        visualizing human intracranial electrode data. *J Neurosci Methods* **281**, 40-48,
840        doi:10.1016/j.jneumeth.2017.01.022 (2017).

841    48    Desikan, R. S. *et al.* An automated labeling system for subdividing the human cerebral
842        cortex on MRI scans into gyral based regions of interest. *Neuroimage* **31**, 968-980 (2006).

843    49    Tager-Flusberg, H. & Seery, A. in *Neural circuit development and function in the brain*
844        (eds JLR Rubenstein & P Rakic) 315-330 (Elsevier, 2013).

845    50    Quian Quiroga, R., Reddy, L., Kreiman, G., Koch, C. & Fried, I. Invariant visual
846        representation by single neurons in the human brain. *Nature* **435**, 1102-1107,
847        doi:10.1038/nature03687 (2005).

848    51    Quian Quiroga, R., Kraskov, A., Koch, C. & Fried, I. Explicit encoding of multimodal
849        percepts by single neurons in the human brain. *Current biology : CB* **19**, 1308-1313,
850        doi:10.1016/j.cub.2009.06.060 (2009).

851    52    Quian Quiroga, R., Kreiman, G., Koch, C. & Fried, I. Sparse but not 'Grandmother-cell'
852        coding in the medial temporal lobe. *Trends in Cognitive Science* **12**, 87-91 (2008).

853    53    Geva-Sagiv, M. *et al.* Augmenting hippocampal-prefrontal neuronal synchrony during
854        sleep enhances memory consolidation in humans. *Nat Neurosci* **26**, 1100-1110,
855        doi:10.1038/s41593-023-01324-5 (2023).

856    54    Wojtasik, M. *et al.* Cytoarchitectonic Characterization and Functional Decoding of Four
857        New Areas in the Human Lateral Orbitofrontal Cortex. *Front Neuroanat* **14**, 2,
858        doi:10.3389/fnana.2020.00002 (2020).

859    55    Kringelbach, M. L. The human orbitofrontal cortex: linking reward to hedonic experience.
860        *Nat Rev Neurosci* **6**, 691-702, doi:10.1038/nrn1747 (2005).

861    56    Ongur, D. & Price, J. L. The organization of networks within the orbital and medial
862           prefrontal cortex of rats, monkeys and humans. *Cereb Cortex* **10**, 206-219,
863           doi:10.1093/cercor/10.3.206 (2000).
864    57    Xiao, Y. *et al.* Integration of recognition, episodic, and associative memories during
865           complex human behavior. *Biorxiv* **2023.03.27.534384**,
866           doi:https://doi.org/10.1101/2023.03.27.534384 (2023).
867    58    Hunt, L. T. *et al.* Triple dissociation of attention and decision computations across
868           prefrontal cortex. *Nat Neurosci* **21**, 1471-1481, doi:10.1038/s41593-018-0239-5 (2018).
869    59    Nogueira, R. *et al.* Lateral orbitofrontal cortex anticipates choices and integrates prior with
870           current information. *Nat Commun* **8**, 14823, doi:10.1038/ncomms14823 (2017).
871    60    Noonan, M. P. *et al.* Separate value comparison and learning mechanisms in macaque
872           medial and lateral orbitofrontal cortex. *Proc Natl Acad Sci U S A* **107**, 20547-20552,
873           doi:10.1073/pnas.1012246107 (2010).
874    61    de Araujo, I. E., Rolls, E. T., Kringelbach, M. L., McGlone, F. & Phillips, N. Taste-
875           olfactory convergence, and the representation of the pleasantness of flavour, in the human
876           brain. *Eur J Neurosci* **18**, 2059-2068, doi:10.1046/j.1460-9568.2003.02915.x (2003).
877    62    Dronkers, N. F., Wilkins, D. P., Van Valin, R. D., Jr., Redfern, B. B. & Jaeger, J. J. Lesion
878           analysis of the brain areas involved in language comprehension. *Cognition* **92**, 145-177,
879           doi:10.1016/j.cognition.2003.11.002 (2004).
880    63    Mesulam, M. M. *et al.* Neuropathological fingerprints of survival, atrophy and language in
881           primary progressive aphasia. *Brain* **145**, 2133-2148, doi:10.1093/brain/awab410 (2022).
882    64    Mesulam, M. M. *et al.* Primary progressive aphasia and the evolving neurology of the
883           language network. *Nat Rev Neurol* **10**, 554-569, doi:10.1038/nrneurol.2014.159 (2014).
884    65    Fried, I., Cerf, M., Rutishauser, U. & Kreiman, G. *Single neuron studies of the human*
885           *brain. Probing cognition.*, 408 (MIT Press, 2014).
886    66    Mukamel, R. & Fried, I. Human intracranial recordings and cognitive neuroscience.
887           *Annual review of psychology* **63**, 511-537, doi:10.1146/annurev-psych-120709-145401
888           (2012).
889    67    Dijksterhuis, D. E. *et al.* Pronouns reactivate conceptual representations in human
890           hippocampal neurons. *bioRxiv*, 2024.2006.2023.600044, doi:10.1101/2024.06.23.600044
891           (2024).
892    68    Wang, J., Tao, A., Anderson, W. S., Madsen, J. R. & Kreiman, G. Mesoscopic
893           physiological interactions in the human brain reveal small-world properties. *Cell Rep* **36**,
894           109585, doi:10.1016/j.celrep.2021.109585 (2021).
895    69    Dale, A. M., Fischl, B. & Sereno, M. I. Cortical surface-based analysis. I. Segmentation
896           and surface reconstruction. *Neuroimage* **9**, 179-194 (1999).
897    70    Joshi, A. *et al.* Unified framework for development, deployment and robust testing of
898           neuroimaging algorithms. *Neuroinformatics* **9**, 69-84, doi:10.1007/s12021-010-9092-8
899           (2011).
900    71    Brainard, D. The Psychophysics Toolbox. *Spatial Vision* **10**, 433-436 (1997).
901    72    Mitra, P. & Bokil, H. *Observed brain dynamics* (Oxford University Press, 2008).
902    73    Bansal, A., Golby, A., Madsen, J. & Kreiman, G. in *COSYNE.*
903
904

905    **Figure Captions**

906

907    **Figure 1. Task schematic, electrode locations, and multimodal responses. a**. Task
908    schematic. Two words were sequentially presented either in visual modality or auditory modality.
909    Participants indicated whether the two words were the same (e.g., "apple apple" or "green green",
910    8% of trials of each type) or different (e.g., "green apple" or "apple green": 42% of trials of each
911    type, **Methods**). In the 84% of trials where the two-words were different, there was an adjective
912    followed by a noun or a noun follower by an adjective. **b-f.** Location of all electrodes overlayed on
913    the Desikan-Killiany Atlas shown with different views. Each white circle shows one electrode. **b.**
914    Left lateral view (n=693), **c.** Left medial view (n=693), **d.** Superior, whole brain view (n=1,801), **e.**
915    Inferior, whole brain view (n=1,801), **f.** Right lateral view (n=1108) **g.** Right medial view (n=1108).
916    **h**. Trial-averaged ($\pm$ SEM) gamma power for responses to auditory (light grey) or visual (black)
917    presentations for an example electrode in the left rostral middle frontal gyrus (electrode location
918    shown in **k**). Responses are aligned to word onset (vertical dashed line). The arrows indicate the
919    half-maximum time. **i**, **j**. Raster plots showing each individual trial for the same electrode for each
920    of the 1,496 words for auditory (**i**) and visual (**j**) presentations (see color scale on right).

921

922    **Figure    2.    Neural    signals    distinguish    between    different    parts    of    speech.**
923    **a-d.** Trial-averaged normalized gamma-band power of responses from an example electrode in
924    the left lateral orbitofrontal cortex (see location in **i**) to nouns (blue) or adjectives (red) during
925    presentation of auditory stimuli (**a, b,** n=435 grammatical and 432 ungrammatical trials) or visual
926    stimuli (**c**, **d,** n=435 grammatical and 432 ungrammatical trials) aligned to the onset (vertical
927    dashed line) of the first word (**a**, **c**) or second word (**b**, **d**). Shaded areas denote s.e.m. Horizontal
928    gray lines denote windows of statistically significant differences between responses to nouns
929    versus adjectives (t-test $p<0.05$, Benjamini-Hochberg false detection rate, $q<0.05$).
930    **e-h**. Raster plots showing the responses in each individual trial (see color scale on bottom right).
931    The red and blue curves in **a-d** correspond to the averages of noun and adjective trials,
932    respectively, in **e-h.**
933    **i.** Location of the example electrode in the left lateral orbitofrontal cortex.
934    **j.** Z-scored $\beta$ coefficients for Generalized Linear Model used to predict area under the curve
935    between 200 ms and 800 ms post word onset, using four task predictors: Noun versus Adjectives,
936    Grammatically correct versus incorrect, number of syllables (auditory presentation) and word
937    length (visual presentation). Asterisks denote statistically significant coefficients, corrected for
938    multiple comparisons (**Methods**).

**k.** Inferior axial view of both hemispheres showing electrodes that revealed statistically significant differences between nouns and adjectives for both audio and visual presentation (orange circles, n=13 electrodes). All the electrodes whose responses were significantly explained *only* by the Nouns versus Adjective task predictor in the GLM are included in this plot.

**l, m.** All electrodes from **k** projected onto the left hemisphere are shown on the frontal plane (**l**) and the axial plane (**m,** same plane as **k**). Electrodes that respond more strongly to nouns, i.e., Nouns versus Adjectives β>0 (n=10 electrodes), are shown in blue and electrodes that responded more strongly to adjectives (β<0, n=3 electrodes), are shown in red. All units are in MNI305 coordinates. Linear support vector machines separating these electrodes are shown with a thick black line. Kernel density curves (bandwidth 2) outline the marginal distributions of noun-preferring (blue) and adjective-preferring (red) electrodes along the lateral-medial axis (**l,m:** top x-axis), ventral-dorsal axis (**l:** right z-axis) and anterior-posterior axis (**m:** right y-axis). P-values indicate significant differences between the coordinates for noun- and adjective-preferring electrodes (ranksum test).

**Fig. 3 | Neural signals distinguishing nouns and adjectives in single trials generalize across Word1 and Word2.**

**a, b, d, e, g, h.** Average cross-validated performance of a support vector machine classifier (SVM, 80% training/20% test) decoding nouns versus adjectives for all electrodes in the left lateral orbitofrontal cortex (LOF) **(a, d, g)** or the right LOF **(b, e, h)**. The dotted horizontal black line shows the chance level. Shaded areas denote s.e.m. Solid horizontal black bar shows time points where performance significantly differed from chance (100 random shuffles, ranksum test, p<0.01). The inputs to the SVM included the top-N principal components of the electrode response that explained >70% variance for the training data at each time bin (**Methods**). **a, b**: Features from auditory and visual responses were combined and used for training and testing on a dataset of both Word1 and Word2 trials. **c, d**: Generalization across word order was evaluated on a dataset where Word1 trials were used for training and word2 trials were used for testing. **g, h**: Training on Word2 and testing on Word1. Black: original labels; Gray: shuffled labels (see **Figure S9** for decoding performance when the number of electrodes was same across all regions and both hemispheres).

**c, f, i.** Summary of average of max-decoding performance for distinguishing nouns versus adjectives in each hemisphere (dark: left; white: right)) for different brain regions. Bottom asterisks denote regions with significant decoding performance with respect to chance and performance from the real and null distribution do not overlap within 3 standard deviations of each other

973    (p<0.01, ranksum test, corrected for multiple comparisons, **Methods**). Shaded box: maximum of

974    the mean ± SD. for the null distribution across all regions. Top asterisks with a U-bracket denote

975    significant differences between decoding accuracy of the left versus the right hemisphere (p<0.01,

976    ranksum test, corrected for multiple comparisons). Regions are sorted in descending order of

977    performance in panel **c**. **c**: Classifiers were trained and tested with features from both Word1 and

978    Word2 trials. **f**: Classifiers were trained on Word1 trials and tested on Word2 trials. **i**: Classifiers

979    were trained on Word2 trials and tested on Word1 trials (see **Figure S10** for controls on word1-

980    only, word2-only, audio-only, visual-only, audio-to-vision and vision-to-audio performance).

981

982    **Figure 4. Neural signals in left LOF generalize across languages in a bilingual participant**

983    **and in monolingual participants**.

984    **a-h**. Trial averaged responses of an electrode in the left lateral orbitofrontal cortex from a bilingual

985    patient. The format follows **Fig. 2a-d**. (**a-d**) English words (audio: n=190 grammatical and 185

986    ungrammatical trials; vision: n=189 grammatical and 191 ungrammatical trials). (**e-h**) Spanish

987    words (audio: n=184 grammatical and ungrammatical trials; vision: 184 grammatical and 186

988    ungrammatical trials). Auditory responses (**a**, **b**, **e**, **f**). Visual responses (**c**, **d**, **g**, **h**). Word 1 (**a**, **c**,

989    **e**, **g**) and Word 2 (**b**, **d**, **f**, **h**).

990    **i, j**. Z-scored β coefficients for Generalized Linear Model to predict area under the curve (AUC)

991    for the English experiment (**i**) and for the Spanish experiment (**j**). The AUC computed between

992    200 ms and 800 ms post word onset using four task predictors: Noun versus Adjectives,

993    Grammatical versus Ungrammatical, number of syllables (auditory presentation) and word length

994    (visual presentation). Asterisks denote statistically significant coefficients corrected for multiple

995    comparisons (**Methods**). The word order for grammatically correct trials in English is an adjective

996    followed by a noun, such as "green apple". This word order gets flipped in grammatically correct

997    Spanish trials.

998    **k**. Inferior view of the 9 out of 38 electrodes (8 audiovisual: **Figure 2k,** 1 visual-only: **Figure S4,**

999    see **Table S4** and **S7**) in the left lateral orbitofrontal cortex that showed noun versus adjective

1000    differences across different languages in which the experiment was conducted (significant Nouns

1001    versus Adjectives β, p<0.01 corrected for multiple comparisons). These electrodes come from 4

1002    different participants. Electrodes from the bilingual patient are in green with a black arrow

1003    indicating the example electrode. Electrodes from one monolingual English participant are in pink

1004    and those from 2 monolingual Taiwanese participants are in brown.

1005

1006    **Fig. 5 | Neural signals distinguish between different parts of speech in sentences.**

**a**. Task Schematic. Sentences comprising four words sequentially presented either in visual or auditory modality were followed by an image. The sentences were either semantic (50% S sentences, e.g., "the girls ate cakes") or non-semantic (50% NS sentences, e.g., "the cakes ate girls"). Participants were instructed to indicate via a button press whether the sentence described the image accurately or not (**Methods**).

**b,c**. Trial-averaged normalized gamma-band power of responses from an example electrode in the pars triangularis (see electrode location in **f**) to nouns (blue) or verbs (black) during presentation of auditory stimuli (**b**, n=628 nouns and 314 verbs) or visual stimuli (**c**, n=628 nouns and 314 verbs) aligned to word onset (vertical dashed line). Shaded areas denote s.e.m. Horizontal gray lines denote windows of statistically significant differences between responses to nouns versus verb (t-test $p<0.05$, Benjamini-Hochberg false detection rate, $q<0.05$).

**d,e**. Raster plots showing the responses in each individual trial (see color scale on bottom right). The blue and black curves in **b,c** correspond to the averages of nouns and verbs, respectively, in **d,e**.

**f**. Location of the example electrode in the pars triangularis (see **Figure S11** for electrode coverage).

**g**. Z-scored β coefficients for the Generalized Linear Model used to predict area under the curve between 200 ms and 800 ms post word onset, using four task predictors: nouns versus verbs, semantic versus non-semantic, number of syllables (auditory presentation) and word length (visual presentation). Asterisks denote statistically significant coefficients, corrected for multiple comparisons (**Methods**).

**h,i**. Lateral view of left (**h**) and right (**i**) hemispheres showing electrodes that revealed statistically significant differences between nouns and verbs for both audio and visual presentation (orange circles, n=41 electrodes, 27 left). Electrodes whose responses were significantly explained only by the nouns versus verbs task predictor in the GLM are included in this plot.

**j,k**. All electrodes from **h,i** projected onto the left hemisphere are shown on the lateral plane (j) and the axial plane (k). All the electrodes that respond more strongly to nouns, i.e., nouns versus verbs β>0 (n=23 electrodes), are shown in blue and electrodes that responded more strongly to verbs (β<0, n=18 electrodes), are shown in black. All units are in MNI305 coordinates. Kernel density curves (bandwidth 2) outline the marginal distributions of noun-preferring (blue) and verb-preferring (black) electrodes along the anterior-posterior axis (**j,k**: y-axis), ventral-dorsal axis (**j**: left z-axis) and lateral-medial axis (**k**: left x-axis, zero being more medial). P-values indicate significant differences between the coordinates for noun- and verb-preferring electrodes (ranksum test).
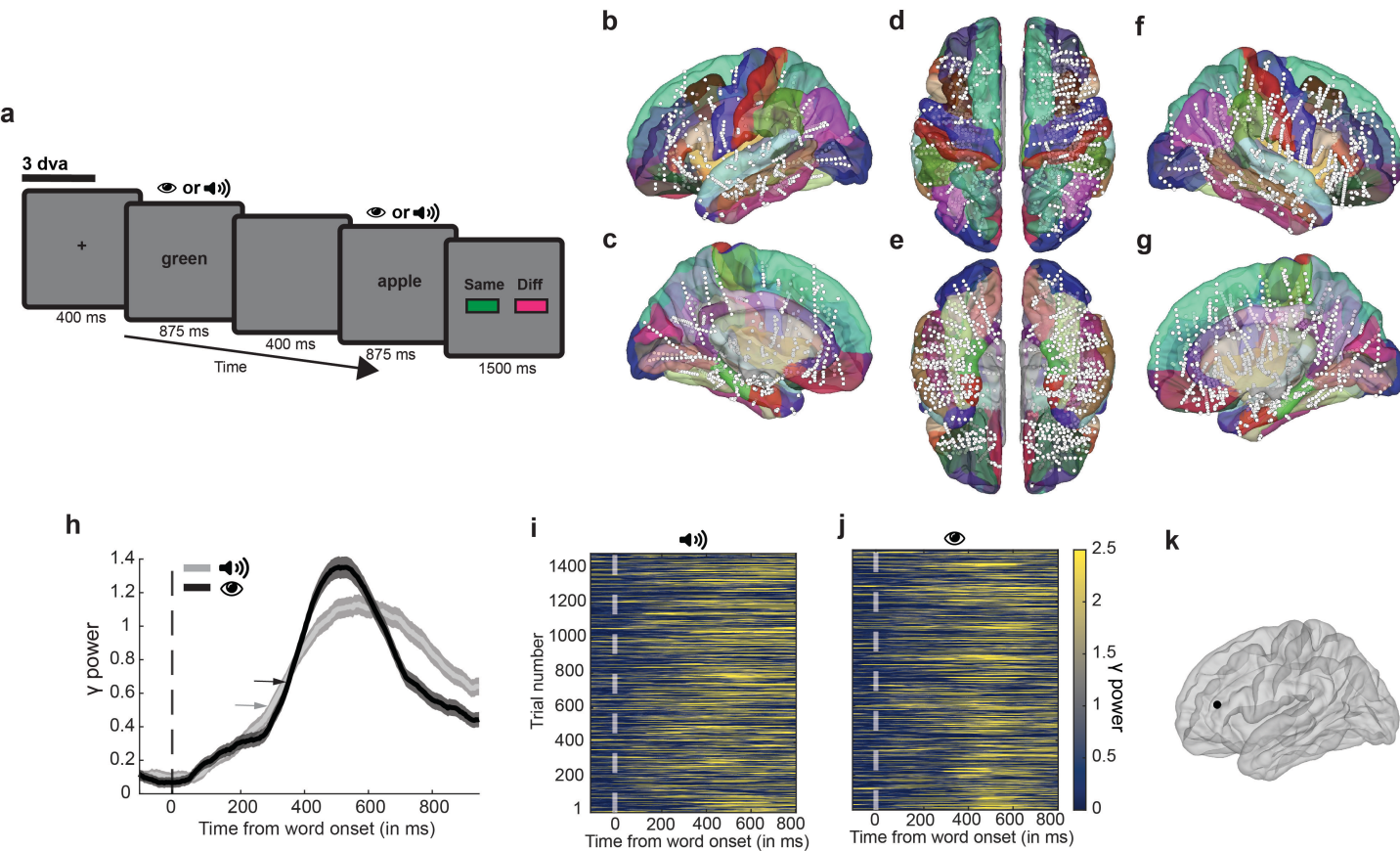
# Figure 1



**Fig. 1 | Task design, electrode locations and multimodal responses.**
**a. Task schematic**. Two words were sequentially presented either in visual modality or auditory modality. Participants indicated whether the two words were the same (e.g., "apple apple" or "green green", 8% of trials of each type) or different (e.g., "green apple" or "apple green": 42% of trials of each type, Methods). In the 84% of trials where the two-words were different, there was an adjective followed by a noun or a noun follower by an adjective. **b-g**. Location of all electrodes overlayed on the Desikan-Killiany Atlas shown with different views. Each white circle shows one electrode. **b**. Left lateral view (n=693), **c**. Left medial view (n=693), **d**. Superior, whole brain view (n=1,801), **e**. Inferior, whole brain view (n=1,801), **f**. Right lateral view (n=1108) **g**. Right medial view (n=1108). **h**. Trial-averaged (± SEM) gamma power for responses to auditory (light grey) or visual (black) presentations for an example electrode in the left rostral middle frontal gyrus (electrode location shown in **k**). Responses are aligned to word onset (vertical dashed line). The arrows indicate the half-maximum time. **i**, **j**. Raster plots showing each individual trial for the same electrode for each of the 1,496 words for auditory (**i**) and visual (**j**) presentations (see color scale on right).
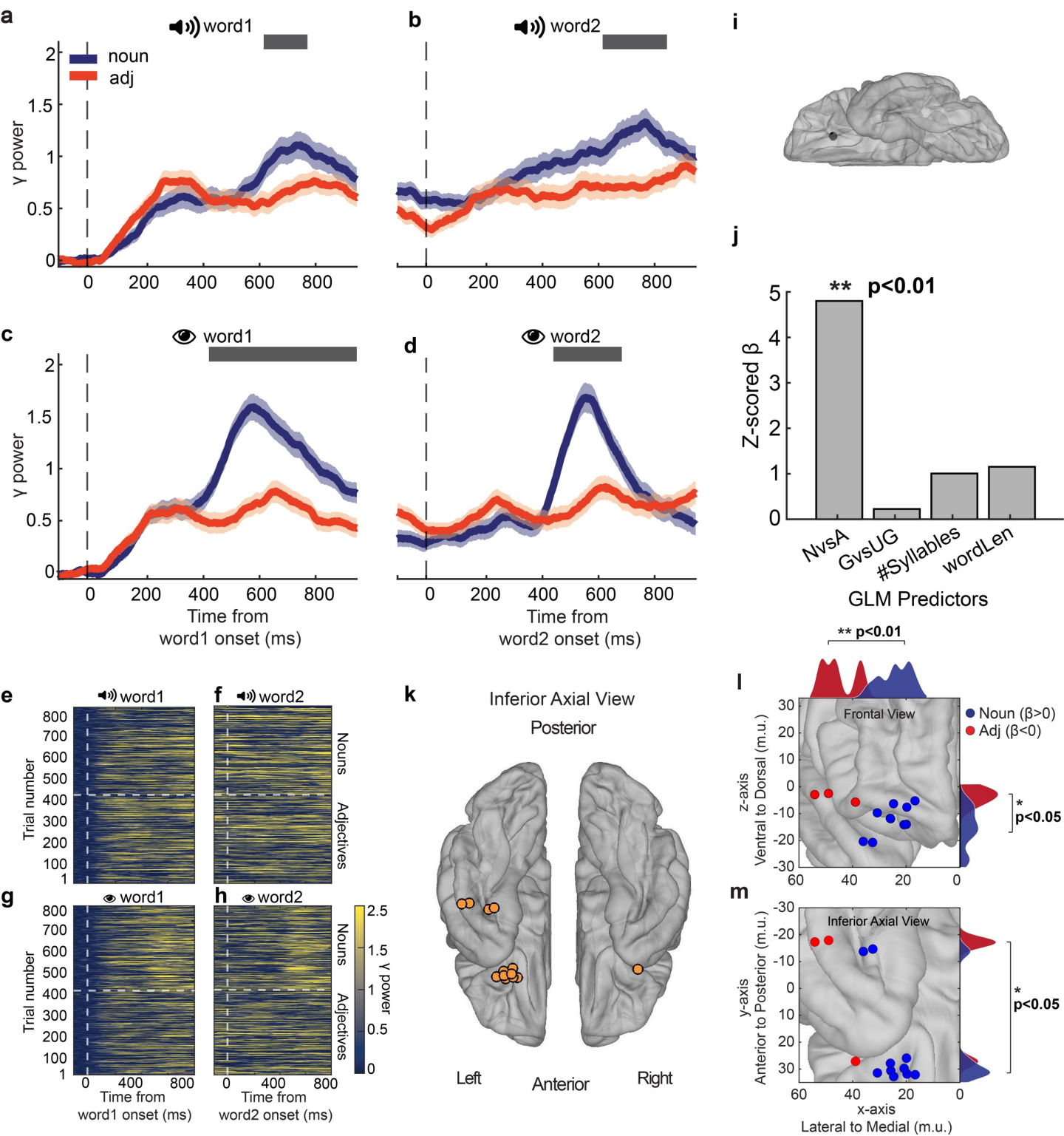
# Figure 2

**Fig. 2 | Neural signals distinguish between different parts of speech.**
**a-d.** Trial-averaged normalized gamma-band power of responses from an example electrode in the left lateral orbitofrontal cortex (see location in **i**) to nouns (blue) or adjectives (red) during presentation of auditory stimuli (**a, b,** n=435 grammatical and 432 ungrammatical trials) or visual stimuli (**c, d,** n=435 grammatical and 432 ungrammatical trials) aligned to the onset (vertical dashed line) of the first word (**a**, **c**) or second word (**b**, **d**). Shaded areas denote s.e.m. Horizontal gray lines denote windows of statistically significant differences between responses to nouns versus adjectives (t-test $p < 0.05$, Benjamini-Hochberg false detection rate, $q < 0.05$).

**e-h**. Raster plots showing the responses in each individual trial (see color scale on bottom right). The red and blue curves in **a-d** correspond to the averages of noun and adjective trials, respectively, in **e-h.**

**i.** Location of the example electrode in the left lateral orbitofrontal cortex.

**j.** Z-scored $\beta$ coefficients for Generalized Linear Model used to predict area under the curve between 200 ms and 800 ms post word onset, using four task predictors: Noun versus Adjectives, Grammatically correct versus incorrect, number of syllables (auditory presentation) and word length (visual presentation). Asterisks denote statistically significant coefficients, corrected for multiple comparisons (**Methods**).

**k.** Inferior axial view of both hemispheres showing electrodes that revealed statistically significant differences between nouns and adjectives for both audio and visual presentation (orange circles, n=13 electrodes). All the electrodes whose responses were significantly explained only by the Nouns versus Adjective task predictor in the GLM are included in this plot.

**l, m.** All electrodes from **k** projected onto the left hemisphere are shown on the frontal plane (**l**) and the axial plane (**m,** same plane as **k**). All the electrodes that respond more strongly to nouns, i.e., Nouns versus Adjectives $\beta > 0$ (n=10 electrodes), are shown in blue and electrodes that responded more strongly to adjectives ($\beta < 0$, n=3 electrodes), are shown in red. All units are in MNI305 coordinates. Kernel density curves (bandwidth 2) outline the marginal distributions of noun-preferring (blue) and adjective-preferring (red) electrodes along the lateral-medial axis (**l,m:** x-axis, zero being more medial), ventral-dorsal axis (**l:** right z-axis) and anterior-posterior axis (**m:** right y-axis). P-values indicate significant differences between the coordinates for noun- and adjective-preferring electrodes (ranksum test).
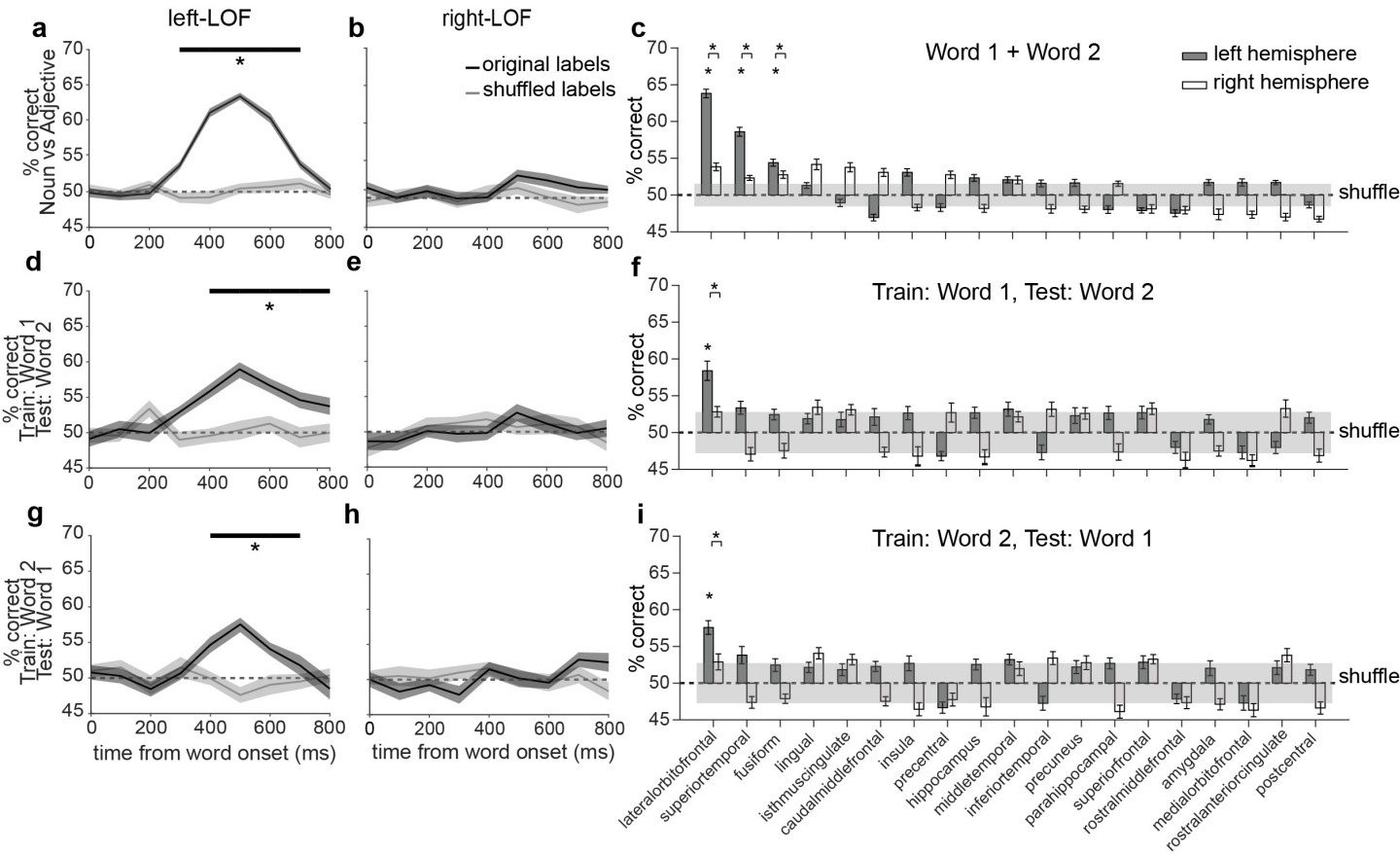
# Figure 3



**Fig. 3 | Neural signals distinguishing nouns and adjectives in single trials generalize across Word1, and Word2.**

**a, b, d, e, g, h.** Average cross-validated performance of a support vector machine classifier (SVM, 80% training/20% test) decoding nouns versus adjectives for all electrodes in the left lateral orbitofrontal cortex (LOF) **(a, d, g)** or the right LOF **(b, e, h)**. The dotted horizontal black line shows the chance level. Shaded areas denote s.e.m. Solid horizontal black bar shows time points where performance significantly differed from chance (100 random shuffles, ranksum test, p<0.01). The inputs to the SVM included the top-N principal components of the electrode response that explained >70% variance for the training data at each time bin (**Methods**). **a, b**: Features from auditory and visual responses were combined and used for training and testing on a dataset of both Word1 and Word2 trials. **c, d**: Generalization across word order was evaluated on a dataset where Word1 trials were used for training and word2 trials were used for testing. **g, h**: Training on Word2 and testing on Word1. Black: original labels; Gray: shuffled labels (see **Figure S9** for decoding performance when the number of electrodes was same across all regions and both hemispheres).

**c, f, i.** Summary of average of max-decoding performance for distinguishing nouns versus adjectives in each hemisphere (dark: left; white: right) for different brain regions. Bottom asterisks denote regions with significant decoding performance with respect to chance and performance from the real and null distribution do not overlap within 3 standard deviations of each other (p<0.01, ranksum test, corrected for multiple comparisons, **Methods**). Shaded box: maximum of the mean ± SD. for the null distribution across all regions. Top asterisks with a U-bracket denote significant differences between decoding accuracy of the left versus the right hemisphere (p<0.01, ranksum test, corrected for multiple comparisons). Regions are sorted in descending order of performance in panel **c**. **c**: Classifiers were trained and tested with features from both Word1 and Word2 trials. **f**: Classifiers were trained on Word1 trials and tested on Word2 trials. **i**: Classifiers were trained on Word2 trials and tested on Word1 trials (see **Figure S10** for controls on word1-only, word2-only, audio-only, visual-only, audio-to-vision and vision-to-audio performance).
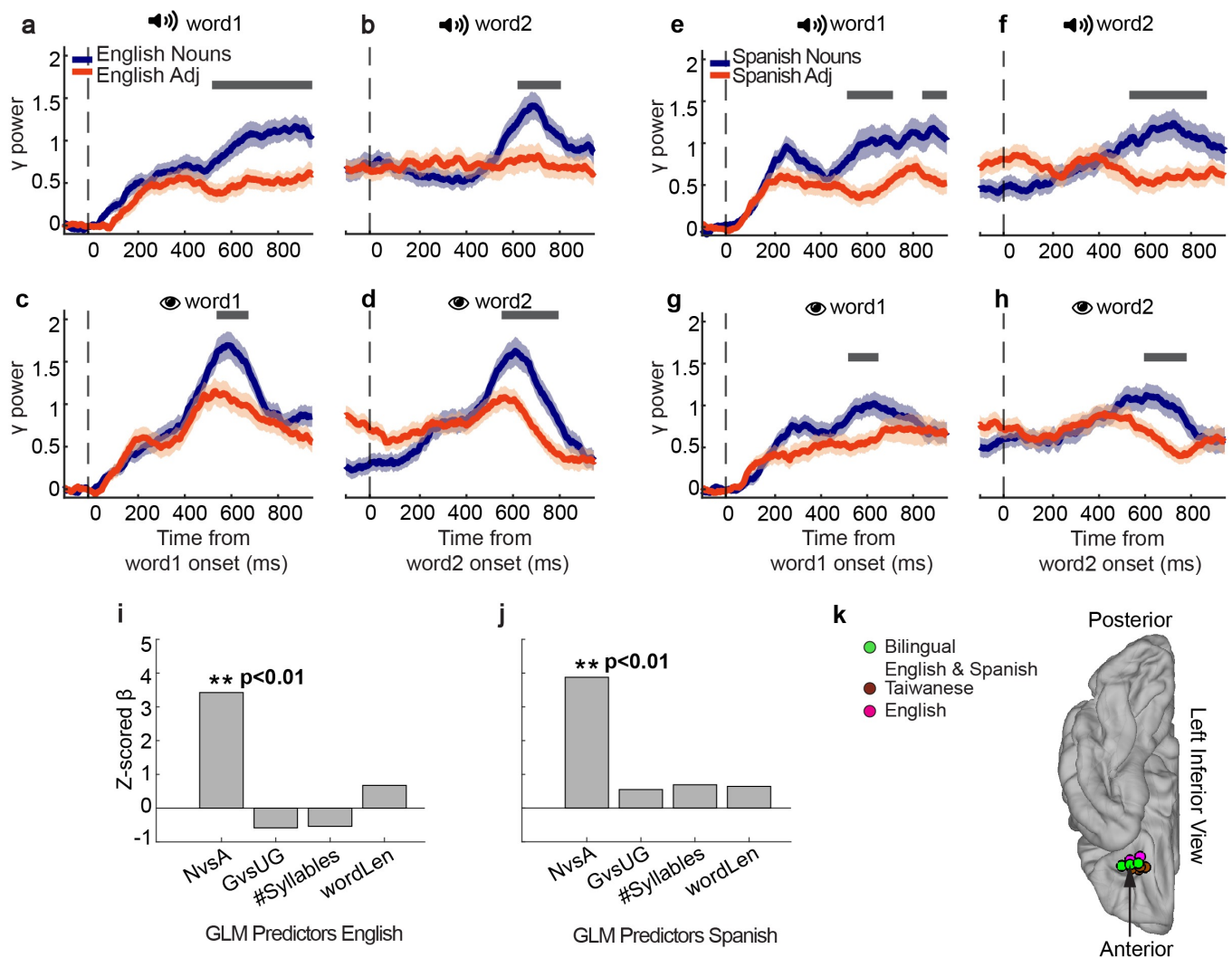
# Figure 4



**Fig. 4 | Neural signals in left LOF generalize across languages in a bilingual participant and in monolingual participants. a-h.** Trial averaged responses of an electrode in the left lateral orbitofrontal cortex from a bilingual patient. The format follows **Fig. 2a-d**. (**a-d**) English words (audio: n=190 grammatical and 185 ungrammatical trials; vision: n=189 grammatical and 191 ungrammatical trials). (**e-h**) Spanish words (audio: n=184 grammatical and ungrammatical trials; vision: 184 grammatical and 186 ungrammatical trials). Auditory responses (**a, b, e, f**) . Visual responses (**c, d, g, h**) . Word 1 (**a, c, e, g**) and Word 2 (**b, d, f, h**) . **i,j.** Z-scored β coefficients for Generalized Linear Model to predict area under the curve (AUC) for the English experiment (**i**) and for the Spanish experiment (**j**)**.** The AUC computed between 200 ms and 800 ms post word onset using four task predictors: Noun versus Adjectives, Grammatical versus Ungrammatical, number of syllables (auditory presentation) and word length (visual presentation). Asterisks denote statistically significant coefficients corrected for multiple comparisons **(Methods)**. The word order for grammatically correct trials in English is an adjective followed by a noun, such as "green apple". This word order gets flipped in grammatically correct Spanish trials. **k.** Inferior view of all the 9 out of 38 electrodes (8 audiovisual: **Figure 2k,** 1 visual-only: **Figure S4,** see **Table S4** and **S7**) in the left lateral orbitofrontal cortex that showed noun versus adjective differences across different languages in which the experiment was conducted (significant Nouns versus Adjectives β, p<0.01 corrected for multiple comparisons). These electrodes come from 4 different participants. Electrodes from the bilingual patient are in green with a black arrow indicating the example electrode. Electrodes from one monolingual English participant are in pink and those from 2 monolingual Taiwanese participants are in brown.
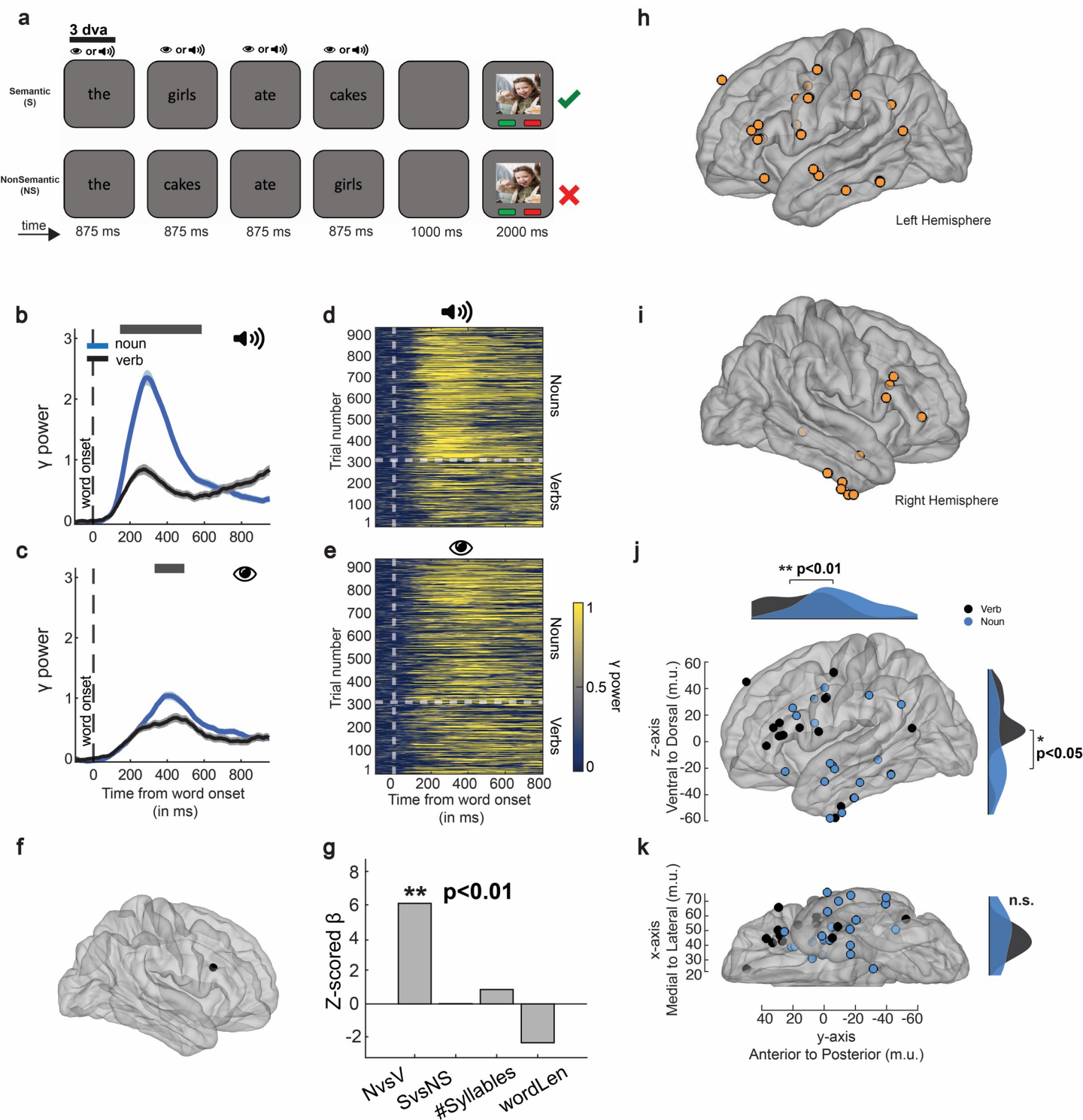
# Figure 5



**a**

| | | | | | |
|---|---|---|---|---|---|
| Semantic (S) | the | girls | ate | cakes | ✅ |
| NonSemantic (NS) | the | cakes | ate | girls | ❌ |

3 dva

time →  875 ms  875 ms  875 ms  875 ms  1000 ms  2000 ms

**b** — noun / verb — γ power vs Time from word onset (in ms) 🔊

**c** — γ power vs Time from word onset (in ms) 👁

**d** — Trial number vs Time from word onset (in ms) 🔊 — Nouns / Verbs

**e** — Trial number vs Time from word onset (in ms) 👁 — Nouns / Verbs — γ power

**f**

**g** — Z-scored β — ** p<0.01 — NvsV, SvsNS, #Syllables, wordLen

**h** — Left Hemisphere

**i** — Right Hemisphere

**j** — ** p<0.01 — Verb / Noun — z-axis Ventral to Dorsal (m.u.) — * p<0.05

**k** — x-axis Medial to Lateral (m.u.) — n.s. — y-axis Anterior to Posterior (m.u.)

**Fig. 5 | Neural signals distinguish between different parts of speech in sentences.**

**a. Task Schematic.** Sentences comprising four words sequentially presented either in visual or auditory modality were followed by an image. The sentences were either semantic (50% **S sentences,** e.g., "the girls ate cakes") or non-semantic (50% **NS sentences,** e.g., "the cakes ate girls"). Participants were instructed to indicate via a button press whether the sentence described the image accurately or not **(Methods).**

**b,c.** Trial-averaged normalized gamma-band power of responses from an example electrode in the pars triangularis (see eletrode location in **f**) to nouns (blue) or verbs (black) during presentation of auditory stimuli (**b,** n=628 nouns and 314 verbs) or visual stimuli (**c,** n=628 nouns and 314 verbs) aligned to word onset (vertical dashed line). Shaded areas denote s.e.m. Horizontal gray lines denote windows of statistically significant differences between responses to nouns versus verb (t-test $p < 0.05$, Benjamini-Hochberg false detection rate, $q < 0.05$).

**d,e.** Raster plots showing the responses in each individual trial (see color scale on bottom right). The blue and black curves in **b,c** correspond to the averages of noun and verb trials, respectively, in **d,e.**

**f.** Location of the example electrode in the pars triangularis (see **Figure S11** for electrode coverage).

**g.** Z-scored $\beta$ coefficients for Generalized Linear Model used to predict area under the curve between 200 ms and 800 ms post word onset, using four task predictors: Noun versus Verbs, Semantically correct versus incorrect, number of syllables (auditory presentation) and word length (visual presentation). Asterisks denote statistically significant coefficients, corrected for multiple comparisons (**Methods**).

**h,i.** Lateral view of left (**h**) and right (**i**) hemispheres showing electrodes that revealed statistically significant differences between nouns and verbs for both audio and visual presentation (orange circles, n=41 electrodes, 27 left). Electrodes whose responses were significantly explained only by the Nouns versus Verbs task predictor in the GLM are included in this plot.

**j,k.** All electrodes from **h,i** projected onto the left hemisphere are shown on the lateral plane (**j**) and the axial plane (**k**). All the electrodes that respond more strongly to nouns, i.e., Nouns versus Verbs $\beta > 0$ (n=23 electrodes), are shown in blue and electrodes that responded more strongly to verbs ($\beta < 0$, n=18 electrodes), are shown in black. All units are in MNI305 coordinates. Kernel density curves (bandwidth 2) outline the marginal distributions of noun-preferring (blue) and verb-preferring (black) electrodes along the anterior-posterior axis (**j,k:** y-axis), ventral-dorsal axis (**j:** left z-axis) and lateral-medial axis (**k:** left x-axis, zero being more medial). P-values indicate significant differences between the coordinates for noun- and verb-preferring electrodes (ranksum test).