

1 **The role of co-speech gestures in retrieval and prediction during naturalistic**
2 **multimodal narrative processing**

3 Sergio Osorio^{a,b,*,1}, Benjamin Straube^c, Lars Meyer^{d,e,+}, & Yifei He^{c,+}

4

5 ^a Laboratory for Cognitive and Evolutionary Neuroscience, School of Medicine, Pontificia
6 Universidad Católica de Chile, Santiago de Chile, Chile

7 ^b Interdisciplinary Centre for Neuroscience, Pontificia Universidad Católica de Chile,
8 Santiago de Chile, Chile

9 ^c Translational Neuroimaging Group, Department of Psychiatry and Psychotherapy
10 Philipps-Universität Marburg, Marburg, Germany

11 ^d MPRG Language Cycles, Max Planck Institute for Human Cognitive and Brain Sciences,
12 Leipzig, Germany

13 ^e Clinic for Phoniatics and Pedaudiology, University Hospital Münster, Germany

14

15 * Corresponding author: srosorio@uc.cl

16 + Shared senior authors

17

18 ¹ Current affiliation: Department of Neurology and Athinoula A. Martinos Center for
19 Biomedical Imaging, Massachusetts General Hospital, Harvard Medical School, Boston, MA.

20 **Abstract**

21 During daily communication, visual cues such as gestures accompany the speech signal and
22 facilitate semantic processing. However, how gestures impact lexical retrieval and semantic
23 prediction, especially in a naturalistic setting, remains unclear. Here, participants watched a
24 naturalistic multimodal narrative, where an actor narrated a story and spontaneously produced
25 co-speech gestures. For all content words, word frequency and lexical surprisal were
26 regressed against the EEG using temporal response functions (TFRs), which were fitted
27 separately, additively, and interactively for words accompanied and not accompanied by
28 gestures. Results from our analyses suggest a robust modulation effect of gesture on the
29 frequency-dependent regression N400. Besides, we also observed some evidence of
30 modulative effect of gesture on the surprisal-N400 effect based on the single-predictor model.
31 Our finding thus suggests that, on a neural level, the presence of co-speech gestures facilitates
32 lexical retrieval and potentially semantic prediction during the processing of naturalistic
33 multimodal stimuli.

34

35 **Keywords:** multimodality, co-speech gestures, surprisal, word frequency, regression ERP,
36 lexical retrieval, semantic prediction, N400

37

38

39 **Introduction**

40 Most human communication is essentially multimodal—besides speech, we use a range of
41 visual signals such as eye-gaze, body orientation, and hand gestures, to convey meaning and
42 social communicative intent. Multimodality has been proposed to facilitate information
43 transmission (Holler & Levinson, 2019; McNeill, 2008). Amongst the modalities, hand
44 gestures (hereafter referred to as gestures) stand out as particularly unique and effective
45 means of expressing complex meaning together with accompanied speech (He et al., 2015;
46 Kelly et al., 2010; Özyürek, 2014; Özyürek et al., 2007). More importantly, they are reported
47 to have a positive impact on speech perception and auditory sentence comprehension, by
48 facilitating speech processing at various levels (Alibali & Kita, 2010; Bosker & Peeters,
49 2021; Cuevas et al., 2019; Drijvers & Özyürek, 2017; Holle et al., 2012; Kelly et al., 2010; Y.
50 Zhang et al., 2021).

51 Recent electrophysiological studies suggest that gestures modulate amplitudes of
52 evoked activities for both speech perception and sentence processing. At lower perceptual
53 levels, co-speech gestures modulate the early N1–P2 components when single words are
54 being processed (Kelly et al., 2004; Sun et al., 2021). At the semantic level, humans
55 automatically integrate gesture and speech semantics during online processing, as reflected in
56 the N400 component (Fabbri-Destro et al., 2015; Kelly et al., 2004; Özyürek et al., 2007;
57 Willems et al., 2007; Wu & Coulson, 2005). Thus, increasing evidence suggests that gestures
58 modulate N400 amplitude when word or sentence semantics are being processed (He et al.,
59 2020; Holle & Gunter, 2007; Morett et al., 2020; Wang & Chu, 2013; Y. Zhang et al., 2021).

60 However notably, these previous studies have predominantly made use of factorial
61 designs where the brain response to single words and sentences is investigated under
62 carefully controlled conditions. Although such designs have long been the basis for our

63 understanding of speech processing and sentence comprehension, questions have been raised
64 about whether observed effects actually hold during ecologically valid communication
65 scenarios (Hamilton & Huth, 2020; Kandylaki & Bornkessel-Schlesewsky, 2019; Meyer et
66 al., 2020; Willems et al., 2020). More recent studies on auditory language processing have
67 therefore investigated how different lexical and contextual features modulate the N400 by
68 using experimental paradigms that employ naturalistic speech stimuli such as long auditory
69 narratives (Alday et al., 2017; Broderick et al., 2018; Goldstein et al., 2022; Sassenhagen,
70 2019; Yan & Jaeger, 2020). In line with the classic sentence processing literature using
71 factorial experiments (Kutas & Federmeier, 2011; Van Petten & Kutas, 1990), these studies
72 suggest that the N400 is modulated by either the lexical frequency of words that reflects an
73 automatic, bottom-up retrieval mechanism (Sassenhagen, 2019), or by metrics that measure
74 more context-dependent predictive mechanisms, such as semantic similarity and lexical
75 surprisal (Broderick et al., 2018; Goldstein et al., 2022).

76 In the EEG literature on the semantic processing of co-speech gestures, in contrast,
77 despite heterogeneity in design, most studies either directly compared the N400 amplitude of
78 single words as accompanied by different types of gestures (e.g., Wang & Chu, 2013), or
79 measured the semantic N400 effect in response to match/mismatch between speech and
80 gestures (e.g., He et al., 2020; Özyürek et al., 2007). In this regard, these studies provide us
81 with insights into either semantic processing in general (as reflected in the N400 amplitude),
82 or rather the interplay between semantic prediction and integration as reflected in the N400
83 effect arisen from semantic mismatch (Lau et al., 2008; Nieuwland et al., 2020). As a result,
84 it remains unclear how different sub-stages of semantic processes such as lexical retrieval or
85 semantic prediction is influenced by gestures. Also importantly, the naturalistic approach as
86 implemented in the auditory speech processing literature has not yet been adapted to shed
87 light on these issues.

88 Here, we conducted an EEG study employing a naturalistic paradigm to investigate
89 whether and how co-speech gestures modulate the N400 as reflecting word-by-word lexical
90 retrieval and semantic prediction. To this end, we obtained measures of frequency (which
91 quantifies lexical retrieval) and surprisal (which quantifies a word’s unpredictability as a
92 measure of semantic prediction) for each word in a multimodal narrative. We then isolated
93 brain responses to individual content words while controlling for overlapping responses to
94 neighboring lexical items using multivariate time-resolved regression (Crosse et al., 2016;
95 Ehinger & Dimigen, 2019; Sassenhagen, 2019). For lexical retrieval, it has been shown that
96 words that occur more frequently in daily communication generally elicit less pronounced
97 N400s than rare words (Sassenhagen, 2019; Van Petten & Kutas, 1990). Word
98 unpredictability, which can be measured using human ratings (i.e., cloze probability) or
99 computational measures of conditional probability derived from Probabilistic Language
100 Models (PLMs), has in turn been also extensively used as a predictor of cognitive load during
101 semantic processing (Frank et al., 2015; Huizeling et al., 2022; Kutas & Federmeier, 2011;
102 Kutas & Hillyard, 1980; Monsalve et al., 2012; Rayner & Duffy, 1986). Here, we quantified
103 words’ unpredictability with lexical surprisal, an information-theoretic measure derived from
104 a deep neural network (GPT-2), as a measure of lexicosemantic pre-activation during
105 naturalistic language processing. GPT-2 is a pre-trained transformer-based model that can
106 estimate the unexpectedness of a lexical item considering its previous context, based on a
107 large training database of written text. GPT-2 outperforms other types of language models in
108 generating upcoming words (i.e., lexical-semantic prediction), including simple embedding
109 and recurrent-neural network models. It was also found that metabolic activity during
110 language processing is fit well by GPT-2’s lexical-semantic predictions (Schrimpf et al.,
111 2021).

112 Lexical surprisal is defined as the negative conditional log-probability of a word given
113 its preceding context. It quantifies the cognitive demand associated with the construction of
114 semantic representations on the basis of incremental probabilistic disambiguation (Hale,
115 2001; Levy, 2008). Lexical surprisal values predict reaction times during self-paced reading
116 of single sentences (Monsalve et al., 2012) and gaze duration during natural reading
117 (Goodkind & Bicknell, 2018). Neurobiologically, the N400—as a marker of cognitive effort
118 during lexical-semantic prediction—increases in amplitude with higher lexical surprisal
119 (Frank et al., 2015; Frank & Willems, 2017; Yan & Jaeger, 2020). More recently, GPT-2-
120 derived measures of surprisal have also been shown to be associated with an N400-like
121 regression effect during natural language comprehension (Goldstein et al., 2022; Heilbron et
122 al., 2022).

123 While it is well-established that both lexical frequency and surprisal correlate with the
124 amplitude of N400 component, whether and how multimodal cues can also modulate this
125 electrophysiological response during naturalistic speech processing is less clear. In a recent
126 study (Y. Zhang et al., 2021), by investigating the processing of isolated multimodal
127 sentences, the authors showed that the regression-based N400 associated with surprisal may
128 be modulated by co-speech gestures, among other multimodal cues (e.g., prosody). Taking a
129 step further, we aimed at extending this work from several novel perspectives: (1) We tested
130 if the facilitative effect of co-speech gestures can be observed in more ecologically valid
131 narratives (Willems et al., 2020). (2) We employed a time-resolved regression technique
132 (multivariate Temporal Response Functions, mTRFs) instead of conventional ERPs, which is
133 tailored to analyze EEG recordings from the processing of continuous speech (Crosse et al.,
134 2016). (3) We examined the potential facilitative effects of gestures on both lexical retrieval
135 (word frequency) and semantic prediction (GPT2-surprisal). Notably, although both measures
136 have been linked to the N400, they are usually highly correlated and may interact under

137 certain contexts (Halgren et al., 2002; Huizeling et al., 2022; Kretzschmar et al., 2015).
138 Consequently, to what extent the frequency–derived or surprisal-derived N400s may be
139 affected by gestures remains unknown.

140 We hypothesized that, during the processing of a naturalistic multimodal narrative, the
141 presence of gestures would result in a decreased amplitude of the N400 response to both
142 frequency and surprisal, analogous to the dampening of the N400 in the presence of co-
143 speech gestures in factorial experiments. Twenty participants were presented with
144 audiovisual clips in which an actor narrated a story presented in German using co-speech
145 gestures spontaneously (figure 1b). Content words in the speech stimuli were coded
146 according to whether their meanings were either repeated, emphasized, or complemented by
147 accompanying co-speech gestures (*Gesture present*) or not (*Gesture absent*), for all types of
148 gestures that are produced by the actor (see methods). Word frequency was included as the
149 corpus rank bin count from the Projekt Deutscher Wortschatz (Goldhahn et al., 2012), with
150 higher ranks (1–24) indicating less frequent words. Words not found in the corpus were
151 coded as 2 + the highest rank (Sassenhagen, 2019). The GPT-2 transformer model was used
152 to estimate surprisal values. Responses associated with lexical retrieval and semantic
153 predictions were obtained through mTRFs for frequency and surprisal. We estimated model
154 fit for two single-predictor models (one for surprisal and the other one for frequency), as well
155 as for an additive model (frequency + surprisal) and an interaction model (frequency +
156 surprisal + frequency * surprisal), with and without gestures as a categorical predictor. We
157 found that adding gesture information improved model fit regardless of model type, but the
158 additive model showed the best fit to the data. Then, by hypothesis, we set up separate
159 models for frequency and surprisal for gesture and no-gesture words to investigate how the
160 EEG correlates of retrieval and semantic prediction error scale with the presence of co-speech
161 gestures. Individual mTRF responses were modelled as the sum of the product of unknown

162 beta coefficients with individual values of frequency or surprisal per word (see figure 1a for a
163 visual summary of our analytical protocol). Results are in line with our hypotheses: words
164 that were not accompanied by gestures during a naturalistic narrative showed both
165 frequency– and surprisal-related N400 responses, which were reduced for words that were
166 accompanied by gestures. Further analysis including both frequency and surprisal as
167 regressors in the same additive model only revealed significant modulation effect of gestures
168 on frequency-dependent N400 responses. To summarize, our findings provide—for the first
169 time—clear evidence that co-speech gestures are associated with modulated frequency-
170 dependent N400 effects, thus indicating a potentially facilitative role of gesture on bottom-up
171 lexical retrieval. For surprisal, results from our surprisal-only model agree with prior
172 observations showing the facilitative effect of gesture on probabilistic disambiguation during
173 lexical-semantic prediction. However, to what extent this effect can be dissociated from
174 lexical retrieval requires further examination.

175 (Figure 1 here)

176

177 **Materials and methods**

178 *Participants*

179 Participants ($n = 20$, *mean age* = 24.1 years, *range* = 19-34 years, 14 females) were native
180 German speakers recruited from the Marburg-Giessen area, Germany. Participants were
181 right-handed, had normal hearing and normal or corrected to normal vision. None of them
182 reported any medical or neuropsychiatric condition. Participants read and signed an informed
183 consent before participating in the study. They were compensated with seven euros per hour
184 for participation. The research protocol and procedures were approved by the Ethics
185 Committee at Phillips University Marburg and were conducted in accordance with the
186 Declaration of Helsinki.

187

188 *Stimuli*

189 Participants watched 16 video clips (individual clips lasting between 1:02 and 3:31 min; story
190 duration = 32:12 minutes) of a professional male actor narrating an adapted version of the
191 story *Der Kuli Kimgun* as naturally as possible. Consent was obtained from the actor to use
192 his image for research and publication purposes. Foreign words in the original story were
193 replaced by German synonyms. We analyzed EEG responses to all content words in the story.
194 Firstly, content words whose semantic representations that were either associated,
195 emphasized, or complemented by gestures were coded as *Gesture present*; all other content
196 words were coded as *Gesture absent*. For example, in figure 1, the screenshots depict two
197 consecutive beat gestures that were used by the actor to emphasize the two “three” in “Samir
198 stayed three days and three nights on the hills”). In this case, even if the gesture onsets from
199 “stay” and offsets around “on”, we only coded the two “three”s as words in the *Gesture*

200 *present* condition, and all other words were coded as *Gesture absent*. This coding applies
201 analogously for non-referential beat gestures. For iconic, metaphoric, and emblems, it is then
202 their lexical affiliates that were coded as gesture present. The coding of the *Gesture*
203 present/absent conditions were double-checked by two independent expert coders. From the
204 total number of content words in the story ($n = 3582$), 466 were coded as *Gesture present*.
205 The actor was free to decide when and how to make use of gestures. Throughout the story, a
206 total of 493 hand gestures were conducted by the actor. Table 1 reports the frequency of each
207 type of gestures. For a sample description (screenshot and transcript) of the multimodal story,
208 please refer to Cuevas et al., 2019. In supplementary figure s1, we also illustrate a few
209 additional examples of the hand gestures and the corresponding *Gesture present* words.
210 Notably, despite prior research showing mixed effect of beat gestures on semantic processing
211 and potentially differential effect of iconic gestures and other gestures (Hintz et al., 2022;
212 Morett et al., 2020; Wang & Chu, 2013; Y. Zhang et al., 2021), for the purpose of testing the
213 facilitative role of gestures in general on semantic processing, we collapsed across all gesture
214 types for all analyses for the main analysis (see Cuevas et al., 2019 for an fMRI study using
215 the same stimuli). We nevertheless reported specific effects of both iconic and beat gestures
216 in the supplement.

217 (Table 1 about here)

218 Word frequency was included as the corpus rank bin count from the Projekt
219 Deutscher Wortschatz (Goldhahn et al., 2012), with higher ranks (1–24) indicating less
220 frequent words. Words not found in the corpus were coded as 2 + the highest rank. To model
221 word-by-word lexical-semantic processing demands, we employed a deep neural network
222 transformer-based model (GPT-2). Whereas earlier language models (e.g., n-gram models,
223 recurrent neural networks) derive surprisal from serial cumulative calculations on prior word

246 Electrophysiological data was acquired using a Brain Products 32-channel EEG system
247 (Brain Products GmbH, Gilching, Germany). Electrodes were positioned according to the 10-
248 20 international standard. EEG data was collected at sampling rate of 250Hz without any on-
249 line filters, referenced to FCz. The EEG dataset has been reported by Sassenhagen (2019) for
250 a different research question.

251

252 *Data analysis*

253 *EEG data preprocessing*

254 Data were preprocessed using a modified version of the Harvard Automated Preprocessing
255 Pipeline (Gabard-Durnam et al., 2018) together with the EEGLab toolbox (Delorme &
256 Makeig, 2004). EEG data were re-referenced to the average of electrodes TP9/10 (mastoids).
257 Bad channels were removed for later interpolation based on joint probability (Delorme &
258 Makeig, 2004). Line noise was removed with ZapLine (de Cheveigné, 2019), and data were
259 lowpass-filtered using an 10-Hz one-pass Hamming sinc FIR filter. For artifact detection and
260 rejection, data underwent Independent Component Analysis (Delorme et al., 2007).
261 Components underwent frequency-domain thresholding (Castellanos & Makarov, 2006).
262 Artifact components were then selected automatically by ICLABEL (Pion-Tonachini et al.,
263 2019), MARA ($p < 0.05$; Winkler et al., 2011) and ADJUST (Mognon et al., 2011). On
264 average across participants, 2.85 components ($SD = 1.89$) were rejected. Bad channels
265 identified initially were interpolated. Linear detrending was applied. Any remaining data
266 exceeding $50 \mu V$ with a moving window of 2 seconds were excluded from statistical analysis.
267 On average, across participants, this resulted in 3.59 % of words ($SD = 5.08$) being excluded.
268 We defined a group of representative centro-parietal electrodes (electrodes Cz, Pz, C3, C4,
269 P3, P4, CP1, CP2) as our region of interest (ROI), and a 300-500 ms time window as our time

270 window of interest. This is consistent with the topographical features and latency of the N400
271 response previously reported elsewhere in the literature.

272

273 *Multivariate time-resolved regression*

274 For multivariate time-resolved regression analyses, we used the mTRF toolbox
275 (Crosse et al., 2016). This method models brain responses using ridge-regression by fitting a
276 multivariate temporal response function (mTRF) to brain signals, which allows mapping
277 between stimulus features and neural activity. For this reason, multivariate time-resolved
278 regression is particularly well-suited to investigate brain responses to continuous, naturalistic
279 stimuli. Another important advantage of ridge-regression is that it is robust against
280 collinearity. This is particularly important because lexical surprisal and lexical frequency
281 values are highly correlated ($r = 0.64$, $p < 0.001$).

282 Brain responses were modelled for each subject as the sum of the product of unknown
283 beta coefficients with individual values per content word for word onset, lexical frequency,
284 GPT2-derived surprisal, and the interaction between frequency and surprisal, estimated as the
285 product between these two predictors. We analyzed content words only, as function words do
286 not induce much word-by-word N400 (Frank et al., 2015). All predictors except word onset
287 were z normalized before mTRF estimation. We implemented four separate encoding models,
288 one for surprisal using the formula $y \sim \text{word onset} + \text{surprisal}$, another one for frequency
289 using the formula $y \sim \text{word onset} + \text{frequency}$, one additive model for the combined effect of
290 surprisal and frequency ($y \sim \text{word onset} + \text{surprisal} + \text{frequency}$) and one last model for the
291 interaction between surprisal and frequency ($y \sim \text{word onset} + \text{surprisal} + \text{frequency} +$
292 $\text{surprisal} * \text{frequency}$). Here, importantly, the first two single-predictor models evaluate the
293 effects of frequency and surprisal separately, and how these effects were modulated by

294 gestures (see below). The latter two models consider additionally the situations where both
295 the effect of frequency and surprisal maybe dependent on, and interact with each other
296 (Dambacher et al., 2006; Huizeling et al., 2022; Payne et al., 2015; Van Petten & Kutas,
297 1990). For model fit evaluation, the four mTRF models were obtained from all content words
298 collapsed regardless of whether they were accompanied by gestures or not. Next, the model's
299 prediction accuracy was evaluated via cross-validation (see next section). Then, the four
300 models were estimated and evaluated again after adding a categorical predictor indicating the
301 presence or absence of gestures. A linear mixed effect model was used to test the effect of
302 model type and gesture information as an additional regressor on model fit.

303

304 *Model optimization and evaluation*

305 Model optimization and model evaluation were conducted via a *leave-one-out* 10-fold
306 cross-validation procedure using the mTRF toolbox (Crosse et al., 2016). For single predictor
307 models, the lambda (i.e., ridge) parameter was consistently set to zero. For all the other
308 models, optimal lambda parameters were identified by evaluating model fit for a logarithmic
309 space of 31 ridge values between 0.01 and 10 and between 100-600ms after word onset. This
310 resulted in a *partition-by-lambda-by-sensor* matrix of correlation coefficients. The optimal
311 lambda parameter was programmatically set as the mean spearman r value across the
312 electrodes within the pre-defined ROI that maximized model fit (supplementary table s1). For
313 model evaluation, the predictive power of the trained data after cross-validation was tested
314 against the test data partition, which returns a set of 27 correlation coefficients, one for each
315 channel. The mean spearman r value was then obtained for the set of electrodes within our
316 predefined ROI (supplementary table s2).

317

318 *Pair-wise comparison between gesture absent and gesture present conditions*

319 In a next step, we separately obtained mTRFs for gesture absent and gesture present
320 content words based on the four models above. This set of analyses, in addition to the model
321 fit comparisons, provides additionally information on if the presence of gesture *enhances* or
322 *modulates* the frequency– or surprisal-dependent N400s. For statistical analyses, we
323 separately obtained each subject’s median beta values for gesture present and gesture absent
324 mTRF models in the group of representative centro-parietal electrodes within our ROI, and
325 within a predefined time-window of 300-500ms. This is consistent with the topographical
326 features and latency of the N400 response previously reported elsewhere in the literature.
327 Within each model, we compared the extracted frequency– and surprisal-N400 beta
328 amplitudes by means of a Wilcoxon signed-rank test. We used median values and a non-
329 parametric test because Kolmogorov-Smirnov normality tests indicated that not all
330 distributions were normal (e.g., surprisal-dependent N400s in the gesture absent, $p = 0.02$).
331 Given our strong hypothesis for a decreased N400 amplitude in the presence of a co-speech
332 gesture for both frequency and surprisal, we opted for one-tailed testing, which increases
333 statistical power (Cho & Abe, 2013).

334

335

336 **Results**

337 To investigate the overall effect of gestures on mTRFs, we modelled brain responses to all
338 content words. We computed four different models: two single-predictor models using lexical
339 frequency and surprisal separately, one additive model for the combined effect of frequency
340 and surprisal, and a model testing the interaction of the two regressors. We then repeated this
341 procedure after adding to all models an additional categorical predictor indicating the
342 presence or absence of gestures. Because model fit varies numerically from one run to
343 another, we implemented these analyses 10 times and averaged model fit values across runs.
344 We report the model fit averaged across all subjects in table 3. Results show that model fit
345 improves when gestures are included as an additional regressor, regardless of model type.
346 Among the four models that include gestures as a categorical predictor, the additive model
347 shows the best fit. For statistical comparison between model fits, we conducted a linear
348 mixed effect model for model type and gesture information, controlling for the within-subject
349 nature of the design by including random effects for subject and the interaction of subject
350 with model type and subject with gesture information. Results revealed a statistically
351 significant main effect of gesture information ($F = 33.718$, $p = 1.36e-05$) and model type ($F =$
352 9.68 , $p = 2.94e-05$) in predicting a better model fit. No statistically significant effect was
353 found for the model type by gesture interaction. For the main effect of model type, post hoc
354 Tukey contrasts revealed significant better model fit for the additive model in comparison to
355 the frequency-only ($z = 2.97$, $p = 0.015$) and the surprisal-only ($z = 3.17$, $p = 8.25e-03$)
356 models, with no difference between the additive and the interaction models ($z = 0.425$ $p =$
357 0.974).

358 (Table 3 here)

359

383 The mTRF results for GPT-2-surprisal are illustrated in figure 3. Similar to the
384 frequency-dependent results, beta values for *Gesture absent* words are more negative than
385 beta values for *Gesture present* words within a time-window that is consistent with the
386 latency and the topography of the classic N400 responses. Again, we extracted the median
387 mTRF beta values between 300–500 ms for the same *a priori* defined centro-parietal
388 electrodes. Given our directional hypothesis, we expected the mTRF response for the *Gesture*
389 *absent* words to be more negative than mTRF for *Gesture present* words in a time window
390 consistent with the N400. A Wilcoxon signed-rank test ($z = -2.07$, signed rank = 49, $p =$
391 0.019 , one-tailed) showed that the median beta values between 300 and 500 ms for *Gesture*
392 *absent* words ($n = 20$, median = -2.43 , $SD = 3.96$) were significantly more negative than the
393 median beta values for *Gesture present* words ($n = 20$, median = 0.085 , $SD = 10.15$), with a
394 moderate effect size ($r = 0.46$). Similar to lexical frequency, conventional ERP effects for
395 high– vs. low-surprisal words (with median split), both collapsed across and within the
396 *Gesture present* and *Gesture absent* conditions, are reported in the supplementary figure s3.

397 (Figure 3 here)

398

399 Notably, within the same time window and for identical electrodes, the modulation
400 effects for the *Gesture absent* and *Gesture present* comparisons are statistically significant for
401 both frequency and GPT-2 surprisal, although the beta coefficients appear to be larger for the
402 frequency-dependent effect by visual inspection. Thus, we compared if gestures elicit a
403 stronger modulatory effect for the regression-based N400 for lexical frequency than for
404 surprisal. To this end, we conducted a Wilcoxon signed-rank test to directly compare the
405 median difference between *Gesture present* and *Gesture absent* responses for frequency
406 (median = 9.54 , $SD = 9.77$) and surprisal (median = 5.77 , $SD = 9.19$). We found that this

407 difference is indeed significant ($z = 2.20$, signed rank = 164, $p = 0.028$, two-tailed). Thus, it
408 appears that gesture modulates the frequency-derived N400 more strongly than the surprisal-
409 derived N400.

410 In a further step, based on the fact that the additive model outperforms the single
411 predictor models, we additionally compared the regression coefficient between the Gesture
412 present and absent conditions for both frequency and surprisal, when they were both entered
413 as regressors in the additive model. Results of this analysis are illustrated in figure 4. Here,
414 beta values differ significantly between *Gesture present* and *Gesture absent* conditions for
415 lexical frequency ($z = -2.967$, signed rank = 25, $p = 0.0015$) with a moderate effect size ($r =$
416 0.66), whereas no significant difference was observed for surprisal ($z = 0.952$, signed rank =
417 130, $p = 0.829$, $r = 0.21$). This analysis corroborates the effect of gesture on frequency-
418 dependent N400s. Here, however, in comparison to the single predictor effects, the scalp
419 distribution of the frequency effect appears to be more anterior. This pattern seems to be
420 reminiscent to the report on the effect of gesture on speech processing being more anterior
421 (Kandana Arachchige et al., 2021). Therefore, based on the additive model together with the
422 gesture information (with gesture) binary regressor, we additionally compared the model fit
423 between our N400 parietal ROI and a set of frontal electrodes (F1/2/3/4/z, FC1/2). A t-test for
424 the difference between beta values extracted from the parietal and frontal ROIs showed that
425 this difference was not statistically significant ($t = 2.06$, $p = 0.053$).

426 (Figure 4 here)

427

428 Further, even though the model fit did not differ between the additive and the
429 interaction models, we nevertheless conducted a control analysis to investigate the effect of
430 gesture on a potential interaction between frequency and lexical frequency (figure s4,

431 supplementary materials). Results of the mTRF analysis for the interaction model suggest an
432 effect of frequency ($z = -2.86$, signed rank = 28, $p = 0.002$) for a moderate effect size ($r =$
433 0.64), where words not accompanied by gestures ($z = -3.863$, $SD = 7.25$) are associated to
434 more negative beta coefficients than words accompanied by gestures ($z = 7.70$, $SD = 17.11$).
435 However, the gesture effect on the interaction term in this model was not statistically
436 significant ($z = 1.47$, signed rank = 144, $p = 0.93$, $r = 0.33$, figure s3).

437 We also explored if iconic and beat gestures may have differential impacts of either
438 frequency-dependent or surprisal-dependent N400. This comparison, however, needs to be
439 treated with caution given the much lower numbers of data points for each gesture type. We
440 evaluated both gesture type's potential impacts also in two steps. Firstly, for model
441 comparison, for both frequency-only and surprisal-only models separately, we included a
442 binary regressor regarding the presence of either an iconic or a beat gesture, and evaluated if
443 the inclusion of iconic/beat gesture regressor improves the model fit, and if they differ from
444 each other. For the frequency model, results showed that including coding for both types of
445 gestures generally improves model fit (iconic: $t = 3.79$, $p = 0.02$, beat: $t = 2.29$, $p = 0.024$).
446 However, no difference between both types of gestures were observed ($t = 1.57$, $p = 0.137$).
447 For surprisal, including coding for both types of gestures also improved model fit (iconic: $t =$
448 2.29 , $p = 0.038$, beat: $t = 1.74$, $p = 0.015$). We also observed no difference between both
449 types of gestures ($t = 1.75$, $p = 0.096$).

450 To investigate the effect of gesture types on frequency-derived N400s, we estimated
451 two separate models for frequency, based on words that were marked as gesture present with
452 iconic gestures and beat gestures respectively. We then analyzed their corresponding mTRFs
453 to gesture present and gesture absent content words for each gesture type. Results (figure s5,
454 supplementary materials) indicated that N400s are similarly modulated by both gesture types,

455 as both iconic gesture absent words (median = -5.37, SD = 3.96) and beat gesture absent
456 words (median = -5.37, SD = 3.96) are associated with more negative beta coefficients than
457 their gesture present counterparts (iconic: median = 7.03, SD = 13.33; beat: median = -0.12,
458 SD = 15.30). Both effects were statistically significant (iconic: $z = -2.93$, signed rank = 26, p
459 = 0.002, $r = 0.66$; beat: $z = -2.40$, signed rank = 40, p value = 0.008, $r = 0.54$, figure s4). We
460 repeated this analysis for surprisal derived N400s. For surprisal (figure s6, supplementary
461 materials), only beat gestures show a statistical effect ($z = -2.632$, signed rank = 34, $r =$
462 0.589) where gesture absent words (median = -2.43, SD = 3.96) are statistically more
463 negative than gesture present words (median = 4.50, SD = 13.24). In contrast, iconic gestures
464 do not have a statistically significant effect on the surprisal derived N400 amplitude ($z = -$
465 0.690, signed rank = 86, $p = 0.245$, $r = 0.154$).

466

467 Discussion

468 We investigated whether the presence of gestures modulates the cognitive demands
469 associated with lexical retrieval and semantic prediction during the processing of naturalistic
470 multimodal stimuli. With a set of mTRF analyses, we found that providing the gesture coding
471 significantly improves the mTRF model fit. Most importantly, extending a prior study
472 (Sassenhagen, 2019), across all models, we observed robust evidence that co-speech gestures
473 reduced the amplitude of the frequency-dependent N400, thus suggesting a facilitative role of
474 gestures on lexical retrieval. This finding significantly elaborates the semantic processing
475 literature of multimodal language: prior studies typically employ classic semantic violation
476 paradigms (He et al., 2020; Morett et al., 2020; Wang & Chu, 2013), or disambiguation
477 paradigms (Holle & Gunter, 2007) to derive the N400 effect. Consequently, although results
478 from these studies speak strongly for a facilitative role of gesture, it is, more specifically,
479 semantic *integration* that benefits from the visual modality. Here, building on the well-
480 established link between lexical frequency and the N400, we showed that bottom-up lexical
481 retrieval, as indexed by the N400, benefits from a complementary visual modality, just as
482 how it interacts with sentence context (Van Petten & Kutas, 1990). Our results also align with
483 an extensive line of literature on the facilitative effect of gesture on lexical retrieval during
484 production (Hadar et al., 1998; Hadar & Butterworth, 1997; Lanyon & Rose, 2009). To our
485 knowledge, this is the first time that an interaction between frequency and co-speech gestures
486 was observed, although directly from a naturalistic paradigm.

487 Regarding semantic prediction, GPT-2 surprisal indeed models an apparent N400
488 during the processing of a naturalistic multimodal narrative. This is in line with studies
489 showing a surprisal N400 effect using auditory-only stimuli (Frank et al., 2015; Frank &
490 Willems, 2017; Yan & Jaeger, 2020) and corroborate a similar effect of surprisal during the

491 processing of multimodal stimuli (Y. Zhang et al., 2021). We extended these findings by
492 using speech stimuli from a long and continuous multimodal narrative rather than single
493 sentences, thus establishing the modulatory effect of co-speech gestures on the N400 during
494 naturalistic comprehension. Similar to more recent studies, we also implemented a time-
495 resolved version of regression-based ERPs, namely multivariate Temporal Response
496 Functions, which allowed us to unmix the ERP responses to words of interest from that to
497 preceding and succeeding words in the continuous EEG signal (Crosse et al., 2016). With this
498 in mind, our findings imply that surprisal can quantify a word's unpredictability not only at
499 the level of single lexical items or sentences (Hale, 2001; Levy, 2008), but also during word-
500 by-word processing of naturalistic narratives that are multimodal in nature.

501 However notably, our analyses in the additive and interactive models only showed a
502 significant effect of gesture on frequency-dependent N400, but not on the surprisal-dependent
503 N400. This finding may be considered evidence of the modulation effect of gestures on
504 surprisal being potentially dependent on frequency. In the literature on unimodal (visual or
505 auditory) naturalistic language processing, word frequency is commonly, but not always
506 input as a covariate when modelling the effect of context-dependent measures such as
507 surprisal or semantic similarity (Armeni et al., 2019; Weissbart et al., 2019; Willems et al.,
508 2016, but see Broderick et al., 2018). In the multimodal language processing literature, to
509 date, the effect of gestures on semantic prediction (as indexed by surprisal) has been
510 investigated in only one recent study (Y. Zhang et al., 2021), but there the potential effect of
511 word frequency was not controlled for. In our study, the null effect for surprisal in the
512 *frequency + surprisal* model does not necessarily imply no effect of gestures on semantic
513 prediction: for example, the absence of a significant effect of surprisal in the additive model
514 may be alternatively explained by the fact that frequency captures the shared variance
515 between frequency and surprisal; and after all, in the single-predictor model, we still

516 observed significant effect of gesture on the surprisal-N400. Thus, although our results are in
517 accordance with prior studies in suggesting a modulatory effect of gestures on the N400s
518 (which may reflect semantic processing in general), they would need further validation,
519 especially regarding how lexical retrieval, semantic prediction, and semantic integration are
520 interactively affected by gestures.

521 An important contribution of the current study is that the effect of gestures on the
522 frequency- and/or surprisal-dependent N400 is observed in a multimodal narrative (Willems
523 et al., 2020). This extends previous results from factorial studies (Fabbri-Destro et al., 2015;
524 Kelly et al., 2004; Özyürek et al., 2007; Willems et al., 2007; Wu & Coulson, 2005) to the
525 processing of more naturalistic stimuli. A possible reason for this facilitative effect is that
526 gesture onset in the current experiment preceded the onset of critical words, preactivating
527 semantic representations or facilitating lexical retrieval, and thus alleviating the cognitive
528 demand associated to the decoding of meaning (He et al., 2020; Maess et al., 2016; Szewczyk
529 & Schriefers, 2018; ter Bekke et al., 2020; Y. Zhang et al., 2021) . This effect could be
530 especially highlighted in a more naturalistic setting where the actor of the video was able to
531 freely produce the spontaneous gestures during recording. Thus, it could be that there is a
532 bias to use gestures for words that are contextually more salient (Pouw et al., 2021; Trujillo et
533 al., 2021), or that gestures co-occur with enhanced articulation, or even with a slow-down of
534 the speaking rate, both of which could potentially facilitate semantic processing (e.g.
535 Broderick et al., 2018). Consequently, when perceived by the comprehender, the onset of a
536 gesture would automatically signal the ease of processing of the upcoming key word and
537 would jointly reduce the effort of lexical retrieval and/or semantic prediction *together* with
538 more enhanced articulation. However, as we did not control for these factors (e.g., artificially
539 cover the actor's mouth) for the purpose of maintaining maximal naturalness of the stimuli,
540 these potential confounds would need to be further assessed in more controlled conditions.

541 Alternatively, from a neurobiological perspective, it has been suggested that a left-lateralized,
542 modality-independent system exists in anterior and posterior temporal regions that maps
543 semantic information into common conceptual representations (Andric et al., 2013; Straube et
544 al., 2012, 2013; Xu et al., 2009, but see Jouravlev et al., 2019 for an alternative view).
545 Therefore, it could also be possible that the facilitative effect of gestures during
546 lexicosemantic retrieval/prediction indexed by the attenuation of the frequency and surprisal
547 N400 responses reflects the integration of matching acoustic and visual symbolic
548 representations encoded by multisensory neuronal populations in this supramodal semantic
549 system. This would be similar to what has been documented for low-level perceptual features
550 during audiovisual speech perception (Park et al., 2018) and to additive and supra-additive
551 effects during multisensory integration (Stein & Stanford, 2008). Both scenarios, however,
552 highlight the relevance of speech gestures in human communication and raise interesting
553 questions about the evolutionary origins and relevance of their facilitative effect.

554 Our study has a number of limitations. First, for practical considerations, we did not
555 set up an auditory-only condition that could potentially serve as unimodal baseline. Although
556 this approach is used by a number of recent naturalistic language processing studies (S.
557 Zhang et al., 2022), it is nevertheless vulnerable to the lack of control of stimulus-related
558 features between the *Gesture present* and *Gesture absent* conditions. Secondly, in the current
559 study, even if we have compared iconic and beat gestures on their effects on frequency and
560 surprisal, and have found potentially differential effects on the surprisal-dependent N400s
561 (see figure s5, supplement), the interpretation of this set of analysis should be cautious
562 because of the low number of datapoints available in the current study. Here, we found that
563 both iconic and beat gestures consistently modulate the frequency-dependent N400; but for
564 surprisal, although the mTRF model fit generally improved when both iconic and beat gesture
565 coding were additionally included as regressors, the pairwise comparison between gesture

566 present and gesture absent words only showed a significant modulation effect of beat gesture
567 on the surprisal-dependent N400. Other recent literature has provided evidence of dissociable
568 effects of gesture types on electrophysiological responses using both naturalistic and factorial
569 approaches, with iconic (i.e., meaningful) gestures being associated to a reduction in the
570 amplitude of the N400 and beat gestures being associated to increased N400 effects as
571 derived from semantic violation or surprisal (Hintz et al., 2022; Y. Zhang et al., 2021). On the
572 other hand, there is another line of literature showing that the N400 amplitude of single words
573 in a sentence may still be modulated by the presence of beat gestures, suggesting a potentially
574 facilitative effect of beat gestures (Morett et al., 2020; Wang & Chu, 2013). Clearly, given
575 the mixture of current literature, future experiments are necessary to shed light on this issue.
576 Moreover, in this study we used GPT-2 for analyzing surprisal. While it was informative,
577 subsequent large language models may outperform GPT-2 due to their expanded training
578 datasets and refined architectures, and may provide better fit to the EEG data (Digutsch &
579 Kosinski, 2023; Mahowald et al., 2023; Michaelov et al., 2023). Further, in the current study
580 we interpolated bad electrodes based on a relatively sparse 32-channel system, this potential
581 distortion may be best handled with a conceptual replication with an independent dataset.
582 Lastly, like many other studies that investigate semantic processing using narratives, we did
583 not employ any behavioral tasks during the experiential procedure (e.g., Broderick et al.,
584 2018; Goldstein et al., 2022; Willems et al., 2016). As a result, the interpretation of any
585 degree of facilitation, either of lexical retrieval or semantic prediction, can be interpreted on a
586 neural level at best. For this reason, a naturalistic paradigm that most optimally combines
587 behavior and its neural substrates becomes imperative for further research (Gratton et al.,
588 2022).

589

590 **Acknowledgments**

591 This project was funded by the Deutsche Forschungsgemeinschaft (DFG), funding number
592 HE8029/2-1, the von Behring-Röntgen-Stiftung (funding number 59-0002, 64-0001), and the
593 Excellence Program ‘The Adaptive Mind’ of the Hessian Ministry of Higher Education. SO
594 received funding from Agencia Nacional de Investigación y Desarrollo (ANID), national grant
595 for doctoral studies N° 21181786.

596

597 **Competing interests**

598 The authors declare no competing interests.

599

600 **Data availability**

601 Data will be made available upon request by contacting YH at yifei.he@staff.uni-marburg.de.

602

603 **Author contributions**

604 SO implemented EEG preprocessing and data analysis scripts, implemented multivariate
605 regression analyses, created figures, interpreted results, and wrote and edited the manuscript.

606 BS designed the experiment, acquired funding, and reviewed the manuscript. LM

607 implemented multivariate regression analyses, GPT2 lexical-surprisal analyses, interpreted
608 results, wrote, reviewed, and edited the manuscript. YH designed the experiment, acquired

609 data, created figures, interpreted results, and wrote, reviewed, and edited the manuscript.

610 **References**

- 611 Alday, P. M., Schlesewsky, M., & Bornkessel-Schlesewsky, I. (2017). Electrophysiology
612 reveals the neural dynamics of naturalistic auditory language processing: Event- related
613 potentials reflect continuous model updates. *ENeuro*, 4(6).
614 <https://doi.org/10.1523/ENEURO.0311-16.2017>
- 615 Alibali, M. W., & Kita, S. (2010). Gesture highlights perceptually present information for
616 speakers. *Gesture*, 10(1), 3–28. <https://doi.org/10.1075/gest.10.1.02ali>
- 617 Andric, M., Solodkin, A., Buccino, G., Goldin-meadow, S., Rizzolatti, G., & Small, S. L.
618 (2013). Brain function overlaps when people observe emblems, speech, and grasping.
619 *Neuropsychologia*, 51(8), 1619–1629.
620 <https://doi.org/10.1016/j.neuropsychologia.2013.03.022>
- 621 Armeni, K., Willems, R. M., van den Bosch, A., & Schoffelen, J. M. (2019). Frequency-
622 specific brain dynamics related to prediction during language comprehension.
623 *NeuroImage*, 198(September 2018), 283–295.
624 <https://doi.org/10.1016/j.neuroimage.2019.04.083>
- 625 Bosker, H. R., & Peeters, D. (2021). Beat gestures influence which speech sounds you hear.
626 *Proceedings of the Royal Society B: Biological Sciences*, 288(1943).
627 <https://doi.org/10.1098/rspb.2020.2419>
- 628 Broderick, M. P., Anderson, A. J., Di Liberto, G. M., Crosse, M. J., & Lalor, E. C. (2018).
629 Electrophysiological correlates of semantic dissimilarity reflect the comprehension of
630 natural, narrative speech. *Current Biology*, 28(5), 803-809.e3.
631 <https://doi.org/10.1016/j.cub.2018.01.080>

632 Castellanos, N. P., & Makarov, V. A. (2006). Recovering EEG brain signals: Artifact
633 suppression with wavelet enhanced independent component analysis. *Journal of*
634 *Neuroscience Methods*, 158(2), 300–312.
635 <https://doi.org/10.1016/j.jneumeth.2006.05.033>

636 Cho, H. C., & Abe, S. (2013). Is two-tailed testing for directional research hypotheses tests
637 legitimate? *Journal of Business Research*, 66(9), 1261–1266.
638 <https://doi.org/10.1016/j.jbusres.2012.02.023>

639 Crosse, M. J., Liberto, G. M. Di, Bednar, A., & Lalor, E. C. (2016). The Multivariate
640 Temporal Response Function (mTRF) toolbox : A MATLAB toolbox for relating
641 neural signals to continuous stimuli. *Frontiers in Human Neuroscience*, 10(November),
642 1–14. <https://doi.org/10.3389/fnhum.2016.00604>

643 Cuevas, P., Steines, M., He, Y., Nagels, A., Culham, J., & Straube, B. (2019). The facilitative
644 effect of gestures on the neural processing of semantic complexity in a continuous
645 narrative. *NeuroImage*, 195(March), 38–47.
646 <https://doi.org/10.1016/j.neuroimage.2019.03.054>

647 Dambacher, M., Kliegl, R., Hofmann, M., & Jacobs, A. M. (2006). Frequency and
648 predictability effects on event-related potentials during reading. *Brain Research*,
649 1084(1), 89–103. <https://doi.org/10.1016/j.brainres.2006.02.010>

650 de Cheveigné, A. (2019). ZapLine: a simple and effective method to remove power line
651 artifacts. *NeuroImage*, 1(1), 1–13. <https://doi.org/http://dx.doi.org/10.1101/782029>

652 Delorme, A., & Makeig, S. (2004). EEGLAB: An open source toolbox for analysis of single-
653 trial EEG dynamics including independent component analysis. *Journal of Neuroscience*
654 *Methods*, 134(1), 9–21. <https://doi.org/10.1016/j.jneumeth.2003.10.009>

655 Delorme, A., Sejnowski, T., & Makeig, S. (2007). Enhanced detection of artifacts in EEG
656 data using higher-order statistics and independent component analysis. *NeuroImage*,
657 34(4), 1443–1449. <https://doi.org/10.1016/j.neuroimage.2006.11.004>

658 Digutsch, J., & Kosinski, M. (2023). Overlap in meaning is a stronger predictor of semantic
659 activation in GPT-3 than in humans. *Scientific Reports*, 13(1), 1–7.
660 <https://doi.org/10.1038/s41598-023-32248-6>

661 Drijvers, L., & Özyürek, A. (2017). Visual context enhanced: The joint contribution of iconic
662 gestures and visible speech to degraded speech comprehension. *Journal of Speech,*
663 *Language, and Hearing Research*, 60(January), 212–222.
664 https://doi.org/10.1044/2016_JSLHR-H-16-0101

665 Ehinger, B. V., & Dimigen, O. (2019). Unfold: An integrated toolbox for overlap correction,
666 non-linear modeling, and regression-based EEG analysis. *PeerJ*, 2019(10), 1–33.
667 <https://doi.org/10.7717/peerj.7838>

668 Fabbri-Destro, M., Avanzini, P., De Stefani, E., Innocenti, A., Campi, C., & Gentilucci, M.
669 (2015). Interaction between words and symbolic gestures as revealed by N400. *Brain*
670 *Topography*, 28, 591–605. <https://doi.org/10.1007/s10548-014-0392-4>

671 Frank, S. L., Otten, L. J., Galli, G., & Vigliocco, G. (2015). The ERP response to the amount
672 of information conveyed by words in sentences. *Brain and Language*, 140, 1–11.
673 <https://doi.org/10.1016/j.bandl.2014.10.006>

674 Frank, S. L., & Willems, R. M. (2017). Word predictability and semantic similarity show
675 distinct patterns of brain activity during language comprehension. *Language, Cognition*
676 *and Neuroscience*, 32(9), 1192–1203. <https://doi.org/10.1080/23273798.2017.1323109>

- 677 Gabard-Durnam, L. J., Leal, A. S. M., Wilkinson, C. L., & Levin, A. R. (2018). The harvard
678 automated processing pipeline for electroencephalography (HAPPE): Standardized
679 processing software for developmental and high-artifact data. *Frontiers in Neuroscience*,
680 *12*, 97. <https://doi.org/10.3389/fnins.2018.00097>
- 681 Goldhahn, D., Eckart, T., & Quasthoff, U. (2012). Building large monolingual dictionaries at
682 the leipzig corpora collection: From 100 to 200 languages. *Proceedings of the 8th*
683 *International Conference on Language Resources and Evaluation, LREC 2012*, 759–
684 765.
- 685 Goldstein, A., Zada, Z., Buchnik, E., Schain, M., Price, A., Aubrey, B., Nastase, S. A., Feder,
686 A., Emanuel, D., Cohen, A., Jansen, A., Gazula, H., Choe, G., Rao, A., Kim, C., Casto,
687 C., Fanda, L., Doyle, W., Friedman, D., ... Hasson, U. (2022). Shared computational
688 principles for language processing in humans and deep language models. *Nature*
689 *Neuroscience*, *25*(3), 369–380. <https://doi.org/10.1038/s41593-022-01026-4>
- 690 Goodkind, A., & Bicknell, K. (2018). Predictive power of word surprisal for reading times is
691 a linear function of language model quality. *Proceedings of the 8th Workshop on*
692 *Cognitive Modeling and Computational Linguistics*, 10–18.
693 <https://doi.org/10.18653/v1/w18-0102>
- 694 Gratton, C., Nelson, S. M., & Gordon, E. M. (2022). Brain-behavior correlations: Two paths
695 toward reliability. *Neuron*, *110*(9), 1446–1449.
696 <https://doi.org/10.1016/j.neuron.2022.04.018>
- 697 Hadar, U., & Butterworth, B. (1997). Iconic gesture, imagery and word retrieval in speech.
698 *Semiotica*, *115*, 147–172.
- 699 Hadar, U., Wenkert-Olenik, D., Krauss, R., & Soroker, N. (1998). Gesture and the processing

700 of speech: Neuropsychological evidence. *Brain and Language*, 62(1), 107–126.

701 Hale, J. T. (2001). A probabilistic earley parser as a psycholinguistic model. *Proceedings of*
702 *the Second Meeting of the North American Chapter of the Association for*
703 *Computational Linguistics on Language Technologies*, 1–8.
704 <https://doi.org/10.3115/1073336.1073357>

705 Halgren, E., Dhond, R. P., Christensen, N., Van Petten, C., Marinkovic, K., Lewine, J. D., &
706 Dale, A. M. (2002). N400-like magnetoencephalography responses modulated by
707 semantic context, word frequency, and lexical class in sentences. *NeuroImage*, 17(3),
708 1101–1116. <https://doi.org/10.1006/nimg.2002.1268>

709 Hamilton, L. S., & Huth, A. G. (2020). The revolution will not be controlled: natural stimuli
710 in speech neuroscience. *Language, Cognition and Neuroscience*, 35(5), 573–582.
711 <https://doi.org/10.1080/23273798.2018.1499946>

712 He, Y., Gebhardt, H., Steines, M., Sammer, G., Kircher, T., Nagels, A., & Straube, B. (2015).
713 The EEG and fMRI signatures of neural integration: An investigation of meaningful
714 gestures and corresponding speech. *Neuropsychologia*, 72, 27–42.
715 <https://doi.org/10.1016/j.neuropsychologia.2015.04.018>

716 He, Y., Luell, S., Muralikrishnan, R., Straube, B., & Nagels, A. (2020). Gesture’s body
717 orientation modulates the N400 for visual sentences primed by gestures. *Human Brain*
718 *Mapping*, 41(17), 4901–4911. <https://doi.org/10.1002/hbm.25166>

719 Heilbron, M., Armeni, K., Schoffelen, J. M., Hagoort, P., & De Lange, F. P. (2022). A
720 hierarchy of linguistic predictions during natural language comprehension. *Proceedings*
721 *of the National Academy of Sciences of the United States of America*, 119(32), 1–12.
722 <https://doi.org/10.1073/pnas.2201968119>

- 723 Hintz, F., Khoe, Y. H., Strauß, A., Psomakas, A., & Holler, J. (2022). Electrophysiological
724 evidence for the enhancement of gesture-speech integration by linguistic predictability
725 during multimodal discourse comprehension. *PsyArXiv*.
726 <https://doi.org/https://doi.org/10.31234/osf.io/avudx>
- 727 Holle, H., & Gunter, T. C. (2007). The role of iconic gestures in speech disambiguation: ERP
728 evidence. *Journal of Cognitive Neuroscience*, *19*(7), 1175–1192.
729 <https://doi.org/10.1162/jocn.2007.19.7.1175>
- 730 Holle, H., Obermeier, C., Schmidt-Kassow, M., Friederici, A. D., Ward, J., & Gunter, T. C.
731 (2012). Gesture facilitates the syntactic analysis of speech. *Frontiers in Psychology*,
732 *3*(March), 1–12. <https://doi.org/10.3389/fpsyg.2012.00074>
- 733 Holler, J., & Levinson, S. C. (2019). Multimodal language processing in human
734 communication. *Trends in Cognitive Sciences*, *23*(8), 639–652.
735 <https://doi.org/10.1016/j.tics.2019.05.006>
- 736 Huizeling, E., Arana, S., Hagoort, P., & Schoffelen, J. M. (2022). Lexical Frequency and
737 Sentence Context Influence the Brain’s Response to Single Words. *Neurobiology of*
738 *Language*, *3*(1), 149–179. https://doi.org/10.1162/nol_a_00054
- 739 Jouravlev, O., Zheng, D., Balewski, Z., Le Arnz Pongos, A., Levan, Z., Goldin-Meadow, S.,
740 & Fedorenko, E. (2019). Speech-accompanying gestures are not processed by the
741 language-processing mechanisms. *Neuropsychologia*, *132*(August 2018), 107132.
742 <https://doi.org/10.1016/j.neuropsychologia.2019.107132>
- 743 Kandana Arachchige, K. G., Simoes Loureiro, I., Blekic, W., Rossignol, M., & Lefebvre, L.
744 (2021). The Role of Iconic Gestures in Speech Comprehension: An Overview of
745 Various Methodologies. *Frontiers in Psychology*, *12*(April), 1–15.

- 746 <https://doi.org/10.3389/fpsyg.2021.634074>
- 747 Kandylaki, K. D., & Bornkessel-Schlesewsky, I. (2019). From story comprehension to the
748 neurobiology of language. *Language, Cognition and Neuroscience*, 34(4), 405–410.
749 <https://doi.org/10.1080/23273798.2019.1584679>
- 750 Kelly, S. D., Kravitz, C., & Hopkins, M. (2004). Neural correlates of bimodal speech and
751 gesture comprehension. *Brain and Language*, 89(1), 253–260.
752 [https://doi.org/10.1016/S0093-934X\(03\)00335-3](https://doi.org/10.1016/S0093-934X(03)00335-3)
- 753 Kelly, S. D., Özyürek, A., & Maris, E. (2010). Two sides of the same coin: Speech and
754 gesture mutually interact to enhance comprehension. *Psychological Science*, 21(2), 260–
755 267. <https://doi.org/10.1177/0956797609357327>
- 756 Kretzschmar, F., Schlewsky, M., & Staub, A. (2015). Dissociating word frequency and
757 predictability effects in reading: Evidence from coregistration of eye movements and
758 EEG. *Journal of Experimental Psychology: Learning Memory and Cognition*, 41(6),
759 1648–1662. <https://doi.org/10.1037/xlm0000128>
- 760 Kutas, M., & Federmeier, K. D. (2011). Thirty years and counting: Finding meaning in the
761 N400 component of the event-related brain potential (ERP). *Annual Review of*
762 *Psychology*, 62, 621–647. <https://doi.org/10.1146/annurev.psych.093008.131123>
- 763 Kutas, M., & Hillyard, S. (1980). Reading senseless sentences: Brain potentials reflect
764 semantic incongruity. *Science*, 207(4427), 203–205.
- 765 Lanyon, L., & Rose, M. L. (2009). Do the hands have it? The facilitation effects of arm and
766 hand gesture on word retrieval in aphasia. *Aphasiology*, 23(7–8), 809–822.
- 767 Lau, E. F., Phillips, C., & Poeppel, D. (2008). A cortical network for semantics:

- 768 (De)constructing the N400. *Nature Reviews Neuroscience*, 9(12), 920–933.
769 <https://doi.org/10.1038/nrn2532>
- 770 Levy, R. P. (2008). Expectation-based syntactic comprehension. *Cognition*, 106(3), 1126–
771 1177. <https://doi.org/10.1016/j.cognition.2007.05.006>
- 772 Maess, B., Mamashli, F., Obleser, J., Helle, L., & Friederici, A. D. (2016). Prediction
773 signatures in the brain: Semantic pre-activation during language comprehension.
774 *Frontiers in Human Neuroscience*, 10(November), 1–11.
775 <https://doi.org/10.3389/fnhum.2016.00591>
- 776 Mahowald, K., Ivanova, A. A., Blank, I. A., Kanwisher, N., Tenenbaum, J. B., & Fedorenko,
777 E. (2023). Dissociating language and thought in large language models: a cognitive
778 perspective. *ArXiv Preprint ArXiv:2301.06627*.
- 779 McNeill, D. (2008). *Gesture and thought*. Chicago University Press.
- 780 Meyer, L., Sun, Y., & Martin, A. E. (2020). Synchronous, but not entrained: exogenous and
781 endogenous cortical rhythms of speech and language processing. *Language, Cognition
782 and Neuroscience*, 35(9), 1089–1099. <https://doi.org/10.1080/23273798.2019.1693050>
- 783 Michaelov, J. A., Bardolph, M. D., Van Petten, C. K., Bergen, B. K., & Coulson, S. (2023).
784 Strong Prediction: Language Model Surprisal Explains Multiple N400 Effects.
785 *Neurobiology of Language*, 1–29. https://doi.org/10.1162/nol_a_00105
- 786 Mognon, A., Jovicich, J., Bruzzone, L., & Buiatti, M. (2011). ADJUST: An automatic EEG
787 artifact detector based on the joint use of spatial and temporal features.
788 *Psychophysiology*, 48(2), 229–240. <https://doi.org/10.1111/j.1469-8986.2010.01061.x>
- 789 Monsalve, I. F., Frank, S. L., & Vigliocco, G. (2012). Lexical surprisal as a general predictor

790 of reading time. *EACL 2012 - 13th Conference of the European Chapter of the*
791 *Association for Computational Linguistics, Proceedings*, 398–408.

792 Morett, L. M., Landi, N., Irwin, J., & McPartland, J. C. (2020). N400 amplitude, latency, and
793 variability reflect temporal integration of beat gesture and pitch accent during language
794 processing. *Brain Research*, 1747(August), 147059.
795 <https://doi.org/10.1016/j.brainres.2020.147059>

796 Nieuwland, M. S., Barr, D. J., Bartolozzi, F., Busch-Moreno, S., Darley, E., Donaldson, D. I.,
797 Ferguson, H. J., Fu, X., Heyselaar, E., Huettig, F., Husband, E. M., Ito, A., Kazanina,
798 N., Kogan, V., Kohút, Z., Kulakova, E., Mézière, D., Politzer-Ahles, S., Rousselet, G.,
799 ... Von Grebmer Zu Wolfsturn, S. (2020). Dissociable effects of prediction and
800 integration during language comprehension: Evidence from a largescale study using
801 brain potentials. *Philosophical Transactions of the Royal Society B: Biological Sciences*,
802 375(1791). <https://doi.org/10.1098/rstb.2018.0522>

803 Özyürek, A. (2014). Hearing and seeing meaning in speech and gesture: Insights from brain
804 and behaviour. *Philosophical Transactions of the Royal Society B: Biological Sciences*,
805 369(1651). <https://doi.org/10.1098/rstb.2013.0296>

806 Özyürek, A., Willems, R. M., Kita, S., & Hagoort, P. (2007). On-line integration of semantic
807 information from speech and gesture: Insights from event-related brain potentials.
808 *Journal of Cognitive Neuroscience*, 19(4), 605–616.
809 <https://doi.org/10.1162/jocn.2007.19.4.605>

810 Park, H., Ince, R. A. A., Schyns, P. G., Thut, G., & Gross, J. (2018). Representational
811 interactions during audiovisual speech entrainment: Redundancy in left posterior
812 superior temporal gyrus and synergy in left motor cortex. *PLoS Biology*, 16(8), 1–26.

813 <https://doi.org/10.1371/journal.pbio.2006558>

814 Payne, B. R., Lee, C. L., & Federmeier, K. D. (2015). Revisiting the incremental effects of
815 context on word processing: Evidence from single-word event-related brain potentials.
816 *Psychophysiology*, 52(11), 1456–1469.

817 Pion-Tonachini, L., Kreutz-Delgado, K., & Makeig, S. (2019). ICLabel: An automated
818 electroencephalographic independent component classifier, dataset, and website.
819 *NeuroImage*, 198, 181–197.

820 Pouw, W., Wit, J. D., Bögels, S., Rasenberg, M., Milivojevic, B., & Ozyurek, A. (2021).
821 Semantically related gestures move alike: Towards a distributional semantics of gesture
822 kinematics. *International Conference on Human-Computer Interaction*, 269–287.

823 Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2019). Language
824 models are unsupervised multitask learners. *OpenAI Blog*, 1(8).
825 <http://arxiv.org/abs/2007.07582>

826 Rayner, K., & Duffy, S. A. (1986). Lexical complexity and fixation times in reading: Effects
827 of word frequency, verb complexity, and lexical ambiguity. *Memory & Cognition*, 14(3),
828 191–201. <https://doi.org/10.3758/BF03197692>

829 Ryu, S. H., & Lewis, R. L. (2021). Accounting for agreement phenomena in sentence
830 comprehension with transformer language models: Effects of similarity-based
831 interference on surprisal and attention. *CMCL 2021 - Workshop on Cognitive Modeling
832 and Computational Linguistics, Proceedings*, 61–71.
833 <https://doi.org/10.18653/v1/2021.cmcl-1.6>

834 Sassenhagen, J. (2019). How to analyse electrophysiological responses to naturalistic

835 language with time-resolved multiple regression. *Language, Cognition and*
836 *Neuroscience*, 34(4), 474–490. <https://doi.org/10.1080/23273798.2018.1502458>

837 Schrimpf, M., Blank, I. A., Tuckute, G., Kauf, C., Hosseini, E. A., Kanwisher, N.,
838 Tenenbaum, J. B., & Fedorenko, E. (2021). The neural architecture of language:
839 Integrative modeling converges on predictive processing. *Proceedings of the National*
840 *Academy of Sciences of the United States of America*, 118(45).
841 <https://doi.org/10.1073/pnas.2105646118>

842 Stein, B. E., & Stanford, T. R. (2008). Multisensory integration: Current issues from the
843 perspective of the single neuron. *Nature Reviews Neuroscience*, 9(4), 255–266.
844 <https://doi.org/10.1038/nrn2331>

845 Straube, B., Green, A., Weis, S., & Kircher, T. (2012). A Supramodal Neural Network for
846 Speech and Gesture Semantics: An fMRI Study. *PLoS ONE*, 7(11).
847 <https://doi.org/10.1371/journal.pone.0051207>

848 Straube, B., He, Y., Steines, M., Gebhardt, H., Kircher, T., Sammer, G., & Nagels, A. (2013).
849 Supramodal neural processing of abstract information conveyed by speech and gesture.
850 *Frontiers in Behavioral Neuroscience*, 7(September), 1–14.
851 <https://doi.org/10.3389/fnbeh.2013.00120>

852 Sun, J., Wang, Z., & Tian, X. (2021). Manual gestures modulate early neural responses in
853 loudness perception. *Frontiers in Neuroscience*, 15(September), 1–17.
854 <https://doi.org/10.3389/fnins.2021.634967>

855 Szewczyk, J. M., & Schriefers, H. (2018). The N400 as an index of lexical preactivation and
856 its implications for prediction in language comprehension. *Language, Cognition and*
857 *Neuroscience*, 33(6), 665–686. <https://doi.org/10.1080/23273798.2017.1401101>

858 ter Bekke, M., Drijvers, L., & Holler, J. (2020). The predictive potential of hand gestures
859 during conversation: An investigation of the timing of gestures in relation to speech.
860 *PsyArXiv*. <https://doi.org/10.31234/osf.io/b5zq7>

861 Trujillo, J., Özyürek, A., Holler, J., & Drijvers, L. (2021). Speakers exhibit a multimodal
862 Lombard effect in noise. *Scientific Reports*, *11*(1), 1–12. [https://doi.org/10.1038/s41598-](https://doi.org/10.1038/s41598-021-95791-0)
863 [021-95791-0](https://doi.org/10.1038/s41598-021-95791-0)

864 Van Petten, C., & Kutas, M. (1990). Interactions between sentence context and word
865 frequency-related brainpotentials. *Memory & Cognition*, *18*(4), 380–393.
866 <https://doi.org/10.3758/BF03197127>

867 Wang, L., & Chu, M. (2013). The role of beat gesture and pitch accent in semantic
868 processing: an ERP study. *Neuropsychologia*, *51*(13), 2847–2855.
869 <https://doi.org/10.1016/j.neuropsychologia.2013.09.027>

870 Weissbart, H., Kandylaki, K. D., & Reichenbach, T. (2019). Cortical tracking of surprisal
871 during continuous speech comprehension. *Journal of Cognitive Neuroscience*, *32*(1),
872 155–166. https://doi.org/10.1162/jocn_a_01467

873 Willems, R. M., Frank, S. L., Nijhof, A. D., Hagoort, P., & Van Den Bosch, A. (2016).
874 Prediction during natural language comprehension. *Cerebral Cortex*, *26*(6), 2506–2516.
875 <https://doi.org/10.1093/cercor/bhv075>

876 Willems, R. M., Nastase, S. A., & Milivojevic, B. (2020). Narratives for neuroscience.
877 *Trends in Neurosciences*, *43*(5), 271–273. <https://doi.org/10.1016/j.tins.2020.03.003>

878 Willems, R. M., Özyürek, A., & Hagoort, P. (2007). When language meets action: The neural
879 integration of gesture and speech. *Cerebral Cortex*, *17*(10), 2322–2333.

- 880 <https://doi.org/10.1093/cercor/bhl141>
- 881 Winkler, I., Haufe, S., & Tangermann, M. (2011). Automatic classification of artifactual
882 ICA-components for artifact removal in EEG signals. *Behavioral and Brain Functions*,
883 7(1), 30. <https://doi.org/10.1186/1744-9081-7-30>
- 884 Wu, Y. C., & Coulson, S. (2005). Meaningful gestures: Electrophysiological indices of iconic
885 gesture comprehension. *Psychophysiology*, 42(6), 654–667.
886 <https://doi.org/10.1111/j.1469-8986.2005.00356.x>
- 887 Xu, J., Gannon, P. J., Emmorey, K., Smith, J. F., & Braun, A. R. (2009). Symbolic gestures
888 and spoken language are processed by a common neural system. *Proceedings of the*
889 *National Academy of Sciences of the United States of America*, 106(49), 20664–20669.
890 <https://doi.org/10.1073/pnas.0909197106>
- 891 Yan, S., & Jaeger, T. F. (2020). (Early) context effects on event-related potentials over
892 natural inputs. *Language, Cognition and Neuroscience*, 35(5), 658–679.
893 <https://doi.org/10.1080/23273798.2019.1597979>
- 894 Zhang, S., Jixing, L., Yang, Y., & Hale, J. (2022). Decoding the silence: Neural bases of zero
895 pronoun resolution in Chinese. *Brain and Language*, 224(November 2021), 105050.
896 <https://doi.org/10.1016/j.bandl.2021.105050>
- 897 Zhang, Y., Frassinelli, D., Tuomainen, J., Skipper, J. I., & Vigliocco, G. (2021). More than
898 words: Word predictability, prosody, gesture and mouth movements in natural language
899 comprehension. *Proceedings of the Royal Society B: Biological Sciences*, 288(1955).
900 <https://doi.org/10.1098/rspb.2021.0500>

901

902 Table 1. Summary of gesture occurrences

Gesture Type	Occurrence	%
Iconic	181	36.71
Metaphoric	85	17.24
Emblematic	19	3.85
Beat	148	30.02
Deictic	60	12.17
Total	493	100

903

904

905

906

907

908

909

910

911

912

913

914 Table 2. Summary of word statistics

	N° of Words	Word duration				GPT-2 Surprisal				Word frequency			
		Mean	Median	SD	Range	Mean	Median	SD	Range	Mean	Median	SD	Range
Gesture Present	466	0.545	0.510	0.199	1.080	13.714	12.679	7.331	47.010	12.695	11	5.628	27
Gesture Absent	1491	0.491	0.460	0.195	2.810	13.056	11.538	7.932	53.972	11.674	11	5.586	27

915

916

917

918

919

920

921

922

923

924

925

926

927

928 Table 3. Model fit comparison

Model	Fit (spearman r)			
	<i>Without gesture as regressor</i>		<i>With gesture as regressor</i>	
	Mean	SD	Mean	SD
<i>Surprisal</i>	0.058	0.013	0.069	0.016
<i>Frequency</i>	0.056	0.015	0.069	0.013
<i>Additive</i>	0.066	0.017	0.076	0.015
<i>Interaction</i>	0.065	0.015	0.075	0.016

929

930

931

932

933

934

935

936

937

938

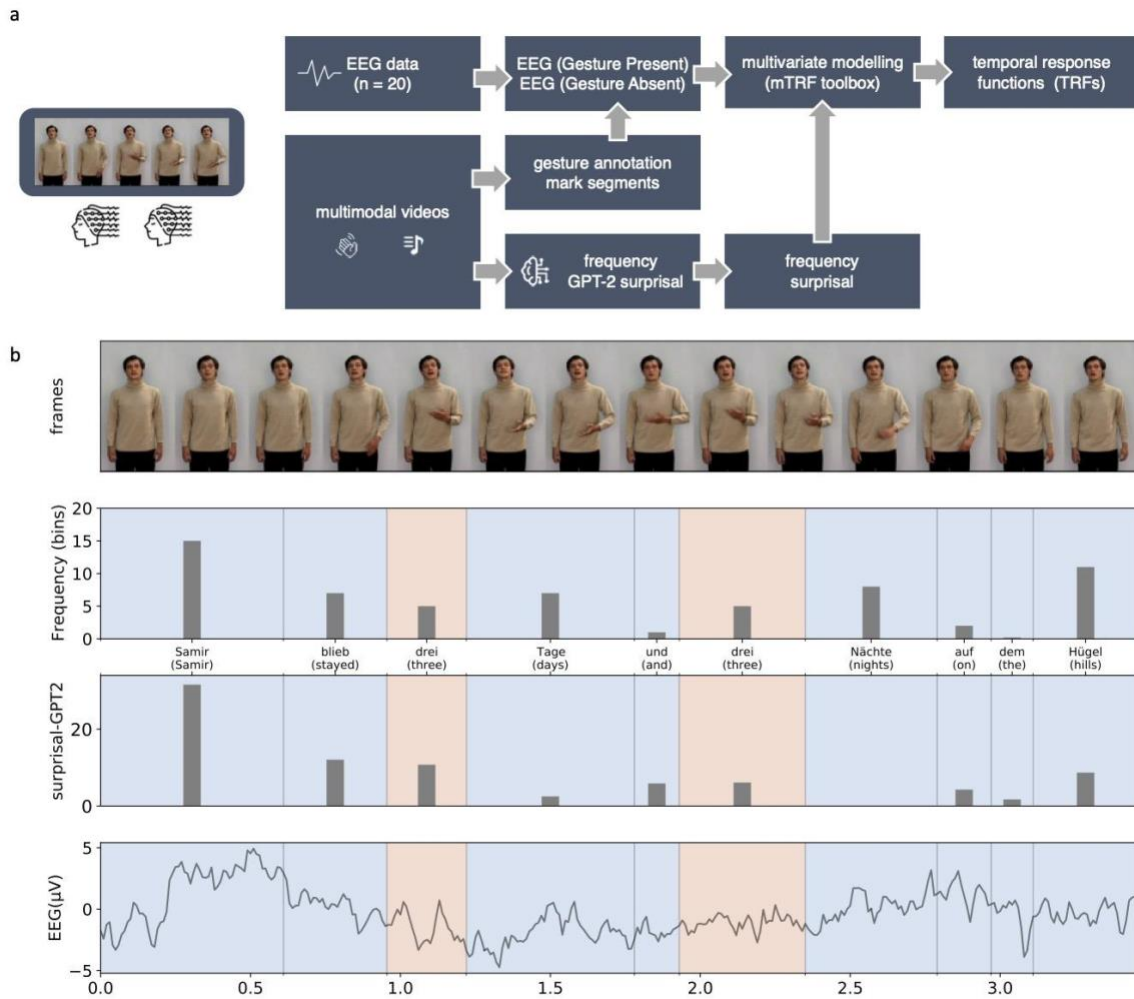
939

940

941

942 **Figures**

943 **Figure 1**



944

945

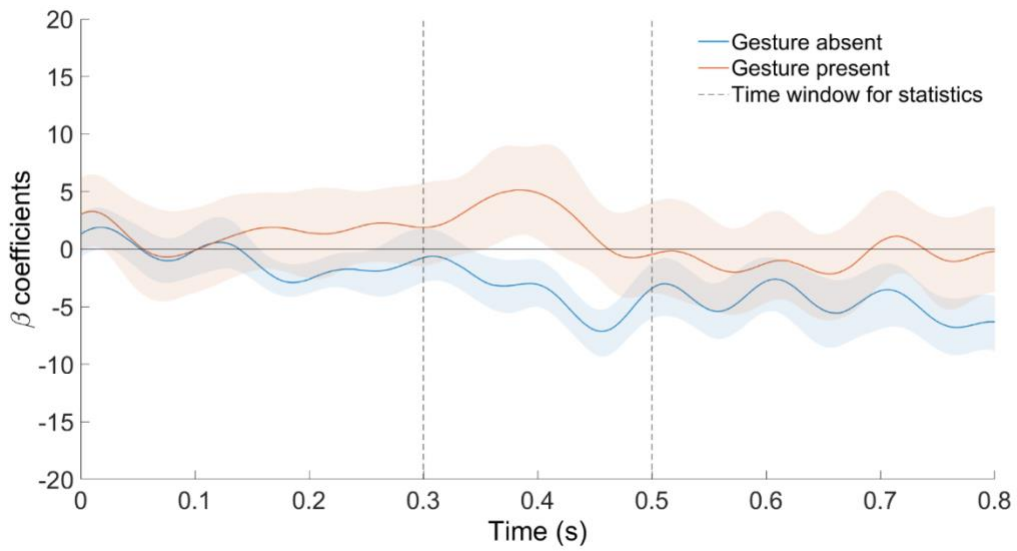
946

947

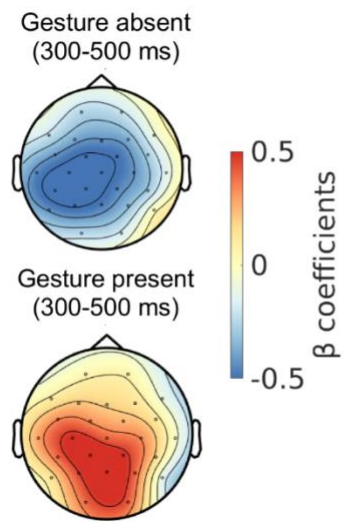
948

949

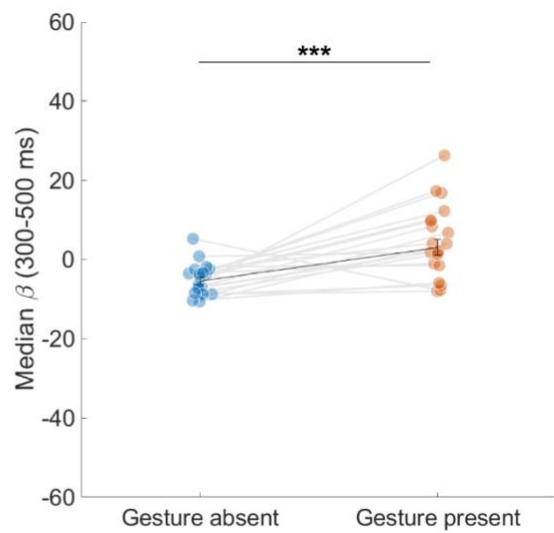
a



b



c



951

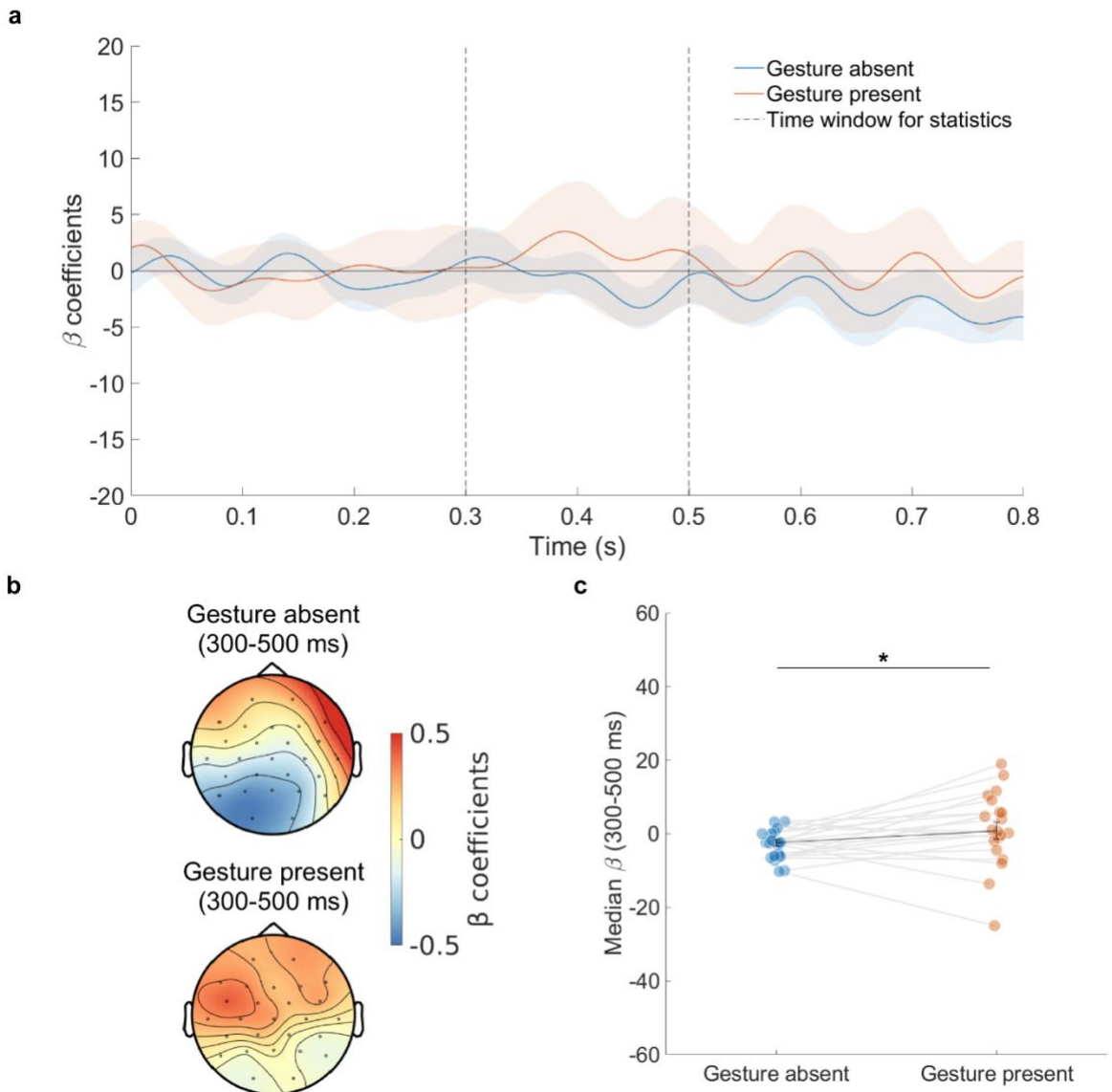
952

953

954

955

956



958

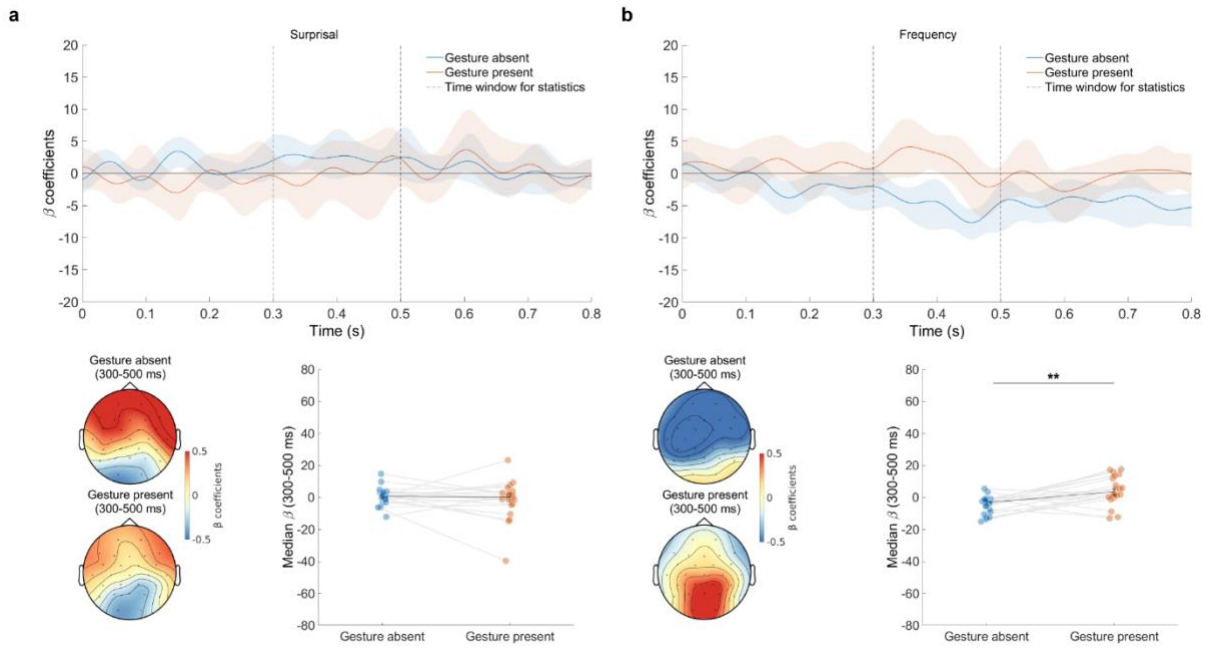
959

960

961

962

963



965

966

967

968

969

970

971

972

973

974

975

976 **Figure captions**

- 977 • **Figure 1.** Schematic representation of experimental and analytic protocol. **a.** EEG
978 data was acquired while participants were exposed to multimodal naturalistic speech
979 stimuli. Besides lexical frequency, word-by-word surprisal was obtained using the
980 GPT-2 model. Multimodal stimuli were annotated for the presence of absence of
981 gestures. EEG data were then modelled using the mTRF toolbox to obtain regression-
982 based ERPs (i.e., mTRFs) for gesture present and gesture absent data separately. **b.**
983 Sample frames from the multimodal video stimuli (top) along with corresponding
984 word-by-word frequency and GPT-2-derived surprisal values (middle) and the
985 average EEG signal for a group of centro-parietal electrodes (bottom). Words marked
986 as *Gesture present* are marked with red background.
- 987 • **Figure 2.** Results of TRF analyses for lexical frequency. **a.** Mean mTRF to *Gesture*
988 *absent* (blue) and *Gesture present* (orange) words obtained from the average of our
989 electrodes of interest. Coloured shades show the 95% bootstrapped confidence
990 intervals. Dotted lines represent our time-window of interest for statistical analyses
991 (300–500 ms). **b.** Topographical distributions of mean mTRF responses between 300–
992 500 ms. **c.** Individual and group median mTRF beta values between 300–500 ms. The
993 line and asterisks represent a statistically significant effect ($p < 0.001$, one-tailed).
- 994 • **Figure 3.** Results of TRF analyses for GPT-2 surprisal. **a.** Mean mTRF to *Gesture*
995 *absent* (blue) and *Gesture present* (orange) words obtained from the average of a
996 group of centro-parietal electrodes. Coloured shades show the 95% bootstrapped
997 confidence intervals. Dotted lines represent our time-window of interest for statistical
998 analyses (300–500 ms). **b.** Topographical distributions of mean mTRF responses
999 between 300–500 ms. **c.** Individual and group median mTRF beta values between

1000 300–500 ms. The line and asterisk represent a statistically significant effect ($p < 0.05$,
1001 one-tailed).

1002 • **Figure 4.** Results of TRF analyses for the additive model **a.** Effect of GPT-2
1003 surprisal. **b.** Effect of lexical frequency. For both panels, mTRF for *Gesture absent*
1004 (blue) and *Gesture present* (orange) words obtained as the average of an *a priori*
1005 defined group of centro-parietal electrodes. Coloured shades show the 95%
1006 bootstrapped confidence intervals. Dotted lines represent our time-window of interest
1007 for statistical analyses (300–500 ms). Topographical distributions are based on the
1008 mean mTRF responses between 300–500 ms. Point and line-plots show Individual
1009 and group median mTRF beta values between 300–500 ms. The line and asterisks
1010 represent a statistically significant effect ($p < 0.01$, one-tailed).

Supplementary materials

The role of co-speech gestures in retrieval and prediction during naturalistic multimodal narrative processing

This PDF includes:

Tables s1 and s2

Figures s1, s2, s3, s4, s5, s6

Table s1. Optimal Lamba values for additive and interaction models after cross-validation

Subject	Additive		Interaction	
	<i>Gesture absent</i>	<i>Gesture present</i>	<i>Gesture absent</i>	<i>Gesture present</i>
1	0.631	0.040	0.794	0.200
2	0.794	10.000	1.995	10.000
3	0.010	0.794	0.158	0.010
4	0.010	0.200	0.316	0.251
5	0.010	0.010	1.585	0.631
6	0.398	2.512	1.000	0.010
7	1.000	1.995	1.259	0.794
8	6.310	1.259	5.012	0.794
9	0.010	0.501	1.000	0.501
10	0.501	0.010	0.794	0.010
11	0.794	0.631	0.398	0.200
12	0.010	10.000	0.010	1.000
13	0.010	1.259	0.010	0.010
14	7.943	0.032	2.512	0.079
15	2.512	0.010	0.010	0.251
16	0.010	0.398	0.251	0.200
17	0.010	1.995	0.631	0.010
18	0.200	0.010	0.794	0.020
19	1.000	0.010	1.259	0.032
20	0.010	10.000	0.794	10.000

Table s2. Model fit values (spearman r) for the four models

Subject	Surprisal		Frequency		Additive		Interaction	
	<i>Gesture absent</i>	<i>Gesture present</i>	<i>Gesture absent</i>	<i>Gesture present</i>	<i>Gesture absent</i>	<i>Gesture present</i>	<i>Gesture absent</i>	<i>Gesture present</i>
1	0,088	-0,013	0,103	0,042	0,030	0,048	0,108	0,037
2	0,037	0,007	0,051	0,031	0,044	0,041	0,039	0,056
3	0,032	0,047	0,028	0,048	0,089	0,039	0,103	0,046
4	0,073	0,053	0,059	0,061	0,082	0,053	0,033	0,048
5	0,024	0,042	0,047	0,020	0,026	0,051	0,037	0,053
6	0,078	0,054	0,116	0,033	0,042	0,083	0,084	0,062
7	0,045	0,040	0,034	0,036	0,066	0,044	0,053	0,010
8	0,055	0,099	0,086	0,043	0,026	0,072	0,028	0,035
9	0,047	0,039	0,092	0,055	0,042	0,074	0,066	0,056
10	0,066	0,039	0,028	0,037	0,042	0,023	0,026	0,070
11	0,050	0,010	0,070	0,059	0,036	0,087	0,050	0,003
12	0,049	0,020	0,078	0,035	0,033	0,021	0,029	0,036
13	0,046	-0,008	0,064	0,054	0,036	0,004	0,093	0,048
14	0,046	0,064	0,068	0,045	0,100	0,070	0,100	0,069
15	0,057	0,019	0,039	0,050	0,075	0,028	0,065	0,102
16	0,123	0,092	0,090	0,100	0,089	0,078	0,084	0,082
17	0,041	0,032	0,033	0,045	0,135	0,039	0,110	0,067
18	0,067	0,035	0,021	0,054	0,086	0,046	0,046	0,045
19	0,022	0,052	0,072	0,060	0,039	0,075	0,074	0,066
20	0,075	0,028	0,043	0,004	0,080	0,021	0,050	0,033



Seine Arbeit war, auf die hohen Dorfpalmen zu klettern, die wie Mastbäume am Strand aufragen...

His job was to climb the tall village palms that towered like mast trees on the beach...



...oben in die Blattsprossen der Palmen ein paar Kerbschnitte zu ritzen und kleine Blechgefäße darunter zu hängen...

...carving a few notch cuts in the top of the leaf shoots of the palm trees and hanging small tin containers underneath...



...glitt Samir von seinem Palmenschaft herunter, bis der Landungssteg ausgelegt wurde und er mit den anderen Jungen Säcke voll Reis...

...Samir slid down from his palm shaft until the jetty was laid out and he and the other boys were carrying sacks full of rice...

Figure s1. Example of the gestures as produced in the multimodal narrative (left panel).

The middle panel shows the corresponding original German sentence and the *Gesture present* words (underlined). The right panel shows the English translation of the sentences and the *Gesture present* words.

Conventional ERP effects for high- vs. low-frequency words

We conducted a median-split analysis based on lexical frequency for all content words and compared the conventional ERP effects between high- and low-frequency words for (1) all words and (2) words within the *Gesture present* and the *Gesture absent* conditions. ERPs were all baseline-corrected based on the -200-0 ms time-locked to the onset of each word. We plot the ERP waveforms in figure s1. Analogous to the mTRF results, we directly extracted the ERP amplitudes between 300–500ms for centro-parietal electrodes and entered these values into a 2 x 2 parametric repeated-measures ANOVA with the factors FREQUENCY (high vs. low) and GESTURE (present vs. absent), as interaction is unable to be tested via non-parametric tests. ANOVA revealed a main effect of GESTURE ($F_{(1,19)} = 41.90, p < 0.00003$) and no main effect of frequency ($F_{(1,19)} = 1.04, p = 0.32$). The interaction between the two factors was significant ($F_{(1,19)} = 5.92, p < 0.00003$). We then conducted Wilcoxon signed-rank tests for the effect of frequency within *Gesture present* and *Gesture absent*. No effect of frequency was observed in the *Gesture absent* condition ($z = 2.25$, signed rank = 165, $p = 0.99$), but low-frequency content words (above median) showed more positive N400 amplitudes in comparison to high-frequency words ($z = -2.14$, signed rank = 47, $p = 0.015$) in the *Gesture present* condition. This effect, besides being significant, is not in accordance with the hypothesized pattern that low-frequency words should be more negative in the N400 window. Overall, the hypothesized effect of frequency was not reliably observed in the ERP results.

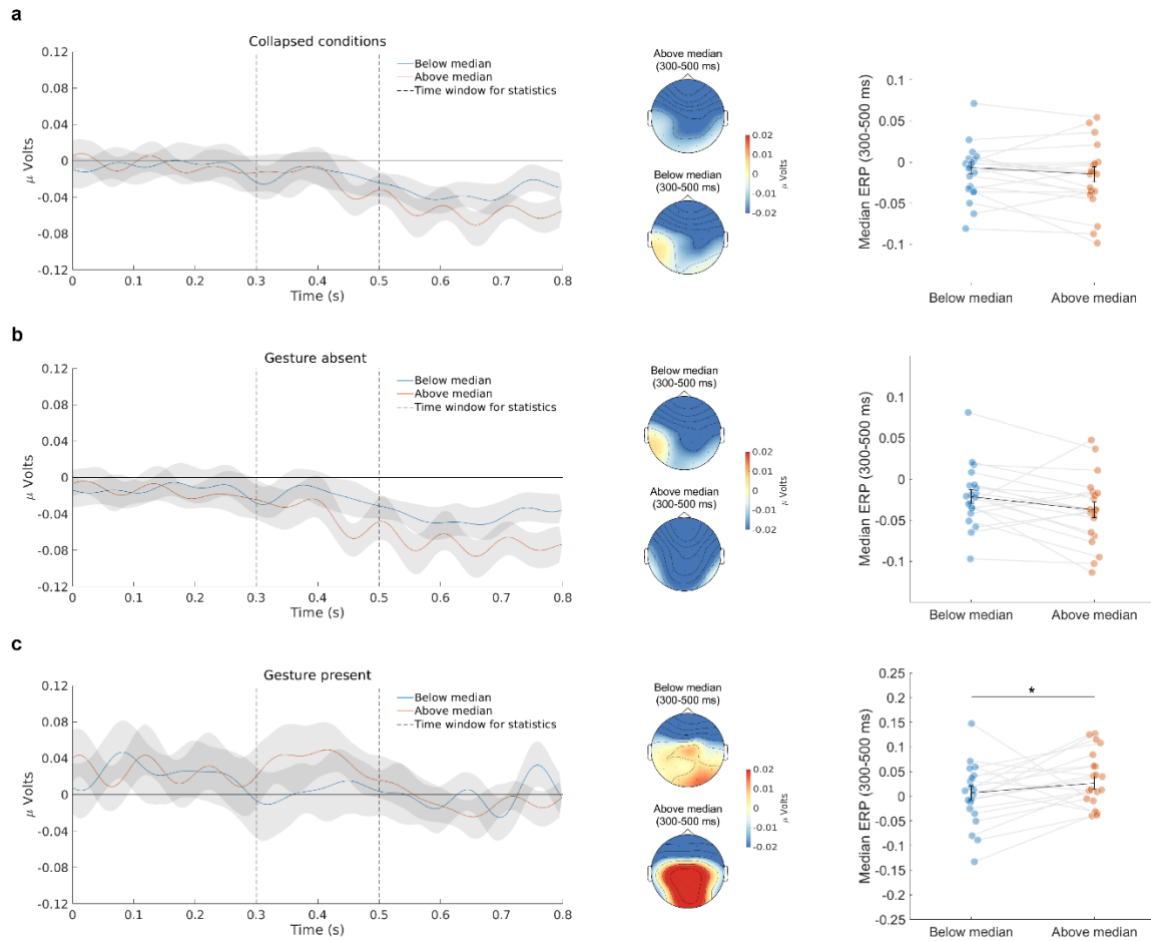


Figure s2. Results of the ERP analyses (median split) for lexical frequency for **a.** All words collapsed across *Gesture absent* and *Gesture present* conditions. **b.** Words in the *Gesture absent* condition, and **c.** Words in the *Gesture present* condition. For all panels, waveforms of low-frequency words (Above median) are depicted in red, and high-frequency words (Below median) are depicted in blue. ERPs were obtained from the average of all electrodes. Coloured shades show the 95% bootstrapped confidence intervals. Dotted lines represent our time-window of interest for statistical analyses (300-500 ms). Topographical distributions show ERP responses between 300-500 ms. Point- and line-plots show individual and group median ERP amplitudes between 300-500 ms. The line and asterisk represent a statistically significant effect ($p < 0.05$, one-tailed).

Conventional ERP effects for high- vs. low-surprisal words

For surprisal, we also conducted a median-split ERP analysis based on GPT-2 surprisal for all content words and compared the conventional ERP effects between high- and low-surprisal words for 1) all words and 2) words within the *Gesture present* and the *Gesture absent* conditions. ERPs were all baseline-corrected based on the -200-0 ms time-locked to the onset of each word. We plot the ERP waveforms in Figure s2. Analogous to the mTRF results, we directly extracted the ERP amplitudes between 300-500ms for centro-parietal electrodes and entered these values into a 2 x 2 parametric repeated-measures ANOVA with the factors SURPRISAL (high vs. low) and GESTURE (present vs. absent). ANOVA revealed a main effect of GESTURE ($F_{(1,19)} = 52.212$, $p = 7.34e-07$), but no main effect of SURPRISAL ($F_{(1,19)} = 0.078$, $p = 0.79$). The interaction between the factors was also non-significant ($F_{(1,19)} = 0.861$, $p = .365$). Non-parametric Wilcoxon signed-rank tests also showed no effect of SURPRISAL within the *Gesture absent* ($z = 1.4$, signed rank = 142, $p = 0.919$), and the *Gesture present* conditions ($z = -0.35$, signed rank = 95, $p = 0.361$). Similar to lexical frequency, ERPs on high vs. low surprisal did not exhibit interpretable patterns.

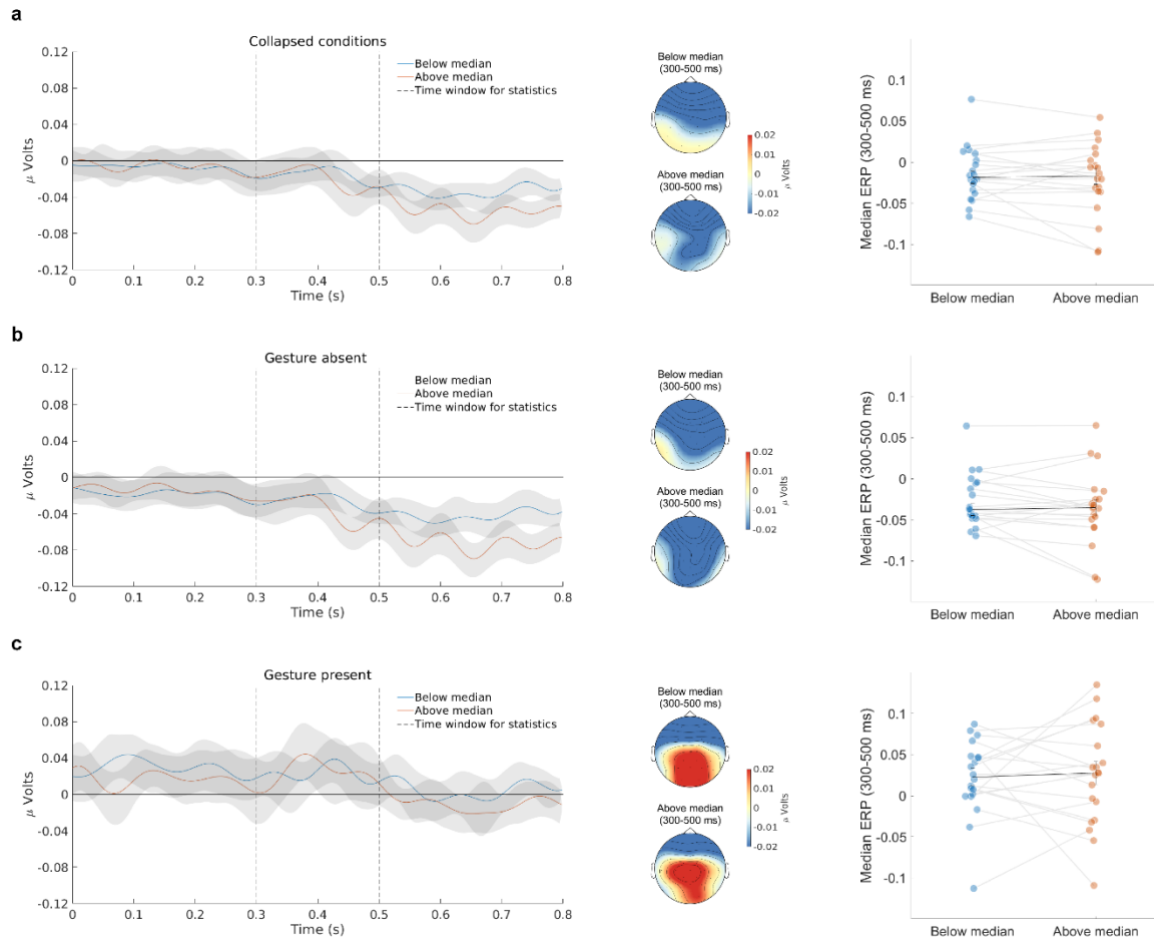


Figure s3: Results of the ERP analyses for GPT-2 surprisal for **a.** All words collapsed across *Gesture absent* and *Gesture present* conditions. **b.** Words in the *Gesture absent* condition. and **c.** Words in the *Gesture present* condition. For all panels, waveforms of high-surprisal words (Above median) are depicted in red, and low-surprisal words (Below median) are depicted in blue. ERPs were obtained from the average of all electrodes. Coloured shades show the 95% bootstrapped confidence intervals. Dotted lines represent our time-window of interest for statistical analyses (300–500 ms). Topographical distributions show ERP responses between 300–500 ms. Point- and line-plots show individual and group median ERP amplitudes between 300–500 ms.

Discussion on the median-split ERPs

Our conventional median-split ERP analyses revealed no systematic and interpretable effects for either frequency or surprisal. Our findings, at first glance, might be at odds with prior publications showing either an N400 effect of high vs. low surprisal using unimodal language stimuli e.g., visually presented sentences (Frank et al., 2015). Notably, in this study, the stimuli were presented in an RSVP manner, and only one critical word within one sentence were analysed for the ERPs. Regarding frequency, despite a long list of literature employing factorial ERP design (e.g., Van Petten and Kutas, 1990), no prior literature has used even unimodal but naturalistic stimuli with median-split to examine how the N400 ERPs vary as a function of word frequency. On the other hand, in the multimodal language processing literature, a number of studies have reported an N400 effect using conventional ERP methods (Kelly et al., 2004; Morett et al., 2020; Wang and Chu, 2013; Wu and Coulson, 2005). Similar to the study from Frank and colleagues (2015), the gesture-speech N400 studies also focused on a single critical word within sentences, so that there is no overlap of EEG time-windows between words. As a result, conventional ERP analysis focusing on the N400 window has led to reliable and interpretable data patterns. This is, however, not the case for the current study, as shown from our results. Here, the validity, and the comparability to prior studies when using conventional ERP methods is undermined by two important confounds: the temporal segments of words being analysed overlap with each other; and all these words were multimodal in nature (with facial expressions, lip movements, and where applicable, hand gestures). Thus, we argue that the conventional ERP analysis (here using the median-split approach) has limitations to answer our research questions; and the effect of

surprisal and word frequency, and their interaction with co-speech gestures is best examined via the temporally-resolved regression approach, as we reported in the main text.

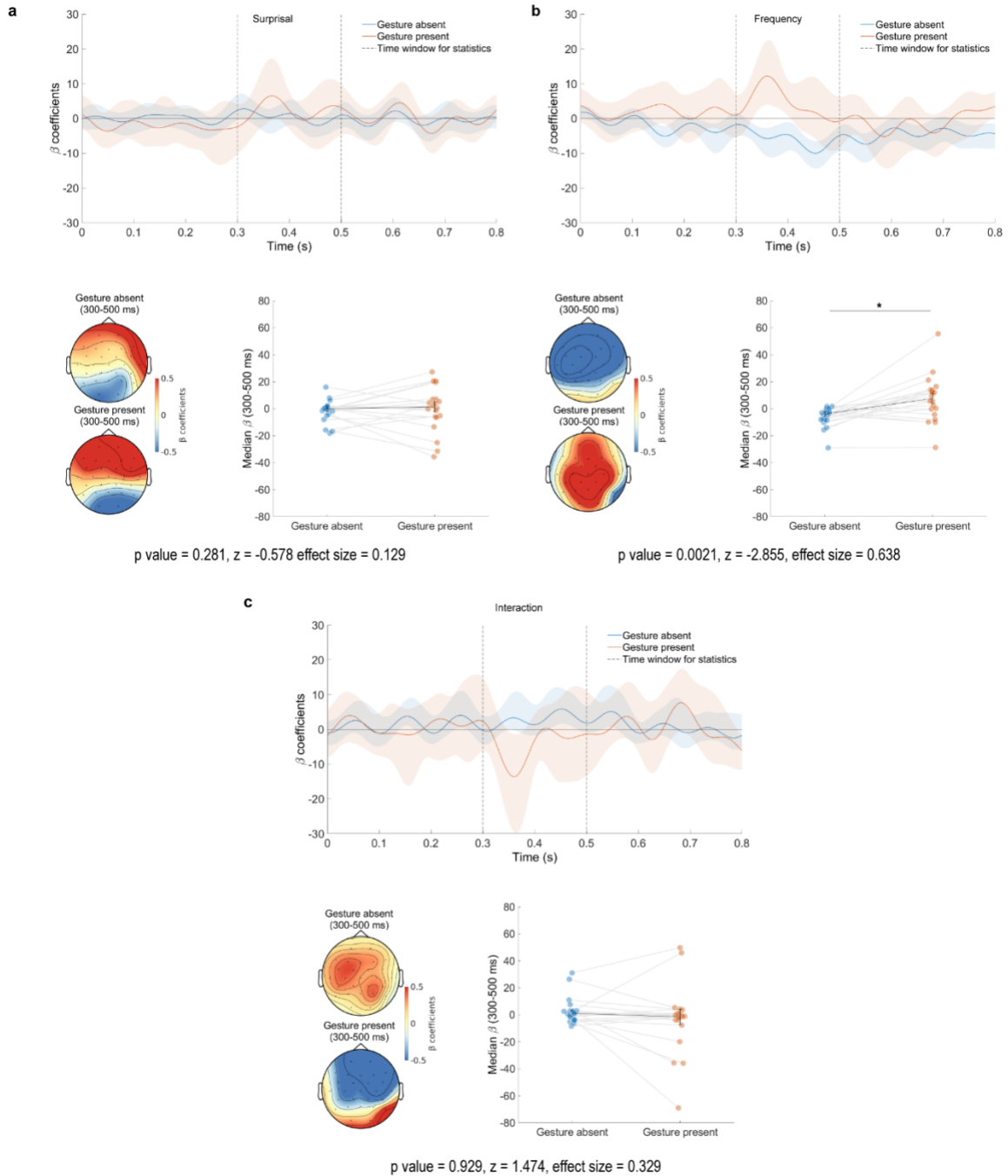


Figure s4. Results of mTRF analyses for surprisal (a), frequency (b) and their interaction (c). Mean mTRF to *Gesture absent* (blue) and *Gesture present* (orange) words obtained from the average of a group of centro-parietal electrodes with their 95% bootstrapped confidence intervals. Topographical distributions are based on the mean mTRF responses between 300–

500 ms. Point and line-plots show Individual and group median mTRF beta values between 300–500 ms. The line and asterisk represent a statistically significant effect ($p < 0.05$, one-tailed).

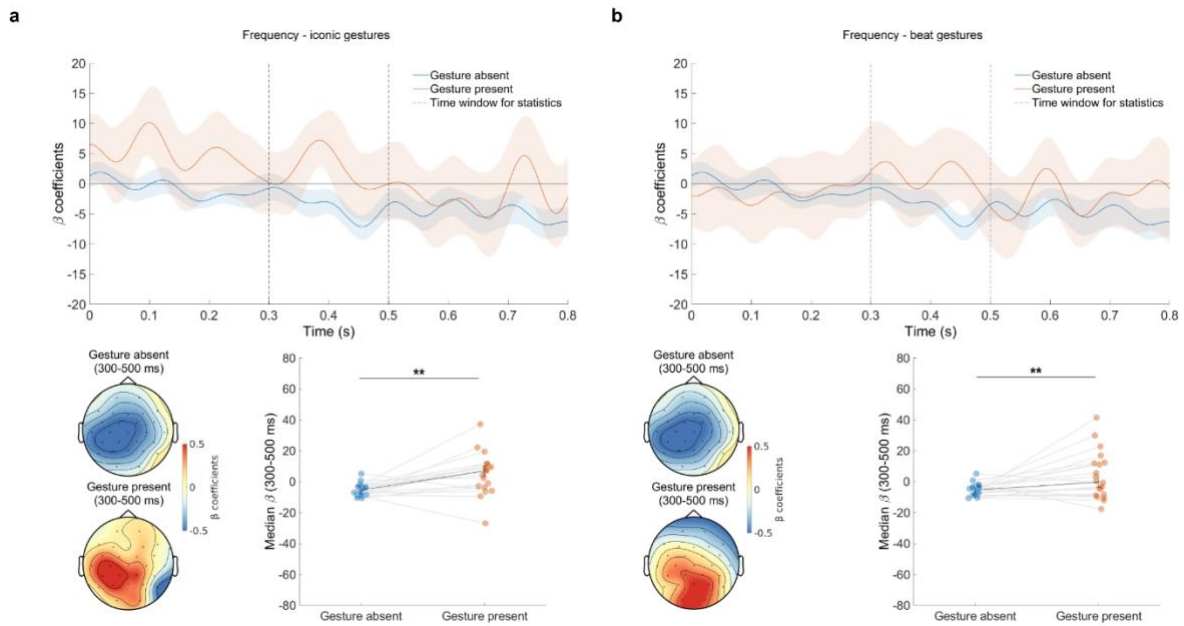


Figure s5. Results of mTRF analyses for frequency using iconic gestures (a) and beat gestures (b) only. Mean mTRF to *Gesture absent* (blue) and *Gesture present* (orange) words obtained from the average of a group of centro-parietal electrodes with their 95% bootstrapped confidence intervals. Topographical distributions are based on the mean mTRF responses between 300–500 ms. Point and line-plots show Individual and group median mTRF beta values between 300–500 ms. The line and asterisks represent a statistically significant effect ($p < 0.01$, one-tailed).

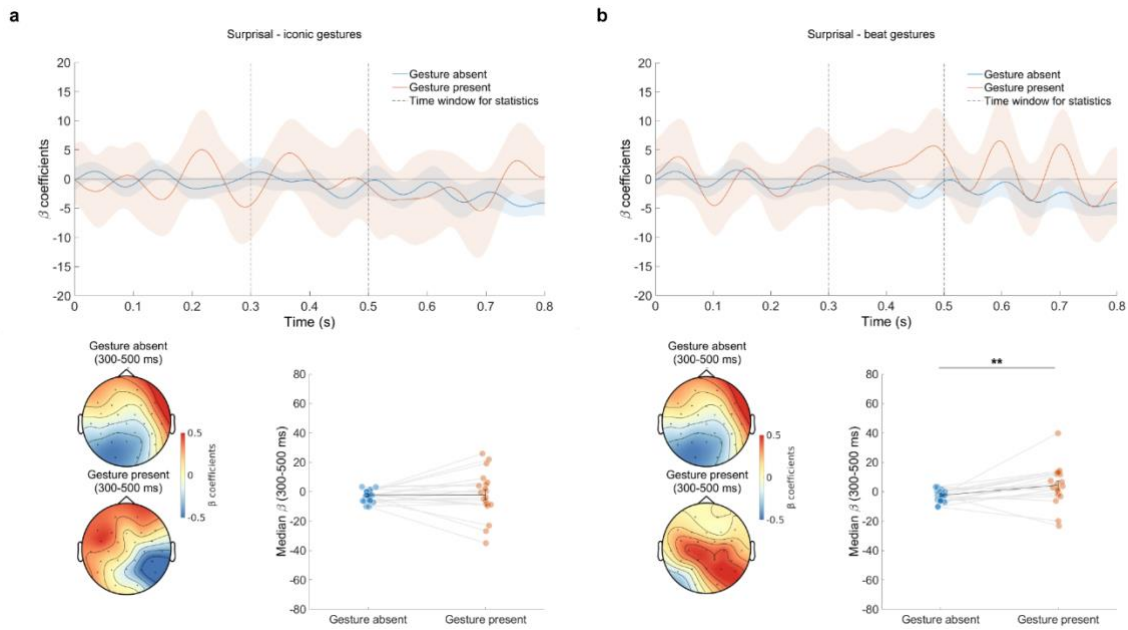


Figure s6. Results of mTRF analyses for surprisal using iconic gestures (a) and beat gestures (b) only. Mean mTRF to *Gesture absent* (blue) and *Gesture present* (orange) words obtained from the average of a group of centro-parietal electrodes with their 95% bootstrapped confidence intervals. Topographical distributions are based on the mean mTRF responses between 300–500 ms. Point and line-plots show Individual and group median mTRF beta values between 300–500 ms. The line and asterisks represent a statistically significant effect ($p < 0.01$, one-tailed).