

1 **Design and evaluation of the effectiveness of a corpus of congruent and incongruent English**  
2 **sentences for the study of event related potentials**

3 Joaquin T. Valderrama <sup>a,b,\*</sup> Elizabeth F. Beach <sup>a</sup>, Mridula Sharma <sup>b</sup>, Shivali Appaiah-Konganda <sup>b</sup>,  
4 Elaine Schmidt <sup>b,c</sup>.

5 <sup>a</sup> *National Acoustic Laboratories, Australian Hearing Hub. 16 University Avenue, Macquarie*  
6 *University, New South Wales, 2109, Sydney, Australia.*

7 <sup>b</sup> *Department of Linguistics, Australian Hearing Hub. 16 University Avenue, Macquarie University,*  
8 *New South Wales, 2109, Sydney, Australia.*

9 <sup>c</sup> *Department of Theoretical and Applied Linguistics, University of Cambridge. 9 West Road,*  
10 *Cambridge, CB3 9DP, United Kingdom.*

11 \* Corresponding author

12 Joaquin T. Valderrama  
13 National Acoustic Laboratories  
14 Australian Hearing Hub  
15 Level 5, 16 University Avenue  
16 Macquarie University NSW 2109  
17 Sydney, Australia  
18 Phone: +61 2 9412 6878  
19 Email address: [joaquin.valderrama@nal.gov.au](mailto:joaquin.valderrama@nal.gov.au), [joaquin.valderrama@mq.edu.au](mailto:joaquin.valderrama@mq.edu.au).  
20 ORCID: 0000-0002-5529-8620

21 Word count: 6,601 words.

22

23 **Abstract**

24 Objective: To design and evaluate the effectiveness of a stimulus material in eliciting the N400  
25 event related potential (ERP).

26 Design: A set of 700 semantically congruent and incongruent sentences was developed in  
27 accordance with current linguistic norms, and validated with an electroencephalography (EEG)  
28 study, in which the influence of age and gender on the N400 ERP magnitude was analysed.

29 Study sample: Forty-five normal-hearing subjects (19-57 years, 21 females) participated in the  
30 EEG study.

31 Results: The stimulus material used in the EEG study elicited a robust N400 ERP, with a  
32 morphology consistent with the literature. Results also showed no statistically significant effect  
33 of age or gender on the N400 magnitude.

34 Conclusions: The material presented in this paper constitutes the largest complete stimulus set  
35 suitable for both auditory and text-based N400 experiments. This material may help facilitate  
36 the efficient implementation of future N400 ERP studies, as well as promote standardization and  
37 consistency across studies.

38 **Keywords**

39 N400; Speech Perception; Language-Related ERPs; Semantic Violation.

40

## 41 1. Introduction

42 The analysis of event-related potentials (ERPs) has been fundamental for understanding the  
43 neural basis of language encoding in the human brain. In particular, the N400 ERP has been  
44 identified as an index of lexical and semantic processing (Kutas and Hillyard, 1980; Osterhout  
45 and Holcomb, 1992; Brown and Hagoort, 1993, Hagoort, 2008), and is widely used in the design  
46 of experiments aiming to evaluate language comprehension (Boudewyn et al., 2012;  
47 Federmeier, 2007; Kutas and Federmeier, 2000; Van Petten and Luka, 2012). The N400 is  
48 characterized as a negative deflection of the average ERP waveform that peaks around 400 ms  
49 after a semantic violation, presenting a larger magnitude at midline central-parietal sites (Kutas,  
50 1993; Duncan et al., 2009; Jerger and Martin, 2009).

51 The large body of literature around the N400 has confirmed that this ERP is associated with  
52 language processing (Kutas and Federmeier, 2000; Federmeier, 2007). In particular, a review of  
53 studies using electrophysiology, functional magnetic resonance imaging,  
54 magnetoencephalography, and intracranial recordings showed that the middle temporal gyrus  
55 is the main area of the brain involved in generation of the N400 (Lau et al., 2008), an area  
56 involved in long-term storage of lexico-semantic information (Hickok and Poeppel, 2004, 2007;  
57 Gitelman et al., 2005; Martin, 2007).

58 A typical N400 ERP study involves the auditory and/or visual presentation of a number of  
59 congruent and incongruent sentences<sup>1</sup>, and the recording of the subject's associated neural  
60 response through surface electrodes placed on the head. The N400 is estimated by comparing  
61 the average ERP waveforms corresponding to the congruent and incongruent sentences. It is  
62 well known that the creation of the stimulus material is a complex and arduous process that  
63 includes: (i) ideation of a set of congruent and incongruent sentences with a consistent morpho-

---

<sup>1</sup> The semantic violation that elicits the N400 ERP can also be primed by a list of words that do not form sentences or discourse (Titone and Salisbury, 2004; Romei et al. 2011; Brown and Hagoort, 1993).

64 syntactic structure; (ii) evaluation of the sentences by independent reviewers; and in case of  
65 auditory stimuli, (iii) recording of the sentences by a trained speaker using sophisticated  
66 equipment such as a studio microphone, a high-fidelity soundcard, and a sound-proof booth;  
67 and (iv) advanced post-processing of the audio files (Duncan et al., 2009; Swaab et al., 2011).

68 Most N400 studies use stimulus materials that have been developed in-house, and often the  
69 sentences are provided as an appendix (e.g. Litcofsky and Van Hell, 2017; Holt et al., 2018).

70 However, sentence development is very time-consuming and with multiple researchers  
71 repeating the process, it is not only highly inefficient, it makes comparison of results difficult.

72 Varying degrees of semantic violation in sentences can lead to variation the amplitude of the  
73 elicited N400 response (Kutas and Hillyard, 1984), which further adds to the difficulty of

74 comparing results between different studies that use different stimuli. If researchers had access

75 to a standardized set of sentences, this would reduce the need for repeated stimulus

76 development by different research groups and also facilitate comparisons between different

77 studies (Bradshaw, 1984). In 1980, Bloom and Fischler published one of the most widely used

78 sentence sets. This study provided 329 sentences with different levels of predictability. A subset

79 of sentences from this study was used by Kutas and Hillyard (1984) to confirm that the

80 magnitude of the N400 was inversely correlated with the degree of expectation associated with

81 the final word of the sentences. Although Bloom and Fischler (1980) provided 329 sentences,

82 studies that evaluate multiple different test conditions often require a larger number of stimuli.

83 Block and Baldwin (2010) extended Bloom and Fischler's set by adding 398 new sentences that

84 followed a similar format. Although this expanded sentence set is likely to be large enough for

85 visual-only studies in which sentences are presented as text only, each N400 study that

86 presented auditory or auditory-visual stimuli would still need to make a recording of their

87 selected sentences for use in their study. To the best of our knowledge, no researchers have

88 published all of the necessary stimulus material to run an N400 ERP study with auditory stimuli.

89 This paper provides a large set of congruent and incongruent English sentences, details the  
90 creation process, evaluates the quality of the sentences through a subjective evaluation of their  
91 meaningfulness, and provides the stimulus material – in text format and auditory recordings. In  
92 addition, the paper presents the results of an electrophysiology validation study which assessed  
93 the appropriateness of the stimulus material for N400 ERP studies. This involved evaluation of  
94 the N400 through grand-average ERP signals, analysis of the scalp-distribution, characterization  
95 of the individual variability of its magnitude, and an investigation of the influence of age and  
96 gender in a large set of normal hearing subjects.

## 97 **2. Stimulus material**

### 98 2.1. Congruent and incongruent sentences

99 A set of 350 congruent and 350 incongruent sentences was created according to the following  
100 structure: <<The + [1<sup>st</sup> noun: 2 syllables] + [verb: 1 syllable] + the + [2<sup>nd</sup> noun: 2 syllables] +  
101 [complement: 3 syllables]>>. An example of a sentence with this structure could be ‘The toddler  
102 likes the biscuits with some milk’. The first noun and the verb defined the context of the  
103 sentence, whereas the congruency of the sentence was determined by the second noun, i.e. the  
104 ‘critical word’. Incongruent sentences were those in which the critical word was not coherent  
105 with the context. The final complement provided the sentence with continuity following the  
106 potential congruency violation in order to avoid an overlap of N400 and wrap-up effects  
107 (Hagoort and Brown, 2000). The preposition <<the>> before the critical word aimed to provide  
108 an identical preceding phonetic context for all critical words in order to standardize the  
109 assimilation effect (i.e. a change in the pronunciation of a phoneme due to an adjacent sound)  
110 across the entire set of sentences (Ohala, 1988).

111 [Table 1]

112 Congruent sentences can be recorded naturally, however, reading aloud incongruent sentences  
113 can lead to acoustic confounds derived from exaggerated or unnatural prosody introduced by

114 the speaker who is aware that what they are saying is somehow 'odd' (Dimitrova et al., 2012;  
115 Meulman et al., 2014). Thus, it is conceivable that listeners might have access to prosodic cues  
116 indicating that a violation is to come before it has actually occurred. This expectancy would skew  
117 any N400 effects. To overcome this confound, each incongruent sentence was constructed by  
118 combining two naturally spoken congruent sentences using a cross-splicing procedure  
119 (Steinhauer et al., 2010; Meulman et al., 2014). For each incongruent sentence, an additional  
120 sentence was created by replacing the incongruent critical word with a different noun that was  
121 congruent with the context of the sentence. Table 1 shows an example of the three types of  
122 sentences created. In this example, sentences 1a and 1c are congruent sentences that were  
123 recorded naturally. Sentence 1b is an incongruent sentence that was constructed by cross-  
124 splicing the critical word from sentence 1c ('glasses') and replacing it with 'petrol' from 1a. To  
125 facilitate the cross-splicing procedure, we chose critical words with phonemes that have clear  
126 onsets following a brief silence in the otherwise continuous speech stream (stops: [p], [b], [t],  
127 [k], [g]; affricates: [tʃ] or [dʒ]); and avoided words that started with vowels.

128 All words used in the sentences were contained in SUBTLEX-UK, a word-frequency database for  
129 British English based on subtitles of British television programmes (Van Heuven et al., 2014). In  
130 order to ensure that the words of the sentences were familiar to the subjects, only words with  
131 a *Zipf* value<sup>2</sup> between 4 and 7 were selected. This value is an indication of medium- to high- word  
132 frequency (Monsell et al., 1989; Van Heuven et al., 2014).

133 [Figure 1, single column]

134 The meaningfulness of the congruent and incongruent sentences was evaluated by young adult  
135 native English-speaking students from Macquarie University (Sydney, Australia). The

---

<sup>2</sup> The *Zipf scale* takes the log<sub>10</sub> of the frequency per billion words (Van Heuven et al., 2014). Thus a *Zipf* value equal to 4 indicates a frequency per million words (*fpmw*) of 10; and a *Zipf* equal to 7 corresponds to a *fpmw* of 10,000. This scale comes from the American linguist George Kingsley Zipf, who was the first to formulate a law about the regularities of word frequency distribution (Zipf, 1949).

136 participants read each sentence and provided a rating between 1 and 6, where 1 – ‘The sentence  
137 makes complete sense’ and 6 – ‘The sentence makes no sense at all’. Each sentence was  
138 evaluated by at least five evaluators. The presentation order of the 700 sentences was  
139 randomized, and a final score was obtained for each sentence by estimating the mean of the  
140 evaluations. Sentences with a mean score closer to 1 were considered very congruent, while  
141 sentences with a mean score closer to 6 were considered very incongruent. During a final review  
142 of the list of sentences, five incongruent sentences were found to be potentially congruent, and  
143 therefore, they were assigned a congruency score of 3.5. Figure 1 shows the histogram of the  
144 subjective meaningfulness ratings of the congruent and incongruent sentences. This figure  
145 shows that the distribution is scattered towards the extreme values, indicating a predominance  
146 of sentences rated by the evaluators as either very congruent or very incongruent.

## 147 2.2. Questions and fillers

148 A set of questions that focused on the content of the sentences was also created. The questions  
149 occurred randomly during stimulus presentation to help sustain the attention of the  
150 participants. This set consisted of 200 *wh*-questions with an equal number focused on the  
151 congruent and incongruent sentences. For the sentences shown in table 1, the associated  
152 questions (Q) and responses (R) were [congruent sentence] Q: *Where does the driver put the*  
153 *petrol?* R: *In the car*; [incongruent sentence] Q: *Where does the mother break the petrol?* R: *On*  
154 *the shelf*. In addition, 130 filler sentences were devised to reduce predictability. These were  
155 similar in duration to the stimulus sentences, but differed in structure. Some examples of these  
156 fillers were “*Every Sunday the father goes to church*” and “*The package was sent by express*  
157 *post*”.

## 158 2.3. Audio recordings

159 The sentences, questions and fillers were recorded by a trained female native Australian English  
160 speaker in a sound-proof recording studio using a C535-EB vocal microphone (AKG Acoustics  
161 GmbH, Vienna, Austria), a StudioLive 16.4.2 audio mixer (PreSonus, Baton Rouge, LA), and a

162 sampling rate of 48 kHz. The congruent sentences (i.e. XXXa) and the additional congruent  
163 sentences necessary to form the incongruent sentences (i.e. XXXc) were recorded consecutively  
164 at least twice. To facilitate cross-splicing, the speaker was instructed to emphasize the critical  
165 word and to give a brief pause of about 0.5 seconds after the end of the critical word.

166 The recorded audio files were processed offline using Praat (Boersma, 2001; Boersma and  
167 Weenink, 2016). Processing consisted of (1) extracting the sentences from the continuous audio  
168 files; (2) selecting the best candidate from the recorded options for each sentence based on  
169 intonation, creakiness, clarity, intensity, and ease for cross-splicing; (3) cross-splicing the critical  
170 word to construct the incongruent sentences (i.e., XXXb); (4) adjusting the intensity of all  
171 sentences according to their root-mean square (RMS) value; and (5) setting the time points at  
172 which relevant language components occurred (markers) in order to identify the onset of the  
173 associated language-related ERPs during subsequent data processing. Markers were placed at  
174 (i) the onset of the initial <<The>>, (ii) the onset of the <<the>> preceding the critical word, (iii)  
175 the onset of the critical word, (iv) the onset of the complement, and (v) the end of the sentence.  
176 All cross-splicing edits were performed at zero-crossing points to avoid undesired audible  
177 artefacts like clicks or pops. The quality of the final audio files and the absence of imperfections  
178 in the cross-splicing process was validated by author 5.

179 The full list of congruent and incongruent sentences, questions, and fillers, along with the mean  
180 subjective evaluation of their meaningfulness, are provided as supporting material in appendix  
181 A. This appendix also includes a description of the rationale for the excluded sentences, i.e. those  
182 with a congruency rating set to 3.5. The raw audio files and markers are also provided as  
183 supplementary material (appendix B).

### 184 **3. Experimental validation**



185 The feasibility of the proposed set of sentences to evoke the N400 ERP was evaluated with an  
186 electroencephalography (EEG) study, in which the influence of age and gender on the N400  
187 magnitude was analysed.

### 188 3.1. Methods

#### 189 3.1.1. Ethics

190 The experimental protocol followed in this study was in accordance with the National  
191 Statements on Ethical Conduct in Human Research and was approved by the Human Research  
192 Ethics Committees of Macquarie University and Australian Hearing (Refs 5201400862;  
193 AHHREC2014-5).

#### 194 3.1.2. Participants

195 Forty-five participants (aged 19-57, mean = 38.78 years, SD = 11.22 years, 21 females) were  
196 recruited from the general community and Macquarie University. The inclusion criteria required  
197 that participants had English as a first language and normal or near-normal pure-tone hearing  
198 thresholds in both ears in the typical range of frequencies evaluated in the clinic (Dillon, 2012;  
199 Katz, 2014). Normal hearing was defined as a hearing loss  $\leq 20$  dB hearing level (HL) at 0.25 – 6  
200 kHz; and near-normal thresholds were considered as  $\leq 25$  dB HL up to 2 kHz,  $\leq 30$  dB HL at 3 kHz,  
201  $\leq 35$  dB HL at 4 kHz, and  $\leq 40$  dB HL at 6 kHz (Moore et al., 2012). All participants gave written  
202 consent to participate, and received \$40 after completing the study.

#### 203 3.1.3. Auditory stimulus

204 The auditory stimuli consisted of 80 congruent sentences, 80 incongruent sentences, 14 fillers,  
205 and 24 questions taken from the recorded materials described in section 2. Eighty highly  
206 congruent and 80 highly incongruent sentences were chosen randomly from the 160 sentences  
207 with the lowest and highest meaningfulness score, i.e. the 160 most- and 160 least-congruent  
208 sentences respectively. All of the 160 most congruent sentences were rated as 1.0, and the  
209 ratings for the 160 least congruent sentences were between 5.8 and 6.0. The 14 fillers were

210 selected randomly from the list of fillers and presented randomly every 6 to 10 sentences, i.e.  
211 with a probability of occurrence of 8.7%. The 24 questions were equally distributed between the  
212 congruent and incongruent sentences, and were presented randomly every 2 to 11 sentences,  
213 i.e. with a probability of occurrence of 15%. A short 1 kHz tone (or ‘beep’) was presented 1  
214 second before every question to inform the participants that they were about to hear a question  
215 about the preceding sentence and that a brief oral response was expected from them. The time  
216 between the end of one sentence and the onset of the following sentence was 4 seconds, except  
217 for questions, in which a time period of 8 seconds was provided to allow participants time to  
218 respond. The sentences, fillers and questions were distributed in four blocks of about 7 minutes  
219 each.

#### 220 3.1.4. Stimulus presentation

221 The auditory stimuli were presented to the subjects diotically at 60 dB sound-pressure level (SPL)  
222 through ER-3A insert earphones (Etymotic Research Inc., Elk Grove Village, IL) placed in the ear  
223 canal after otoscopic examination. The insert earphones were connected to a Fireface UCX audio  
224 soundcard (RME Audio, Haimhausen, Germany). Stimulus level was calibrated in a type HA2  
225 artificial ear 2-cc acoustic coupler, connected to a type 4144 pressure microphone, which was  
226 connected to a type 2636 measuring amplifier through a type 2639 preamplifier cable (Brüel &  
227 Kjær Sound & Vibration Measurement A/S, Nærum, Denmark).

#### 228 3.1.5. EEG recording

229 The EEG recording sessions took place in an electromagnetically shielded booth at the Speech  
230 and Hearing Clinic of Macquarie University. During the session, participants were seated in a  
231 comfortable armchair and were instructed to remain still during the test, to leave their neck and  
232 shoulder muscles relaxed, and to respond orally to the questions that followed a beep sound.

233 The recording of the neural response was carried out by 64 surface electrodes placed on the  
234 head using the TC64 EasyCap EEG recording cap (EASYCAP GmbH, Herrsching, Germany), in

235 which the electrode positions were distributed according to the International 10-20 system  
236 convention (Klem et al., 1999). The ground electrode was placed at AFz, and all active electrodes  
237 were referenced to the left mastoid (i.e., Tp9). The impedances of the electrodes with the skin  
238 were kept below 5 k $\Omega$  in all recordings. The EEGs were acquired by a Neuroscan SynAmps RT  
239 recording system, which was controlled by Curry 7 software (Compumedics Limited, Abbotsford,  
240 Australia). The recording sampling rate was 1 kHz; and the bandpass of the analogue filters was  
241 [0.01 - 300] Hz.

### 242 3.1.6. Data analysis

243 EEG data were analysed with custom scripts developed in Matlab (The Mathworks Inc., Natick,  
244 MA). The EEG recordings were processed in a 4-step procedure: (1) re-referencing, (2) digital  
245 filtering, (3) blink-artifact suppression, and (4) segments averaging. First, EEGs were re-  
246 referenced to the combined mastoid by subtracting from each raw EEG channel (i.e., [XX-Tp9])  
247 the Tp10 EEG channel divided by 2, i.e.  $[Tp10 - Tp9]/2$ . This way, all EEG channels were referenced  
248 to the combined mastoid, as  $[XX - Tp9] - [Tp10 - Tp9]/2 = [XX - (Tp9 + Tp10)]/2$ . The re-referenced  
249 EEGs were digitally filtered by a zero-phase 4<sup>th</sup> order Butterworth filter with bandpass cut-off  
250 frequencies [0.05 - 20] Hz. Blink artifacts were suppressed using iterative template matching and  
251 suppression (ITMS: Valderrama et al., 2018), a technique that allows blink-artifact suppression  
252 from single-channel EEG recordings. In the present multichannel application, ITMS was first  
253 applied to FP1 (i.e., an EEG channel situated in the vertical of the left eye) in order to facilitate  
254 detection of blink events. Once blink events were detected on the FP1 channel, a simplified  
255 version of ITMS was applied to the remaining EEG channels. The simplified ITMS procedure  
256 consisted of the following processes [described in detail in (Valderrama et al., 2018)]: (1)  
257 template estimation, (2) amplitudes estimation, (3) blink-artifact model estimation, and (4)  
258 blink-artifact model suppression. The Matlab functions that implement the original and  
259 simplified versions of ITMS are provided as supporting material (appendix C). Finally, the EEG

260 segments corresponding to 1 second pre- and 3 seconds post- critical word onset were averaged  
261 to obtain the ERPs associated with congruent and incongruent sentences.

262 The N400 ERP was quantified according to the area under the curve (AuC), which was estimated  
263 as the area (in  $s \cdot \mu V$ ) between the congruent and incongruent ERPs in the [0.4 – 0.8] s time  
264 interval (Swaab et al., 2011). Grand-average ERP signals were obtained by averaging the ERPs  
265 associated with congruent and incongruent sentences across participants. For this analysis, the  
266 10% of the ERPs with highest root-mean square (RMS) value in each channel (i.e., those most  
267 contaminated by noise) were discarded from the average in order to improve the quality of the  
268 ERPs (Thornton, 2007). Scalp topographic maps were also created using functions from the  
269 FieldTrip Matlab toolbox (Oostenveld et al., 2011). In addition, a cluster analysis was carried out  
270 to evaluate the N400 ERP in the frontal, central, and parietal-occipital areas of the brain. In each  
271 area, clusters were formed by averaging the ERPs corresponding to the following channels:  
272 frontal [AF3, AF1, AF2, AF4, F3, F1, Fz, F2, F4, FC3, FC1, FCz, FC2, FC4]; central [FC3, FC1, FCz,  
273 FC2, FC4, C3, C1, Cz, C2, C4, CP3, CP1, CPz, CP2, CP4]; and parietal-occipital [CP3, CP1, CPz, CP2,  
274 CP4, P3, P1, Pz, P2, P4, PO3, PO4, O1, Oz, O2].

275 The cluster analysis evaluated the presence of the N400 ERP, and assessed the influence of  
276 gender and age on its magnitude. The presence of the N400 ERP was evaluated by a one-sample  
277  $t$ -test, in which the null hypothesis was that the mean of the AuC distribution was zero (i.e., no  
278 N400 ERP was obtained). The gender analysis was based on a two-sample  $t$ -test, in which the  
279 null hypothesis was that the mean of the AuC distribution for males and females were equal  
280 (i.e., males and females showed similar AuCs). The influence of age on the AuC was evaluated  
281 using a linear regression test, in which the null hypothesis was that the slope of the correlation  
282 was equal to zero (i.e., AuC was not influenced by the age of the participants).

283 Statistical analysis was conducted using functions from the ‘Statistics and Machine Learning’  
284 Matlab toolbox. The normality assumption was evaluated in all distributions according to the

285 Lilliefors test of normality; and a conservative alpha-value of 0.01 was chosen as the criterion  
286 for statistical significance.

## 287 3.2. Results

### 288 3.2.1. Hearing thresholds

289 [Figure 2, double column]

290 Figure 2 shows the pure-tone audiometric threshold distributions for the left and right ears at  
291 the test frequencies. Boxplots represent the quartiles of the distribution. This figure shows that  
292 of the 45 participants, 44 presented with normal or near-normal hearing according to the study  
293 criteria. The remaining participant met the near-normal criteria for all but one of their  
294 thresholds, which was slightly outside the inclusion criteria.

### 295 3.2.2. N400 ERP

296 [Figure 3, double column]

297 Figure 3.A shows the grand-average ERPs associated with the onset of the congruent (in blue)  
298 and incongruent (in brown) critical words at selected EEG channels. The N400 ERP is shown in  
299 this figure as a negative voltage deviation of the incongruent sentence compared to the  
300 congruent sentence in the time interval [0.4 – 0.8] s after the critical word onset. The  
301 topographic maps presented in figure 3.B show that the N400 ERP magnitude is highest in the  
302 central area of the brain around the aforementioned time interval.

303 [Figure 4, double column]

304 Figure 4 presents an analysis of the N400 AuC at the frontal, central, and parietal-occipital areas.  
305 The first column shows the grand-average ERPs corresponding to each cluster of EEG channels.  
306 The second column shows the AuC histogram, the fitted distribution, and the probability (p) that  
307 the AuC distribution comes from a normal distribution with mean equal to 0. The results  
308 presented in the three evaluated areas show that all the p-values were statistically significant,

309 thus indicating that the set of sentences used in this study successfully evoked the N400 ERP.  
310 However, it is also noteworthy that a significant portion of the participants presented negative  
311 AuC values in each of the three areas [frontal: 14 (31%), central: 10 (22%), parietal-occipital: 9  
312 (20%)]. The fitted distributions for males and females presented in the third column show that  
313 females tended to show greater AuC, although this difference was not significant in any of the  
314 three topographic areas. Finally, the age plots shown in the fourth column indicate no effect of  
315 age on the AuC estimate.

#### 316 **4. Discussion**

317 This paper is an important contribution to the N400 literature, providing researchers with access  
318 to a new verified set of 700 congruent and incongruent English sentences and their audio  
319 recordings, appropriate for use in N400 ERP studies. To our knowledge, this constitutes the  
320 largest complete stimulus set suitable for both auditory and text-based experiments.

321 The subjective evaluation of the meaningfulness of the sentences showed that most sentences  
322 were classified as either highly congruent (on a 1-to-6 scale, 304/350 sentences had a mean  
323 score less than or equal to 1.5) or highly incongruent (250/350 sentences had a mean score  
324 greater than or equal to 5.5). This large set of highly congruent and incongruent sentences may  
325 be useful when testing subjects in multiple sessions or in different stimulus conditions [e.g.  
326 different SNRs, transducers (insert vs loudspeakers), etc.].

327 The ability of the sentences to elicit the N400 ERP was evaluated in a sample of 45 subjects with  
328 good hearing. In order to test a significant portion of the proposed set, this experiment used 80  
329 sentences selected randomly from the subset of the 160 most congruent sentences, and 80  
330 sentences from the 160 most incongruent sentences. This experiment demonstrated the  
331 suitability of the selected test sentences for evoking N400 ERP. Consistent with the majority of  
332 published N400 studies (Osterhout and Nicol, 1999; Swaab et al., 2011; Boudewyn et al., 2012),  
333 the grand-average ERPs and the topographic maps presented in figure 3 showed a dominant

334 central distribution of the N400 magnitude. In addition, the incongruous sentences produced an  
335 ERP that was more positive at the peak that followed the N400 in the parietal area. This effect  
336 is known as 'post-N400 positivity' (Matsumoto et al., 2005; Van Petten and Luka, 2006;  
337 Federmeier et al., 2007), and is associated with an impossible or semantically anomalous  
338 interpretation (Van Petten and Luka, 2012).

339 The distributions of the individual AuC estimates at the midline frontal, central and parietal areas  
340 had statistically significant positive means, thus providing evidence of the elicitation of the N400  
341 ERP at a group level. However, these distributions also showed that a significant portion of the  
342 individuals [frontal: 14 subjects (31%); central: 10 subjects (22%); parietal: 9 subjects (20%)]  
343 presented negative AuC values, consistent with other studies showing only limited reliability of  
344 N400 ERP studies when conducting single-subject analyses (Cruse et al., 2014).

345 The analysis also showed that although females presented fitted AuC distributions with a higher  
346 mean, this effect was not statistically significant. The few previous studies that have analysed  
347 the influence of gender on the N400 ERP have reported conflicting results. On the one hand,  
348 Tsolaki et al. (2015) showed no gender-related differences; whereas several other studies have  
349 found that females presented earlier and larger N400 ERP components (Daltrozzo et al., 2007;  
350 Proverbio et al., 2010; Steffensen et al., 2008; Wirth et al., 2007). Whether or not males and  
351 females process lexical violations differently is still under debate. With regard to age, we  
352 observed no effect on the AuC in any of the midline clusters. These results are consistent with  
353 recent literature, in which no significant N400 changes with age have been reported (Federmeier  
354 et al., 2003; Grieder et al., 2012; Komes et al., 2014; Tsolaki et al., 2015; Wilkinson et al., 2013).  
355 It is noteworthy that the age and gender analyses conducted in this paper were supported by a  
356 gender-balanced distribution of the participants (21 females against 24 males), and by a uniform  
357 distribution of age across a wide range (from 19 to 57 years).

358 Attached to this paper, the full text of all stimulus sentences, filler sentences, and questions are  
359 provided as supporting material (appendix A). In addition, this paper also provides the processed

360 audio files (appendix B). Access to this stimulus material is intended to facilitate the efficient  
361 implementation of future N400 ERP studies, and promote standardization and comparisons  
362 across studies.

363 Two important limitations should be considered when using the stimulus material presented in  
364 this paper for other studies. The first limitation is that the audio files were recorded by an  
365 Australian-English speaker, which could affect the results if a study is conducted in a population  
366 that uses a different English dialect. In this case, it is advised that the sentences are recorded by  
367 a trained speaker of the same dialect as the study population, following the steps described in  
368 section 2.3. The second limitation is that the subjective rating of the congruency of the sentences  
369 was carried out by young adult university students only, which may introduce bias since their  
370 evaluations are not being representative of the general population (which included individuals  
371 with a lower education level and from different age groups). In cases where the study population  
372 differs significantly from the university-educated young adults who rated the congruency of the  
373 sentences in this study, the authors suggest that the congruency of the sentences is re-evaluated  
374 by people with a similar profile as the study population.

### 375 **Acknowledgements**

376 The authors gratefully acknowledge Ms Lorna Betts for her help with the audio recording of the  
377 sentences; Mr Greg Stewart (NAL: National Acoustic Laboratories, Sydney, Australia) for his help  
378 with the calibration of the stimuli; and the Macquarie University students who participated in  
379 the subjective evaluation of the sentences. This work was supported by the Australian  
380 Government Department of Health. No potential conflict of interest was reported by the  
381 authors.

### 382 **Appendix**

383 Supplementary material associated with this article can be found at [URL]. Appendix A presents  
384 a table of (1) the congruent and incongruent sentences; (2) the associated individual and mean



385 subjective congruency ratings; (3) the list of questions and expected responses; (4) the list of  
386 filler sentences; and (5) the rationale for excluding the ambiguous sentences. Appendix B  
387 includes the processed audio files of the recorded sentences and markers. Appendix C includes  
388 the Matlab functions and associated files necessary to implement the blink-artifact removal  
389 technique 'iterative template matching and suppression (ITMS)' in a multichannel EEG  
390 configuration.

391

392 **References**

- 393 Block, C.K., Baldwin, C.L. (2010). Cloze probability and completion norms for 498 sentences:  
394 Behavioral and neural validation using event-related potentials. *Behavior Research*  
395 *Methods* 42, 665-670.
- 396 Bloom, P.A., Fischler, I. (1980). Completion norms for 329 sentence contexts. *Memory &*  
397 *Cognition* 8, 631-642.
- 398 Boersma, P. (2001). Praat, a system for doing phonetics by computer. *Glott International* 5, 341-  
399 345.
- 400 Boersma, P., Weenink, D. (2016). Praat: doing phonetics by computer [Computer program].  
401 Version 6.0.19, retrieved 13 June 2016 from <http://www.praat.org/>
- 402 Bradshaw, J.L. (1984). A guide to norms, ratings, and lists. *Memory & Cognition* 12, 202-206.
- 403 Brown, C., Hagoort, P.J. (1993). The processing nature of the N400: evidence from masked  
404 priming. *Journal of Cognitive Neuroscience* 5, 34-44.
- 405 Boudewyn, M.A., Gordon, P.C., Long, D., Polse, L., Swaab, T.Y. (2012). Does discourse  
406 congruence influence spoken language comprehension before lexical association?  
407 Evidence from event-related potentials. *Language and Cognitive Processes* 27, 698-  
408 733.
- 409 Cruse, D., Beukema, S., Chennu, S., Malins, J.G., Owen, A.M., McRae, K. (2014). The reliability of  
410 the N400 in single subjects: Implications for patients with disorders of consciousness.  
411 *NeuroImage: Clinical* 4, 788-799.
- 412 Daltrozzo, J., Wioland, N., Kotchoubey, B. (2007). Sex differences in two event-related potential  
413 components related to semantic priming. *Archives of Sexual Behavior* 36, 555-568.
- 414 Dillon, H. (2012). *Hearing Aids* (Thieme Publishers, New York), 608 p.

- 415 Dimitrova, D.V., Stowe, L.A., Redeker, G., Hoeks, J.C. (2012). Less is not more: neural responses  
416 to missing and superfluous accents in context. *Journal of Cognitive Neuroscience* 24,  
417 2400-2418.
- 418 Duncan, C.C., Barry, R.J., Connolly, J.F., Fischer, C., Michie, P.T., Näätänen, R., Polich, J., Reinvang,  
419 I., Van Petten, C. (2009). Event-related potentials in clinical research: Guidelines for  
420 eliciting, recording, and quantifying mismatch negativity, P300, and N400. *Clinical*  
421 *Neurophysiology* 120, 1883-1908.
- 422 Federmeier, K.D. (2007). Thinking ahead: the role and roots of prediction in language  
423 comprehension. *Psychophysiology* 44, 491-505.
- 424 Federmeier, K.D., Van Petten, C., Schwartz, T.J., Kutas, M. (2003). Sounds, words, sentences:  
425 age-related changes across levels of language processing. *Psychology and Aging* 18,  
426 858-872.
- 427 Federmeier, K.D., Wlotko, E.W., de Ochoa-Dewald, E., Kutas, M. (2007). Multiple effects of  
428 sentential constraint on word processing. *Brain Research* 1146, 75-84.
- 429 Gitelman, D.R., Nobre, A.C., Sonty, S., Parrish, T.B., Mesulam, M.M. (2005). Language network  
430 specializations: an analysis with parallel task designs and functional magnetic  
431 resonance imaging. *Neuroimage* 26, 975-985.
- 432 Grieder, M., Crinelli, R.M., Koenig, T., Wahlund, L.-O., Dierks, T., Wirth, M. (2012).  
433 Electrophysiological and behavioural correlates of stable automatic semantic retrieval  
434 in aging. *Neuropsychologia* 50, 160-171.
- 435 Hagoort, P. (2008). The fractionation of spoken language understanding by measuring electrical  
436 and magnetic brain signals. *Philosophical Transactions of the Royal Society of London*  
437 – Series B: Biological Sciences 363, 1055-1069.
- 438 Hagoort, P., Brown, C. M. (2000). ERP effects of listening to speech: semantic ERP  
439 effects. *Neuropsychologia* 38, 1518-1530.

- 440 Hickok, G., Poeppel, D. (2004). Dorsal and ventral streams: a framework for understanding  
441 aspects of the functional anatomy of language. *Cognition* 92, 67-99.
- 442 Hickok, G., Poeppel, D. (2007). The cortical organization of speech processing. *Nature Reviews*  
443 *Neuroscience* 8, 393-402.
- 444 Holt, R., Kung, C., Demuth, K. (2018). Listener characteristics modulate the semantic processing  
445 of native vs. foreign-accented speech. *PLoS ONE* 13(12): e0207452.
- 446 Jerger, J., Martin, J. (2009). Some effects of aging on event-related potentials during a linguistic  
447 monitoring task. *International Journal of Audiology* 44, 321-330.
- 448 Klem, G.H., Lüders, H.O., Jasper, H.H., Elger, C. (1999). The ten-twenty electrode system of the  
449 International Federation. *The International Federation of Clinical Neurophysiology.*  
450 *Electroencephalography and Clinical Neurophysiology* 52, 3-6.
- 451 Katz, J. (2014). *Handbook of clinical audiology*, edited by Marshall Chasin, Kristina English, Linda  
452 J. Hood, and Kim L. Tillery (Wolters Kluwer Health, Philadelphia), 927 p.
- 453 Komes, J., Schweinberger, S.R., Wiese, H. (2014). Fluency affects source memory for familiar  
454 names in younger and older adults: evidence from event-related brain potentials.  
455 *NeuroImage* 92C, 90-105.
- 456 Kutas, M. (1993). In the company of other words: electrophysiological evidence for single-word  
457 and sentence context effects. *Language and Cognitive Process* 8, 533-572.
- 458 Kutas, M., Federmeier, K.D. (2000). Electrophysiology reveals semantic memory use in language  
459 comprehension. *Trends in Cognitive Sciences* 4, 463-470.
- 460 Kutas, M., Hillyard, S.A. (1980). Reading senseless sentences: brain potentials reflect semantic  
461 incongruity. *Science* 207, 203-205.
- 462 Kutas, M., Hillyard, S.A. (1984). Brain potentials during reading reflect word expectancy and  
463 semantic association. *Nature* 307, 161-163.

- 464 Litcofsky, K.A., Van Hell, J.G. (2017). Switching direction affects switching costs: Behavioral, ERP  
465 and time-frequency analyses of intra-sentential codeswitching. *Neuropsychologia* 97,  
466 112-139.
- 467 Martin, A. (2007). The representation of object concepts in the brain. *Annual Review of*  
468 *Psychology* 58, 25-45.
- 469 Matsumoto, A., Iidaka, T., Haneda, K., Okada, T., Sadato, N. (2005). Linking semantic priming  
470 effect in functional MRI and event-related potentials. *Neuroimage* 24, 624-634.
- 471 Meulman, N., Stowe, L.A., Sprenger, S.A., Bresser, M., Schmid, M.S. (2014). An ERP study on L2  
472 syntax processing: When do learners fail? *Frontiers in Psychology* 5, art, 1072, 17p.
- 473 Monsell, S., Doyle, M.C., Haggard, P.N. (1989). Effects of frequency on visual word recognition  
474 tasks - Where are they? *Journal of Experimental Psychology: General* 118, 43-71.
- 475 Moore, B.C.J., Creeke, S., Glasberg, B.R., Stone, M.A., Sek, A. (2012). A version of the TEN Test  
476 for use with ER-3A insert earphones. *Ear and Hearing* 33, 554-557.
- 477 Ohala, J.J. (1988). "The phonetics and phonology of aspects of assimilation," in *Papers in*  
478 *Laboratory Phonology I – Between the Grammar and Physics of Speech*, edited by  
479 Kingston, J., Beckman, M.E. (Cambridge University Press, Cambridge, United Kingdom),  
480 258-275.
- 481 Oostenveld, R., Fries, P., Maris, E., Schoffelen, J.M. (2011). FieldTrip: Open Source Software for  
482 Advanced Analysis of MEG, EEG, and Invasive Electrophysiological Data.  
483 *Computational Intelligence and Neuroscience* 2011, art. 156869, 9 p.
- 484 Osterhout, L., Nicol, J. (1999). On the distinctiveness, independence, and time course of the  
485 brain responses to syntactic and semantic anomalies. *Language and Cognitive*  
486 *Processes* 14, 283-317.
- 487 Osterhout, L., Holcomb, P.J. (1992). Event-related brain potentials elicited by syntactic anomaly.  
488 *Journal of Memory and Language* 31, 785-806.

- 489 Proverbio, A.M., Riva, F., Zani, A. (2010). When neurons do not mirror the agent's intentions:  
490 sex differences in neural coding of goal-directed actions. *Neuropsychologia* 48, 1454-  
491 1463.
- 492 Romei, L., Wambacq, I.J.A., Besing, J., Koehnke, J., Jerger, J. (2011). Neural indices of spoken  
493 word processing in background multi-talker babble. *International Journal of Audiology*  
494 50, 321-333.
- 495 Steffensen, S.C., Ohran, A.J., Shipp, D.N., Hales, K., Stobbs, S.H., Fleming, D.E. (2008). Gender-  
496 selective effects of the P300 and N400 components of the visual evoked potential.  
497 *Vision Research* 48, 917-925.
- 498 Steinhauer, K., Abada, S.H., Pauker, E., Itzhak, I., Baum, S.R. (2010). Prosody-syntax interactions  
499 in aging: Event-related potentials reveal dissociations between on-line and off-line  
500 measures. *Neuroscience Letters* 472, 133-138.
- 501 Swaab, T.Y., Ledoux, K., Camblin, C.C., Boudewyn, M.A. (2011). Language-related ERP  
502 components. In: Kappenman, E.S., Luck, S.J. (Eds.), *The Oxford Handbook of Event-  
503 Related Potential Components*. Oxford Handbooks Online, Oxford, United Kingdom, 49  
504 p.
- 505 Thornton, A.R.D. (2007). Instrumentation and recording parameters. In: Burkard, R., Don, M.,  
506 Eggermont, J. (Eds.), *Auditory Evoked Potentials: Basic Principles and Clinical  
507 Application*. Lippincott William & Wilkins, Baltimore, MD, pp. 73-101.
- 508 Titone, D.A., Salisbury, D.F. (2004). Contextual modulation of N400 amplitude to lexically  
509 ambiguous words. *Brain and Cognition* 55, 470-478.
- 510 Tsolaki, A., Kosmidou, V., Hadjileontiadis, L., Kompatsiaris, I.Y., Tsolaki, M. (2015). Brain source  
511 localization of MMN, P300 and N400: Aging and gender differences. *Brain Research*  
512 1603, 32-49.

- 513 Van Heuven, W.J.B., Mandera, P., Keuleers, E., Brysbaert, M. (2014). SUBTLEX-UK: A new and  
514 improved word frequency database for British English. *The Quarterly Journal of*  
515 *Experimental Psychology* 67, 1176-1190.
- 516 Van Petten, C., Luka, B.J. (2006). Neural localization of semantic context effects in  
517 electromagnetic and hemodynamic structures. *Brain and Language* 97, 279-293.
- 518 Van Petten, C., Luka, B.J. (2012). Prediction during language comprehension: benefits, costs, and  
519 ERP components. *International Journal of Psychophysiology* 83, 176-190.
- 520 Valderrama, J.T., de la Torre, A., Van Dun, B. (2018). An automatic algorithm for blink-artifact  
521 suppression based on iterative template matching: application to single-channel  
522 recording of cortical auditory evoked potentials. *Journal of Neural Engineering* 15, art.  
523 016008, 16 p.
- 524 Wilkinson, A.J., Yang, L., Dyson, B.J. (2013). Modulating younger and older adult's performance  
525 in ignoring pictorial information during a word matching task. *Brain and Cognition* 83,  
526 351-359.
- 527 Wirth, M., Horn, H., König, T., Stein, M., Federspiel, A., Meier, B., Michel, C.M., Strik, W. (2007).  
528 Sex differences in semantic processing event-related brain potentials distinguish  
529 between lower and higher order semantic analysis during word reading. *Cerebral*  
530 *Cortex* 17, 1987-1997.
- 531 Zipf, G.K. (1949). *Human behaviour and the principle of least effort*. Oxford, England: Addison-  
532 Wesley Press, 573 p.
- 533

534 **Figure legends**

- 535 • Figure 1. (Color online) Histogram of the subjective meaningfulness ratings for the list of  
536 congruent and incongruent sentences.
- 537 • Figure 2. Pure-tone audiometric threshold distributions for left and right ears. Boxplots  
538 indicate the minimum, 1<sup>st</sup> quartile, median, 3<sup>rd</sup> quartile, and maximum values of the  
539 distributions. Straight and dashed lines indicate the limits for normal and near-normal  
540 hearing, respectively.
- 541 • Figure 3. (Color online) [A] Grand-Average N400 ERPs for selected EEG channels. [B]  
542 Topographic plots show the scalp distribution of the area under the curve estimated from  
543 the Grand-Average ERPs at different time intervals.
- 544 • Figure 4. (Color online) Cluster analysis for the area under the curve (AuC) in frontal, central,  
545 and parietal-occipital areas: [1<sup>st</sup> column] Grand-average cluster ERPs; [2<sup>nd</sup> column] raw AuC  
546 distribution, fitted distribution, and probability ( $p$ ) that the distribution comes from a  
547 normal distribution with mean equal to 0; [3<sup>rd</sup> column] fitted distribution for males and  
548 females,  $p$ -value indicating the probability that the two distributions present the same  
549 mean; [4<sup>th</sup> column] effect of age on AuC evaluated through a linear regression analysis.
- 550



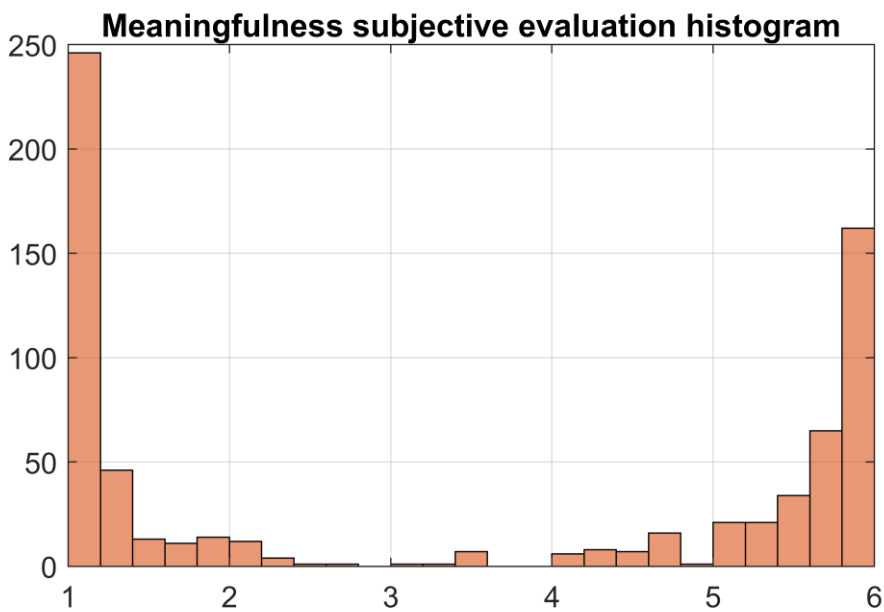
551 **Tables**

- 552 • Table 1. Example of the three types of sentences. Sentences 1a and 1c are congruent
- 553 sentences that were recorded naturally. Sentence 1b is an incongruent sentence artificially
- 554 constructed from sentences 1a and 1c by cross-splicing the critical word. The critical words
- 555 in these sentences are highlighted in bold.

1a	The driver puts the <b>petrol</b> in the car
1b	The mother breaks the <b>petrol</b> on the shelf
1c	The mother breaks the <b>glasses</b> on the shelf

556

557

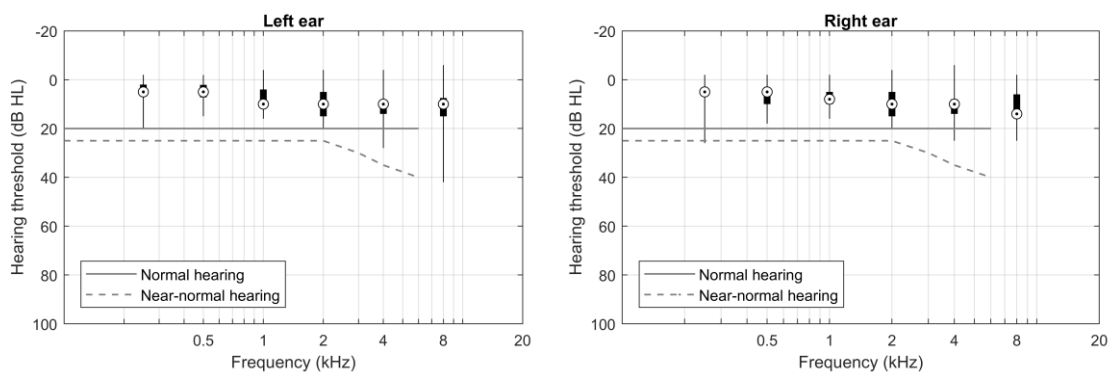


558

Totally congruent

Not congruent at all

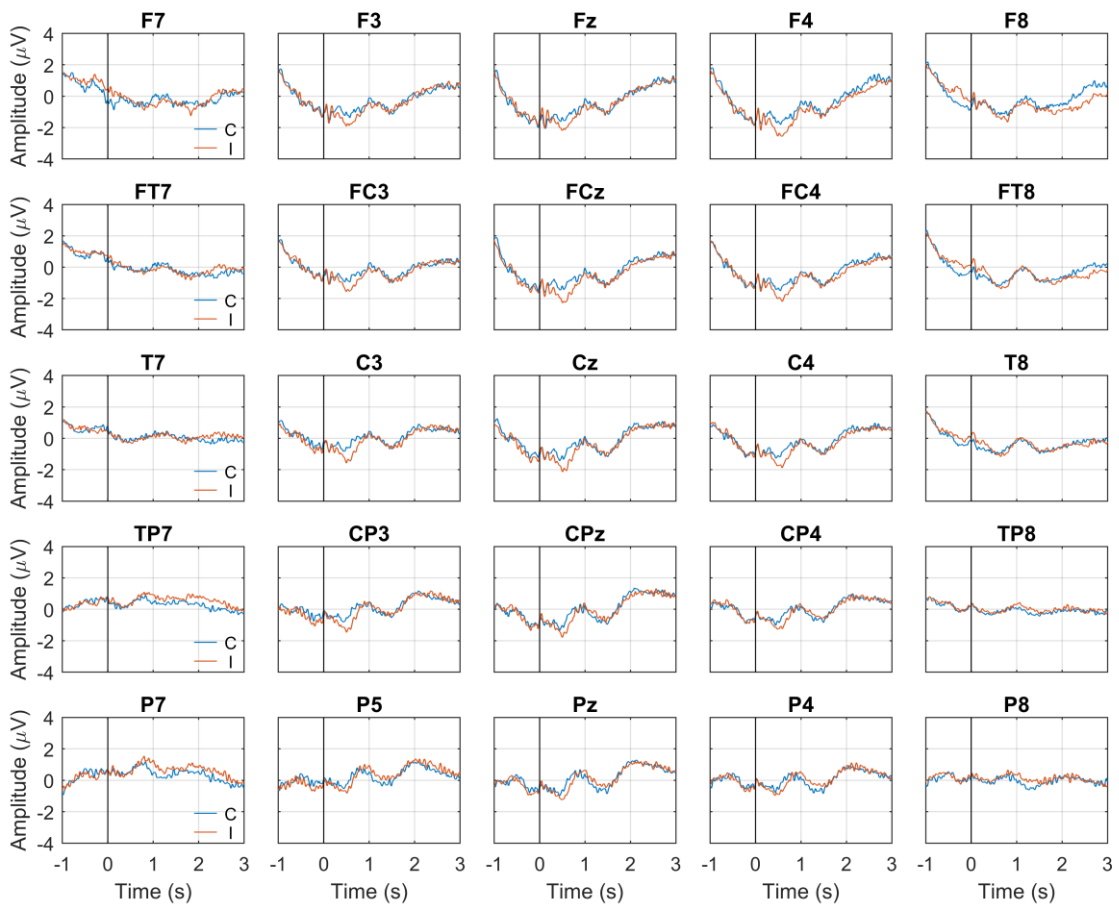
559



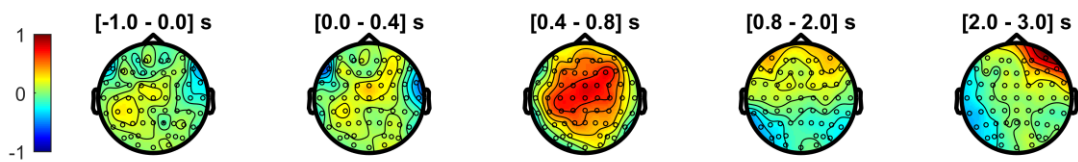
560

561

**A. Grand-Average ERPs**

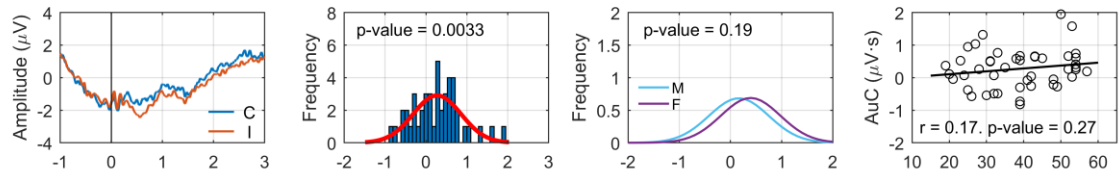
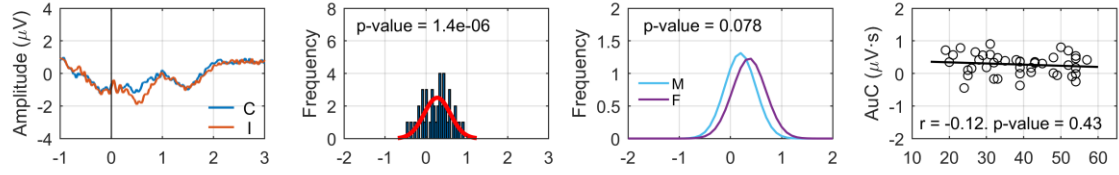


**B. Topographic maps**



562

563

**A. Frontal****B. Central****C. Parietal-Occipital**