



Kate Stone | Bruno Nicenboim | Shravan Vasishth | Frank Rösler

Understanding the effects of constraint and predictability in ERP

Suggested citation referring to the original publication:

Neurobiology of Language 4 (2022) 2, pp. 1 - 71

DOI: https://doi.org/10.1162/nol_a_00094

ISSN: 2641-4368

Journal article | Version of record

Secondary publication archived on the Publication Server of the University of Potsdam:

Zweitveröffentlichungen der Universität Potsdam : Humanwissenschaftliche Reihe 829

ISSN: 1866-8364

URN: <https://nbn-resolving.org/urn:nbn:de:kobv:517-opus4-587594>

DOI: <https://doi.org/10.25932/publishup-58759>

Terms of use:

This work is licensed under a Creative Commons License. This does not apply to quoted content from other authors. To view a copy of this license visit

<https://creativecommons.org/licenses/by/4.0/>.

Abstract

14

15 Intuitively, strongly constraining contexts should lead to stronger probabilistic
16 representations of sentences in memory. Encountering unexpected words could therefore be
17 expected to trigger costlier shifts in these representations than expected words. However,
18 psycholinguistic measures commonly used to study probabilistic processing, such as the
19 N400 event-related potential (ERP) component, are sensitive to word predictability but
20 not to contextual constraint. Some research suggests that constraint-related processing
21 cost may be measurable via an ERP positivity following the N400, known as the anterior
22 post-N400 positivity (PNP). The PNP is argued to reflect update of a sentence
23 representation and to be distinct from the posterior P600, which reflects conflict detection
24 and reanalysis. However, constraint-related PNP findings are inconsistent. We sought to
25 conceptually replicate Federmeier et al. (2007) and Kuperberg et al. (2020), who observed
26 that the PNP, but not the N400 or the P600, was affected by constraint at unexpected but
27 plausible words. Using a pre-registered design and statistical approach maximising power,
28 we demonstrated a dissociated effect of predictability and constraint: strong evidence for
29 predictability but not constraint in the N400 window, and strong evidence for constraint
30 but not predictability in the later window. However, the constraint effect was consistent
31 with a P600 and not a PNP, suggesting increased conflict between a strong representation
32 and unexpected input rather than greater update of the representation. We conclude that
33 either a simple strong/weak constraint design is not always sufficient to elicit the PNP, or
34 that previous PNP constraint findings could be an artifact of smaller sample size.

35 *Keywords:* N400, anterior PNP, posterior P600, probabilistic processing, constraint,
36 predictability, entropy

Understanding the effects of constraint and predictability in ERP

37
38 Readers can use contextual cues from words and sentences to construct a mental
39 representation of an event. This representation can be viewed as probabilistic, with
40 plausible upcoming words and sentence structures preactivated in anticipation of their
41 appearance (Kuperberg et al., 2020; Kuperberg & Jaeger, 2016; Kutas & Federmeier,
42 2011). Assuming that readers generate such a representation, its probabilistic strength
43 should depend on how constraining the sentential context is. For example, in sentence (1)a,
44 the strong constraint of the context makes the word *true* highly predictable, whereas in
45 (1)b, the weak contextual constraint means no specific word is predictable (Federmeier
46 et al., 2007):

- 47 (1) a. *Strongly constraining:*
48 Sam could not believe her story was... true/published
49 b. *Weakly constraining:*
50 I was impressed by how much he... knew/published

51 The reader's probabilistic representation should therefore be stronger in (1)a than
52 (1)b, so that encountering the low-predictable word *published* is more unexpected (in the
53 sense that the reader expected a different event) in (1)a, even though *published* is equally
54 unpredictable in both contexts (according to a cloze test; Federmeier et al., 2007).
55 Nonetheless, psycholinguistic measures typically used to study probabilistic
56 processing—including the N400 event-related potential (ERP) component—have been
57 found to correspond only to the matched predictability of *published* between (1)a and (1)b,
58 and not the mismatch in constraint (Federmeier et al., 2007; Kuperberg et al., 2020; Kutas
59 & Hillyard, 1984; Van Petten & Luka, 2012). Instead, an anteriorly distributed positive
60 deflection in the ERP after the N400, the post-N400 positivity (PNP), may hold the key to
61 measuring the constraint/predictability dissociation (Brothers et al., 2020; Federmeier
62 et al., 2007; Kuperberg et al., 2020). However, empirical findings involving the PNP are

63 inconsistent (Federmeier & Kutas, 1999; Frank et al., 2015; Lai et al., 2021; Szewczyk &
64 Schriefers, 2013; Thornhill & Van Petten, 2012; Wlotko & Federmeier, 2007). Given the
65 potential importance of the PNP in studying reader's probabilistic representations, in this
66 registered report, we addressed possible sample size concerns in previous studies by testing
67 the PNP in a confirmatory study with a larger sample size.

68 **The post-N400 positivity (PNP)**

69 An incidental finding in many studies of the N400 has been that of a late positivity
70 beginning at around 600 ms in the anterior scalp region. This anterior positivity appears to
71 be spatially and functionally distinct from the more well-known posterior P600 (Kuperberg
72 et al., 2020). The P600 has been variously linked to conflict detection and repair processes
73 in a fronto-temporal cortical circuit (Bornkessel-Schlesewsky & Schlewsky, 2008; Brouwer
74 et al., 2017; Brouwer & Hoeks, 2013; Fitz & Chang, 2019; Kim & Osterhout, 2005;
75 Kuperberg et al., 2003; Meerendonk et al., 2009; Metzner et al., 2017; Osterhout &
76 Holcomb, 1992). In contrast, the anterior PNP has been linked to the update of event
77 representations, possibly involving the inhibition of representations falsified by unexpected
78 input via left prefrontal cortex (Kutas, 1993). Extending this characterisation, recent
79 research has suggested that the PNP is only elicited when unexpected input is still
80 plausible in the given context (DeLong et al., 2014; Kuperberg et al., 2020). For example,
81 in (2) below, *swimmers* is the most expected continuation, while *trainees* and *drawer* are
82 both low probability. However, *trainees* is still plausible in the context, while *drawer* is not.
83 A PNP and P600 were elicited by *trainees* relative to the expected *swimmers*, but not by
84 *drawer*, which only elicited a P600 (DeLong et al., 2014):

85 (2) The lifeguards received a report of sharks right near the beach [...] Hence they
86 cautioned the swimmers/trainees/drawer

87 The fact that only the plausible *trainees* and not the implausible *drawer* elicited the
88 PNP has led some to hypothesise that the PNP reflects a change in activity associated

89 with *successfully* updating the mental representation of an event, which may include the
90 inhibition of previous representations (Kuperberg et al., 2020; Kutas, 1993; Ness &
91 Meltzer-Asscher, 2018). Under this assumption and the assumption that the P600 reflects
92 reanalysis (Kim and Osterhout, 2005; Kuperberg et al., 2003; Osterhout and Holcomb,
93 1992, cf. Bornkessel-Schlesewsky and Schlewsky, 2008; Brouwer et al., 2017; Fitz and
94 Chang, 2019), Kuperberg et al. (2020) have proposed that an unexpected word (in this
95 example *trainees*) triggers a large but successful update of the readers' representation of
96 the event, including suppression of the more predictable event *caution the swimmers*. The
97 magnitude of this update is reflected by the presence of a PNP. According to Kuperberg et
98 al. (2020), the unexpected word also engages reanalysis processes during attempts to
99 accommodate it, which are reflected in the presence of a P600. In contrast, the implausible
100 *drawer* triggers no change in the existing event representation (PNP absent), even though
101 reanalysis processes may be engaged (P600 present).

102 More importantly for research on probabilistic processing, the PNP also appears to
103 be sensitive to contextual constraint. Like the N400, the PNP has been found to be larger
104 for low vs. high probability words (Brothers et al., 2017; Brothers et al., 2020; DeLong
105 et al., 2014; DeLong et al., 2011; Federmeier et al., 2007; Kuperberg et al., 2020; Ness &
106 Meltzer-Asscher, 2018; Thornhill & Van Petten, 2012); but unlike the N400, the PNP
107 appears to be larger for low probability words in strongly vs. weakly constraining contexts
108 (Brothers et al., 2020; Federmeier et al., 2007; Kuperberg et al., 2020). Returning to the
109 example in (1) above, Federmeier et al. (2007) found that the unexpected word *published*
110 elicited a larger PNP in the strongly constraining (1)a than in the weakly constraining
111 (1)b, even though their cloze probabilities and corresponding N400 amplitudes were the
112 same. The PNP would therefore appear to suggest that a stronger probabilistic
113 representation was built in (1)a than in (1)b, and that the stronger representation was
114 more costly to update.

115 However, not all studies eliciting the PNP involve a constraint manipulation

116 (Van Petten & Luka, 2012), and thus it is difficult to attribute the PNP exclusively to the
117 manipulation of contextual constraint, rather than to part of a biphasic response to low
118 probability words following the N400. Furthermore, not all studies manipulating constraint
119 show consistent effects on the PNP. Contrary to Federmeier et al. (2007) and Kuperberg
120 et al. (2020), Federmeier and Kutas (1999) found that *expected* words elicited a larger PNP
121 than unexpected words, and only in low constraint sentences. It should be noted that
122 expected words in the Federmeier and Kutas (1999) “low” constraint condition had a mean
123 cloze probability of 0.59 with a range 0.17 to 0.78; nonetheless, the direction of the PNP
124 constraint effect was the opposite of that described elsewhere. In high constraint sentences,
125 no difference in the PNP was observed between expected and unexpected words. More
126 recently, Szewczyk and Schriefers (2013) noted a larger, centrally distributed post-N400
127 positivity for unexpected vs. expected words, but in both high- and low-constraint
128 contexts. Moreover, the effect was found in only two of four conditions involving
129 unexpected words, despite all unexpected words being plausible.

130 Not only is there inconsistency in how constraint affects the PNP, sometimes
131 constraint-based effects are not elicited at all. In an experiment using the same materials
132 as Federmeier et al. (2007), Wlotko and Federmeier (2007) did not find any evidence of an
133 effect of constraint on the PNP. The lack of a constraint effect on the PNP was perhaps
134 particularly surprising given that constraint was found to affect the earlier P2 component.
135 This dissociation is interesting given that early and late positivities may share a neural
136 generative process, although this is the subject of much debate (Coulson et al., 1998;
137 Osterhout, 1999; Osterhout et al., 1996; Sassenhagen & Fiebach, 2019). If the PNP does
138 indeed share a generative process with the P2, it is therefore surprising that the effect of
139 constraint was not observed in both.

140 In a study more specifically investigating the PNP, Thornhill and Van Petten (2012)
141 also failed to find any constraint-related difference in PNP amplitude. The authors raise
142 the possibility that the concept of “weak expectation” may need close attention in

143 designing low-constraint experimental stimuli. Low constraint is typically measured using
144 cloze probability; however, the authors suggest that low cloze probability may sometimes
145 reflect a lack of agreement between cloze test participants on the best way to continue a
146 sentence, rather than a “weak” mental representation of the event. More recently, it has
147 been suggested that the *richness* of the mental representation may also determine whether
148 the PNP is seen at an unexpected word (Brothers et al., 2020). For example, in (3)a below,
149 expectation for the upcoming word can only be derived from the three words immediately
150 preceding it. In contrast, in (3)b, a richer context is built across the whole of the preceding
151 sentence. A constraint effect on the PNP was only seen at the unexpected word in (3)b and
152 not in (3)a, suggesting that the richer context allowed a more committed event
153 representation in (3)b, which required a greater update in order to accommodate the
154 unexpected word (Brothers et al., 2020):

155 (3)

156 a. *Locally constraining:*

157 He was thinking about what needed to be done on his way home. He finally arrived.

158 James unlocked the door/laptop

159 b. *Globally constraining:*

160 Tim really enjoyed baking apple pie for his family. He had just finished mixing the

161 ingredients for the crust. To proceed, he flattened the dough/foil

162 One possible explanation for the inconsistency among studies observing a PNP is
163 that its temporal proximity to the N400 makes it susceptible to component overlap
164 (DeLong et al., 2011; Luck, 2005a). Depending on the study design, this may mean that a
165 difference in the PNP is simply the result of an earlier difference in the N400. Other
166 explanations for the inconsistency are that the PNP is simply a broadly distributed P600,
167 or even a methodological artifact. One further complication is that the PNP may have a

168 relationship with the P3 family of components which is as yet unclear (Coulson et al.,
169 1998; Garnsey, 1993; Kuperberg et al., 2020; Kutas & Hillyard, 1980; Osterhout, 1999;
170 Osterhout et al., 1996; Sassenhagen & Fiebach, 2019; Van Petten & Luka, 2012). With
171 these issues in mind, in the present study we treat the N400 and PNP—with temporal and
172 spatial signatures defined by previous research—as distinct measures that can be used to
173 disentangle the influence of contextual constraint. Crucially, the PNP effect should be
174 manipulated by constraint while the N400 should not. Even if the N400 and PNP do arise
175 from generators that exhibit variable latency, finding evidence that they are affected
176 differentially by constraint will still allow conclusions about the usefulness of the PNP in
177 investigating readers' probabilistic representations. On the other hand, variable latency
178 may obscure any true effect and we may find no support for our hypotheses. In this case, a
179 null result would provide a starting point for future designs or analyses to more explicitly
180 address the contribution of latency variation. With this in mind, we make no claims about
181 the possibility of component overlap or latency variation with respect to the current study.

182 To summarise, while there is evidence to suggest that the PNP may be sensitive to
183 the strength of readers' probabilistic sentence representations, there is still inconsistency
184 within the PNP literature. The operationalisation of contextual constraint may also
185 require more careful consideration. Providing strong evidence for an association between
186 the PNP and contextual constraint, and thus a link between the PNP and representation
187 strength, would provide a crucial tool for future research into understanding how
188 probabilistic representations are built, and how readers' expectations about the upcoming
189 sentence influences their processing of incoming language input.

190 Moreover, providing further evidence for the PNP establishes a basis with which to
191 investigate the neurobiology of post-N400 positive deflections, including the P600. For
192 example, the link between the PNP and “suppression” (Kuperberg et al., 2020) or
193 “inhibition” (Kutas, 1993; Ness & Meltzer-Asscher, 2018) suggests engagement of executive
194 processes in the prefrontal cortex (e.g. Hagoort, 2013). These executive processes are

195 proposed to have a distinct cortical location and function from the types of processes to
196 which the P600 is sensitive (Hagoort, 2013; Hagoort & Indefrey, 2014). The P600 is instead
197 proposed to index involvement of circuits between the left inferior prefrontal cortex and the
198 temporal lobe as information from memory is retrieved and integrated during attempts to
199 revise a disconfirmed sentence representation (Brouwer et al., 2017; Brouwer & Hoeks,
200 2013). Strong evidence for the PNP would aid future investigations in this direction.

201 **The current study**

202 Recent research efforts have highlighted the fact that one of the critical findings in
203 research on probabilistic preactivation is difficult to replicate (Nieuwland et al., 2018) and
204 that the effect sizes of this predictability manipulation is likely much smaller than thought
205 (Nicenboim et al., 2020). Overestimated effect sizes and/or effects in an unexpected
206 direction can be the result of Type M(agnitude) and S(ign) errors in underpowered study
207 designs with too few participants and/or too few experimental items (Gelman & Carlin,
208 2014). ERP experiments are particularly susceptible to being underpowered given that
209 they are costly, both in terms of time, labour, equipment maintenance, and replacement of
210 disposable elements. Resource constraints therefore may prevent the recruitment of
211 sufficient number of participants to offset the high level of signal-to-noise ratio inherent in
212 ERP data (Luck, 2005a; Luck & Gaspelin, 2016). Many ERP studies also involve the
213 comparison of ERP components at target words that are not identical, which may
214 introduce additional noise through variability in frequency and lexical representations.
215 Investigation of the PNP would therefore greatly benefit from a confirmatory study using a
216 large number of participants.

217 We expected to show a dissociated effect of constraint on the N400 and PNP in a
218 relatively large number of participants (see *Participants* section below). The key findings
219 that we wished to replicate were those of Federmeier et al. (2007) and Kuperberg et al.
220 (2020), who found that only the PNP and not the N400 was affected by constraint. We
221 extended the design of Federmeier et al. by measuring PNP and N400 effects at matching

222 words with matching pre-critical regions, eliminating any potential lexical- or
223 frequency-based variation. Kuperberg et al. (2020) also measured ERPs at matching
224 words, but we extended their design by operationalising contextual constraint as the
225 continuous variable “entropy”. Entropy is a measure of uncertainty at the target word that
226 takes into account how the context of a sentence has affected the distribution of probable
227 words at that position (see the section *Cloze test* below for a more detailed definition). In
228 addition, we used constraint (entropy) and word predictability (log cloze probability) as
229 continuous rather than categorical predictors in the statistical analysis, which maximises
230 statistical power (Cohen, 1983). A discussion of the use of log cloze probability can be
231 found in *Section 2.6* on statistical analyses. A successful replication would make a solid
232 contribution to evidence that the PNP will be of great value in future investigations of
233 probabilistic processing.

234 Methods

235 The *Introduction* and *Methods* sections of this manuscript received Stage 1 approval
236 as a registered report and were pre-registered at
237 https://osf.io/bxg3n/?view_only=bf5946cadb3f47ccb44ad284e0ca9ec6.

238 Participants

239 In total, EEG was recorded from 74 participants. Seven participants were excluded
240 due to software problems during the recording and three because >75% of their EEG was
241 affected by artefact. This left a final sample size of 64. The participant sample size was
242 determined via a stopping rule based on the inference criteria used in our statistical
243 analysis (the Bayes factor), as well as time and resource limitations. We planned to recruit
244 participants either until we reached a Bayes factor of 10 in favour of the null or the
245 alternative hypotheses, or until we reached 150 participants, whichever came first. 150
246 participants was thought to be the maximum feasible number that we could collect data
247 from given limited resources and time. However, a major protocol deviation was made with
248 the approval of the editor and reviewers: A Bayes factor of 10 was exceeded for the PNP

249 constraint effect at 40 participants, but the Bayes factor for the N400 constraint effect
250 remained stable at approximately 1, regardless of sample size. Due to the difficulty in
251 recruiting participants during the Covid-19 pandemic and because it seemed unlikely that
252 the Bayes factor for the N400 constraint effect would reach 10 even with 150 participants,
253 we ceased recruitment early. We discuss the inconclusive Bayes factor further in the *Results*
254 section and present a design analysis which suggests that even over 150 participants would
255 not have been sufficient to reach the pre-registered Bayes factor threshold.

256 More detail on the statistical analysis is provided below, but support for our
257 hypotheses was assessed using Bayes factors for the effect of entropy (PNP prior: a
258 truncated normal distribution $N_-(0, 0.2)$; N400 prior: a normal distribution $N(0, 0.2)$), and
259 cloze probability (PNP prior: a truncated normal distribution $N_-(0, 0.2)$; N400 prior: a
260 truncated normal distribution $N_+(0, 0.2)$). Section *2.6.1 Statistical models and predictions*
261 provides further detail and motivates the use of truncated prior distributions.

262 Even with the protocol deviation, to our knowledge, the sample size is the largest
263 amount of data to date on this topic and we reached strong evidence (a Bayes factor of at
264 least 10, in line with Jeffreys, 1939) in favour of two pre-registered hypotheses without
265 reaching the maximum of 150 participants. For the hypotheses for which even 150
266 participants would not have yielded strong evidence, the experiment is still informative
267 because the estimates from our data can be used in a future meta-analysis in order to
268 synthesise the evidence available so far. For examples illustrating the importance of
269 evidence synthesis in psycholinguistics, see Bürki-Foschini et al. (2022), Jäger et al. (2017),
270 Nicenboim et al. (2020), and Vasishth and Engelmann (2022).

271 The inclusion criteria for participants in the study were: native German speakers
272 with no other language acquired before age 6, no history of developmental or acquired
273 reading, production, or hearing disorder, no history of developmental or acquired
274 neurological disorder, and no current need for or intake of psychopharmaceutical
275 medication. All participants' vision was normal or corrected to normal. Participants were

276 excluded from the final analysis if there were technical problems with the EEG recording, if
 277 more than 75% of EEG segments were badly affected by artifact, or if the attention check
 278 was failed (post-stimulus questions answered with an accuracy of less than 70%).

279 **Materials**

280 Each experimental item consisted of four sentences. An example item is below. In
 281 the example, target nouns for the respective analyses are in bold face:

282 (4)

283 **Strong constraint, high cloze probability noun:**

284 a. Auf Annetts Terrasse schien im Sommer zu viel Sonne, um noch draußen sitzen
 On Annett's terrace shone in summer too much sun in order outside sit
 285 zu können. Daher kaufte sie sich einen großen **Schirm** und...
 to be able. Therefore bought she herself a.MASC large.MASC **umbrella.MASC** and...

286 **Strong constraint, low cloze probability noun:**

287 b. Auf Annetts Terrasse schien im Sommer zu viel Sonne, um noch draußen sitzen
 On Annett's terrace shone in summer too much sun in order outside sit
 288 zu können. Daher kaufte sie sich einen großen **Hut** und...
 to be able. Therefore bought she herself a.MASC large.MASC **hat.MASC** and...

289 **Weak constraint, low cloze probability noun:**

290 c. Annett mag es gerne gemütlich, wenn sie etwas Zeit für sich findet. Daher
 Annett likes it really cozy when she some time for herself finds. Therefore
 291 kaufte sie sich einen großen **Schirm** und...
 bought she herself a.MASC large.MASC **umbrella.MASC** and...

292 **Weak constraint, low cloze probability noun:**

293 d. Annett mag es gerne gemütlich, wenn sie etwas Zeit für sich findet. Daher
 Annett likes it really cozy when she some time for herself finds. Therefore
 294 kaufte sie sich einen großen **Hut** und...
 bought she herself a.MASC large.MASC **hat.MASC** and...

295 *Cloze test*

296 To assess noun predictability, native German speakers completed sentences
 297 truncated after the determiner before the target noun. For the strongly constraining
 298 conditions, we used the publicly available stimuli from Nicenboim et al. (2020) and so the
 299 cloze procedure for the strongly constraining condition is as reported in that paper. For the
 300 weakly constraining condition, 60 new participants completed truncated sentences
 301 presented in Ixex (Drummond, 2016) either in the lab, or online via Prolific
 302 (www.prolific.co). Plural and singular forms of the same word were collapsed, as were
 303 nouns with the same stem (e.g. *Schirm* “umbrella” and *Sonnenschirm* “sun umbrella” or
 304 “parasol”). The cloze probability of the target noun in each condition was computed as the
 305 proportion of participants who gave that word or word stem out of the total number of
 306 participants.

307 To assess the contextual constraint of our conditions, we calculated entropy at the
 308 noun site. Entropy is a measure of uncertainty in terms of how the probability mass of
 309 cloze test responses is distributed. For example, in a strong constraint context, nine cloze
 310 test completions may be the word “umbrella” and one may be “hat”. Probability mass is
 311 therefore concentrated on “umbrella” and entropy is low (high constraint). In a weak
 312 constraint context, the cloze completions may be ten different words; now probability mass
 313 is evenly distributed and entropy is high (low constraint). We quantified Entropy (H) as
 314 the negative sum of cloze probabilities (P) for all nouns provided by participants for a
 315 particular sentence in the cloze test, multiplied by their respective logs: $H = - \sum_{i=1}^n P_i \log P_i$.
 316 For example, if nine cloze completions were “umbrella” and one was “hat” then:
 317 $H = -(P_{umbrella} \cdot \log P_{umbrella} + P_{hat} \cdot \log P_{hat}) = -(0.9 \cdot \log 0.9 + 0.1 \cdot \log 0.1) = 0.47$.
 318 Summary statistics for cloze probability and entropy are reported in Table 1 as well as in
 319 Appendix B, Figure ??.

Condition	log ₂ cloze probability		Proportion target word (%)		Entropy (bits)	
	Mean	95% range	Mean	95% range	Mean	95% range
a) Strong constraint, high predictable noun	-0.40	-1.00, -0.07	79.60	50.00, 100.00	0.68	0.00, 1.59
b) Strong constraint, low predictable noun	-3.71	-4.58, -2.50	5.47	4.17, 14.60	0.68	0.00, 1.59
c) Weak constraint, low predictable noun	-4.09	-5.09, -1.51	7.49	2.94, 34.20	2.44	1.47, 3.12
d) Weak constraint, low predictable noun	-4.46	-5.09, -2.34	4.93	2.94, 17.80	2.44	1.47, 3.12

Table 1

Cloze probability and entropy descriptive statistics. log₂ cloze probability is presented, as log₂ cloze probability will be used in the statistical model. Since cloze probability can only range between zero and one, log₂ cloze probability values will range between minus infinity and zero. The 95% range refers to the 2.5th and 97.5th percentiles of the data. Proportion target word refers to the raw percentage of cloze completions where the target word was given. Entropy reflects contextual constraint, where low values indicate strong constraint (low variety of completions given), and high values weak constraint (high variety of low probability completions given).

320 Design

321 Sentences were constructed in quartets, although the experimental design was
 322 non-factorial, with conditions a) and b), and b) and d) being collapsed in two respective
 323 analyses. Condition c) was presented for lexical balance:

324 a) Strong constraint, high predictable noun

325 b) Strong constraint, low predictable noun

326 c) Weak constraint, low predictable noun

327 d) Weak constraint, low predictable noun

328 Stimuli were presented in a Latin square design such that all participants saw only
329 one sentence from each item. There were 224 items in total. The collapsed conditions
330 meant that in each analysis, each participant would contribute data from 112 items. Since
331 all sentences were grammatical and plausible, filler sentences were not used.

332 Procedure

333 Participants were tested in a single session. For the EEG recording, participants
334 were seated in a shielded EEG cabin at distance of approximately 60 cm from a 56 cm
335 presentation screen. The experimental presentation paradigm was built using OpenSesame
336 (Mathôt et al., 2012). Each experimental session began with instruction screens advising
337 participants that they would read two related sentences for each trial: the first sentence
338 was presented several words at a time and the second (the critical sentence) was presented
339 word-by-word. Participants were advised that after some sentences, they must answer a
340 question as quickly and accurately as possible. Each experimental session began with five
341 practice trials.

342 Each trial in the experiment began with a 500 ms fixation cross in the centre of the
343 screen followed by a blank screen jittered with a mean of 1000 ms and standard deviation
344 of 250 ms. Each sentence was presented word-by-word for a duration of 190 ms per word
345 plus 20 ms for each letter. The target word, however, was presented for 700 ms regardless
346 of length so that the segment of EEG on which we conduct our analysis would not include
347 the onset of the following word. The inter-stimulus interval was 300 ms. After 50% of the
348 sentences, a yes/no comprehension question appeared; for example, *Hat Annett eine*
349 *Terrasse?* (Does Annett have a terrace?). Answering the question via a video game
350 controller triggered the beginning of the next trial. The order of presentation of sentences

351 within each list was fully randomised by the presentation software. Breaks were offered
352 after every 30 sentences.

353 Before starting the EEG experiment, participants performed a stop signal task
354 (Lappin & Eriksen, 1966; Logan & Cowan, 1984) that closely followed the design of
355 Verbruggen et al. (2008). The purpose of the stop signal task was to measure individual
356 differences in the ability to stop an action (a button press) once they had already initiated
357 it. This information was correlated with participants' PNP responses, with the hypothesis
358 that poorer performance on the stop signal task may correlate with smaller
359 constraint-related differences in the PNP; that is, if the PNP is related to suppressing the
360 mental representation of a sentence that has been falsified by unexpected input, people
361 who are better at inhibiting responses on the stop signal task might also show larger PNP
362 constraint effects. However, this was an exploratory analysis and we pre-registered no
363 specific analysis plan here. The testing session including EEG setup lasted approximately
364 three hours.

365 **EEG recording parameters and preprocessing pipeline**

366 EEG was recorded from 32 scalp sites by means of AgAgCl active electrodes
367 mounted in an elastic electrode cap at the standard 10-20 system (Jasper, 1958). Eye
368 movements and blinks were monitored with bipolar electrodes next to the left and right
369 outer canthus as well as below and above the right eye. EEG and EOG was recorded with
370 a TMSi Refa amplifier with active shielding at a sampling rate of 512 Hz and a low-pass
371 filter of 138 Hz, in line with manufacturer recommendations. Recordings were initially
372 referenced to the left mastoid and re-referenced offline to the average of the left and right
373 mastoid channels.

374 EEG was filtered offline using zero phase FIR filters with a bandpass of 0.01 – 30 Hz
375 on whole, unsegmented EEG blocks (i.e. continuous blocks recorded between participants'
376 breaks). The width of the transition band at the low cut-off frequency was 0.01 Hz and at
377 the high cut-off frequency, 7.5 Hz. Data was then segmented into whole sentences and

378 blinks and eye movements corrected using independent component analysis (ICA; Jung
379 et al., 2001) with the Fast ICA algorithm (Hyvärinen et al., 2001). ICA components were
380 inspected for each participant and removed if they strongly correlated with the ocular
381 channels. The data were then further segmented to extract the target region, and segments
382 were rejected if they contained a voltage difference of over 100 μV in a time window of 150
383 ms or containing a voltage step of over 50 $\mu\text{V}/\text{ms}$. In total, this pipeline resulted in the
384 rejection of 16% of the target noun segments, leaving approximately 3000 target segments
385 per condition. Corrected signal was then segmented and baseline-corrected relative to a
386 200 ms interval preceding the stimulus.

387 **Analyses**

388 The dependent variables in our planned analyses were:

- 389 • N400: Average ERP amplitude (μV) over electrodes Cz, CP1, CP2, P3, Pz, P4, and
390 POz in the window 300-500 ms following target word onset.
- 391 • PNP: Average ERP amplitude (μV) over electrodes Fpz, Fp1, Fp2, F3, Fz, F4 in the
392 window 600-1000 ms following target word onset.

393 As mentioned above, constraint was operationalised as entropy, where increasing
394 entropy reflected decreasing constraint. Noun predictability was operationalised as
395 smoothed cloze probability transformed to \log_2 . Additive smoothing was used with
396 pseudocounts set to one to avoid taking the log of zero (Laplace or Lindstone smoothing;
397 Chen & Goodman, 1999; Lindstone, 1920). The log transformation reflected the
398 assumption that the effect of cloze probability on N400 amplitude is continuous and
399 non-linear. In other words, changes in cloze probability at the upper end of the probability
400 scale will not affect N400 amplitude as much as changes at the lower end of the scale.
401 Thus, the model will estimate the same average change in amplitude for a difference in
402 cloze probability of 0.09 to 0.26 as for a change of 0.26 to 0.74, even though the latter
403 represents a larger change in raw cloze probability. Log transformed cloze probability has

404 previously been demonstrated to give a better fit to ERP data (Delaney-Busch et al., 2019;
405 Frank et al., 2015; Nicenboim et al., 2020), as well as to reading time data (Hale, 2001;
406 Levy, 2008; Smith & Levy, 2013), is consistent with Pareto and Zipf distributions of word
407 frequency (Baayen, 2001), and with scaling laws in other areas of cognitive research (Kello
408 et al., 2010).

409 Both entropy and log cloze probability were centred according to the mean of the
410 conditions included in the model (see below), such that the model estimated the one-unit
411 change in ERP amplitude at average values of log cloze probability and entropy (average
412 values are in Table 1 above).

413 *Statistical models and predictions*

414 Linear mixed effects models with correlated by-item intercept estimates and full
415 variance-covariance matrices for by-subject random effects were fit in the *rstan/Stan*
416 wrapper *brms* (Buerkner, 2018) in *R* (R Core Team, 2020).¹ Only random intercepts were
417 estimated for items because once the conditions were collapsed to treat entropy and cloze
418 probability as continuous predictors, there were only two entropy/cloze values per item
419 (corresponding to each sentence context). Since this was unlikely to be sufficient to
420 precisely calculate by-item random slopes, to reduce computation time we included by-item
421 intercepts only.

422 Our priors for the models were informed by the model estimates of previous
423 Bayesian ERP analyses, which suggested that intercept variability was higher than
424 individual variability between participants and items (Nicenboim et al., 2020). Using prior

¹ The complete list of software used for this paper is the following: *R* (Version 3.6.3; R Core Team, 2020) and the R-packages *bayesplot* (Version 1.8.1; Gabry et al., 2019), *brms* (Version 2.16.3; Buerkner, 2018), *eeguana* (Version 0.1.8.9001; Nicenboim, 2018), *job* (Version 0.3.0; Lindeløv, 2021), *lme4* (Version 1.1-30; Bates et al., 2015), *LSAfun* (Version 0.6.3; Günther et al., 2015), *patchwork* (Version 1.1.1; Pedersen, 2022), *rstan* (Version 2.21.3; Stan Development Team, 2020), *tidybayes* (Version 3.0.2; Kay, 2022), *tidyverse* (Version 1.3.1; Wickham et al., 2019)

425 predictive checks against simulated data, we then calibrated the priors so that they were in
 426 line with previous findings, but not strictly informative. These regularising priors were
 427 used to ensure stable and psycholinguistically plausible estimates (Chung et al., 2015;
 428 Gelman et al., 2008; Gelman et al., 2017). We confirmed that the joint behaviour of these
 429 priors in the model would generate plausible estimates using prior predictive checks
 430 (Gelman et al., 2017; Schad et al., 2020); see Figure 3. The priors were:

$$\begin{aligned} \textit{intercept} &\sim \textit{Normal}(0, 5) \\ \beta_{\textit{predictability}} &\sim \textit{Normal}(0, 1) \\ \beta_{\textit{constraint}} &\sim \textit{Normal}(0, 1) \\ \sigma_{\textit{subject,item}} &\sim \textit{Normal}_+(0, 0.5) \\ \sigma_{\textit{residual}} &\sim \textit{Normal}_+(8, 2) \\ \rho &\sim \textit{LKJ}(2) \end{aligned}$$

431 Models for estimation were fit with 50,000 iterations, including a warmup of 1000
 432 iterations. Model convergence was assessed by ensuring that the number of bulk and tail
 433 effective samples for every parameter estimate was at least 2000 and that \hat{R} values—the
 434 correlations of between- and within-chain variance—did not exceed 1.01. If these checks
 435 were violated, the number of iterations for each model was increased, or sampler behaviour
 436 modified, as indicated by warning messages from *brms*.

437 Support for our specific hypotheses (detailed below) was assessed using Bayes
 438 factors. As we had very specific, pre-registered hypotheses about the direction of these
 439 effects, the priors used for the Bayes factor analysis were truncated such that they
 440 constitute one-sided tests. As discussed above, conclusions about evidence for or against
 441 our hypotheses was based on Bayes factors computed using priors of $\textit{Normal}_-(0, 0.2)$ for
 442 the effect of entropy (constraint) and cloze probability (predictability) on the PNP, and
 443 $\textit{Normal}(0, 0.2)$ for the effect of entropy (constraint) and $\textit{Normal}_+(0, 0.2)$ for the effect of

444 cloze probability (predictability) on the N400, according to which of the questions (see
445 Sections 2.6.1.1 and 2.6.1.2) was being tested. These truncated priors were used for
446 hypothesis testing, but exploratory analyses with two-sided tests was also used to assess
447 evidence for non-hypothesised effects.

448 Models for the Bayes factor analyses were fit with 50,000 iterations in line with
449 Buerkner (2018) recommendations, including a warmup of 1000 iterations. Convergence
450 was assessed as for the estimation models—at least 2000 bulk and tail effective samples for
451 each parameter estimate, and $\hat{R} \leq 1.01$. Bayes factors were calculated using bridge
452 sampling (Bennett, 1976; Gronau et al., 2017; Meng & Wong, 1996). The strength of
453 evidence for or against our hypotheses was assessed with reference to Jeffreys (1939) scale,
454 where a Bayes factor indicating evidence at a ratio of 3:1 in favour of an effect is considered
455 the minimum meaningful support for that effect, and only 10:1 or larger values are
456 considered strong evidence. Given the sensitivity of the Bayes factor to the choice of prior
457 (Lee & Wagenmakers, 2014), we also computed Bayes factors for a range of different priors
458 on the effects of constraint (entropy) or predictability (cloze probability) while holding all
459 other priors (e.g. intercept, random effects) constant as defined above. The priors for these
460 sensitivity analyses ranged from $Normal(0, 0.2)$ to $Normal(0, 2)$, both truncated and
461 non-truncated.

462 **Effect of low predictability at the noun under differing constraint.** Our
463 main comparison of interest concerned the effect of constraint when noun predictability
464 was low. With respect to the N400, in line with previous research we expected that words
465 with similar cloze probabilities elicit N400s with similar amplitudes, regardless of how
466 constraining their context was. With respect to the PNP, if it is the case that the PNP
467 reflects the cost of revising a probabilistic event representation (Kuperberg et al., 2020),
468 then we should expect that low cloze probability words elicit a PNP that is larger in
469 contexts that are strongly constraining than in contexts that are weakly constraining.

470 For this comparison, we took sentences from conditions (b) and (d), which both had

471 low cloze probability nouns but varied in entropy (high entropy = weak constraint, low
 472 entropy = strong constraint); this can be seen in Figure 1A. Conditions (b) and (d) were
 473 collapsed together and ERP amplitude analysed as a function of continuous entropy.
 474 Although noun cloze probability in both conditions was low, there was some variability due
 475 to the differing contexts and thus log cloze probability was added as a continuous nuisance
 476 predictor in the models. In short, Figure 1A shows our predictions that when cloze
 477 probability is low:

- 478 • the N400 would be of equally high (negative) amplitude regardless of entropy
 479 (constraint). There may be a small effect of cloze probability;
- 480 • the PNP would become more positive as entropy decreases (i.e. as constraint
 481 increases). There may be a small effect of cloze probability.

482 Note that cloze probability and entropy are somewhat correlated (see Appendix B,
 483 Figure ??). This is because it is difficult to build stimuli that hold cloze probability
 484 constant while systematically varying entropy. However, our pre-registered hypotheses do
 485 not concern the effect of an interaction, and adding an interaction term to the model may
 486 only estimate variance otherwise explained by entropy (or cloze probability). For this
 487 reason, we chose to omit an interaction from the model.

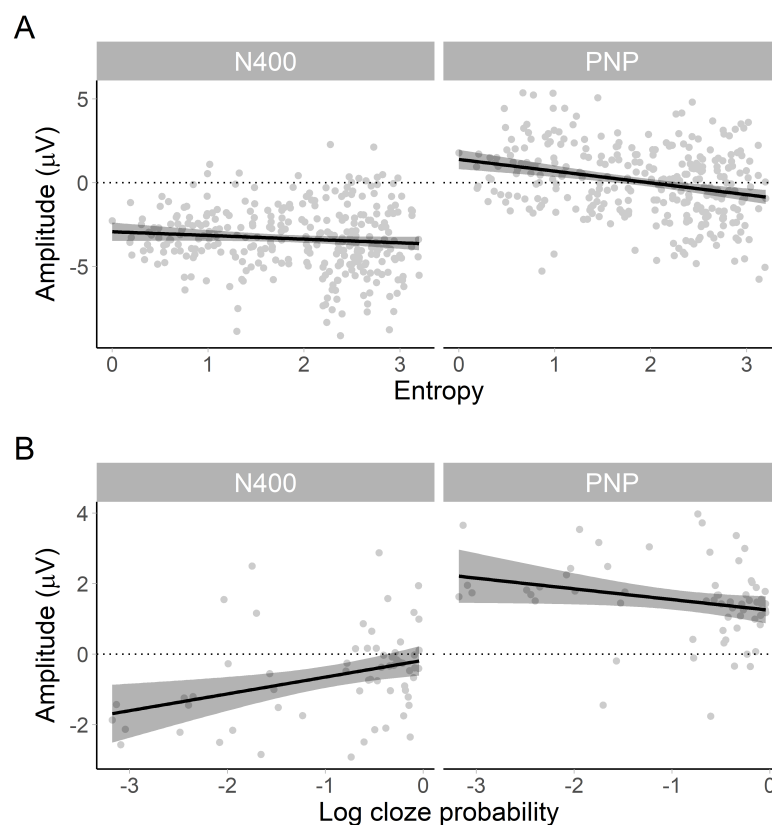
488 R *brms* model specification:

$$N400 \sim \text{constraint} + \text{predictability} + (1|\text{item}) + (1 + \text{constraint} + \text{predictability}|\text{subj})$$

$$PNP \sim \text{constraint} + \text{predictability} + (1|\text{item}) + (1 + \text{constraint} + \text{predictability}|\text{subj})$$

Figure 1

Simulated direction of the effect of constraint and predictability on average amplitude in the N400 and PNP time windows. A. In our first analysis, we collapsed conditions (b) and (d) such that predictability (cloze probability) was low but constraint (entropy) varied. Increasing entropy means decreasing constraint. Thus, as entropy increase on the x-axis, PNP amplitude should become less positive. In other words, the PNP at unexpected words should be more positive at low values of entropy (high constraint) than at high values of entropy (low constraint). N400 amplitude should not be affected by constraint, but may be sensitive to small differences in cloze probability between conditions (b) and (d). This was accounted for in the statistical analysis by adding cloze probability as a nuisance variable. B. In our second analysis, we collapsed conditions (a) and (b) such that constraint was high (low entropy), but predictability (cloze probability) varied. Cloze probability values are negative due to the log transformation. As cloze probability increases toward zero on the x-axis, the N400 becomes less negative and the PNP less positive. In other words, as predictability increases, the size of both the N400 and the PNP decrease.



489 **Effect of differing predictability at the noun under strong constraint.** As

490 a sanity check, we also compared conditions (a) and (b). It is well-established that
 491 decreasing cloze probability should increase amplitude of the N400 (i.e. make it more
 492 negative; Kutas & Federmeier, 2011) and of the PNP (i.e. make it more positive;
 493 Federmeier et al., 2007; Kuperberg et al., 2020). Under this assumption, when constraint
 494 was matched, we expected a larger N400 and PNP for low vs. high cloze probability words.
 495 For this comparison, we took sentences from conditions (a) and (b), which both had strong
 496 constraint but varied in cloze probability; see Figure 1B. Thus, conditions (a) and (b) were
 497 collapsed and ERP amplitude analysed as a function of continuous log cloze probability. As
 498 can be seen in Figure 1B, we expected that when constraint was strong:

- 499 • the N400 would become more negative as cloze probability decreases;
- 500 • the PNP would become more positive as cloze probability decreases.

501 R *brms* model specification:

$$N400 \sim \text{predictability} + (1|item) + (1 + \text{predictability}|subj)$$

$$PNP \sim \text{predictability} + (1|item) + (1 + \text{predictability}|subj)$$

502 ***Prior distributions and predictive check for the statistical models***

503 As an additional check that our prior specification would result in sensible estimates
 504 for our models, we conducted a prior predictive check (Gelman et al., 2017; Schad et al.,
 505 2020). In Figure 2, we show the prior distributions for each parameter in our statistical
 506 models. In Figure 3, we show the posterior distributions of a model simulating the
 507 predicted effect of entropy on the PNP and the N400 using only the priors. The estimated
 508 effect of entropy based on the priors (light blue lines) is plausible with respect to the effect
 509 based on simulated data (dark blue line), confirming that the joint behaviour of our priors
 510 in the model did not lead to implausible parameter estimates.

Figure 2

Prior distributions for the model parameters.

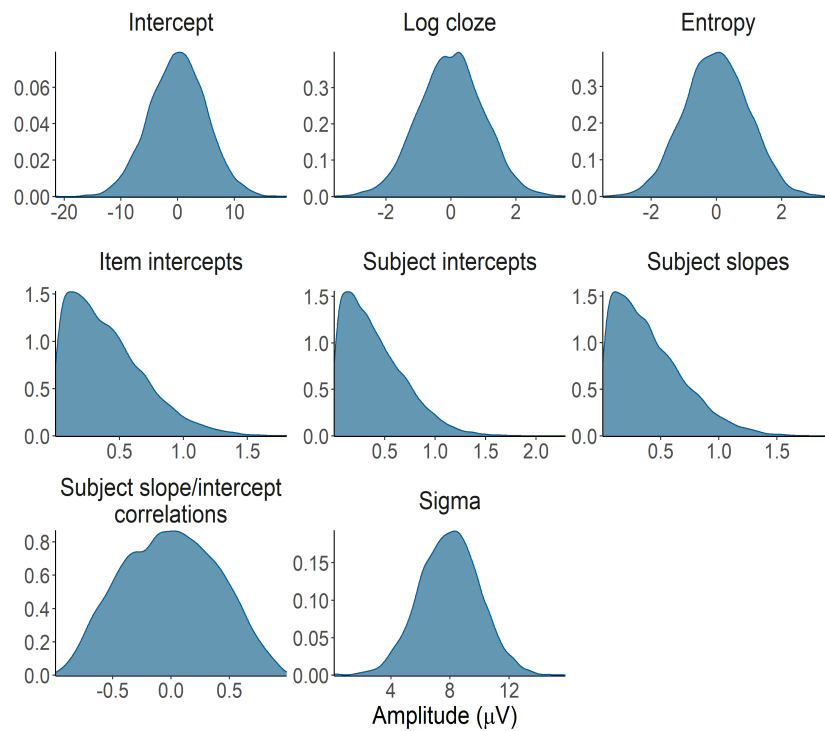
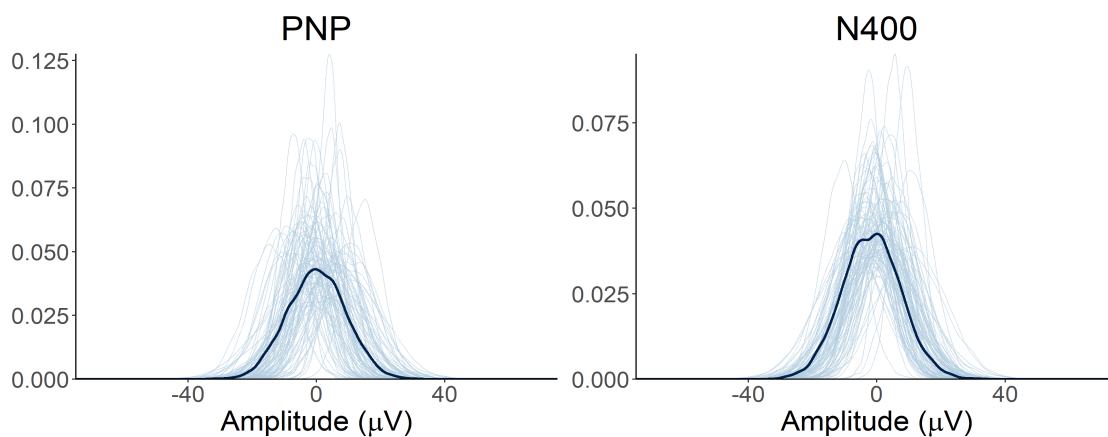


Figure 3

Prior predictive check. Prior predictive distributions for the effect of entropy on the PNP and N400 (light blue lines) based on the model priors suggests the priors generate plausible estimates consistent with simulated data (dark blue lines).



Results

511

512 In the following sections we report first the results of the pre-registered analyses,
513 then the results of our exploratory analyses. Data and code for all analyses are available at
514 https://osf.io/fndk5/?view_only=43f02800be0f4bd0b9309da36350778d.

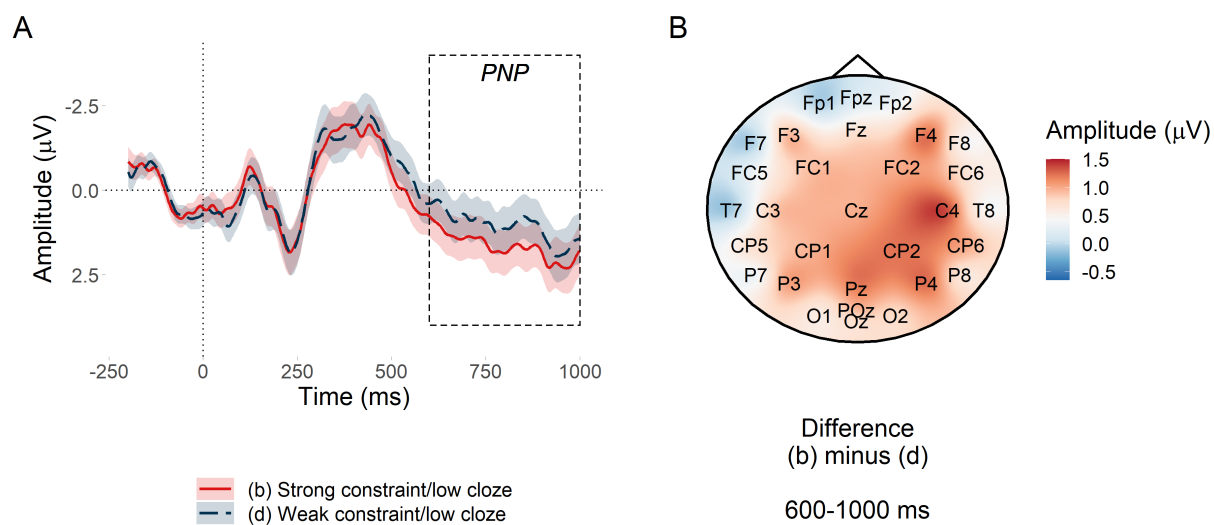
515 Pre-registered analyses

516 *Effect of low predictability at the noun under differing constraint*

517 **PNP window.** Figure 4A plots mean amplitude at the target word in the anterior
518 region of interest. The PNP was most positive for low probability words in low entropy
519 (strongly constraining) contexts and became less positive as entropy increased (constraint
520 weakened) by a estimated mean amplitude of $-0.26\mu V$ per bit of entropy, with a 95%
521 credible interval of $[-0.48, -0.05]\mu V$. Credible intervals reported throughout the
522 manuscript are quantile-based. The Bayes factor indicated strong evidence for H_1 over H_0 ,
523 $BF_{10} = 17.17$, consistent with Federmeier et al. (2007) and Kuperberg et al. (2020).
524 However, those studies predicted that the effect would be centred over anterior electrodes,
525 whereas Figure 4B suggests that in the current study, the scalp distribution of the
526 constraint effect was centred over posterior electrodes; we return to this in the exploratory
527 analyses. Sensitivity analyses testing the sensitivity of the Bayes factor to the choice of
528 prior for all pre-registered analyses are presented in Appendix C.

Figure 4

PNP constraint effect at low predictability nouns. **A.** Mean amplitude at the target low probability noun in the anterior region of interest. Since constraint in the statistical analysis was represented by the continuous predictor entropy, conditions (b) and (d) are divided by the median split of their entropy values. Ribbons indicate 95% confidence intervals. **B.** Subtraction plot of mean amplitude at low predictability target words between high and low median split entropy.

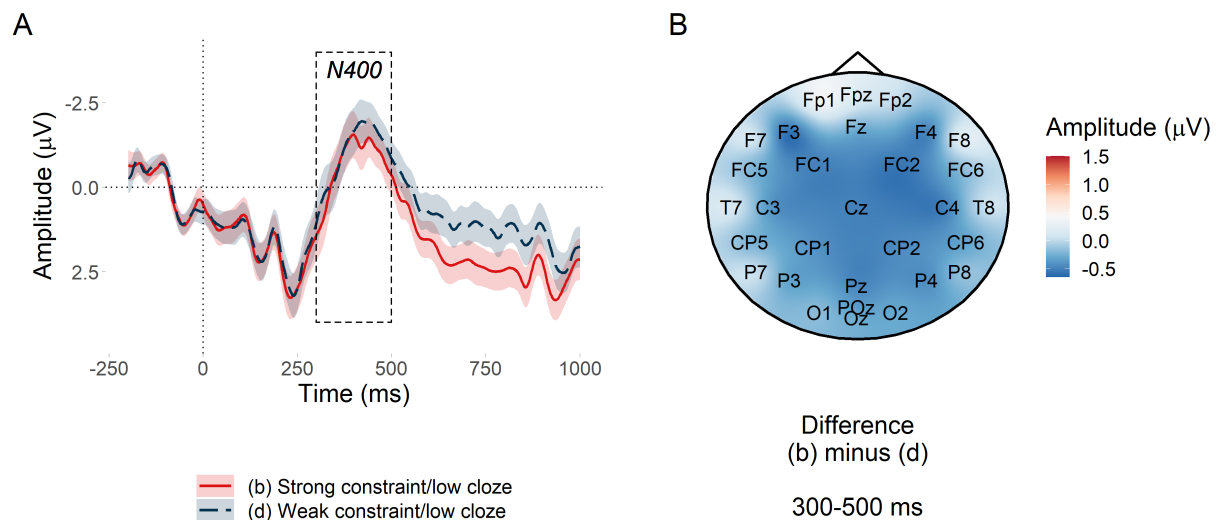


529 **N400 window.** Our pre-registered analysis yielded inconclusive evidence about the
 530 effect of constraint in the N400 window, $\hat{\beta} = -0.09[-0.30, 0.12]\mu V$, $BF_{10} = 0.76$. We
 531 attribute the inconclusive result to what appears to be between-condition differences in the
 532 behaviour of the N400 prior to and after its peak amplitude, as can be seen in Figure 5A.
 533 Prior to the peak, there was no visible effect of constraint. Past the peak however, from
 534 about 400 ms, there appeared to be a small constraint effect, which could be consistent
 535 with the beginning of post-N400 processing. Alternatively, it could reflect differences in
 536 mean latency of the N400 between the two conditions, with one condition peaking slightly
 537 later and thus having a higher amplitude for longer (we thank a reviewer for this

538 suggestion). Figure 5B shows a very small difference between high and low entropy in the
 539 N400 window with a topographic distribution typical of the N400.

Figure 5

N400 constraint effect at low predictability nouns. A. Mean amplitude at the target low probability noun in the posterior region of interest. Conditions (b) and (d) are divided by the median split of their entropy values. Ribbons indicate 95% confidence intervals. B. Subtraction plot of mean amplitude between the high and low constraint low predictability target words. Conditions (b) and (d) are divided by the median split of their entropy values.

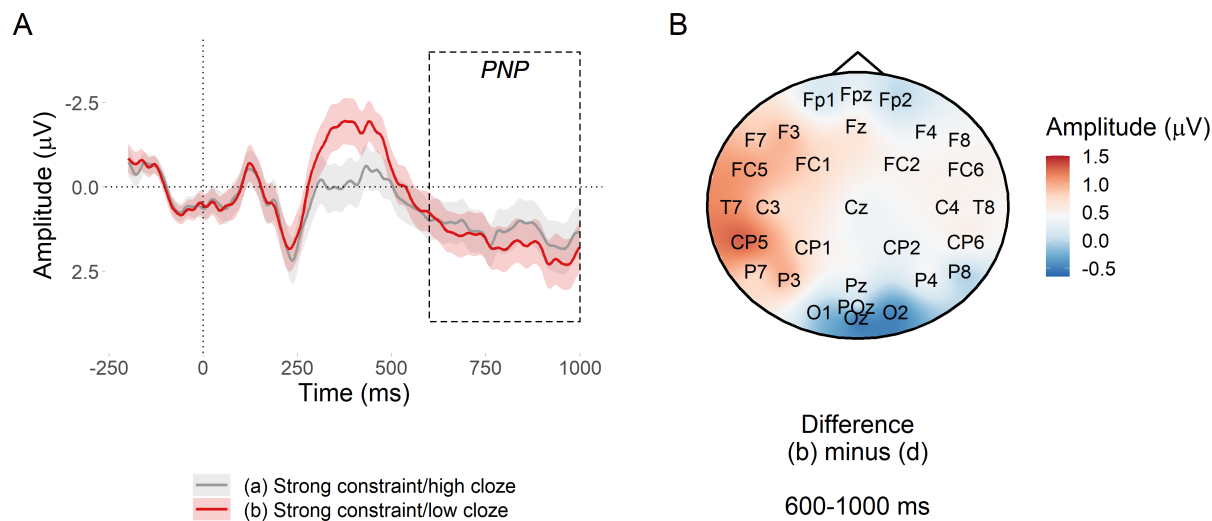


540 **Effect of differing predictability at the noun under strong constraint**

541 **PNP window.** Figure 6A suggests a small predictability effect in the expected
 542 direction with respect to Kuperberg et al. (2020), but the evidence was inconclusive,
 543 $\hat{\beta} = -0.11[-0.24, -0.01]\mu V, BF_{10} = 1.67$. However, Figure 6B suggests that there may
 544 have been a more left lateralised predictability effect; a similar predictability effect was also
 545 observed in Kuperberg et al. (2020) but was not analysed separately.

Figure 6

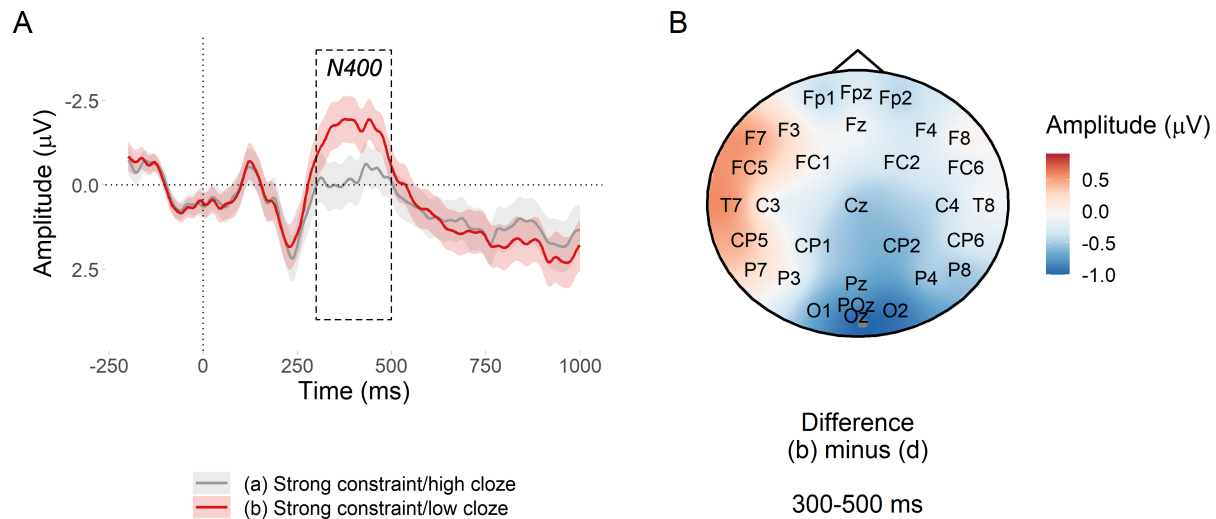
PNP predictability effect at nouns in strongly constraining contexts. A. Mean amplitude at the target noun in the posterior region of interest. Ribbons indicate 95% confidence intervals. B. Subtraction plot of mean amplitude between the high and low predictability target words.



546 **N400 window.** Low probability words in strongly constraining contexts elicited a
 547 large N400 in comparison to high probability words (Figure 7). There was extremely strong
 548 evidence for the effect, $\hat{\beta} = 0.56[0.41, 0.71]\mu V, BF_{10} > 20^7$.

Figure 7

N400 predictability effect at nouns in strongly constraining contexts. A. Mean amplitude at the target noun in the posterior region of interest. Ribbons indicate 95% confidence intervals. B. Subtraction plot of mean amplitude between the high and low predictability target words.



Discussion

549

550 Using the pre-registered analysis plan, we observed strong evidence that low
 551 probability words elicited more positive amplitude in the post-N400 window in strongly
 552 versus weakly constraining contexts. The direction of this effect was in line with previous
 553 research (Federmeier et al., 2007; Kuperberg et al., 2020), but its scalp distribution was
 554 consistent with a posterior P600 and not an anterior PNP. The effect of predictability in
 555 the PNP window was inconclusive, which contradicts Kuperberg et al. (2020). The N400
 556 window was more consistent with previous research: Although between-condition
 557 differences in the behaviour of the N400 before and after its peak amplitude were apparent
 558 in the latter part of the window, it did not appear that constraint affected the N400
 559 (Federmeier & Kutas, 1999; Federmeier et al., 2007; Kuperberg et al., 2020; Lai et al.,

2021; Szewczyk & Schriefers, 2013; Thornhill & Van Petten, 2012) and there was strong evidence for the standard N400 predictability effect (Kutas & Federmeier, 2011).

These findings support our hypotheses only partially. In support of our hypotheses, the constraint effect was apparent in the post-N400 window and not in the N400 window. This demonstrates a dissociated effect of probabilistic representation strength as processing progresses over time: It does not appear to affect initial semantic processing in 300-500 ms window (Kutas & Federmeier, 2011; Rabovsky et al., 2018), but it does appear to affect the downstream consequences of this processing in the 600-1000 ms window. Contrary to our hypotheses, the topography of the late positive effect was more consistent with a P600 than with the PNP reported in the literature. The P600 has been associated with conflict monitoring and syntactic reanalysis—a different type of processing than that proposed for the PNP (Bornkessel-Schlesewsky & Schlewsky, 2008; Brouwer et al., 2017; Fitz & Chang, 2019; Kim & Osterhout, 2005; Kuperberg et al., 2003; Osterhout & Holcomb, 1992).

Since a constraint effect on the P600 was unexpected in the current design, in the following section we first establish statistical evidence for the effect. We also examine whether word predictability affected the P600, since it was shown to affect the PNP in the previous research we had been trying to replicate. We then present a number of exploratory analyses probing different factors that could have resulted in the observed constraint effect being posterior (P600) rather than anterior (PNP).

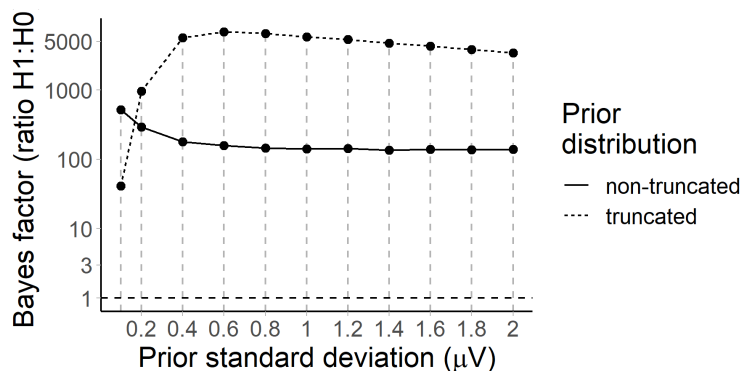
In other exploratory analyses, we examine the two effects for which we did not find conclusive evidence—the PNP predictability effect and the N400 constraint effect—and simulate datasets with larger sample sizes to determine what a sufficient sample size would have to be to yield conclusive evidence. Finally, we analyse the Stop Signal task to determine whether participants who were better at suppressing motor responses also showed larger constraint-based PNPs or P600s. We turn now to these exploratory analyses.

585 **Exploratory analyses**586 ***Statistical evidence for the P600 constraint effect***

587 We analysed average amplitude in the 600-1000 ms across the posterior region of
 588 interest (electrodes Cz, CP1, CP2, P3, Pz, P4, and POz). The model was that used for the
 589 PNP, but since we did not have a priori hypotheses about the direction or magnitude of
 590 the constraint effect, we examined a range of priors. Figure 8 suggests that there was
 591 strong evidence (BF_{10} from 41 to 5472) that low probability words elicited a more positive
 592 P600 in strong versus weak constraint regardless of prior, although the Bayes factor peaked
 593 around a prior standard deviation of $0.6\mu V$ (truncated to assume a negative effect),
 594 $\hat{\beta} = -0.60[-0.86, -0.34]\mu V$.

Figure 8

Bayes factors for the P600 constraint effect under a range of priors. The dashed line at a Bayes factor of 1 indicates equivalent evidence for H_1 and H_0 . Bayes factors above this line indicate evidence in favour of H_1 , with Bayes factors of over 10 generally considered to indicate strong evidence (Jeffreys, 1939).

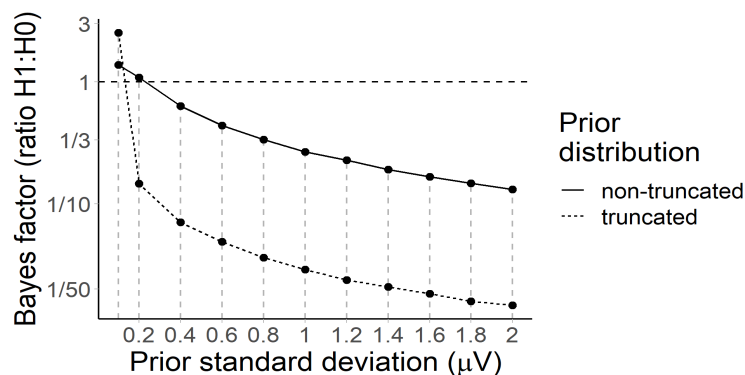
595 ***Predictability and the posterior P600***

596 In a previous study, both contextual constraint and word predictability affected the
 597 PNP (Kuperberg et al., 2020). Assuming that a similar underlying process drove the P600
 598 constraint effect in the current study, we additionally tested the effect of predictability in
 599 the 600-1000 ms window. We fit the same model as used to test the PNP predictability

600 effect, but used mean amplitude across posterior electrodes Cz, CP1, CP2, P3, Pz, P4, and
 601 POz. We used a range of priors and computed a Bayes factor for each. Figure 9 suggests
 602 that for prior standard deviations of $0.2\mu V$ or more that assumed a negative effect, there
 603 was strong evidence against a predictability effect, $\hat{\beta} = -0.11[-0.24, -0.01]\mu V$, prior:
 604 $\beta \sim Normal_-(0, 0.2)$. For priors that made no assumption about the direction of the
 605 effect, evidence against a predictability effect was weaker, but tended in the same direction
 606 as for priors assuming a positive effect.

Figure 9

Bayes factors for the P600 predictability effect under a range of priors. The horizontal dashed line at a Bayes factor of 1 indicates equivocal evidence for H_1 and H_0 . Above this line, evidence increases for H_1 , below this line, for H_0 . Evidence above 10 for H_1 or below $1/10$ for H_0 is generally considered to be strong. The plot panels show the estimated ratio of evidence for H_1 over H_0 (BF_{10}).



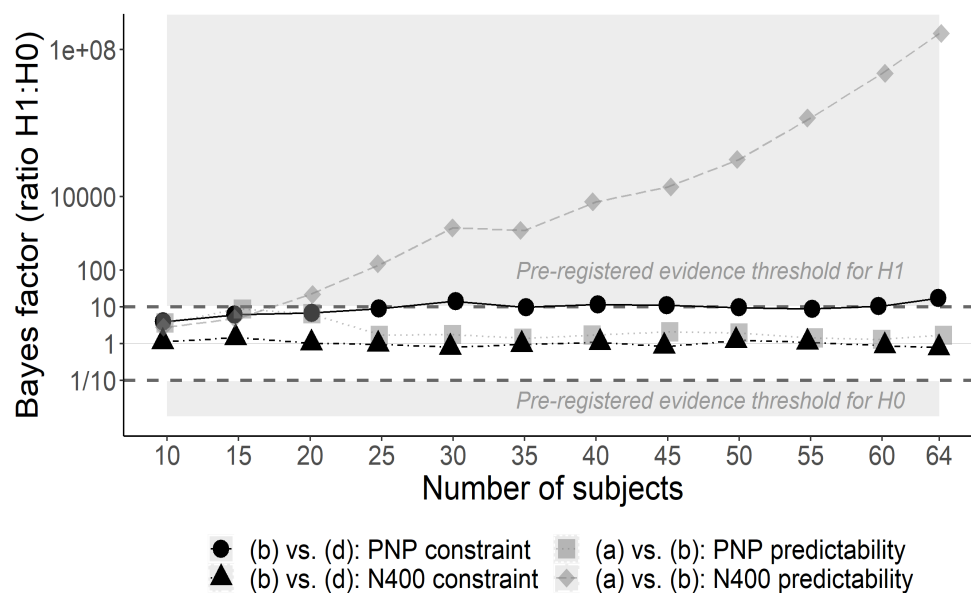
607 ***How many subjects would have been needed to yield conclusive evidence?***

608 Using our pre-registered analysis plan, we were unable to find conclusive evidence
 609 for two of our four pre-registered hypotheses. Figure 10 plots the Bayes factor for each of
 610 our four comparisons as sample size increased. Our two key comparisons are highlighted in
 611 black. Despite the Bayes factor remaining inconclusive for one of these key
 612 comparisons—the N400 constraint effect—we ceased recruitment due to the difficulty in
 613 recruiting participants during the Covid-19 pandemic. The post-peak N400

614 constraint-related differences may also have prevented the Bayes factor from ever being
 615 able to distinguish between null and alternative hypotheses, even if we had reached our
 616 pre-registered cap of 150 participants, which would have been infeasible given the poor
 617 recruitment rate.

Figure 10

Ratio of evidence for H1:H0 (Bayes factor) as sample size increases. The key contrasts regarding the effect of constraint on the PNP and N400 are in black.



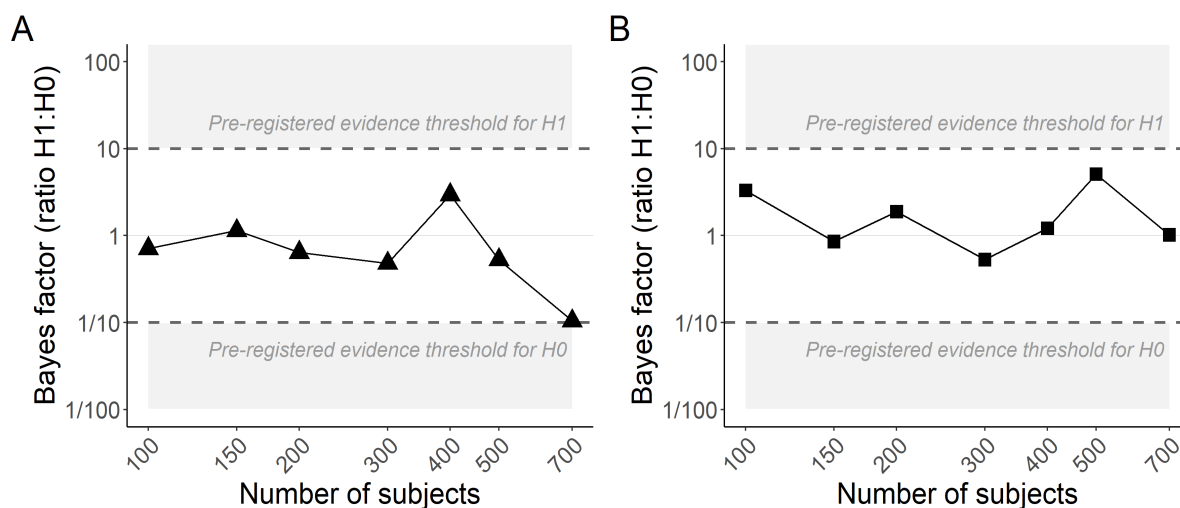
618 We therefore conducted a design analysis (Gelman & Carlin, 2014) to determine
 619 how many participants would be needed in a future experiment to yield conclusive evidence
 620 for the null hypothesis. We assumed that the estimates from the final sample of 64
 621 participants reflected true values and used them to simulate new datasets for between 100
 622 and 700 participants. A Bayes factor for the N400 constraint effect was computed for each
 623 sample size. Figure 11A suggests that even with the pre-registered cap of 150 participants,
 624 we would not have furnished strong evidence against the constraint effect on the N400
 625 using our pre-registered analysis plan. The analysis suggested that, assuming that the
 626 estimates obtained from the present data are indeed the true values, at least 700
 627 participants would be needed to demonstrate strong evidence against a constraint effect

628 using the current experimental design.

629 Since our secondary hypothesis about the PNP predictability effect also yielded
 630 inconclusive evidence with 64 participants, we repeated the same design analysis and noted
 631 that again, assuming our parameter estimates reflected the ground truth, the pre-registered
 632 cap of 150 participants would not have yielded conclusive evidence using the current
 633 design. Figure 13B suggests that if there were a true predictability effect, not even 700
 634 participants would have been sufficient to yield conclusive evidence for it.

Figure 11

Bayes factors at simulated sample sizes. A. *N400 constraint effect: One dataset was simulated for each sample size to which the pre-registered analysis model was fit. Each point in the plot reflects the Bayes factor for that sample size. B.* *PNP predictability effect: Each point reflects the Bayes factor for a pre-registered analysis applied to a simulated dataset.*



635 **Factors that could have changed the scalp appearance of the constraint effect**
 636 **or its underlying cognitive process**

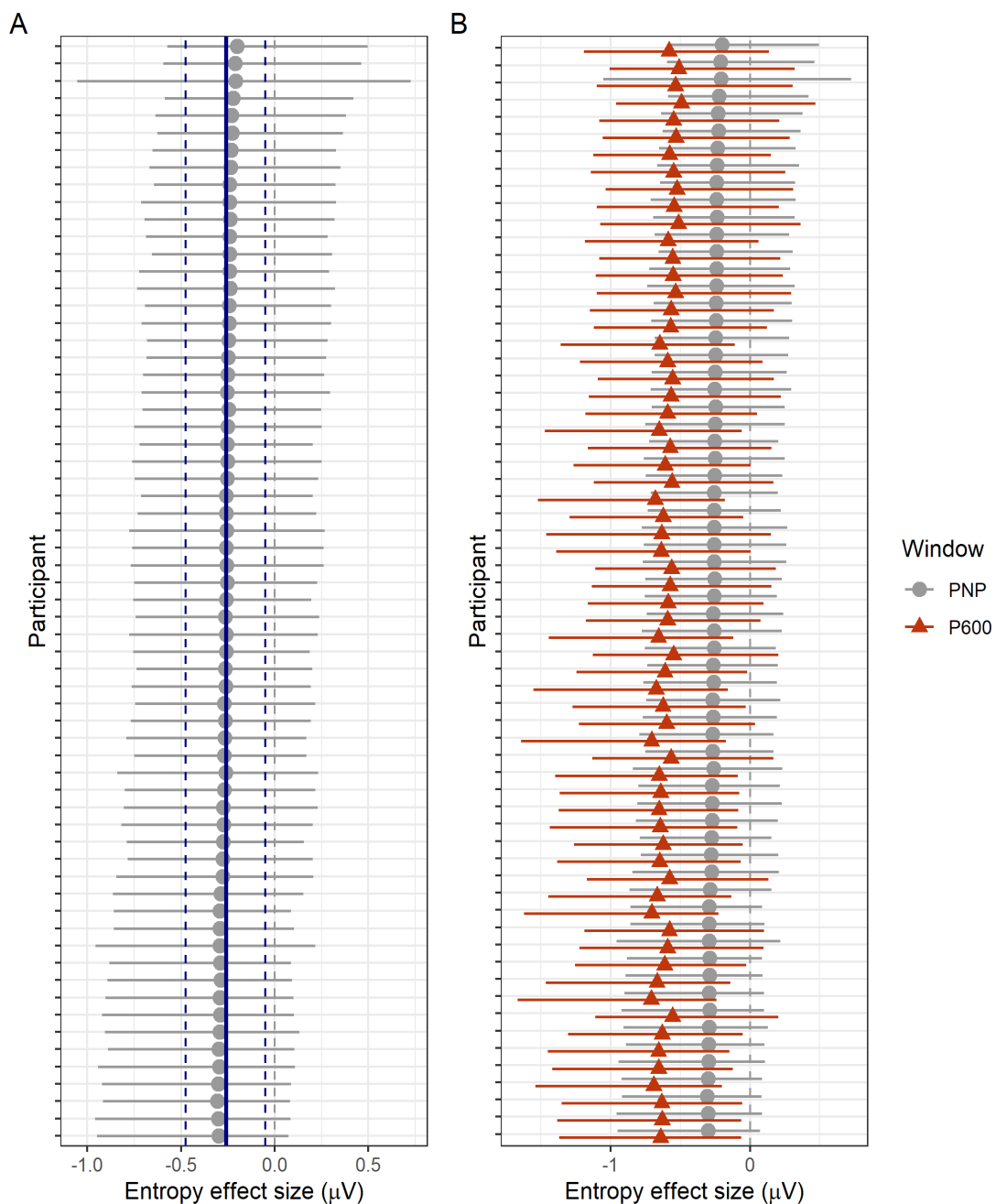
637 **Individual variability.** The scalp topography of an averaged ERP can be affected
 638 by factors such as variability in cortical folding and skull thickness between participants
 639 (Luck, 2005a). We examined individual variability by plotting posterior estimates of the

640 entropy effect by participant for the PNP Figure 12A and P600 Figure 12B. However,
641 individual estimates largely reflected the group mean with no obvious systematic outliers.

642 Another possibility is that individual participants differed in their response to the
643 unexpected word: some may have suppressed the disconfirmed sentence parse (PNP), while
644 others attempted to reanalyse the sentence (P600). If this were the case and we simply had
645 more P600-type processors in our participant pool, one could expect a crossover effect
646 where participants with smaller PNP constraint effects showed larger P600 constraint
647 effects, and vice versa. Individual PNP estimates are plotted against P600 estimates in
648 Figure 12B, but do not suggest a crossover effect. To quantify the relationship between
649 individual PNP and P600 constraint effects, we fit a multivariate linear mixed effects
650 model with the same form as the constraint models above, except that there were two
651 response variables: mean amplitude in the PNP and in the P600 windows/regions. A prior
652 for the correlation of the PNP and P600 constraint effects was also added: $LKJ(2)$. A
653 crossover between the PNP and P600 constraint effects would yield a negative correlation
654 estimate; instead, the model suggested a positive correlation, $\hat{\rho} = 0.61[0.60, 0.63]$. In other
655 words, participants with larger PNPs also tended to exhibit larger P600s.

Figure 12

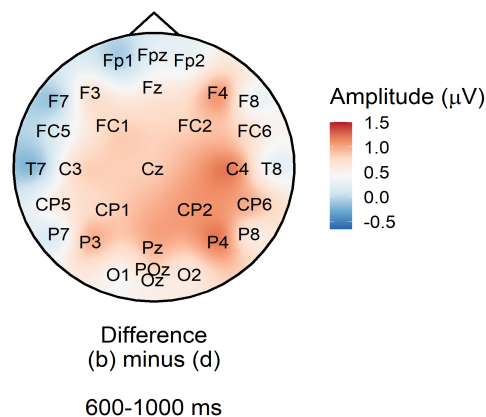
Individual posterior estimates of the effect of entropy in the post-N400 window. A. Individual posteriors from the pre-registered model of the anterior PNP (grey) are plotted against the group estimate (blue). Points show posterior means and errorbars the 95% credible intervals. B. Individual posterior estimates for the PNP (grey) are overlaid with estimates from the model fit to P600 amplitudes at the top of this section (orange).



656 **The operationalisation of constraint as entropy.** A major difference between
 657 the current study and Kuperberg et al. (2020) and Federmeier et al. (2007) is the use of
 658 entropy as a continuous measure of constraint. Instead, as in those studies, we could have
 659 used cloze probability of the most often given response, which, in the high constraint
 660 condition (b) was 0.80, 95% range = [0.50, 1.00] and in the low constraint condition (d),
 661 0.10, 95% range = [0.06, 0.50]. To determine whether a categorical definition of constraint
 662 would have changed the topography of the constraint effect, we re-plotted Figure 4B by
 663 subtracting condition (b) from condition (d) as defined by their category, rather than by a
 664 median split of entropy values. As can be seen in Figure 13, the distribution of the
 constraint effect was still posteriorly focused and was actually lower in magnitude.

Figure 13

Subtraction plot of mean amplitude at low predictability target words between high and low constraint as defined by category rather than entropy.



665

666 **Semantic association of target nouns with their context.** Another difference
 667 between the current study and Kuperberg et al. (2020) is that there was a semantic
 668 association between the target noun and its preceding context. Kuperberg et al. (2020)
 669 deliberately kept semantic association low. Assuming that low semantic association means
 670 weaker preactivation of the target word by the context, it could be that readers in
 671 Kuperberg et al. had to work harder to update their sentence representation at the

672 unexpected noun than participants in the current study, and that this extra work was
 673 necessary to elicit a detectable PNP constraint effect. If so, we could expect that low
 674 semantic association is a necessary condition for eliciting the constraint effect. In Table 2
 675 below, we compare semantic association of target nouns and their contexts across three
 676 studies: the current study, Kuperberg et al. (2020) and Federmeier et al. (2007). For our
 677 own stimuli, we computed cosine similarity using the *LSAfun* package in R (Günther et al.,
 678 2015). We used a pretrained German latent semantic analysis (LSA) space with 300
 679 dimensions (Günther, 2022) created from the 1.7 billion-word deWaC corpus (Baroni et al.,
 680 2009). Kuperberg et al. (2020) also computed cosine similarities using LSA and we present
 681 the values reported in their paper. For Federmeier et al. (2007), we computed cosine
 682 similarities using *LSAfun* and a pretrained English LSA space with 300 dimensions
 683 (Günther, 2022) created using the British National Corpus, the ukWaC corpus (Baroni
 684 et al., 2009), and a 2009 Wikipedia dump (we thank Kara Federmeier for providing the
 685 stimuli).

Table 2

Cosine similarity of target nouns with their context. Conditions names for all studies are presented in line with condition names from the current study.

Condition	Current study		Kuperberg et al. (2020)		Federmeier et al. (2007)	
	Mean	95% range	Mean	95% CI	Mean	95% range
a) Strong constraint, high cloze	0.40	0.17, 0.61	0.18	0.10, 0.26	0.40	0.18, 0.64
b) Strong constraint, low cloze	0.36	0.17, 0.58	0.01	-0.01, 0.03	0.33	0.17, 0.52
c) Weak constraint, low cloze	0.34	0.13, 0.54	-	-	0.36	0.14, 0.59
d) Weak constraint, low cloze	0.33	0.15, 0.56	0.01	-0.01, 0.03	0.34	0.12, 0.56

686

While semantic association in the current study was notably higher than in

687 Kuperberg et al., it was comparable with Federmeier et al., and yet Federmeier et al. saw a
688 distinct PNP constraint effect and no associated P600 effect. The degree of semantic
689 association between target noun and context thus may not explain our findings.

690 In the current experiment, it was possible to quantify whether cosine similarity
691 affected whether a constraint-based anterior PNP or posterior P600 effect was seen using
692 our model of the potential crossover effect above. We fit the same multivariate linear mixed
693 effects model with the two response variables mean amplitude in the PNP and P600
694 windows/regions, but added scaled and centred cosine similarity as a predictor interacting
695 with entropy. The main effect of cosine similarity was not consistent with a change in
696 amplitude, $\hat{\beta}_{PNP} = 0.10[-0.11, 0.32]\mu V$, $\hat{\beta}_{P600} = -0.12[-0.35, 0.11]\mu V$, nor was its
697 interaction with entropy, $\hat{\beta}_{PNP} = 0.02[-0.21, 0.24]\mu V$, $\hat{\beta}_{P600} = -0.05[-0.26, 0.17]\mu V$. As
698 before, the model yielded a strong positive correlation between the PNP and the P600,
699 $\hat{\rho} = 0.61[0.59, 0.63]$, suggesting that readers who exhibited larger PNPs still exhibited
700 larger P600s, even after semantic relatedness was taken into consideration.

701 **Task-related effects.** One of the factors that may play a role in the topography of
702 positive components in the post-N400 window is the type of task (Friederici et al., 2002;
703 Kuperberg & Brothers, 2019). During our experiment, participants answered a yes/no
704 question after 50% of sentences (28 sentences per condition). In Figure 14 below, we
705 compare topography and mean ERP amplitude in the late window between target nouns
706 that appeared in a sentence directly following a sentence that was one of the 50% of
707 question trials (Figure 14A and B), with nouns that appeared after a sentence with no
708 question (Figure 14C and D). Conditions b) and d) have been collapsed and split into high
709 and low constraint by their median entropy value. The posterior P600 effect is markedly
710 smaller in trials following a question (Figure 14B versus Figure 14D), suggesting readers
711 behaved differently when they may have expected another question versus when they did
712 not.

713 Participants' expectations with respect to an upcoming question could have had

714 various effects on their processing. For example, although questions were randomly
715 distributed, participants may have thought that question trials were more likely to appear
716 immediately after no-question trials and focussed more on the sentences, enhancing their
717 conflict-detection response and eliciting the P600 constraint effect after no-question trials
718 (Figure 14D). Alternatively, participants may have been primed to expect another question
719 trial if they had just seen one, and engaged a more PNP-type of processing such as
720 suppressing information not relevant to answering the question. This could explain the
721 absence of the P600 in post-question trials, although there was no suggestion of a PNP in
722 Figure 14B. Using the same model and priors as for the pre-registered PNP constraint
723 analysis, there was only inconclusive statistical evidence for the anterior PNP constraint
724 effect in the post-question trials, $\hat{\beta} = -0.23[-0.49, -0.02]\mu V$, $BF_{10} = 4$. When compared
725 with the strong evidence for the same effect when all trials were used (see main
726 pre-registered analysis), this finding does not suggest a functional dissociation between the
727 PNP and P600 on post-question and post-no-question trials.

Figure 14

Comparison of post-N400 amplitude at target nouns on trials after a trial

where a question was asked versus where no question was asked. A. ERP

amplitude in the anterior and posterior scalp regions on trials following a question. B.

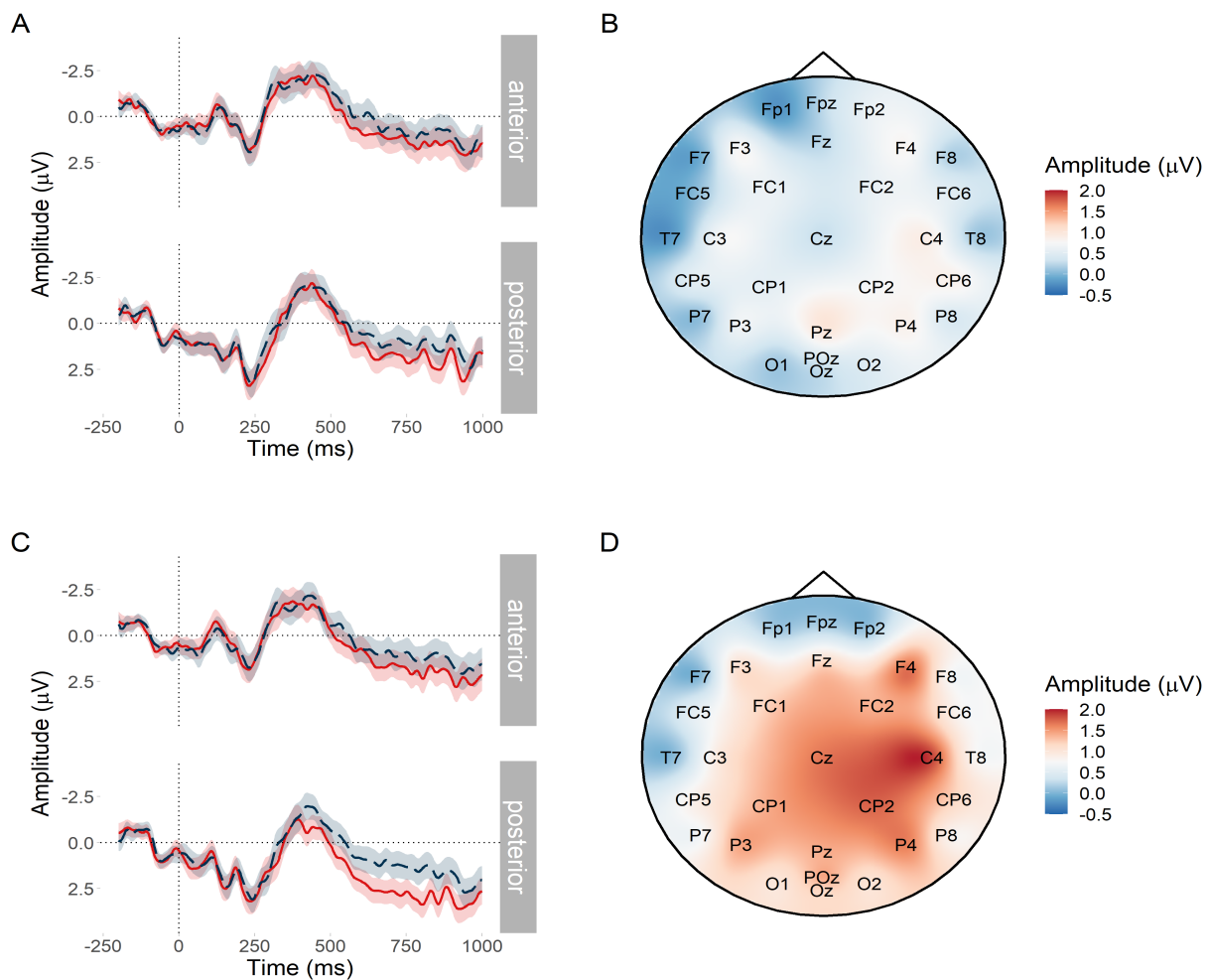
Topography of the constraint comparison on trials following a question (strong minus weak

constraint via median split entropy). C. ERP amplitude in the anterior and posterior scalp

regions on trials following a no-question trial. D. Topography of the constraint comparison

on trials following a no-question trial (strong minus weak constraint via median split

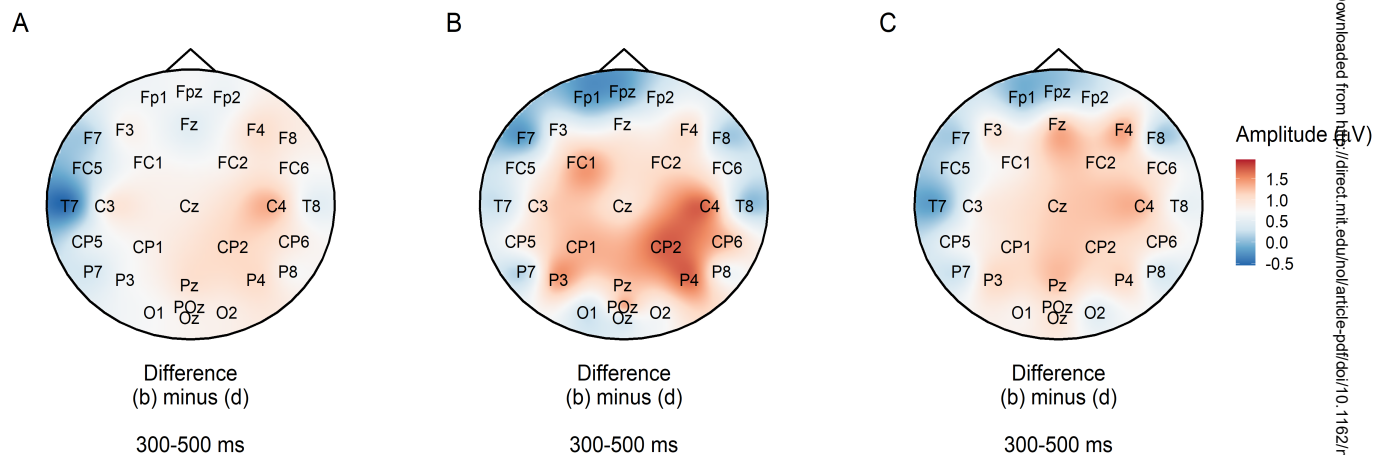
entropy).



728 **Trial order effects.** One reason for the absence of an anterior PNP may have been
729 due to participants not having engaged in predictive processing once they got used to or
730 guessed the purpose of the experiment. If so, this may have been visible across the
731 experiment, e.g. with an anterior PNP early on when participants were still predicting, and
732 a posterior P600 later as prediction stopped. Figure 15 suggests this was not the case, with
733 no PNP apparent at any stage of the experiment. We quantified a trial order effect by
734 adding trial number as an interaction with entropy to our pre-registered constraint model.
735 We fit two separate models, one of amplitude in the anterior region of interest (PNP) and
736 one of amplitude in the posterior region (P600). There appeared to be a main effect of trial
737 order in the anterior region, with amplitude becoming less positive as the experiment
738 progressed, $\hat{\beta} = -0.14[-0.36, 0.07]$, but this did not interact with entropy,
739 $\hat{\beta} = 0.005[-0.25, 0.27]$. In other words, there was no suggestion that a constraint effect on
740 the PNP differed across the experiment. In the posterior region, there appeared to be
741 neither a main effect of trial order, $\hat{\beta} = -0.04[-0.26, 0.16]$, nor an interaction of trial order
742 with entropy, $\hat{\beta} = 0.11[-0.15, 0.37]$.

Figure 15

Comparison of post-N400 amplitude at target nouns in different stages of the experiment. A. First third of the experiment. B. Middle third of the experiment. C. Final third of the experiment.

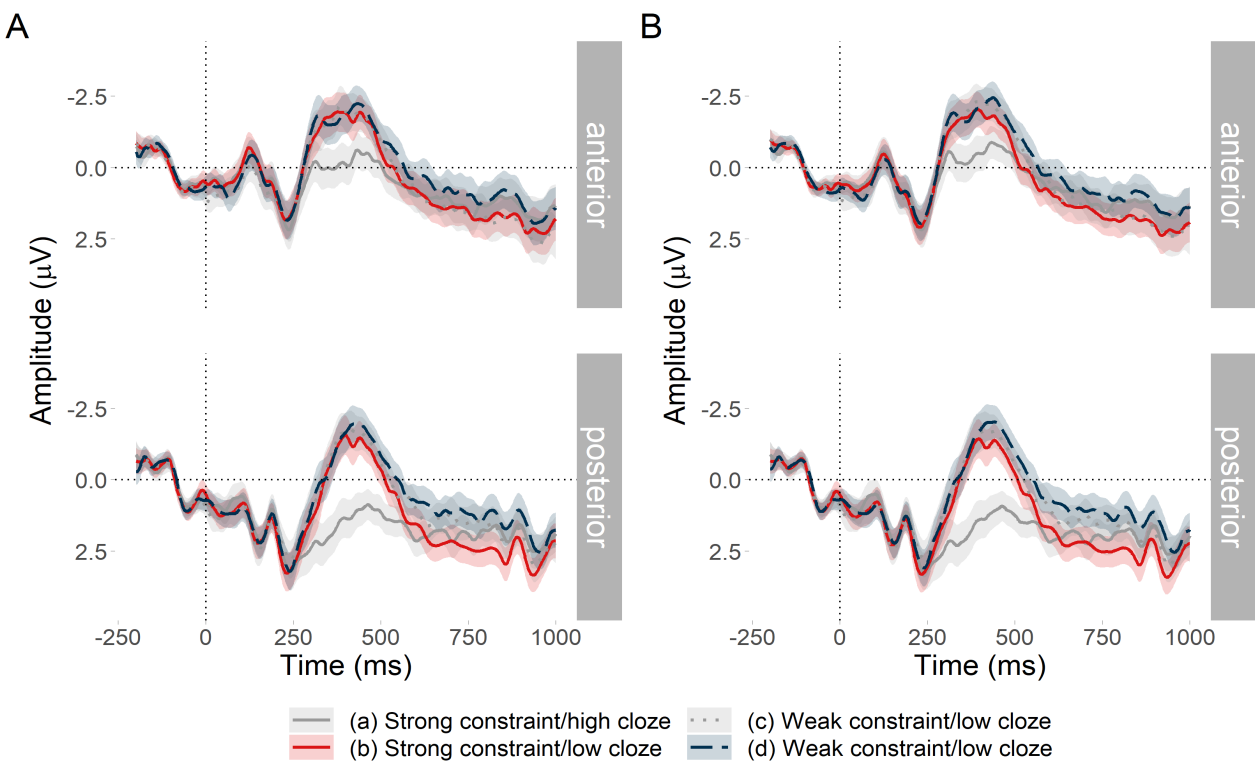


743 **The choice of temporal filter.** One ERP preprocessing step that can potentially
 744 alter the appearance of ERP components is the choice of filter (Luck, 2005a; Tanner et al.,
 745 2015; Vanrullen, 2011). Filter choice can create artificial differences, usually in the
 746 temporal appearance of ERP components, where amplitude from one time window is
 747 “smeared” into another as an artifact of the filtering process. The degree of smear depends
 748 on various filter settings and filter types, and can affect things like component overlap,
 749 which may have been present in our N400 window. Although smearing is more likely to
 750 affect the magnitude of an effect rather than its topography, we compared two different
 751 filter types. For our pre-registered preprocessing pipeline, we used finite impulse response
 752 (FIR) filters, but another common choice is infinite impulse response (IIR) filters. We
 753 re-processed the data using a Butterworth zero-phase (two-pass forward and reverse)
 754 non-causal IIR filter with filter order 16 (effective, after forward-backward) and cut-offs at
 755 0.01 and 30 Hz, (-6.02 dB).

756 ERPs after both types of preprocessing are plotted in Figure 16. Figure 16A shows
 757 the ERP using FIR filters (Figure 4B in the main text) and Figure 16B the ERP using IIR
 758 filters. We observed small differences in the amplitude of the ERP signal in each of our
 759 analysis windows, but nothing of a degree that would have changed our conclusions.

Figure 16

Comparison of FIR and IIR filters on the entropy effect in the post-N400 window. A. Mean amplitude over time at target words after preprocessing using FIR filters. B. Mean amplitude over time after preprocessing using IIR filters.



760 *Correlation of post-N400 amplitude with the stop signal task*

761 In a final exploratory analysis, we examined whether performance on a response
 762 inhibition task would predict the magnitude of the PNP constraint effect, with the
 763 hypothesis that better inhibitors might elicit larger PNPs. Before undergoing the EEG
 764 recording, participants completed a stop signal task. Participants saw either a circle or a

765 square on a screen and were instructed to press the “J” key on a keyboard as soon as they
766 saw a circle and the “F” key as soon as they saw a square, unless they heard a tone
767 presented via headphones, in which case they should not press anything. Our exploratory
768 hypothesis was that participants who performed better at suppressing their responses after
769 stop signals might also show larger PNP effects, if in fact the PNP were related to
770 suppression. The stop signal tone was a 750 Hz sine wave tone presented for 75 ms with no
771 attack or decay. The stop signal varied in its delay after the visual presentation,
772 determined via a tracking procedure: The starting delay was 250 ms and 50 ms was
773 subtracted after unsuccessful stop trials (i.e. trials where a response was made despite
774 hearing the tone), and 50 ms added after successful stop trials. The minimum stop signal
775 delay was 50 ms and the maximum, 1000 ms. The mean stop signal delay was 525 ms, 95%
776 CI = [511, 539] ms (see Table 3 for further descriptive statistics).

777 Participants were given four practice trials. The main experiment contained eight
778 trials per block and three blocks. Each block contained four circles and four squares
779 presented in random order. Stop signals were presented after one of the squares and one of
780 the circles. Each trial began with a fixation dot presented for 250 ms, followed by the
781 visual presentation. A keyboard response to the visual presentation triggered a blank
782 screen of 500 ms duration and the beginning of the next trial. If no response was made, the
783 next trial began after a timeout of 1250 ms. At the end of each block, participants were
784 given feedback about their proportions of incorrect responses, missed responses, and
785 correctly suppressed responses, as well as their average reaction time. The duration of the
786 feedback screen was determined by participants. The task was presented using
787 OpenSesame (Mathôt et al., 2012) on a 56 cm monitor in a sound-insulated cabin.

788 Of the 64 participants whose EEG was recorded, 59 had useable stop signal data:
789 one participant was excluded as they were unable to understand the stop signal task and
790 two participants’ stop signal data were not saved in error. Two further participants were
791 excluded because their mean response time on go trials was more than two standard

Table 3

Stop signal task descriptive statistics. Means and 95% confidence intervals are presented for the probability of (incorrectly) not responding on a go trial, the probability of (incorrectly) making any response on a stop trial, stop signal delay after visual presentation (SSD), stop signal reaction time (SSRT), reaction time (RT) of any response on go trials, and reaction time (RT) of any response on stop trials.

Measure	Mean	95% CI
Probability of no response on go trial	0.06	0.04, 0.08
Probability of response on stop trial	0.20	0.18, 0.21
Mean SSD	527	514, 541
SSRT	245	230, 260
RT on go trials	896	872, 920
RT on stop trials	781	754, 809

792 deviations faster or slower than stop trials, violating the assumptions of the stop signal
 793 reaction time calculation (Verbruggen et al., 2019). Stop signal reaction time (SSRT) was
 794 calculated via the integration method in Verbruggen et al. (2019).

795 We used SSRT as a predictor of amplitude in two separate models, one for the PNP
 796 and one for the P600. We used the same model specification as for the main analysis, but
 797 added log transformed SSRT as a continuous predictor interacting with entropy. All
 798 predictors were scaled and centred. Since there was only one SSRT observation per
 799 participant, random slopes were not estimated. With respect to the prior, we had no a
 800 priori expectation about the direction in which SSRT would affect amplitude: faster SSRTs

801 (better response inhibition) could hypothetically result in either a more marked inhibitory
802 response to unexpected input and higher amplitude, or a more efficient inhibitory response
803 and lower amplitude. We also did not expect that the effect of SSRT would be any larger
804 than that of entropy or cloze probability. We therefore used the same prior for SSRT as for
805 entropy and cloze probability, only non-truncated: $Normal(0, 0.2)\mu V$. Due to the mix of
806 truncated and non-truncated priors on the predictors, which *brms* did not allow at the time
807 of analysis, the model was fit in the *RStan* R package (Stan Development Team, 2018,
808 2020).

809 The posterior estimates of the interaction of entropy and SSRT on both PNP and
810 P600 amplitude were both centred around zero, $\hat{\beta}_{PNP} = 0.09[-0.12, 0.32]\mu V$ and
811 $\hat{\beta}_{P600} = 0.05[-0.18, 0.28]\mu V$, which was not consistent with faster SSRTs being predictive
812 of amplitude, regardless of constraint. Estimates were also not consistent with a main
813 effect of SSRT on amplitude in either the anterior, $\hat{\beta}_{PNP} = 0.11[-0.20, 0.41]\mu V$, or the
814 posterior scalp region, $\hat{\beta}_{P600} = -0.03, [-0.30, 0.25]\mu V$. In sum, the data were not
815 suggestive that faster performance on the stop signal task was associated with either PNP
816 or P600 amplitude. However, accuracy on the stop signal task was too high according to
817 guidelines set out by Verbruggen et al. (2019), which violates some assumptions in
818 computing SSRT. More specifically, the probability of responding after a stop signal should
819 be around 0.50, or at least between 0.25 and 0.75; our participants had a mean probability
820 of 0.20. The finding should thus be taken with caution.

821 General discussion

822 Our study addressed the idea that encountering a low predictability noun in a
823 context where a different noun was highly predictable should trigger greater processing cost
824 than a low probability noun in a context where no particular noun was predictable. We set
825 out to conceptually replicate the finding that a contextual constraint-based processing cost
826 at unexpected but still plausible words is reflected by an increase in anterior PNP
827 amplitude (Federmeier et al., 2007; Kuperberg et al., 2020). Using an experimental design

828 that maximised our ability to detect constraint effects and a sample size determined by
829 reaching a threshold for strong evidence, we were able to partially replicate previous
830 findings. We observed strong evidence that low probability words elicited more positive
831 amplitude in the post-N400 window in strongly versus weakly constraining contexts, but
832 the scalp distribution of this positivity was consistent with a posterior P600 and not an
833 anterior PNP. Also in contrast with previous findings (Kuperberg et al., 2020), the effect of
834 predictability in the post-N400 window was inconclusive, both for the PNP and the P600.
835 This suggests that the critical factor in determining processing at the target noun was not
836 how predictable that specific noun was, but rather how strongly the preceding context had
837 driven expectations about the event as a whole in which the target noun, and also other
838 words or concepts, might be expected. Findings in the N400 window were highly consistent
839 with previous research: constraint did not appear to affect the N400 (Federmeier & Kutas,
840 1999; Federmeier et al., 2007; Kuperberg et al., 2020; Lai et al., 2021; Szewczyk &
841 Schriefers, 2013; Thornhill & Van Petten, 2012) and there was strong evidence for the
842 standard N400 predictability effect (Kutas & Federmeier, 2011).

843 **Is the PNP affected by contextual constraint?**

844 The anterior PNP is proposed to be a distinct ERP phenomenon reflecting the cost
845 of shifting the interpretation of a sentence after unexpected input, becoming larger when
846 the preceding context increases certainty about a particular interpretation (Federmeier
847 et al., 2007; Kuperberg et al., 2020). We note here an assumption: that increased ERP
848 amplitude in one condition relative to another can be interpreted as increased processing
849 cost in the higher amplitude condition. However, a cost-amplitude association may not
850 reflect the true state of affairs since latency variability can create the appearance of
851 artificial amplitude differences (Luck, 2005b). The precise link between ERPs and neuronal
852 activity is also still unclear. However, for the purposes of this paper, we assume a
853 cost-amplitude link, based on the typical pattern that more “difficult” tasks (like dealing
854 with semantically unexpected words or odd syntax) reliably increase the amplitude of at

855 least the N400 and late positive components.

856 The mechanism underlying the PNP is proposed to be separate from that of another
857 post-N400 positive component—the posterior P600—since in two previous studies only the
858 PNP was affected by a constraint manipulation at plausible but unexpected words and not
859 the P600 (Federmeier et al., 2007; Kuperberg et al., 2020). In one of these studies, the
860 reverse observation was made for words that were anomalous in their contexts: constraint
861 affected the P600 but not the PNP (Kuperberg et al., 2020). Together, these findings have
862 been taken to suggest that the PNP reflects the successful update of a sentence
863 representation with plausible input and the P600 an error signal triggered by implausible
864 input. The current findings contrast with Kuperberg et al. and Federmeier et al. in two
865 ways: first, we did not observe a constraint effect for plausible words in the anterior PNP
866 but rather in the posterior P600, and second, the effect on the P600 was elicited by
867 plausible unexpected words. In this section we examine a number of possible explanations
868 for the contrasting findings.

869 With respect to the posterior rather than the anterior distribution of the constraint
870 effect, we ruled out with exploratory analyses that the difference was related to our
871 definition of constraint, or to individual variability in constraint effects. Since the type of
872 filter used during EEG preprocessing can also alter at least the temporal appearance of
873 ERPs (Luck, 2005a; Tanner et al., 2015; Vanrullen, 2011), we additionally re-processed
874 the data using a different filter, but the topography of the constraint effect remained
875 posterior. The combination of filter settings and the choice of baseline can create artificial
876 differences in ERP topography (Tanner et al., 2016): We used average amplitude over a
877 pre-stimulus period of 200 ms as a baseline and a bandpass filter of 0.01-30 Hz. Of the
878 previous studies in which constraint was examined, all used 100 or 200 ms pre-stimulus
879 baselines (100 ms for all but two studies), with which effects on the PNP both were and
880 were not observed; that is, there was no systematic effect of the baseline duration on
881 whether or not a PNP constraint effect was observed. Almost every study used different

882 bandpass filter settings which—while of concern for ERP research more broadly—again
883 does not suggest a systematic effect on the appearance of the PNP (although we did not
884 manipulate these settings directly and so cannot rule it out).

885 The type of task that participants do during the EEG recording can also affect the
886 appearance, including the topography, of late positive components (Friederici et al., 2002;
887 Kuperberg & Brothers, 2019; Sassenhagen & Bornkessel-Schlesewsky, 2015; Sassenhagen
888 et al., 2014), so we reviewed task types among previous studies. Participants in the current
889 study answered yes/no comprehension questions after 50% of sentences. In previous studies
890 where a constraint effect on the anterior PNP was observed (but not on the posterior
891 P600), participants had to judge whether each sentence “made sense” and additionally
892 answered yes/no questions about filler sentences (Kuperberg et al., 2020), or had no task
893 during the experiment but were informed they would complete a word recognition task
894 after the experiment (Federmeier et al., 2007). Of the previous studies that have observed
895 no or contrasting effects of constraint on the PNP/P600, participants either had to indicate
896 after each sentence whether a probe word appeared in that sentence (Thornhill &
897 Van Petten, 2012), or were informed they would complete a word recognition task after the
898 experiment (Federmeier & Kutas, 1999; Lai et al., 2021; Szewczyk & Schriefers, 2013;
899 Wlotko & Federmeier, 2007). Thus, there did not appear to be systematic differences in
900 task type between studies. In addition, we did not observe statistical evidence that the
901 presence of absence of a question influenced whether participants exhibited a PNP or P600
902 in the current study. Future studies directly manipulating the effect of task type on
903 eliciting the PNP versus the P600 would better address this question, however.

904 With respect to the P600 being elicited by plausible words, this is somewhat
905 unusual since the target noun and context were also syntactically well-formed and the P600
906 has traditionally been associated with reanalysis after syntactic violations (Hagoort et al.,
907 1993; Osterhout & Holcomb, 1992). However, P600s are also reliably observed at the verb
908 in role-reversal sentences which are syntactically well-formed, just semantically odd, e.g.

909 *the dog that the man bit* (Kim & Osterhout, 2005; Kuperberg et al., 2007; Kuperberg et al.,
910 2003). Van Petten and Luka (2012) also note a number of predictability studies where a
911 P600 was elicited by plausible unexpected words that did not involve a role reversal. Thus
912 rather than being limited to reanalysis after syntactic violations, the P600 has been
913 proposed to signal a more general conflict detection or integration process recruiting the
914 left inferior frontal gyrus (Brouwer et al., 2017; Brouwer & Hoeks, 2013; Fitz & Chang,
915 2019; van de Meerendonk et al., 2011). In our case, it likely reflects the conflict between
916 readers' strong event representation and the low probability input (Kuperberg et al., 2020;
917 Laszlo & Federmeier, 2009; Vissers et al., 2006).

918 The combination of strong constraint and high semantic relationship between target
919 words and their contexts in the current study are known to increase the likelihood of the
920 P600's appearance in syntactically well-formed sentences (Kuperberg & Brothers, 2019).
921 Since semantic association was higher in our study than in Kuperberg et al. (2020), we
922 reasoned that this could have contributed to the difference in topography. For example,
923 high semantic association would mean that lexical preactivation of the presented target
924 word by the context is stronger than when semantic association is weak, even in the low
925 predictability conditions. Stronger semantic association and stronger preactivation in our
926 study may not have required the engagement of PNP-related resources when a low
927 probability word triggered a shift in interpretation. In contrast, weaker semantic
928 association and weaker preactivation in Kuperberg et al. (2020) may have made the shift
929 costlier and the PNP more pronounced. However, we compared semantic association
930 between target words and their contexts in the current study against Kuperberg et al.
931 (2020) and Federmeier et al. (2007; Table 2) and noted that semantic relationship in
932 Federmeier et al.'s stimuli was comparable with our study—yet they observed a PNP and
933 not a P600. Future experiments comparing plausible, low probability words with strong
934 and weak semantic association with their contexts may yield further insights.

935 One likely factor contributing to the difference between the current and previous

936 studies is that of statistical power: fewer participants and/or fewer critical trials in
937 previous studies may have led to a lower signal-to-noise ratio in the EEG recordings. It is a
938 known issue in ERP research that if the signal-to-noise ratio is not sufficiently high, scalp
939 topography can be misleading and statistical false positives can occur (Luck, 2005a; Luck
940 & Gaspelin, 2016). False positives occur when low power leads to an overestimate of the
941 effect size or a type M (magnitude) error (Gelman & Carlin, 2014). Type S (sign) errors
942 may also result, explaining why at least one previous study reports a PNP constraint effect
943 in the opposite direction (Federmeier & Kutas, 1999).

944 The current study therefore raises the possibility that the PNP constraint effect
945 observed in previous studies may actually be part of a broad P600 response where lower
946 sample size has contributed to Type M and S errors in the anterior region of the scalp.
947 This would account for the anterior PNP constraint effect's inconsistent appearance in
948 previous studies despite similar experimental designs. If true, then our findings also suggest
949 that the processing cost of strong probabilistic representations does not always result from
950 having to update an interpretation or suppress disconfirmed interpretations after receiving
951 conflicting input, but can instead be the cost of detecting the conflict itself.

952 **Why was a constraint-based P600 effect not observed in previous studies?**

953 If the anterior PNP constraint effect really is just the edge of a P600 constraint
954 effect, then one would expect to see a P600 constraint effect in at least some previous
955 studies. One previous study did in fact observe a P600 constraint effect, but only at
956 anomalous (implausible) words (Kuperberg et al., 2020). For anomalous words, the P600
957 became more positive for anomalous words in highly constraining contexts. This is
958 consistent with the P600 constraint effect elicited by syntactic violations (Gunter et al.,
959 2000; Hoeks et al., 2004); although in Hoeks et al. (2004) the effect was in the opposite
960 direction and statistical evidence was not strong. One possibility therefore is that the
961 anomalous sentences in Kuperberg et al. (2020) encouraged participants to treat
962 unexpected-but-still-plausible words differently to the “real” conflict posed by anomalous

963 words (as Kuperberg et al. hypothesised it would). In the absence of anomalous words in
964 the current study, participants may have responded to low probability words in the same
965 way as if they were errors. However, this would not account for why a P600 constraint
966 effect was not observed in Federmeier et al. (2007)—who also did not have an anomalous
967 condition—nor in other previous studies without anomalous conditions who observed
968 contrasting or no PNP effects. This may again be due to a power issue, but we have also
969 made suggestions above as to future research that could help to disentangle the PNP and
970 P600.

971 Aside from the absence of anomalous words, another possibility is that features of
972 the current study design encouraged conflict monitoring rather than prediction in
973 participants. Generating predictions while reading is thought to be one of the necessary
974 conditions for eliciting the PNP (Federmeier, 2022). It is possible that the large number of
975 sentences and simple manipulation in the current design meant participants stopped
976 predicting once they got used to (or even guessed the purpose of) the experiment. If this
977 were the case, one might expect a constraint-based PNP early in the experiment and a
978 constraint-based P600 later; we examined trial order effects and while the P600 constraint
979 effect was visually most pronounced in the middle of the experiment, no PNP constraint
980 effect was obvious either visually or statistically. Moreover, in order for readers to have
981 shown a larger P600 in the strong constraint condition at the low predictability target, with
982 all else being equal, the strong constraint of the context must have been used to generate
983 some degree of expectation for a different upcoming word relative to the weak constraint
984 context. This would suggest that readers were indeed predicting upcoming words. One
985 hypothesis for a future experiment is that there is a difference between passive expectations
986 when participants settle into a long experiment, and active predictions in more challenging
987 experimental designs. One could imagine that the former encourages conflict-monitoring
988 and thus a P600 response and the latter, suppression of previous predictions and a PNP
989 response. There is some evidence that conscious prediction strategies modulate the PNP

990 (Brothers et al., 2017), though not to the point of eliciting a P600 instead.

991 **The effect of probabilistic strength on processing cost**

992 Topography aside, the firm conclusion from the current and previous studies is that
993 the effect of probabilistic representation strength on processing cost only becomes
994 observable in the time window after the N400. The lack of a constraint effect in the N400
995 window is consistent with existing accounts of the N400 suggesting that the underlying
996 cognitive processes are seated in the medial temporal gyral and posterior temporal areas of
997 the ventral stream, at a time where phonetic and orthographic activation gives way to
998 lexical retrieval and semantic unification (Friederici, 2012; Hagoort, 2013; Hickok &
999 Poeppel, 2007; Lau et al., 2008). Retrieval and unification generate a probabilistic
1000 representation of the sentence, which in turn influences the activation of related words and
1001 concepts. Under these accounts, the N400 is only sensitive to the level of activation in this
1002 area, such that two words with the same activation level will elicit the same amplitude
1003 N400, regardless of how they came to be activated (Fitz & Chang, 2019; Hagoort et al.,
1004 2009; Kutas & Federmeier, 2011; Lau et al., 2008; Rabovsky et al., 2018). In our study, the
1005 low probability target in a strongly constraining context would have had low activation
1006 because the context suggested it was not a likely next word, whereas the low probability
1007 target in a weakly constraining context would have had low activation because the context
1008 did not suggest any particular next word. Their respective N400s were therefore of a
1009 similar amplitude.

1010 Further down the ventral stream, in the post-N400 processing time window, is
1011 where we observed the consequences of a strong probabilistic representation. In the current
1012 study, a strong representation increased sensitivity to input that conflicted with
1013 expectations (assuming a conflict-based function of the P600). Interestingly, low
1014 predictability lexical items seen in strongly constraining contexts did not elicit conclusive
1015 differences in P600 amplitude relative to high predictability lexical items, suggesting that
1016 conflict was driven by the semantic richness of the preceding context rather than a simple

1017 unexpectedness detection. This indicates a change in processing with respect to the
1018 previous N400 window, where word predictability was important.

1019 Source localisation of processing associated with the P600 has proven difficult
1020 (Friederici, 2011), however the P600 has been associated with a left inferior
1021 prefrontal-temporal cortical circuit (Brouwer et al., 2017; Brouwer & Hoeks, 2013) which
1022 also includes areas of the frontal inferior gyrus thought to mediate suppression of previous
1023 interpretations and possibly hints at the involvement of executive control (Hagoort, 2013;
1024 Kutas, 1993). Thus while we interpret our P600 constraint effect as a conflict signal, we do
1025 not believe our findings rule out that a shift in interpretation or suppression of previous
1026 representations occurs: we simply did not observe evidence that such a process is
1027 inevitably engaged by manipulating contextual strength, or that it is mappable to a single
1028 ERP phenomenon (for a discussion of the difficult “mapping problem” in behavioural
1029 neuroscience see Rösler, 2012). Indeed, if both processes involve the same cortical circuit at
1030 the same time, they may be difficult to disentangle without experimental methodologies
1031 better suited to spatial mapping such as MEG or fMRI.

1032 **Predictability and the PNP**

1033 In contrast with constraint, word predictability is a more reliable factor in eliciting
1034 the PNP, with low probability words triggering more positive amplitudes than high
1035 probability words (Brothers et al., 2017; Brothers et al., 2020; DeLong et al., 2014; DeLong
1036 et al., 2011; Federmeier et al., 2007; Hodapp & Rabovsky, 2021; Kuperberg et al., 2020;
1037 Ness & Meltzer-Asscher, 2018; Thornhill & Van Petten, 2012). It was therefore surprising
1038 that the current study did not find stronger evidence of a predictability effect in the
1039 anterior scalp region. However, we did see a left-lateralised effect (Figure 6). Among
1040 previous studies reporting an anterior PNP predictability effect, several observed this to be
1041 distributed across frontal and/or left lateral electrodes (DeLong et al., 2014; DeLong et al.,
1042 2011; Federmeier et al., 2007; Hodapp & Rabovsky, 2021; Kuperberg et al., 2020; Szewczyk
1043 & Schriefers, 2013). One possibility is that the left-lateralisation of the predictability effect

1044 is somehow related to the presence of a constraint manipulation; however, a left-lateralised
1045 effect appears to be evenly distributed across previous studies both with and without
1046 constraint manipulations. We thus refrain from interpreting the finding, but make note of
1047 it as being potentially in need of future characterisation.

1048 **Reflections on sample size and the sequential Bayes factor design**

1049 A major concern in ERP research is how to balance the labour and financial cost of
1050 EEG recordings with statistical power. The sequential Bayes factor design revealed that
1051 some research questions may be answerable with relatively small samples. For example,
1052 Figure 12 indicates that there was already strong evidence for the standard N400 high vs.
1053 low cloze probability effect at a sample size of 20 participants. However, here we urge
1054 caution: This was a large effect size that had a clear, directional, a priori hypothesis which
1055 we encoded into the statistical model using a truncated prior. A truncated prior will yield
1056 strong evidence more quickly for such a large effect, but a truncated prior must be carefully
1057 theoretically motivated a priori. Truncated priors will not be suitable for all types of
1058 research questions and should be interpreted with a higher threshold for evidence.
1059 However, designing informative priors for effects of interest based on previous data may be
1060 useful for keeping sample size within practical limits.

1061 Sample size must of course also be large enough to sufficiently account for the effects
1062 of interindividual variability and prevalence (i.e. some subjects may be “non responders”).
1063 ERP research is particularly sensitive to interindividual effects given the limitations of the
1064 EEG method (i.e. cortical and skull differences, low signal-to-noise ratio), and such effects
1065 are difficult to characterise in small samples (we thank an anonymous reviewer for this
1066 note). One approach to deciding whether a given sample is sufficiently large when it has
1067 been determined via a stopping rule with a narrow, truncated prior is to examine posterior
1068 estimates under a range of priors, both truncated and non-truncated, to see how well an
1069 estimated effect “holds up” under different assumptions (prior sensitivity analyses should
1070 be conducted regardless, but may be additionally useful for this question).

1071 Nonetheless, for our research question regarding constraint, we were able to provide
1072 strong evidence of an effect with considerably fewer participants than we had anticipated.
1073 For those effects that remained inconclusive at our final sample size, there were reasons we
1074 had not anticipated at the design stage of the study (e.g. a pandemic) and we were able to
1075 demonstrate using a design analysis that we would not have found strong evidence even
1076 with an infeasibly large sample. We were thus able to cut our losses and conserve
1077 resources. A sequential Bayes factor design may therefore be an efficient method of sample
1078 size determination for future EEG research.

1079

Conclusions

1080 In a relatively high-powered experimental design, we confirm previous research
1081 demonstrating a dissociated effect of contextual constraint on the ERP, in which the
1082 strength of a probabilistic representation affects processing in the post-N400 but not the
1083 N400 window. We also demonstrate a dissociated effect of word predictability on the ERP,
1084 in which there is a clear effect of predictability in the N400 but not the post-N400 window.
1085 Together these findings suggest that N400 amplitude is more sensitive to individual word
1086 predictability than context, whereas context is more important than predictability to the
1087 processes associated with the post-N400 window. We conclude that in the current study,
1088 the processing cost of stronger probabilistic expectations in the post-N400 window resulted
1089 from greater conflict between expectations and input, rather than from a greater shift in
1090 interpretation or suppression of previous representations. We base this conclusion on our
1091 observation of a posterior P600 rather than an anterior PNP. While a shift in
1092 interpretation or suppression could have occurred, these processes may not be the
1093 inevitable result of strong contextual constraint and may not be mappable to a unique
1094 ERP phenomenon. We propose either that eliciting a constraint effect in the anterior PNP
1095 requires a more complex experimental design than a straightforward strong/weak
1096 constraint comparison, or that the constraint-related PNP effect observed in previous
1097 studies could even be an artifact of low sample size.

Acknowledgements

1098

1099 We thank Johanna Thieke, Romy Leue, and Lisa Plagemann for their help with
1100 stimuli development, and for data collection under difficult circumstances during the
1101 Covid-19 pandemic. We also thank Trevor Brothers and Gina Kuperberg for sharing data,
1102 and Kara Federmeier for sharing stimuli and feedback on the results. Finally, we thank two
1103 anonymous reviewers and editor Kate Watkins for their help and feedback at all stages of
1104 the review process.

1105

Funding information

1106 Shravan Vasishth, Volkswagen Foundation grant 89953 and Deutsche
1107 Forschungsgemeinschaft (German Research Foundation) project number 317633480 – SFB
1108 1287, Project Q. Shravan Vasishth and Frank Rösler, Deutsche Forschungsgemeinschaft
1109 grant VA482/8–1. Article processing charges, Deutsche Forschungsgemeinschaft project
1110 number 491466077.

References

- 1111
1112 Baayen, R. H. (2001). *Word Frequency Distributions* (Vol. 18). Springer Science & Business
1113 Media.
- 1114 Baroni, M., Bernardini, S., Ferraresi, A., & Zanchetta, E. (2009). The WaCky wide web: A
1115 collection of very large linguistically processed web-crawled corpora. *Language*
1116 *Resources and Evaluation*, 43(3), 209–226.
1117 <https://doi.org/10.1007/s10579-009-9081-4>
- 1118 Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting Linear Mixed-Effects
1119 Models Using lme4. *Journal of Statistical Software*, 67(1), 1–48.
1120 <https://doi.org/10.18637/jss.v067.i01>
- 1121 Bennett, C. H. (1976). Efficient estimation of free energy differences from Monte Carlo
1122 data. *Journal of Computational Physics*, 22(2), 245–268.
1123 [https://doi.org/10.1016/0021-9991\(76\)90078-4](https://doi.org/10.1016/0021-9991(76)90078-4)
- 1124 Bornkessel-Schlesewsky, I., & Schlesewsky, M. (2008). An alternative perspective on
1125 “semantic P600” effects in language comprehension. *Brain Research Reviews*, 59(1),
1126 55–73. <https://doi.org/10.1016/J.BRAINRESREV.2008.05.003>
- 1127 Brothers, T., Swaab, T. Y., & Traxler, M. J. (2017). Goals and strategies influence lexical
1128 prediction during sentence comprehension. *Journal of Memory and Language*, 93.
1129 <https://doi.org/10.1016/j.jml.2016.10.002>
- 1130 Brothers, T., Wlotko, E. W., Warnke, L., & Kuperberg, G. (2020). Going the Extra Mile:
1131 Effects of Discourse Context on Two Late Positivities During Language
1132 Comprehension. *Neurobiology of Language*, 1(1), 135–160.
1133 https://doi.org/10.1162/nol_a_00006
- 1134 Brouwer, H., Crocker, M. W., Venhuizen, N. J., & Hoeks, J. C. J. (2017). A
1135 Neurocomputational Model of the N400 and the P600 in Language Processing.
1136 *Cognitive Science*, 41(S6), 1318–1352. <https://doi.org/10.1111/cogs.12461>

- 1137 Brouwer, H., & Hoeks, J. C. J. (2013). A time and place for language comprehension:
1138 Mapping the N400 and the P600 to a minimal cortical network. *Frontiers in Human*
1139 *Neuroscience*, 7. <https://doi.org/10.3389/fnhum.2013.00758>
- 1140 Buerkner, P.-C. (2018). brms: An R package for Bayesian multilevel models using Stan.
1141 *Journal of Statistical Software*, 80(1). <https://doi.org/10.18637/jss.v080.i01>
- 1142 Bürki-Foschini, A., Alario, F.-X., & Vasishth, S. (2022). EXPRESS: When words collide:
1143 Bayesian meta-analyses of distractor and target properties in the picture-word
1144 interference paradigm. *Quarterly Journal of Experimental Psychology*,
1145 17470218221114644. <https://doi.org/10.1177/17470218221114644>
- 1146 Chen, S. F., & Goodman, J. (1999). An empirical study of smoothing techniques for
1147 language modeling. *Computer Speech & Language*, 13(4), 359–394.
1148 <https://doi.org/10.1006/csla.1999.0128>
- 1149 Chung, Y., Gelman, A., Rabe-Hesketh, S., Liu, J., & Dorie, V. (2015). Weakly Informative
1150 Prior for Point Estimation of Covariance Matrices in Hierarchical Models. *Journal*
1151 *of Educational and Behavioral Statistics*, 40(2), 136–157.
1152 <https://doi.org/10.3102/1076998615570945>
- 1153 Cohen, J. (1983). The cost of dichotomization. *The cost of dichotomization*, 7(3), 249–253.
- 1154 Coulson, S., King, J. W., & Kutas, M. (1998). Expect the Unexpected: Event-related Brain
1155 Response to Morphosyntactic Violations. *Language and Cognitive Processes*, 13(1),
1156 21–58. <https://doi.org/10.1080/016909698386582>
- 1157 Delaney-Busch, N., Morgan, E., Lau, E., & Kuperberg, G. (2019). Neural evidence for
1158 Bayesian trial-by-trial adaptation on the N400 during semantic priming. *Cognition*,
1159 187, 10–20. <https://doi.org/10.1016/j.cognition.2019.01.001>
- 1160 DeLong, K. A., Quante, L., & Kutas, M. (2014). Predictability, plausibility, and two late
1161 ERP positivities during written sentence comprehension. *Neuropsychologia*, 61(1).
1162 <https://doi.org/10.1016/j.neuropsychologia.2014.06.016>

- 1163 DeLong, K. A., Urbach, T. P., Groppe, D. M., & Kutas, M. (2011). Overlapping dual ERP
1164 responses to low cloze probability sentence continuations. *Psychophysiology*, *48*(9),
1165 1203–1207. <https://doi.org/10.1111/j.1469-8986.2011.01199.x>
- 1166 Drummond, A. (2016). *Ibex: Software for psycholinguistic experiments*.
1167 <https://github.com/addrummond/ibex>
- 1168 Federmeier, K. D. (2022). Connecting and considering: Electrophysiology provides insights
1169 into comprehension. *Psychophysiology*, *59*(1), e13940.
1170 <https://doi.org/10.1111/psyp.13940>
1171 _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/psyp.13940>
- 1172 Federmeier, K. D., & Kutas, M. (1999). A rose by any other name: Long-term memory
1173 structure and sentence processing. *Journal of Memory and Language*, *41*(4),
1174 469–495. <https://doi.org/10.1006/jmla.1999.2660>
- 1175 Federmeier, K. D., Wlotko, E. W., Ochoa-Dewald, E. D., & Kutas, M. (2007). Multiple
1176 effects of sentential constraint on word processing. *Brain Research*, (1146), 75–84.
1177 <https://doi.org/10.1016/j.brainres.2006.06.101>
- 1178 Fitz, H., & Chang, F. (2019). Language ERPs reflect learning through prediction error
1179 propagation. *Cognitive Psychology*, *111*, 15–52.
1180 <https://doi.org/10.1016/j.cogpsych.2019.03.002>
- 1181 Frank, S. L., Otten, L. J., Galli, G., & Vigliocco, G. (2015). The ERP response to the
1182 amount of information conveyed by words in sentences. *Brain and Language*, *140*,
1183 1–11. <https://doi.org/10.1016/j.bandl.2014.10.006>
- 1184 Friederici, A. D. (2011). The brain basis of language processing: From structure to
1185 function. *Physiological Reviews*, *91*(4), 1357–1392.
1186 <https://doi.org/doi:10.1152/physrev.00006.2011>
- 1187 Friederici, A. D. (2012). The cortical language circuit: From auditory perception to
1188 sentence comprehension. *Trends in Cognitive Sciences*, *16*(5), 262–268.
1189 <https://doi.org/10.1016/j.tics.2012.04.001>

- 1190 Friederici, A. D., Hahne, A., & Saddy, D. (2002). Distinct Neurophysiological Patterns
1191 Reflecting Aspects of Syntactic Complexity and Syntactic Repair. *Journal of*
1192 *Psycholinguistic Research*, 31(1), 45–63. <https://doi.org/10.1023/A:1014376204525>
- 1193 Gabry, J., Simpson, D., Vehtari, A., Betancourt, M., & Gelman, A. (2019). Visualization in
1194 Bayesian workflow. *Journal of the Royal Statistical Society Series A (Statistics in*
1195 *Society)*, 182(2). <http://arxiv.org/abs/1709.01449>
- 1196 Garnsey, S. M. (1993). *Event-related brain potentials in the study of language: An*
1197 *introduction: Language and Cognitive Processes: Vol 8, No 4*. Retrieved October 21,
1198 2019, from <https://www.tandfonline.com/doi/abs/10.1080/01690969308407581>
- 1199 Gelman, A., & Carlin, J. (2014). Beyond Power Calculations: Assessing Type S (Sign) and
1200 Type M (Magnitude) Errors. *Perspectives on Psychological Science*, 9(6), 641–651.
1201 <https://doi.org/10.1177/1745691614551642>
- 1202 Gelman, A., Jakulin, A., Pittau, M. G., & Su, Y.-S. (2008). A weakly informative default
1203 prior distribution for logistic and other regression models. *Annals of Applied*
1204 *Statistics*, 2(4), 1360–1383. <https://doi.org/10.1214/08-AOAS191>
- 1205 Gelman, A., Simpson, D., & Betancourt, M. (2017). The Prior Can Often Only Be
1206 Understood in the Context of the Likelihood. *Entropy*, 19(10), 555.
1207 <https://doi.org/10.3390/e19100555>
- 1208 Gronau, Q. F., Sarafoglou, A., Matzke, D., Ly, A., Boehm, U., Marsman, M., Leslie, D. S.,
1209 Forster, J. J., Wagenmakers, E.-J., & Steingroever, H. (2017). A tutorial on bridge
1210 sampling. *Journal of Mathematical Psychology*, 81, 80–97.
1211 <https://doi.org/10.1016/j.jmp.2017.09.005>
- 1212 Gunter, T. C., Friederici, A. D., & Schriefers, H. (2000). Syntactic Gender and Semantic
1213 Expectancy: ERPs Reveal Early Autonomy and Late Interaction. *Journal of*
1214 *Cognitive Neuroscience*, 12(4), 556–568. <https://doi.org/10.1162/089892900562336>

- 1215 Günther, F. (2022). *Homepage of Fritz Günther - Semantic Spaces*. Retrieved June 26,
1216 2022, from
1217 https://sites.google.com/site/fritzgntr/software-resources/semantic_spaces
- 1218 Günther, F., Dudschig, C., & Kaup, B. (2015). LSAfun - An R package for computations
1219 based on Latent Semantic Analysis. *Behavior Research Methods*, *47*(4), 930–944.
1220 <https://doi.org/10.3758/s13428-014-0529-0>
- 1221 Hagoort, P. (2013). MUC (Memory, Unification, Control) and beyond. *Frontiers in*
1222 *Psychology*, *4*. <https://doi.org/10.3389/fpsyg.2013.00416>
- 1223 Hagoort, P., Baggio, G., & Willems, R. M. (2009). Semantic unification. In M. S. Gazzaniga
1224 (Ed.), *The Cognitive Neurosciences* (4th ed., pp. 819–836). MIT Press. https://pure.mpg.de/pubman/faces/ViewItemOverviewPage.jsp?itemId=item_64579
- 1225 https://pure.mpg.de/pubman/faces/ViewItemOverviewPage.jsp?itemId=item_64579
- 1226 Hagoort, P., Brown, C., & Groothusen, J. (1993). The syntactic positive shift (sps) as an
1227 erp measure of syntactic processing. *Language and Cognitive Processes*, *8*(4),
1228 439–483. <https://doi.org/10.1080/01690969308407585>
- 1229 Hagoort, P., & Indefrey, P. (2014). The Neurobiology of Language Beyond Single Words.
1230 *Annual Review of Neuroscience*, *37*(1), 347–362.
1231 <https://doi.org/10.1146/annurev-neuro-071013-013847>
- 1232 Hale, J. (2001). A probabilistic Earley parser as a psycholinguistic model. *NAACL '01:*
1233 *Second meeting of the North American Chapter of the Association for*
1234 *Computational Linguistics on Language technologies 2001*, 1–8.
1235 <https://doi.org/10.3115/1073336.1073357>
- 1236 Hickok, G., & Poeppel, D. (2007). The cortical organization of speech processing. *Nature*
1237 *Reviews Neuroscience*, *8*(5), 393–402. <https://doi.org/10.1038/nrn2113>
- 1238 Hodapp, A., & Rabovsky, M. (2021). The N400 ERP component reflects an error-based
1239 implicit learning signal during language comprehension. *European Journal of*
1240 *Neuroscience*, *54*(9), 7125–7140. <https://doi.org/10.1111/ejn.15462>

- 1241 Hoeks, J. C. J., Stowe, L. A., & Doedens, G. (2004). Seeing words in context: The
1242 interaction of lexical and sentence level information during reading. *Cognitive Brain*
1243 *Research*, 19(1), 59–73. <https://doi.org/10.1016/j.cogbrainres.2003.10.022>
- 1244 Hyvärinen, A., Karhunen, J., & Oja, E. (2001). Adaptive and learning systems for signal
1245 processing, communications, and control. In *Independent component analysis*
1246 (pp. 11–14). John Wiley & Sons, Inc. <https://doi.org/10.1002/0471221317>
- 1247 Jäger, L. A., Engelmann, F., & Vasishth, S. (2017). Similarity-based interference in
1248 sentence comprehension: Literature review and Bayesian meta-analysis. *Journal of*
1249 *Memory and Language*, 94, 316–339. <https://doi.org/10.1016/j.jml.2017.01.004>
- 1250 Jasper, H. (1958). The 10/20 international electrode system. *EEG and Clinical*
1251 *Neurophysiology*, 10(2), 370–375.
- 1252 Jeffreys, H. (1939). *Theory of Probability*. Oxford University Press.
- 1253 Jung, T.-p., Makeig, S., Westerfield, M., Townsend, J., Courchesne, E., & Sejnowski, T. J.
1254 (2001). Analyzing and Visualizing Single-Trial Event-Related Potentials. *Human*
1255 *Brain Mapping*, 185, 166–185. <https://doi.org/10.1002/hbm.1050>
- 1256 Kay, M. (2022). *{Tidybayes}: Tidy Data and Geoms for {Bayesian} Models*.
1257 <http://mjskay.github.io/tidybayes/>
1258 10.5281/zenodo.1308151
- 1259 Kello, C. T., Brown, G. D. A., Ferrer-i-Cancho, R., Holden, J. G., Linkenkaer-Hansen, K.,
1260 Rhodes, T., & Van Orden, G. C. (2010). Scaling laws in cognitive sciences. *Trends*
1261 *in Cognitive Sciences*, 14(5), 223–232. <https://doi.org/10.1016/j.tics.2010.02.005>
- 1262 Kim, A., & Osterhout, L. (2005). The independence of combinatory semantic processing:
1263 Evidence from event-related potentials. *Journal of Memory and Language*, 52(2),
1264 205–225. <https://doi.org/10.1016/j.jml.2004.10.002>
- 1265 Kuperberg, G., & Brothers, T. (2019). *P600 Task Discussion*. Retrieved May 22, 2022,
1266 from [https://projects.iq.harvard.edu/kuperberglab/additional-materials/P600-task-](https://projects.iq.harvard.edu/kuperberglab/additional-materials/P600-task-discussion)
1267 [discussion](https://projects.iq.harvard.edu/kuperberglab/additional-materials/P600-task-discussion)

- 1268 Kuperberg, G., Brothers, T., & Wlotko, E. W. (2020). A Tale of Two Positivities (and the
1269 N400): Distinct neural signatures are evoked by confirmed and violated predictions
1270 at different levels of representation. *Journal of Cognitive Neuroscience*, *32*(1),
1271 12–35. https://doi.org/10.1162/jocn_a_01465
- 1272 Kuperberg, G., & Jaeger, T. F. (2016). What do we mean by prediction in language
1273 comprehension? *Language Cognition & Neuroscience*, *31*(1).
1274 <https://doi.org/10.1080/23273798.2015.1102299>
- 1275 Kuperberg, G., Kreher, D. A., Sitnikova, T., Caplan, D. N., & Holcomb, P. J. (2007). The
1276 role of animacy and thematic relationships in processing active English sentences:
1277 Evidence from event-related potentials. *Brain and Language*, *100*(3), 223–237.
1278 <https://doi.org/10.1016/j.bandl.2005.12.006>
- 1279 Kuperberg, G., Sitnikova, T., Caplan, D., & Holcomb, P. J. (2003). Electrophysiological
1280 distinctions in processing conceptual relationships within simple sentences. *Cognitive*
1281 *Brain Research*, *17*(1), 117–129. [https://doi.org/10.1016/S0926-6410\(03\)00086-7](https://doi.org/10.1016/S0926-6410(03)00086-7)
- 1282 Kutas, M. (1993). In the company of other words: Electrophysiological evidence for
1283 single-word and sentence context effects. *Language and Cognitive Processes*, *8*(4),
1284 533–572. <https://doi.org/10.1080/01690969308407587>
- 1285 Kutas, M., & Federmeier, K. D. (2011). Thirty Years and Counting: Finding Meaning in the
1286 N400 Component of the Event-Related Brain Potential (ERP). *Annual Review of*
1287 *Psychology*, *62*(1), 621–647. <https://doi.org/10.1146/annurev.psych.093008.131123>
- 1288 Kutas, M., & Hillyard, S. A. (1980). Reading between the lines: Event-related brain
1289 potentials during natural sentence processing. *Brain and Language*, *11*(2), 354–373.
1290 [https://doi.org/10.1016/0093-934X\(80\)90133-9](https://doi.org/10.1016/0093-934X(80)90133-9)
- 1291 Kutas, M., & Hillyard, S. A. (1984). Brain potentials during reading reflect word
1292 expectancy and semantic association. *Nature*, *307*(5947), 161–163.
1293 <https://doi.org/10.1038/307161a0>

- 1294 Lai, M. K., Rommers, J., & Federmeier, K. D. (2021). The fate of the unexpected:
1295 Consequences of misprediction assessed using ERP repetition effects. *Brain*
1296 *Research*, 147290. <https://doi.org/10.1016/j.brainres.2021.147290>
- 1297 Lappin, J. S., & Eriksen, C. W. (1966). Use of a delayed signal to stop a visual
1298 reaction-time response. *Journal of Experimental Psychology*, 72(6), 805–811.
1299 <https://doi.org/10.1037/h0021266>
- 1300 Laszlo, S., & Federmeier, K. D. (2009). A Beautiful Day in the Neighborhood: An
1301 Event-Related Potential Study of Lexical Relationships and Prediction in Context.
1302 *Journal of memory and language*, 61(3), 326–338.
1303 <https://doi.org/10.1016/j.jml.2009.06.004>
- 1304 Lau, E., Phillips, C., & Poeppel, D. (2008). A cortical network for semantics:
1305 (de)constructing the N400. *Nature Reviews Neuroscience*, 9(12), 920–933.
1306 <https://doi.org/10.1038/nrn2532>
- 1307 Lee, M., & Wagenmakers, E.-J. (2014). *Bayesian Cognitive Modeling: A Practical Course*.
1308 Cambridge University Press. <https://doi.org/10.1017/CBO9781139087759>
- 1309 Levy, R. (2008). Expectation-based syntactic comprehension. *Cognition*, 106(3),
1310 1126–1177. <https://doi.org/10.1016/j.cognition.2007.05.006>
- 1311 Lindeløv, J. K. (2021). *Job* (Version 0.3.0). <https://lindeloev.github.io/job>
- 1312 Lindstone, G. J. (1920). Note on the General Case of the Bayes-Laplace Formula for
1313 Inductive or a Posteriori Probabilities. *Transactions of the Faculty of Actuaries*, (8),
1314 182–192.
- 1315 Logan, G. D., & Cowan, W. B. (1984). On the ability to inhibit thought and action: A
1316 theory of an act of control. *Psychological Review*, 91(3), 295–327.
1317 <https://doi.org/10.1037/0033-295X.91.3.295>
- 1318 Luck, S. J. (2005a). *The Event-Related Potential Technique in Cognitive Neuroscience*.
1319 MIT Press.

- 1320 Luck, S. J. (2005b). Ten Simple Rules for Designing and Interpreting ERP Experiments. In
1321 T. C. Handy (Ed.), *Event-related Potentials: A Methods Handbook* (pp. 17–32). MIT
1322 press.
- 1323 Luck, S. J., & Gaspelin, N. (2016). How to Get Statistically Significant Effects in Any ERP
1324 Experiment (and Why You Shouldn't). *Psychophysiology*, *44*(24).
1325 <https://doi.org/10.1111/psyp.12639>
- 1326 Mathôt, S., Schreij, D., & Theeuwes, J. (2012). OpenSesame: An open-source, graphical
1327 experiment builder for the social sciences. *Behavior Research Methods*, *44*(2),
1328 314–324. <https://doi.org/10.3758/s13428-011-0168-7>
- 1329 Meerendonk, N. V. D., Kolk, H. H. J., Chwilla, D. J., & Vissers, C. T. W. M. (2009).
1330 Monitoring in Language Perception. *Language and Linguistics Compass*, *3*(5),
1331 1211–1224. <https://doi.org/10.1111/j.1749-818X.2009.00163.x>
- 1332 Meng, X.-L., & Wong, W. H. (1996). Simulating ratios of normalizing constants via a
1333 simple identity: A theoretical exploration. *Statistica Sinica*, *6*(4), 831–860.
1334 <https://doi.org/https://www.jstor.org/stable/24306045>
- 1335 Metzner, P., von der Malsburg, T., Vasishth, S., & Rösler, F. (2017). The Importance of
1336 Reading Naturally: Evidence From Combined Recordings of Eye Movements and
1337 Electric Brain Potentials. *Cognitive Science*, *41*. <https://doi.org/10.1111/cogs.12384>
- 1338 Ness, T., & Meltzer-Asscher, A. (2018). Lexical inhibition due to failed prediction:
1339 Behavioral evidence and ERP correlates. *Journal of Experimental Psychology:*
1340 *Learning, Memory, and Cognition*, *44*(8), 1269–1285.
1341 <https://doi.org/10.1037/xlm0000525>
- 1342 Nicenboim, B. (2018). *Eeguana: A package for manipulating EEG data in R*.
1343 <https://github.com/bnicenboim/eeguana>
- 1344 Nicenboim, B., Vasishth, S., & Rösler, F. (2020). Are words pre-activated probabilistically
1345 during sentence comprehension? Evidence from new data and a Bayesian

- 1346 random-effects meta-analysis using publicly available data. *Neuropsychologia*,
1347 107427. <https://doi.org/10.1016/j.neuropsychologia.2020.107427>
- 1348 Nieuwland, M. S., Politzer-Ahles, S., Heyselaar, E., Segaert, K., Darley, E., Kazanina, N.,
1349 Von Grebmer Zu Wolfsthurn, S., Bartolozzi, F., Kogan, V., Ito, A., Mézière, D.,
1350 Barr, D. J., Rousselet, G. A., Ferguson, H. J., Busch-Moreno, S., Fu, X.,
1351 Tuomainen, J., Kulakova, E., Husband, E. M., ... Huettig, F. (2018). Large-scale
1352 replication study reveals a limit on probabilistic prediction in language
1353 comprehension (B. G. Shinn-Cunningham, Ed.). *eLife*, 7, e33468.
1354 <https://doi.org/10.7554/eLife.33468>
- 1355 Osterhout, L. (1999). A Superficial Resemblance Does Not Necessarily Mean You Are Part
1356 of the Family: Counterarguments to Coulson, King and Kutas (1998) in the
1357 P600/SPS-P300 Debate. *Language and Cognitive Processes*, 14(1), 1–14.
1358 <https://doi.org/10.1080/016909699386356>
- 1359 Osterhout, L., & Holcomb, P. J. (1992). Event-related brain potentials elicited by syntactic
1360 anomaly. *Journal of Memory and Language*, 31(6), 785–806.
1361 [https://doi.org/10.1016/0749-596X\(92\)90039-Z](https://doi.org/10.1016/0749-596X(92)90039-Z)
- 1362 Osterhout, L., McKinnon, R., Bersick, M., & Corey, V. (1996). On the Language
1363 Specificity of the Brain Response to Syntactic Anomalies: Is the Syntactic Positive
1364 Shift a Member of the P300 Family? *Journal of Cognitive Neuroscience*, 8(6),
1365 507–526. <https://doi.org/10.1162/jocn.1996.8.6.507>
- 1366 Pedersen, T. L. (2022). *Patchwork: The Composer of Plots*.
1367 <https://patchwork.data-imaginist.com>
- 1368 R Core Team. (2020). *R: A Language and Environment for Statistical Computing*. Vienna,
1369 Austria. <https://www.R-project.org>
- 1370 Rabovsky, M., Hansen, S. S., & McClelland, J. L. (2018). Modelling the N400 brain
1371 potential as change in a probabilistic representation of meaning. *Nature Human*
1372 *Behaviour*, 2(9), 693. <https://doi.org/10.1038/s41562-018-0406-4>

- 1373 Rösler, F. (2012). Some unsettled problems in behavioral neuroscience research.
1374 *Psychological Research*, 76(2), 131–144. <https://doi.org/10.1007/s00426-011-0408-6>
- 1375 Sassenhagen, J., & Bornkessel-Schlesewsky, I. (2015). The P600 as a correlate of ventral
1376 attention network reorientation. *Cortex*, 66, A3–A20.
1377 <https://doi.org/10.1016/j.cortex.2014.12.019>
- 1378 Sassenhagen, J., & Fiebach, C. J. (2019). Finding the P3 in the P600: Decoding shared
1379 neural mechanisms of responses to syntactic violations and oddball targets.
1380 *NeuroImage*, 200, 425–436. <https://doi.org/10.1016/j.neuroimage.2019.06.048>
- 1381 Sassenhagen, J., Schlewsky, M., & Bornkessel-Schlesewsky, I. (2014). The P600-as-P3
1382 hypothesis revisited: Single-trial analyses reveal that the late EEG positivity
1383 following linguistically deviant material is reaction time aligned. *Brain and*
1384 *Language*, 137, 29–39. <https://doi.org/10.1016/j.bandl.2014.07.010>
- 1385 Schad, D. J., Betancourt, M., & Vasishth, S. (2020). Toward a principled Bayesian
1386 workflow: A tutorial for cognitive science. *Psychological Methods*.
1387 <https://doi.org/10.1037/met0000275>
- 1388 Smith, N. J., & Levy, R. (2013). The effect of word predictability on reading time is
1389 logarithmic. *Cognition*, 128(3). <https://doi.org/10.1016/j.cognition.2013.02.013>
- 1390 Stan Development Team. (2018). *The Stan Core Library* (Version Version 2.18.0).
1391 <http://mc-stan.org>
- 1392 Stan Development Team. (2020). *RStan: The R interface to Stan* (Version 2.21.2).
1393 <http://mc-stan.org/>
- 1394 Szewczyk, J. M., & Schriefers, H. (2013). Prediction in language comprehension beyond
1395 specific words: An ERP study on sentence comprehension in Polish. *Journal of*
1396 *Memory and Language*, 68(4), 297–314. <https://doi.org/10.1016/j.jml.2012.12.002>
- 1397 Tanner, D., Morgan-Short, K., & Luck, S. J. (2015). How inappropriate high-pass filters can
1398 produce artifactual effects and incorrect conclusions in ERP studies of language and
1399 cognition. *Psychophysiology*, 52(8), 997–1009. <https://doi.org/10.1111/psyp.12437>

- 1400 Tanner, D., Norton, J. J. S., Morgan-Short, K., & Luck, S. J. (2016). On high-pass filter
1401 artifacts (they're real) and baseline correction (it's a good idea) in ERP/ERMF
1402 analysis. *Journal of Neuroscience Methods*, *266*, 166–170.
1403 <https://doi.org/10.1016/j.jneumeth.2016.01.002>
- 1404 Thornhill, D. E., & Van Petten, C. (2012). Lexical versus conceptual anticipation during
1405 sentence processing: Frontal positivity and N400 ERP components. *International*
1406 *Journal of Psychophysiology*, *83*(3), 382–392.
1407 <https://doi.org/10.1016/j.ijpsycho.2011.12.007>
- 1408 van de Meerendonk, N., Indefrey, P., Chwilla, D. J., & Kolk, H. H. J. (2011). Monitoring in
1409 language perception: Electrophysiological and hemodynamic responses to spelling
1410 violations. *NeuroImage*, *54*(3), 2350–2363.
1411 <https://doi.org/10.1016/j.neuroimage.2010.10.022>
- 1412 Van Petten, C., & Luka, B. J. (2012). Prediction during language comprehension: Benefits,
1413 costs, and ERP components. *International Journal of Psychophysiology*, *83*(2),
1414 176–190. <https://doi.org/http://dx.doi.org/10.1016/j.ijpsycho.2011.09.015>
- 1415 Vanrullen, R. (2011). Four Common Conceptual Fallacies in Mapping the Time Course of
1416 Recognition. *Frontiers in Psychology*, *2*. Retrieved June 24, 2022, from
1417 <https://www.frontiersin.org/article/10.3389/fpsyg.2011.00365>
- 1418 Vasishth, S., & Engelmann, F. (2022). *Sentence Comprehension as a Cognitive Process:*
1419 *{A} Computational Approach*. Cambridge University Press.
1420 <https://doi.org/10.1017/9781316459560.007>
- 1421 Verbruggen, F., Aron, A. R., Band, G. P., Beste, C., Bissett, P. G., Brockett, A. T.,
1422 Brown, J. W., Chamberlain, S. R., Chambers, C. D., Colonius, H., Colzato, L. S.,
1423 Corneil, B. D., Coxon, J. P., Dupuis, A., Eagle, D. M., Garavan, H., Greenhouse, I.,
1424 Heathcote, A., Huster, R. J., ... Boehler, C. N. (2019). A consensus guide to
1425 capturing the ability to inhibit actions and impulsive behaviors in the stop-signal

- 1426 task (M. J. Frank, D. Badre, T. Egner, & D. Swick, Eds.). *eLife*, 8, e46323.
1427 <https://doi.org/10.7554/eLife.46323>
- 1428 Verbruggen, F., Logan, G. D., & Stevens, M. A. (2008). STOP-IT: Windows executable
1429 software for the stop-signal paradigm. *Behavior Research Methods*, 40(2), 479–483.
1430 <https://doi.org/10.3758/BRM.40.2.479>
- 1431 Vissers, C. T. W. M., Chwilla, D. J., & Kolk, H. H. J. (2006). Monitoring in language
1432 perception: The effect of misspellings of words in highly constrained sentences.
1433 *Brain Research*, 1106(1), 150–163. <https://doi.org/10.1016/j.brainres.2006.05.012>
- 1434 Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L. D., François, R.,
1435 Grolemund, G., Hayes, A., Henry, L., Hester, J., Kuhn, M., Pedersen, T. L.,
1436 Miller, E., Bache, S. M., Müller, K., Ooms, J., Robinson, D., Seidel, D. P.,
1437 Spinu, V., ... Yutani, H. (2019). Welcome to the Tidyverse. *Journal of Open Source*
1438 *Software*, 4(43), 1686. <https://doi.org/10.21105/joss.01686>
- 1439 Wlotko, E. W., & Federmeier, K. D. (2007). Finding the right word: Hemispheric
1440 asymmetries in the use of sentence context information. *Neuropsychologia*, 45(13),
1441 3001–3014. <https://doi.org/10.1016/j.neuropsychologia.2007.05.013>