

## FOR PEER REVIEW - CONFIDENTIAL

**Large-scale replication study reveals a limit on probabilistic prediction in language comprehension**

Tracking no: 09-11-2017-RA-eLife-33468R1

Mante Nieuwland (MPI for Psycholinguistics), Stephen Politzer-Ahles (The Hong Kong Polytechnic University), Evelien Heyselaar (University of Birmingham), Katrien Segaert (University of Birmingham), Emily Darley (University of Bristol), Nina Kazanina (University of Bristol), Sarah Von Grebmer Zu Wolfsturn (University of Bristol), Federica Bartolozzi (University of Edinburgh), Vita Kogan (University of Edinburgh), Aine Ito (University of Edinburgh), Diane Mézière (University of Edinburgh), Dale Barr (University of Glasgow), Guillaume Rousselet (University of Glasgow), Heather Ferguson (University of Kent), Simon Busch-Moreno (University College London), Xiao Fu (University College London), Jyrki Tuomainen (University College London), Eugenia Kulakova (University College London), E. Husband (University of Oxford), David Donaldson (University of Stirling), Zdenko Kohút (University of York), Shirley-Ann Rueschemeyer (University of York), and Falk Huettig (Max Planck Institute for Psycholinguistics)

**Abstract:**

Do people routinely pre-activate the meaning and even the phonological form of upcoming words? The most acclaimed evidence for phonological prediction comes from a 2005 *Nature Neuroscience* publication by DeLong, Urbach and Kutas, who observed a graded modulation of electrical brain potentials (N400) to nouns and preceding articles by the probability that people use a word to continue the sentence fragment ('cloze'). In our direct replication study spanning 9 laboratories ( $N=334$ ), pre-registered replication-analyses and exploratory Bayes Factor analyses successfully replicated the noun-results but, crucially, not the article-results. Pre-registered single-trial analyses also yielded a statistically significant effect for the nouns but not the articles. Exploratory Bayesian single-trial analyses showed that the article-effect may be non-zero but is likely far smaller than originally reported and too small to observe without very large sample sizes. Our results do not support the view that readers routinely pre-activate the phonological form of predictable words.

**Impact statement:** Large-scale replication study with brain potentials challenges the view that people routinely predict the phonological form of a predictable word during language comprehension

**Competing interests:** No competing interests declared

**Author contributions:**

Mante Nieuwland: Conceptualization; Data curation; Software; Formal analysis; Supervision; Validation; Investigation; Visualization; Methodology; Writing—original draft; Project administration; Writing—review and editing Stephen Politzer-Ahles: Software; Formal analysis; Validation; Investigation; Visualization; Methodology; Writing—original draft; Project administration; Writing—review and editing Evelien Heyselaar: Investigation; Writing—review and editing Katrien Segaert: Resources; Supervision; Writing—review and editing Emily Darley: Investigation; Nina Kazanina: Resources; Supervision; Writing—original draft; Writing—review and editing Sarah Von Grebmer Zu Wolfsturn: Investigation; Federica Bartolozzi: Investigation; Vita Kogan: Investigation; Aine Ito: Investigation; Writing—review and editing Diane Mézière: Investigation; Dale Barr: Software; Formal analysis; Writing—review and editing Guillaume Rousselet: Resources; Formal analysis; Supervision; Methodology; Writing—review and editing Heather Ferguson: Resources; Supervision; Funding acquisition; Writing—review and editing Simon Busch-Moreno: Software; Formal analysis; Investigation; Xiao Fu: Investigation; Jyrki Tuomainen: Resources; Supervision; Eugenia Kulakova: Software; Formal analysis; Investigation; E. Husband: Resources; Supervision; Writing—review and editing David Donaldson: Resources; Supervision; Writing—review and editing Zdenko Kohút: Investigation; Shirley-Ann Rueschemeyer: Supervision; Writing—review and editing Falk Huettig: Conceptualization; Funding acquisition; Writing—review and editing

**Funding:**

European Research Council: Heather J. Ferguson, ERC Starting grant 636458 The funders had no role in study design, data collection and interpretation, or the decision to submit the work for publication.

**Datasets:**

Datasets Generated: Replication Recipe Analysis plan: Mante Nieuwland, 2018, <https://osf.io/eyzaq/>, Available at the Open Science Framework  
Reporting Standards: N/A

**Ethics:**

Human Subjects: Yes Ethics Statement: All participants were informed about the procedure of the experiment and then gave informed consent to use the data for research and dissemination/publication purpose. Ethical approval for EEG experimentation was obtained at each involved institution, according to custom guidelines of the ethics committee at each institution. Clinical Trial: No Animal Subjects: No

**Author Affiliation:**

Mante Nieuwland(Neurobiology of Language,MPI for Psycholinguistics,Netherlands) Stephen Politzer-Ahles(Department of Chinese and Bilingual Studies,The Hong Kong Polytechnic University,Hong Kong) Evelien Heyselaar(School of Psychology,University of Birmingham,United Kingdom) Katrien Segaert(School of Psychology,University of Birmingham,United Kingdom) Emily Darley(School of Experimental Psychology,University of Bristol,United Kingdom) Nina Kazanina(,University of Bristol,United Kingdom) Sarah Von Grebmer Zu Wolfsthurn(School of Experimental Psychology,University of Bristol,United Kingdom) Federica Bartolozzi(School of Philosophy, Psychology and Language Sciences,University of Edinburgh,United Kingdom) Vita Kogan(School of Philosophy, Psychology and Language Sciences,University of Edinburgh,United Kingdom) Aine Ito(School of Philosophy, Psychology and Language Sciences,University of Edinburgh,United Kingdom) Diane Mézière(School of Philosophy, Psychology and Language Sciences,University of Edinburgh,United Kingdom) Dale Barr(Institute of Neuroscience and Psychology,University of Glasgow,United Kingdom) Guillaume Rousselet(Institute of Neuroscience and Psychology,University of Glasgow,United Kingdom) Heather Ferguson(School of Psychology,,University of Kent,United Kingdom) Simon Busch-Moreno(Division of Psychology and Language Sciences,University College London,United Kingdom) Xiao Fu(Division of Psychology and Language Sciences,University College London,United Kingdom) Jyri Tuomainen(Division of Psychology and Language Sciences,University College London,United Kingdom) Eugenia Kulakova(Institute of Cognitive Neuroscience,University College London,United Kingdom) E. Husband(Faculty of Linguistics, Philology & Phonetics,University of Oxford,United Kingdom) David Donaldson(Psychology, Faculty of Natural Sciences,University of Stirling,United Kingdom) Zdenko Kohút(Department of Psychology,University of York,United Kingdom) Shirley-Ann Rueschemeyer(Department of Psychology,University of York,United Kingdom) Falk Huettig(Psychology of Language,Max Planck Institute for Psycholinguistics,Netherlands)

**Dual-use research:** No

**Permissions:** Have you reproduced or modified any part of an article that has been previously published or submitted to another journal? No

1                   **Large-scale replication study reveals a limit on probabilistic prediction in**  
2                   **language comprehension**

3

4       Mante S. Nieuwland<sup>1,5</sup>, Stephen Politzer-Ahles<sup>2,10</sup>, Evelien Heyselaar<sup>3</sup>, Katrien Segaert<sup>3</sup>,  
5       Emily Darley<sup>4</sup>, Nina Kazanina<sup>4</sup>, Sarah Von Grebmer Zu Wolfsturn<sup>4</sup>, Federica Bartolozzi<sup>5</sup>,  
6       Vita Kogan<sup>5</sup>, Aine Ito<sup>5,10</sup>, Diane Mézière<sup>5</sup>, Dale J. Barr<sup>6</sup>, Guillaume Rousselet<sup>6</sup>, Heather J.  
7       Ferguson<sup>7</sup>, Simon Busch-Moreno<sup>8</sup>, Xiao Fu<sup>8</sup>, Jyrki Tuomainen<sup>8</sup>, Eugenia Kulakova<sup>9</sup>, E.  
8       Matthew Husband<sup>10</sup>, David I. Donaldson<sup>11</sup>, Zdenko Kohút<sup>12</sup>, Shirley-Ann Rueschemeyer<sup>12</sup>,  
9                   Falk Huettig<sup>1</sup>

10

11                   <sup>1</sup> Max Planck Institute for Psycholinguistics, Nijmegen, The Netherlands

12                   <sup>2</sup> Department of Chinese and Bilingual Studies, The Hong Kong Polytechnic University, Kowloon, Hong Kong

13                   <sup>3</sup> School of Psychology, University of Birmingham, Birmingham, United Kingdom

14                   <sup>4</sup> School of Experimental Psychology, University of Bristol, Bristol, United Kingdom

15                   <sup>5</sup> School of Philosophy, Psychology and Language Sciences, University of Edinburgh, Edinburgh, United Kingdom

16                   <sup>6</sup> Institute of Neuroscience and Psychology, University of Glasgow, Glasgow, United Kingdom

17                   <sup>7</sup> School of Psychology, University of Kent, Canterbury, United Kingdom

18                   <sup>8</sup> Division of Psychology and Language Sciences, University College London, London, United Kingdom

19                   <sup>9</sup> Institute of Cognitive Neuroscience, University College London, London, United Kingdom

20                   <sup>10</sup> Faculty of Linguistics, Philology & Phonetics; University of Oxford, Oxford, United Kingdom

21                   <sup>11</sup> Psychology, Faculty of Natural Sciences, University of Stirling, Stirling, United Kingdom

22                   <sup>12</sup> Department of Psychology, University of York, York, United Kingdom

23

24                   **Corresponding author:**

25       Mante S. Nieuwland, Max Planck Institute for Psycholinguistics, Wundtlaan 1, 6525 XD  
26       Nijmegen, The Netherlands. E-mail: [mante.nieuwland@mpi.nl](mailto:mante.nieuwland@mpi.nl), phone: +31-24-3521911

27     **ABSTRACT**

28         Do people routinely pre-activate the meaning and even the phonological form of  
29         upcoming words? The most acclaimed evidence for phonological prediction comes from a  
30         2005 *Nature Neuroscience* publication by DeLong, Urbach and Kutas, who observed a  
31         graded modulation of electrical brain potentials (N400) to nouns and preceding articles by the  
32         probability that people use a word to continue the sentence fragment ('cloze'). In our direct  
33         replication study spanning 9 laboratories ( $N=334$ ), pre-registered replication-analyses and  
34         exploratory Bayes Factor analyses successfully replicated the noun-results but, crucially, not  
35         the article-results. Pre-registered single-trial analyses also yielded a statistically significant  
36         effect for the nouns but not the articles. Exploratory Bayesian single-trial analyses showed  
37         that the article-effect may be non-zero but is likely far smaller than originally reported and  
38         too small to observe without very large sample sizes. Our results do not support the view that  
39         readers routinely pre-activate the phonological form of predictable words.

40

41 INTRODUCTION

42 In the last decades, the idea that people routinely and implicitly predict upcoming  
43 words during language comprehension turned from a highly controversial hypothesis to a  
44 widely accepted assumption. Initial objections to prediction in language were based on a lack  
45 of empirical support (e.g., Zwitserlood, 1989), incompatibility with traditional bottom-up  
46 models and contemporary interactive models of language comprehension (e.g., Kintsch,  
47 1988; Marslen-Wilson & Tyler, 1988), and the purported futility of prediction in a generative  
48 system where sentences can continue in infinitely many different ways (Jackendoff, 2002).

49 Current theories of language comprehension, however, reject such objections and posit  
50 prediction as an integral and inevitable mechanism by which comprehension proceeds  
51 quickly and incrementally (e.g., Altmann & Mirkovic, 2009; Dell & Chang, 2014; Pickering  
52 & Garrod, 2013). Prediction, i.e., context-based pre-activation of an upcoming linguistic  
53 input, is thought to occur at all levels of linguistic representation (semantic, morpho-syntactic  
54 and phonological/orthographic) and serves to facilitate the integration of newly available  
55 bottom-up information into the unfolding sentence- or discourse-representation. In this line of  
56 thought, language is yet another domain in which the brain acts as a prediction machine  
57 (Clark, 2013; Van Berkum, 2010; see also Friston, 2005, 2010; Summerfield & De Lange,  
58 2014), hard-wired to continuously match sensory inputs with top-down, grammatical or  
59 probabilistic expectations based on context and memory.

60 What promoted linguistic prediction from outlandish and deeply contentious to  
61 ubiquitous and somewhat anodyne? One of the key and most acclaimed pieces of empirical  
62 evidence for linguistic prediction to date comes from a landmark *Nature Neuroscience*  
63 publication by DeLong, Urbach and Kutas (2005), whose approach exploited a phonological  
64 rule of English whereby the indefinite article is realized as *a* before consonant-initial words  
65 and as *an* before vowel-initial words. In their experiment, participants read sentences of

66 varying degree of contextual constraint that led to expectations for a particular consonant- or  
67 vowel-initial noun. This expectation was operationalized as a word's cloze probability  
68 (cloze), calculated in a separate, non-speeded sentence completion task as the percentage of  
69 continuations of a sentence fragment with that word (Taylor, 1953). For example, the  
70 sentence fragment "The day was breezy so the boy went outside to fly..." is continued with  
71 'a' by 86% of participants, and "The day was breezy so the boy went outside to fly a..." is  
72 continued with 'kite' by 89% of participants. In the main experiment, word-by-word sentence  
73 presentation enabled DeLong and colleagues to examine electrical brain activity elicited by  
74 articles that were concordant with the highly expected but yet unseen noun ('a', followed by  
75 'kite'), or by articles that were incompatible with the highly expected noun and heralded a  
76 less expected one<sup>1</sup> ('an', followed by 'airplane'). The dependent measure was the amplitude  
77 of the N400<sup>2</sup> event-related potential (ERP), a negative ERP deflection that peaks  
78 approximately 400 ms after word onset and is maximal at centroparietal electrodes (Kutas &  
79 Hillyard, 1980). The N400 is elicited by every word of an unfolding sentence and its  
80 amplitude is smaller (less negative) with increasing ease of semantic processing (Kutas &  
81 Hillyard, 1984). DeLong et al. found that the N400 amplitude for a given word decreased as a  
82 function of increasing cloze probability, both for nouns and, critically, for articles. DeLong et  
83 al. presented the systematic and graded N400 modulation by article-cloze as strong evidence  
84 that participants activated the nouns and articles in advance of their appearance, and that the

---

<sup>1</sup> Of note, an unexpected like 'a/an' does not rule out that the expected noun appears, just that it appears as the immediately following word (e.g., 'an old kite'), we return to this issue in the Discussion.

<sup>2</sup> In this article, we use "N400 amplitude" as a shorthand for "ERP amplitude in the time window associated with the N400"; this ERP amplitude is actually a sum of the N400 ERP component and other ERP components (reflecting other aspects of cognition) that overlap with it in time and space.

85 disconfirmation of this prediction by the less-expected articles resulted in processing  
86 difficulty (higher N400 amplitude at the article).

87 The results obtained with this elegant design warranted a much stronger conclusion  
88 than related results available at the time. Previous studies that employed a visual-world  
89 paradigm had revealed listeners' anticipatory eye-movements towards visual objects on the  
90 basis of probabilistic or grammatical considerations (e.g., Altmann & Kamide, 1999).  
91 However, predictions in such studies are scaffolded onto already-available visual context, and  
92 therefore do not measure purely pre-activation, but perhaps re-activation of word information  
93 previously activated by the visual object itself (Huettig, 2015). DeLong and colleagues  
94 examined brain responses to information associated with concepts that were not pre-specified  
95 and had to be retrieved from long-term memory 'on-the-fly'. Furthermore, DeLong and  
96 colleagues were the first to muster evidence for highly specific pre-activation of a word's  
97 phonological form, rather than merely its semantic (e.g., Federmeier & Kutas, 1999) or  
98 morpho-syntactic features (e.g., Van Berkum, Brown, Zwitserlood, Kooijman & Hagoort,  
99 2005; Wicha, Moreno & Kutas, 2004). Crucially, as their demonstration involved  
100 semantically identical articles (function words) rather than nouns or adjectives (content  
101 words) that are rich in meaning, the observed N400 modulation by article-cloze is unlikely to  
102 reflect difficulty interpreting the articles themselves. Most notably, DeLong and colleagues  
103 were the first to examine brain activity elicited by a range of more- or less-predictable  
104 articles, not simply most- versus least-expected. Based on the observed correlation, they  
105 argued that pre-activation is not all-or-none and limited to highly constraining contexts, but  
106 occurs in a graded, probabilistic fashion, with the strength of a word pre-activation  
107 proportional to its cloze probability. Moreover, they concluded that prediction is an integral  
108 part of real-time language processing and, most likely, a mechanism for propelling the  
109 comprehension system to keep up with the rapid pace of natural language.

110 DeLong et al.'s study has had an immense impact on the field of psycholinguistics,  
111 neurolinguistics and beyond. It is cited by authoritative reviews (e.g., Altmann & Mirkovic,  
112 2009; Hagoort, 2017; Lau, Phillips & Poeppel, 2008; Pickering & Clark, 2014; Pickering &  
113 Garrod, 2007) as delivering decisive evidence for probabilistic prediction of words all the  
114 way up to their phonological form. Moreover, as a demonstration of pre-activation of  
115 phonological form (sound) during reading, it is sometimes cited as evidence for 'prediction  
116 through production' (e.g., Pickering & Garrod, 2013), the hypothesis that linguistic  
117 predictions are implicitly generated by the language production system. To date, DeLong et  
118 al. has received a total of 757 citations (Google Scholar), averaging to more than 1 citation  
119 per week over the past decade, with an increasing number of citations in each subsequent  
120 year. The results also played an important role in settling an ongoing debate in the  
121 neuroscience of language. It provided the clearest evidence that the N400 component, which  
122 some researchers had long taken to directly index the high-level compositional processes by  
123 which people integrate a word's meaning with its context (Brown & Hagoort, 1993; Chwilla,  
124 Brown & Hagoort, 1995; Connolly & Phillips, 1994; Friederici, Steinhauer & Frisch, 1999;  
125 Van Berkum, Hagoort & Brown, 1999; Van Petten, Coulson, Rubin, Plante, & Parks, 1999),  
126 actually reflected non-compositional processes by which word information is accessed as a  
127 function of context (e.g., Kutas & Hillyard, 1984).

128 But how robust are gradient effects of form prediction? In over a decade that has  
129 passed since the publication by DeLong and colleagues, there is still no published study that  
130 directly replicates their graded pattern of results (for an overview, see Ito, Martin &  
131 Nieuwland, 2017b). DeLong and colleagues also performed an alternative analysis of the  
132 same data, using cloze as a categorical variable instead of a continuous variable. This analysis  
133 did not yield a statistically significant result (p.59 in DeLong, 2009), and was not mentioned  
134 in the published report. In at least three other unpublished data sets (DeLong, 2009;

135 Miyamoto, 2016), DeLong and colleagues did not find a significant correlation between  
136 article-N400 and cloze probability. Martin, Thierry, Kuipers, Boutonnet, Foucart and Costa  
137 (2013) claimed a successful conceptual replication in native speakers of English but not in  
138 bilinguals. However, their study did not test for a graded effect of cloze, and differed from  
139 the original in many crucial aspects of the experimental design, data-preprocessing and  
140 statistical analysis, clouding both a qualitative and quantitative comparison to the original  
141 results. Moreover, two attempts to replicate the Martin et al. results in English monolinguals  
142 failed to yield a reliable effect of cloze on article-ERPs (Ito, Martin & Nieuwland, 2017b; for  
143 results that combined data from monolinguals and bilinguals, see Ito, Martin & Nieuwland,  
144 2017a).

145 As the tremendous scientific impact of the DeLong et al. findings is at odds with the  
146 apparent lack of replication attempts, we report here a direct replication study. Inspired by  
147 recent demonstrations for the need for large subject-samples in psychology and neuroscience  
148 research (Button et al., 2013; Open Science Collaboration, 2015), our replication spanned 9  
149 laboratories each with a sample size equal to or greater than that of the original. In addition to  
150 duplicating the original analysis, our replication attempt also seeks to improve upon DeLong  
151 et al.'s data analysis. DeLong et al.'s original analysis reduced an initial pool of 2560 data  
152 points (32 subjects who each read 80 sentences) to 10 grand-average values, by averaging  
153 N400 responses over trials within 10 cloze probability decile-bins (cloze 0-10, 11-20, et  
154 cetera), per participant and then averaging over participants, even though these bins held  
155 greatly different numbers of observations (for example, the 0-10 cloze bin contained 37.5%  
156 of all data, whereas the 90-100 cloze bin contained only about 4%, which means that the  
157 reliability of the estimates per bin greatly differ, increasing the likelihood of obtaining  
158 spurious results; for additional discussion see Ito et al., 2017b). These 10 values were  
159 correlated with the average cloze value per bin, yielding numerically high correlation

160 coefficients with large confidence intervals (for example, the Cz electrode showed a  
161 statistically significant  $r$ -value of 0.68 with a 95% confidence interval ranging from 0.09 to  
162 0.92). However, this analysis potentially compromises power by discretizing cloze  
163 probability into deciles and not distinguishing various sources of subject-, item-, bin-, and  
164 trial-level variation. Furthermore, treating subjects as a fixed rather than random factor  
165 potentially inflates false positive rates, since the overall cloze effect is confounded with by-  
166 subject variation in the effect (Barr, Levy, Scheepers & Tily, 2013; Clark, 1973).

167 In our replication study, we followed two pre-registered analysis routes: a *replication*  
168 *analysis* that duplicated the DeLong et al. analysis, and a *single-trial analysis* that modelled  
169 variance at the level of item and subject (with a linear mixed-effects model), which offers  
170 better control over false-positives than the replication analysis when analyzing effects of the  
171 continuous predictor cloze probability. The effect of cloze on noun-elicited N400s (DeLong  
172 et al., 2015; Kutas & Hillyard, 1984) is a necessary but not sufficient evidence for the claim  
173 on pre-activation in language processing (as it is also compatible with the view that the  
174 noun's cloze probability correlates with the ease of integration of that noun into the context).  
175 It serves as a manipulation check to ensure that the experiment is able to successfully detect  
176 graded variation in N400 amplitude, but does not provide strong evidence for the prediction  
177 of phonological form. That evidence would come from the ERPs elicited by articles.  
178 Observing a reliable effect of cloze on article-elicited N400s in the replication analysis and,  
179 in particular, in the single-trial analysis, would constitute powerful evidence for the pre-  
180 activation of phonological form during reading.

181

## 182 **RESULTS**

183 We first obtained offline cloze probabilities for all target articles and nouns from a  
184 group of native English speakers. These values closely resembled those of the original study

185 (see Methods for details). In the subsequent ERP experiment, a different group of participants  
186 ( $N=334$ ) read the sentences word-by-word from a computer display at a rate of 2 words per  
187 second while we recorded their electrical brain activity at the scalp. The replication analysis  
188 and single-trial analysis described below were each pre-registered at <https://osf.io/eyzaq/>.

189 **Replication analysis**

190 We sorted the articles and nouns into 10 bins based on each word's cloze probability  
191 (e.g., items with 0-10% cloze were put in one bin, 10-20% in another, etc.). For each  
192 laboratory, we averaged ERPs per bin first within, then across, participants. No baseline  
193 correction was used, following the procedure described in the Methods section in DeLong et  
194 al (2005). We then correlated the averaged cloze values per bin with mean ERP amplitude in  
195 the N400 time window (200-500 ms) elicited by the nouns (for the noun analysis) or articles  
196 (for the article analysis) from the corresponding bin, yielding a Pearson correlation  
197 coefficient ( $r$ -value) per EEG channel. This analysis yielded a very different pattern than  
198 DeLong et al. observed (Fig. 1). In no laboratory did article-N400 amplitude at centro-  
199 parietal sites become significantly smaller (less negative) as article-cloze probability  
200 increased (in fact, in most laboratories the pattern went into the opposite direction). Only in  
201 one laboratory (Lab 2) did the correlation coefficient have a  $p$ -value below .05 in the  
202 predicted direction (positive) at any electrode (uncorrected for multiple comparisons), but this  
203 effect was observed at a few left-frontal electrodes, not at the central-parietal electrodes  
204 where DeLong et al. found their N400 effects. Moreover, in two laboratories (Labs 3 and 5),  
205 a statistically significant effect was observed in the opposite direction, larger (more negative)  
206 article-N400 amplitude for articles with increasing cloze probability. For the nouns, the  
207 pattern was more similar to the DeLong et al. results. In six laboratories (Lab 2, 3, 4, 6, 7, and  
208 9), noun-N400 amplitude for nouns at central-parietal or parietal-occipital electrodes became

209 smaller with increasing noun-cloze, and in two other laboratories (Lab 5 and 8) the effects  
210 clearly went in the expected direction without reaching statistical significance.

211 <INSERT FIG 1>

212 DeLong et al. recently mentioned using a 500 ms baseline correction procedure that  
213 was not mentioned in the published study (personal communication by DeLong, March  
214 2017). In an exploratory analysis, we therefore recomputed the correlations based on data  
215 pooled from all laboratories using this baseline correction procedure (Fig. 2). This analysis  
216 also showed a lack of statistically significant positive correlations for the articles, but  
217 statistically significant positive correlations for the nouns. In exploratory Bayesian analyses  
218 reported below, we perform an analysis to establish whether these results are consistent with  
219 the size and direction of the effects reported by DeLong et al., regardless of statistical  
220 significance.

221 <INSERT FIG 2>

222

### 223 **Single-trial analysis**

224 We first performed baseline correction by subtracting the average amplitude in the 100  
225 ms time window before word onset. Baseline-corrected ERPs for relatively expected and  
226 unexpected words and difference waveforms are shown in Fig. 3. Then, for the data pooled  
227 across all laboratories, we used linear mixed effects models to regress the N400 amplitude (in  
228 a spatiotemporal region of interest selected *a priori* based on the DeLong et al. results) on  
229 cloze probability. For the articles, the effect of cloze was not statistically significant at the  
230  $\alpha=.05$  level,  $\beta = .29$ , CI [-.08, .67],  $\chi^2(1) = 2.31$ ,  $p = .13$  (see Fig. 4, left panel)<sup>3</sup>, with  $\beta$   
231 referring to the N400 difference in microvolts associated with stepping from 0% to 100%

---

<sup>3</sup> Unless otherwise indicated,  $p$ -values are two-tailed, and CIs are two-tailed 95% confidence intervals.

232 cloze. The effect of cloze on N400 amplitude at the article did not significantly differ  
233 between laboratories,  $\chi^2(8) = 7.90, p = .44$ . For the nouns, however, higher cloze values were  
234 strongly associated with smaller N400s,  $\beta = 2.22$ , CI [1.76, 2.69],  $\chi^2(1) = 56.50, p < .001$  (see  
235 Fig. 4, right panel). This pattern did not significantly differ between laboratories,  $\chi^2(8) =$   
236 11.59,  $p = .17$ . The effect of cloze on noun-N400s was statistically different from its effect on  
237 article-N400s,  $\chi^2(1) = 31.38, p < .001$ .

238 <INSERT FIG 3 & 4>

239 *Exploratory (i.e., not pre-registered) single-trial analyses*

240 The effect of article-cloze did not significantly vary as a function of subject  
241 comprehension question accuracy,  $\chi^2(1) = 0.45, p = .50$ . In addition, the effect of article-cloze  
242 was also not statistically significant when subject comprehension accuracy was included in  
243 the analysis (100 ms baseline:  $\beta = .24$ , CI [-.17, .64],  $\chi^2(1) = 1.27, p = .26$ ).

244 In our dataset, an analysis in the 500 to 100 ms time window *before* article-onset revealed  
245 a non-significant effect of cloze that resembled the pattern observed *after* article-onset,  $\beta =$   
246 .16, CI [-.07, .39],  $\chi^2(1) = 1.82, p = .18$  (Fig. 5). Because the sentence context of each item  
247 was identical for the expected and unexpected article, effects in the pre-article window cannot  
248 be meaningfully related to the appearance of the article. Effects in this window must  
249 therefore be due to a spurious mix of ‘residual EEG background noise’ (activity that differed  
250 between expected and unexpected conditions but was unrelated to actual expectancy) with  
251 EEG activity associated with the specific word appearing before the article (which varied  
252 between items in terms of lexical characteristics, contextual constraint, and sentence  
253 position). The observed result in this time window therefore suggests that a 500 ms baseline  
254 correction procedure, which was used but not reported in DeLong et al. (2005), would better  
255 correct for pre-article voltage-levels. We repeated our analysis with the 500 ms baseline  
256 correction procedure. Compared to the article-cloze effect observed in the pre-registered

257 analysis, the observed effect with the new baseline procedure (Fig. 5) was numerically  
258 smaller and yielded a higher  $p$ -value ( $\beta = .14$ , CI [-.25, .53],  $\chi^2(1) = 0.46$ ,  $p = .50$ ).

259 <INSERT FIG 5>

260 Upon request of reviewers for this journal, we also performed an additional exploratory  
261 analysis with cloze as a dichotomous variable (based on a medium-split, thus disregarding the  
262 known variability in cloze values). We note that this type of analysis was not reported in  
263 DeLong et al. (2005), although it was reported in the corresponding thesis chapter (DeLong,  
264 2009) and did not yield a statistically significant effect of cloze on article-elicited ERPs. We  
265 performed this analysis for articles (100 and 500 ms baseline correction) and nouns. The  
266 results did not change substantially, and, in fact, each analysis yielded a lower  $\chi^2$  value (and  
267 higher  $p$ -value) for the cloze variable than the corresponding analysis with cloze as a  
268 continuous predictor. The results can be reproduced from our online dataset and code.

269

## 270 **Exploratory Bayesian analyses**

271 For the articles, our pre-registered replication analyses yielded non-significant  $p$ -  
272 values, indicating failure to reject the null-hypothesis that cloze has no effect on N400  
273 activity. To better adjudicate between the null-hypothesis ( $H_0$ ) and an alternative hypothesis  
274 ( $H_r$ ), we performed an exploratory replication Bayes factor analysis for correlations  
275 (Wagenmakers, Verhagen & Ly 2016). The obtained replication Bayes factor quantifies the  
276 evidence that there is an effect in the size and direction reported by DeLong et al. (see Fig. 6).  
277 For the articles, this yielded strong to extremely strong evidence for the null hypothesis that  
278 the effect of cloze is zero, with  $BF_{0r}$  values up to 154 (at the Cz electrode depicted by  
279 DeLong et al.,  $BF_{0r} = 77$ ), and strongest evidence at the posterior channels. For the nouns, we  
280 obtained extremely strong evidence for the alternative hypothesis that the effect is non-zero,  
281 particularly at posterior channels, with  $BF_{10}$  values up to 9,163,515 (at Cz,  $BF_{r0} = 10,725$ ).

282 The pattern of results was similar when the 500 ms pre-stimulus baseline correction was  
283 applied.

284 < INSERT FIG 6 >

285 Next, we computed Bayesian mixed-effect model estimates ( $\beta$ ) and 95% credible  
286 intervals (CrI) for our single-trial analyses, using priors based on the results from DeLong et  
287 al. In both of our article-analyses credible intervals included zero (100 ms baseline:  $\beta = .31$ ,  
288 CrI [-.06 .69]; 500 ms baseline:  $\beta = .17$ , CrI [-.22 .55]). For the nouns, zero was not within  
289 the credible interval:  $\beta = 2.24$ , CrI [1.77 2.70]. The analyses suggest that the data (combined  
290 with prior assumptions about the effect) are not very consistent with the hypothesis that the  
291 article-effect is zero (further information and posterior summaries are available in Fig. 7), but  
292 also are extremely inconsistent with the hypothesis that the article-effect is as big as that  
293 observed by DeLong and colleagues (2005). The data are most consistent with an effect that  
294 is more likely to be positive than zero or negative, but is very small (so small that it was not  
295 detected at traditional significance levels in this large-scale experiment with substantially  
296 higher power than previous experiments).

297 < INSERT FIG 7 >

298

## 299 **Control experiment**

300 Lack of a statistically significant, article-elicited prediction effect could reflect a  
301 general insensitivity of our participants to the phonologically conditioned variation of the  
302 English indefinite article, i.e., *a/an* alternation. We ruled out this alternative explanation in an  
303 additional experiment that followed the replication experiment as part of the same  
304 experimental session. Participants read 80 short sentences containing the same nouns as the  
305 replication experiment, preceded by a phonologically licit or illicit article (e.g., “David found  
306 a/an apple...”), presented in the same manner as before. In each laboratory, nouns following

307 illicit articles elicited a late positive-going waveform compared to nouns following licit  
308 articles (see Fig. 8), starting at about 500 ms after word onset and strongest at parietal  
309 electrodes. This standard P600 effect (Osterhout & Holcomb, 1992) was confirmed in a  
310 single-trial analysis,  $\chi^2(1) = 83.09, p < .001$ , and did not significantly differ between labs,  
311  $\chi^2(8) = 8.98, p = .35$ .

312 <INSERT FIG 8>

313

314 **DISCUSSION**

315 In a landmark study, DeLong, Urbach and Kutas observed a statistically significant,  
316 graded modulation of article- and noun-elicited electrical brain potentials (N400) by the pre-  
317 determined probability that people continue a sentence fragment with that word (cloze). They  
318 concluded that people routinely and probabilistically pre-activate upcoming words to a high  
319 level of detail, including whether a word starts with a consonant or vowel. Our *direct*  
320 *replication* study spanning 9 laboratories found a statistically significant effect of cloze on  
321 noun-elicited N400 activity but, critically, no significant effect of cloze on article-elicited  
322 N400 activity. This pattern was observed in a pre-registered replication analysis that  
323 duplicated the original study's analysis, and a pre-registered single-trial analysis that  
324 modelled variance at the level of item and subject. Exploratory Replication Bayes Factor  
325 analyses confirmed that we successfully replicated the direction and size of the correlations  
326 reported by Delong et al. for the nouns, but not for the articles. Exploratory Bayesian mixed-  
327 effects model analyses suggested that, while there is some evidence that the true population-  
328 level effect may be in the direction reported by DeLong and colleagues, the effect is likely far  
329 smaller than what they reported. In fact, the effect is likely too small to be meaningfully  
330 observed without very large sample sizes. Finally, a control experiment confirmed that our

331 participants did respect the phonological alternation *a/an* of the article with nouns used in the  
332 replication experiment.

333 Our findings thus challenge one empirical cornerstone of the ‘strong prediction view’  
334 held by current theories of language comprehension (e.g., Altmann & Mirkovic, 2009;  
335 Pickering & Garrod, 2013). The strong prediction view entails two key claims. The first is  
336 that people pre-activate words at all levels of representation in a routine and implicit (i.e.,  
337 non-strategic) fashion. Pre-activation is not limited to a word’s meaning, but includes its  
338 grammatical features and even its orthographic and/or phonological form. This would put  
339 language on a par with other cognitive systems such as visual perception, wherein higher-  
340 level brain regions attempt to predict lower-level inputs (Friston, 2005, 2010; Summerfield &  
341 De Lange, 2014). The second claim is that pre-activation occurs at all levels of contextual  
342 support and gradually increases in strength with the level of contextual support. When  
343 contextual support for a specific word is high, like at a 100% cloze value, the word’s form  
344 and meaning is strongly pre-activated. When contextual support for a word is low, like when  
345 it is one amongst 20 words each with a 5% cloze value, pre-activation is distributed across  
346 multiple potential continuations. However, even then, a word’s form and meaning are pre-  
347 activated, just weakly so. The strength of pre-activation is probabilistic, that is, linked to  
348 estimated probability of occurrence.

349 DeLong and colleagues, and subsequently other scientists (e.g., Dell & Chang, 2014;  
350 Pickering & Clark, 2013), took their results as the evidence to support both these claims.  
351 DeLong et al. (2005) was – and still is - the only study to date that measured pre-activation at  
352 the prenominal articles *a* and *an* that do not differ in their semantic or grammatical content,  
353 and that observed a graded relationship between cloze and N400 activity across a range of  
354 low- and high-cloze words, rather than merely a difference between low- and high-cloze  
355 words. Given that the use of these articles depends on whether the next word starts with a

356 vowel or consonant, their results were considered as powerful evidence that participants  
357 probabilistically pre-activated the initial sound of upcoming nouns.

358 However, we show that there is no statistically significant effect of cloze on article-  
359 elicited N400 activity, using a sample size more than ten times that of the original, and a  
360 statistical analysis that better accounts for sources of non-independence than the original  
361 averaging-based correlation approach. If an effect of cloze on article-N400s exists at all, its  
362 true effect size is so small that it cannot be reliably detected even in an expansive multi-  
363 laboratory approach, let alone in the typical sample size in psycholinguistic and  
364 neurolinguistic experiments (roughly,  $N= 30$ ). This means that even if article-cloze is  
365 associated with a graded modulation of N400 amplitudes, this effect seems to be so small that  
366 it cannot be reliably measured with small samples, and thus the previous studies may not  
367 have contributed much reliable information to our understanding of this effect. Moreover, it  
368 is also possible that the effect is sensitive to specifics of the experimental procedure and  
369 context such that it lacks generalizability. Current theoretical positions thus either require  
370 new strong evidence for phonological pre-activation or require revision. In particular, one  
371 claim from the strong prediction view, namely that pre-activation routinely occurs across all –  
372 including phonological – levels (Pickering & Garrod, 2013), can no longer be viewed as  
373 having strong empirical support. Our work impels the field to think differently about what  
374 constitutes strong evidence within a theory, but also highlights the need for a theory of  
375 linguistic prediction to formulate quantitative predictions about the effect-size of to-be-  
376 observed effects (for discussion, see also Vasishth, Mertzen, Jäger & Gelman, 2018).

377 By contrast, we observed a strong and statistically significant effect of cloze on noun-  
378 elicited activity in the majority of our analyses. Although three of the nine laboratories did  
379 not show statistically significant correlations between noun-cloze and N400s, data pooled  
380 across all laboratories showed a strong and statistically significant noun-cloze effect, and our

381 Replication Bayes Factor analysis overwhelmingly replicated the direction and size of the  
382 noun-cloze effect of DeLong et al. Moreover, our single-trial analysis revealed a significant  
383 noun-cloze effect in each of the laboratories, further demonstrating that our single-trial  
384 analysis is a more powerful approach than the averaged-based correlation approach of  
385 DeLong et al. These results are therefore consistent with the handful of studies that reported a  
386 graded relationship between noun-cloze and noun-N400s (DeLong et al., 2005; Kutas &  
387 Hillyard, 1984; Wlotko & Federmeier, 2012).

388 Where do our results leave the strong prediction view? Following the experimental logic  
389 of DeLong et al, we do not have sufficient evidence to conclude that people routinely pre-  
390 activate the initial phoneme of an upcoming noun, or perhaps any other word form  
391 information. Without pre-activation of the initial phoneme, the specific instantiation of the  
392 article does not cause people to revise their prediction about the meaning of the upcoming  
393 noun, thus lacking any impact on processing. Crucially, this conclusion is incompatible with  
394 the strong prediction view, because it suggests that pre-activation does not occur at the level  
395 of detail that is often assumed. Our results are also incompatible with an alternative  
396 interpretation of the DeLong et al. findings that people predict the article itself together with  
397 the noun (Ito, Corley, Pickering, Martin & Nieuwland, 2016; Van Petten & Luka, 2012), and  
398 they pose a serious challenge to the theory that comprehenders predict upcoming words,  
399 including their initial phonemes, through implicit production (Pickering & Garrod, 2013).  
400 Crucially, the idea that prediction is probabilistic, rather than all-or-none, is now  
401 questionable, given that there is no other published report of a pre-activation gradient (also,  
402 see Van Petten & Luka, 2012, for a critique of the DeLong et al. conclusions that graded  
403 effects evidence graded pre-activation). Although other studies have claimed prediction of  
404 form (Ito et al., 2016) or a prediction gradient (Smith & Levy, 2013), no study has  
405 indisputably demonstrated graded pre-activation, i.e., graded effects occurring *before* the

406 noun. Effects that are observed upon, rather than before the noun, do not purely index pre-  
407 activation but can index a mixture of memory retrieval and semantic integration processes  
408 instigated by the noun itself (Baggio & Hagoort, 2011; Lau, Namyst, Fogel & Delgado, 2016;  
409 Nieuwland et al., 2018; Otten & Van Berkum, 2008; Steinhauer, Royle, Drury & Fromont,  
410 2017). Therefore, there is currently no clear evidence to support routine probabilistic pre-  
411 activation of a noun's phonological form during sentence comprehension.

412 Our results, however, should not be taken as evidence against prediction in language  
413 processing more generally, and we believe that prediction could play an important role in  
414 language comprehension. In addition, our results do not necessarily exclude phonological  
415 form pre-activation, and we temper our conclusion with a caveat stemming from the *a/an*  
416 manipulation. For this manipulation to 'work', people must specifically predict the initial  
417 phoneme of the next word, and revise this prediction when faced with an unexpected article.  
418 However, because articles are only diagnostic about the next word within the noun phrase,  
419 rather than about the head noun itself, an unexpected article does not refute the upcoming  
420 noun, it merely signals that another word would come first (e.g., 'an old kite'). This opens up  
421 explanations for why the *a/an* manipulation 'fails' (see also Ito et al., 2017a,b). In addition,  
422 comprehenders may not predict the noun to follow immediately, but at a later point; the  
423 unexpected article then does not evoke a change in prediction. Predictions about a specific  
424 position may be disconfirmed too often in natural language to be viable. This idea is  
425 supported by corpus data (Corpus of Contemporary American English and British National  
426 Corpus), showing a mere 33% probability that *a/an* is directly followed by a noun.  
427 Alternatively, people predict the noun to come next, but only revise their prediction about its  
428 linear position while retaining the prediction about its meaning. So perhaps a revision of the  
429 predicted meaning, not the position, is required to trigger differential ERPs. In both of these

430 hypothetical scenarios, people do not revise their prediction about the upcoming noun's  
431 meaning unless they must.

432 Our results can be straightforwardly reconciled with effects reported for other pre-  
433 nominal manipulations, such as those of Dutch or Spanish article-gender (e.g., Van Berkum  
434 et al., 2005; Otten, Nieuwland & Van Berkum, 2008; Otten & Van Berkum, 2009; Wicha et  
435 al., 2004). Unlike a/an articles, gender-marked articles can immediately disconfirm the noun,  
436 because article- and noun-gender agrees regardless of intervening words (e.g., the Spanish  
437 article 'el' heralds a masculine noun). Revising the prediction about the noun presumably  
438 results in a semantic processing cost, thereby modulating N400 activity (e.g., Kochari &  
439 Flecken, 2017; Otten & Van Berkum, 2009). Although gender-marked articles do not  
440 consistently incur the exact same type of effect (for a recent review, see Kochari & Flecken,  
441 2017) and have only been observed at very high cloze values, previous studies suggest that a  
442 noun's grammatical gender can be pre-activated along with its meaning. Compared to this  
443 gender-manipulation, DeLong et al.'s study based on the English a/an manipulation claimed a  
444 stronger version of the prediction view, namely that people predict which word comes next  
445 up to its phonological form *and*, make backwards prediction as to the phonological form of  
446 the preceding linguistic material even on the basis of probabilistic, graded information.

447 What do our results say about prediction during natural language processing? Like the  
448 conclusions by DeLong et al., ours are limited by the generalization from language  
449 comprehension in a laboratory setting. On one hand, a rich conversational or story context  
450 may enhance predictions of upcoming words, and listeners may be more likely to pre-activate  
451 the phonological form of upcoming words than readers. On the other hand, our laboratory  
452 setting offered particularly good conditions for prediction of the next word's initial sound to  
453 occur. Each article was always immediately followed by a noun, unlike in natural language.  
454 Moreover, our word presentation rate was slow compared to natural reading rates, which may

455 facilitate predictive processing (Ito et al., 2016; Wlotko & Federmeier, 2015). In natural  
456 reading, articles are hardly fixated and often skipped (e.g., O'Regan 1979). In short,  
457 arguments can be made both for and against phonological form prediction in natural language  
458 settings, and novel avenues of experimentation are needed to settle this issue.

459 DeLong and colleagues recently stated an omission in the description of their data  
460 analysis, i.e., a baseline procedure was applied to the data but inadvertently omitted from the  
461 description (DeLong et al., 2005). We have shown that our conclusions hold regardless of the  
462 baseline procedure. In a recent commentary, Delong, Urbach, and Kutas (2017) also  
463 described filler-sentences in their experiment, which were omitted from their original report,  
464 and were neither provided nor mentioned to us by the authors upon our request for the  
465 stimuli. DeLong et al. used the existence of these filler-sentences to dismiss an alternative  
466 explanation of their original findings, namely that an unusual experimental context wherein  
467 every sentence contains an article-noun combination leads participants to strategically predict  
468 upcoming nouns. Following this logic, our results were obtained *despite* an experimental  
469 context that could inadvertently encourage strategic prediction (for demonstrations of  
470 experimental context boosting predictive processing, see Brothers, Swaab & Traxler, 2017;  
471 Lau, Holcomb & Kuperberg, 2013). Therefore, the presence of fillers in their experiment  
472 versus absence in ours cannot straightforwardly explain the different results, and may even  
473 strengthen our conclusions.

474 Since becoming publicly available as a pre-print (Nieuwland et al., 2017), our study has  
475 been simultaneously criticized for being not a sufficiently direct replication (due to the  
476 differences in fillers and baseline procedure; DeLong, Urbach & Kutas, 2017; Yan,  
477 Kuperberg & Jaeger, 2017) and for being a too direct replication (because we base our  
478 analysis on the same theoretical assumptions as the original study, rather than applying an ad-  
479 hoc transformation or different kind of analysis that might 'reveal' the effect; e.g., Yan et al.,

480 2017). As an example of the latter, an unpublished commentary by Yan et al. (2017) raises an  
481 interesting point that cloze probability should be log-transformed to better approximate their  
482 suggested index of probabilistic semantic prediction, the Bayesian surprise over the noun  
483 semantics upon encountering the article. Yan et al. describe a number of exploratory  
484 reanalyses of our single-trial data with the log-transform, and one of those exploratory  
485 analyses yields a small but statistically significant effect of article-cloze ( $p=.015$ ). Ultimately,  
486 however, their conclusion is not that different from ours, namely that there is some evidence  
487 in our data that the effect is non-zero. More importantly, their commentary demonstrates that  
488 our dataset, like any complex EEG dataset, can be analyzed in many different ways, which  
489 can lead to different outcomes. However, even if alternative analyses are well-motivated after  
490 the fact, the problem remains that they are contingent on the data, and the accompanying  
491 researcher degrees of freedom lead to a multiple comparison problem (e.g., Gelman & Loken,  
492 2013; Luck & Gaspelin, 2017). We pre-registered our main analyses and none of these  
493 allowed us conclude that the DeLong et al. study replicated. Yan et al. present an alternative  
494 analysis that is exploratory and that itself requires further replication. Moreover, their  
495 analysis also raises a novel set of important concerns. For example, log-transformation of  
496 cloze also boosts the effect in the pre-article time window ( $p=.058$ ), where there cannot be a  
497 meaningful effect, possibly because it amplifies ‘noise’ (between-item differences at the low  
498 end of the cloze-scale that have nothing to do with prediction of the article). Furthermore,  
499 log-transformation does not yield a significant effect with the original baseline procedure of  
500 DeLong et al., and it strongly boosts the impact of items with zero cloze, i.e., the items that  
501 are problematic because their predictability cannot be accurately estimated (of note, without  
502 zero-cloze values in their analysis, higher cloze leads to more negative, not positive voltage).  
503 Yan et al. report that log-transformation yields somewhat higher  $t$ -values of cloze in this  
504 dataset and changes our non-significant effect into a significant effect, but it remains unclear

505 whether log-transformation is indeed ‘better’. Crucially, the difference between significant  
506 and not-significant itself may not be significant (Gelman & Stern, 2006), log-transformation  
507 does not yield higher t-values consistently across laboratories, does not necessarily improve  
508 model fit, and does not yield higher t-values or improve model fit in another large dataset  
509 (collapsed data from Ito et al., 2017a; Nieuwland, 2016; Nieuwland & Martin, 2012; total  $N =$   
510 124). Finally, it is unknown whether log-transformation weakens rather than strengthens the  
511 effect of the original study. Details of these and further concerns are available on  
512 <https://osf.io/mb2ud>. In sum, these concerns merely add to our main point, namely that even  
513 if analysis decisions are justifiable in retrospect, a flexible analysis practice can result in  
514 capitalizing on noise (Gelman & Loken, 2013).

515 To conclude, we failed to replicate the main result of DeLong et al., a landmark study  
516 published more than ten years ago that has not been directly replicated since. Our results  
517 suggest that, if there is an effect of article-cloze probability on the amplitude of the N400, it  
518 is too small and/or too sensitive to unknown experimental design factors to have been  
519 meaningfully measured in previous small-sample-size experiments. Our findings thus do not  
520 lend clear support the ‘strong prediction view’ in which people routinely and probabilistically  
521 pre-activate information at all levels of linguistic representation, including phonological form  
522 information such as the initial phoneme of an upcoming noun. Consequently, there is  
523 currently no convincing evidence that people routinely pre-activate the phonological form of  
524 an upcoming noun during written sentence comprehension. In addition, our findings further  
525 highlight the importance of direct replication, large sample size studies, transparent reporting  
526 and of pre-registration to advance reproducibility and replicability in the neurosciences.

527 MATERIALS AND METHODS

528           **Experimental design and materials.** Nieuwland requested all original materials from  
529 DeLong et al., including the questions and norms, with the stated purpose of direct replication  
530 (personal communication, November 4 and 19, 2015), upon which DeLong et al. made  
531 available the 80 sentences described in the original study. These sentences were then adapted  
532 from American to British spelling and underwent a few minor changes to ensure their  
533 suitability for British participants. The complete set of materials and the list of changes to the  
534 original materials are available online (Supplementary Table 1 and 2). The materials were 80  
535 sentence contexts with two possible continuations each: a more or less expected indefinite  
536 article + noun combination. The noun was followed by at least one subsequent word. All  
537 article + noun continuations were grammatically correct. Each article + noun combination  
538 served once as the more expected continuation and the other time as the less expected  
539 continuation, in different contexts. We divided the 160 items in two lists of 80 sentences such  
540 that each list contained each noun only once. Each participant was presented with only one  
541 list (thus, each context was seen only once). One in four sentences was followed by a yes/no  
542 comprehension question, which yielded a mean response accuracy of 95% (after taking into  
543 account ambiguity in three of the questions, see Supplemental Table 2 and 3). While this  
544 percentage is very similar to that reported by DeLong et al., we note that this cannot be  
545 directly compared to the accuracy reported in DeLong et al., because we had to create new  
546 comprehension questions in the absence of the original ones. Regardless, because Delong et  
547 al. suggested that our results were due to poor language comprehension (DeLong, Urbach &  
548 Kutas, 2017), we describe an exploratory analysis in which we attempt to account for  
549 variation in response accuracy in the statistical model.

550           We obtained article cloze and noun cloze ratings from a separate group of native  
551 speakers of English who were students at the University of Edinburgh and did not participate

552 in the ERP experiment. They were instructed to complete the sentence fragment with the best  
553 continuation that comes to mind (Taylor, 1953). We obtained article cloze ratings from 44  
554 participants for 80 sentence contexts truncated before the critical article. Noun cloze ratings  
555 were obtained by first truncating the sentences after the critical articles, and presenting two  
556 different, counterbalanced lists of 80 sentences to 30 participants each, such that a given  
557 participant only saw each sentence context with the expected or the unexpected article. The  
558 obtained values closely resemble those of the original study, with the same range (0-100% for  
559 articles and nouns), slightly lower median values (for articles and nouns, 29% and 40%,  
560 compared to 31% and 46% in the original study), but slightly higher mean values (for articles  
561 and nouns, 41% and 46%, compared to 36% and 44%). Because the sentence materials we  
562 used describe common situations that can be understood by any English speaker, and because  
563 students at the University of Edinburgh come from across the whole of the UK, we had no *a*  
564 *priori* expectation that cloze ratings would differ substantially across laboratories, and thus  
565 we did not obtain cloze norms from other sites. Consistent with this assumption, nothing in  
566 our results suggests stronger cloze effects in University of Edinburgh students compared to  
567 other students, suggesting that our cloze norms are sufficiently representative for the other  
568 universities. The raw cloze responses are available on our OSF page.

569 **Participants.** Participants were students from the University of Birmingham, Bristol,  
570 Edinburgh, Glasgow, Kent, Oxford, Stirling, York, or volunteers from the participant pool of  
571 University College London or Oxford University, who received cash or course credit for  
572 taking part in the ERP experiment. Participant information and EEG recording information  
573 per laboratory is available online (Supplementary Table 3). We pre-registered a target sample  
574 size of 40 participants per laboratory, which was thought to give at least 32 participants (the  
575 sample size of DeLong et al.) per laboratory after accounting for data loss, as was later  
576 confirmed. Due to logistic constraints, not all laboratories reached an N of 40. Because in two

577 labs corruption of data was incorrectly assumed before computing trial loss, these  
578 laboratories tested slightly more than 40 participants. All participants ( $N = 356$ ; 222 women)  
579 were right-handed, native English speakers with normal or corrected-to-normal vision,  
580 between 18–35 years (mean, 19.8 years), free from any known language or learning disorder.  
581 Eighty-nine participants reported a left-handed parent or sibling.

582 **Procedure.** After giving written informed consent, participants were tested in a single  
583 session. Sentences were presented visually in the center of a computer display, one word at a  
584 time (200 ms duration, followed by a blank screen of 300 ms duration<sup>4</sup>). Participants were  
585 instructed to read sentences for comprehension and answer yes/no comprehension questions  
586 by pressing hand-held buttons. The electroencephalogram (EEG) was recorded from at least  
587 32 electrodes.

588 The replication experiment was followed by a control experiment, which served to  
589 detect sensitivity to the correct use of the a/an rule in our participants. Participants read 80  
590 relatively short sentences (average length 8 words, range 5-11) that contained the same  
591 critical words as the replication experiment, preceded by a correct or incorrect article. As in  
592 the replication experiment, each critical word was presented only once, and was followed by  
593 at least one more word. All words were presented at the same rate as the replication  
594 experiment. There were no comprehension questions in this experiment. After the control  
595 experiment, participants performed a Verbal Fluency Test and a Reading Span test; the

---

<sup>4</sup> Due to a programming error, in four labs (1, 3, 5 and 8, which used E-prime scripts) the critical articles and nouns, but not other words, were followed by a 380 ms blank instead of the intended 300 ms. This delay is unlikely to have affected the results because if it was noticed at all, which is unlikely, it could only be noticed 500 ms after the article, i.e., after the N400 window associated with the article. Of note, the pattern of the results from the pre-registered single-trial analysis did not change when we removed these labs from the analysis.

596 results from these tests are not discussed here. All stimulus presentation scripts are publicly  
597 available in two different software packages (E-Prime and Presentation) on  
598 <https://osf.io/eyzaq>.

599         **Data processing.** Data processing was performed in BrainVision Analyzer 2.1 (Brain  
600 Products, Germany). We performed one pre-registered replication analysis that followed the  
601 DeLong et al. analysis as closely as possible and one pre-registered single-trial analysis  
602 (Open Science Framework, <https://osf.io/eyzaq>). All non-pre-registered analyses are  
603 considered as exploratory. First, we interpolated bad channels from surrounding channels,  
604 and downsampled to a common set of 22 EEG channels per laboratory which were similar in  
605 scalp location to those used by DeLong et al. One laboratory did not have 12 of the selected  
606 22 channels in its EEG channel montage, and we matched the full 22-channel layout used for  
607 other laboratories by creating 12 virtual channels from neighbouring channels using  
608 topographic interpolation by spherical splines. We then applied a 0.01-100 Hz digital band-  
609 pass filter (including 50 Hz Notch filter), re-referenced all channels to the average of the left  
610 and right mastoid channels (in a few participants with a noisy mastoid channel, only one  
611 mastoid channel was used), and segmented the continuous data into epochs from 500 ms  
612 before to 1000 ms after word onset. We then performed visual inspection of all data segments  
613 and rejected data with amplifier blocking, movement artifacts, or excessive muscle activity.  
614 Subsequently, we performed independent component analysis (Jung et al., 2000) on a 1-Hz  
615 high-pass filtered version of the data, and applied the obtained weightings to the original data  
616 to correct for blinks, eye movements or steady muscle artefacts. After this, we automatically  
617 rejected segments containing a voltage difference of over 120  $\mu$ V in a time window of 150  
618 ms or containing a voltage step of over 50  $\mu$ V/ms. Participants with fewer than 60/80 article  
619 trials or 60/80 noun trials were removed from the analysis, leaving a total of 334 participants

620 (range across laboratories 32-42, and therefore each lab had a sample size at least as large as  
621 DeLong et al.). On average, participants had 77 article trials and 77 noun trials.

622 **Pre-registered replication analysis.** We applied a 4<sup>th</sup>-order Butterworth band-pass  
623 filter at 0.2-15 Hz to the segmented data, averaged trials per participant within 10% cloze  
624 bins (0-10, 11-20, etc. until 91-100), and then averaged the participant-wise averages  
625 separately for each laboratory. Because the bins did not contain equal numbers of trials (the  
626 intermediate bins contained fewest trials), like in DeLong et al., not all participants  
627 contributed a value for each bin to the grand average per laboratory. For nouns and articles  
628 separately, and for each EEG channel, we computed the correlation between ERP amplitude  
629 in the 200-500 ms time window per bin with the average cloze probability per bin.

630 **Pre-registered single-trial analysis.** In this analysis, we did not apply the 0.2-15 Hz  
631 band-pass filter, which carries the risk of inducing data distortions (Luck, 2014; Tanner,  
632 Morgan-Short & Luck, 2015). However, we deemed it necessary to perform a baseline  
633 correction of the data. This procedure corrects for spurious voltage differences before word  
634 onset, increasing confidence that observed effects are elicited by the word rather than  
635 differences in brain activity that already existed before the word and is a standard procedure  
636 in ERP research (Luck, 2014). DeLong et al. (2005) did not report a baseline correction, nor  
637 did any of the related work from DeLong and colleagues that was reported in DeLong (2009).  
638 Yet baseline correction has been used in many other publications from the Kutas Cognitive  
639 Electrophysiology Lab. We chose a 100 ms pre-stimulus baseline as the most frequently used  
640 one both in other studies from the Kutas lab and in similar studies from other labs. For each  
641 trial, we performed baseline correction by subtracting the mean voltage of the -100 to 0 ms  
642 time window from each data point in the epoch.

643 Instead of averaging N400 data across trials and participants for subsequent statistical  
644 analysis, we performed linear mixed effects model analysis (Baayen, Davidson & Bates,

645 2008) of the single-trial N400 data, using the “lme4” package (Bates, Maechler, Bolker &  
646 Walker, 2014) in the R software (R CoreTeam, 2014). This approach simultaneously models  
647 variance associated with each subject and with each item. Especially when analyzing effects  
648 of a continuous predictor variable such as cloze probability, LMER offers better control over  
649 false-positive results than the averaged-based correlation analysis of the original. Using a  
650 spatiotemporal region-of-interest approach based on the DeLong et al. results, our dependent  
651 measure (N400 amplitude) was the average voltage across 6 centro-parietal channels  
652 (Cz/C3/C4/Pz/P3/P4) in the 200-500 ms window for each trial. Analysis scripts and data to  
653 run these scripts are publicly available on <https://osf.io/eyzaq>.

654 For articles and nouns separately, we used a maximal random effects structure as justified  
655 by the design (Barr et al., 2013), which did not include random effects for ‘laboratory’ as  
656 there were only 9 laboratories. Z-scored cloze was entered in the model as a continuous  
657 variable, and laboratory was entered as a deviation-coded nuisance predictor. We tested the  
658 effects of ‘laboratory’ and ‘cloze’ through model comparison with a  $\chi^2$  log-likelihood test.  
659 We tested whether the inclusion of a given fixed effect led to a significantly better model fit.  
660 The first model comparison examined laboratory effects, namely whether the cloze effect  
661 varied across laboratories (cloze-by-laboratory interaction) or whether the N400 magnitudes  
662 varied over laboratory (laboratory main effect). If laboratory effects were not significant, we  
663 dropped them from the analysis because they were not of theoretical interest. For the articles  
664 and nouns separately, we compared the subsequent models below. Each model included the  
665 random effects associated with the fixed effect ‘cloze’ (see Barr et al., 2014). All output  $\beta$   
666 estimates and 95% confidence intervals (CI) were transformed from z-scores back to raw  
667 scores, and then back to the 0-100% cloze range, so that the voltage estimates represent the  
668 change in voltage associated with a change in cloze probability from 0 to 100.  
669 Model 1: N400 ~ cloze \* laboratory + (cloze | subject) + (cloze | item)

670 Model 2: N400 ~ cloze + laboratory + (cloze | subject) + (cloze | item)  
671 Model 3: N400 ~ cloze + (cloze | subject) + (cloze | item)  
672 Model 4: N400 ~ (cloze | subject) + (cloze | item)

673 In an analysis that included the data from both articles and nouns, we also tested the  
674 differential effect of cloze on article ERPs and on noun ERPs by comparing models with and  
675 without an interaction between cloze and the deviation-coded factor ‘wordtype’  
676 (article/noun). Random correlations were removed for the models to converge.

677 Model 1: N400 ~ cloze \* wordtype + (cloze \* wordtype || subject) + (cloze \* wordtype ||  
678 item)  
679 Model 2: N400 ~ cloze + wordtype + (cloze \* wordtype || subject) + (cloze \* wordtype ||  
680 item)

681 **Exploratory correlation analysis.** Of note, DeLong et al. have recently described  
682 using a 500 ms baseline correction procedure that they failed to mention in DeLong et al.  
683 (2005). Using this baseline correction procedure, we recomputed the correlations that we  
684 obtained in our Replication analysis (Fig. 2). To compare our results most directly with those  
685 reported in Figure 1C of DeLong el al. (2005), we pooled data from all the laboratories to  
686 obtain a single  $r$ -value for each EEG-channel. Data were pooled after computing bin-  
687 averages per laboratory as in the original study, treating the laboratories as multiple  
688 observations of each bin-average.

689 **Exploratory single-trial analyses.** We performed an exploratory analysis in the 500  
690 to 100 ms time window *before* the article, using the originally (-100 to 0 ms) baselined data,  
691 using Model 3 and 4 from the article analysis. This window covers the first 400 ms of the  
692 word that preceded the article. Analysis in this window yielded a similar pattern as in the pre-  
693 registered analysis, which indicates that a baseline correction procedure covering the entire  
694 500 ms pre-stimulus window would account better for pre-article voltage levels. We

695 performed this additional analysis, the results of which did not change our conclusions and  
696 are shown in Figure 5.

697 We also performed an exploratory analysis in which we control for a potential  
698 influence of response accuracy, taken as a proxy for the subject's attention to the task, on  
699 predictive processing of linguistic input. We entered the (z-transformed) average response  
700 accuracy of each subject in our model, and compared the models below. Comparison of  
701 Model 1 and 2 tested whether the effect of cloze on the article-N400s depended on subject  
702 accuracy. Comparison of Model 2 and 3 tested whether there was a significant effect of cloze  
703 on article-N400s when subject accuracy was included in the model.

704 Model 1:  $N400 \sim \text{accuracy} * \text{cloze} + (\text{cloze} | \text{subject}) + (\text{cloze} | \text{item})$

705 Model 2:  $N400 \sim \text{accuracy} + \text{cloze} + (\text{cloze} | \text{subject}) + (\text{cloze} | \text{item})$

706 Model 3:  $N400 \sim \text{accuracy} + (\text{cloze} | \text{subject}) + (\text{cloze} | \text{item})$

707 **Exploratory Bayesian analyses.** Supplementing the Replication analysis, we  
708 performed a Replication Bayes factor analysis for correlations (Wagenmakers et al., 2016)  
709 using as prior the size and direction of the effect reported in the original study. We performed  
710 this test for each electrode separately, after collapsing the data points from the different  
711 laboratories. Because we had no articles in the 40-50 % cloze bin, there was a total of 9 and  
712 10 data points per laboratory for the articles and nouns, respectively. Our analysis used priors  
713 estimated from the DeLong et al. results, matched as closely as possible to our electrode  
714 locations. A Bayes factor between 3 and 10 is considered moderate evidence, between 10-30  
715 is considered strong evidence, 30-100 is very strong evidence, and values over 100 are  
716 considered extremely strong evidence (Jeffreys, 1961). In addition to using a 100 ms pre-  
717 stimulus baseline, we also computed the replication Bayes factors using the 500 ms pre-  
718 stimulus time window for baseline correction. Results are shown in Figure 5.

719           Supplementing the pre-registered single-trial analyses, we performed an exploratory  
720       Bayesian mixed-effects model analysis using the *brms* package for R (Buerkner, 2016),  
721       which fits Bayesian multilevel models using the Stan programming language (Stan  
722       Development Team, 2016). Nieuwland requested to use the results of a mixed-effects model  
723       reanalysis of the DeLong et al. data as an appropriate prior (personal communication from  
724       Nieuwland, November 14 and 22 2017); this request was declined by DeLong and  
725       colleagues. We were therefore limited to using a prior centered on a point estimate based on  
726       the Delong et al. correlation analysis, namely our estimate of the observed effect size at Cz  
727       for a difference between 0% cloze and 100% cloze ( $1.25 \mu\text{V}$  and  $3.75 \mu\text{V}$  for articles and  
728       nouns, respectively, based on visual inspection of the graphs) and a prior centered on zero for  
729       the intercept. Both priors had a normal distribution and a standard deviation of 0.5 (given the  
730       a priori expectation that average ERP voltages in this window generally fluctuate on the order  
731       of a few microvolts; note that these units are expressed in terms of the *z*-scored cloze values,  
732       rather than the original cloze values, such that  $\mu$  for the cloze prior was 0.45, which  
733       corresponds to a raw cloze effect of 1.25). We computed estimates and 95% credible intervals  
734       for each of the mixed-effects models we tested, and transformed these back into raw cloze  
735       units. The credible interval is the range of values such that one can be 95% certain that it  
736       contains the true effect, given the data, priors and the model. The results from these analyses  
737       are shown in Figure 7; the analyses suggest that, while there may be a small positive  
738       association between article cloze and ERP amplitude elicited by the articles, the effect is  
739       substantially smaller than that estimated by Delong and colleagues (2005) and likely is too  
740       small to be observed without very large sample sizes.

741           **Control experiment.** Analysis of the control experiment involved a comparison  
742       between a model with the categorical factor ‘grammaticality’ (grammatical/ungrammatical)  
743       and a model without. Our dependent measure (P600 amplitude; Osterhout & Holcomb, 1992)

744 was the average voltage across 6 centro-parietal channels (Cz/C3/C4/Pz/P3/P4) in the 500-  
745 800 ms window for each trial. Results are shown in Figure 8.

746 Model 1: P600 ~ grammaticality + (grammaticality | subject) + (grammaticality | item)

747 Model 2: P600 ~ (grammaticality | subject) + (grammaticality | item)

748    **Acknowledgements**

749        This work was partly funded by ERC Starting grant 636458 to H.J.F. We thank Matt  
750        Davis for his comments on a previous draft of this work. We thank Alexander Ly and Eric-  
751        Jan Wagenmakers for their support in computing the Replication Bayes Factors.

752

753    **Competing financial interests**

754        The authors declare no competing financial interests.

755

756 **REFERENCES**

- 757 Altmann, G. T., & Kamide, Y. (1999). Incremental interpretation at verbs: restricting the  
758 domain of subsequent reference. *Cognition*, 73(3), 247-264.
- 759 Altmann, G. T., & Mirkovic, J. (2009). Incrementality and Prediction in Human Sentence  
760 Processing. *Cogn Sci*, 33(4), 583-609. doi: 10.1111/j.1551-6709.2009.01022.x
- 761 Baayen, R. H., Davidson, D. J., & Bates, D. M. (2008). Mixed-effects modeling with crossed  
762 random effects for subjects and items. *Journal of Memory and Language*, 59(4), 390-  
763 412. doi: 10.1016/j.jml.2007.12.005
- 764 Baggio, G., & Hagoort, P. (2011). The balance between memory and unification in  
765 semantics: A dynamic account of the N400. *Language and Cognitive Processes*,  
766 26(9), 1338-1367. doi: 10.1080/01690965.2010.542671
- 767 Barr, D. J., Levy, R., Scheepers, C., & Tily, H. J. (2013). Random effects structure for  
768 confirmatory hypothesis testing: Keep it maximal. *Journal of Memory and Language*,  
769 68(3), 255-278. doi: 10.1016/j.jml.2012.11.001
- 770 Brothers, T., Swaab, T. Y., & Traxler, M. J. (2017). Goals and strategies influence lexical  
771 prediction during sentence comprehension. *Journal of Memory and Language*, 93,  
772 203-216. doi: 10.1016/j.jml.2016.10.002
- 773 Brown, C., & Hagoort, P. (1993). The processing nature of the n400: evidence from masked  
774 priming. *J Cogn Neurosci*, 5(1), 34-44. doi: 10.1162/jocn.1993.5.1.34
- 775 Brown, C. M., Hagoort, P., & Chwilla, D. J. (2000). An event-related brain potential analysis  
776 of visual word priming effects. *Brain Lang*, 72(2), 158-190. doi:  
777 10.1006/brln.1999.2284
- 778 Button, K. S., Ioannidis, J. P., Mokrysz, C., Nosek, B. A., Flint, J., Robinson, E. S., &  
779 Munafo, M. R. (2013). Power failure: why small sample size undermines the  
780 reliability of neuroscience. *Nat Rev Neurosci*, 14(5), 365-376. doi: 10.1038/nrn3475
- 781 Chwilla, D. J., Brown, C. M., & Hagoort, P. (1995). The N400 as a function of the level of  
782 processing. *Psychophysiology*, 32(3), 274-285.
- 783 Clark, A. (2013). Whatever next? Predictive brains, situated agents, and the future of  
784 cognitive science. *Behav Brain Sci*, 36(3), 181-204. doi:  
785 10.1017/S0140525X12000477
- 786 Clark, H. H. (1973). Language as Fixed-Effect Fallacy - Critique of Language Statistics in  
787 Psychological Research. *Journal of Verbal Learning and Verbal Behavior*, 12(4),  
788 335-359. doi: Doi 10.1016/S0022-5371(73)80014-3
- 789 Connolly, J. F., & Phillips, N. A. (1994). Event-related potential components reflect  
790 phonological and semantic processing of the terminal word of spoken sentences. *J  
791 Cogn Neurosci*, 6(3), 256-266. doi: 10.1162/jocn.1994.6.3.256
- 792 Dell, G. S., & Chang, F. (2014). The P-chain: relating sentence production and its disorders  
793 to comprehension and acquisition. *Philos Trans R Soc Lond B Biol Sci*, 369(1634),  
794 20120394. doi: 10.1098/rstb.2012.0394
- 795 DeLong, K. A. (2009). *Electrophysiological explorations of linguistic pre-activation and its  
796 consequences during online sentence processing* (Doctoral dissertation). Cognitive  
797 Science, UC San Diego. b6301658. Retrieved from:  
798 <http://escholarship.org/uc/item/4q7520sb>
- 799 DeLong, K. A., Urbach, T. P., & Kutas, M. (2005). Probabilistic word pre-activation during  
800 language comprehension inferred from electrical brain activity. *Nat Neurosci*, 8(8),  
801 1117-1121. doi: 10.1038/nn1504
- 802 DeLong, K.A., Urbach, T.P., & Kutas, M. (2017). Is there a replication crisis? Perhaps. Is this  
803 an example? No: a commentary on Ito, Martin, and Nieuwland (2016), *Language,  
804 Cognition and Neuroscience*, DOI: 10.1080/23273798.2017.1279339

- 805 Federmeier, K. D., & Kutas, M. (1999). A rose by any other name: Long-term memory  
806 structure and sentence processing. *Journal of Memory and Language*, 41(4), 469-495.  
807 doi: DOI 10.1006/jmla.1999.2660
- 808 Friederici, A. D., Steinhauer, K., & Frisch, S. (1999). Lexical integration: sequential effects  
809 of syntactic and semantic information. *Mem Cognit*, 27(3), 438-453.
- 810 Friston, K. (2005). A theory of cortical responses. *Philos Trans R Soc Lond B Biol Sci*,  
811 360(1456), 815-836. doi: 10.1098/rstb.2005.1622
- 812 Friston, K. (2010). The free-energy principle: a unified brain theory? *Nat Rev Neurosci*,  
813 11(2), 127-138. doi: 10.1038/nrn2787
- 814 Gelman A., Loken E. (2013). The garden of forking paths: Why multiple comparisons can be  
815 a problem, even when there is no “fishing expedition” or “p-hacking” and the research  
816 hypothesis was posited ahead of time. Retrieved from  
817 [http://www.stat.columbia.edu/~gelman/research/unpublished/p\\_hacking.pdf](http://www.stat.columbia.edu/~gelman/research/unpublished/p_hacking.pdf)
- 818 Gelman, A., & Stern, H. (2006). The difference between “significant” and “not significant” is  
819 not itself statistically significant. *The American Statistician*, 60(4), 328-331.
- 820 Hagoort, P. (2017). The core and beyond in the language-ready brain. *Neurosci Biobehav  
821 Rev*. doi: 10.1016/j.neubiorev.2017.01.048
- 822 Hauk, O., Davis, M. H., Ford, M., Pulvermuller, F., & Marslen-Wilson, W. D. (2006). The  
823 time course of visual word recognition as revealed by linear regression analysis of  
824 ERP data. *Neuroimage*, 30(4), 1383-1400. doi: 10.1016/j.neuroimage.2005.11.048
- 825 Huettig, F. (2015). Four central questions about prediction in language processing. *Brain Res*,  
826 1626, 118-135. doi: 10.1016/j.brainres.2015.02.014
- 827 Ito, A., Corley, M., Pickering, M. J., Martin, A. E., & Nieuwland, M. S. (2016). Predicting  
828 form and meaning: Evidence from brain potentials. *Journal of Memory and  
829 Language*, 86, 157-171. doi: 10.1016/j.jml.2015.10.007
- 830 Ito, A., Martin, A. E., & Nieuwland, M. S. (2017a). How robust are prediction effects in  
831 language comprehension? Failure to replicate article-elicited N400 effects. *Language,  
832 Cognition and Neuroscience*, 32(8), 954-965. doi: 10.1080/23273798.2016.1242761
- 833 Ito, A., Martin, A. E., & Nieuwland, M. S. (2017b). Why the A/AN prediction effect may be  
834 hard to replicate: a rebuttal to Delong, Urbach, and Kutas (2017). *Language,  
835 Cognition and Neuroscience*, 32(8), 974-983. doi: 10.1080/23273798.2017.1323112
- 836 Jackendoff, R. (2002). *Foundations of language : brain, meaning, grammar, evolution*.  
837 Oxford ; New York: Oxford University Press.
- 838 Jung, T. P., Makeig, S., Humphries, C., Lee, T. W., McKeown, M. J., Iragui, V., &  
839 Sejnowski, T. J. (2000). Removing electroencephalographic artifacts by blind source  
840 separation. *Psychophysiology*, 37(2), 163-178.
- 841 Kintsch, W. (1988). The role of knowledge in discourse comprehension: a construction-  
842 integration model. *Psychol Rev*, 95(2), 163-182.
- 843 Kutas, M., & Federmeier, K. D. (2011). Thirty years and counting: finding meaning in the  
844 N400 component of the event-related brain potential (ERP). *Annu Rev Psychol*, 62,  
845 621-647. doi: 10.1146/annurev.psych.093008.131123
- 846 Kutas, M., & Hillyard, S. A. (1980). Reading senseless sentences: brain potentials reflect  
847 semantic incongruity. *Science*, 207(4427), 203-205.
- 848 Kutas, M., & Hillyard, S. A. (1984). Brain potentials during reading reflect word expectancy  
849 and semantic association. *Nature*, 307(5947), 161-163.
- 850 Lau, E. F., Holcomb, P. J., & Kuperberg, G. R. (2013). Dissociating N400 Effects of  
851 Prediction from Association in Single-word Contexts. *J Cogn Neurosci*, 25(3), 484-  
852 502.

- 853 Lau, E., Namyst, A., Fogel, A., & Delgado, T. (2016). A direct comparison of N400 effects  
854 of predictability and incongruity in adjective-noun combination. *Collabra: Psychology*, 2(1).
- 855
- 856 Lau, E. F., Phillips, C., & Poeppel, D. (2008). A cortical network for semantics:  
857 (de)constructing the N400. *Nat Rev Neurosci*, 9(12), 920-933. doi: 10.1038/nrn2532
- 858 Luck, S. J., & Gaspelin, N. (2017). How to get statistically significant effects in any ERP  
859 experiment (and why you shouldn't). *Psychophysiology*, 54(1), 146-157.
- 860 Marslen-Wilson, W., & Tyler, L. K. (1980). The temporal structure of spoken language  
861 understanding. *Cognition*, 8(1), 1-71.
- 862 Martin, C. D., Thierry, G., Kuipers, J. R., Boutonnet, B., Foucart, A., & Costa, A. (2013).  
863 Bilinguals reading in their second language do not predict upcoming words as native  
864 readers do. *Journal of Memory and Language*, 69(4), 574-588. doi:  
865 10.1016/j.jml.2013.08.001
- 866 Miyamoto, K. (2016). Hemispheric Differences in Linguistic Prediction Given High  
867 Constraint Contexts. Presentation given at the Kutas Cognitive Electrophysiology Lab  
868 April 23rd, 2016. Retrieved from <https://www.slideshare.net/KianaMiyamoto/kutas-lab-latart>
- 869
- 870 Nieuwland, M. S. (2016). Quantification, prediction, and the online impact of sentence truth-  
871 value: Evidence from event-related potentials. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 42(2), 316-334.
- 872
- 873 Nieuwland, M. S., & Martin, A. E. (2012). If the real world were irrelevant, so to speak: The  
874 role of propositional truth-value in counterfactual sentence comprehension. *Cognition*,  
875 122(1), 102-109.
- 876 Nieuwland, M., Politzer-Ahles, S., Heyselaar, E., Segaert, K., Darley, E., Kazanina, N., ... &  
877 Mézière, D. (2017). Limits on prediction in language comprehension: A multi-lab  
878 failure to replicate evidence for probabilistic pre-activation of phonology. *BioRxiv*,  
879 111807.
- 880 Open Science, C. (2015). PSYCHOLOGY. Estimating the reproducibility of psychological  
881 science. *Science*, 349(6251), aac4716. doi: 10.1126/science.aac4716
- 882 O'Regan, K. (1979). Saccade size control in reading: evidence for the linguistic control  
883 hypothesis. *Percept Psychophys*, 25(6), 501-509.
- 884 Osterhout, L., & Holcomb, P. J. (1992). Event-Related Brain Potentials Elicited by Syntactic  
885 Anomaly. *Journal of Memory and Language*, 31(6), 785-806. doi: Doi 10.1016/0749-  
886 596x(92)90039-Z
- 887 Otten, M., Nieuwland, M. S., & Van Berkum, J. J. (2007). Great expectations: specific lexical  
888 anticipation influences the processing of spoken language. *BMC Neurosci*, 8, 89. doi:  
889 10.1186/1471-2202-8-89
- 890 Otten, M., & Van Berkum, J. (2008). Discourse-Based Word Anticipation During Language  
891 Processing: Prediction or Priming? *Discourse Processes*, 45(6), 464-496. doi: Pii  
892 90598764910.1080/01638530802356463
- 893 Otten, M., & Van Berkum, J. J. A. (2009). Does working memory capacity affect the ability  
894 to predict upcoming words in discourse? *Brain Res*, 1291, 92-101. doi:  
895 10.1016/j.brainres.2009.07.042
- 896 Pickering, M. J., & Clark, A. (2014). Getting ahead: forward models and their place in  
897 cognitive architecture. *Trends Cogn Sci*, 18(9), 451-456. doi:  
898 10.1016/j.tics.2014.05.006
- 899 Pickering, M. J., & Garrod, S. (2013). An integrated theory of language production and  
900 comprehension. *Behav Brain Sci*, 36(4), 329-347. doi: 10.1017/S0140525X12001495
- 901 Smith, N. J., & Levy, R. (2013). The effect of word predictability on reading time is  
902 logarithmic. *Cognition*, 128(3), 302-319. doi: 10.1016/j.cognition.2013.02.013

- 903 Steinhauer, K., Royle, P., Drury, J. E., & Fromont, L. A. (2017). The priming of priming:  
904 Evidence that the N400 reflects context-dependent post-retrieval word integration in  
905 working memory. *Neurosci Lett*, 651, 192-197. doi: 10.1016/j.neulet.2017.05.007
- 906 Summerfield, C., & de Lange, F. P. (2014). Expectation in perceptual decision making:  
907 neural and computational mechanisms. *Nat Rev Neurosci*, 15(11), 745-756. doi:  
908 10.1038/nrn3838
- 909 Tanner, D., Morgan-Short, K., & Luck, S. J. (2015). How inappropriate high-pass filters can  
910 produce artifactual effects and incorrect conclusions in ERP studies of language and  
911 cognition. *Psychophysiology*, 52(8), 997-1009. doi: 10.1111/psyp.12437
- 912 Taylor, W. L. (1953). "Cloze Procedure": A New Tool For Measuring Readability.  
913 *Journalism Quarterly*, 30(4), 415-433.
- 914 Van Berkum, J. J. (2010). The brain is a prediction machine that cares about good and bad-  
915 any implications for neuropragmatics?. *Italian Journal of Linguistics*, 22, 181-208.
- 916 Van Berkum, J. J., Brown, C. M., Zwitserlood, P., Kooijman, V., & Hagoort, P. (2005).  
917 Anticipating upcoming words in discourse: evidence from ERPs and reading times. *J  
918 Exp Psychol Learn Mem Cogn*, 31(3), 443-467. doi: 10.1037/0278-7393.31.3.443
- 919 Van Berkum, J. J., Hagoort, P., & Brown, C. M. (1999). Semantic integration in sentences  
920 and discourse: evidence from the N400. *J Cogn Neurosci*, 11(6), 657-671.
- 921 Van Petten, C., Coulson, S., Rubin, S., Plante, E., & Parks, M. (1999). Time course of word  
922 identification and semantic integration in spoken language. *Journal of Experimental  
923 Psychology-Learning Memory and Cognition*, 25(2), 394-417. doi: Doi  
924 10.1037/0278-7393.25.2.394
- 925 Van Petten, C., & Luka, B. J. (2012). Prediction during language comprehension: benefits,  
926 costs, and ERP components. *Int J Psychophysiol*, 83(2), 176-190. doi:  
927 10.1016/j.ijpsycho.2011.09.015
- 928 Vasishth, S., Mertzen, D., Jäger, L.A. & Gelman, A. (2018). The statistical significance filter  
929 leads to overoptimistic expectations of replicability. Psyarxiv
- 930 Wagenmakers, E. J., Verhagen, J., & Ly, A. (2016). How to quantify the evidence for the  
931 absence of a correlation. *Behav Res Methods*, 48(2), 413-426. doi: 10.3758/s13428-  
932 015-0593-0
- 933 Wicha, N. Y., Moreno, E. M., & Kutas, M. (2004). Anticipating words and their gender: an  
934 event-related brain potential study of semantic integration, gender expectancy, and  
935 gender agreement in Spanish sentence reading. *J Cogn Neurosci*, 16(7), 1272-1288.  
936 doi: 10.1162/0898929041920487
- 937 Wlotko, E. W., & Federmeier, K. D. (2012). So that's what you meant! Event-related  
938 potentials reveal multiple aspects of context use during construction of message-level  
939 meaning. *Neuroimage*, 62(1), 356-366. doi: 10.1016/j.neuroimage.2012.04.054
- 940 Wlotko, E. W., & Federmeier, K. D. (2015). Time for prediction? The effect of presentation  
941 rate on predictive sentence comprehension during word-by-word reading. *Cortex*, 68,  
942 20-32.
- 943 Yan, S., Kuperberg, G. R., & Jaeger, T.F. (2017). Prediction (Or Not) During Language  
944 Processing. A Commentary On Nieuwland et al.(2017) And Delong et al.(2005).  
945 *bioRxiv*, 143750.
- 946 Zwitserlood, P. (1989). The locus of the effects of sentential-semantic context in spoken-  
947 word processing. *Cognition*, 32(1), 25-64.
- 948
- 949
- 950

951 **FIGURE CAPTIONS**  
952

953 **Figure 1. Replication analysis.** Correlations between N400 amplitude and article/noun cloze  
954 probability per laboratory. N400 amplitude is the mean voltage in the 200-500 ms time  
955 window after word onset. A positive value corresponds to the canonical finding that N400  
956 amplitude became smaller (less negative—more positive) with increasing cloze probability.  
957 Here and in all further plots, negative voltages are plotted upwards. Upper graph: Scatter  
958 plots showing the correlation between cloze and N400 activity at electrode Cz, for each lab.  
959 The position of Cz and the other electrodes is displayed in the head plot in between the upper  
960 and lower graph. Lower graph: Scalp distribution of the *r*-values for each lab. Asterisks (\*)  
961 indicate electrodes that showed a statistically significant correlation (two-tailed  $p < 0.05$ , not  
962 corrected for multiple comparisons). Exact *r*- and *p*-values for each laboratory and EEG  
963 channel are available as source data (Figure 1-source data 1-4) and on <https://osf.io/eyzaq>.  
964

965 **Figure 2. Replication analysis.** Scalp distribution and *r*-values at each channel based on data  
966 pooled from all laboratories, using a 500 ms baseline correction procedure as used by  
967 DeLong et al (2005). Data were pooled after computing bin-averages per laboratory as in the  
968 original study, treating the laboratories as multiple observations of each bin-average.  
969 Asterisks (\*) indicate electrodes that showed a statistically significant correlation (two-tailed,  
970 not corrected for multiple comparisons). Exact *r*- and *p*-values for EEG channel are available  
971 as source data (Figure 2-source data 1-4).  
972

973 **Figure 3. Single-trial analysis.** Grand-average ERPs elicited by relatively expected and  
974 unexpected words (cloze higher/lower than 50%) and the associated difference waveforms  
975 (low minus high cloze) at electrode Cz. Dotted lines indicate 1 standard deviation above or  
976 below the grand average.  
977

978 **Figure 4. Single-trial analysis.** Relationship between cloze and ERP amplitude for articles  
979 and nouns in the N400 spatiotemporal window, as illustrated by the mean ERP values per  
980 cloze value (number of observations reflected in circle size), along with the regression line  
981 and 95% confidence interval. A change in article cloze from 0 to 100 is associated with a  
982 change in amplitude of 0.296  $\mu$ V (95% confidence interval: -.08 to .67). A change in noun-  
983 cloze from 0 to 100 is associated with a change in amplitude of 2.22  $\mu$ V (95% confidence  
984 interval: 1.75 to 2.69). The data for these analyses were pooled across all 9 labs.  
985

986 **Figure 5. Exploratory single-trial analyses.** The relationship between cloze and ERP  
987 amplitude as illustrated by the mean ERP values per cloze value (number of observations  
988 reflected in circle size), along with the regression line and 95% confidence interval, from two  
989 exploratory analyses. We performed a test which used a longer baseline time window (500  
990 ms, left panel) to better control for pre-article voltage levels. This test reduced the initially  
991 observed effect of article-cloze,  $\beta = .14$ , CI [-.25, .53],  $\chi^2(1) = 0.46$ ,  $p = .50$ . An analysis in  
992 the 500 to 100 ms time window *before* article-onset (right panel) revealed a non-significant  
993 effect of cloze that resembled the pattern observed *after* article-onset,  $\beta = .16$ , CI [-.07, .39],

994  $\chi^2(1) = 1.82, p = .18$ , shedding doubt on the conclusion that the observed results are due to  
995 the presentation of the articles.

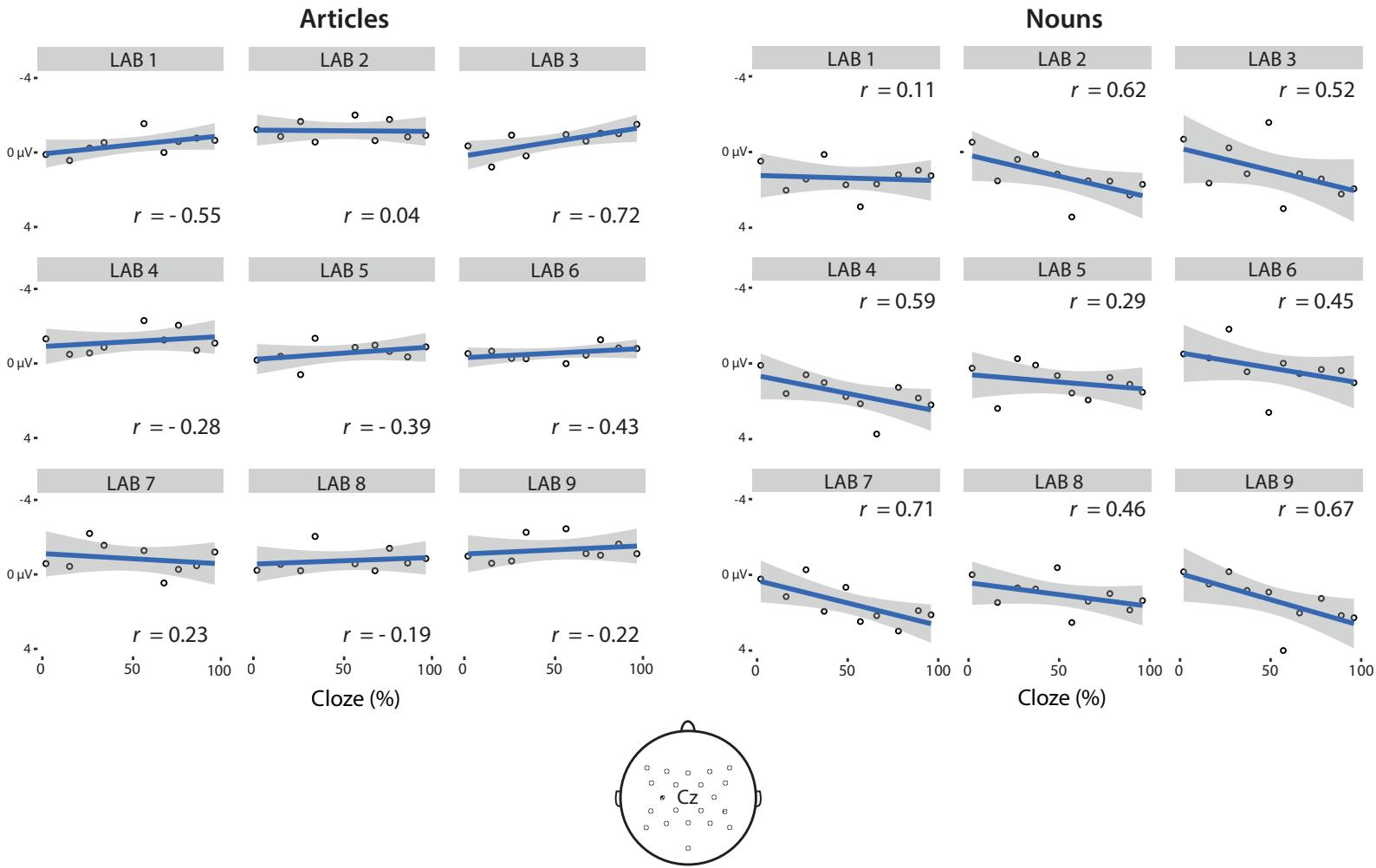
996  
997 **Figure 6. Exploratory replication Bayes factor analysis.** This analysis quantifies the  
998 obtained evidence for the null hypothesis ( $H_0$ ) that N400 is not impacted by cloze, or for the  
999 alternative hypothesis ( $H_1$ ) that N400 is impacted by cloze with the direction *and* size of  
1000 effect reported by DeLong et al. Scalp maps show the common logarithm of the replication  
1001 Bayes factor for each electrode, capped at log(100) for presentation purposes. Electrodes that  
1002 yielded at least moderate evidence for or against the null hypothesis (Bayes factor of  $\geq 3$ ) are  
1003 marked by an asterisk. At posterior electrodes where DeLong et al. found their effects, our  
1004 article data yielded strong to extremely strong evidence for the null hypothesis, whereas our  
1005 noun data yielded extremely strong evidence for the alternative hypothesis (upper graphs).  
1006 These results were obtained with the procedure described in DeLong et al. (no baseline  
1007 correction), and with a 500 ms pre-word baseline correction (lower graphs), the procedure  
1008 later described by DeLong and colleagues.  
1009

1010 **Figure 7. Exploratory Bayesian mixed-effects model analyses.** Posterior distributions for  
1011 the effect of cloze on ERP amplitudes in the N400 window. The x-axis shows cloze effect  
1012 sizes (i.e., changes in microvolts associated with an increase from 0% cloze probability to  
1013 100% cloze probability). The black line indicates the posterior distribution of effects; higher  
1014 values of the posterior density at a given effect size indicate higher probability that this is the  
1015 true effect size in the population. The peak of the posterior distribution roughly corresponds  
1016 to the point estimate of the effect size (the regression coefficient) fitted from the Bayesian  
1017 mixed effect model, i.e., the most likely value of the true effect size. The middle 95% of the  
1018 posterior distribution, shaded in pink, corresponds to a two-tailed 95% credible interval for  
1019 the effect size—i.e., an interval that we can be 95% confident contains the true effect. The  
1020 green dotted line indicates the prior distribution (i.e., our expectation about where the true  
1021 effect would lie before the data were collected). For the articles, this prior is centred on  
1022 1.25 $\mu$ V, an approximation of the effect observed by Delong and colleagues (2005), and for  
1023 the nouns it is centred on 3.5 $\mu$ V. The black connected dots illustrate the ratio between the  
1024 posterior and prior distribution (i.e., the Bayes Factor) at the effect size of 0 $\mu$ V; for example,  
1025 a Bayes Factor of 4 suggests we can be 4 times more certain that the true effect is zero after  
1026 having conducted this experiment than before, or, in other words, that the data increased our  
1027 confidence in the null effect of zero fourfold. We performed these analyses for each of the  
1028 linear mixed-effects model analyses we performed. We note that in all the article-analyses,  
1029 the posterior probability of the estimated effect being greater than zero is around 80 or 90%,  
1030 although this is also true for the pre-stimulus variable, shedding doubt that the observed  
1031 results are due to presentation of the articles. In none of our article-analyses did zero lie  
1032 outside the obtained credible interval, whereas for the nouns, zero lay outside the credible  
1033 interval. These results are consistent with a failure to replicate the size of the article-effect  
1034 reported by DeLong et al. and a successful replication of the noun-effect.  
1035

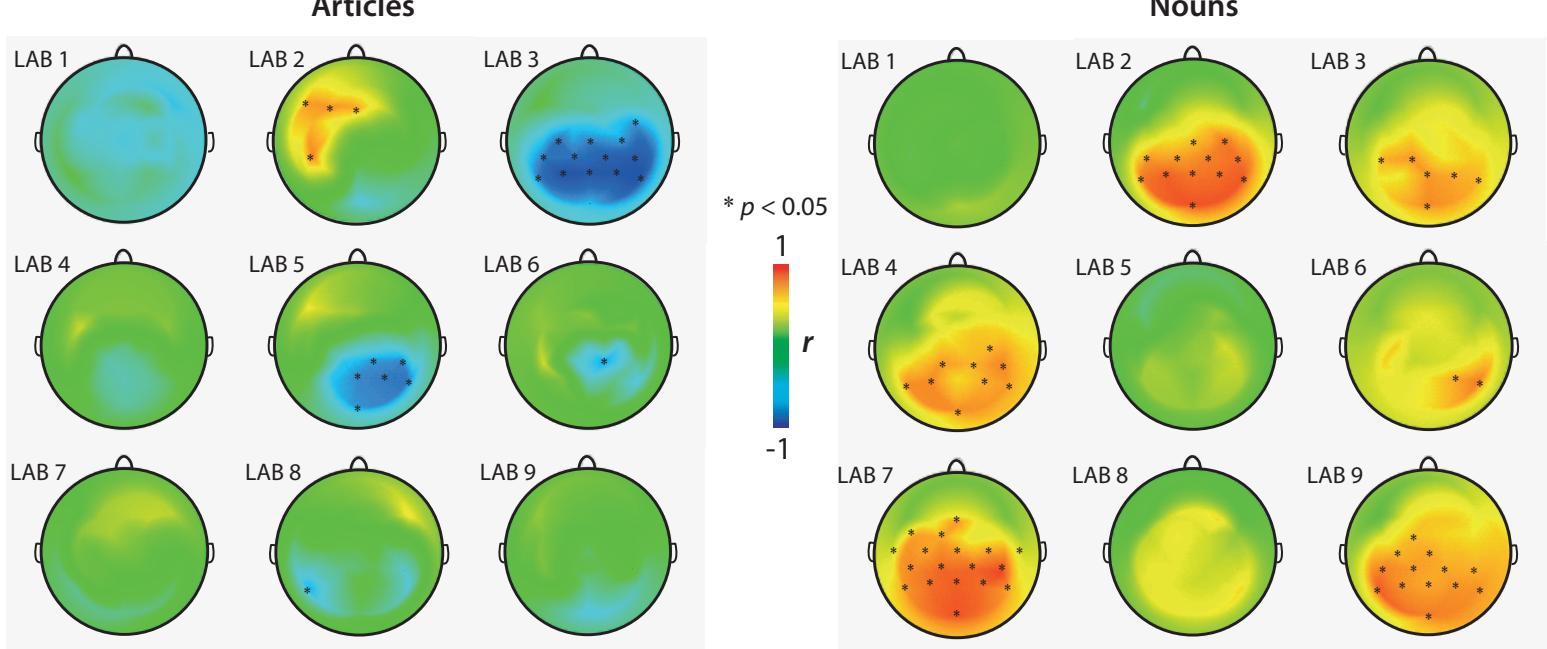
1036 **Figure 8. Control experiment.** P600 effects at electrode Pz per lab associated with flouting  
1037 of the English a/an rule. Plotted ERPs show the grand-average difference waveform and

1038 standard deviation for ERPs elicited by ungrammatical expressions ('an kite') minus those  
1039 elicited by grammatical expressions ('a kite').  
1040  
1041

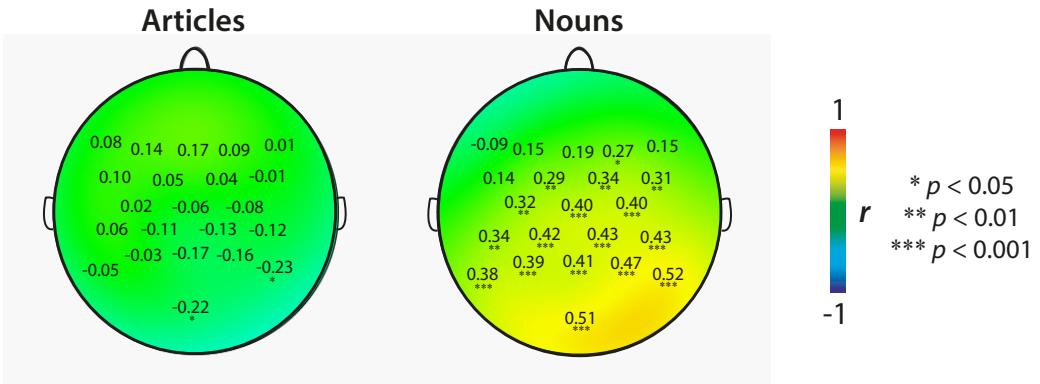
## Replication analysis: Correlation results at Cz



**r-values at all electrode positions**



### r-values at all electrode positions

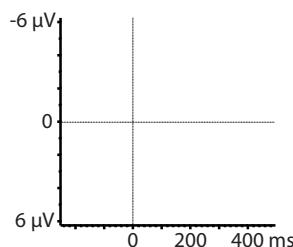
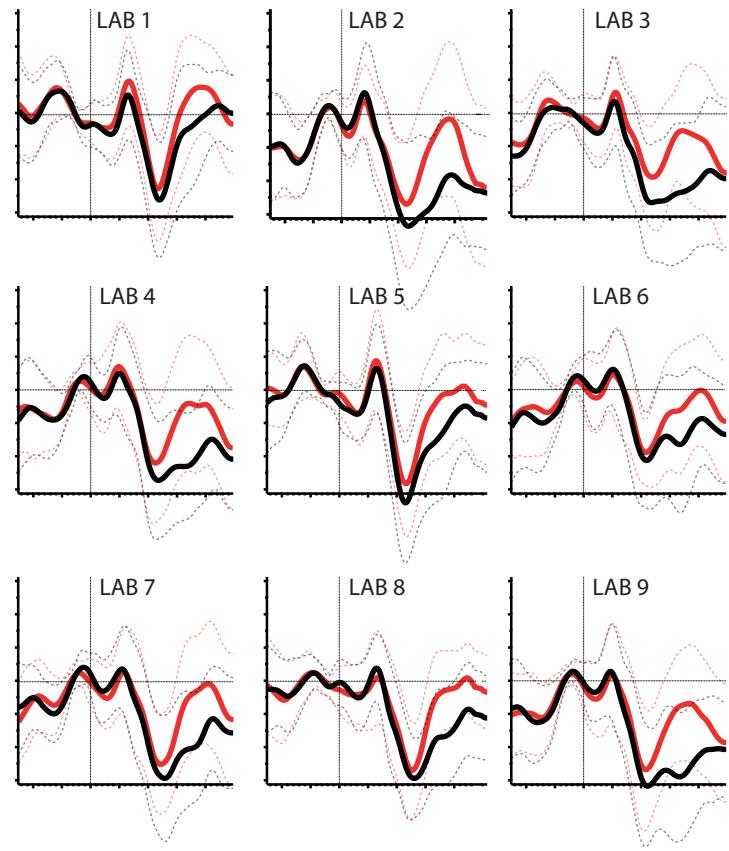
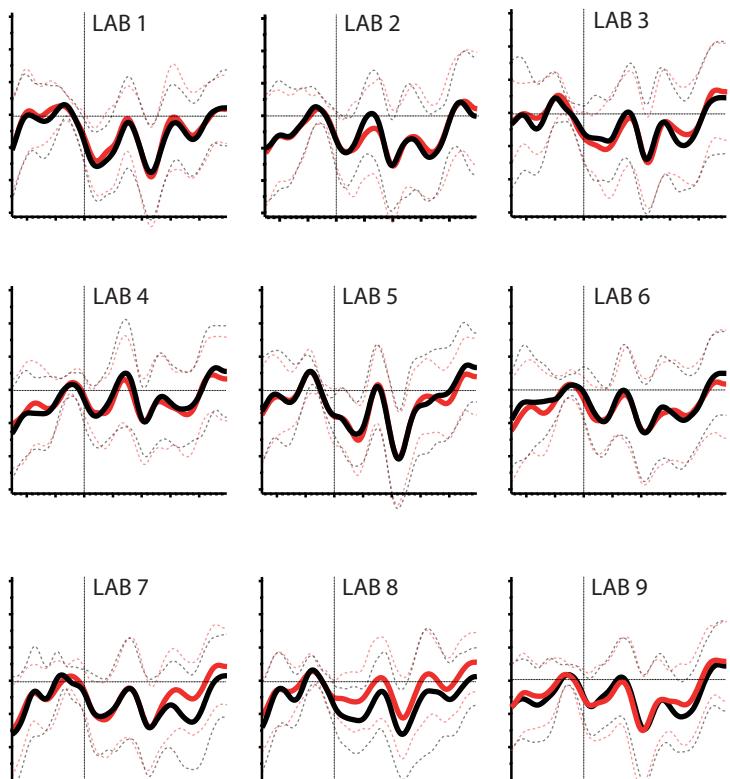


**Single-trial analysis:  
Grand-average ERPs (Cz)**

**Articles**

— Low Cloze  
— High Cloze

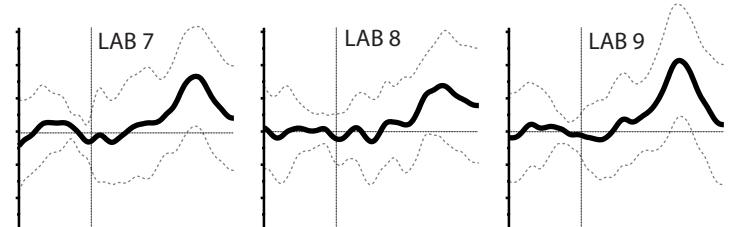
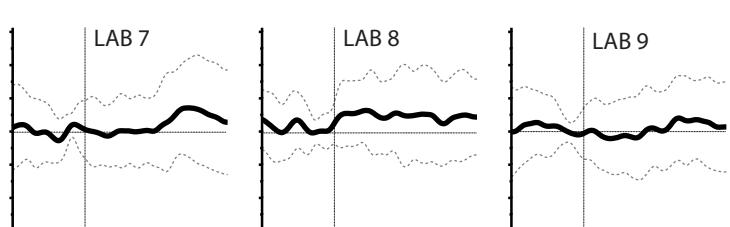
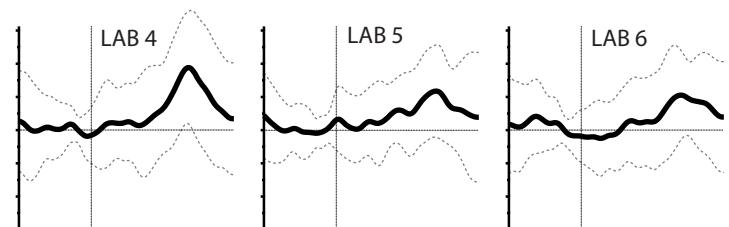
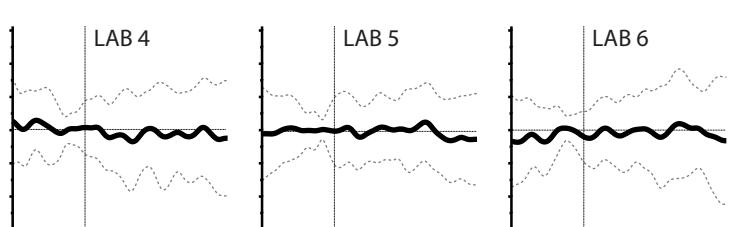
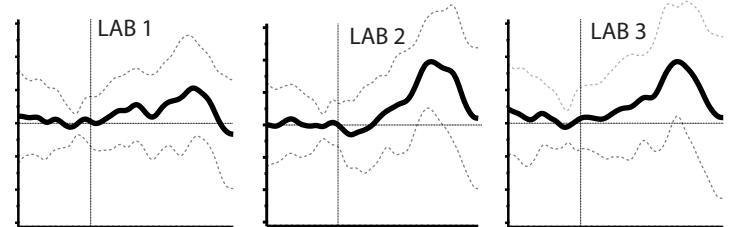
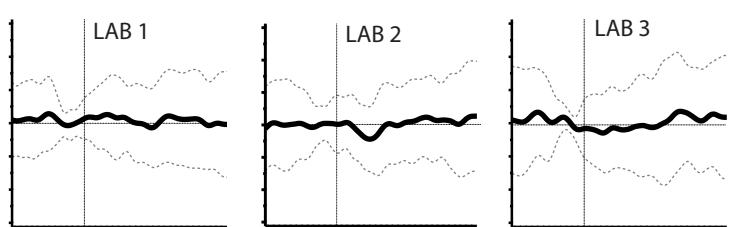
**Nouns**

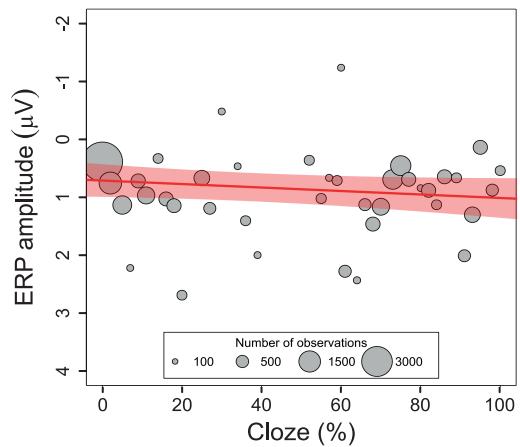
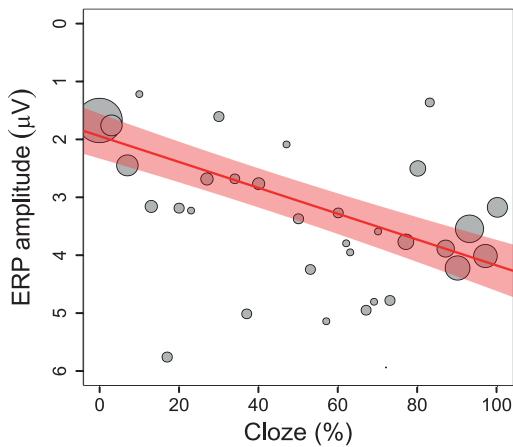


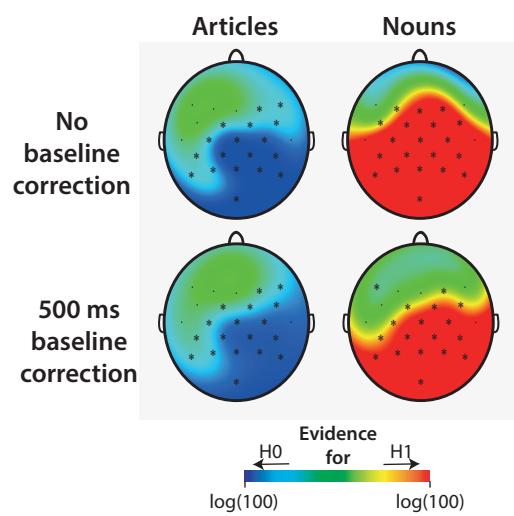
**Articles**

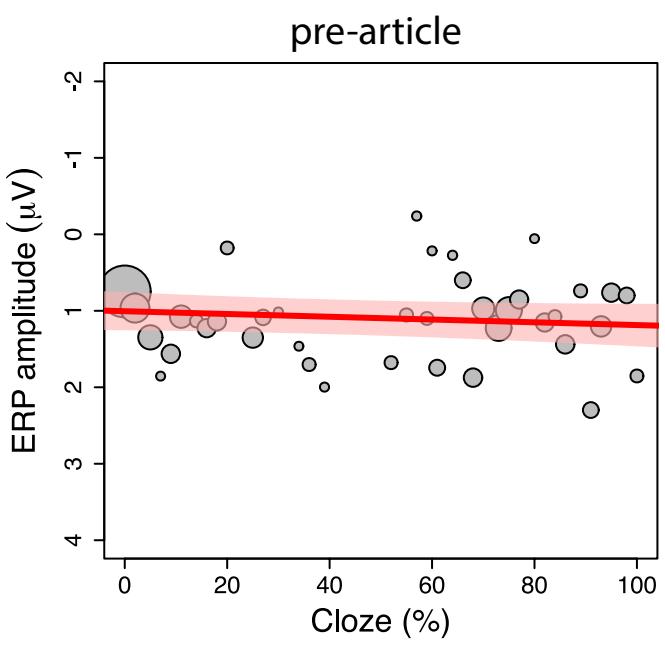
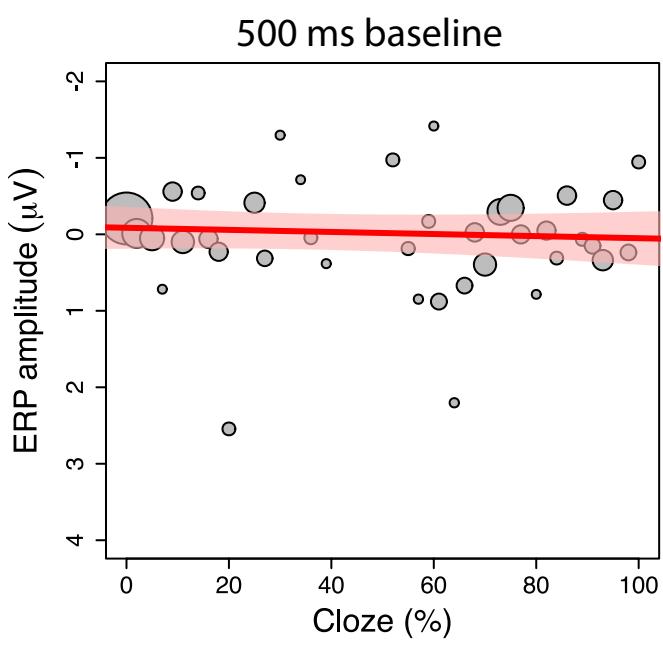
**Difference ERPs**

**Nouns**



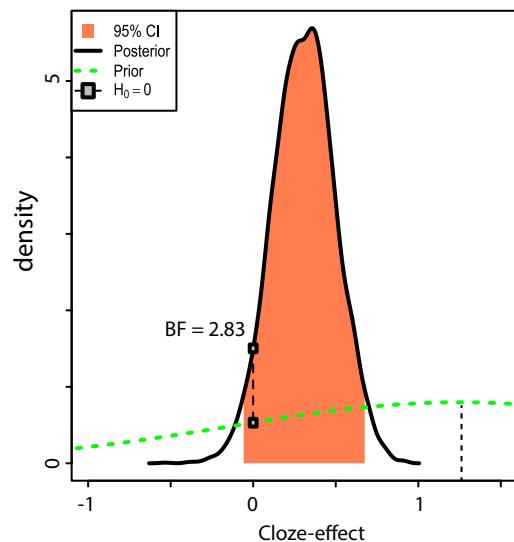
**Articles****Nouns**



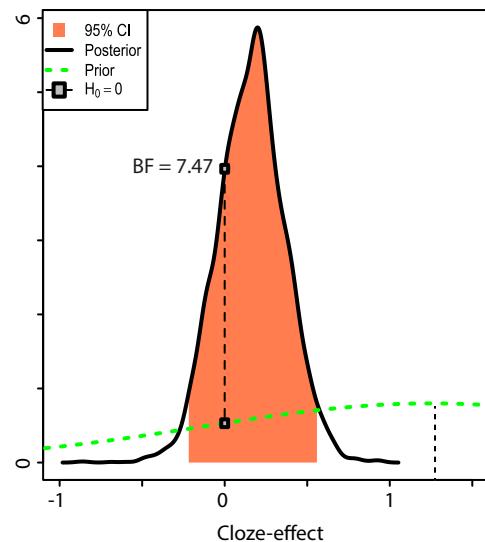


## ARTICLES

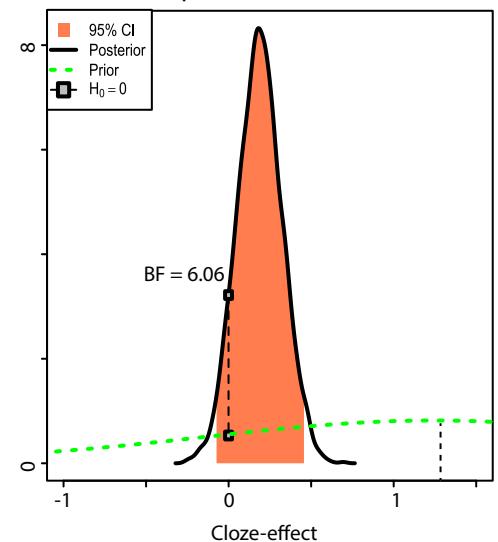
100 ms baseline



500 ms baseline

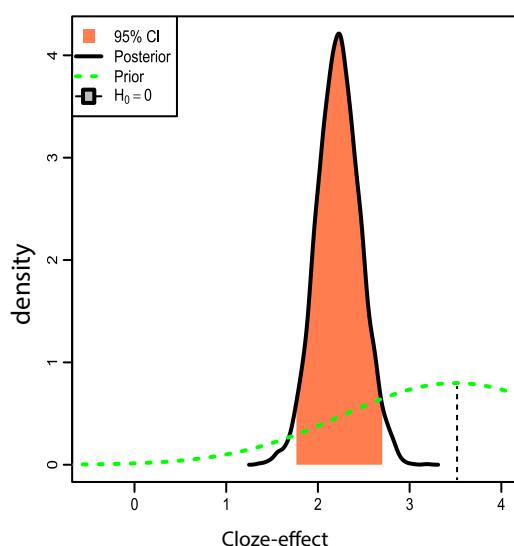


pre-article



## NOUNS

100 ms baseline



**Control experiment:**  
**Ungrammatical - Grammatical 'P600 effect' (Pz)**

