

# An exploratory data analysis of word form prediction during word-by-word reading

# Thomas P. Urbach<sup>a,1</sup>, Katherine A. DeLong<sup>a</sup>, Wen-Hsuan Chan<sup>a</sup>, and Marta Kutas<sup>a,b</sup>

<sup>a</sup>Department of Cognitive Science, University of California San Diego, La Jolla, CA, 92037; and <sup>b</sup>Department of Neurosciences, University of California San Diego, La Jolla, CA, 92037

Edited by Gary S. Dell, University of Illinois at Urbana-Champaign, Champaign, IL, and approved July 13, 2020 (received for review December 22, 2019)

In 2005, we reported evidence indicating that upcoming phonological word forms-e.g., kite vs. airplane-were predicted during reading. We recorded brainwaves (electroencephalograms [EEGs]) as people read word-by-word and then correlated the predictability in context of indefinite articles that preceded nouns (a kite vs. an airplane) with the average event-related brain potentials (ERPs) they elicited [K. A. DeLong, T. P. Urbach, M. Kutas, Nat. Neurosci. 8, 1117-1121 (2005)]. Amid a broader controversy about the role of word-form prediction in comprehension, those findings were recently challenged by a failed putative direct replication attempt [M. S. Nieuwland et al., eLife 7, e33468 (2018); nine labs, one experiment, and 2.6e4 observations]. To better understand the empirical justification for positing an association between prenominal article predictability and scalp potentials, we conducted a wide-ranging exploratory data analysis (EDA), pooling our original data with extant data from two followup studies (one lab, three experiments, and 1.2e4 observations). We modeled the time course of article predictability in the single-trial data by fitting linear mixed-effects regression (LMER) models at each time point and scalp location spanning a 3-s interval before, during, and after the article. Model comparisons based on Akaike information criteria (AIC) and slope-regression ERPs [rERPs; N. J. Smith, M. Kutas, Psychophysiology 52, 157-168 (2015)] provide substantial empirical support for a small positive association between article predictability and scalp potentials approximately 300 to 500 ms after article onset, predominantly over bilateral posterior scalp. We think this effect may reasonably be attributed to prediction of upcoming word forms.

language | prediction | EEG | rERP | EDA

P sycholinguistic theories of language comprehension gener-ally endorse the near-immediate, "incremental" construction of structured representations of meaning, as words, phrases, sentences, and discourses rapidly unfold over time (1). New information must be integrated with this evolving semantic representation, and some accounts further posit predictive or preparatory mechanisms that facilitate processing and help the system keep up with the input (2-4). The hypothesis that the comprehension system actively predicts is difficult to test experimentally. The challenge is to find evidence of predictive processing that cannot plausibly be attributed to rapid integration. For instance, given a sentence context like, The day was breezy so the boys went outside to fly ..., knowledge of the world and English make some continuations more predictable (a kite) and others less so (an airplane). It is possible that the supporting context leads the processor to predict (anticipate or expect) the word *kite* before it arrives, in which case online measures sensitive to experimental manipulations of processing difficulty-e.g., self-paced reading times, eye movements, event-related brain potentials (ERPs), and event-related magnetic fields (ERFs)-might show an experimental effect in the expected direction-i.e., faster reading times, shorter gaze durations, or reduced negative deflection of event-related potential around 400 ms (N400) ERPs/ERFsfor kite vs. airplane. However, if the effects observed at these nouns could, with equal justification, be attributed to violated

predictions or integration difficulty (or both), these findings are compatible with, but do not constitute strong evidence for, prediction, and parsimony favors integration mechanisms alone, which are necessary on any account.

The crux of the experimental challenge is time: Strong tests that information is *pre*-dicted come from measurements made before it actually arrives. Seminal laboratory studies measuring eye movements while listening to meaningful sentences in a controlled visual environment (5-7) found that people tended to glance at mentioned objects quickly or even prior to hearing a likely word, indicating rapid language-driven anticipation of upcoming semantic or conceptual content. To date, the clearest evidence for prediction of specifically linguistic information comes from paradigms that recruit sequential dependencies, wherein one type of grammatical element, such as a word or morphological marking, regularly co-occurs with another element. The seminal ERP studies (8, 9), were conducted by Wicha, Bates, Moreno, and Kutas using grammatical gender agreement between indefinite articles and nouns in Spanish—e.g., feminine una canasta ("a basket") vs. masculine un costal ("a sack"). If a Spanish sentence is likely to continue about a basket, the corresponding indefinite article is likely to be *una*, not *un*, and vice versa if the likely continuation is about a sack. Since the two forms of the indefinite article have the same meaning ("some singular thing"), they should be equally easy or difficult to integrate. Wicha et al. (8, 9) recorded electrical brain potentials at the scalp (electroencephalograms [EEGs]) as people read sentences word-by-word on a computer screen and found small

# Significance

Complex biological systems do not merely react, they anticipate. In 2005, the human-language comprehension system was considered an exception. We concluded not, based on our recordings of electrical brain activity measured before the critical words arrived during sentence reading, described in a now widely cited report. This, and the emergence of the "statistical crisis" in psychology, led to a large-scale replication attempt that failed. This prompted us to revisit the issue by analyzing our original data and two replication extensions with an exploratory data analysis (EDA) approach, enabled by advances in scientific computing technology. Our original conclusion was supported: Brains can anticipate specific upcoming words. We offer this as a case study in EDA for cognitive neurophysiology, more generally.

Author contributions: T.P.U., K.A.D., and M.K. designed research; K.A.D. performed research; T.P.U. and W.-H.C. contributed new reagents/analytic tools; T.P.U. and W.-H.C. analyzed data; T.P.U. wrote the paper; T.P.U., K.A.D., W.-H.C., and M.K. contributed substantive editorial revisions; and T.P.U. suggested the design of experiment 1.

The authors declare no competing interest.

This article is a PNAS Direct Submission.

This article contains supporting information online at https://www.pnas.org/lookup/suppl/doi:10.1073/pnas.1922028117/-/DCSupplemental.

Published under the PNAS license

<sup>&</sup>lt;sup>1</sup>To whom correspondence may be addressed. Email: turbach@ucsd.edu.

differences between the average ERPs elicited by articles that were compatible vs. incompatible with the grammatical gender of the likely continuation. These effects varied with the particulars of the experimental design: Incompatible articles elicited an N400-like relative negativity when the referent of the likely noun was depicted with a line drawing (8, 10) and a relative positive deflection around 500 to 700 ms when the continuations were orthographic words (9). With other lexical variables controlled by the experimental design, the difference between *un* and *una* is plausibly attributed to a mismatch between the grammatical gender of the article and the gender of the likeliest continuation, indicating that the continuation had been predicted before it was encountered.

Subsequent studies have used related sequential dependency designs to probe other languages for evidence of prediction, e.g., via case-marking in Dutch (11), grammatical gender in Dutch (refs. 12-14, but see ref. 15), Polish (16), and German (17). For these types of experimental designs, the nature of the linguistic dependency constrains the inferences that can be drawn about what information is anticipated (discussed in refs. 3 and 18). English does not mark grammatical gender or case agreement on nouns, but does attest a phonological dependency between alternate forms of the indefinite article a, which precedes consonant-sound-initial words, and an, which precedes vowel-sound-initial words: a kite vs. an airplane. We recruited this sequential dependency in previous work (ref. 19; hereafter, DUK05), recording scalp potentials while people read sentences like, The day was breezy so the boys went outside to fly [a kite/an airplane] in the park, one word at a time on a computer screen. We observed a positive correlation between the predictability, in context, of the indefinite articles that preceded the nouns a kite vs. an airplane and the average ERPs they elicited 200 to 500 ms over bilateral central and posterior scalp. Since the a/analternation depends on the initial speech sound of the next word, we took the systematic association between the ERP amplitude and offline article cloze probability to suggest "that individuals can use linguistic input to pre-activate representations of upcoming words in advance of their appearance" (ref. 19, p. 1119), and "Our observation of an ERP expectancy effect at the article leads us to conclude that predictions can be for specific phonological forms-words beginning with either vowels or consonants. In this sense, we propose that prediction can be highly specific, at least under some circumstances" (ref. 19, pp. 1119-1120).

Controversy has emerged recently regarding the strength of evidence for word-form prediction in variations of the a/an design. For instance, we did not observe the effect in younger adults with sentences at a faster presentation rate (ref. 20, experiment 2, 3.3 words per second) or in older adults at two words per second (21), and other groups have reported statistically reliable (22), marginal (ref. 23, experiment 2), and null results (ref. 23, experiment 1). A recent large-scale study by Nieuwland et al. (ref. 24; hereafter, NIET18) proposed to resolve the question by reusing the experimental materials and design of the original DUK05 a/an study (healthy younger adults reading two words per second in central vision) and analyzing EEG data collected from nine laboratories around Great Britain. That report makes four main points: 1) It is important to replicate experimental findings; 2) the prenominal article correlation with grand average ERPs reported in DUK05 could be a spurious statistical result; 3) with the same stimuli, generally similar procedures, more participants (n = 338), and more appropriate statistical analyses, they failed to observe a reliable effect at the prenominal article with either the potentially problematic average ERP correlation analysis or planned and posthoc single-trial, linear mixed-effect regression (LMER) model analyses; and 4) if there is such an effect, it is relatively small. We concur. The value of replication is uncontroversial, although rather than simply running the same experiment over and over, there may be more to learn from replication and extension, as illustrated by the followup studies DeLong conducted in the laboratory between 2005 and 2010 and that we have analyzed anew for this report (SI Appendix, Table S1, data acquisition and reporting timeline and references therein). We recognize the limitations of inferences drawn from correlations between averages and, thus, analyze single-trial EEG data with LMER models for this report. It is also clear that NIET18 failed to observe an effect of prenominal article predictability with the preregistered LMER analysis of scalp potentials averaged across six scalp locations and a 300ms poststimulus interval. However, when the existence of such an effect is in question, there seems little reason to suppose that the most informative general answer is to be had by selecting one temporal interval and a small set of scalp locations in advance and drawing inferences about what is or is not going on throughout the brain, as comprehension processes evolve from the analysis of this aggregated snapshot. In what follows, we propose alternatives that build on the strengths of the NIET18 analysis and aim to overcome some of its limitations.

The key empirical premise in the argument for word-form prediction based on the *a*/*an* experimental design is that indefinite article predictability, operationalized as cloze probability, is positively associated with the amplitude of scalp potentials elicited by the articles around 400 ms poststimulus over central and posterior scalp-i.e., that article N400 ERP amplitude correlates inversely with cloze probability. Accordingly, we investigated this association in three EEG datasets recorded in a/an-design experiments previously conducted in our laboratory: the original DUK05 experiment and two replication-extension experiments that revised and extended the stimulus materials and experimental conditions. In all three experiments, healthy young adults read sentences two words per second in central vision as in the original DUK05 report and NIET18. In contrast with the absence of evidence reported in NIET18, our exploratory LMER modeling of the single-trial EEG data moment-by-moment at 26 scalp locations finds empirical support for the hypothesized association, which, in turn, may reasonably be attributed to prediction of upcoming word forms.

Exploratory EEG Data Analysis with Regression ERPs. The data from these three experiments have already been analyzed in a number of other ways, published and unpublished (SI Appendix, Table S1), and the results are known. These circumstances rightly prompt concern about circular analyses, multiple comparisons, and *p*-hacking when choosing which and how among the many available hypotheses to test with confirmatory null hypothesis tests (e.g., refs. 25–28). Since accept-or-reject-at- $\alpha$  confirmatory null-hypothesis testing is not appropriate, we present a series of data-driven exploratory analyses, along with what Tukey terms rough confirmatory assessments of strength of evidencei.e., a flexible data investigation in the sense he contrasts with the rigid steps of *data processing* and confirmatory hypothesis tests (29-31). Consequently, in concept and execution, the analyses reported herein have more in common with the iterative phases of model development, diagnosis, evaluation, and selection found in applied statistical modeling than in boilerplate data processing that passes from EEG recordings to results through a predetermined sequence of steps and declares victory by rejecting (or failing to reject; refs. 15, 23, and 24) a null hypothesis at P < 0.05. Researchers intrigued or outraged by this approach will find an engaging manifesto in Tukey's "Badmandments" (ref. 32, prologue), a clear overview for psychologists in Behrens (33), and methodological guidance in standard texts, e.g., Cohen, et al. (ref. 34, chapters 4 and 10), Fox (ref. 35, Data Craft: chapters 2-4), and Kutner et al. (ref. 36, chapters 9 and 10 and figure 9.1).

Our exploratory analyses used the same class of LMER models as NIET18 and differed primarily in that we evaluated a greater variety of models and modeled the data at a higher spatial and temporal resolution in the regression ERP (rERP) framework recently described and motivated by Smith and Kutas (ref. 37 and references therein for related approaches). For these analyses, we sweep an LMER model across the single-trial EEG and fit the data for all subjects and items at each time point of the digital recording. As Smith and Kutas point out, modeling the EEG data in this manner is a generalization of conventional sum-and-divide time-domain averaging. For a set of nsingle-trial EEG epochs (segments of the recording), each timealigned to an experimental event of interest, the time-domain average  $ERP(t) = \frac{1}{n} \sum_{i=1}^{n} EEG_i(t)$  at time t is mathematically identical to the estimated intercept,  $\hat{\beta}_0$ , of an intercept-only linear model of the same data,  $EEG(t) = \beta_0 + \epsilon$ , fit by ordinary least-squares regression. This means plotting, measuring, analyzing, and interpreting time-domain average ERP waveforms and the time series of estimated linear model intercepts,  $\hat{\beta}_0(t)$ , are literally one and the same. This approach generalizes to more complex models, notably, multiple-regression models that may include continuous and categorical predictor variables, and other classes of models including linear mixedeffects models. For models with multiple predictor variables, e.g.,  $EEG(t) = \beta_0 + \beta_1 X_1 + \dots + \beta_J X_J + \epsilon$ , fitting the model yields a time series of estimated coefficients,  $\hat{\beta}_{i}(t)$ , for each regressor,  $X_i$ , the waveforms that Smith and Kutas dubbed rERPs. Furthermore, besides the estimated model parameters, fitting a model at each time point also yields the corresponding time series of residual errors and derived quantities, such as error variance, coefficient SEs and CIs, and goodness-of-fit measures. Modeling time-series data is nothing new; the key insight of the rERP framework is that the logic of conventional event-related timedomain averaging extends to event-related time-domain modeling more generally, and thereby to the investigation of eventrelated brain activity by methods and procedures from applied statistical data modeling developed to fit, diagnose, compare, and interpret different models. The endgame is to determine which model(s), among the many possible, are likely to better or best account for systematic relationships between predictor and response variables, i.e., between experimental variables and event-related brain activity. Determining the existence and form of these associations is the first (though not last) step in causal inference.

# Approach

To investigate the association, if any, between the predictability of articles and the brain responses they elicit during word-byword reading, we swept LMER models across single-trial EEG recordings before, during, and after the onset of articles that vary in cloze probability. We make inferences based on the time course and scalp distribution of model goodness-of-fit measures and rERPs. Details and further discussion appear in *Materials and Methods* and *SI Appendix*. The analysis reproduction recipe, open-source scripts, and additional figures are available online at the Open Science Foundation (OSF): UDCK (https://osf.io/ tksur/) (38).

**EEG Data: Three Experiments.** After the original study reported in DUK05, DeLong and colleagues continued to investigate aspects of predictive processing in younger and older adults. For this report, we selected two additional studies conducted between 2005 and 2010 that incorporated the *a/an* prenominal indefinite article manipulation and extended the original study design with additional conditions and materials (see *SI Appendix*, Table S1 for a summary and references). The rationale for selecting these particular studies is that they tested healthy young adults reading two words per second in central vision, which affords a close comparison between and across the original DUK05 and NIET18 studies. Furthermore, the additional materials devel-

#### Table 1. EEG experiment participants, items, and article cloze

E	Ρ	I	Ν	М	SD	Range
1	32	80	2,136 (0.16)	0.38	0.35	0.0 to 0.97
2	32	160	4,668 (0.07)	0.44	0.41	0.0 to 1.0
3	24	240	5,232 (0.08)	0.39	0.38	0.0 to 1.0
All	88	320 <sup>†</sup>	1,2043 (0.10)	0.408	0.389	0.0 to 1.0

E, EEG experiment. I, number of items in the experimental design for modeling item as a random variable. Each item corresponds to the context prior to the critical article and provides one cloze value for a and one for an (see *SI Appendix* for article-cloze distributions and data exclusions). N, number of single trials analyzed after excluding EEG artifacts (proportions in parentheses) and stimulus irregularities (0.01). P, number of participants. The observed article-cloze mean (*M*) and SD (*SD*) on each row are computed for the single-trial data on that row and may be used to transform estimated regression coefficients for standardized article cloze back to the original cloze scale of 0 to 1. <sup>†</sup>Experiment 3 used 160 of the same prearticle item contexts as experiment 2 and added 80 new ones, 80 + 160 + 80 = 320 distinct items. Modeling item random effects takes this into account (*SI Appendix, Stimulus and Item Coding*).

oped by revising and extending the DUK05 materials fill in gaps in the distribution of contextually supported noun and the corresponding prenominal article cloze values in the DUK05 materials. This makes the pooled datasets appropriate for modeling article-cloze probability as a continuous predictor. So, for this report, we pooled the data from these three studies and modeled approximately 12,000 single-trial epochs (Table 1), recorded at 26 scalp locations spanning the interval from about 1.5 s before to 1.5 s after the critical article (*Materials and Methods* and *SI Appendix*, *EEG Experimental Procedures*).

Modeling: Linear Mixed-Effects rERPs. To characterize the time course and scalp distribution of article-cloze effects in the rERP framework, we swept each of the LMER models in Table 2 across the single-trial EEG data and computed the lme4::lmer() profiled maximum likelihood (ML) fit for the 1.2e4 observations at each time point and each channel (39). For exposition, Table 2 presents the models in the formula language of lme4, which specifies LMER models in two parts: the "fixed effect" predictor terms and the "random effect" terms enclosed in parentheses. This syntax aligns with a matrix equation specification of the model,  $y = X\beta + Zb + \epsilon$ , that shows the observed response variable y modeled in two parts as the sum of  $\beta$ -weighted regressors for fixed effects  $(X\beta)$  and b-weighted regressors for random effects (*Zb*). For an introduction to LMER modeling in psychology experiments, see the development of equation 9 in ref. 40, and see ref. 39 for a formal treatment of the model and fitting algorithms.

To highlight the approach in this report, we can unpack  $X\beta$  as the column vectors,  $X = [1, x_{cloze}]$ , a column of ones and the per-item article cloze values, and the scalar coefficients,  $\beta = [\beta_0, \beta_{cloze}]$  for the intercept and article cloze:

$$EEG = \beta_0 \mathbf{1} + \beta_{cloze} \mathbf{x}_{cloze} + \mathbf{Z}\mathbf{b} + \boldsymbol{\epsilon}.$$
 [1]

The analyses that follow map neatly onto the terms of Eq. 1. First, to select random effects for subjects, items, and experiments, we compared models with different Zb (Fig. 1). Second, to evaluate evidence for an association between article cloze and scalp potentials, we compared (full) models like Eq. 1 that include the article-cloze regressor,  $x_{cloze}$ , with corresponding (reduced) models that do not (Fig. 2). Third, the *LMER ERP* (ImerERP) waveforms are the estimated coefficients for the intercept,  $\hat{\beta}_0$ , and article cloze  $\hat{\beta}_{cloze}$  over time for each EEG channel (Fig. 3).

Fo	rm	ul	а
----	----	----	---

Random effects	
Maximal	
M0	cloze + (cloze   expt) + (cloze   subject) + (cloze   item)
Drop 1 slope	
M1	cloze + (cloze   expt) + (cloze   subject) + (1   item)
M2	cloze + (cloze   expt) + (1   subject) + (cloze   item)
M3	cloze + (1   expt) + (cloze   subject) + (cloze   item)
Drop 2 slopes	
M4	cloze + (cloze   expt) + (1   subject) + (1   item)
M5	cloze + (1   expt) + (cloze   subject) + (1   item)
M6	cloze + (1   expt) + (1   subject) + (cloze   item)
Drop 3 slopes	
M7	cloze + (1   expt) + (1   subject) + (1   item)
Drop 1 random term	
M8	cloze + (1   subject) + (1   item)
M9	cloze + (1   expt) + (1   subject)
M10	cloze + (1   expt) + (1   item)
Article cloze fixed-effect	
comparisons	
KIM	
M5	cloze + (1   expt) + (cloze   subject) + (1   item)
M5r	(1   expt) + (cloze   subject) + (1   item)
KIP	
M7	cloze + (1   expt) + (1   subject) + (1   item)
M7r	(1   expt) + (1   subject) + (1   item)
Experiment as a fixed effect	
KIM	
M11	cloze + expt + (cloze   subject) + (1   item)
M11r	expt + (cloze   subject) + (1   item)
KIP	
M12	cloze + expt + (1   subject) + (1   item)
M12r	expt + (1   subject) + (1   item)
Experiments 1, 2, and 3	
modeled separately	
KIM	
M13	cloze + (cloze   subject) + (1   item)
M13r	(cloze   subject) + (1   item)
KIP	
M14	cloze + (1   subject) + (1   item)
M14r	(cloze   subject) + (1   item)

Fixed and random intercepts are implicit and modeled by default.

Model Evaluation: Akaike Information Criterion and  $\Delta_i$ . To have the same metric for comparing larger sets of models en masse and model pairs (41), we evaluated models on estimated Akaike information criterion (AIC). In outline, the general form of the AIC =  $-2\log(\mathcal{L}) + 2K$  rewards goodness-of-fit through the maximized likelihood,  $\mathcal{L}$ , of the model given the data, while penalizing model complexity in proportion to the number of model parameters, K. Better-fitting models of the same data have larger likelihoods, hence smaller  $-2\log(\mathcal{L})$  (deviance). Simpler models have fewer parameters, i.e., smaller K. So, among a set of models of the same data, the better-fitting, simpler model(s),  $M_i$ , have lower AIC values than worse-fitting and/or more complex models. We evaluated the degree of empirical support for models in a set according to Burnham and Anderson's heuristics for  $\Delta_i = AIC_i - AIC_{min}$ , the difference between the AIC for model,  $M_i$ , and the minimum AIC among models being compared: "models having  $\Delta_i \leq 2$  have substantial support (evidence), those in which  $4 \le \Delta_i \le 7$  have considerably less support, and models having  $\Delta_i > 10$  have essentially no support" (ref. 42, pp. 270-271). Critically, these heuristics treat AIC differences less than two as meaningless for model selection-i.e., they characterize evidential ties and begin to look for AIC differences around four or greater to differentiate alternative models. Taken together, the AIC and heuristics comprise a practical general framework for investigating—comparing and selecting among—sets and pairs of models with fixed and random effects (*SI Appendix, AIC Model Selection*).

**Random-Effects Selection.** There is some debate in the recent mixed-effects modeling literature about whether maximal or parsimonious random effects are appropriate for hypothesis testing with LMER models (43, 44). The debate turns, in part, on how the decision to include, e.g., random slopes in addition to random intercepts impacts the rate of incorrect null hypothesis rejections (type I errors) vs. loss of power and failure to reject the null hypothesis (type II errors). We took the present project as an opportunity to evaluate the consequences of the decision as a case study of exploratory data analysis. Specifically, among the 11 candidate models with random effects ranging from maximal to minimal,  $M0, \ldots, M10$  (Table 2), we selected two for further investigation according to different decision rules: "Keep It Maximal" (KIM), select the maximal random effects for which the



Random effects model selection:  $\Delta_{Mi} = AIC_{Mi} - AIC_{min}$ 

Fig. 1. The time course and scalp distribution of AIC  $\Delta_{Mi}$  comparisons among models in the set { $M0, \ldots, M10$ } (Table 2). Each panel,  $\Delta_{Mi}$ , indicates how the AIC for model M<sub>i</sub> compares with the best-supported model (minimum AIC) among the 11 candidates at each time point and channel:  $\Delta_{Mi} = AIC_{Mi} - AIC_{min}$ . Since there is always some minimum AIC, somewhere among the models,  $\Delta_{Mi} = 0$ . As the panels show, this varies by time point and channel. The x axis is time in milliseconds; vertical lines indicate stimulus word onsets, and critical-article onset is at 0. The rainbow line plots show the time course of  $\Delta_{Mi}$  (y axis) for each channel in colors given by the channel legend; horizontal lines indicate the Burnham and Anderson (42)  $\Delta_i$  heuristic intervals bounded by 2, 4, 7, and 10. A few values for M9 and most for M10 are above 50 and not shown. The adjacent blue and red raster plots show the same data: Darker colors correspond to larger  $\Delta_{Mi}$  values; shading levels correspond to the heuristic intervals. EEG channels are arrayed on the y axis in the order given by the channel color legend: The top 11 rows are the left hemiscalp, the next four are midline, and the bottom 11 rows are the right hemiscalp. At a glance, the lightest patches among the raster plots indicate the best-supported (or equally well-supported) model(s) in the set ( $0 \le \Delta_{Mi} \le 2$ ), and darker patches indicate that the model is less well supported than an alternative ( $\Delta_{Mi} > 2$ ). Times and channels where Ime4::Imer() fitting generated a warning are indicated with red. Models M5 and M7 were selected for further investigation, based on the KIM and KIP selection rules, respectively. These results are for models fit to approximately 1.2e4 single-trial observations at 8-ms intervals and 26 EEG channels (Table 1).

model converges reliably; and "Keep It Parsimonious" (KIP), select the simplest random effects for which the model converges reliably and does not have substantially less support than the alternatives ( $\Delta_{M_i} \ge 4$ ).

Evidence for an Article-Cloze Effect:  $\Delta_M$  and ImerERPs. The critical empirical question is whether there is an association between article cloze and scalp potentials generated by brain activity in response to encountering those articles. We approached this in two ways, based on fitting the models selected by the KIM and KIP decision rules: 1) We computed  $\Delta_{M}$  and  $\Delta_{Mr}$  for the full and corresponding reduced model pairs, taking  $\Delta_{Mr} > 4$  as indicative of substantially less support for the reduced model; and 2) we examined the magnitude and CIs of the article-cloze (slope) rERPs for the full model.

The possible outcomes and interpretations of this rERP modeling are straightforward. If the article cloze and scalp potentials are unrelated, including article cloze in the model should have little impact on the goodness-of-fit, and  $\Delta_M$  for the full vs. reduced model should be around two because of the AIC penalty for the additional parameter. And in this same case, the articlecloze (slope) rERP waveforms should tend to be around zero plus or minus random variation, i.e., the x-y trend line for article cloze (x) vs. EEG (y) at each point in time should tend to be flat. Alternatively, if there is an approximately linear association between article-cloze probability and scalp potentials, the deviance term of the AIC for the full model should be smaller. In this case, the extent to which  $\Delta_{Mr}$  for the reduced model is greater than two indicates the degree to which the full model is better supported by the data after adjusting for its increased complexity, with  $\Delta_{Mr} > 4$  indicating a substantial difference in support. Furthermore, the time course and scalp distribution of the  $\Delta_{Mr}$  values and lmerERPs are important. To support the inference that the potentials are generated by a brain response to the article, an AIC  $\Delta_{Mr}$  effect should be evident in the interval after article onset, and not before. Likewise, the article-cloze (slope) rERP waveforms should tend to hover around zero prior to article onset and then deviate from zero afterward, with the polarity of the deviation, positive or negative, indicating the direction of the association (correlation).

Taken together, the full vs. reduced model pair  $\Delta_i$  values and the magnitude of the ImerERPs relative to their CIs are the basis of our evaluation of the strength of evidence for an article-cloze effect, the rough confirmatory analysis in Tukey's sense. In Tukey's view (ref. 31, p. 24), strong confirmatory null hypothesis testing requires designing, executing, and analyzing an experiment to ask and answer one question, thereby reducing the entire project to a single bit of information-one or zero, significant or not (ref. 32, p. 277). By contrast, our exploratory modeling aims to gauge where and when and to what extent-if any-there is evidence to support a linear approximating model of the relationship between article cloze and scalp potentials.

# Results

The following summarizes the main findings in the critical interval from 1.5 s before article onset up to the onset of the following word. Note that Figs. 1-3 display the 3 s of data modeled, which spans the two words after the article.

Random-Effects Selection. The LMER models M0, M1, ..., M10 (Table 2) hold constant the intercept and fixed effect of article and vary the random effects. Fig. 1 shows there is no unique best supported model with minimum AIC at all time points



**Fig. 2.** AIC  $\Delta_M$  pairwise full vs. reduced LMER model comparisons. (A) KIM: full (M5) vs. reduced (M5r). (B) KIP: full (M7) vs. reduced (M7r). Axes, scales, and data are as in Fig. 1. A and B, Top and Middle show AIC  $\Delta_M$  and  $\Delta_{Mr}$  for the full and reduced models, respectively, across the 3-s epoch, article onset at zero. The insets in A and B, Bottom zoom in to show AIC  $\Delta_M$  for the reduced model at the critical prenominal article in more detail. For both comparisons, during the 1.5-s interval preceding the critical article, the full and reduced models are equally supported,  $\Delta_M$  and  $\Delta_{Mr} < 2$ , with a few idiosyncratic exceptions. During the interval around 300 to 500 ms following the article onset (highlighted in magenta), the reduced models are substantially and systematically less supported at bilateral posterior scalp locations,  $\Delta_{M5r}$  and  $\Delta_{M7r} > 4$ , as indicated in by traces above four in the rainbow line plots and darker blue bands in raster plots.

and EEG channels—i.e., no single model where  $\Delta_{Mi} = AIC_{Mi}$  –  $AIC_{min} = 0$ . However, some models were much less supported than others in the 1.5-s prearticle to 0.5-s postarticle interval, and we selected two for further investigation. First, in accord with both decision rules, we ruled out models with substantial numbers of fitting warnings (M0, M1, M2, M3, M4, and M6), each of which included item or experiment random slopes for article cloze. Of those remaining, in accord with the KIM decision rule, we selected M5 with random intercepts for experiment, subject, and item and a random slope for subjects as the model with the maximal random effects that reliably converged, KIM M5: cloze + (1 | expt) + (cloze | subject) + (1 | item). We examined the remaining models with simpler random effects and, unsurprisingly, found intervals of substantially less support  $(\Delta_{Mi} > 4)$  for models that dropped any one of the experiment, subject, or item random variables entirely (M8, M9, or M10). Consequently, in accord with the KIP rule, we selected model M7 with random intercepts for experiment, subject, and item as the model with the most parsimonious random effects that was well-supported by the design and the data, KIP M7: cloze + (1 | expt) + (1 | subject) + (1 | item). Neither the KIM (M5) nor KIP (M7) models were entirely free of fitting warnings, but these were scattered irregularly across the times and channels and few in number, especially during the interval of interest. Although KIM and KIP decision rules may represent different extremes, in this particular instance, the models selected, M5 and M7, differed only in whether or not to include an article random slope for subjects.

**Evidence for an Article Cloze Effect.** With the KIM (M5) and KIP (M7) models selected for further investigation, we turned to the research question of primary interest: Is there evidence of an

association between article predictability and scalp potentials? We addressed this by pairwise AIC model comparisons between the full and reduced KIM (M5 and M5r) and KIP (M7 and M7r) models (Fig. 2) in conjunction with the values of the estimated coefficients for the article-cloze predictor in the full models,  $\hat{\beta}_{cloze}$ —i.e., the article-cloze rERPs (Fig. 3*B*).

We note first that  $\Delta_{M5}$  and  $\Delta_{M7}$  for the full KIM and KIP models, respectively, accord with the definitions of AIC and  $\Delta_M$ . These values range between zero and two at all times and channels (Fig. 2, *Top*), except for a few anomalous values, where the fitting failed to converge for the maximal model M5. These expected results support the face validity of the AIC estimates and  $\Delta_M$  calculations, which appear to be generally well-behaved for these models and data.

The key evidence for an article-cloze effect is observed at those scalp locations and times where the reduced models  $\Delta_{MSr}$  and  $\Delta_{M7r}$  values are >4, indicating a substantial decrease in goodness-of-fit when the article-cloze predictor is omitted from the model. For these reduced models (Fig. 2, *Middle* and *Bottom*), there are two intervals of immediate interest: the prestimulus interval (-1.5 to 0 s) and the critical article (0 to 0.5 s). The interval spanning the words immediately following the article (0.5 to 1.5 s), is relevant as well, albeit less directly, as we touch on in *Discussion*.

**Prestimulus**  $\Delta_{M}$ . During the 1.5 s preceding the onset of the critical article, values for the reduced KIM model range between zero and two (Fig. 24,  $\Delta_{M5r}$ ), with occasional irregular values above two (indicated by the darker blue speckles) and, again, a few anomalously large AIC values coincident with model-fitting warnings. The findings for the reduced KIP model with the parsimonious random effects are similar (Fig. 2*B*,  $\Delta_{M7r}$ ), except that



LMER model M5 fixed-effect regression ERPs

**B** Article cloze slope  $\widehat{\beta}_{cloze}$ 



**Fig. 3.** Model M5 linear mixed-effects rERPs (3 s, 26 channels). Solid lines plot the estimated regression parameter over time (milliseconds) relative to critical article onset at zero; bands indicate 95% CIs, and positive values are plotted up. Anterior to posterior scalp locations are arrayed from top to bottom in each panel. (A) Intercept ImerERPs ( $\hat{\beta}_0$ ) are analogs of grand mean average ERPs and show the characteristic morphology of visual evoked-potential responses, sharply defined transient peaks and troughs, especially prominent over the lateral occipital scalp. (*B*) Article-cloze ImerERPs ( $\hat{\beta}_{cloze}$ ) characterize the slope of the straight-line relationship between standardized article cloze and scalp potentials as it evolves over time. The *y* axis is  $\mu$ V per unit standardized cloze. The cloze ImerERPs show a transient positive response, predominantly over the bilateral posterior scalp, around 300 to 500 ms after article onset (magenta highlight) and not before, indicating a positive association between cloze probability and scalp potentials in response to the critical prenominal articles.

there are fewer fitting warnings and no anomalous  $\Delta_{M7r}$  excursions. Since  $\Delta_M \leq 2$  for the most part during the prestimulus interval, and rarely >4, we conclude that support for the full and reduced models does not differ substantially in this interval for either the KIM or KIP random effects. This evidential tie in the prestimulus interval is instructive for what it does not show. Given the design of the experiment, and the epoch centered on the entire 1.5-s prestimulus baseline, an effect of article predictability should be evident upon encountering the article, but not before. If the modeling showed an article-cloze effect prior to article onset, it could indicate something amiss in the design or execution of the experiments, the model specification or fitting, or the model comparison metric. In so far as we can determine with the present approach, examination of the 1.5 s of prestimulus activity for the 26 scalp locations at 8-ms intervals reveals no clear indication of these potential defects. Consequently, we suppose that article cloze effects observed in the interval following article onset may reasonably be attributed to a brain response to the article.

**Critical Article**  $\Delta_M$ . Following the onset of the critical *a*/*an* indefinite articles, the AIC differences between the full and reduced models do not appear to be dramatically different from those in the prestimulus interval until about 300 ms poststimulus. Then, between around 300 ms and the onset of the next word, AIC values for the reduced models, M5r and M7r, are systematically larger, predominantly over the bilateral posterior scalp, peaking around 400 ms (Fig. 2A and B, Bottom, magenta highlight). This increase was not observed over the anterior scalp. The results for the KIM and KIP models are similar: The KIP model  $\Delta_{M7r}$  values are slightly larger in some cases, there are fewer fitting warnings, and no anomalously large AIC values. For both the KIM and KIP comparisons, there appears to be an oscillation around 10 Hz in the reduced models ( $\Delta_{M5r}$ ,  $\Delta_{M7r}$ ) during the interval 300 to 500 ms poststimulus, and perhaps earlier, over the posterior scalp. These oscillations may indicate residual alpha-band noise EEG, though the possibility of an event-related 10-Hz amplitude modulation should not be overlooked. These oscillations make evaluation of the time course of AIC differences on a scale below about a 10th of a second precarious, but the slower phasic response is evident with or without the oscillations. We interpret this phasic increase in  $\Delta_{M5r}$  and  $\Delta_{M7r}$  above four for the KIM and KIP pairwise model comparisons as empirical supportrough confirmation-of a systematic association between article cloze and scalp potentials 300 to 500 ms over the posterior scalp. This effect is the crux of the argument for word-form prediction.

Article-Cloze ImerERPs. Whereas the full vs. reduced model AIC comparisons indicate when (around 300 to 500 ms poststimulus) and where (bilateral posterior scalp) there is evidence of an article-cloze effect, the magnitude and polarity of the estimated rERP slope coefficients characterize the magnitude and direction of the association under the assumption of a linear relationship. We found that the magnitude and CIs for the KIM and KIP intercept ( $\hat{\beta}_0$ ) and article cloze ( $\hat{\beta}_{cloze}$ ) ImerERPs are essentially indistinguishable over the entire 3-s epoch (*SI Appendix*, Fig. S4), and we present results here for the KIM model only (Fig. 3).

The model intercept lmerERPs  $(\hat{\beta}_0)$  are the rERP analog of grand mean average ERPs. These show the morphology characteristic of visual evoked potentials, a series of six transient responses to the six words presented two per second over the 3-s epoch (Fig. 3*A*). For the critical article cloze lmerERPs ( $\hat{\beta}_{cloze}$ ), we found that, prior to the onset of the article, they hovered around zero, and the 95% CIs for the point estimates generally span zero (Fig. 3*B*). Then, following the onset of the critical

PSYCHOLOGICAL AND COGNITIVE SCIENCES

article, we observed a biphasic positive response. The first phase began around 300 ms after the article, was larger predominantly over the posterior scalp, increased to a peak around 400 ms, and then decreased until shortly after the onset of the following word. The polarity of this deflection indicates a positive association, i.e., as cloze probability of an article increases, scalp potentials over the posterior scalp become more positive. This interval, about 300 to 500 ms postarticle, was the first time in the epoch where the lower bound of the 95% CI for the article-cloze rERP was above zero for sustained periods. A second, larger phasic positive deflection was observed, peaking around 400 ms after the word following the article, with a time course and scalp distribution corresponding to the larger second phase of increased AIC  $\Delta_{Mr}$  for the reduced KIM and KIP models that emerged after the onset of the word following the article (Fig. 2A and B, Middle).

In sum, we observed what appears to be a systematic, eventrelated lmerERP response to the article with a polarity, latency, and scalp distribution that coincide with previously reported reductions in N400 ERP amplitude with increasing cloze probability. We interpret this as direct evidence that the brain response to the article systematically covaries with the predictability of the indefinite articles *a* and *an*. To the extent the predictability of the article is dependent on the predictability of the not-yet-presented noun and its initial speech sound, the positive-going phasic article cloze lmerERP response is reasonably interpreted as indirect evidence for word form prediction.

# **Interim Summary**

When we modeled about 12,000 EEG single trials moment by moment at 26 scalp locations with appropriate linear mixedeffects models, we found that models that included article cloze probability as a predictor variable did a substantially better job accounting for the variability in potentials recorded over the posterior scalp around 300 to 500 ms after the onset of the article. The face validity of the modeling generally, and pairwise AIC model comparison results in particular, are bolstered by the facts that 1)  $\Delta_{\rm M} \leq 2$  for the full models are in line with theory; 2) the full and reduced models are equally supported during the prestimulus interval when no difference is expected; and 3) the direction of the observed positive association between articlecloze probability and scalp potentials characterized by the slope rERPs agrees with the reported reductions in average N400 ERP amplitude with increasing cloze probability (19, 45). As best we could determine, for these data, the perhaps-contentious choice to fit models with maximal or parsimonious random effects made little difference for characterizing the time course, scalp distribution, or strength of empirical support for the article-cloze effect based on model comparisons or for estimating the fixed effect of article cloze-i.e., the magnitude and precision of the ImerERP estimates.

# **Followup Analyses**

Since exploratory data investigation arrives at conclusions through an iterative process of evaluating assumptions and alternatives, we conducted a number of followup analyses, summarized briefly here (see *SI Appendix* for further details and discussion).

**Influential Data Diagnosis.** A general issue for the interpretation of estimated regression model coefficients is whether subsets of extreme or outlying observations exert a disproportionate influence on estimates and exaggerate (or obscure) patterns seen in the bulk of the data. For modeling the time course of the article-cloze effect, this question is whether the morphology of the lmerERP waveforms, in particular, is driven by a subset of unrepresentative data. Mixed-effects modeling is computationally intensive, and influence diagnostics based on model refitting are intractable for data on the scale of this analysis at present, so we fell back to ordinary least squares (SI Appendix, Influential Data Diagnosis). We identified and excluded a subset of about 5% of the single-trial epochs that contained the highest proportion of potentially influential observations. We then refit the KIP and KIM models to this trimmed dataset and computed how much the amplitude of the intercept and article cloze ImerERPs changed as a consequence of the trimming-i.e., we computed a version of the scaled deleted observation difference in beta (DFBETAS) data diagnostic, adapted for rERPs. We assumed that article-cloze DFBETAS  $\pm 2$  would indicate an unusually large change in the rERP estimate based on a large n Student's t distribution (35). We found that there were few DFBETAS excursions of that magnitude, and those that occur do so at the peaks and troughs of approximately 10-Hz oscillations (SI Appendix, Fig. S6). This oscillation suggests that the epochs identified and excluded contained high-amplitude alpha-band activity. Crucially, the time course and distribution of  $\Delta_{Mr}$  values for the reduced KIM (M5r) and KIP (M7r) models of the trimmed data still showed the phasic increase over posterior scalp around 300 to 500 ms, and the article-cloze Imer-ERPs ( $\beta_{cloze}$ ) showed the corresponding positive deflection (SI Appendix, Fig. S7). So it appears that the article-cloze effects observed in the initial analysis were not driven entirely by this subset of potentially influential trials.

Modeling Experiment as a Fixed Effect. The designs and procedures of EEG experiments 1, 2, and 3 are sufficiently similar to justify pooling the data for purposes of modeling the brain response to the critical indefinite articles, provided that systematic variation between the experiments is also accounted for. Since, for our purposes, systematic differences between the experiments is nuisance variation and the different numbers of trials in the three experiments make the design substantially unbalanced, we modeled experiment as a random variable. However, views may differ on the appropriate treatment of categorical variables as fixed vs. random, and the consequences for drawing model-based inferences, particularly when the number of levels is small (for discussion, see ref. 46, p. 20ff and ref. 47, pp. 246 and 275ff). So, we investigated the question by modeling the single-trial EEG with article cloze and experiment as fixed effects, retaining the KIM and KIP random effects for subjects and items; Table 2 KIM (M11 and M11r) and KIP (M12 and M12r). We found that fitting full and reduced models with experiment as a fixed effect converged reliably, and the pattern of AIC  $\Delta_{\rm M}$  and  $\Delta_{\rm Mr}$  for the pairwise full vs. reduced model comparisons, and article-cloze rERPs and their CIs were essentially the same as for models with a random intercept for experiment (SI Appendix, Fig. S8). So, in this instance, the choice of fixed vs. random effect for the experiment variable was immaterial for inferences about the article-cloze effects.

Modeling Experiments 1, 2, and 3 Separately. To assess whether the article-cloze effect observed for the data pooled across the three experiments was representative of each experiment individually, we split the data by experiment and fit the full and reduced model pairs in Table 2: KIM (M13 and M13r) and KIP (M14 and M14r). For each experiment, we examined AIC  $\Delta_{\rm M}$ and  $\Delta_{Mr}$  measures and the article lmerERPs (experiment 1, SI Appendix, Fig. S9; experiment 2, SI Appendix, Fig. S10; and experiment 3, SI Appendix, Fig. S11). The results were mixed for the AIC model comparisons and somewhat more consistent for the article-cloze lmerERPs. For the experiment 1 data, fitting the full and reduced models with KIM random effects had considerable difficulty converging. Fitting the full and reduced KIP models converged reliably with irregular intervals of  $\Delta_{M14r} > 4$  throughout the 3-s epoch and no clear break in the pattern between the prearticle and postarticle interval that suggests an event-related

brain response to the article. So, the AIC model comparisons did not provide clear evidence for a relationship between article cloze and an event-related EEG response in experiment 1. For the experiment 2 data, the KIM and KIP models converged reliably with only a modest increase in convergence failures for the KIM models. Overall, the time course and scalp distributions were generally similar to those for models of the data pooled across all three experiments, with scattered idiosyncratic  $\Delta_{M13r}$  > 4 in the prestimulus interval and a systematic onset and offset around 300 and 500 ms postarticle, respectively. For the experiment 3 data, there are slightly more convergence failures for the KIM models, and prestimulus AIC differences for the reduced model are evident, more so for the KIP comparison, though not to the extent observed for experiment 1. In the critical interval around 300 to 500 ms postarticle, AIC differences larger than in experiment 1 and smaller than experiment 2 rise and fall. In all three experiments, the article-cloze lmerERPs tended to vary around zero prior to the critical article onset, after which they showed a small positive deflection followed by a larger one over the bilateral posterior scalp. The onset of this rERP response in experiment 1 appears to be perhaps 100 to 200 ms later than in experiments 2 and 3, though the timing in experiment 1 was obscured by a pronounced oscillation around 10 Hz. In sum, the AIC  $\Delta_M$  results observed for the data pooled across the experiments appeared to be more representative of experiments 2 and 3 than experiment 1. The pattern of article-cloze slope lmerERPs was more consistent, and all three experiments showed a similar, albeit more variable, biphasic positive response following the article, similar to that observed for the pooled data.

LMER Modeling Interval Mean Amplitude. Whereas the rERP analyses described thus far model the moment-by-moment time

course of the article-cloze effect from 1.5 s before to 1.5 s after the article, experimental EEG studies using event-related designs, including DUK05 and NIET18, often base inferences about event-related brain responses on measurements of scalp potentials aggregated over a specific time interval-e.g., mean amplitude between 200 or 300 and 500 ms poststimulus-relative to mean amplitude in a specified prestimulus baseline intervale.g., 100, 200, or 500 ms. To compare the LMER rERP results with interval mean amplitude analyses, we reduced the singletrial EEG time-series data to four sets of summary measures: mean amplitude in two poststimulus intervals (200 to 500 ms and 300 to 500 ms), each measured relative to a baseline of mean amplitude in two intervals (100 and 500 ms prestimulus). We then modeled these single-trial time-averaged mean amplitude measurements by fitting the KIM (M5 and M5r) and KIP (M7 and M7r) model pairs at each of the 26 EEG channels separately (c.f., NIET18 LMER analyses of mean potentials aggregated in the interval 200 to 500 ms poststimulus across six centro-parietal scalp locations).

Consistent with the ImerERP time-course analysis, modeling the potentials averaged across these temporal intervals also found a positive association between article cloze, with a posterior scalp distribution (Fig. 4). Across the different combinations of model random effects, baseline intervals, and N400 measurement intervals, only the poststimulus measurement interval had much impact on the results (OSF: udck19\_pipeline\_5. html; https://osf.io/hbgfs/). Regardless of the random effects or prestimulus baseline interval, the magnitudes of the estimated article-cloze coefficients for the longer and earlier 200- to 500-ms poststimulus interval measurements tend to be around  $\frac{1}{3}$  smaller than for the measurements made 300 to 500 ms poststimulus (Fig. 4A vs. Fig. 4B). Attenuated article effects in the 200- to



**Fig. 4.** Comparison of KIM models M5 and M5r of single-trial mean EEG amplitude measured in a longer, earlier-starting interval 200 to 500 ms poststimulus (*A*) and a shorter, later-starting interval 300 to 500 ms poststimulus (*B*). *A* and *B*, *Left* show the AIC  $\Delta_{M5r}$  values for the pairwise full (M5) vs. reduced (M5r) KIM model comparison.  $\Delta_{M5}$  for the full model (not shown) was between zero and two, as expected for this comparison. A and *B*, *Right* show the magnitude of the estimated fixed-effect coefficient for article cloze,  $\hat{\beta}_{cloze}$ , with positive values in red and with filled circles only at locations where the 95% CI for the estimate did not include zero. Like the temporally fine-grained rERP models, this single-trial LMER modeling indicates a positive association between article cloze and potentials over the bilateral posterior scalp around 400 ms poststimulus, albeit more robust for the shorter and later interval 300 to 500 ms post-stimulus. Results in this figure are for poststimulus potentials measured relative to mean amplitude in a 500-ms prestimulus baseline; results for measurements relative to a 100-ms prestimulus baseline were similar. See (OSF: udck19\_pipeline\_5.html; https://osf.io/hbgfs/) for these and additional analyses.

500-ms postarticle interval are consistent with the time-course rERP modeling, which found no clear evidence of the article effect before 300 ms poststimulus.

Lurking Variables and Spurious ImerERPs. Another general issue for the interpretation of an estimated regression model coefficient is the spurious effect that can result from a "lurking" variable-i.e., a variable that is causally related to the response variable and correlated with the predictor, but omitted from the model (for discussion, see SI Appendix, pp. 6 to 8). If the articlecloze lmerERPs in Fig. 3 are driven purely by correlation with some causal factor unrelated to the form of the indefinite article, interpreting them as support for word-form prediction would be unwarranted. The impact of a lurking variable on a regression coefficient can be quantified as the omitted variable bias (e.g., ref. 35, pp. 111-112), which we used to investigate the impact of a variable known to be correlated with article cloze, but unrelated to the form of the indefinite article.\* Since our normative stimulus testing was free response, the proportion of indefinite articles goes down as the proportion of nonarticle responses (e.g., bare plurals, adjectives, or definite articles) goes up. The article and nonarticle cloze probabilities are negatively correlated (r = -0.264, P < 0.0001; *SI Appendix*, Fig. S12). We modeled the nonarticle-cloze rERP (SI Appendix, Fig. S13) and found that, despite this correlation, the omitted variable bias does not account for the article-cloze ImerERP (SI Appendix, Fig. S14). Numerous variables are associated with article cloze and scalp potentials to some degree. However, unless the correlations are strong and the omitted variable rERPs are large, the bias is small and, thus, unlikely to account for the article-cloze effect.

# Discussion

The project reported herein aims to shed light on the recent theoretical controversy about whether the human-language comprehension mechanism anticipates the phonological form of upcoming words. The crucial empirical question is whether processing at the prenominal articles *a/an* varies with their predictability since, other things equal, the factor responsible for the form of the indefinite article is the initial speech sound of a not-yet-encountered word. Because of this phonological dependency, direct evidence of an effect of predict-*ability* at the article may be reasonably interpreted as indirect evidence that, by then, upcoming noun word forms were predict-*ed*.

To investigate the time course of the electrical brain activity, we modeled single-trial EEGs recorded before, during, and after presentation of pre-nominal indefinite articles (a/an) in three experiments that manipulated the predictability (cloze probability) of nouns in sentence contexts read by healthy younger adults at two words per second in central vision. Our interim conclusion was that models that include article-cloze probability as a continuous predictor do a substantially better job accounting for the variability in potentials recorded over the bilateral posterior scalp around 300 to 500 ms after the onset of the article than do models that omit this variable. Since this was not the case during the 1.5 s prior to the article, we interpreted these results as evidence of a systematic association between article-cloze probability and scalp potentials generated by the brain response to the article. The latency, polarity, and scalp distribution of this article-cloze effect is generally consistent with the association between cloze probability and scalp potentials (19, 45).

Exploratory investigation of alternatives indicated that evidence for the association does not appear to depend on the choice of maximal (KIM) or parsimonious (KIP) random effects, to be driven by the influence of a subset of unrepresentative data, or to depend on whether the experiment variable is modeled as a fixed or random effect. That said, the article-cloze effect appears to be markedly smaller (less variability accounted for and lower-amplitude slope lmerERPs) than a corresponding effect at the following word (Figs. 2 and 3, immediately after the magenta highlight). In this experimental design (... a kite ...), article-cloze probability is correlated, though not perfectly, with noun-cloze probability. The larger  $\Delta_{Mr}$  and lmerERP effects for the article-cloze predictor variable on the following word are likely a consequence of this relationship, but cannot be strictly attributed to the contextually supported nouns because in a subset of materials in experiment 2, a phonologically legal adjective is interposed between the article and noun, an orange kite. Given the high proportion of nouns relative to adjectives in the combined data, it is reasonable to suppose that modeling potentials elicited by the nouns with noun cloze as a predictor variable would find similar, if not larger, effects, but testing this speculation is tangential to the present aims and beyond the scope of this report. Although the comparison is imperfect, in all of the models investigated, the magnitude of the transient article-cloze rERP response at the article was smaller than at the following word. In this respect, the pattern is consistent with other studies that recruit sequential-dependency experimental designs to test for prediction in language comprehension and report relatively small and variable ERP effects at the probe word (8, 9, 11-14, 16, 17).

LMER modeling the single-trial data for each experiment separately found that article-cloze slope lmerERPs for all three experiments showed a biphasic positive response following the article, similar to that observed for the pooled data, albeit more variable. The AIC  $\Delta_{\rm M}$  patterns for the individual experiment pairwise model comparisons were similar to the pooled data for two of the datasets, experiment 2 and experiment 3 to a lesser extent, but not experiment 1. This is not entirely surprising, since there are roughly twice as many single-trial observations in experiments 2 and 3 as in experiment 1 (Table 1). It may be that the two-part stimulus presentation procedure and/or the additional materials developed for experiments 2 and 3 afford a better opportunity to observe a small article-cloze effect with a single-trial LMER analysis than do the procedures and materials used for the DUK05 study. While the rERP modeling does not show clear evidence of an article effect for the experiment 1 data considered on their own, the findings are consistent with the stronger support provided by the replication and extension studies that followed. We also modeled single-trial mean amplitude in the postarticle intervals 200 to 500 ms and 300 to 500 ms with the same KIM and KIP LMER models used for the time-course modeling. The choice of KIM vs. KIP model and choice of measurement relative to a shorter (100 ms) vs. longer (500 ms) prestimulus baseline interval had a negligible impact on the results, but in all cases, the magnitude of the article-cloze effect was markedly smaller for the 200- to 500-ms poststimulus interval.

Taken together, this pattern of findings may be relevant to understanding the failure to observe an effect of article cloze reported in NIET18. That study tested only the smaller set of *a/an* items and single-sentence rapid serial visual presentation (RSVP) used for the study reported in DUK05 (experiment 1 in this report), whereas we found that the article-cloze effects may be more readily observed in the followup experiments 2 and 3 with the expanded sets of items and two-part stimulus presentation. The LMER analyses reported in NIET18 were conducted on single-trial mean amplitudes in the interval 200- to 500-ms postarticle, averaged over six centro-parietal electrode locations, whereas our time-course modeling at each scalp location found the article-cloze effect to have a more posterior distribution and somewhat later onset (Figs. 2 and 4). The LMER model pairs

<sup>\*</sup>We thank an anonymous reviewer for suggesting this example.

compared in NIET18 for the likelihood-ratio tests of the null hypothesis assumed maximal random effects with correlated random intercepts and slopes for subjects and items, whereas we found that in pairwise AIC model comparisons, the article-cloze effect was, at times, slightly attenuated for the maximal relative to parsimonious model (Fig. 2,  $\Delta_{M5r}$  vs.  $\Delta_{M7r}$ ). So, although the decisions made in conducting and analyzing the study reported in NIET18 are defensible for purposes of conducting a direct replication of the DUK05 study, they may be suboptimal for answering the scientific question of interest about word-form prediction.

The failure of the NIET18 report to observe a prenominal cloze probability effect in a much larger data sample with generally similar design parameters as the DUK05 report raised the possibility that there is no such systematic relationship between prenominal article cloze and electrical brain activity at all. This is the primary research question that our project was designed to address, using an exploratory-analysis approach. To answer this specific question, we selected data from experiments similar to both DUK05 and NIET18: a/an designs testing young adults reading at two words per second in central vision. This selection affords meaningful comparisons among the studies, but it also means that the results do not answer looming secondary questions about how various experimental variables, such as presentation rate or age, among others, might impact the model fits and morphology of article-cloze rERP waveforms. Still less does the analysis answer broader questions about the generalizability of the findings in the way a meta-analysis might. Although we pooled data across multiple studies, ours is a forensic EEG data investigation, not a meta-analysis. And, considered in its entirety, the pattern of results from the lmerERP modeling we conducted does appear to provide direct evidence of an association (quantitative relationship) between prenominal article cloze and scalp potentials. Of course, the time course, scalp distribution, and polarity of article-cloze slope lmerERPs-i.e., the estimated  $\hat{\beta}_{cloze}$  coefficients—are key to this interpretation. And, of course, if a model omits (any) relevant predictor variables, estimates of the coefficients for variables that are included may be biased, and, in turn, inferences drawn from the model may be wrong; we never know with certainty whether a model omits relevant predictors. Interpreting our findings as evidence of a structural relation between the predictability of the stimulus and the brain response it elicits requires the stronger assumption that there are no serious lurking variables. This caveat applies to all regression modeling. All of the more reason to systematically explore the data, "look for what can be seen, even if not anticipated." (ref. 48, p. 24).

# Conclusions

In contrast with the large-scale null result reported in NIET18, our moderately large-scale LMER modeling of single-trial EEG moment-by-moment at 26 scalp locations finds direct empirical support for an association between the predictability of prenominal indefinite articles and the brain's response to encountering them in word-by-word reading. This effect may reasonably be attributed to prediction of upcoming word forms in answer to the question of scientific interest. The exploratory modeling reported herein illustrates an approach to experimental EEG data analysis that may prove a useful complement to confirmatory null hypothesis testing.

# **Materials and Methods**

**Methods.** All normative stimulus testing and EEG studies were conducted under human-subjects research protocols approved by the University of California, San Diego Institutional Review Board. Volunteers were recruited by flyer and through the campus subject pool. Upon their arrival at the laboratory, the experimental procedures were explained verbally, and participants were presented with a printed consent form describing the procedures and

potential risks. Individuals who elected to participate in the study provided their written informed consent and received 2 h of course credit, cash payment, or a combination, at their discretion. The normative predictability of the critical prenominal indefinite articles and nouns was operationally defined as the relative frequency of production in a sentence-fragmentcompletion task (cloze probability) in separate testing with individuals who did not participate in the EEG experiments. Participants in the EEG studies were healthy, young-adult, right-handed native English speakers. Salient differences between the EEG experiments included the number of participants and experimental items (Table 1), the presentation mode (one vs. two sentences per trial), experimental conditions ( $\pm$  prenominal adjectives,  $\pm$ filler items), counterbalancing scheme, the distribution of cloze probabilities, and normative plausibility of critical nouns (SI Appendix, Table S1). In all three EEG experiments, sentences containing the critical prenominal articles were read word-by-word at a fixed rate of approximately two per second, and the EEG data-acquisition and data-processing procedures were the same (SI Appendix, EEG Recording and Data Processing). Prior to modeling, the EEG data were visually screened for artifacts, smoothed (25-Hz low-pass phase-compensated finite impulse response filter), downsampled to 125 samples per second, centered by subtracting the mean of the 1,496ms prestimulus interval for each channel, and rescreened for EEG artifacts by computer algorithm (see SI Appendix, EEG Experimental Procedures for details and (OSF: udck19 pipeline 1.html; https://osf.io/y2wa3/) for exclusions tabulated by experiment, participant, and item).

LMER Model Fitting. For the data pooled across the three experiments, each observation was coded for the experiment, subject, and stimulus item. Each item corresponded to the context prior to the critical article and provided one cloze value for a and one for an (see SI Appendix, Fig. S2 for the distributions of article cloze across and within each design). Prior to modeling the EEG, the article-cloze predictor variable was scaled from proportions of response (0.0 to 1.0) to standardized units ("z scores") by centering and dividing by the SD. The 1.2e4 screened single-trial EEG epochs were stacked into a dataframe (4.5e6 rows = 1.2e4 epochs  $\times$  375 samples per epoch), with each row indexed for epoch and time-stamped relative to article onset, with the experiment, subject, item, standardized cloze values, and the 26 EEG channels in columns. To model these single-trial data, we used fitgrid (49), an open-source Python package we developed in the laboratory that implements mixed-effects model fitting via the pymer4 (50) interface to the ImerTest (51) and Ime4 (39) R packages (52). With fitgrid, we swept each LMER model in Table 2 across 3-s epochs of data with the critical article in the middle (375 time points, 8-ms intervals = 125 samples per second: 26 electrode locations spaced about 5 cm apart) and collected the lme4::lmer() profiled ML fits (REML = FALSE) in a tabular grid. From this grid of model fits, we extracted summary measures returned by ImerTest::Imer() for the fit at each time and channel, including AIC,  $\hat{\beta}_i$  estimates for the intercept and article-cloze ImerERPs and their 95% Wald Cls, and fitting algorithm warnings (ref. 49; fitgrid.lmer). The  $\hat{\beta}_{\text{cloze}}$  lmerERPs in Fig. 3B and interval mean amplitude coefficients in Fig. 4 for standardized cloze may be converted to coefficients  $\hat{B}_{cloze}$  on the original cloze scale ( $\mu V$ /cloze) as  $\hat{B}_{cloze} = \hat{\beta}_{cloze} / SD_{cloze}$  with the article cloze SD values in Table 1.

**Data Availability.** Stimulus materials, aggregated behavioral and EEG data, summary measures, data-analysis software, and reproduction recipe are deposited in the Open Science Foundation repository OSF: UDCK (38) and licensed under the Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License, which may be viewed here: http://creativecommons.org/licenses/by-nc-nd/4.0. Behavioral and EEG data related to an identifiable natural person are maintained under control of the principal investigator M.K. and co-investigators. Contact the corresponding author for information about further use of the research materials or for access to privacy-sensitive data under a written data-sharing agreement.

ACKNOWLEDGMENTS. We thank Anna Stoermann for her painstaking work preparing the normative and EEG data and report; Andrey Portnoy for his key original contributions designing and implementing the fitgrid software used for these analyses; Lauren Liao for her work on early prototypes; Nathaniel Smith for sharing his code for reading the laboratory's EEG data; Melissa Troyer, Anna Stoermann, Reina Mizrahi, Seana Coulson, Megan Bardolph, and Emily Provenzano for critical discussion of EEG data analysis, modeling, and testing; and Gregory A. Miller and Kara D. Federmeier for their comments on earlier drafts. This research was supported by NIH Grants SR01HD022614-27 and 1R01AG048252-01A1. This report describes de novo secondary analyses of extant behavioral and electrophysiological data; these findings have not previously been published in any form.

- M. C. MacDonald, Y. Hsiao, "Sentence comprehension," in The Oxford Handbook of Psycholinguistics, S. A. Rueschemeyer, M. G. Gaskell, Eds. (Oxford University Press, Oxford, UK, ed. 2, 2018).
- G. T. M. Altmann, J. Mirkovic, Incrementality and prediction in human sentence processing. Cognit. Sci. 33, 583–609 (2009).
- G. R. Kuperberg, T. F. Jaeger, What do we mean by prediction in language comprehension?. Lang. Cognit. Neuros. 31, 32–59 (2016).
- M. Kutas, K. A. DeLong, N. J. Smith, "A look around at what lies ahead: Prediction and predictability in language processing" in *Predictions in the Brain: Using Our Past* to Generate a Future, M. Bar, Ed. (Oxford University Press, New York, NY, 2011), pp. 190–207.
- R. M. Cooper, Control of eye fixation by meaning of spoken language: New methodology for real-time investigation of speech perception, memory, and language processing. *Cognit. Psychol.* 6, 84–107 (1974).
- M. K. Tanenhaus, M. J. Spiveyknowlton, K. M. Eberhard, J. C. Sedivy, Integration of visual and linguistic information in spoken language comprehension. *Science* 268, 1632–1634 (1995).
- G. T. M. Altmann, Y. Kamide, Incremental interpretation at verbs: Restricting the domain of subsequent reference. *Cognition* 73, 247–264 (1999).
- N. Y. Y. Wicha, E. A. Bates, E. M. Moreno, M. Kutas, Potato not pope: Human brain potentials to gender expectation and agreement in Spanish spoken sentences. *Neurosci. Lett.* 346, 165–168 (2003).
- N. Y. Y. Wicha, E. M. Moreno, M. Kutas, Anticipating words and their gender: An event-related brain potential study of semantic integration, gender expectancy, and gender agreement in Spanish sentence reading. J. Cognit. Neurosci. 16, 1272–1288 (2004).
- N. Y. Y. Wicha, E. M. Moreno, M. Kutas, Expecting gender: An event related brain potential study on the role of grammatical gender in comprehending a line drawing within a written sentence in Spanish. *Cortex* 39, 483–508 (2003).
- J. J. A. Van Berkum, C. M. Brown, P. Zwitserlood, V. Kooijman, P. Hagoort, Anticipating upcoming words in discourse: Evidence from ERPs and reading times. J. Exp. Psychol. Learn. Mem. Cogn. 31, 443–467 (2005).
- M. Otten, M. S. Nieuwland, J. J. A. van Berkum, Great expectations: Specific lexical anticipation influences the processing of spoken language. *BMC Neurosci.* 8, 89 (2007).
- M. Otten, J. Van Berkum, Discourse-based word anticipation during language processing: Prediction or priming? *Discourse Process* 45, 464–496 (2008).
- 14. M. Otten, J. J. A. Van Berkum, Does working memory capacity affect the ability to predict upcoming words in discourse? *Brain Res.* **1291**, 92–101 (2009).
- A. R. Kochari, M. Flecken, Lexical prediction in language comprehension: A replication study of grammatical gender effects in Dutch. *Lang. Cognit. Neurosci.* 34, 239–253 (2019).
- J. M. Szewczyk, H. Schriefers, Prediction in language comprehension beyond specific words: An ERP study on sentence comprehension in Polish. J. Mem. Lang. 68, 297–314 (2013).
- B. Nicenboim, S. Vasishth, F. Rösler, Are words pre-activated probabilistically during sentence comprehension? Evidence from new data and a Bayesian random-effects meta-analysis using publicly available data. *Neuropsychologia* 142, 107427 (2020).
- K. A. DeLong, W. H. Chan, M. Kutas, Similar time courses for word form and meaning preactivation during sentence comprehension. *Psychophysiology* 56, 14 (2019).
- K. A. DeLong, T. P. Urbach, M. Kutas, Probabilistic word pre-activation during language comprehension inferred from electrical brain activity. *Nat. Neurosci.* 8, 1117–1121 (2005).
- K. A. DeLong, "Electrophysiological explorations of linguistic pre-activation and its consequences during online sentence processing," Doctoral dissertation (University of California, San Diego, La Jolla, CA, 2009).
- K. A. DeLong, D. M. Groppe, T. P. Urbach, M. Kutas, Thinking ahead or not? Natural aging and anticipation during reading. *Brain Lang.* 121, 226–239 (2012).
- C. D. Martin et al., Bilinguals reading in their second language do not predict upcoming words as native readers do. J. Mem. Lang. 69, 574–588 (2013).
- A. Ito, A. E. Martin, M. S. Nieuwland, How robust are prediction effects in language comprehension? Failure to replicate article-elicited N400 effects. *Lang. Cognit. Neurosci.* 32, 954–965 (2017).

- M. S. Nieuwland et al., Large-scale replication study reveals a limit on probabilistic prediction in language comprehension. eLife 7, e33468 (2018).
- M. L. Head, L. Holman, R. Lanfear, A. T. Kahn, M. D. Jennions, The extent and consequences of p-hacking in science. *PLoS Biol.* 13, e1002106 (2015).
- L. D. Nelson, J. Simmons, U. Simonsohn, "Psychology's renaissance" in Annual Review of Psychology, S. T. Fiske, Ed. (Annual Reviews, Palo Alto, CA, 2018), vol. 69, pp. 511– 534.
- 27. A. Gelman, E. Loken, The statistical crisis in science. Am. Sci. 102, 460-465 (2014).
- N. Kriegeskorte, W. K. Simmons, P. S. F. Bellgowan, C. I. Baker, Circular analysis in systems neuroscience: The dangers of double dipping. *Nat. Neurosci.* 12, 535–540 (2009).
- J. W. Tukey, Analyzing data: Sanctification or detective work. Am. Psychol. 24, 83–91 (1969).
- J. W. Tukey, Data analysis, computation and mathematics. Q. Appl. Math. 30, 51–65 (1972).
- J. W. Tukey, We need both exploratory and confirmatory. Am. Statistician 34, 23–25 (1980).
- J. W. Tukey, "Data analysis and behavioral science or learning to bear the quantitative man's burden by shunning Badmandments" in *The Collected Works of John W. Tukey: Philosophy and Principles of Data Analysis 1949-1964*, L. V. Jones, Ed. (CRC Press, Boca Raton, FL, 1986) Vol. III, pp. 187–389.
- J. T. Behrens, Principles and procedures of exploratory data analysis. *Psychol. Methods* 2, 131–160 (1997).
- J. Cohen, P. Cohen, S. West, L. Aiken, Applied Multiple Regression/Correlation Analysis for the Behavioral Sciences (Lawrence Erlbaum Associates, Mahwah, NJ, ed. 3, 2003).
- J. Fox, Applied Regression Analysis and Generalized Linear Models (Sage Publications, Thousand Oaks, CA, 2008).
- M. H. Kutner, C. J. Nachtsheim, J. Neter, W. Li, Applied Linear Statistical Models (McGraw-Hill Irwin, Boston, MA, 2005), Vol. 5.
- N. J. Smith, M. Kutas, Regression-based estimation of ERP waveforms: I. The rERP framework. *Psychophysiology* 52, 157–168 (2015).
- T. P. Urbach, K. DeLong, W. Chan, M. Kutas, UDCK2020 An exploratory data analysis of word form prediction during word-by-word reading. OSF project repository https://doi.org/10.17605/OSF.IO/TKSUR (2020).
- D. Bates, M. Mächler, B. Bolker, S. Walker, Fitting linear mixed-effects models using lme4. 67, 48 (2015).
- R. H. Baayen, D. J. Davidson, D. M. Bates, Mixed-effects modeling with crossed random effects for subjects and items. J. Mem. Lang. 59, 390–412 (2008).
- K. Burnham, D. Anderson, Model Selection and Multimodel Inference: A Practical Information-Theoretic Approach (Springer-Verlag, Berlin, Germany, ed. 2, 2002).
- K. P. Burnham, D. R. Anderson, Multimodel inference: Understanding AIC and BIC in model selection. Socio. Methods Res. 33, 261–304 (2004).
- D. J. Barr, R. Levy, C. Scheepers, H. J. Tily, Random effects structure for confirmatory hypothesis testing: Keep it maximal. J. Mem. Lang. 68, 255–278 (2013).
- H. Matuschek, R. Kliegl, S. Vasishth, H. Baayen, D. Bates, Balancing type I error and power in linear mixed models. J. Mem. Lang. 94, 305–315 (2017).
- M. Kutas, S. A. Hillyard, Brain potentials during reading reflect word expectancy and semantic association. *Nature* 307, 161–163 (1984).
- A. Gelman, Analysis of variance: Why it is more important than ever. Ann. Stat. 33, 1–31 (2005).
- A. Gelman, J. Hill, Data Analysis Using Regression and Multilevel/Hierarchical Models (Cambridge University Press, Cambridge, UK, 2007).
- J. Tukey, Exploratory Data Analysis (Addison-Wesley Publishing Company, Reading, MA, 1977).
- T. Urbach, A. S. Portnoy, Fitgrid 0.4.8, https://doi.org/10.5281/zenodo.3581504 (2019).
  E. Jolly, Pymer4: Connecting R and Python for linear mixed modeling. J. Open Source Software 3, 862 (2018).
- A. Kuznetsova, P. B. Brockhoff, R. H. B. Christensen, ImerTest package: Tests in linear mixed effects models. J. Stat. Software 82, 1–26 (2017).
- R Core Team, R: A Language and Environment for Statistical Computing (R Foundation for Statistical Computing, Vienna, Austria, 2019).