

# A Tale of Two Positivities and the N400: Distinct Neural Signatures Are Evoked by Confirmed and Violated Predictions at Different Levels of Representation

Gina R. Kuperberg<sup>1,2</sup>, Trevor Brothers<sup>1,2</sup>, and Edward W. Wlotko<sup>1,3</sup>

## Abstract

■ It has been proposed that hierarchical prediction is a fundamental computational principle underlying neurocognitive processing. Here, we ask whether the brain engages distinct neurocognitive mechanisms in response to inputs that fulfill versus violate strong predictions at different levels of representation during language comprehension. Participants read three-sentence scenarios in which the third sentence constrained for a broad event structure, for example, {*Agent caution animate–Patient*}. *High constraint* contexts additionally constrained for a specific event/lexical item, for example, a two-sentence context about a beach, lifeguards, and sharks constrained for the event, {*Lifeguards cautioned Swimmers*}, and the specific lexical item *swimmers*. *Low constraint* contexts did not constrain for any specific event/lexical item. We measured ERPs on critical nouns that fulfilled and/or violated each of these constraints. We found clear, dissociable effects to fulfilled semantic predictions (a reduced

N400), to event/lexical prediction violations (an increased *late frontal positivity*), and to event structure/animacy prediction violations (an increased *late posterior positivity/P600*). We argue that the late frontal positivity reflects a large change in activity associated with successfully updating the comprehender's current situation model with new unpredicted information. We suggest that the late posterior positivity/P600 is triggered when the comprehender detects a conflict between the input and her model of the communicator and communicative environment. This leads to an initial failure to incorporate the unpredicted input into the situation model, which may be followed by second-pass attempts to make sense of the discourse through reanalysis, repair, or reinterpretation. Together, these findings provide strong evidence that confirmed and violated predictions at different levels of representation manifest as distinct spatiotemporal neural signatures. ■

## INTRODUCTION

The goal of language comprehension is to extract the communicator's intended message. This is challenging. Linguistic inputs unfold quickly, they are often ambiguous, and our communicative environments are noisy. It therefore helps if we can use context to mobilize our stored linguistic (and nonlinguistic) knowledge to predict upcoming information. If incoming information matches these predictions, its processing will be facilitated. On the other hand, there will be times when we encounter inputs that violate strong predictions. In this study, we use ERPs to ask whether the brain engages distinct neurocognitive mechanisms in response to inputs that fulfill and inputs that violate strong predictions at different levels of representation.

ERPs have provided some of the strongest evidence that the brain is sensitive to predictive processes during language comprehension. The ERP component that is primarily sensitive to *fulfilled* semantic predictions is the N400—a centroparietally distributed negative-going

waveform that is largest between 300 and 500 msec after the onset of an incoming word. The N400 is highly sensitive to the probability of encountering a word's semantic features given its preceding context (Kuperberg, 2016; Kutas & Federmeier, 2011; Kutas & Hillyard, 1984).<sup>1</sup> If this word's features match semantic features that have already been predicted by a highly constraining context (e.g., as in “birthday” following context (1)), then it will evoke a smaller (less negative) N400 than a word appearing in a low constraint context (e.g., “collection” following context (2)). Importantly, however, the N400 is not a direct index of a lexical prediction violation: Its amplitude is just as large to “collection” following a low constraint context, as in (2), as to “collection” when this violates a highly lexically constraining context, as in (3) (Federmeier, Wlotko, De Ochoa-Dewald, & Kutas, 2007; Kutas & Hillyard, 1984).

- (1). He bought her a pearl necklace for her birthday.
- (2). He looked worried because he might have broken his collection.
- (3). He bought her a pearl necklace for her collection.

In contrast to the N400, a set of later positive-going ERP components, visible at the scalp surface between

<sup>1</sup>Tufts University, <sup>2</sup>Harvard Medical School, <sup>3</sup>Moss Rehabilitation Research Institute, Elkins Park, PA

approximately 500 and 1000 msec, do appear to be differentially sensitive to words that violate strong contextual constraints. The initial ERP research characterizing these late positivities came from different research groups working within different theoretical frameworks. One set of studies focused on a late *posteriorly* distributed positivity (maximal at parietal and occipital sites), otherwise known as the P600. This *late posterior positivity/P600* was initially characterized as a response produced by syntactic anomalies or syntactically dispreferred continuations (Hagoort, Brown, & Groothusen, 1993; Osterhout & Holcomb, 1992), but it was subsequently noted that, under certain conditions, it is also evoked by semantic incongruities—the “semantic P600” (e.g., “Every morning at breakfast the eggs would \*eat...”; Kuperberg, Sitnikova, Caplan, & Holcomb, 2003; for a review, see Kuperberg, 2007).<sup>2</sup> Another set of studies focused on a late *frontally* distributed positivity (maximal at prefrontal and frontal sites), which is typically evoked by words that violate strong lexical constraints, such as in sentence (3) above (e.g., Van Petten & Luka, 2012; Federmeier et al., 2007; see also Kutas, 1993).

The studies that characterized these late positive effects used different sets of stimuli with different syntactic structures, different positions of critical words, and different tasks. This made it difficult to compare the effects across different studies. However, based on a review of the early literature, Van Petten and Luka (2012) noted that, although both the late posterior and frontal positivity effects were associated with unexpected linguistic input, the main factor that distinguished them was the plausibility of the resulting interpretation. The late posterior positivity/P600 was produced by highly implausible words. In fact, as noted by Kuperberg (2007, Section 3.4, p. 32), it is typically produced by semantically *anomalous* words that result in an *impossible* interpretation (e.g., Paczynski & Kuperberg, 2012; van de Meerendonk, Kolk, Vissers, & Chwilla, 2010; Kuperberg, Kreher, Sitnikova, Caplan, & Holcomb, 2007; Kuperberg et al., 2003). In contrast, the *late frontal positivity* is produced by unexpected but plausible critical words. This distinction was confirmed in recent studies showing that, in lexically constraining contexts, the same individuals who produced a late frontal positivity effect on unexpected plausible critical words produced a late posterior positivity effect on unexpected highly implausible critical words (in English by DeLong, Quante, & Kutas, 2014; in German by Quante, Bölte, & Zwitserlood, 2018; in Hebrew by Ness & Meltzer-Asscher, 2018).

Although the research characterizing the late posterior positivity/P600 effect and the late frontal positivity effect has proceeded somewhat independently, the mechanisms proposed to underlie each of these effects are quite similar. They include the detection of conflict between the predicted and bottom-up input (late frontal positivity: DeLong, Urbach, Groppe, & Kutas, 2011; late posterior positivity/P600: van de Meerendonk, Kolk, Chwilla, &

Vissers, 2009; Kuperberg, 2007), the suppression of incorrectly predicted information and enhancement of activity associated with the bottom-up input (late frontal positivity: Federmeier et al., 2007; Kutas, 1993; late posterior positivity/P600: van de Meerendonk et al., 2009), and prolonged attempts to integrate the bottom-up input to reach a new higher level interpretation (late frontal positivity: Brothers, Swaab, & Traxler, 2015; DeLong et al., 2014; Federmeier, Kutas, & Schul, 2010; late posterior positivity/P600: Brouwer, Fitz, & Hoeks, 2012; Kuperberg, 2007).

This leaves open many questions. For example, precisely what representations must be predicted and violated to trigger each effect? And exactly why or how are such prediction violations linked to comprehenders’ interpretation of the input as plausible or impossible? The primary goal of this study was to begin to address these questions and to bring the literatures discussing the late frontal positivity, the late posterior positivity/P600, and the N400 effects together under a common theoretical umbrella. We aimed to dissociate the two late positivity effects from one another and from the N400 effect by explicitly manipulating predictions that are fulfilled and violated at different grains and levels of representation in a single experiment.

In explaining the logic of our design, we assume that the comprehender’s draws upon a *communication model*—a set of linguistic and non-linguistic knowledge that describes her beliefs about the communicator and the broader communicative environment. We assume that this communication model is of a literal English speaker who is communicating information about events and states that are possible in the real world (cf. Degen, Tessler, & Goodman, 2015; Frank & Goodman, 2012). We focus on three hierarchical levels of representation within this communication model. At the top of the hierarchy is the *situation model*—a high-level representation of meaning, established during deep comprehension, that describes the full set of events, actions, and characters being communicated (Zwaan & Radvansky, 1998; Van Dijk & Kintsch, 1983). Below the situation model is the *event level*, which represents information about the *event structures* (sets of events that are compatible with the communication model) and the specific events or states that are currently being communicated. The third “semantic feature” level represents the semantic features and properties of the individual words and sets of words associated with these events and event structures.

We assume that these three levels are distinguished in at least two ways. The first is by the timescale of the linguistic input with which they interact. In general, it requires more time (and linguistic information) to build a rich situation model than to infer a single event; similarly, it requires more time/linguistic information to infer a whole event than the meaning of a single word. Because of this, information that is fully encoded at a higher level of the hierarchy will subsume information encoded at a lower level. Second, we

assume that information represented at each of these different levels is encoded in a different form and therefore that it interacts with different *types* of linguistic information. For example, the semantic features level can receive input from phonological or orthographic inputs that are used to decode the semantic features associated with specific words. The event level can interact with sets of semantic features (e.g., animacy information) and certain syntactic cues, which provide information about an event structure, as well as with finer-grained features that provide information about a specific event or state.

We assume that, during language comprehension, as new linguistic information becomes available and incrementally decoded, it is passed up the hierarchy in a bottom-up fashion. In addition, information can flow down the hierarchy in a top-down fashion. Specifically, the situation model influences the probability/degree of activation over specific events. Similarly, the probability/degree of activation of specific events can influence the probability/degree of activation over semantic features represented at the lower semantic features level. Within this dynamic framework, if information from a higher level reaches and changes the state of information represented at a lower level before new bottom-up linguistic information arrives and is decoded at that lower level, we refer to this as top-down *prediction*. In other words, we use the term *prediction* both in a spatial sense (top-down effects from a higher to a lower level of representation) and in a temporal sense (input at time point  $t$  changes the state of activity at  $t + 1$ ).

We conceptualize the N400 as reflecting access to the semantic features associated with new bottom-up inputs that have not already been predicted—that is, as changes in activity at the semantic feature level that are induced by new inputs, with smaller changes associated with inputs whose semantic features have already been predicted by the prior context (see Kuperberg, 2016). We hypothesize that the late frontal positivity reflects a large change in activity associated with successfully updating the comprehender's prior situation model. We suggest that the late posterior positivity/P600 is evoked when new bottom-up input *conflicts* with the constraints of the communication model itself, which prevents the input from being initially incorporated into the current situation model. This conflict and resulting initial interpretive failure may lead to additional second-pass attempts to make sense of the input (see Shetreet, Alexander, Romoli, Chierchia, & Kuperberg, 2019, for recent discussion).

### Design of the Current Study

As a step toward testing this theory, we created a set of well-controlled three-sentence discourse scenarios. These contexts were written so that, just before the onset of a critical word, they constrained for different types of

information at the levels of representation described above. This is illustrated schematically in Figure 1.

1A. *Low constraint context*: Eric and Grant received the news late in the day. They mulled over the information and decided it was better to act sooner rather than later. Hence, they cautioned the... (trainees/drawer).

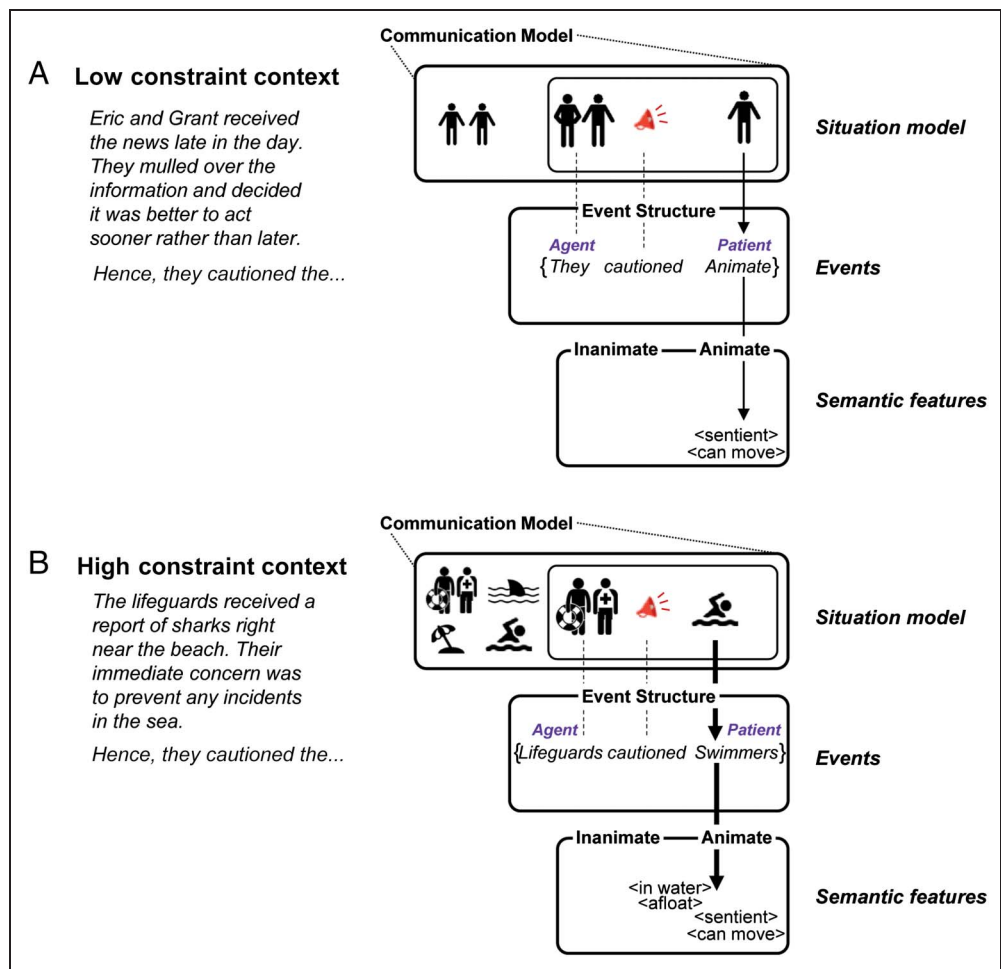
1B. *High constraint context*: The lifeguards received a report of sharks right near the beach. Their immediate concern was to prevent any incidents in the sea. Hence, they cautioned the... (swimmers/trainees/drawer).

*Low constraint* contexts established a relatively simple situation model. For example, in 1A, the first two sentences simply introduce two unknown characters, Eric and Grant, who are considering taking an unknown action (see Figure 1A, left side of situation model). The first few words of the third sentence (“They cautioned the...”) establishes that Eric and Grant are about to caution an unknown person/people (see Figure 1A, right side of situation model). Although this situation model constrains strongly for a particular event structure,  $\{Agent\} cautioned animate-Patient$ , it does not constrain for any specific event.<sup>3</sup> The  $\{Agent\} cautioned animate-Patient$  event structure, in turn, constrains for semantic features that are typically associated with animate entities (e.g.,  $\langle sentient \rangle$ ,  $\langle can\ move \rangle$ ).

*High constraint* contexts established a rich situation model. For example, in 1B, the first two sentences establish the presence of lifeguards and sharks in a beach scene, which is likely to include swimmers (see Figure 1B, left side of situation model). This means that the same first few words of the third sentence (“They cautioned the...”) establish that the lifeguards are highly likely to caution a group of swimmers (Figure 1B, right side of situation model). Thus, in these high constraint contexts, the situation model constrains not only for a particular event structure,  $\{Agent\} cautioned animate-Patient$ , but also for a specific event,  $\{Lifeguards\} cautioned Swimmers$ . This specific event, in turn, constrains not only for semantic features that are typically associated with animate entities (e.g.,  $\langle sentient \rangle$ ,  $\langle can\ move \rangle$ ), but also for features that are more specifically associated with the lexical item, *swimmers* (e.g.,  $\langle in\ water \rangle$ ,  $\langle afloat \rangle$ ).

Following each context, we introduced critical words in the direct object position of the third sentence, which either confirmed or violated each of these constraints. This yielded five types of discourse scenarios (see Table 1 for examples). In the *high constraint expected* scenarios, the critical word (“swimmers”) was highly predictable because it satisfied all top-down constraints at both the event level and the semantic features level. In all four other conditions, the critical words were unpredictable, but each for a different reason. In the *low constraint unexpected* scenarios, the critical word (“trainees”) was

**Figure 1.** Schematic illustration of the state of the language comprehension system at three hierarchical levels of representation after reading a low constraint context and a high constraint context, but before encountering the upcoming critical noun. Please see Design of the Current Study section for full explanation. Icons used in this and Figures 7 and 8 were obtained from thenounproject.com.



plausible but unexpected because the discourse context did not constrain for this item or any other single event or lexical item. The critical word also did not violate event structure/animacy constraints of the prior verb. In the *high constraint unexpected* scenarios, the critical word (“trainees”) was plausible but unexpected because it violated constraints for a specific event, {*Lifeguard cautioned Swimmers*}, and for the specific semantic features associated with the expected word (“swimmers”). Again, it did not violate event structure/animacy constraints. In the *low constraint anomalous* scenarios, the critical word (“drawer”) was anomalous because it violated event structure/animacy constraints—its inanimate features were incompatible with an animate Patient. However, it did not violate constraints for a specific event or for the semantic features associated with a specific lexical item because the discourse context did not constrain strongly for any single event/lexical item. Finally, in the *high constraint anomalous* scenarios, the critical word (“drawer”) produced a “double violation,” violating both event structure/animacy constraints as well as constraints for a specific event/lexical item.

Critical words in the four unpredictable scenarios were matched on their general semantic relatedness with the preceding content words in their contexts, as assessed

using the latent semantic analysis (LSA: Landauer, Foltz, & Laham, 1998; Landauer & Dumais, 1997). Participants read and monitored the coherence of the scenarios, and we measured ERPs on the critical words.

Our first set of questions concerned the N400. Based on numerous previous studies, we expected that the N400 would be smaller on critical words in the high constraint expected scenarios than in all four types of unpredictable scenarios. Based on the work by Kutas and Hillyard (1984) and Federmeier et al. (2007), we also predicted that, in the plausible scenarios, the N400 evoked by unpredictable critical words would be insensitive to the lexical constraint of the context (no difference between the high constraint unexpected and low constraint unexpected conditions). In addition, our design allowed us to ask whether the amplitude of the N400 evoked by semantically anomalous critical words would also be insensitive to the lexical constraint of the context (predicting no difference between the high constraint anomalous and low constraint anomalous conditions). Finally, we were interested in whether the N400 would be larger on the anomalous than the unexpected but plausible critical words. If so, this would provide evidence that the N400 evoked by a given word can be sensitive to its



**Table 1.** Examples of Five Experimental Conditions Created around the Same Verb (Here, “Cautioned”)

<i>Scenario Type</i>	<i>Example</i>	<i>Lexical Constraint of Context</i>	<i>Cloze of Critical Word</i>	<i>SSV<sup>a</sup> between Critical Word and Preceding Context</i>
1. High constraint expected	The lifeguards received a report of sharks right near the beach. Their immediate concern was to prevent any incidents in the sea. Hence, they cautioned the <u>swimmers</u> ...	71% (13%)	71% (13%)	0.18 (0.18)
2. Low constraint unexpected	Eric and Grant received the news late in the day. They mulled over the information and decided it was better to act sooner rather than later. Hence, they cautioned the <u>trainees</u> ...	20% (9%)	0.2% (0.8%)	0.01 (0.05)
3. High constraint unexpected	The lifeguards received a report of sharks right near the beach. Their immediate concern was to prevent any incidents in the sea. Hence, they cautioned the <u>trainees</u> ...	71% (13%)	0.1% (0.5%)	0.01 (0.05)
4. Low constraint anomalous	Eric and Grant received the news late in the day. They mulled over the information and decided it was better to act sooner rather than later. Hence, they cautioned the <u>drawer</u> ...	20% (9%)	0% (0%)	0.01 (0.06)
5. High constraint anomalous	The lifeguards received a report of sharks right near the beach. Their immediate concern was to prevent any incidents in the sea. Hence, they cautioned the <u>drawer</u> ...	71% (13%)	0% (0%)	0.01 (0.04)

The critical word in each of the example final sentences is underlined (although this was not the case in the experiment itself). The final sentences continued with additional words, as indicated by the three dots. Means are shown with standard deviations in parentheses.

<sup>a</sup>SSV: Semantic similarity values. The cosine similarity between vector corresponding to the critical word and the vector corresponding to the preceding context, extracted using LSA using the pairwise comparison tool ([term-to-document](http://lsa.colorado.edu/)) at [lsa.colorado.edu/](http://lsa.colorado.edu/) (possible range of SSVs: +1 to -1).

implausibility (here, operationalized by whether that word matched or mismatched the verb’s animacy selection restrictions), even when its cloze probability and its semantic relatedness with its preceding content words are matched across conditions (see Paczynski & Kuperberg, 2011, 2012).

Our second set of questions concerned the late positivities. We hypothesized that we would be able to dissociate the scalp distribution of these positivities based on whether they violated strong predictions for a single event/lexical item or for a whole event structure. Specifically, we hypothesized that the high constraint unexpected critical words, which violated strong constraints for a specific event/lexical item, would selectively evoke a late frontal positivity (larger than that produced by critical words in any other condition), whereas the low constraint anomalous critical words, which violated strong constraints for an event structure/animacy features, would selectively evoke a late posterior positivity/P600 (larger than that produced by critical words in any of the

plausible scenarios). This would provide strong evidence that distinct neurocognitive mechanisms can be triggered by violations at these different grains of representation.

Finally, our design allowed us to ask not only whether the neurocognitive mechanisms underlying the late frontal positivity and the late posterior positivity/P600 are distinct, but also whether they are independent of one another. This would, in turn, constrain our understanding of the precise levels of representation at which violations must be detected to trigger each of these effects (see Discussion). The key condition to address this question was the doubly violated scenarios (high constraint anomalous) in which the critical words violated both constraints for a specific event/lexical item as well as for a broader event structure. One possibility is that the neurocognitive processes underlying the late frontal positivity and the late posterior positivity/P600 are independent of one another. For example, if the late frontal positivity only reflects the detection of a violated prediction at the level of semantic features whereas the late posterior positivity/P600 reflects

the detection of a violated event structure, then, assuming additivity of independent ERP effects (see Osterhout & Nicol, 1999), one might expect the doubly violated critical words to evoke both a late frontal positivity and a late posterior positivity/P600. If, on the other hand, the neurocognitive processes underlying the late frontal positivity and the late posterior positivity/P600 are interdependent, then we should see non-additive effects. For example, if the detection of an event structure violation interferes with processes reflected by the late frontal positivity, then we should see a late posterior positivity/P600 but no late frontal positivity on the double violations.

## METHODS

### Development and Norming of Materials

Participants read the five types of scenarios discussed above (see Table 1). In all scenarios, the first two sentences introduced a context. The third sentence began with an adjunct phrase of one to four words, followed by a pronominal subject that referred back to the first two sentences, followed by a verb, a determiner, a “critical word” (always a direct object noun), and then three additional words. The scenarios varied by the constraint of the context (i.e., the combination of the first two sentences and the third sentence until just before the critical word) and by whether the critical word matched or violated the event/lexico-semantic constraints of the high constraint contexts and/or the event structure/animacy constraints defined by the thematic-semantic properties of the previous verb.

To construct these scenarios, we began with a set of 100 preferentially transitive verbs (50% selecting for animate direct objects, 50% for inanimate direct objects), which, in the absence of a discourse context, were not highly predictive of any particular upcoming event or noun. The constraints of these verbs in minimal context were assessed in an offline cloze norming study (see below). We then wrote high and low constraint two-sentence discourse contexts around each verb (matching the average number of words across the two levels of constraint) and carried out a second cloze norming study of these discourse contexts (see below). Based on the results of this cloze norming, we created the five scenario types. To create the high constraint expected scenarios, each high constraint context was paired with the noun with the highest cloze probability for that context. To create the high constraint unexpected scenarios, each high constraint context was paired with a direct object noun of zero (or very low) cloze probability, but that was still plausible in relation to this context. To create the low constraint unexpected scenarios, the same unexpected noun was paired with the low constraint context. To create the anomalous scenarios, each high constraint and low constraint context was paired with a noun that violated the animacy selectional restrictions of the verb.

Thus, one of the five conditions had predictable critical words, and four had unpredictable critical words.

In constructing these scenarios, we used LSA (Landauer et al., 1998; Landauer & Dumais, 1997) to match the semantic similarity between the critical words and the entirety of the prior contexts across the four unpredictable conditions. To carry out this analysis, we used the pairwise comparison tool from [lsa.colorado.edu/](http://lsa.colorado.edu/) with default values (local and global weighting functions, 300 latent semantic dimensions). We extracted pairwise term-to-document semantic similarity values (SSVs)—cosine similarities between the vectors corresponding to our critical words and “pseudo-document” vectors corresponding to each of our contexts. All of our critical words appeared in the original term–document matrix, except for five for which we substituted close synonyms. The mean SSVs for each of the five types of scenario are shown in Table 1. SSVs showed no significant difference across the four unpredictable conditions,  $F(3, 396) = 0.173$ ,  $p = .914$ . Unsurprisingly, the SSVs were significantly greater in the predictable scenarios than in each of the four unpredictable scenarios (all pairwise comparisons,  $ps < .001$ ,  $ts > 9$ ).

### Cloze Norming Studies

As noted above, we carried out two cloze norming studies to construct the stimuli. For both studies, participants were recruited through the crowd-sourcing platform, Amazon Mechanical Turk. They were asked to complete each context with the first word that came to mind (Taylor, 1953) and, in an extension of the standard cloze procedure, to also provide two additional words that could complete the sentence (see Federmeier et al., 2007; Schwanenflugel & LaCount, 1988). Responses were excluded from participants who reported that their first language learned was not English or if they reported any psychiatric or neurological disorders. Responses were also excluded from any participants who failed to follow instructions (“catch” questions were used as periodic attention checks).

*Cloze Norming Study 1: Selection of non-constraining verbs.* We started with a set of 600 transitively biased verbs, compiled from various sources including Levin (1993) and materials from previous studies conducted in our laboratory (Paczynski & Kuperberg, 2011, 2012). Verbs with log Hyperspace Analogue to Language frequency (Lund & Burgess, 1996) of two standard deviations below the mean (based on English Lexicon Project database; Balota et al., 2007) were excluded. For each verb, we constructed a simple active, past tense sentence stem that consisted of only a proper name, the verb, and a determiner (e.g., “Harry explored the...”). These sentences were divided into six lists to decrease the time demands on any individual participant during cloze norming. Between 89 and 106 participants

(depending on list) who met inclusionary criteria provided completions for each verb.

The lexical constraint of each verb in these minimal contexts was calculated by identifying the most common completion across participants and tallying the proportion of participants who provided this completion. The set of verbs used in the final set of stimuli had an average constraint of 16.4% ( $SD = 0.7\%$ ).

*Cloze Norming Study 2: Selection and characterization of the final set of discourse stimuli.* We began with an initial set of 198 pairs of high constraint and low constraint contexts (the combination of the first two sentences and the third sentence, including the verb and the determiner). These were pseudorandomly divided into two lists, such that each list contained only one of the two contexts associated with each verb. The two lists were then divided into thirds to decrease time demands on any individual participant during cloze norming. Between 51 and 69 participants who met inclusionary criteria provided completions for each scenario.

Cloze probabilities were calculated based on the percentage of respondents providing the critical noun used in the experiment. Alternate word forms (e.g., singular/plural) were collapsed, but synonyms or lexical alternatives were not collapsed (e.g., couch/sofa). The lexical constraints of the full discourse contexts were calculated as described above. The mean and standard deviation of the cloze probability and the lexical constraint for each scenario type is given in Table 1.

#### *Counterbalancing across Lists*

The final set of 500 scenarios (100 sets of five) was rotated across the five lists, with 20 items per condition in each list. Counterbalancing ensured that the same adjunct phrase plus verb in each of the final sentences appeared in all conditions across the five lists, and the same critical words were counterbalanced across the four types of unpredictable scenarios (low constraint unexpected, high constraint unexpected, low constraint anomalous, high constraint anomalous). The critical words in the high constraint expected scenarios differed from those in the four types of unpredictable scenarios (see Supplementary Materials, Table S1 for information about the lexical characteristics of these critical words).<sup>4</sup> The same high constraint contexts were counterbalanced across the high constraint expected, high constraint unexpected, and high constraint anomalous scenarios, and the same low constraint contexts were counterbalanced across the low constraint unexpected and low constraint anomalous scenarios.

To each list, we then added an additional 20 high constraint anomalous, 20 low constraint unexpected, and 20 low constraint anomalous filler scenarios so that each participant saw 160 scenarios in total. This ensured that each participant saw an equal number of high and low constraint contexts and that, at each level of

contextual constraint (high constraint vs. low constraint), half of the scenarios were plausible and the other half were anomalous.

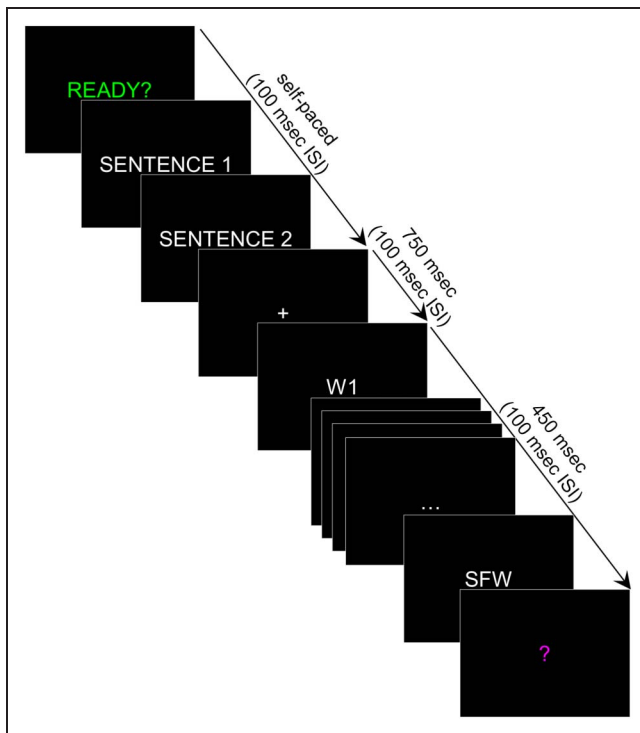
#### **Participants**

We report data from 39 native English speakers (age = 18–32 years, mean = 21.6,  $SD = 3.6$ , 21 men). Forty participants were originally recruited, but one failed to complete the session. Participants were recruited from Tufts University and the surrounding communities. They were screened on the basis of the following exclusion criteria: significant exposure to any language other than English before the age of 5 years, history of psychiatric or neurological diagnoses or injury, and use of psychoactive medication within the preceding 6 months. Participants were right-handed, as assessed by the Edinburgh Handedness Inventory (Oldfield, 1971). They provided written informed consent and were paid for their time, and all protocols were approved by Tufts University Social, Behavioral, and Educational Research Institutional Review Board.

#### **Stimulus Presentation and Task**

Participants sat in a dimly lit room while stimuli were presented approximately 150 cm from the computer screen. Stimuli were presented using PsychoPy 1.83 software (Peirce, 2007) and were displayed on an LCD monitor using white Arial font set to 0.1 of the screen height on a black background. Details of the structure of each trial are given in Figure 2 and legend.

Participants' task was to press one of two buttons after seeing a "?" cue to indicate whether they judged each scenario to "make sense" or not. This task encouraged active coherence monitoring during online comprehension and was intended to prevent participants from completely disregarding the anomalies (see Sanford, Leuthold, Bohan, & Sanford, 2011, for evidence that detecting anomalies is necessary to produce a late positivity/P600 effect at all, and see [projects.iq.harvard.edu/kuperberg/lab/additional-materials/P600-task-discussion](http://projects.iq.harvard.edu/kuperberg/lab/additional-materials/P600-task-discussion) for a discussion of the use of this task in ERP studies). In addition, following approximately 32 trials of the 160 trials (all fillers), participants answered a yes/no comprehension question about the preceding scenario. For example, the scenario, "Damien was not at all like himself that day. Something in the atmosphere made him act differently, and everyone could tell. Relentlessly, he ridiculed the patients until they cried." was followed by the question, "Is this behavior typical of Damien?" This encouraged participants to engage in deep comprehension, attending to the scenarios as a whole rather than just the third sentence in which the anomalies appeared. The experiment was divided into four blocks, with block order randomized across participants. Participants were given 10 practice trials before the main experiment.



**Figure 2.** Presentation of stimuli. Each trial began with the prompt “READY?” (centered in green font). When participants initiated the trial with a button press, the first two sentences were presented successively in full on the screen. Participants advanced through these sentences at their own pace using a button press. Then a centered fixation cross appeared (750 msec; 100 msec ISI), followed by the third sentence, which was presented word-by-word (each word 450 msec; 100 msec ISI). After the sentence-final word (SFW), a “?” appeared (centered, magenta font), prompting participants to make an acceptability judgment about whether or not the preceding scenario “made sense” or not. The “?” remained on the screen until the response was registered.

### EEG Recording and Processing

EEG was recorded using a Biosemi ActiveTwo acquisition system from 32 active electrodes in a modified 10/20 system montage. Signals were digitized at 512 Hz and a pass-band of DC at 104 Hz. All processing was carried out using EEGLab (Delorme & Makeig, 2004) and ERPLab (Lopez-Calderon & Luck, 2014) analysis packages in the MATLAB environment. After importing the data for processing, the EEG was referenced offline to the average of the left and right mastoid channels. Both a high-pass (0.05 Hz) and a low-pass (30 Hz) 24 dB per octave filter were applied offline to the continuous data.

The EEG was then segmented into initial epochs spanning from  $-300$  msec before until  $+1400$  msec after the presentation of critical words. Independent components analysis (using the extended infomax algorithm implemented in EEGLab) was then applied to the epoched data to correct for blinks and other eye movement artifact. Then, trials with any residual artifacts were rejected using a semiautomated procedure with participant-specific artifact detection thresholds. After artifact rejection, there

remained on average 18.6 trials per condition ( $SD = 1.3$ ). Artifact rejection rates did not differ significantly across the five conditions,  $F(4, 152) = 1.91, p = .13$ . Our epoch of interest was  $-100$  to  $1000$  msec.

We then extracted single trial artifact-free ERP data, using a baseline of  $-100$  to  $0$  msec, by averaging across electrode sites and time windows that defined specific spatiotemporal ROIs that we selected *a priori* to operationalize each of our three ERP components of interest (see Figure 3). Based on numerous previous studies, we operationalized the N400 as the average voltage within the central region between 300 and 500 msec. In the prior literature, the two late positivities have shown variable time courses (e.g., DeLong et al., 2014; Paczynski & Kuperberg, 2011, 2012; Federmeier et al., 2007; Kuperberg et al., 2003), and so, conservatively, we operationalized both of these components as the average voltage across a long 600–1000 msec time window. For the late frontal positivity, we averaged across electrode sites within the prefrontal region (based on Federmeier et al., 2007), and for the late posterior positivity/P600, we averaged across electrode sites within the posterior region (based on Paczynski & Kuperberg, 2011, 2012; Kuperberg et al., 2003).

### ERP Statistical Analysis

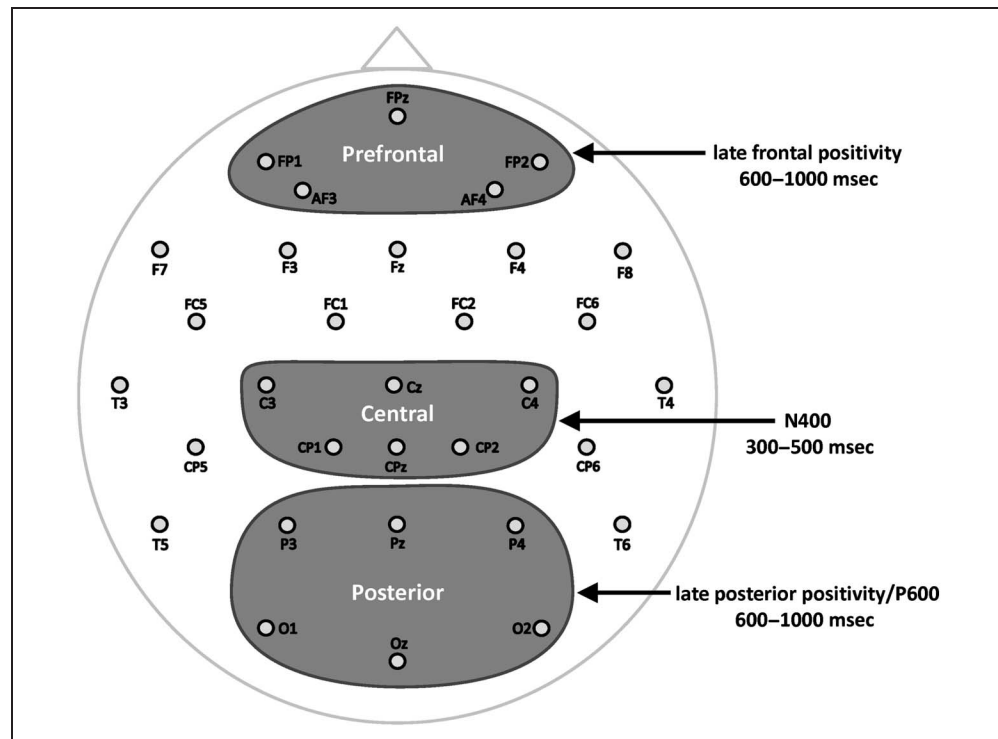
We report the results of a series of linear mixed-effect regression models, fit in R Version 3.2.4 (R Core Team, 2016) using the *lme4* package Version 1.1-11 (see Bates, Mächler, Bolker, & Walker, 2015). Our dependent measure was the average trial-level ERP amplitude within each of the spatiotemporal ROIs described above.<sup>5</sup> The maximal random effects structure was used (Barr, Levy, Scheepers, & Tily, 2013) with by-item and by-subject random intercepts and by-item and by-subject random slopes for each fixed effect of interest. For each contrast of interest,  $p$  values were estimated using a Satterthwaite approximation, as implemented by the *lmerTest* package Version 2.0-30 (Kuznetsova, Brockhoff, & Christensen, 2015).

On the N400, we were primarily interested in whether, across the four unpredictable conditions, there were main effects of Constraint, Plausibility, and/or an interaction between these two factors, and so we constructed a model that crossed two levels of Constraint (high constraint, low constraint) and two levels of Plausibility (plausible, anomalous). We also carried out pairwise comparisons between each of the four unpredictable conditions and the high constraint expected condition to confirm the presence of significant N400 effects.

On the late positivities, we were primarily interested in the contrasts between specific unpredictable conditions and all other conditions. We therefore proceeded straight to planned pairwise comparisons. For the late frontal positivity, we were interested in the comparison between the high constraint unexpected critical words and each of the four other conditions, and for the late posterior



**Figure 3.** Spatiotemporal ROIs used for analysis. The N400 was operationalized as the average voltage between 300 and 500 msec across all electrode sites within the central region. The late frontal positivity was operationalized as the average voltage (600–1000 msec) in the prefrontal region. The late posterior positivity/P600 was operationalized as the average voltage (600–1000 msec) in the posterior region.



positivity/P600, we were interested in the comparison between the anomalous critical words (both high constraint anomalous and low constraint anomalous) and each of the three plausible conditions. For both late positivities, we carried out additional pairwise comparisons between all remaining conditions to determine the specificity of any effects.

## RESULTS

### Behavioral Findings

On average, participants correctly judged the acceptability of the scenarios on 89.6% of trials. Participants were most accurate in responding “YES” to the high constraint expected scenarios (97.6%), less accurate in responding “NO” to the high constraint anomalous scenario (91.0%) and the low constraint anomalous scenario (89.8%), and least accurate in responding “YES” to the high constraint unexpected and low constraint unexpected scenarios (both 86.3%). Participants correctly answered 86.4% of the comprehension questions on average (27.6 of 32 questions,  $SD = 2.0$ ), indicating that they attended to the full contexts of the scenarios, rather than just to the final sentences.

### ERP Findings

#### N400

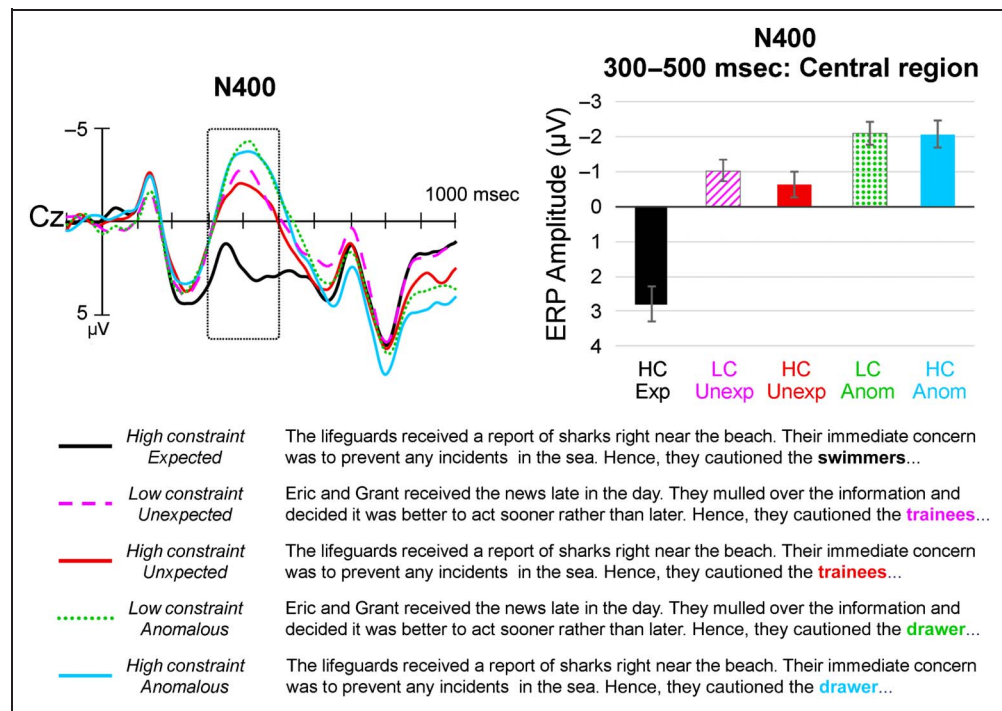
Figure 4 shows grand-averaged ERP waveforms at electrode site, Cz, for each of the five conditions. It also

shows the mean voltage (averaged across the central 300–500 msec spatiotemporal region used to operationalize the N400) for each condition. As expected, the N400 was smaller (less negative) on critical words in the high constraint expected scenarios than on critical words in each of the four types of unpredictable scenario (all  $t_s > 6.0$ , all  $p_s < .001$ ). Across the four types of unpredictable scenarios, a model that crossed Constraint (high constraint, low constraint) and Plausibility (plausible, anomalous) showed no main effect of Constraint ( $t = -0.54$ ,  $p = .59$ ), no interaction between Constraint and Plausibility ( $t = 0.31$ ,  $p = .75$ ), but a significant main effect of Plausibility ( $t = -3.08$ ,  $p = .003$ ) due to a slightly larger (more negative) N400 on the anomalous than the plausible critical words. Voltage maps for all pairwise comparisons are shown in Supplementary Materials, Figure S1, and the full set of statistical comparisons between all pairs of conditions is reported in Supplementary Materials, Table S2.

#### Late Frontal Positivity

Figure 5 (top) shows grand-averaged ERP waveforms at electrode site FPz in each of the five conditions. It also shows the mean voltage (averaged across the prefrontal 600–1000 msec spatiotemporal region used to operationalize the late frontal positivity) for each condition. As hypothesized, the late frontal positivity was larger (more positive) to critical words that violated strong event/lexical constraints (high constraint unexpected scenarios) than to critical words in all other conditions (all  $t_s > 2.3$ ,

**Figure 4.** N400. Left: Grand-averaged ERP waveforms at electrode Cz in each of the five conditions. Right: Mean N400 amplitude in each of the five conditions, with the N400 operationalized as the average voltage across all time points between 300 and 500 msec across all electrode sites within the central ROI. Negative voltage is plotted upward. Error bars represent  $\pm 1 SEM$ , calculated within subjects (Morey, 2008).



all  $ps < .02$ ). This effect was selective: Pairwise contrasts that did not include the high constraint unexpected condition failed to show any significant effects on the late frontal positivity (all  $ts < 1.7$ , all  $ps > .11$ ). Of particular note, no significant late frontal positivity effect was produced by the critical words in the high constraint anomalous scenarios, which violated both event/lexical and event structure/animacy constraints, relative to the low constraint anomalous, the high constraint expected, or the low constraint unexpected critical words (all  $ts < 1.0$ ,  $ps > .49$ ). Voltage maps that illustrate the scalp distribution of the late frontal positivity effects are shown in Figure 6A. The full set of statistical comparisons between all pairs of conditions is reported in Supplementary Materials, Table S3.

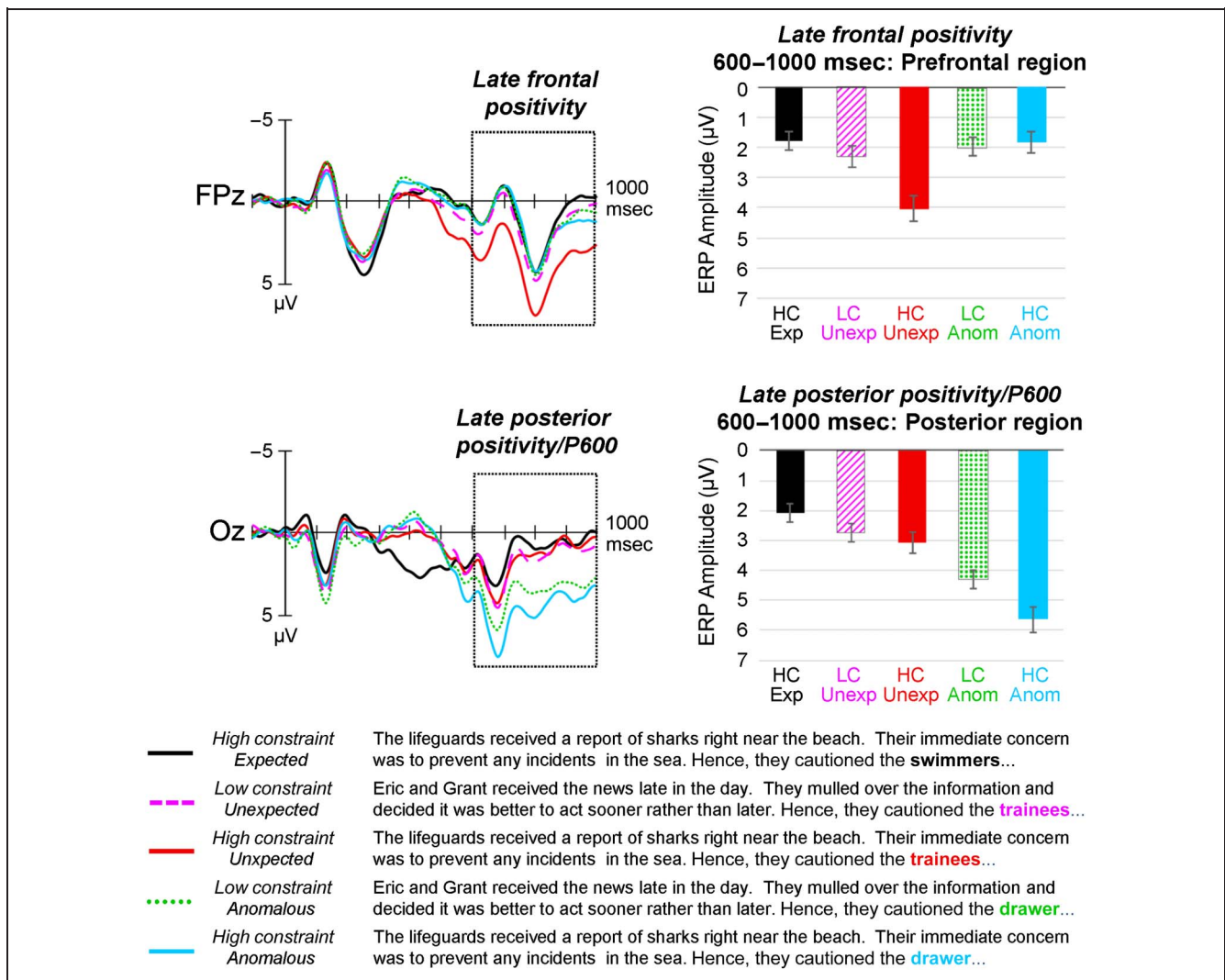
#### Late Posterior Positivity/P600

Figure 5 (bottom) shows grand-averaged ERP waveforms at electrode site Oz in each of the five conditions. This figure also shows the mean voltage (averaged across the posterior 600–1000 msec spatiotemporal region used to operationalize the late posterior positivity/P600) for each condition. As hypothesized, the late posterior positivity/P600 was larger (more positive) to anomalous critical words that violated event structure/animacy constraints (both high constraint anomalous and low constraint anomalous) than to plausible critical words (high constraint expected, low constraint unexpected, and high constraint unexpected critical words), all  $ts > 3.7$ , all  $ps < .001$ . Again, this effect was selective: Pairwise comparisons that did not include the anomalous conditions failed to reveal any significant effects on the

late posterior positivity/P600 (all  $ts < 1.83$ , all  $ps > .07$ ). Of particular note, the high constraint anomalous scenarios, which violated both event/lexical and event structure/animacy constraints, produced a larger late posterior positivity/P600 than the low constraint anomalous scenarios, which violated just event structure/animacy constraints ( $t = 2.31$ ,  $p = .03$ ). Voltage maps that illustrate the scalp distribution of the late posterior positivity/P600 effects are shown in Figure 6B. The full set of statistical comparisons between all pairs of conditions is reported in Supplementary Materials, Table S4.

## DISCUSSION

We used a single set of well-controlled stimuli to examine neural activity in response to words that confirmed and violated constraints at different levels of representation during language comprehension. Our findings were clear. Words that confirmed predictions of semantic features produced an attenuated N400; words that violated event/lexical constraints produced an enhanced late frontal positivity, whereas words that violated event structure/animacy constraints elicited an enhanced late posterior positivity/P600. These findings suggest that the brain recognized and distinguished the different levels and grains of representation that were predicted and violated as a consequence of our experimental manipulations. In the sections below, we discuss these findings in more detail and how they can be accommodated within a “hierarchical generative framework” of language comprehension (see Kuperberg, 2016; Kuperberg & Jaeger, 2016; see also Xiang & Kuperberg, 2015). We then



**Figure 5.** The late positivities. Top: The late frontal positivity. Top left: Grand-averaged ERP waveforms at electrode FPz in each of the five conditions. Top right: Mean late frontal positivity amplitude in each of the five conditions, with the late frontal positivity operationalized as the average voltage across all time points between 600 and 1000 msec across all electrode sites within the prefrontal ROI. Bottom: The late posterior positivity/P600. Bottom left: Grand-averaged ERP waveforms at electrode Oz in each of the five conditions. Bottom right: Mean late posterior positivity/P600 amplitude in each of the five conditions, with the late posterior positivity/P600 operationalized as the average voltage across all time points between 600 and 1000 msec across all electrode sites within the posterior ROI. Negative voltage is plotted upward. Error bars represent  $\pm 1$  SEM, calculated within subjects (Morey, 2008).

return to the functional significance of each of these ERP components and discuss the broader theoretical implications of our findings.

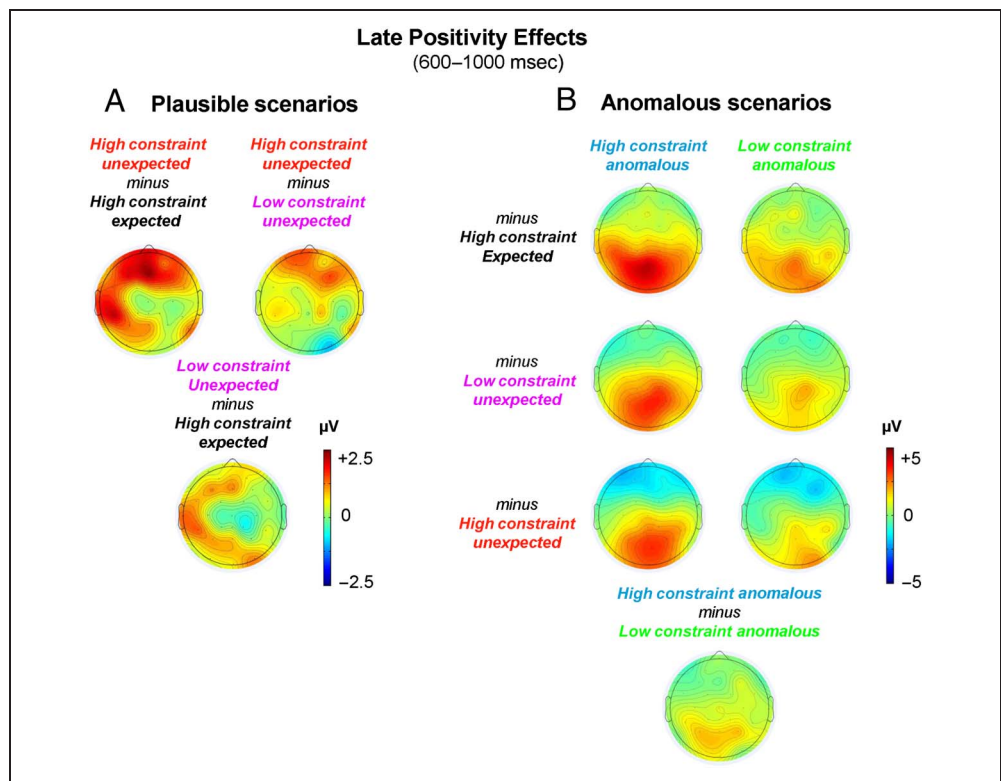
### The N400: Confirmed Predictions at the Level of Semantic Features

As expected, the N400 was significantly smaller to semantically predictable critical words in the high constraint expected scenarios than to unpredictable words in each of the four other conditions. Importantly, however, the N400 evoked by the four types of unpredictable words was not sensitive to the lexical constraint of the preceding context. Replicating the studies of Kutas and Hillyard (1984) and Federmeier et al. (2007) in a new stimulus set,

the amplitude of the N400 evoked by critical words in the high constraint unexpected and low constraint unexpected scenarios was indistinguishable. We further show that this insensitivity of the N400 to lexical constraint extended to anomalous critical words: The amplitude of the N400 to critical words in the high constraint anomalous and low constraint anomalous scenarios was also indistinguishable. Taken together, these findings support the hypothesis that the N400 primarily reflects the degree of match between semantic features that are predicted by a context and the semantic features associated with an incoming word, rather than the detection or recovery from violations of strong lexico-semantic predictions.

In addition, we show that critical words whose semantic features mismatched animacy constraints of their

**Figure 6.** Voltage maps showing the scalp topographies of the late positivity effects. Voltages are averaged across the 600–1000 msec time window for all pairwise contrasts. (A) Left: Voltage maps contrasting each of the plausible conditions with one another. These illustrate the frontal scalp distribution of the late positivity effect produced by the high constraint unexpected critical words, relative to the high constraint expected critical words and relative to the low constraint unexpected critical words. The scalp distribution of the effect produced by the low constraint unexpected, relative to the high constraint expected critical words, is similar, but the effect was not statistically significant. (B) Right: Voltage maps contrasting the anomalous critical words with critical words in each of the plausible conditions. These illustrate the posterior scalp distribution of the late positivity effect evoked by the anomalous critical words. The effects produced by the high constraint anomalous critical words were larger than those produced by the low constraint anomalous critical words. Note that the late frontal positivity effects (left) and the late posterior positivity/P600 effects (right) are shown at different voltage scales to better illustrate the scalp distribution of each effect.



preceding verbs (both in the high constraint anomalous and the low constraint anomalous scenarios) elicited a slightly larger N400 amplitude than unpredicted but plausible critical words that were consistent with these animacy constraints (low constraint unexpected and high constraint unexpected scenarios).<sup>6</sup> This finding shows that the amplitude of the N400 can be sensitive to when lexical probability (as operationalized by cloze) is low, and over and above the effects of simple semantic relationships between the critical word and the prior context (as operationalized by LSA). Rather than reflecting a post-lexical integrative mechanism, however, we interpret this N400 plausibility effect as arising from the same semantic predictive mechanism that gave rise to the N400 effects of cloze probability—a point that we return to below.

### Late Positivities: Violations of Strong Predictions

In contrast to the N400, the late positive components were selectively modulated by violations of strong predictions. Moreover, the scalp topographies of these positivities differed, depending on the grain of representation that was violated: Critical words that only violated specific event/lexical constraints without violating event structure/animacy constraints (high constraint unexpected) selectively evoked a late frontal positivity, which was larger than that produced by critical words in any

other condition. In contrast, critical words that only violated event structure/animacy constraints (low constraint anomalous) evoked a late posterior positivity/P600, which was larger than that produced by critical words in any of the three plausible scenarios. These findings suggest that the underlying neural sources contributing to each of these effects are at least partially distinct.

Importantly, our findings on the double violations, which violated both event structure/animacy constraints and finer-grained event/lexical constraints, suggest that, although distinct, the neurocognitive mechanisms underlying the late frontal positivity and the late posterior positivity/P600 are interdependent: The doubly violated high constraint anomalous critical words did not produce a late frontal positivity effect and instead only produced a late posterior positivity/P600 effect, which was larger than that produced by the low constraint anomalous critical words. This suggests a trade-off relationship between these two components, which has important implications for their functional interpretation.

With regard to the late frontal positivity, the absence of an effect on the doubly violated critical words suggests that this component cannot simply reflect the detection of a lexical prediction violation. In the high constraint anomalous scenarios, the critical word “drawer” also violated strong lexical constraints. If the late frontal positivity only reflected the detection of a discrepancy between the



predicted and encountered lexical item, or a competitive process operating purely at the level of semantic features, then a frontal positivity should have also been evoked in this condition. The fact that no such effect was observed suggests that the late frontal positivity reflected a higher level process that was engaged only in the high constraint unexpected scenarios. Specifically, we propose that this component indexes a large change in activity associated with successfully updating the higher level situation model—the process of shifting from the prior situation model to a new situation model on the basis of new bottom-up input. As we discuss further below, we also suggest that this higher level shift entailed top-down feedback suppression of the incorrectly predicted semantic features and selection of the correct semantic features.

With regard to the late posterior positivity/P600, which was selectively produced by the anomalous critical words, we similarly suggest that, rather than simply reflecting the detection of an animacy violation, this effect also reflected higher level activity. Specifically, we suggest that it was triggered when the bottom-up input conflicted with the constraints of the existing communication model. This resulted in an (initial) failure to successfully update the current situation model with the new input (see Shetreet et al., 2019, for recent discussion), thereby “blocking” any late frontal positivity effect. Our finding that the late posterior positivity/P600 was larger on the doubly violated critical words (high constraint anomalous > low constraint anomalous) suggests that the comprehender’s original high-certainty prediction for a specific event made it easier to detect conflict (previous work suggests that the detection of anomalies is critical to producing the late posterior positivity/P600; see Sanford et al., 2011). This finding is consistent with previous reports that the late posterior positivity/P600 evoked by syntactic violations is also influenced by the lexical constraint of the preceding context: Its amplitude is larger to syntactic anomalies within high constraint than low constraint contexts (Gunter, Friederici, & Schriefers, 2000).

As we discuss next, this pattern of ERP findings can be understood within a hierarchical generative framework of language comprehension (Kuperberg, 2016; Kuperberg & Jaeger, 2016). Below, we outline this framework and offer functional interpretations of the N400, the late frontal positivity, and the late posterior positivity/P600 components within this framework. We then discuss these interpretations in relation to previous interpretations of these components, previous models of language comprehension, and more general accounts of predictive processing in the brain.

### **A Hierarchical Generative Framework of Language Comprehension**

Within a hierarchical generative framework of language comprehension (see Kuperberg & Jaeger, 2016, Section 5, p. 16; Kuperberg, 2016), the agent draws upon an internal generative model<sup>7</sup>—a hierarchical network of stored

linguistic and nonlinguistic representations that she believes are relevant to her communicative goals and her beliefs about the communicator and the broader communicative environment. In the present study, we assume that the comprehender’s goal is deep comprehension and that she believes that the communicator is an English speaker who is communicating literally about events that are possible in the real world (cf. Degen et al., 2015; Frank & Goodman, 2012). Her generative model will therefore reflect the subset of all her linguistic and non-linguistic knowledge that is consistent with these beliefs—her communication model. At the top of this model, she will construct a “situation model” of the current discourse. Because the linguistic input unfolds linearly over time, and because it takes time to pass new information up from lower to higher levels of representation, the comprehender cannot achieve the goal of comprehension (i.e., construct the situation model) all at once. However, if she has hypotheses at a higher level of representation, then she can test them by propagating probabilistic predictions down to lower levels of the hierarchy, thereby pre-activating information at these lower levels. Then, as new bottom-up information becomes available to each level of representation, any information that confirms these top-down predictions is “explained away” (its processing is facilitated), whereas any information that remains unexplained is propagated up the generative hierarchy where it can be used to update hypotheses represented at higher levels.

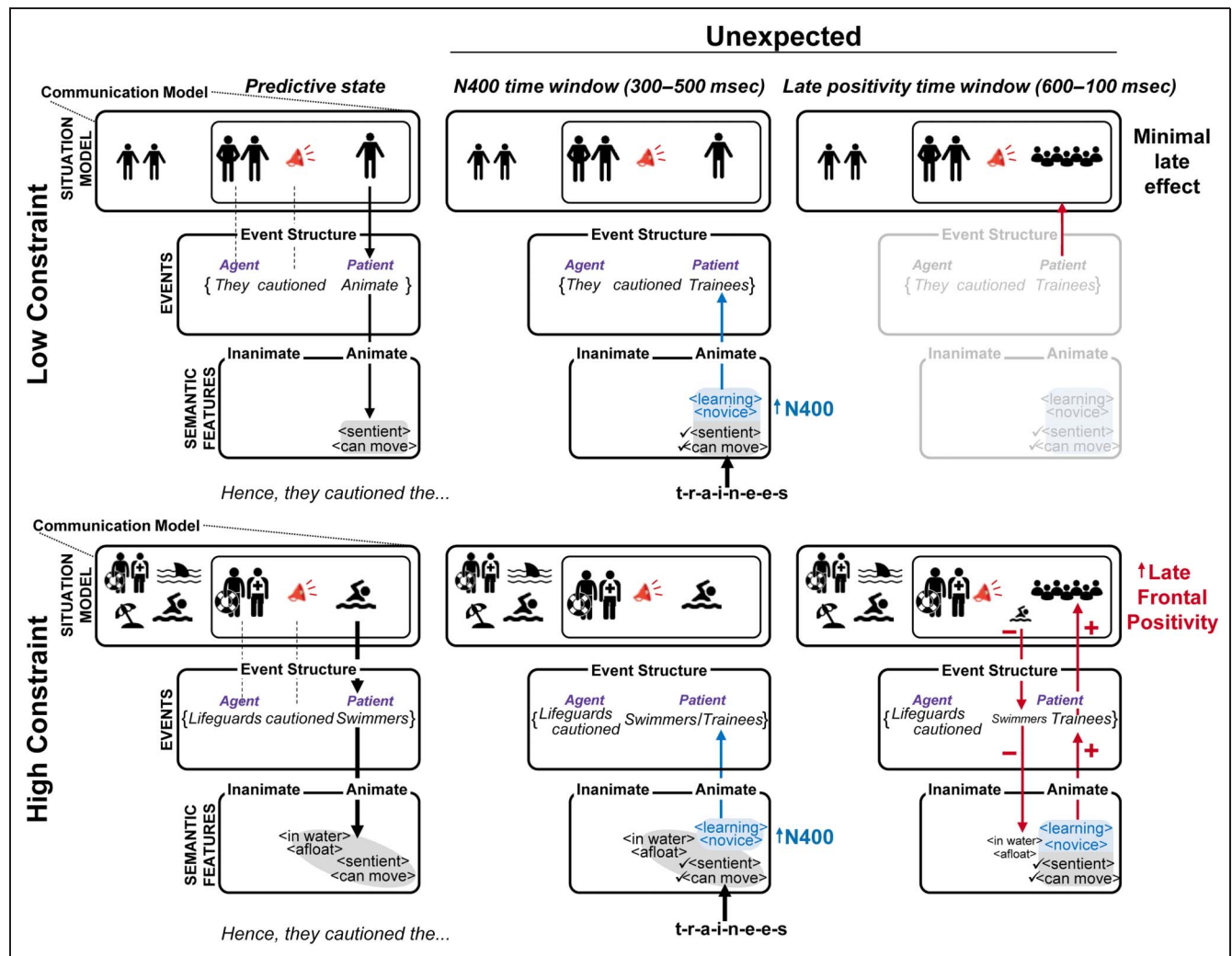
In most situations, at any given point in time, the arrival of unexplained input at the high-level situation model will induce only a small change in its state; larger changes are accumulated gradually over time. Sometimes, however, a particular input may induce a large change in state in the situation model. For example, if at an earlier point in time, the comprehender had already settled, with high certainty, on a particular high-level message that she believed the communicator intended to convey, then the arrival of new unexplained information at the situation model will lead to a large change in state, corresponding to the large shift in the comprehender’s interpretation.

So, as long as the bottom-up input is compatible with the constraints of the communication model, then, through incremental cycles of prediction and hypothesis updating, the comprehender should be able to “home in” on the communicator’s intended message with increasing certainty, as the linguistic input progressively unfolds over time. However, if the comprehender receives unpredicted information that conflicts with the constraints of the communication model itself, then these cycles of updating the situation model come to a halt, leading to a temporary comprehension failure. This may trigger prolonged second-pass attempts to make sense of the discourse scenarios through reanalysis, attempts to repair the input, and/or attempts to revise or adapt the communication model itself. We now illustrate these principles with the example stimuli used in this study.

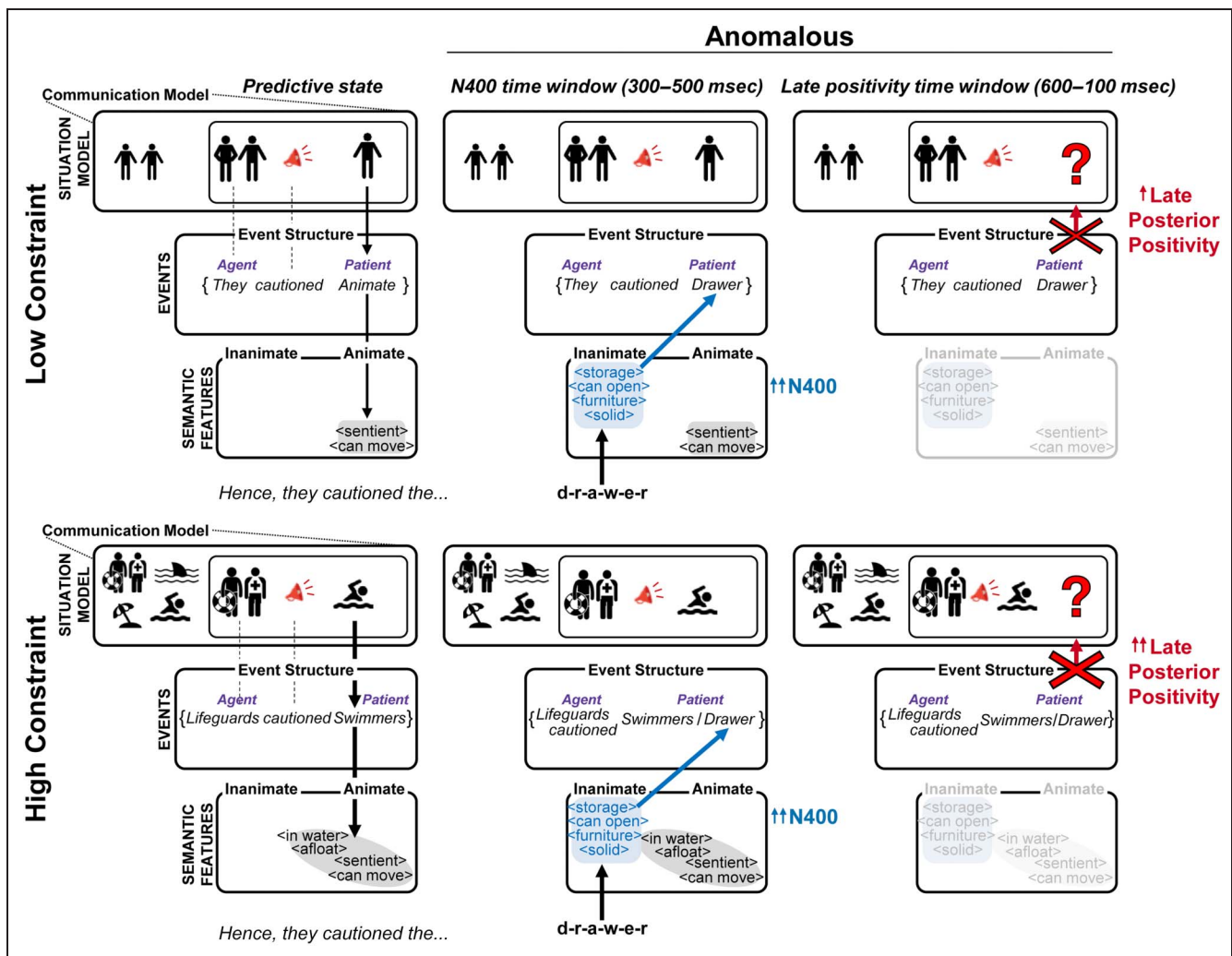
As shown in the top left panels of Figures 7 and 8, we assume that, after reading a low constraint context but before encountering the critical word, the comprehender has established a simple situation model that constrains for the event structure {Agent cautioned animate-Patient} (see also Altmann & Mirkovic, 2009). This leads to the prediction of semantic features that are characteristic of animate entities (e.g., ⟨sentient⟩, ⟨can move⟩). As shown in the bottom left panels of Figures 7 and 8, we assume that, after reading a high constraint context, the comprehender has built a rich situation model, leading her to strongly predict the specific event, {Lifeguards cautioned Swimmers} (see also Altmann & Mirkovic, 2009; Hare, Jones, Thomson, Kelly, & McRae, 2009; McRae & Matsuki, 2009). This, in turn, leads her to predict not only semantic features associated with animate entities (e.g., ⟨sentient⟩, ⟨can move⟩) but also additional

semantic features associated with “swimmers” (e.g., ⟨in water⟩, ⟨afloat⟩).<sup>8</sup>

Within this framework, the amplitude of the N400 evoked by the incoming critical word can be naturally understood as reflecting the retrieval of, or access to, semantic features that have not already been predicted by the context (see Kuperberg, 2016, for further discussion). In the high constraint expected scenarios, the critical word “swimmers” offers no new semantic information, and so it evokes a small amplitude N400. In the middle panel of Figure 7, in both the low constraint unexpected and the high constraint unexpected scenarios, the semantic features associated with the critical word “trainees” match these pre-activated semantic features that characterize its animacy (⟨sentient⟩, ⟨can move⟩), but the comprehender must retrieve additional semantic features that more specifically characterize



**Figure 7.** Schematic illustration of the state of the language comprehension system at three hierarchical levels of representation before and after encountering plausible critical nouns. States associated with the low constraint unexpected scenarios (Condition 2 in Table 1) are illustrated in the top half of the figure, and states associated with the high constraint unexpected scenarios (Condition 3 in Table 1) are illustrated in the bottom half of the figure. Left: the predictive state before encountering the critical word. Middle: the state during the N400 time window (300–500 msec). Right: the state during the late positivity time window (600–1000 msec). Please see the section titled A Hierarchical Generative Framework of Language Comprehension for full explanation.



**Figure 8.** Schematic illustration of the state of the language comprehension system at three hierarchical levels of representation before and after encountering anomalous critical nouns. States associated with the low constraint anomalous scenarios (Condition 4 in Table 1) are illustrated in the top half of the figure, and states associated with the high constraint anomalous scenarios (Condition 5 in Table 1) are illustrated in the bottom half of the figure. Left: the predictive state before encountering the critical word. Middle: the state during the N400 time window (300–500 msec). Right: the state during the late positivity time window (600–1000 msec). Please see the section titled A Hierarchical Generative Framework of Language Comprehension for full explanation.

“trainees” (e.g., ⟨learning⟩, ⟨novice⟩). The amplitude of the N400, which reflects the retrieval of these unpredicted semantic features, is therefore larger (see Figure 7, middle).

Finally, in both the low constraint anomalous and the high constraint anomalous scenarios, no semantic features associated with the critical word “drawer” have been pre-activated, and so when the bottom–up input is encountered, the comprehender must retrieve all its properties, including its mismatching inanimate features (e.g., ⟨storage⟩, ⟨can open⟩, ⟨furniture⟩, ⟨solid⟩; see also Paczynski & Kuperberg, 2011, 2012). The amplitude of the N400 is therefore even larger.

As shown in both Figures 7 and 8 (middle), in all four types of unpredictable scenarios, any unpredicted semantic features are passed up to higher levels of representation where, at a slightly later stage of processing, activity

differs depending on condition. In the high constraint unexpected scenarios, the arrival of unpredicted input at the situation model produces a late frontal positivity effect. We suggest that this effect primarily reflects a large change in activity associated with successfully updating the situation model from its prior state to a new state. We further suggest in these high constraint unexpected scenarios, that this high-level shift is linked to the top–down suppression of competing lower-level incorrect lexico-semantic predictions and the selection of the correct lexico-semantic features. Specifically, we propose that, before encountering the critical word, the comprehender had built a rich situation model and had used this model to strongly predict a specific event, {Lifeguards cautioned Swimmers}, as well as the semantic features associated with ⟨swimmers⟩ (see bottom left panel of Figure 7). Then, as the unpredicted semantic features associated with

the critical word “trainees” (*learning*, *novice*) arrive at the event level, activity begins to shift toward a new event, *{Lifeguards cautioned Trainees}* (see bottom middle panel of Figure 7). This unpredicted event information is, in turn, passed up so that the comprehender begins to update her situation model. However, to actually *complete* updating the situation model, the comprehender must suppress the original incorrectly predicted event *{Lifeguards cautioned Swimmers}* and the incorrectly predicted semantic features associated with “swimmers,” thereby “selecting” the newly inferred event *{Lifeguards cautioned Trainees}* and the semantic features associated with “trainees” (see bottom right panel of Figure 7).

In the low constraint unexpected scenarios, no significant late frontal positivity effect is produced by the critical words. We suggest that this is because any shift in activity induced by the unpredicted semantic features at the level of the situation model is much smaller. As shown in the top middle panel of Figure 7, as soon as the unpredicted semantic features associated with the critical word “trainees” reach the event level, the comprehender is able to shift activity to infer the event *{They cautioned Trainees}*, within the N400 time window. This is because the event structure has already been correctly predicted, and there is no competition at the level of semantic features. In other words, by fully accessing these unpredicted semantic features, the amount of additional work to infer the specific event is minimal. It is possible that there may be some small additional shift as the comprehender updates her situation model in the late positivity time window, but this change in activity is much smaller than in the high constraint unexpected scenarios.<sup>9</sup>

Finally, we suggest that the late posterior positivity/P600 is triggered when new input conflicts with the constraints of the broader communication model, leading to an initial failure to update the situation model. This may be followed by second-pass attempts to make sense of the discourse scenarios through reanalysis, attempts at repair of the input, and/or revision of the communication model. As shown in the left panels of Figure 8, in both the high constraint and low constraint contexts, the situation model constrains for a particular event structure, *{Agent cautioned animate–Patient}*—the set of events that are compatible with the broader communication model. As shown in the middle panels of Figure 8, when unpredicted semantic features associated with “drawer” arrive at the event level, this leads the comprehender to infer a different event structure, *{Agent cautioned inanimate–Patient}*. This unexplained information is passed up to the situation model. However, unlike in the high constraint unexpected scenarios, the input cannot be explained at this level (the conflict cannot initially be resolved), triggering the late posterior positivity/P600. Note that there is greater conflict in the high constraint anomalous scenarios than in the low constraint anomalous scenarios because the original predictions were

stronger following the high constraint than the low constraint contexts.

### **Relationship with Other Functional Interpretations of These ERP Components and Other Models of Language Processing**

#### *The N400*

A classic debate about the N400 is whether it reflects “lexical access” (e.g., Lau, Phillips, & Poeppel, 2008) or combinatorial “integration”<sup>10</sup> (e.g., Hagoort, Baggio, & Willems, 2009). Implicit in this debate is the assumption that “lexical access” and “integration” are often independent of one another and serial in nature, such that the comprehender first needs to achieve full access to a unique lexical item before she can begin to integrate this item into its preceding context.

The generative framework outlined above makes no such assumption. According to this framework, the comprehender can predict an upcoming event structure or a specific event (and its associated semantic features) before a new word becomes available from the bottom–up input (see also Rabovsky, Hansen, & McClelland, 2018; Metusalem et al., 2012). When the new word is encountered, so long as its semantic and syntactic features are consistent with the predicted event structure/event, little additional work is required to “integrate” it. In other words, because semantic predictions are derived from higher level event-level predictions, in a sense, some or all of a predicted word’s integration has already been completed, even before its occurrence (for further discussion of such “preupdating,” see Kuperberg & Jaeger, 2016, Section 4, pp. 14–15). In these cases, accessing the semantic features of the incoming word, as indexed by the N400, effectively completes the integration of that word into its context. In other cases, when a word violates strong prior event/lexical predictions (as in the high constraint unexpected scenarios), the additional work required to integrate it into the current situation model should primarily be reflected in the amplitude of later, post-N400 components, as discussed further below.

One argument that some have raised in favor of an integration account of the N400 is that its amplitude can be sensitive to the implausibility of the resulting interpretation (although the N400 is not always sensitive to plausibility, see Kuperberg, 2016; see also Shetreet et al., 2019, Supplementary Materials Section 3, for a recent review and discussion). The assumption is that, having accessed a lexical item, it takes more combinatorial work to construct an implausible/incoherent proposition than to construct a plausible/coherent proposition (see Nieuwland et al., in press; Lau, Namyst, Fogel, & Delgado, 2016). The amplitude of the N400 is, in part, taken to reflect this additional work. Because effects of plausibility/coherence on the N400 are sometimes observed over and above the effects of lexical probability (as indexed by cloze probability), it has also sometimes been assumed that any effects of implausibility do not reflect the consequences of prior



prediction (e.g., Nieuwland et al., in press; van Berkum, Zwitserlood, Hagoort & Brown, 2003).

Once again, the generative framework described above makes no such assumptions. We suggest that any sensitivity of the N400 to implausibility, over and above the effects of cloze probability, can again reflect the retrieval of semantic features that have not already been predicted. This is because comprehenders can, in some cases, use their situation model to predict event structures and their associated semantic features, even when these predictions do not correspond to the pre-activation of specific individual lexical items. For example, we have argued that the slightly larger N400 produced by the unpredictable critical word, “drawer”, in the anomalous scenarios reflects the additional retrieval of inanimate features that were not pre-activated by the context (see Wang, Jensen, & Kuperberg, 2018; Szewczyk & Schriefers, 2013, for evidence that semantic features linked to animacy can be pre-activated before encountering new input). Of course, it is possible that shifts in activity at higher levels of representation may begin within the N400 time window (e.g., Nieuwland et al., in press; Lau et al., 2016). Our primary point is that any effects of plausibility on the N400 itself can ultimately stem from the same semantic predictive mechanism that gives rise to the effects of cloze probability on the N400.

We also emphasize that the pre-activation of semantic features by a higher level representation cannot simply be explained by passive “resonance” or “priming” between content words, as is assumed in some theories of discourse comprehension (e.g., Myers & O’Brien, 1998), as well as some views of the N400 (e.g., Brouwer et al., 2012). This is because, in this study, semantic relationships between the critical words and the “bag of words” in the preceding context were matched between the unexpected and the anomalous conditions using LSA (see also Kuperberg, Paczynski, & Ditman, 2011 and see also Kuperberg, 2016, for discussion). Rather, we attribute any preactivation of semantic features related to the animacy of the incoming word to the prediction of the event structure, {*Agent cautioned animate–Patient*}. Finally, we emphasize that the preactivation of semantic features by a higher level representation does not necessarily or always correspond to the prediction of specific stored lexical entities (see also Kutas & Federmeier, 2011; Laszlo & Federmeier, 2009; Federmeier & Kutas, 1999).

### *Late Frontal Positivity*

As noted in the Introduction, the late frontal positivity is classically produced by words that violate strong lexical constraints. As such, it has been linked to the inhibition of incorrectly predicted words (Ness & Meltzer-Asscher, 2018; Federmeier et al., 2007; Kutas, 1993) and/or the integration of the violating bottom–up input to reach a new higher level interpretation (Brothers et al., 2015; DeLong et al., 2014; Federmeier et al., 2010). It is sometimes

implicitly assumed that these processes are independent of one another. For example, one possible account is that, having computed a partial high-level representation of the context (e.g., {The little boy went out to fly his...}) and used this to predict a specific upcoming lexical item (“kite”), upon encountering the unexpected word, “plane”, the comprehender must first inhibit “kite” before she can begin to integrate “plane” into her current contextual representation. The late frontal positivity effect would be assumed to reflect either the inhibition and/or the integrative process.

Within the hierarchical generative framework described above, however, these processes are fundamentally interrelated and they proceed in parallel. Specifically, as discussed above, after reading the high constraint context and before encountering the incoming word, the comprehender is assumed to have already updated her situation model such that it includes a representation of the lifeguards warning a group of swimmers (“pre-updating”). Within this framework, the late frontal positivity corresponds to the large shift in activity associated with updating this situation model when new information is encountered (“re-integration”), as well as the top–down feedback suppression of incorrectly predicted semantic features, which is necessary to complete this reupdating process.

### *Late Posterior Positivity/P600*

Debates about the functional significance of the late posterior positivity/P600 have come from many different perspectives.

In terms of its primary trigger, several authors have proposed that, to produce the effect, the comprehender must detect *conflict* between alternative representations that are computed during language comprehension (e.g., van de Meerendonk et al., 2009; Kuperberg, 2007; Kim & Osterhout, 2005). Others have emphasized that, unlike the late frontal positivity, the late posterior positivity/P600 is produced by words that yield highly implausible interpretations (DeLong et al., 2014; Kuperberg, 2013; Van Petten & Luka, 2012) and, in particular, by words that are, at least initially, perceived as impossible/anomalous during online comprehension (for discussion, see Kuperberg, 2007, Section 3.4; e.g., Paczynski & Kuperberg, 2012; van de Meerendonk et al., 2010).

The interpretation that we have offered above attempts to bridge these accounts. We argue that the primary trigger of the late posterior positivity/P600 is the detection of conflict with the comprehender’s broader communication model. As a result, the comprehender cannot initially incorporate the new input into her current situation model, leading to the initial perception of incoherence/impossibility. This account makes two assumptions. The first is that, to produce a late posterior positivity/P600, the comprehender must have previously established a communication model with a situation model at the top of the hierarchy, whose constraints now conflict with the new bottom–up input. For example, in this study, we assume

that the comprehender had already established a model of a literal English speaker who communicates about events that are possible in the real world (cf. Degen et al., 2015; Frank & Goodman, 2012). Therefore, the event, *{Lifeguards cautioned Drawer}*, cannot be interpreted into the situation model because it conflicts with these constraints. If, however, the comprehender had previously established a communication model that allowed for the expression of fantasy world events, in which drawers and other inanimate objects could understand language, then the same event would have been interpreted as plausible, and no late posterior positivity/P600 would have been produced (e.g., see Nieuwland & Van Berkum, 2006).

The second assumption is that, to trigger the late posterior positivity/P600, the initial failure to incorporate the input into the situation model occurs during fast, online comprehension. Although the online detection of incoherence often patterns with offline judgments of impossibility (as in the current study), this may not always be the case. Late posterior positivity/P600 effects are sometimes seen in sentences that are not judged to be impossible offline, but in which some temporary conflict with the existing communication model is detected during online processing. For example, a late posterior positivity/P600 can be evoked by unexpected words in semantically reversible sentences, even though these sentences are generally judged offline to be implausible, but not impossible (e.g., “The restaurant owner forgot which waitress the customer had served...”: Chow, Smith, Lau, & Phillips, 2016; Kolk, Chwilla, van Herten, & Oor, 2003; see Kuperberg, 2016, for discussion). In addition, a late posterior positivity/P600 is sometimes evoked in non-literal constructions such as nominal metaphors (e.g., De Grauwe, Swain, Holcomb, Ditman, & Kuperberg, 2010, Experiment 2) or in association with certain types of metonymic shifts (e.g., Schumacher, 2013), which may initially yield literal impossible interpretations during online comprehension.

In terms of the neurocognitive mechanisms reflected by the late posterior positivity/P600, one possibility is it only reflects the detection of high-level conflict and the resulting initial failure to incorporate the input into the current situation model. Another possibility is that it additionally reflects subsequent prolonged attempts to make sense of the input. These might include a reanalysis of the prior context to check whether or not it was accurately perceived the first time around (van de Meerendonk et al., 2010) and, if necessary, attempts to repair it. Alternatively, it may reflect second-pass attempts to come up with a revised interpretation of the input (Kuperberg, 2007; Kuperberg, Caplan, Sitnikova, Eddy, & Holcomb, 2006; see also Brouwer et al., 2012). Within the generative framework described above, this type of re-interpretation would entail revising the constraints of the communication model itself (modifying its parameters) so that it can accommodate the previously impossible interpretation. For example, revising the communication model may allow the comprehender to accept a fantasy world scenario as

possible, or come to a non-literal interpretation of the input (e.g., Schumacher, 2013; De Grauwe et al., 2010). As we discuss further below, this type of revision of the communication model is closely linked to longer-term learning/adaptation.

From the current study alone, we cannot determine whether the late posterior positivity/P600 effect reflects any of these additional processes. While our comprehension and judgment tasks suggest that participants were reading carefully for comprehension and successfully detecting the presence of semantic anomalies, they do not provide direct evidence that they were additionally engaged in processes related to reanalysis, repair, or re-interpretation. It will therefore be important for future experiments to further investigate the neurocognitive function of the late posterior positivity/P600 by using tasks that explicitly ask participants to repair the input or attempt to create representations of literally impossible events.

Finally, we note that the late posterior positivity/P600 is likely to share basic computational information processing mechanisms with the well-known posterior P300 component (Sassenhagen & Fiebach, 2019; Sassenhagen, Schlesewsky, & Bornkessel-Schlesewsky, 2014; Osterhout, Kim, & Kuperberg, 2012; Coulson, King, & Kutas, 1998). We discuss these relationships in detail in Kuperberg and Brothers (in preparation).

### **Relationship with More General Theories of Prediction, Prediction Error, and Predictive Coding in the Brain**

The generative framework of language comprehension outlined above is keeping with more general theories that emphasize a central role of prediction in the brain (e.g., Clark, 2013), as well as with many previous studies showing that the brain can encode “prediction error”—the difference between a predicted state of activity and a new state after new input is observed (den Ouden, Kok, & de Lange, 2012; Schultz & Dickinson, 2000).

Both probabilistic prediction and prediction error are important components of “predictive coding” (Friston, 2005; Rao & Ballard, 1997, 1999; Mumford, 1992), which has been proposed as an algorithm for carrying out Bayesian inference in the brain (although it is important to note that not all neural prediction error is in the service of Bayesian inference, and predictive coding may not be the only way in which Bayesian inference is carried out in the brain, see Aitchison & Lengyel, 2017). Hierarchical predictive coding involves passing up prediction error from lower to higher levels of the cortical hierarchy, updating beliefs at higher levels, and passing down new predictions to lower levels, with the overarching goal of minimizing prediction error across the entire cortical network (Friston, 2005). Although we think that it is premature to conclude that language comprehension in the brain is carried out through predictive coding, we do note that there are features of our model that are consistent with some of its principles.

First, we have suggested that unpredicted semantic features, reflected by the N400, are passed up to higher levels of representation. This is consistent with the basic principle that prediction error at lower levels is passed up a generative hierarchy and that it can be indexed by ERPs produced by superficial pyramidal cells (Friston, 2005). Importantly, here we assume that the “prediction error,” indexed by the N400, is at the level of semantic features, and that it can be formalized simply as a word’s semantic probability, independent of the probability of any other words that may have also been predicted by the prior context. In other words, when used in relation to the N400, the term “prediction error” does not necessarily correspond to how the term “prediction violation” (or “error”) is typically used in psycholinguistic models.

Second, we have argued that the late frontal positivity reflects a large high-level shift at the level of the situation model, which is induced by the arrival of unpredicted information that cannot be explained at lower levels of the generative hierarchy. This is consistent with the idea that, by passing up lower level prediction error to higher levels of representation, higher level prediction error can be reduced through Bayesian inference. The high-level shift of the situation model reflected by the *late frontal positivity* may correspond to “Bayesian surprise” (Baldi & Itti, 2010)—the difference between the prior and the posterior probability distribution of higher level hypotheses. Moreover, we have also suggested that this high-level shift is accompanied by top–down feedback suppression/selection to lower levels of representation. This is consistent with some versions of predictive coding, which propose that, having corrected higher level hypotheses, earlier incorrect lower level predictions can also be retrospectively corrected through feedback suppression. This, in turn, can lead to a relative enhancement of stimulus-driven activity that is consistent with the new higher level hypotheses (see Spratling, 2008; Lee & Mumford, 2003).

Finally, we have argued that the late posterior positivity/P600 is triggered by a failure to incorporate unpredicted input into the existing situation model because of conflict with the communication model. This can be conceptualized as reflecting unresolved prediction error at the highest level of the generative hierarchy. The prediction error is unresolved because there is a discrepancy between the input and the model that the comprehender had previously been assumed (cf. “puzzlement” surprise; see Faraji, Preuschoff, & Gerstner, 2018). As discussed briefly below, this “unexpected surprise” (Yu & Dayan, 2005) might cue the comprehender to revise the communication model or to switch models (adaptation), again with the broad goal of reducing overall prediction error, explaining the bottom–up input.

## Implications

By associating different ERP components with neural activity at different levels of a hierarchy of representations,

this hierarchical generative framework makes several predictions and has several implications.

First, to the degree that different hierarchical levels of representation map on to neuroanatomically distinct cortical regions, this framework predicts not only a temporal segregation of responses to words that confirm and violate predictions, but also some spatial segregation of neural activity. Consistent with this hypothesis, we have recently carried out a multimodal neuroimaging study, using ERP together with MEG and fMRI, to show that, within the N400 time window, confirmed predictions are associated with reduced activity within the left temporal cortex, whereas prediction violations are associated with additional recruitment of the left inferior prefrontal cortex, with feedback activity to different parts of the temporal cortex depending on the grain of representation that was violated.

Second, this generative hierarchical framework predicts that, under certain circumstances, it should be possible to violate lower level constraints without seeing evidence of a shift at the level of the situation model. For example, if comprehenders fail to engage in deep comprehension and have not established a higher-level situation model, then lexical prediction violations may be associated with a large N400 without triggering a late frontal positivity, and animacy violations may be associated with a large N400 without triggering a late posterior positivity/P600 (see Brothers, Wlotko, Warnke, & Kuperberg, submitted).

Conversely, it should also be possible to see evidence of higher level activity without violating lower level predictions. Some evidence for this claim comes from studies reporting a late frontal positivity in response to words that do not violate strong lexical constraints of their preceding contexts (e.g., Chow, Lau, Wang, & Phillips, 2018; Zirnstein, van Hell, & Kroll, 2018; Freunberger & Roehm, 2016; Thornhill & Van Petten, 2012). These “low constraint” frontal positivities may be produced if a new input to the situation model is particularly informative, triggering a large update, even when the prior situation model had not led to strong lower level predictions of upcoming semantic features. Similarly, a late posterior positivity/P600 can sometimes be evoked by words that do not violate strong lower level animacy or other selection restrictions if the bottom–up input conflicts with presuppositions that have previously been encoded and maintained within the comprehender’s high-level situation model (e.g., Shetreet et al., 2019).

We also emphasize that the precise conditions evoking the late positivities are likely to vary across languages. For example, in some languages (e.g., Spanish, Dutch), a mismatching gender-marked adjective or determiner may provide strong evidence that a predicted specific event will be violated, leading the comprehender to update her situation model and produce a late frontal positivity, even before she receives direct evidence of violating semantic features (e.g., Wicha, Moreno, & Kutas, 2004). Moreover, conditions evoking the late positivities are also likely to vary between individuals, where variability in

linguistic and domain-general cognitive abilities may predict differences in the amplitude of these components (e.g., see Zirmstein et al., 2018, for evidence of individual variability in the late frontal positivity; see Kim, Oines, & Miyake, 2018; Nakano, Saron, & Swaab, 2010, for evidence of individual variability in the late posterior positivity/P600).

Finally, within this generative framework, we can begin to understand *why* the brain should engage different neurocognitive mechanisms in response to predictions that are confirmed versus predictions that are strongly violated. This is because the framework explicitly links the process of language comprehension to language adaptation (see also Kleinschmidt & Jaeger, 2015; Dell & Chang, 2014; Chang, Dell, & Bock, 2006). In the present study, we have focused on comprehension, discussing the late positivities as responses that can potentially help comprehenders to recover meaning when input violates strong predictions by updating the current situation model (the late frontal positivity) or by diagnosing that something is wrong with the input and reanalyzing/repairing/reinterpreting it (the late posterior positivity/P600). It is sometimes assumed that these types of neurocognitive mechanisms reflect unwanted costs of a predictive language comprehension system (e.g., see Kutas, DeLong, & Smith, 2011; Federmeier, 2007, for discussion). However, another way of understanding their broader functional role is as providing strong signals that the statistical structure of the broader communicative environment has changed—the so-called “unexpected surprise” (Yu & Dayan, 2005). This type of signal may cue the comprehender to adapt to her new communicative environment, which is necessary for her to continue predicting efficiently. Within the hierarchical generative framework discussed here, adaptation would involve modifying the parameters of the comprehender’s existing generative network. This may entail revising the existing communication model so that, for example, it can accommodate future fantasy world scenarios or metaphorical interpretations (as discussed above). It may also entail revising assumptions about the reliability of the communicator, leading to a systematic down-weighting of linguistic cues at lower levels of the generative network. Alternatively, it might involve switching to a new generative model (cf. Qian, Jaeger, & Aslin, 2012, 2016; Gallistel, Krishan, Liu, Miller, & Latham, 2014; see Kleinschmidt & Jaeger, 2015, for a detailed discussion). Some existing evidence suggests that the late posterior positivity/P600 evoked by syntactic anomalies adapts over time (Hanulíková, van Alphen, van Goch, & Weber, 2012; Coulson et al., 1998; see also Kuperberg & Brothers, in preparation, for discussion). It will be important for future studies to determine whether the late positivities play a functional role in actually driving such adaptation processes.

## Notes

1. ERP effects of contextual predictability can sometimes appear to begin before 300 msec. It has been hypothesized that

this earlier divergence reflects an ERP component that is distinct from the N400 and that is more sensitive to the effects of predicting phonological/orthographic properties of an incoming word (e.g., Brothers, Swaab, & Traxler, 2015; Lau, Holcomb, & Kuperberg, 2013).

2. In the present study, we focus on the semantic P600. Although some of the ideas we propose are relevant to understanding the syntactic P600 (as well as posteriorly distributed late positivities that are evoked by other types of linguistic and non-linguistic violations), a full discussion of these relationships is outside the scope of this article.

3. In this study, the event structure was defined largely by the thematic properties of the verb in the third sentence — properties that define the semantic roles around specific types of actions or states (Dowty, 1989; Jackendoff, 1987; Fillmore, 1967; Gruber, 1965). However, other types of linguistic cues can also constrain strongly for particular sets of events (event structures), for example, presupposition triggers (see Shetreet, Alexander, Romoli, Chierchia, & Kuperberg, 2019) and concessive discourse connectives (see Xiang & Kuperberg, 2015).

4. The supplementary material for this paper can be retrieved from [https://projects.iq.harvard.edu/files/kuperberglab/files/kuperbergbrotherswlotko\\_jcn\\_2019\\_supp\\_materials.pdf](https://projects.iq.harvard.edu/files/kuperberglab/files/kuperbergbrotherswlotko_jcn_2019_supp_materials.pdf).

5. We initially carried out statistical analyses on trial-averaged data within each of these spatiotemporal ROIs. However, as a reviewer pointed out, this limits our ability to generalize our findings to new sets of items (see Clark, 1973). The pattern of results for these two sets of analyses did not differ. We also carried out an initial omnibus ANOVA analyses on the trial-averaged data, with Scenario Type (five levels corresponding to the five conditions) and Region (five levels corresponding to five regions across the anterior–posterior distribution of the scalp) as within-subject variables. These analyses confirmed that both the N400 and the late positivity effects differed significantly across the five conditions in their scalp distribution (significant Scenario Type  $\times$  Region interactions in both the 300–500 and 600–1000 msec time windows,  $F_s > 9.66$ ,  $p_s < .0001$ ).

6. Similar to previous findings (Paczynski & Kuperberg, 2011), the degree of N400 modulation depended on whether the verb constrained for an animate or inanimate direct object and on the animacy of the argument. A detailed discussion of these interactions is beyond the scope of the current article.

7. Generative models are classically described within Bayesian probabilistic frameworks of cognition (Griffiths, Kemp, & Tenenbaum, 2008) at Marr’s first computational level of analysis. However, they can also be instantiated at Marr’s second algorithmic level (see McClelland, 1998, 2013).

8. A generative framework also assumes that, at least under some circumstances, the comprehender can predict/pre-activate information at other levels of representation. In this study, this would include other lexical properties associated with the upcoming critical word. Specifically, in both the *high constraint* and *low constraint* contexts, the comprehender is likely to have predicted its syntactic representation (a noun phrase). And, in the *high constraint* contexts, she may have also predicted its phonological/orthographic representation (“swimmers”).

9. Although the late frontal positivity effect evoked by the low constraint unexpected (vs. the high constraint expected) critical words was not significant, examination of the voltage map for this contrast (shown in Figure 6, left) suggests that the spatial distribution of this effect was qualitatively similar to the distribution of the effect observed when contrasting the high constraint unexpected and the high constraint expected critical words within this time window.

10. Note that the term integration in relation to the N400 has not always been used to imply combinatorial processing. It has sometimes been used simply to refer to the use of context to facilitate semantic processing of incoming words, only



once some bottom-up information becomes available, without any assumption that this entails combinatorial processes that build new propositional meaning (e.g., Van Petten & Luka, 2012; Chwilla, Hagoort & Brown, 1998).

## Acknowledgments

This work was funded by the National Institute of Child Health and Human Development (R01 HD08252 to G. R. K.). E. W. W. was supported by an Institutional Research and Academic Career Development Award from the National Institute of General Medical Sciences (K12GM074869 to Tufts University, PI C. Moore). We thank Maria Luiza Cunha Lima, Margarita Zeitlin, and Connie Choi for their contributions to constructing the experimental materials, Margarita Zeitlin and Simone Riley for their assistance with data collection, Sophie Greene for her help in statistical analysis, and Lotte Schoot and Lin Wang for their insightful comments on the article. We are very grateful to Arim Choi Perrachione for her tremendous patience and help in making (and remaking) the figures.

Reprint requests should be sent to Gina R. Kuperberg, Department of Psychology, Tufts University, 490 Boston Avenue, Medford, MA 02155, or via e-mail: GKuperberg@mgh.harvard.edu.

## REFERENCES

- Aitchison, L., & Lengyel, M. (2017). With or without you: Predictive coding and Bayesian inference in the brain. *Current Opinion in Neurobiology*, *46*, 219–227.
- Altmann, G. T., & Mirkovic, J. (2009). Incrementality and prediction in human sentence processing. *Cognitive Science*, *33*, 583–609.
- Baldi, P., & Itti, L. (2010). Of bits and wows: A Bayesian theory of surprise with applications to attention. *Neural Networks*, *23*, 649–666.
- Balota, D. A., Yap, M. J., Hutchison, K. A., Cortese, M. J., Kessler, B., Loftis, B., et al. (2007). The English lexicon project. *Behavior Research Methods*, *39*, 445–459.
- Barr, D. J., Levy, R., Scheepers, C., & Tily, H. J. (2013). Random effects structure for confirmatory hypothesis testing: Keep it maximal. *Journal of Memory and Language*, *68*, 255–278.
- Bates, D. M., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting linear mixed-effects models using lme4. *Journal of Statistical Software*, *67*, 1–48.
- Brothers, T., Swaab, T. Y., & Traxler, M. J. (2015). Effects of prediction and contextual support on lexical processing: Prediction takes precedence. *Cognition*, *136*, 135–149.
- Brothers, T., Wlotko, E. W., Warnke, L., & Kuperberg, G. R. (submitted). Going the extra mile: Effects of discourse context on two late positivities during language comprehension.
- Brouwer, H., Fitz, H., & Hoeks, J. (2012). Getting real about semantic illusions: Rethinking the functional role of the P600 in language comprehension. *Brain Research*, *1446*, 127–143.
- Chang, F., Dell, G. S., & Bock, K. (2006). Becoming syntactic. *Psychological Review*, *113*, 234–272.
- Chow, W.-Y., Lau, E. F., Wang, S., & Phillips, C. (2018). Wait a second! Delayed impact of argument roles on on-line verb prediction. *Language, Cognition and Neuroscience*, *33*, 803–828.
- Chow, W.-Y., Smith, C., Lau, E. F., & Phillips, C. (2016). A “bag-of-arguments” mechanism for initial verb predictions. *Language, Cognition and Neuroscience*, *31*, 577–596.
- Chwilla, D. J., Hagoort, P., & Brown, C. M. (1998). The mechanism underlying backward priming in a lexical decision task: Spreading activation versus semantic matching. *The Quarterly Journal of Experimental Psychology: Section A: Human Experimental Psychology*, *51*, 531–560.
- Clark, A. (2013). Whatever next? Predictive brains, situated agents, and the future of cognitive science. *Behavioral and Brain Sciences*, *36*, 181–204.
- Clark, H. H. (1973). The language-as-a-fixed-effect fallacy: A critique of language statistics in psychological research. *Journal of Verbal Learning and Verbal Behavior*, *12*, 335–359.
- Coulson, S., King, J. W., & Kutas, M. (1998). Expect the unexpected: Event-related brain responses to morphosyntactic violations. *Language and Cognitive Processes*, *13*, 21–58.
- De Grauwe, S., Swain, A., Holcomb, P. J., Ditman, T., & Kuperberg, G. R. (2010). Electrophysiological insights into the processing of nominal metaphors. *Neuropsychologia*, *48*, 1965–1984.
- Degen, J., Tessler, M. H., & Goodman, N. D. (2015). Wonky worlds: Listeners revise world knowledge when utterances are odd. Paper presented at the Proceedings of the Thirty-Seventh Annual Conference of the Cognitive Science Society.
- Dell, G. S., & Chang, F. (2014). The P-chain: Relating sentence production and its disorders to comprehension and acquisition. *Philosophical Transactions of the Royal Society of London: Series B: Biological Sciences*, *369*, 20120394.
- DeLong, K. A., Quante, L., & Kutas, M. (2014). Predictability, plausibility, and two late ERP positivities during written sentence comprehension. *Neuropsychologia*, *61*, 150–162.
- DeLong, K. A., Urbach, T. P., Groppe, D. M., & Kutas, M. (2011). Overlapping dual ERP responses to low cloze probability sentence continuations. *Psychophysiology*, *48*, 1203–1207.
- Delorme, A., & Makeig, S. (2004). EEGLAB: An open source toolbox for analysis of single-trial EEG dynamics including independent component analysis. *Journal of Neuroscience Methods*, *134*, 9–21.
- den Ouden, H. E., Kok, P., & de Lange, F. P. (2012). How prediction errors shape perception, attention, and motivation. *Frontiers in Psychology*, *3*, 548.
- Dowty, D. R. (1989). On the semantic content of the notion of thematic role. In G. Cherchia, B. Partee, & R. Turner (Eds.), *Properties, types and meaning* (pp. 69–129). Norwell, MA: Kluwer.
- Faraji, M., Preuschoff, K., & Gerstner, W. (2018). Balancing new against old information: The role of puzzlement surprise in learning. *Neural Computation*, *30*, 34–83.
- Federmeier, K. D. (2007). Thinking ahead: The role and roots of prediction in language comprehension. *Psychophysiology*, *44*, 491–505.
- Federmeier, K. D., & Kutas, M. (1999). A rose by any other name: Long-term memory structure and sentence processing. *Journal of Memory and Language*, *41*, 469–495.
- Federmeier, K. D., Kutas, M., & Schul, R. (2010). Age-related and individual differences in the use of prediction during language comprehension. *Brain and Language*, *115*, 149–161.
- Federmeier, K. D., Wlotko, E. W., De Ochoa-Dewald, E., & Kutas, M. (2007). Multiple effects of sentential constraint on word processing. *Brain Research*, *1146*, 75–84.
- Fillmore, C. J. (1967). The case for case. Paper presented at the Texas Symposium on Language Universals.
- Frank, M. C., & Goodman, N. D. (2012). Predicting pragmatic reasoning in language games. *Science*, *336*, 998.
- Freunberger, D., & Roehm, D. (2016). Semantic prediction in language comprehension: Evidence from brain potentials. *Language, Cognition and Neuroscience*, *31*, 1193–1205.
- Friston, K. (2005). A theory of cortical responses. *Philosophical Transactions of the Royal Society of London: Series B: Biological Sciences*, *360*, 815–836.

- Gallistel, C. R., Krishan, M., Liu, Y., Miller, R., & Latham, P. E. (2014). The perception of probability. *Psychological Review*, *121*, 96–123.
- Griffiths, T. L., Kemp, C., & Tenenbaum, J. B. (2008). Bayesian models of cognition. In R. Sun (Ed.), *The Cambridge handbook of computational psychology* (pp. 59–100). New York: Cambridge University Press.
- Gruber, J. S. (1965). Studies in lexical relations. (Doctoral dissertation), Massachusetts Institute of Technology.
- Gunter, T. C., Friederici, A. D., & Schriefers, H. (2000). Syntactic gender and semantic expectancy: ERPs reveal early autonomy and late interaction. *Journal of Cognitive Neuroscience*, *12*, 556–568.
- Hagoort, P., Baggio, G., & Willems, R. M. (2009). Semantic unification. In M. S. Gazzaniga (Ed.), *The cognitive neurosciences* (4th ed., pp. 819–836). Cambridge, MA: MIT Press.
- Hagoort, P., Brown, C., & Groothusen, J. (1993). The syntactic positive shift (SPS) as an ERP measure of syntactic processing. *Language and Cognitive Processes*, *8*, 439–483.
- Hanulíková, A., van Alphen, P. M., van Goch, M. M., & Weber, A. (2012). When one person's mistake is another's standard usage: The effect of foreign accent on syntactic processing. *Journal of Cognitive Neuroscience*, *24*, 878–887.
- Hare, M., Jones, M., Thomson, C., Kelly, S., & McRae, K. (2009). Activating event knowledge. *Cognition*, *111*, 151–167.
- Jackendoff, R. (1987). The status of thematic relations in linguistic theory. *Linguistic Inquiry*, *18*, 369–411.
- Kim, A., Oines, L., & Miyake, A. (2018). Individual differences in verbal working memory underlie a tradeoff between semantic and structural processing difficulty during language comprehension: An ERP investigation. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *44*, 406–420.
- Kim, A., & Osterhout, L. (2005). The independence of combinatory semantic processing: Evidence from event-related potentials. *Journal of Memory and Language*, *52*, 205–225.
- Kleinschmidt, D. F., & Jaeger, F. T. (2015). Robust speech perception: Recognize the familiar, generalize to the similar, and adapt to the novel. *Psychological Review*, *122*, 148–203.
- Kolk, H. H., Chwilla, D. J., van Herten, M., & Oor, P. J. (2003). Structure and limited capacity in verbal working memory: A study with event-related potentials. *Brain and Language*, *85*, 1–36.
- Kuperberg, G. R. (2007). Neural mechanisms of language comprehension: Challenges to syntax. *Brain Research*, *1146*, 23–49.
- Kuperberg, G. R. (2013). The proactive comprehender: What event-related potentials tell us about the dynamics of reading comprehension. In B. Miller, L. Cutting, & P. McCardle (Eds.), *Unraveling reading comprehension: Behavioral, neurobiological, and genetic components* (pp. 176–192). Baltimore: Brookes.
- Kuperberg, G. R. (2016). Separate streams or probabilistic inference? What the N400 can tell us about the comprehension of events. *Language, Cognition and Neuroscience*, *31*, 602–616.
- Kuperberg, G. R., & Brothers, T. (in preparation). What can the P300 tell us about the P600? Understanding language comprehension within a decision theoretic framework.
- Kuperberg, G. R., Caplan, D., Sitnikova, T., Eddy, M., & Holcomb, P. J. (2006). Neural correlates of processing syntactic, semantic, and thematic relationships in sentences. *Language and Cognitive Processes*, *21*, 489–530.
- Kuperberg, G. R., & Jaeger, T. F. (2016). What do we mean by prediction in language comprehension? *Language, Cognition and Neuroscience*, *31*, 32–59.
- Kuperberg, G. R., Kreher, D. A., Sitnikova, T., Caplan, D. N., & Holcomb, P. J. (2007). The role of animacy and thematic relationships in processing active English sentences: Evidence from event-related potentials. *Brain and Language*, *100*, 223–237.
- Kuperberg, G. R., Paczynski, M., & Ditman, T. (2011). Establishing causal coherence across sentences: An ERP study. *Journal of Cognitive Neuroscience*, *23*, 1230–1246.
- Kuperberg, G. R., Sitnikova, T., Caplan, D., & Holcomb, P. J. (2003). Electrophysiological distinctions in processing conceptual relationships within simple sentences. *Cognitive Brain Research*, *17*, 117–129.
- Kutas, M. (1993). In the company of other words: Electrophysiological evidence for single-word and sentence context effects. *Language and Cognitive Processes*, *8*, 533–572.
- Kutas, M., DeLong, K. A., & Smith, N. J. (2011). A look around at what lies ahead: Prediction and predictability in language processing. In M. Bar (Ed.), *Predictions in the brain: Using our past to generate a future* (pp. 190–207). New York: Oxford University Press.
- Kutas, M., & Federmeier, K. D. (2011). Thirty years and counting: Finding meaning in the N400 component of the event-related brain potential (ERP). *Annual Review of Psychology*, *62*, 621–647.
- Kutas, M., & Hillyard, S. A. (1984). Brain potentials during reading reflect word expectancy and semantic association. *Nature*, *307*, 161–163.
- Kuznetsova, A., Brockhoff, P. B., & Christensen, R. H. B. (2015). *Tests for random and fixed effects for linear mixed effect models* (lmer objects of lme4 package, R package Version 2.0–33). Retrieved from cran.r-project.org/package=lmerTest.
- Landauer, T. K., & Dumais, S. T. (1997). A solution to Plato's problem: The latent semantic analysis theory of acquisition, induction, and representation of knowledge. *Psychological Review*, *104*, 211–240.
- Landauer, T. K., Foltz, P. W., & Laham, D. (1998). An introduction to latent semantic analysis. *Discourse Processes*, *25*, 259–284.
- Laszlo, S., & Federmeier, K. D. (2009). A beautiful day in the neighborhood: An event-related potential study of lexical relationships and prediction in context. *Journal of Memory and Language*, *61*, 326–338.
- Lau, E. F., Holcomb, P. J., & Kuperberg, G. R. (2013). Dissociating N400 effects of prediction from association in single-word contexts. *Journal of Cognitive Neuroscience*, *25*, 484–502.
- Lau, E. F., Namyst, A., Fogel, A., & Delgado, T. (2016). A direct comparison of N400 effects of predictability and incongruity in adjective-noun combination. *Collabra*, *2*, 13.
- Lau, E. F., Phillips, C., & Poeppel, D. (2008). A cortical network for semantics: (De)constructing the N400. *Nature Reviews Neuroscience*, *9*, 920–933.
- Lee, T. S., & Mumford, D. (2003). Hierarchical Bayesian inference in the visual cortex. *Journal of the Optical Society of America A: Optics, Image, Science, and Vision*, *20*, 1434–1438.
- Levin, B. (1993). *English verb classes and alternations: A preliminary investigation*. Chicago: University of Chicago Press.
- Lopez-Calderon, J., & Luck, S. J. (2014). ERPLAB: An open-source toolbox for the analysis of event-related potentials. *Frontiers in Human Neuroscience*, *8*, 213.
- Lund, K., & Burgess, C. (1996). Producing high-dimensional semantic spaces from lexical co-occurrence. *Behavior Research Methods, Instruments, & Computers*, *28*, 203–208.
- McClelland, J. L. (1998). Connectionist models and Bayesian inference. In M. Oaksford & N. Chater (Eds.), *Rational*

- models of cognition* (pp. 21–52). New York: Oxford University Press.
- McClelland, J. L. (2013). Integrating probabilistic models of perception and interactive neural networks: A historical and tutorial review. *Frontiers in Psychology*, 4, 503.
- McRae, K., & Matsuki, K. (2009). People use their knowledge of common events to understand language, and do so as quickly as possible. *Language and Linguistics Compass*, 3, 1417–1429.
- Metusalem, R., Kutas, M., Urbach, T. P., Hare, M., McRae, K., & Elman, J. L. (2012). Generalized event knowledge activation during online sentence comprehension. *Journal of Memory and Language*, 66, 545–567.
- Morey, R. D. (2008). Confidence intervals from normalized data: A correction to Cousineau (2005). *Tutorial in Quantitative Methods for Psychology*, 4, 61–64.
- Mumford, D. (1992). On the computational architecture of the neocortex. II. The role of cortico-cortical loops. *Biological Cybernetics*, 66, 241–251.
- Myers, J. L., & O'Brien, E. J. (1998). Accessing the discourse representation during reading. *Discourse Processes*, 26, 131–157.
- Nakano, H., Saron, C., & Swaab, T. Y. (2010). Speech and span: Working memory capacity impacts the use of animacy but not of world knowledge during spoken sentence comprehension. *Journal of Cognitive Neuroscience*, 22, 2886–2898.
- Ness, T., & Meltzer-Asscher, A. (2018). Lexical inhibition due to failed prediction: Behavioral evidence and ERP correlates. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 44, 1269–1285.
- Nieuwland, M. S., Barr, D. J., Bartolozzi, F., Busch-Moreno, S., Darley, E., Donaldson, D. I., et al. (in press). Dissociable effects of prediction and integration during language comprehension: Evidence from a large-scale study using brain potentials. *Philosophical Transactions of the Royal Society of London: Series B: Biological Sciences*. <https://doi.org/10.1101/267815>.
- Nieuwland, M. S., & Van Berkum, J. J. A. (2006). When peanuts fall in love: N400 evidence for the power of discourse. *Journal of Cognitive Neuroscience*, 18, 1098–1111.
- Oldfield, R. C. (1971). The assessment and analysis of handedness: The Edinburgh inventory. *Neuropsychologia*, 9, 97–113.
- Osterhout, L., & Holcomb, P. J. (1992). Event-related brain potentials elicited by syntactic anomaly. *Journal of Memory and Language*, 31, 785–806.
- Osterhout, L., Kim, A., & Kuperberg, G. R. (2012). The neurobiology of sentence comprehension. In M. Spivey, M. Joannisse, & K. McRae (Eds.), *The Cambridge handbook of psycholinguistics* (pp. 365–389). Cambridge: Cambridge University Press.
- Osterhout, L., & Nicol, J. (1999). On the distinctiveness, independence and time course of the brain responses to syntactic and semantic anomalies. *Language and Cognitive Processes*, 14, 283–317.
- Paczynski, M., & Kuperberg, G. R. (2011). Electrophysiological evidence for use of the animacy hierarchy, but not thematic role assignment, during verb argument processing. *Language and Cognitive Processes*, 26, 1402–1456.
- Paczynski, M., & Kuperberg, G. R. (2012). Multiple influences of semantic memory on sentence processing: Distinct effects of semantic relatedness on violations of real-world event/state knowledge and animacy selection restrictions. *Journal of Memory and Language*, 67, 426–448.
- Peirce, J. W. (2007). PsychoPy—Psychophysics software in Python. *Journal of Neuroscience Methods*, 162, 8–13.
- Qian, T., Jaeger, T. F., & Aslin, R. N. (2012). Learning to represent a multi-context environment: More than detecting changes. *Frontiers in Psychology*, 3, 228.
- Qian, T., Jaeger, T. F., & Aslin, R. N. (2016). Incremental implicit learning of bundles of statistical patterns. *Cognition*, 157, 156–173.
- Quante, L., Bölte, J., & Zwitserlood, P. (2018). Dissociating predictability, plausibility and possibility of sentence continuations in reading: Evidence from late-positivity ERPs. *PeerJ*, 6, e5717.
- R Core Team. (2016). *R: A language and environment for statistical computing*. Vienna, Austria. <http://www.R-project.org/>.
- Rabovsky, M., Hansen, S. S., & McClelland, J. L. (2018). Modelling the N400 brain potential as change in a probabilistic representation of meaning. *Nature Human Behaviour*, 2, 693–705.
- Rao, R. P., & Ballard, D. H. (1997). Dynamic model of visual recognition predicts neural response properties in the visual cortex. *Neural Computation*, 9, 721–763.
- Rao, R. P., & Ballard, D. H. (1999). Predictive coding in the visual cortex: A functional interpretation of some extra-classical receptive-field effects. *Nature Neuroscience*, 2, 79–87.
- Sanford, A. J., Leuthold, H., Bohan, J., & Sanford, A. J. S. (2011). Anomalies at the borderline of awareness: An ERP study. *Journal of Cognitive Neuroscience*, 23, 514–523.
- Sassenhagen, J., & Fiebach, C. J. (2019). Finding the P3 in the P600: Decoding shared neural mechanisms of responses to syntactic violations and oddball targets. *Neuroimage*, 200, 425–436.
- Sassenhagen, J., Schlesewsky, M., & Bornkessel-Schlesewsky, I. (2014). The P600-as-P3 hypothesis revisited: Single-trial analyses reveal that the late EEG positivity following linguistically deviant material is reaction time aligned. *Brain and Language*, 137, 29–39.
- Schultz, W., & Dickinson, A. (2000). Neuronal coding of prediction errors. *Annual Review of Neuroscience*, 23, 473–500.
- Schumacher, P. B. (2013). When combinatorial processing results in reconceptualization: Toward a new approach of compositionality. *Frontiers in Psychology*, 4, 677.
- Schwanenflugel, P. J., & Lacoount, K. L. (1988). Semantic relatedness and the scope of facilitation for upcoming words in sentences. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 14, 344–354.
- Shetreet, E., Alexander, E. J., Romoli, J., Chierchia, G., & Kuperberg, G. R. (2019). What we know about knowing: Presuppositions generated by factive verbs influence downstream neural processing. *Cognition*, 184, 96–106.
- Spratling, M. W. (2008). Predictive coding as a model of biased competition in visual attention. *Vision Research*, 48, 1391–1408.
- Szewczyk, J. M., & Schriefers, H. (2013). Prediction in language comprehension beyond specific words: An ERP study on sentence comprehension in Polish. *Journal of Memory and Language*, 68, 297–314.
- Taylor, W. (1953). “Cloze” procedure: A new tool for measuring readability. *Journalism Quarterly*, 30, 415–433.
- Thornhill, D. E., & Van Petten, C. (2012). Lexical versus conceptual anticipation during sentence processing: Frontal positivity and N400 ERP components. *International Journal of Psychophysiology*, 83, 382–392.
- Van Berkum, J. J., Zwitserlood, P., Hagoort, P., & Brown, C. M. (2003). When and how do listeners relate a sentence to the wider discourse? Evidence from the N400 effect. *Cognitive Brain Research*, 17, 701–718.
- van de Meerendonk, N., Kolk, H. H. J., Chwilla, D. J., & Vissers, C. T. W. M. (2009). Monitoring in language perception. *Language and Linguistics Compass*, 3, 1211–1224.
- van de Meerendonk, N., Kolk, H. H. J., Vissers, C. T. W. M., & Chwilla, D. J. (2010). Monitoring language perception: Mild and strong conflicts elicit different ERP patterns. *Journal of Cognitive Neuroscience*, 22, 67–82.

- van Dijk, T. A., & Kintsch, W. (1983). *Strategies of discourse comprehension*. New York: Academic Press.
- Van Petten, C., & Luka, B. J. (2012). Prediction during language comprehension: benefits, costs, and ERP components. *International Journal of Psychophysiology*, *83*, 176–190.
- Wang, L., Jensen, O., & Kuperberg, G. R. (2018). Neural evidence for prediction of animacy features by verbs during language comprehension: Evidence from MEG and EEG representational similarity analysis. Paper presented at the 10th Annual Meeting of the Society for the Neurobiology of Language, Quebec City, Canada.
- Wicha, N. Y., Moreno, E. M., & Kutas, M. (2004). Anticipating words and their gender: An event-related brain potential study of semantic integration, gender expectancy, and gender agreement in Spanish sentence reading. *Journal of Cognitive Neuroscience*, *16*, 1272–1288.
- Xiang, M., & Kuperberg, G. (2015). Reversing expectations during discourse comprehension. *Language, Cognition and Neuroscience*, *30*, 648–672.
- Yu, A. J., & Dayan, P. (2005). Uncertainty, neuromodulation, and attention. *Neuron*, *46*, 681–692.
- Zirnstien, M., van Hell, J. G., & Kroll, J. F. (2018). Cognitive control ability mediates prediction costs in monolinguals and bilinguals. *Cognition*, *176*, 87–106.
- Zwaan, R. A., & Radvansky, G. A. (1998). Situation models in language comprehension and memory. *Psychological Bulletin*, *123*, 162–185.