

A hierarchy of linguistic predictions during natural language comprehension

Micha Heilbron^{1,2}, Kristijan Armeni¹, Jan-Mathijs Schoffelen¹,
Peter Hagoort^{1,2}, Floris P. de Lange¹

¹Donders Institute, Radboud University, Nijmegen, the Netherlands

²Max Planck Institute for Psycholinguistics, Nijmegen, the Netherlands

m.heilbron@donders.ru.nl

Abstract

Understanding spoken language requires transforming ambiguous acoustic streams into a hierarchy of representations, from phonemes to meaning. It has been suggested that the brain uses prediction to guide the interpretation of incoming input. However, the role of prediction in language processing remains disputed, with disagreement about both the ubiquity and representational nature of predictions. Here, we address both issues by analysing brain recordings of participants listening to audiobooks, and using a deep neural network (GPT-2) to precisely quantify contextual predictions. First, we establish that brain responses to words are modulated by ubiquitous, probabilistic predictions. Next, we disentangle model-based predictions into distinct dimensions, revealing dissociable signatures of syntactic, phonemic and semantic predictions. Finally, we show that high-level (word) predictions inform low-level (phoneme) predictions, supporting hierarchical predictive processing. Together, these results underscore the ubiquity of prediction in language processing, showing that the brain spontaneously predicts upcoming language at multiple levels of abstraction.

INTRODUCTION

Understanding spoken language requires transforming ambiguous stimulus streams into a hierarchy of increasingly abstract representations, ranging from speech sounds to meaning. It is often argued that during this process, the brain relies on prediction to guide the interpretation of incoming information [1, 2]. Such a 'predictive processing' strategy has not only proven effective for artificial systems processing language [3, 4], but has also been found to occur in neural systems in related domains such as perception and motor control and might constitute a canonical neural computation [5, 6].

There is a considerable amount of evidence that appears in line with predictive language processing. For instance, behavioural and brain responses are highly sensitive to violations of linguistic regularities [7, 8] and to deviations from linguistic expectations

more broadly [9-13]. While such effects are well-documented, two important questions about the role of prediction in language processing remain unresolved [14].

The first question concerns the *ubiquity* of prediction. While some models cast prediction as a routine, integral part of language processing [1, 15, 16], others view it as relatively rare, pointing out that apparent widespread prediction effects might instead reflect other processes like semantic integration difficulty [17, 18]; or that such prediction effects might be exaggerated by the use of artificial, prediction-encouraging experiments focussing on highly predictable 'target' words [17, 19]. The second question concerns the representational nature of predictions: Does linguistic prediction occur primarily at the level of syntax [15, 20-22] or rather at the lexical [16, 23], semantic [24, 25] or the phonological

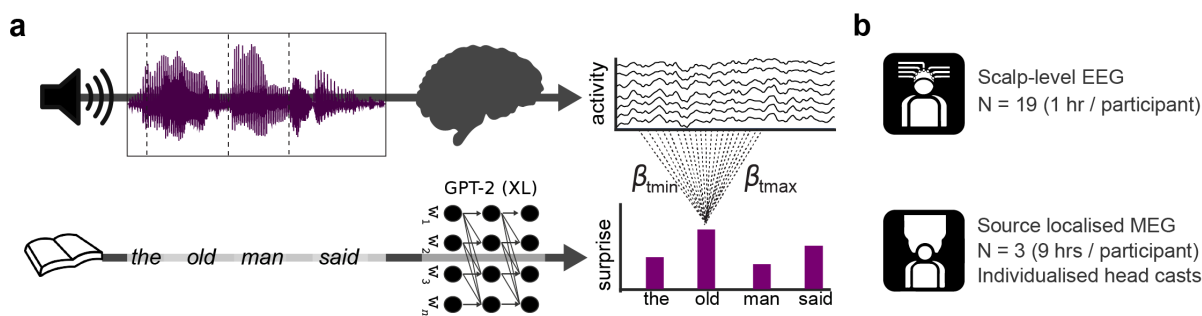


Figure 1: SCHEMATIC OF EXPERIMENTAL AND ANALYTICAL FRAMEWORK **a)** Top row: in both experiments participants listened to continuous recordings from audiobooks while brain activity was recorded. Bottom row: the texts participants listened to were analysed by a deep neural network (GPT-2) to quantify the contextual probability of each word. A regression-based technique was used to estimate the effects of (different levels of) linguistic unexpectedness on the evoked responses within the continuous recordings. **b)** Datasets analysed: one group-level EEG dataset, and one individual subject source-localised MEG dataset.

37 level [13, 26-29]? ERP studies have described brain
38 responses to violations of, and deviations from, both
39 high and low-level expectations, suggesting predic-
40 tion might occur at all levels simultaneously [1, 19],
41 although see [30]. However, it has been disputed
42 whether these findings would generalise to natural
43 language, where violations are rare or absent and
44 with few highly predictable words. In these cases,
45 prediction may be less relevant or might perhaps be
46 limited to the most abstract levels [17, 19, 30].

47 Here, we address both issues, probing the ubiq-
48 uity and nature of linguistic prediction during nat-
49 ural language understanding. Specifically, we anal-
50 ysed brain recordings from two independent exper-
51 iments of participants listening to audiobooks, and
52 use a state-of-the-art deep neural network (GPT-2)
53 to quantify linguistic predictions in a fine-grained,
54 contextual fashion. First, we obtain evidence for pre-
55 dictive processing, confirming that brain responses
56 to words are modulated by *probabilistic* predictions.
57 Critically, the effects of prediction were found over
58 and above those of non-predictive factors such as
59 integration difficulty, and were not confined to a
60 subset of predictable words, but were widespread –
61 supporting the notion of *ubiquitous* prediction. Next,
62 we investigated at which level prediction occurs. To
63 this end, we disentangled the model-based predic-
64 tions into distinct dimensions, revealing dissociable

65 neural signatures of syntactic, phonemic and seman-
66 tic predictions. Finally, we found that higher-level
67 (word) predictions constrain lower-level (phoneme)
68 predictions, supporting hierarchical prediction. To-
69 gether, these results underscore the ubiquity of pre-
70 diction in language processing, and demonstrate
71 that prediction is not confined to a single level of
72 abstraction but occurs throughout the language net-
73 work, forming a hierarchy of predictions across all
74 levels of analysis, from phonemes to meaning.

75 RESULTS

76 We consider data from two independent exper-
77 iments, in which brain activity was recorded while
78 participants listened to natural speech from audio-
79 books. The first experiment is part of a publicly
80 available dataset [31], and contains 1 hour of elec-
81 troencephalographic (EEG) recordings in 19 partici-
82 pants. The second experiment collected 9 hours of
83 magneto-encephalographic (MEG) data in three indi-
84 viduals, using individualised head casts that allowed
85 us to localise the neural activity with high precision.
86 While both experiments had a similar setup (see Fig-
87 ure 1), they yield complementary insights, both at
88 the group level and in three individuals.

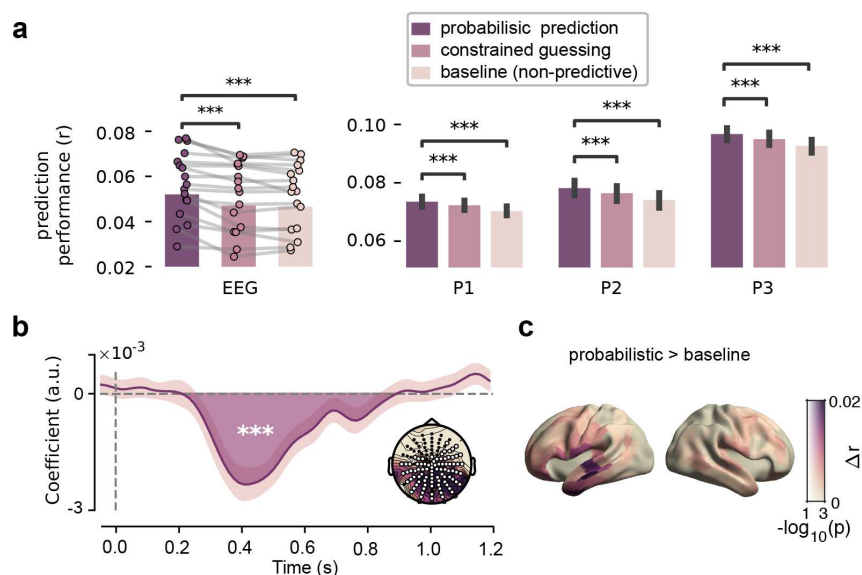


Figure 2: NEURAL RESPONSES ARE MODULATED BY PROBABILISTIC PREDICTIONS

a Model comparison. Cross-validated correlation coefficients for EEG (left) and each MEG participant (right). EEG: dots with connecting lines represent individual participants (averaged over all channels). MEG: bars represent median across runs, bars represent bootstrapped absolute deviance (averaged over language network sources). **b** EEG: coefficients describing the significant effect of lexical surprise (see Figure S3 for the full topography over time). Highlighted area indicates extent of the cluster, shaded error bar indicates bootstrapped SE. Inset shows distribution of absolute t-values and of channels in the cluster. **c** Difference in prediction performance across cortex (transparency indicates FWE-corrected p-values). Significance levels correspond to $P < 0.001$ (***) in a two-tailed one-sample Student's *t* or Wilcoxon sign rank test.

Neural responses to speech are modulated by probabilistic linguistic predictions

We first tested for evidence for linguistic prediction in general. We reasoned that if the brain is constantly predicting upcoming language, neural responses to words should be sensitive to violations of contextual predictions, yielding ‘prediction error’ signals which are considered a hallmark of predictive processing [5]. To this end, we used a regression-based deconvolution approach to estimate the effects of prediction error on evoked responses within the continuous recordings. We focus on this event-related, low-frequency evoked response because it connects most directly to earlier influential neural signatures of prediction in language [7, 30, 32, 33].

To quantify linguistic predictions, we analysed the books participants listened to with a state-of-the-art neural language model: GPT-2 [34]. GPT-2 is

a large transformer-based model that predicts the next word given the previous words, and is currently among the best publicly-available models of its kind. Note that we do not use GPT-2 as a model of human language processing, but purely as a tool to quantify how expected each word is in context.

To test whether neural responses to words are modulated by contextual predictions, we compared three regression models (see S5). The baseline model formalises the hypothesis that natural, passive language comprehension does not invoke prediction. This model did not include regressors related to contextual predictions, but did include several potentially confounding variables (such as word frequency, semantic integration, and acoustics). The *constrained guessing* model formalised the hypothesis that language processing *sometimes* (in constraining contexts) invokes prediction, and that such pre-

125 dictions are an all-or-none phenomenon – together
126 representing how the notion of prediction was classi-
127 cally used in the psycholinguistic literature [33]. This
128 model included all non-predictive variables from the
129 baseline model, plus, in constraining contexts, a lin-
130 ear estimate of word improbability (since all-or-none
131 predictions result in a linear relationship between
132 word probability and brain responses; see meth-
133 ods for details). Finally, the *probabilistic prediction*
134 model included all confounding regressors from the
135 baseline model, plus for every word a logarithmic
136 estimate of word improbability (i.e. *surprise*). This
137 formalises the hypothesis that the brain constantly
138 generates *probabilistic predictions*, as proposed by
139 predictive processing accounts of language [1, 32]
140 and of neural processing more broadly [5, 6].

141 When we compared the ability of these models
142 to predict brain activity using cross-validation, we
143 found that the probabilistic prediction model per-
144 formed better than both other models (see Figure
145 2a). The effect was highly consistent, found in virtu-
146 ally all EEG participants (probabilistic vs constrained
147 guessing, $t_{18} = 5.34$, $p = 4.46 \times 10^{-5}$; probabilistic
148 vs baseline, $t_{18} = 6.43$, $p = 4.70 \times 10^{-6}$) and within
149 each MEG participant (probabilistic vs constrained
150 guessing, all p 's $< 1.54 \times 10^{-6}$; probabilistic vs
151 baseline, all p 's $< 5.17 \times 10^{-12}$).

152 As the *constrained guessing* model differed from
153 the probabilistic model in two ways – by assuming
154 that predictions are (i) categorical and (ii) limited to
155 constraining contexts – we also considered a control
156 model. Like the constrained guessing model, this ex-
157 tended guessing model included a linear estimate of
158 word probability, but for every word rather than only
159 for constraining contexts. Although this model did
160 not outperform the probabilistic prediction model, it
161 did substantially outperform the constrained model
162 (Fig S5). This demonstrates that the effects of pre-
163 diction are not limited to constraining contexts, but
164 apply much more broadly – in line with the idea that
165 predictions are ubiquitous and automatic.

166 Having established that word unexpectedness
167 modulates neural responses, we characterised this
168 effect in space and time. In the MEG dataset, we

169 asked for which neural sources lexical surprise was
170 most important in explaining neural data, by com-
171 paring the prediction performance of the baseline
172 model to the predictive model in a spatially resolved
173 manner. This revealed that overall word unexpect-
174 edness modulated neural responses throughout
175 the language network (see Figure 2c). To investi-
176 gate the temporal dynamics of this effect, we in-
177 spected the regression coefficients, which describe
178 how fluctuations in lexical surprise modulate the
179 neural response at different time lags – together
180 forming a modulation function also known as the *re-*
181 *gression evoked response* [35] or Temporal Response
182 Function (TRF) [27, 36]. When we compared these
183 across participants in the EEG experiment, cluster-
184 based permutation tests revealed a significant effect
185 ($p = 2 \times 10^{-4}$) based on a postero-central cluster
186 with a negative polarity between 0.2 and 0.9 seconds
187 (see Figure 2b and S8). This indicates that surpris-
188 ing words lead to a stronger negative deflection of
189 evoked responses, an effect peaking at 400 ms post
190 word onset and strongly reminiscent of the classic
191 N400 [7, 24, 30]. Coefficients for MEG subjects re-
192 vealed a similar, slow effect at approximately the
193 same latencies (see Fig S4).

194 Together, these results constitute clear evidence
195 for predictive processing by confirming that brain
196 responses to words are modulated by predictions.
197 These modulations are not confined to constraining
198 contexts, occur throughout the language network,
199 evoke an effect reminiscent of the N400, and are
200 best explained by a probabilistic account of predic-
201 tion. This suggests the brain predicts constantly and
202 probabilistically – even when passively listening to
203 natural language.

204 Linguistic predictions are feature-specific

205 The results so far revealed modulations of neural
206 responses by *overall* word unexpectedness. What
207 type of linguistic prediction might be driving these
208 effects? Earlier research suggests a range of possibil-
209 ities, with some proposing that the effect of overall
210 word surprise primarily reflects syntax [15, 20], while
211 others propose that prediction unfolds at the seman-

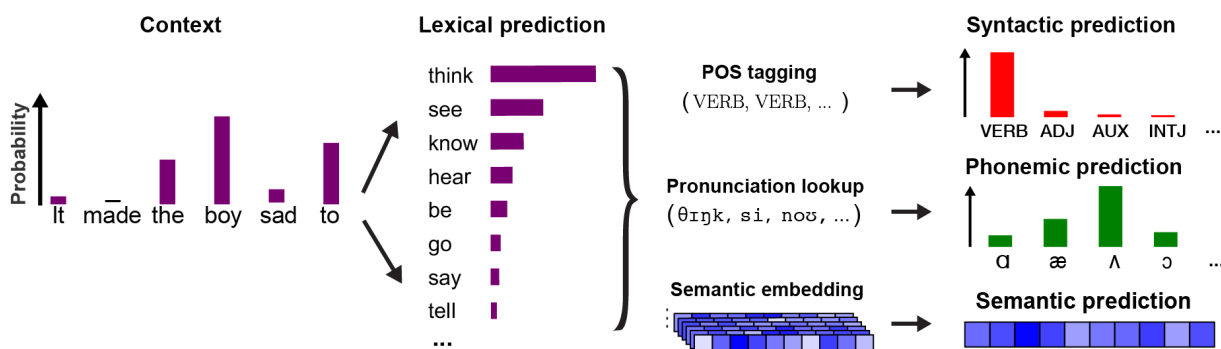


Figure 3: PARTITIONING MODEL-DERIVED PREDICTIONS INTO DISTINCT LINGUISTIC DIMENSIONS.

To disentangle syntactic, semantic and phonemic predictions, the lexical predictions from GPT-2 were analysed. For the syntactic prediction, part-of-speech tagging was performed over all potential sentences (e.g. "It made the boy sad to *think*"). To compute the phonemic prediction, each predicted word was decomposed into its constituent phonemes, and the predicted probabilities were used as a contextual prior in a phoneme model (see Figure 6). For the semantic prediction, a weighted average was computed over the GloVe embeddings of all predicted words.

212 tic [24, 25], or the phonemic level [13, 26, 27] – or at
 213 all levels simultaneously [1].

214 To evaluate these possibilities, we factorised the
 215 aggregate, word-level linguistic predictions from the
 216 artificial neural network into distinct linguistic dimen-
 217 sions (Fig 3). This allows us to derive model-based
 218 estimates of three feature-specific predictions: the
 219 syntactic prediction (defined as the conditional prob-
 220 ability distribution over parts-of-speech, given con-
 221 text), semantic prediction (defined as the predicted
 222 semantic embedding) and phonemic prediction (i.e.
 223 the conditional probability of the next phoneme,
 224 given the phonemes within the word so far and the
 225 prior context). By comparing these predictions to
 226 the presented words, we derived *feature-specific pre-*
 227 *diction errors* which quantified not just the extent to
 228 which a word is surprising overall, but also in what
 229 way: semantically, syntactically or phonemically (see
 230 Methods for definitions).

231 We reasoned that if the brain is generating predic-
 232 tions at a given level (e.g. syntax), then the neural
 233 responses should be sensitive to prediction errors
 234 specific to this level. Moreover, because these differ-
 235 ent features are processed by partly different brain
 236 areas over different timescales, the prediction errors
 237 should be at least partially dissociable. To test this,

238 we formulated a new regression model (Figure S6).
 239 This included all variables from the lexical prediction
 240 model as nuisance regressors, and added three re-
 241 gressors of interest: syntactic surprise (defined for
 242 each word), semantic prediction error (defined for
 243 each content word), and phonemic surprise (defined
 244 for each word-non-initial phoneme).

245 Because these regressors were to some degree
 246 correlated, we first asked whether, and in which
 247 brain area, each of the feature-specific prediction er-
 248 rors explained any unique variance, not explained by
 249 the other regressors. In this analysis, we turn to the
 250 MEG data because of its spatial specificity. As a con-
 251 trol, we first performed the analysis for a predictor
 252 with a known source: the acoustics. This revealed a
 253 clear peak around auditory cortex (Fig S7) especially
 254 in the right hemisphere. This aligns with prior work
 255 [37] and confirms that this approach can localise
 256 which areas are especially sensitive to a given re-
 257 gressor. We then tested the three prediction errors,
 258 finding that each type of prediction error explained
 259 significant unique variance in each individual (Figure
 260 4), except in participant 1 where phonemic surprise
 261 did not survive multiple comparisons correction (but
 262 see Figure 6c and Discussion). This shows that the
 263 brain responds differently to different types of pre-

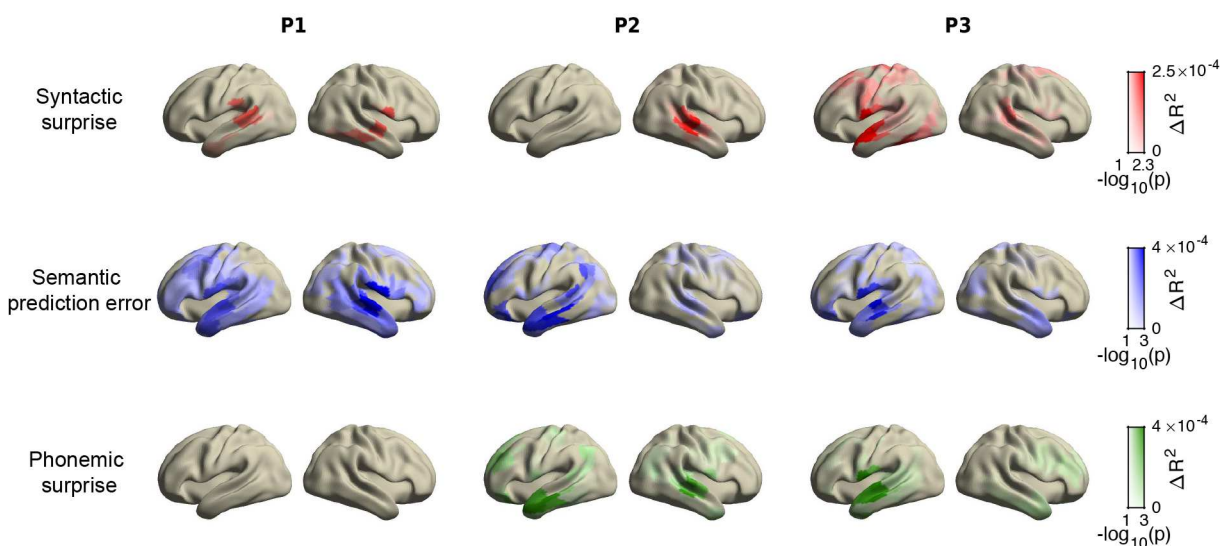


Figure 4: DISSOCIABLE PATTERNS OF EXPLAINED VARIANCE BY SYNTACTIC, SEMANTIC AND PHONEMIC PREDICTIONS.

Unique variance explained by syntactic, semantic and phonemic unexpectedness (quantified via surprise or prediction error) across cortical sources in each MEG participant. In all plots, colour indicates amount of additional variance explained; opacity indicates FWE-corrected statistical significance. Note that $p < 0.05$ is equivalent to $-\log_{10}(p) > 1.3$.

264 diction errors, implying that linguistic predictions
 265 are feature-specific and occur both at high and low
 266 levels of processing simultaneously.

267 Although we observed considerable variation in
 268 lateralisation and exact spatial locations between
 269 individuals, the overall pattern of sources aligned
 270 well with prior research on the neural circuits for
 271 each level. For instance, only for semantic predic-
 272 tion errors we observed a widely distributed set of
 273 neural sources – consistent with the fact that the
 274 semantic (but not the syntactic or phonological) sys-
 275 tem is widely distributed [38, 39]. Moreover, the
 276 temporal areas showing the strongest effect of syntac-
 277 tic surprise are indeed key areas for syntactic
 278 processing [40] and for the posterior temporal areas
 279 predictive syntax in particular [21, 41–43] – though
 280 a clear syntactic effect in the inferior frontal gyrus
 281 (IFG) was interestingly absent. When we compared
 282 the sources of phonemic surprise to those obtained
 283 for lexical surprise, we observed a striking overlap
 284 in all individuals (see Fig. S7, S4 and S13), suggesting
 285 that the phonemic predictions as formalised here

286 mostly relate to predictive (incremental) word recog-
 287 nition at the phoneme level rather than describing
 288 phonological or phonotactic predictions *per se*.

289 Dissociable signatures of syntactic, semantic 290 and phonemic predictions

291 Having established that syntactic, phonemic and
 292 semantic prediction errors independently modu-
 293 lated neural responses in different brain areas, we
 294 further investigated the nature of these effects. This
 295 was done by inspecting the coefficients (or modu-
 296 lation functions), which describe how fluctuations
 297 in a given regressor modulate the response over
 298 time. We first turn to the EEG data because there
 299 the sample size allows for population-level statisti-
 300 cal inference on the coefficients. We fitted the same
 301 integrated model (Figure S6) and performed cluster-
 302 based permutation tests on the modulation func-
 303 tions. This revealed significant effects for each type
 304 of prediction error (Figure 5).

305 First, syntactic surprise evoked an early, positive
 306 deflection ($p = 0.027$) based on a frontal cluster be-

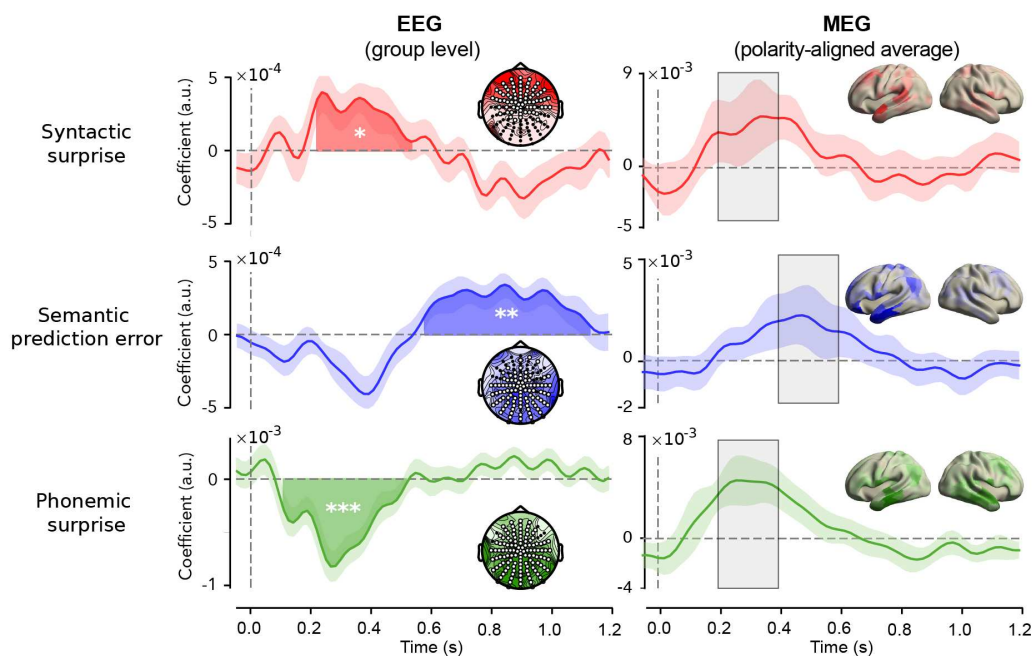


Figure 5: SPATIOTEMPORAL SIGNATURES OF SYNTACTIC, SEMANTIC AND PHONEMIC PREDICTION ERRORS.

Coefficients describing the effects of each prediction-error. EEG (left column): modulation functions averaged across the channels participating for at least one sample in the three main significant clusters (one per predictor). Highlighted area indicates temporal extent of the cluster. Shaded area around waveform indicates bootstrapped standard errors. Stars indicate cluster-level significance; $p < 0.05$ (*), $p < 0.05$ (**), $p < 0.001$ (***). Insets represent selected channels and distribution of absolute t-values. Note that these plots only visualise the effects; for the full topographies of the coefficients and respective statistics, see Figure S8. MEG (right column): polarity aligned responses averaged across the sources with significant explained variance (Figure 4) across participants. Shaded area represents absolute deviation. Insets represent topography of absolute value of coefficients averaged across the highlighted period. Note that due to polarity alignment, sign information is to be ignored for the MEG plots. For average coefficients for each source, see Figure S10 for coefficients of each individual, see Figs S11-S14.

307 tween 200 and 500 ms. This early frontal positivity
 308 converges with two recent studies that investigated
 309 specifically syntactic prediction using models trained
 310 explicitly on syntax [22, 44]. We also observed a late
 311 negative deflection for syntactic surprise ($p = 0.025$;
 312 Figure S9), but this was neither in line with earlier
 313 findings nor replicated in the MEG data. The semantic
 314 prediction error also evoked a positive effect
 315 ($p = 9.1 \times 10^{-3}$) but this was based on a much later,
 316 spatially distributed cluster between 600 and 1100
 317 ms. Although such a late positivity has been prominently
 318 associated with syntactic violations [8], there
 319 is also a considerable body of work reporting such
 320 late positivities for purely semantic anomalies [45]

321 which is more in line with the semantic prediction
 322 error as quantified here (see Discussion). Notably,
 323 we did not find a significant N400-like effect for semantic
 324 prediction error – possibly because this negative
 325 deflection was already explained by the overall
 326 lexical surprise, which was included as a nuisance re-
 327 gressor (Figure S10). Finally, the phonemic surprise
 328 evoked a negative effect ($p = 3 \times 10^{-4}$) based on
 329 an early, distributed cluster between 100 and 500
 330 ms. This effect was similar to the word-level surprise
 331 effect (Figure 2C and S10) but occurred earlier. This
 332 timecourse corresponds to recent studies using simi-
 333 lar regression-based techniques to study (predictive)
 334 phoneme processing in natural listening [13, 28, 46].

335 When we performed the same analysis on the
336 MEG data, we observed striking differences in the
337 exact shape and timing of the modulation functions
338 between individuals (see Figure S11-S14). While
339 this might partly reflect variance in the coefficients
340 due to inherent correlations between the variables,
341 it clearly also reflects true individual differences,
342 demonstrated by one of the strongest and least cor-
343 related regressors (the acoustics) also showing con-
344 siderable variability (see Figure S14). Overall how-
345 ever, we could recover a temporal pattern of effects
346 similar to the EEG results: phonemic and syntactic
347 surprise modulating early responses, and seman-
348 tic prediction error modulating later responses – al-
349 though not as late in the EEG data. This temporal
350 order holds on average (Figures 5 S10) and is espe-
351 cially clear within individuals (Figure S11-S13).

352 Overall, our results (Figure 45) demonstrate that
353 syntactic, phonemic and semantic prediction errors
354 evoke brain responses that are both temporally and
355 spatially dissociable. Specifically, while phonemic
356 and syntactic predictions modulate relatively early
357 neural responses (100-400 ms) in a set of focal
358 temporal (and frontal) areas that are key for syn-
359 tactic and phonetic/phonemic processing, seman-
360 tic predictions modulate later responses (>400 ms)
361 across a widely distributed set of areas across the
362 distributed semantic system. These results reveal
363 that linguistic prediction is not implemented by a
364 single system but occurs throughout the speech and
365 language network, forming a hierarchy of linguistic
366 predictions across all levels of analysis.

367 **Phoneme predictions reveal hierarchical infer-** 368 **ence**

369 Having established that the brain generates lin-
370 guistic predictions across multiple levels of analysis,
371 we finally asked whether predictions at different lev-
372 els might interact. One option is that they are encaps-
373 ulated: Predictions in separate systems might use
374 different information, for instance unfolding over
375 different timescales, rendering them independent.
376 Alternatively, predictions at different levels might in-
377 form and constrain each other, effectively converg-

378 ing into a single multilevel prediction – as suggested
379 by theories of hierarchical cortical prediction [5, 6
380 47].

381 One way to adjudicate between these hypothe-
382 ses is by evaluating different schemes of deriving
383 phoneme predictions. One possibility is that such
384 predictions are only based on information unfold-
385 ing over short timescales. In this scheme, the pre-
386 dicted probability of the next phoneme is derived
387 from the *cohort* of words that are compatible with
388 the phonemes presented so far, with each candi-
389 date word weighted by its overall frequency of oc-
390 currence (see Figure 6A). As such, this scheme pro-
391 poses a *single-level model*: phoneme predictions are
392 based only on information at the level of within-
393 word phoneme sequences unfolding over short
394 timescales, plus a fixed frequency-based prior (cap-
395 turing statistical knowledge of word frequencies
396 within a language).

397 Alternatively, phoneme predictions might not only
398 be based on sequences of phonemes within a word,
399 but also on the longer prior linguistic context. In this
400 case, the probability of the next phoneme would still
401 be derived from the cohort of words compatible with
402 the phonemes presented so far, but now each candi-
403 date word is not weighted by its overall frequency
404 but by its *contextual probability* (Figure 6A). Such a
405 model would be hierarchical, in the sense that pre-
406 dictions are based both – at the first level – on short
407 sequences of phonemes (i.e. of hundreds of millisec-
408 onds long), and on a contextual prior which itself is
409 based – at the higher level – on long sequences of
410 words (i.e. of tens of seconds to minutes long).

411 Here, the first model is more in line with the clas-
412 sic Cohort model of incremental (predictive) word
413 recognition, which suggests that context is only in-
414 tegrated after the selection and activation of lexical
415 candidates [48]. By contrast, the second model is
416 more in line with contemporary theories of hierar-
417 chical predictive processing which propose that high-
418 level cortical predictions (spanning larger spatial or
419 temporal scales) inform and shape low-level predic-
420 tions (spanning finer spatial or temporal scales) [47
421 49]. Interestingly, recent studies of phoneme pre-

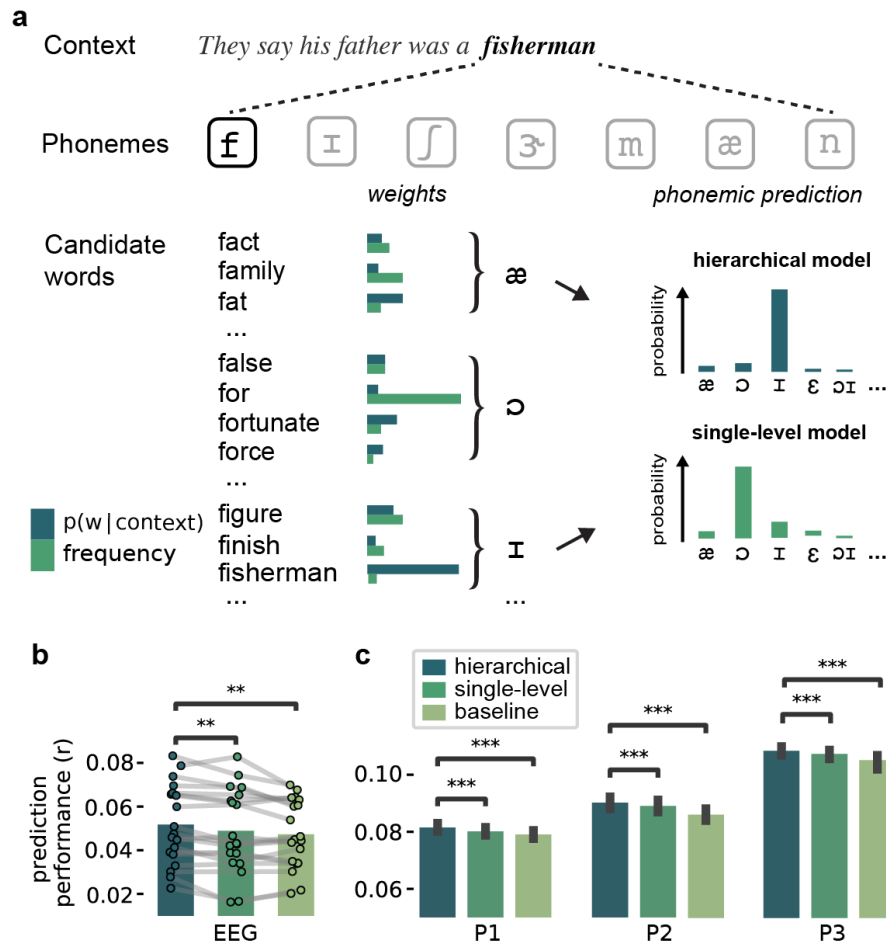


Figure 6: EVIDENCE FOR HIERARCHICAL INFERENCE DURING PHONEME PREDICTION.

a) Two models of phoneme prediction during incremental word recognition. Phonemic predictions were computed by grouping candidate words by their identifying next phoneme, and weighting each candidate word by its prior probability. This weight (or prior) could be either based on a word's overall probability of occurrence (i.e. frequency) or on its conditional probability in that context (from GPT-2). Critically, in the frequency-based model, phoneme predictions are based on a single level: short sequences of within words phonemes (hundreds of ms long) plus a fixed prior. By contrast, in the contextual model, predictions are based not just on short sequences of phonemes, but also on a contextual prior which is itself based on long sequences of prior words (up to minutes long), rendering the model hierarchical (see Methods). **b-c)** Model comparison results in EEG (**b**) and all MEG participants (**c**). EEG: dots with connecting lines represent individual participants (averaged over all channels). MEG: bars represent median across runs, error bars represent bootstrapped absolute deviance (averaged over language network sources). Significance levels correspond to $P < 0.01$ (**) or $P < 0.001$ (***) in a two-tailed paired t or Wilcoxon sign rank test.

422 ditions during natural listening have used both the
423 frequency-based single level model [27, 29] and a
424 context-based (hierarchical) model [13]. However,
425 the models have not been explicitly compared to test
426 which model can best account for prediction-related

427 fluctuations in neural responses to phonemes.

428 To compare these possibilities, we constructed 3
429 phoneme-level regression models (see Figure S15),
430 which all only included regressors at the level of
431 phonemes. First, the baseline model only included

432 non-predictive control variables: phoneme onsets,
433 acoustics, word boundaries and uniqueness points.
434 This can be seen as the phoneme-level equivalent of
435 the baseline model in Figure ???. The baseline model
436 was compared with two regression models which
437 additionally included phoneme surprise. In one of
438 the regression models, this was calculated using a
439 single-level model (with a fixed, frequency-based
440 prior), in the other regression model it was derived
441 from a hierarchical model (with a dynamic, contex-
442 tual prior derived from GPT-2). To improve our ability
443 to discriminate between the hierarchical and single-
444 level model, we not only included surprise but also
445 phoneme entropy (calculated with either model) as
446 a regressor [13].

447 When we compared the cross-validated predictive
448 performance, we first found that in both datasets
449 the predictive model performed significantly better
450 than the non-predictive baseline (Figure 6b-c hierar-
451 chical vs baseline, EEG: $t_{18} = 3.80$, $p = 1.31 \times 10^{-3}$;
452 MEG: all $p's < 5.69 \times 10^{-12}$). This replicates the
453 basic evidence for predictive processing but now
454 at the phoneme rather than word level (Figure ??).
455 Critically, when we compared the two predictive
456 models, we found that the hierarchical model per-
457 formed significantly better, both in EEG ($t_{18} = 3.03$,
458 $p = 7.28 \times 10^{-3}$) and MEG (all $p's < 9.44 \times 10^{-4}$).
459 This suggests that neural predictions of phonemes
460 (based on short sequences of within-word speech
461 sounds) are informed by lexical predictions,
462 effectively incorporating long sequences of prior
463 words as contexts. This is a signature of hierarchi-
464 cal prediction, supporting theories of hierarchical
465 predictive processing.

466 DISCUSSION

467 Across two independent data sets, we combined
468 deep neural language modelling with regression-
469 based deconvolution of human electrophysiological
470 (EEG and MEG) recordings to ask if and how evoked
471 responses to speech are modulated by linguistic
472 expectations that arise naturally while listening to
473 a story. Our results demonstrated that evoked re-
474 sponses are modulated by *probabilistic* predictions.

475 We then introduced a novel technique that allowed
476 us to quantify not just how much a linguistic stimu-
477 lus is surprising, but also at what level – phonemi-
478 cally, syntactically and/or semantically. This revealed
479 dissociable effects, in space and time, of different
480 types of prediction errors: syntactic and phonemic
481 prediction errors modulated early responses in a
482 set of focal, mostly temporal areas, while semantic
483 prediction errors modulated later responses across
484 a widely distributed set of cortical areas. Finally, we
485 found that phonemic prediction error signals were
486 best modelled by a hierarchical model incorporating
487 two levels of context: short sequences of within-
488 word phonemes (up to hundreds of milliseconds
489 long) and long sequences of prior words (up to min-
490 utes long). Together, these results demonstrate that
491 during natural listening, the brain is engaged in pre-
492 diction across multiple levels of linguistic represen-
493 tation, from speech sounds to meaning. The find-
494 ings underscore the ubiquity of prediction during
495 language processing, and fit naturally in predictive
496 processing accounts of language [1, 2] and neural
497 computation more broadly [5, 6, 49, 50].

498 A primary result of this paper is that evoked re-
499 sponses to words are best explained by a predic-
500 tive processing model: regression models including
501 unexpectedness performed better than strong non-
502 predictive baseline models, demonstrating that the
503 effects of prediction on brain responses cannot be
504 reduced to confounding simple features like seman-
505 tic incongruency. This aligns with recent ERP studies
506 aimed specifically at distinguishing prediction from
507 semantic integration [51, 52] and extends those find-
508 ings by analysing not just specific (highly predictable)
509 ‘target’ words, but *all* words in a natural story. In-
510 deed, when we further compared different accounts
511 of prediction, responses were best explained by a
512 regression model casting linguistic predictions as
513 ubiquitous and probabilistic. This supports the no-
514 tion of continuous, graded prediction – as opposed
515 to the classical view of prediction as the all-or-none
516 pre-activation of specific words in highly constrain-
517 ing contexts [33].

518 Because our deconvolution analysis focused on

519 evoked responses, the results can be linked to the
520 rich literature on linguistic violations using tradi-
521 tional ERP methods. This is powerfully illustrated
522 by the modulation function of lexical surprise (Fig-
523 ure 2b) tightly following the N400 modulation effect,
524 one of the first proposed, most robust and most
525 debated ERP signatures of linguistic prediction [7,
526 24, 30]. Similarly, the early negativity we found for
527 phonemic surprise and later positivity for seman-
528 tic prediction error (Fig 5) align well with N200 and
529 the semantic P600 or PNP effects of phonological
530 mismatch and semantic anomaly respectively [33,
531 53]. Unlike most ERP studies, we observed these
532 effects in participants listening to natural stimuli –
533 without any anomalies or violations – not engaged in
534 any task. This critically supports the idea that these
535 responses reflect deviations from *predictions* inher-
536 ent to the comprehension process – rather than
537 reflecting either detection of linguistic anomalies or
538 expectancy effects introduced by the experiment
539 [17, 19].

540 While we found several striking correspondences
541 between the modulation functions recovered from
542 the data and classic effects from the ERP literature,
543 there were also some differences. Specifically, for
544 syntactic surprise, we found neither a late positive
545 effect resembling the syntactic P600 [8] nor an early
546 negative effect akin to the ELAN [54]. One potential
547 explanation for this is that our formalisation (part-of-
548 speech surprise) might not fully capture syntactic vi-
549 olations used in ERP studies. Indeed, a recent paper
550 on syntactic prediction using a similar model-based
551 approach found a P600-like effect not for syntactic
552 surprise but for the number of syntactic reinter-
553 pretation attempts a word induced [22]. Conversely, the
554 early positive effect of syntactic surprise we found
555 – which replicated other model-based findings, de-
556 spite using a different formalisation of syntactic sur-
557 prise [22, 44] – does not have a clear counterpart in
558 the traditional ERP literature. Better understanding
559 such systematic differences between the traditional
560 experimental and model-based approach provides
561 an interesting challenge for future work.

562 Beyond the ERP literature, there has also been

563 earlier model-based work on prediction. How-
564 ever, these studies have mostly quantified feature-
565 unspecific lexical unexpectedness [10, 12, 32, 55,
566 56] or modelled feature-specific predictions at a sin-
567 gle level such as syntax [11, 22, 44, 57], phonemes
568 [13, 27, 28] or semantics [24]. We extend these
569 studies by probing predictions at all these levels si-
570 multaneously. This is important because it allows
571 to control for correlations between levels – since
572 words that are, for instance, syntactically surprising
573 are, on average, also semantically surprising. More-
574 over, prior modelling of feature-specific predictions
575 used domain-specific models that had to be inde-
576 pendently trained, and typically incorporated linguis-
577 tic context in a limited way. By contrast, our method
578 (Figure 3) allows to derive multiple predictions from
579 a single, large pre-trained model (like GPT-2) which
580 has a much deeper grasp of linguistic context. How-
581 ever, a limitation of this method is that the resulting
582 predictions are not independent. Therefore, you
583 cannot test if levels interact without *also* creating a
584 separate, domain-specific model. As such, the disen-
585 tangling approach we used is complementary to the
586 domain-specific modelling approach. Future work
587 could combine the two, for instance to test if the hi-
588 erarchical prediction we observed for phonemes ap-
589 plies to all linguistic levels – or whether predictions
590 at some levels (e.g. syntax) might be independent.

591 In this study, we combined group-level analysis
592 (of the EEG data) and individual-level analysis (of
593 the MEG data). These approaches are complemen-
594 tary. While including more participants allows one to
595 perform population-level inference, acquiring more
596 data per participant allows one to evaluate effects
597 within individuals. By combining both forms of analy-
598 sis, we found that on the one hand, the basic effects
599 of prediction and the comparison of hypotheses
600 about its computational nature (probabilistic predic-
601 tion, hierarchical prediction) were identical within
602 and across each individual (Figure 2, 6, S5). But on
603 the other hand, the exact spatiotemporal character-
604 istics of these effects showed substantial variability
605 (Figure 4, 5, S4, S7-S14). This suggest that while the
606 prediction effects themselves at the EEG group-level

607 are likely present in each individual, the precise spa-
608 tiotemporal signatures (Figure 5) are probably best
609 understood as a statistical average that is not neces-
610 sarily representative of underlying individuals.

611 Because our analysis focused on evoked re-
612 sponses, we chose to probe predictions indirectly:
613 via the neural markers of deviations from these pre-
614 dictions. As such, we cannot rule out that the ef-
615 fects might partly reflect ‘postdiction’. However, a
616 purely postdictive explanation appears unlikely as it
617 implies that after recognition, the brain computes a
618 prediction of the recognised stimulus based on infor-
619 mation available *before* recognition. While the data
620 therefore indirectly support pre-activation, the rep-
621 resentational format of these pre-activations is still
622 an open question. In our analyses – and many theo-
623 retical models [6, 49]) – predictions are formalised
624 as *explicit* probability distributions, but this is almost
625 certainly a simplification. It remains unclear whether
626 the brain represents probabilities implicitly. Alterna-
627 tively, it might use a kind of approximation: graded,
628 anticipatory processing that is perhaps functionally
629 equivalent to probabilistic processing, but avoids
630 having to represent (and compute with) probabili-
631 ties. A potential way to address this question is to
632 try to decode predictions before word onset [58].
633 Interestingly, this approach could be extended to
634 assess whether predicted probabilities are repre-
635 sented before onset at different levels of the linguis-
636 tic hierarchy, to test whether and which predicted
637 distributions are reflected in pre-stimulus activity.

638 Why would the brain constantly predict upcoming
639 language? Three – mutually non-exclusive – func-
640 tions have been proposed. First, predictions can
641 be used for *compression*: if predictable stimuli are
642 represented succinctly, this yields an efficient code
643 [6, 49] – conversely, optimising efficiency can make
644 predictive coding emerge in neural networks [59].
645 A second, perhaps more studied function is that
646 predictions can guide *inference*. Our analysis only
647 probed prediction errors, and hence does not speak
648 directly to such inferential effects of prediction – but
649 earlier work suggests that linguistic context can in-
650 deed enhance neural representations in a top-down

651 fashion [60, 61]; but see [62, 63]. Finally, predictions
652 may guide *learning*: prediction errors can be used to
653 perform error-driven learning without supervision.
654 While learning is perhaps the least-studied function
655 of linguistic prediction in cognitive neuroscience (but
656 see [16]), it is its primary application in Artificial Intel-
657 ligence [64, 65]. In fact, the language model we used
658 (GPT-2) was created to study such predictive learn-
659 ing. These models are trained only to predict words,
660 but learn about language more broadly, and can
661 then be applied to practically any linguistic task [34,
662 65]. Interestingly, models trained with this predic-
663 tive objective also develop representations that are
664 ‘brain-like’, in the sense that they are currently the
665 best encoders of linguistic stimuli to predict brain
666 responses [66–69]. And yet, these predictive mod-
667 els are also brain-unlike in an interesting way – they
668 predict upcoming language only at a single (typically
669 lexical) level.

670 When prediction is used for compression or infer-
671 ence, it seems useful to predict at multiple levels,
672 since redundancies and ambiguities also occur at
673 multiple levels. But if predictions drive learning, why
674 would the brain predict at multiple levels, when ef-
675 fective learning can be achieved using simple, single-
676 level prediction? One fascinating option is that it
677 might reflect the brain’s way to perform credit as-
678 signment within biological constraints. In artificial
679 networks, credit assignment is typically done by first
680 *externally* computing a single, global error term, and
681 then ‘backpropagating’ this error through all levels
682 of the network – but both these steps are biolog-
683 ically implausible [70]. Interestingly, it has been
684 shown that hierarchical predictive coding networks
685 can approximate or even implement classical back-
686 propagation while using only Hebbian plasticity and
687 local error computation [6, 70, 71]. Therefore, if
688 the brain uses predictive error-driven learning, one
689 might expect such prediction to be hierarchical, so
690 error-terms can be locally computed throughout the
691 hierarchy – which is in line with what we find.

692 Beyond the domain of language, there have been
693 other reports of hierarchies of neural prediction, but
694 these have been limited to artificial, predictive tasks

695 or to restricted representational spans, such as suc-
696 cessive stages in the visual system [72-74]. Our re-
697 sults demonstrate that even during passive listening
698 of natural stimuli, the brain is engaged in prediction
699 across disparate levels of abstraction (from speech
700 sounds to meaning) based on timescales separated
701 by three orders of magnitude (hundreds of millisec-
702 onds to minutes). These findings provide important
703 evidence for hierarchical predictive processing in
704 cortex. As such, they highlight how language pro-
705 cessing in the brain is shaped by a domain-general
706 neurocomputational principle: the prediction of per-
707 ceptual inputs across multiple levels of abstraction.

708 METHODS

709 We analysed EEG and source localised MEG data from
710 two experiments. The EEG data is part of a public dataset
711 that has been published about before [27].

712 Participants

713 All participants were native English speakers. In the
714 EEG experiment, 19 subjects (13 male) between 19 and 38
715 years old participated; in the MEG experiment, 3 subjects
716 participated (2 male) aged 35, 30, and 28. Both exper-
717 iments were approved by local ethics committees (EEG:
718 ethics committee of the School of Psychology at Trinity
719 College Dublin; MEG: CMO region Arnhem-Nijmegen).

720 Stimuli and procedure

721 In both experiments, participants were presented con-
722 tinuous segments of narrative speech extracted from au-
723 diobooks. The EEG experiment used a recording of Hem-
724 ington's *The Old Man and the Sea*. The MEG experiment
725 used 10 stories from the *The Adventures of Sherlock Holmes*
726 by Arthur Conan Doyle. In total, EEG subjects listened to
727 ~1 hour of speech (containing ~11,000 words and ~35,000
728 phonemes); MEG subjects listened to ~9 hours of speech
729 (containing ~85,000 words and ~290,000 phonemes).

730 In the EEG experiment, each participants performed
731 only a single session, which consisted of 20 runs of 180s
732 long, amounting to the first hour of the book. Partici-
733 pants were instructed to maintain fixation and minimise
734 movements but were otherwise not engaged in any task.

735 In the MEG experiment, each participant performed a
736 total of ten sessions, each ~1 hour long. Each session was
737 subdivided in 6-7 runs of roughly ten minutes, although

738 the duration varied as breaks only occurred at meaning-
739 ful moments (making sure, for example, that prominent
740 narrative events were not split across runs). Unlike in the
741 EEG experiment, participants in the MEG dataset partici-
742 pants were asked to listen attentively and had to answer
743 questions in between runs: one multiple choice compre-
744 hension question, a question about story appreciation
745 (scale 1-7) and a question about informativeness.

746 MRI acquisition and headcast construction

747 To produce the headcast, we needed to obtain accurate
748 images of the participants's scalp surface, which were
749 obtained using structural MRI scans with a 3T MAGNETOM
750 Skyra MR scanner (Siemens AG). We used a fast low angle
751 shot (FAST) sequence with the following image acquisition
752 parameters: slice thickness of 1 mm; field-of-view of 256
753 \times 256 \times 208 mm along the phase, read, and partition
754 directions respectively; TE/TR = 1.59/4.5 ms.

755 Data acquisition and pre-processing

756 The EEG data were originally acquired using a 128-
757 channel (plus two mastoid channels) using an ActiveTwo
758 system (BioSemi) at a rate of 512 Hz, and downsampled
759 to 128 Hz before being distributed as a public dataset. We
760 visually inspected the raw data to identify bad channels,
761 and performed independent component analysis (ICA) to
762 identify and remove blinks; rejected channels were linearly
763 interpolated with nearest neighbour interpolation using
764 MNE-python.

765 The MEG data were acquired using a 275 axial gra-
766 diometer system at 1200 Hz. For the MEG data, prepro-
767 cessing and source modelling was performed in MATLAB
768 2018b using fieldtrip [75]. We applied notch filtering (But-
769 terworth IIR) at the bandwidth of 49-51, 99-101, and 149-
770 151 Hz to remove line noise. Artifacts related to muscle
771 contraction and squidjumps were identified and removed
772 using fieldtrip's semi-automatic rejection procedure. The
773 data were downsampled to 150 Hz. To identify and re-
774 move eye blink artifacts, ICA was performed using the
775 FastICA algorithm.

776 For both MEG and EEG analyses, we focus on the slow,
777 evoked response and hence restricted our analysis to low-
778 frequency components. To this end, we filtered the data
779 between 0.5 and 8 Hz using a bidirectional FIR bandpass
780 filter. Restricting the analysis to such a limited range of low
781 frequencies (which are known to best follow the stimulus)
782 is common when using regression ERP or TRF analysis,
783 especially when the regressors are sparse impulses [28]

784 [31, 36]. The particular upper bound of 8 Hz is arbitrary but
785 was based on earlier papers using the same EEG dataset
786 to study how EEG tracks acoustic and linguistic content of
787 speech [31, 56, 61].

788 Head and source models

789 The MEG sensors were co-registered to the subjects'
790 anatomical MRIs using position information of three local-
791 ization coils attached to the headcasts. To create source
792 models, FSL's Brain Extraction Tool was used to strip non-
793 brain tissue. Subject-specific cortical surfaces were recon-
794 structed using Freesurfer, and post-processing (downsam-
795 pling and surface-based alignment) of the reconstructed
796 cortical surfaces was performed using the Connectome
797 Workbench command-line tools (v 1.1.1). This resulted
798 in cortically-constrained source models with 7,842 source
799 locations per hemisphere. We created single-shell volume
800 conduction models based on the inner surface of the skull
801 to compute the forward projection matrices (leadfields).

802 Beamformer and parcellation

803 To estimate the source time series from the MEG data,
804 we used linearly constrained minimum variance (LCMV)
805 beamforming, performed separately for each session, us-
806 ing Fieldtrip's `ft_sourceanalysis` routine. To reduce
807 the dimensionality, sources were parcellated, based on a
808 refined version of the Conte69 atlas, which is based on
809 Brodmann's areas. We computed, for each session, parcel-
810 based time series by taking the first principal component
811 of the aggregated time series of the dipoles belonging to
812 the same cortical parcel.

813 Self-attentional language model

814 Contextual predictions were quantified using GPT-2
815 – a large, pre-trained language model [34]. Formally,
816 a language model can be cast as a way of assigning a
817 probability to a sequence of words (or other symbols),
818 (x_1, x_2, \dots, x_n) . Because of the sequential nature of lan-
819 guage, the joint probability, $P(X)$ can, via the chain rule,
820 be factorised as the product of conditional probabilities:

$$\begin{aligned} P(X) &= p(x_1) \times p(x_2 | x_1) \times \dots \times p(x_n | x_{n-1}, \dots, x_1) \\ &= \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1}) \end{aligned} \quad (1)$$

821 Since the advent of neural language models, as op-
822 posed to statistical (Markov) models, methods to compute

823 these conditional probabilities have strongly improved.
824 Improvements have been especially striking in the past
825 two years with the introduction of the *Transformer* [76]
826 architecture, which allows efficient training of very large
827 networks on large, diverse data. This resulted in models
828 that dramatically improved the state-of-the art in language
829 modelling on a range of domains.

830 GPT-2 [34] is one of these large, transformer-based lan-
831 guage models and is currently among the best publicly
832 released models of English. The architecture of GPT-2 is
833 based on the decoder-only version of the transformer. In
834 a single forward pass, it takes a sequence of tokens $U =$
835 (u_1, \dots, u_k) and computes a sequence of conditional
836 probabilities, $(p(u_1), p(u_2|u_1), \dots, p(u_k | u_1, \dots, u_{k-1}))$.
837 Roughly, the full model (see Figure S1) consists of three
838 steps: first, an embedding step encodes the sequence of
839 symbolic tokens as a sequence of vectors which can be
840 seen as the first hidden state h_0 . Then, a stack of trans-
841 former blocks, repeated n times, each apply a series of
842 operations resulting in a new set of hidden states h_l , for
843 each block l . Finally, a (log-)softmax layer is applied to
844 compute (log-)probabilities over target tokens. Formally,
845 then, the model can be summarised in three equations:

$$h_0 = UW_e + W_p \quad (2)$$

$$h_l = \text{transformer_block}(h_{l-1}) \quad \forall l \in [1, n] \quad (3)$$

$$P(u) = \text{softmax}(h_n W_e^T), \quad (4)$$

846 where W_e is the token embedding and W_p is the posi-
847 tion embedding (see below).

848 The most important component of the transformer-
849 block is the *masked multi-headed self-attention* (Fig S1). The
850 key operation is self-attention, a seq2seq operation turn-
851 ing a sequence of input vectors $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k)$ into a
852 sequence of output vectors $(\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_k)$. Fundamen-
853 tally, each output vector \mathbf{y}_i is a weighted average of the
854 input vectors: $\mathbf{y}_i = \sum_{j=1}^k w_{ij} \mathbf{x}_j$. Critically, the weight
855 $w_{i,j}$ is not a parameter but is *derived* from a function over
856 input vectors \mathbf{x}_i and \mathbf{x}_j . The Transformer uses (scaled) *dot*
857 *product attention*, meaning that the function is simply a dot
858 product between the input vectors $\mathbf{x}_i^T \mathbf{x}_j$, passed through
859 a softmax make sure that the weights sum to one, scaled
860 by a constant determined by the dimensionality, $\frac{1}{\sqrt{d_k}}$ (to
861 avoid the dot-products growing too large in magnitude):
862 $w_{ij} = \frac{\exp \mathbf{x}_i^T \mathbf{x}_j / \sum_{j=1}^k \exp \mathbf{x}_i^T \mathbf{x}_j}{\sqrt{d_k}}$.

863 In self-attention, then, each input \mathbf{x}_i is used in three
864 ways. First, it is multiplied by the other vectors to derive

865 the weights for its own output, y_i (as the *query*). Second, it
866 is multiplied by the other vectors to determine the weight
867 for any other output y_j (as the *key*). Finally, to compute
868 the actual outputs it is used in the weighted sum (as the
869 *value*). Different (learned) linear transformations are ap-
870 plied to the vectors in each of these use cases, resulting in
871 the Query, Key and Value matrices (Q, K, V). Putting this
872 all together, we arrive at the following equation:

$$\text{self_attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V, \quad (5)$$

873 where d_k is dimension of the keys/queries. In other words,
874 self_attention simply computes a weighted sum of the
875 values, where the weight of each value is determined by
876 the dot-product similarity of the query with its key. Be-
877 cause the queries, keys and values are linear transforma-
878 tions of the same vectors, the input *attends itself*.

879 To be used as a language model, two elements need
880 to be added. First, the basic self-attention operation is
881 not sensitive to the order of the vectors: if the order of
882 the input vectors is permuted, the output vectors will be
883 identical (but permuted). To make it position-sensitive, a
884 position embedding W_p is simply added during the em-
885 bedding step – see Equation 2. Second, to enforce that
886 the model only uses information from one direction (i.e.
887 left), a mask is applied to the attention weights (before
888 the softmax) which sets all elements above the diagonal
889 to $-\infty$. This makes the self-attention *masked*.

890 To give the model more flexibility, each transformer
891 block actually contains multiple instances of the basic self-
892 attention mechanisms from 5. Each instance (each *head*)
893 applies different linear transformations to turn the same
894 input vectors into a different set of Q, K and V matrices,
895 returning a different set of output vectors. The outputs of
896 all heads are concatenated and then reduced to the initial
897 dimensionality with a linear transformation. This makes
898 the self-attention *multi-headed*.

899 In total, GPT-2 (XL) contains $n = 48$ blocks, with 12
900 heads each; a dimensionality of $d = 1600$ and a context
901 window of $k = 1024$, yielding a total 1.5×10^9 parameters.
902 We used the PyTorch implementation of GPT-2 provided
903 by HuggingFace's *Transformers* package 77.

904 Lexical predictions

905 We passed the raw texts through GPT-2 (Equations 2,4)
906 for each run independently (assuming that listeners' ex-
907 pectations would to some extent 'reset' during the break).

908 This resulted in a (log-)probability distribution over to-
909 kens $P(U)$. Since GPT-2 uses Byte-Pair Encoding, a token
910 can be either punctuation or a word or (for less frequent
911 words) a word-part. How many words actually fit into a
912 context window of length k therefore depends on the text.
913 For words spanning multiple tokens, we computed word
914 probabilities simply as the joint probability of the tokens.
915 For window-placement, we used the constraint that the
916 windows had an overlap of at least 700 tokens, and that
917 they could not start mid-sentence (ensuring that the first
918 sentence of the window was always well-formed).

919 As such, for each word w_i we computed $p(w_i|\text{context})$,
920 where 'context' consisted either of all preceding words in
921 the run, or of a sequence of prior words constituting a
922 well-formed context that was at least 700 tokens long.

923 Syntactic and semantic predictions

924 Feature-specific predictions were computed from the
925 lexical prediction. To this end, we first truncated the un-
926 reliable tail from the distribution using a combination of
927 top-k and nucleus truncation. The nucleus was defined
928 as the "top" k tokens with the highest predicted probabili-
929 ty, where k was set dynamically such that the cumulative
930 probability was at least 0.9. To have enough information
931 also for very low entropy cases (where k becomes small),
932 we forced k to be at least 40.

933 From this truncated distribution, we derived feature-
934 specific predictions by analysing the predicted words. For
935 the syntactic predictions, we performed part of speech
936 tagging on every potential sentence (i.e. the context plus
937 the predicted word) with Spacy to derive the probability
938 distribution over parts-of-speech, from which the syntactic
939 surprise was calculated as the negative log probability of
940 the POS of a word, $-\log(P(\text{POS}_n | \text{context}))$.

941 For the semantic prediction, we took a weighted
942 average of the glove embeddings of the predicted
943 words to compute the expected vector: $\mathbb{E}[G(w_n)] =$
944 $\sum_{i=1}^k P(x_i) G(x_i)$, where $G(w_i)$ is the GloVe embedding
945 for predicted word w_i . From this prediction, we computed
946 the semantic prediction error as the cosine distance be-
947 tween the predicted and observed vector:

$$\text{PE}_{\text{semantic}} = 1 - \frac{\mathbb{E}[G(w_n)] \cdot G(w_n)}{\|\mathbb{E}[G(w_n)]\| \|G(w_n)\|} \quad (6)$$

948 Phonemic predictions

949 Phonemic predictions were formalised in the context
950 of incremental word recognition 27, 29. This process
951 can be cast as probabilistic prediction by assuming that

952 brain is tracking the *cohort* of candidate words consistent
953 with the phonemes so far, each word weighted by its prior
954 probability. We compared two such models that differed
955 only in the prior probability assigned to each word.

956 The first model was the single-level or frequency-
957 weighted model (Fig 6), in which prior probability of words
958 was fixed and defined by a word's overall probability of
959 occurrence (i.e. lexical frequency). The probability of a
960 specific phoneme (A), given the prior phonemes within a
961 word, was then calculated using the statistical definition:

$$P(\varphi_t = A \mid \varphi_{1:t-1}) = \frac{f(C_{\varphi_t=A})}{f(C_{\varphi_{1:t-1}})}. \quad (7)$$

962 Here, $f(C_{\varphi_t=A})$ denotes the cumulative frequency of
963 all words in the remaining cohort of candidate words
964 if the next phoneme were A , and $f(C_{\varphi_{(1:t-1)}})$ denotes
965 the cumulative frequency of all words in the prior cohort
966 (equivalent to $f(C)$ of all potential continuations). If a cer-
967 tain continuation did not exist and the cohort was empty,
968 $f(C_{\varphi_t=A})$ was assigned a laplacian pseudocount of 1. To
969 efficiently compute 7 for every phoneme, we constructed
970 a statistical phonetic dictionary as a digital tree that com-
971 bined frequency information from SUBTLEX database and
972 pronunciation from the CMU dictionary.

973 The second model was equivalent to the first model,
974 except that the prior probability of each word was not de-
975 fined by its overall probability of occurrence, but by its con-
976 ditional probability in that context (based on GPT-2). This
977 was implemented by constructing a separate phonetic dic-
978 tionary for every word, in which lexical frequencies were
979 replaced by implied counts derived from the lexical predic-
980 tion. We truncated the unreliable tail from the distribution
981 and replaced that by a flat tail that assigned each word
982 a pseudocount of 1. This greatly simplifies the problem
983 as it only requires to assign implied counts for the top k
984 predicted words in the dynamic nucleus. Since all counts
985 in the tail are 1, the cumulative implied counts of the nu-
986 cleus is complementary to the the length of the tail, which
987 is simply the difference between the vocabulary size and
988 nucleus size ($V - k$). As such a little algebra reveals:

$$\text{freqs}_n = P_{tr}(w^{(i)}|\text{context}) \frac{V - k}{1 - \sum_{j=1}^k P(w_j^{(i)}|\text{context})}, \quad (8)$$

989 where $P_{tr}(w^{(i)}|\text{context})$ is the truncated lexical lexical
990 prediction, and $P(w_j^{(i)}|\text{context})$ is predicted probability
991 that word i in the text is word j in the sorted vocabulary.

992 Although we computed probabilities using the simple
993 statistical definition of probability, these two ways of as-

994 signing lexical frequencies are equivalent to two kinds of
995 priors in a Bayesian model. Specifically, in the first model
996 the prior over words is the fixed unconditional word proba-
997 bility, while in the second model the prior is the contextual
998 probability, itself based on a higher level (lexical) predic-
999 tion. This makes the second computation *hierarchical* be-
1000 cause phoneme predictions are based on not just (at the
1001 first level) on short sequences of within-word phonemes,
1002 but also on a contextual prior which itself (at the second
1003 level) is based on long sequences of prior words.

1004 Non-predictive control variables

1005 To ensure we were probing effects of predictions, we
1006 had to control for various non-predictive variables: onsets,
1007 acoustics, frequency and semantic congruency. We will
1008 briefly outline our definitions of each.

1009 For speech, it is known that the cortical responses are
1010 sensitive to fluctuations in the envelope – which is specifi-
1011 cally driven by rapid increases of the envelope amplitude
1012 (or ‘acoustic edges’) 78. To capture these fluctuations
1013 in a sparse, impulse-based regressor we quantified the
1014 amplitude of these edges as the variance of the envelope
1015 over each event (e.g. phoneme) following 61. A sec-
1016 ond non-predictive variable is frequency. We accounted
1017 for frequency as the overall base rate or unconditional
1018 probability of a word, defining it similarly to lexical sur-
1019 prise as the unigrams surprise – $\log P(\text{word})$ based on
1020 its frequency of occurrence in subtext.

1021 The final non-predictive variable was semantic congru-
1022 ency or integration difficulty. This speaks to the debate
1023 whether effects of predictability reflect prediction or rather
1024 post-hoc effects arising when integrating a word into the
1025 semantic context. This can be illustrated by considering
1026 a constraining context (‘coffee with milk and ...’). When
1027 we contrast a highly expected word (‘sugar’) and an unex-
1028 pected word (e.g. ‘dog’), the unexpected word is not just
1029 less likely, but also semantically incongruous in the prior
1030 context. As such, the increased processing cost reflected
1031 by effects like N400 increases might not (only) be due to
1032 a violated *prediction* but due to difficulty integrating the
1033 target word (‘dog’) in the semantic context (‘coffee with
1034 milk’) 7 18 51 52. As a proxy for semantic integration
1035 difficulty we computed the semantic congruency of a word
1036 in its context defined as the cosine dissimilarity (see 6)
1037 between the average semantic vector of the prior context
1038 words and the target content word, following 31. This
1039 metric is known to predict N400-like modulations and can
1040 hence capture the extent to which such effects can be

1041 explained by semantic congruency only [31, 52].

1042 Word-level regression models

1043 The word-level models (see Fig S2 for graphical repre-
1044 sentation) captured neural responses to words as a func-
1045 tion of word-level variables. The *baseline* model formalised
1046 the hypothesis that responses to words were not affected
1047 by word unexpectedness but only by the following non-
1048 predictive confounds: word onsets, envelope variability
1049 (acoustic edges), semantic congruency (integration diffi-
1050 culty) and word frequency.

1051 The *probabilistic prediction* model formalised the hy-
1052 pothesis that predictions were continuous and probabilis-
1053 tic. This model was identical to the baseline model plus
1054 the lexical surprise (or negative log probability of a word),
1055 for every word. This was based on normative theories of
1056 predictive processing which state that the brain response
1057 to a stimulus should be proportional to the negative log
1058 probability of that stimulus [6].

1059 The *constrained guessing* model formalised the classical
1060 psycholinguistic notion of prediction as the all-or-none pre-
1061 activation of specific words in specific (highly constraining)
1062 contexts [33]. We translated the idea of all-or-none predic-
1063 tion into a regression model using an insight by Smith and
1064 Levy which implied that all-or-none predictions result in
1065 a linear relationship between word probability and brain
1066 responses [9]. The argument follows from two assump-
1067 tions: (1) all predictions are all-or-none; and (2) incorrect
1068 predictions incur a cost, expressed as a prediction error
1069 brain response (fixed in size because of assumption 1).
1070 For simplicity, we first consider the unconstrained case
1071 (i.e. subjects make a prediction for *every* stimulus), and
1072 we bracket all other factors affecting brain responses by
1073 absolving them into an average brain response, y_{baseline} .
1074 As such, the response to any word is either y_{baseline} (if the
1075 prediction is correct) or $y_{\text{baseline}} + y_{\text{error}}$ (if it was false).
1076 For any individual stimulus, this equation cannot be used
1077 (as we don't know what a subject predicted). But if we
1078 assume that predictions are approximately correct, then
1079 the probability of a given prediction to be incorrect simply
1080 becomes $\sim(1 - p)$. As such, *on average*, the response
1081 becomes $y_{\text{resp}} = y_{\text{baseline}} + (1 - p)y_{\text{error}}$. In other words,
1082 a linear function of word improbability. To extend this
1083 to the constrained case, we only define the improbability
1084 regressor for constraining contexts, and add a constant
1085 to those events to capture (e.g. suppressive) effects of
1086 correct predictions (Figure S2). To identify 'constraining
1087 contexts', we simply took the 10% of words with the lowest

1088 prior lexical entropy. The choice of 10% was arbitrary –
1089 however, using a slightly more or less stringent definition
1090 would not have changed the results because the naive
1091 guessing model (which included linear improbability for
1092 *every* word) performed so much better (see Figure S5).

1093 Integrated regression model

1094 For all analyses on feature-specific predictions, we for-
1095 mulated an integrated regression model with both word-
1096 level and phoneme-level regressors (Figure S6). To avoid
1097 collinearity between word and phoneme level regressors,
1098 phoneme-level regressors were only defined for word-
1099 non-initial phonemes, and word-level regressors were de-
1100 fine for word-onset. As regressors of interest this model
1101 included phonemic surprise, syntactic surprise and se-
1102 mantic prediction error. In principle, we could have also
1103 included phoneme and syntactic entropy rather than just
1104 surprise (e.g. [13]) – however, these were highly corre-
1105 lated with the respective surprise. Since this was already
1106 a complex regression model, including more correlated
1107 regressors would have made the coefficients estimates
1108 less reliable and hence more difficult to interpret. As such,
1109 we did not include both but focussed on surprise because
1110 it has the most direct relation to stimulus evoked effect.

1111 Phoneme-level regression models

1112 To compare different accounts of phoneme prediction,
1113 we formulated three regression models with only regres-
1114 sors at the individual phoneme level (Figure S15). In all
1115 models, following [27] we used separate regressors for
1116 word-initial and word-non-initial phonemes, to account
1117 for juncture phonemes being processed differently. The
1118 baseline model only included non-predictive factors of
1119 word-boundaries, phoneme onsets, envelope variability,
1120 and uniqueness points. The two additional models also
1121 included phoneme surprise and phoneme entropy from
1122 either the hierarchical model or non-hierarchical model.
1123 To maximise our ability to dissociate the hierarchical pre-
1124 diction and non-hierarchical prediction, we included both
1125 entropy and surprise. Although these metrics are corre-
1126 lated, adding both should add more information to the
1127 model-comparison, assuming that there is some effect of
1128 entropy [13]. (Note that here, we were only interested in
1129 model comparison, and not in comparing the coefficients,
1130 which may become more difficult when including both.)

1131 Time resolved regression

1132 As we were interested in the evoked responses, vari-
1133 ables were regressed against EEG data using time-resolved
1134 regression, within a regression ERP/F (or impulse TRF)
1135 framework [31, 35]. Briefly, this involves using impulse
1136 regressors for both constants and covariates defined at
1137 event onsets, and then temporally expanding the design
1138 matrix such that each predictor column C becomes a se-
1139 ries of columns over a range of temporal lags $C_{t_{min}}^{t_{max}} =$
1140 $(C_{t_{min}}, \dots, C_{t_{max}})$. For each predictor one thus estimates
1141 a series of weights $\beta_{t_{min}}^{t_{max}}$ (Fig 1) which can be understood
1142 as the *modulation function* describing how a given regres-
1143 sor modulates the neural response over time, and which
1144 corresponds to the *effective* evoked response that would
1145 have been obtained in a time-locked ERP/ERF design. Here,
1146 we used a range between -0.2 and 1.2 seconds. All data
1147 and regressors were standardised and coefficients were
1148 estimated with ℓ_2 -norm regularised (Ridge) regression:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \|y - X\beta\|_2^2 + \lambda\|\beta\|_2^2, \quad (9)$$

1149 using the scikit learn sparse matrix implementation. In
1150 both datasets, models were estimated by concatenating
1151 the (time-expanded) design matrix across all runs and ses-
1152 sions. Regularisation was set based on leave-one-run-out
1153 R^2 comparison; for inference on the weights in the EEG
1154 data this was done across subjects to avoid doing statistics
1155 over coefficients with different amounts of shrinkage.

1156 Model comparison

1157 In both datasets, model comparison was based on
1158 comparing cross-validated correlation coefficients. Cross-
1159 validation was performed in a leave-one-run-out cross-
1160 validation scheme, amounting to 19-fold cross-validation
1161 in the EEG data and between 63 and 65-fold cross-
1162 validation for the MEG data (in some subjects, some runs
1163 were discarded due to technical problems).

1164 For the EEG data, models' cross-validated prediction
1165 performance was performed across subjects to perform
1166 population-level inference. To this end, we reduced the
1167 scores into a single n_{subs} dimensional vector by taking
1168 the median across folds and the mean across channels.
1169 Critically, we did not select any channels but used the av-
1170 erage across the scalp. For the MEG data, models were
1171 only statistically compared on a within within-subject ba-
1172 sis. Because the MEG data was source localised we could
1173 discard sources known to be of no interest (e.g. early vi-
1174 sual cortex). To this end, we focussed on the language

1175 network, using a rather unconstrained definition encom-
1176 passing all Brodmann areas in the temporal lobe, plus the
1177 temporo-parietal junction, and inferior frontal gyrus and
1178 dorsolateral prefrontal cortex; all bilaterally (see Figure
1179 S16).

1180 Statistical testing

1181 All statistical tests were two-tailed and used an alpha of
1182 0.05. For all simple univariate tests performed to compare
1183 model-performance within and between subjects, we first
1184 verified that the distribution of the data did not violate nor-
1185 mality and was outlier free, determined by the D'Agostino
1186 and Pearson's test implemented in SciPy and the 1.5 IQR
1187 criterion, respectively. If both criteria were met, we used a
1188 parametric test (e.g. paired t-test); otherwise, we resorted
1189 to a non-parametric alternative (e.g. Wilcoxon sign rank).

1190 In EEG, we performed mass-univariate tests on the
1191 coefficients across participants between 0 and 1.2 sec-
1192 onds. This was firstly done using cluster-based permuta-
1193 tion tests [79, 80] to identify clustered significant effects
1194 as in Figure 5 (10,000 permutations per test). Because
1195 the clustered effects as in Figure 5 only provide a partial
1196 view, we also reported more comprehensive picture of the
1197 coefficients across all channels (Figure S3-S8); there, we
1198 also provide multiple-comparison corrected p-values to
1199 indicate statistical consistency of the effects; these were
1200 computed using TFCE. In the MEG, multiple comparison
1201 correction for comparison of explained variance across
1202 cortical areas was done using Threshold Free Cluster En-
1203 hancement (TFCE). In both datasets, mass-univariate test-
1204 ing was performed based on one-sample t-tests plus the
1205 'hat' variance adjustment method with $\sigma = 10^{-3}$.

1206 Polarity-alignment

1207 In the source localised MEG data, the coefficients in
1208 individuals (e.g. Figure S11-S14) are symmetric in polar-
1209 ity, with the different sources in a single response having
1210 an arbitrary sign due to ambiguity of the source polar-
1211 ity. To harmonise the polarities, and avoid cancellation
1212 when visualising the average coefficient, we performed
1213 a polarity-alignment procedure. This was based on first
1214 performing SVD, $\mathbf{A} = \mathbf{A}\mathbf{\Sigma}\mathbf{V}^T$, where \mathbf{A} is the $m \times n$
1215 coefficient matrix, with m being the number of sources
1216 and n the number of regressors; and then multiplying
1217 each row of \mathbf{A} by the sign of the first right singular vector.
1218 Because the right singular vectors (columns of \mathbf{U}) can be
1219 interpreted as the eigen vectors of the source-by-source
1220 correlation matrix, this can be thought of as flipping the

1221 sign of each source as a function of its polarity with respect
1222 to the dominant correlation. This procedure was used for
1223 visualisation purposes only (see Fig S4 and S11 S14).

1224 Data and code availability

1225 Data and code to reproduce all results will be made
1226 public at the Donders Repository. The full MEG dataset
1227 will be made public in a separate resource publication.

1228 Acknowledgements

1229 This work was supported by The Netherlands Organisa-
1230 tion for Scientific Research (NWO Research Talent grant to
1231 M.H.; NWO Vidi grant to F.P.d.L.; NWO Vidi 864.14.011 to
1232 JMS; Gravitation Program Grant Language in Interaction
1233 no. 024.001.006 to P.H.) and the European Union Horizon
1234 2020 Program (ERC Starting Grant 678286, 'Contextvision'
1235 to F.P.d.L.). We wish to thank Michael P Broderick, Giovanni
1236 M. Di Liberto, and colleagues from the Lalor lab for making
1237 the EEG dataset openly available. We thank all the authors
1238 of the open source software we used and apologise for
1239 citation limits that prevent us from citing all tools used.

1240 Contributions

1241 Conceptualisation: MH, FPdL, PH; Formal analysis: MH;
1242 Data collection: KA, JMS; Source modelling: KA, JMS; Original
1243 draft: MH; Final manuscript: MH, FPdL, PH, JMS,KA.

1244 REFERENCES

- 1245 1. Kuperberg, G. R. & Jaeger, T. F. What do we
1246 mean by prediction in language
1247 comprehension? *Language, cognition and*
1248 *neuroscience* **31**, 32–59. ISSN: 2327-3798
1249 (2016).
- 1250 2. Kutas, M., DeLong, K. A. & Smith, N. J. in
1251 *Predictions in the brain: Using our past to*
1252 *generate a future* 190–207 (Oxford University
1253 Press, New York, NY, US, 2011). ISBN:
1254 978-0-19-539551-8. doi:10.1093/acprof:
1255 [oso/9780195395518.003.0065](https://doi.org/10.1093/acprof:oso/9780195395518.003.0065).
- 1256 3. Jelinek, F. *Statistical methods for speech*
1257 *recognition* ISBN: 978-0-262-10066-3 (MIT Press,
1258 Cambridge, MA, USA, 1998).

- 1259 4. Graves, A., Mohamed, A.-r. & Hinton, G. *Speech*
1260 *recognition with deep recurrent neural networks*
1261 *in 2013 IEEE international conference on*
1262 *acoustics, speech and signal processing (IEEE,*
1263 *2013)*, 6645–6649.
- 1264 5. Keller, G. B. & Mrsic-Flogel, T. D. Predictive
1265 Processing: A Canonical Cortical Computation.
1266 en. *Neuron* **100**, 424–435. ISSN: 0896-6273
1267 (October 2018).
- 1268 6. Friston, K. J. A theory of cortical responses.
1269 eng. *Philosophical Transactions of the Royal*
1270 *Society of London. Series B, Biological Sciences*
1271 **360**, 815–836. ISSN: 0962-8436 (April 2005).
- 1272 7. Kutas, M. & Hillyard, S. A. Brain potentials
1273 during reading reflect word expectancy and
1274 semantic association. *Nature* **307**. Place:
1275 United Kingdom Publisher: Nature Publishing
1276 Group, 161–163. ISSN:
1277 1476-4687(Electronic),0028-0836(Print) (1984).
- 1278 8. Hagoort, P., Brown, C. & Groothusen, J. The
1279 syntactic positive shift (SPS) as an ERP
1280 measure of syntactic processing. *Language and*
1281 *Cognitive Processes* **8**, 439–483. ISSN:
1282 1464-0732(Electronic),0169-0965(Print) (1993).
- 1283 9. Smith, N. J. & Levy, R. The effect of word
1284 predictability on reading time is logarithmic.
1285 en. *Cognition* **128**, 302–319. ISSN: 0010-0277
1286 (September 2013).
- 1287 10. Willems, R. M., Frank, S. L., Nijhof, A. D.,
1288 Hagoort, P. & van den Bosch, A. Prediction
1289 During Natural Language Comprehension. en.
1290 *Cerebral Cortex* **26**, 2506–2516. ISSN:
1291 1047-3211 (June 2016).
- 1292 11. Henderson, J. M., Choi, W., Lowder, M. W. &
1293 Ferreira, F. Language structure in the brain: A
1294 fixation-related fMRI study of syntactic
1295 surprisal in reading. eng. *NeuroImage* **132**,
1296 293–300. ISSN: 1095-9572 (2016).
- 1297 12. Armeni, K., Willems, R. M., van den Bosch, A. &
1298 Schoffelen, J.-M. Frequency-specific brain
1299 dynamics related to prediction during
1300 language comprehension. *NeuroImage*. ISSN:

- 1301 1053-8119.
1302 doi:[10.1016/j.neuroimage.2019.04.083](https://doi.org/10.1016/j.neuroimage.2019.04.083).
1303 (Visited on 05/20/2019) (May 2019).
- 1304 13. Donhauser, P. W. & Baillet, S. Two Distinct
1305 Neural Timescales for Predictive Speech
1306 Processing. en. *Neuron* **105**, 385–393.e9. ISSN:
1307 0896-6273 (January 2020).
- 1308 14. Ryskin, R., Levy, R. P. & Fedorenko, E. Do
1309 domain-general executive resources play a
1310 role in linguistic prediction? Re-evaluation of
1311 the evidence and a path forward. en.
1312 *Neuropsychologia* **136**, 107258. ISSN:
1313 0028-3932 (January 2020).
- 1314 15. Levy, R. Expectation-based syntactic
1315 comprehension. en. *Cognition* **106**, 1126–1177.
1316 ISSN: 0010-0277 (March 2008).
- 1317 16. Fitz, H. & Chang, F. Language ERPs reflect
1318 learning through prediction error propagation.
1319 en. *Cognitive Psychology* **111**, 15–52. ISSN:
1320 0010-0285 (June 2019).
- 1321 17. Huettig, F. & Mani, N. Is prediction necessary
1322 to understand language? Probably not.
1323 *Language, Cognition and Neuroscience* **31**,
1324 19–31 (2016).
- 1325 18. Brown, C. & Hagoort, P. The Processing Nature
1326 of the N400: Evidence from Masked Priming.
1327 *Journal of Cognitive Neuroscience* **5**. Publisher:
1328 MIT Press, 34–44. ISSN: 0898-929X (January
1329 1993).
- 1330 19. Nieuwland, M. S. Do ‘early’ brain responses
1331 reveal word form prediction during language
1332 comprehension? A critical review. eng.
1333 *Neuroscience and Biobehavioral Reviews* **96**,
1334 367–400. ISSN: 1873-7528 (2019).
- 1335 20. Hale, J. *A Probabilistic Earley Parser as a*
1336 *Psycholinguistic Model in Second Meeting of the*
1337 *North American Chapter of the Association for*
1338 *Computational Linguistics* (2001). <<https://www.aclweb.org/anthology/N01-1021>>
1339 (visited on 08/21/2020).
1340
- 1341 21. Brennan, J. R., Dyer, C., Kuncoro, A. & Hale, J. T.
1342 Localizing syntactic predictions using
1343 recurrent neural network grammars. en.
1344 *Neuropsychologia*, 107479. ISSN: 0028-3932
1345 (May 2020).
- 1346 22. Hale, J., Dyer, C., Kuncoro, A. & Brennan, J.
1347 *Finding syntax in human encephalography with*
1348 *beam search in Proceedings of the 56th Annual*
1349 *Meeting of the Association for Computational*
1350 *Linguistics (Volume 1: Long Papers)* (Association
1351 for Computational Linguistics, Melbourne,
1352 Australia, July 2018), 2727–2736.
1353 doi:[10.18653/v1/P18-1254](https://doi.org/10.18653/v1/P18-1254). <<https://www.aclweb.org/anthology/P18-1254>>
1354 (visited on 07/28/2020).
1355
- 1356 23. Fleur, D. S., Flecken, M., Rommers, J. &
1357 Nieuwland, M. S. Definitely saw it coming? The
1358 dual nature of the pre-nominal prediction
1359 effect. en. *Cognition* **204**, 104335. ISSN:
1360 0010-0277 (November 2020).
- 1361 24. Rabovsky, M., Hansen, S. S. & McClelland, J. L.
1362 Modelling the N400 brain potential as change
1363 in a probabilistic representation of meaning.
1364 En. *Nature Human Behaviour* **2**, 693. ISSN:
1365 2397-3374 (September 2018).
- 1366 25. Federmeier, K. D. Thinking ahead: the role and
1367 roots of prediction in language
1368 comprehension. eng. *Psychophysiology* **44**,
1369 491–505. ISSN: 0048-5772 (July 2007).
- 1370 26. Gagnepain, P., Henson, R. N. & Davis, M. H.
1371 Temporal Predictive Codes for Spoken Words
1372 in Auditory Cortex. English. *Current Biology* **22**.
1373 Publisher: Elsevier, 615–621. ISSN: 0960-9822
1374 (April 2012).
- 1375 27. Brodbeck, C., Hong, L. E. & Simon, J. Z. Rapid
1376 Transformation from Auditory to Linguistic
1377 Representations of Continuous Speech. eng.
1378 *Current biology: CB* **28**, 3976–3983.e5. ISSN:
1379 1879-0445 (December 2018).
- 1380 28. Di Liberto, G. M., Wong, D., Melnik, G. A. &
1381 de Cheveigne, A. Low-frequency cortical
1382 responses to natural speech reflect

- 1383 probabilistic phonotactics. en. *NeuroImage*
1384 **196**, 237–247. ISSN: 1053-8119 (August 2019).
- 1385 29. Gwilliams, L., Poeppel, D., Marantz, A. &
1386 Linzen, T. *Phonological (un)certainty weights*
1387 *lexical activation in Proceedings of the 8th*
1388 *Workshop on Cognitive Modeling and*
1389 *Computational Linguistics (CMCL 2018)*
1390 (Association for Computational Linguistics,
1391 Salt Lake City, Utah, January 2018), 29–34.
1392 doi:[10.18653/v1/W18-0104](https://doi.org/10.18653/v1/W18-0104) <[https:](https://www.aclweb.org/anthology/W18-0104)
1393 [//www.aclweb.org/anthology/W18-0104](https://www.aclweb.org/anthology/W18-0104)>
1394 (visited on 10/28/2020).
- 1395 30. Nieuwland, M. S. *et al.* Large-scale replication
1396 study reveals a limit on probabilistic prediction
1397 in language comprehension. *eLife* **7** (ed
1398 Shinn-Cunningham, B. G.) e33468. ISSN:
1399 2050-084X (April 2018).
- 1400 31. Broderick, M. P., Anderson, A. J.,
1401 Di Liberto, G. M., Crosse, M. J. & Lalor, E. C.
1402 Electrophysiological correlates of semantic
1403 dissimilarity reflect the comprehension of
1404 natural, narrative speech. *Current Biology* **28**,
1405 803–809 (2018).
- 1406 32. Frank, S. L., Otten, L. J., Galli, G. & Vigliocco, G.
1407 The ERP response to the amount of
1408 information conveyed by words in sentences.
1409 en. *Brain and Language* **140**, 1–11. ISSN:
1410 0093-934X (January 2015).
- 1411 33. Van Petten, C. & Luka, B. J. Prediction during
1412 language comprehension: Benefits, costs, and
1413 ERP components. en. *International Journal of*
1414 *Psychophysiology. Predictive information*
1415 *processing in the brain: Principles, neural*
1416 *mechanisms and models* **83**, 176–190. ISSN:
1417 0167-8760 (February 2012).
- 1418 34. Radford, A. *et al.* Language models are
1419 unsupervised multitask learners. *OpenAI Blog*
1420 **1**, 8 (2019).
- 1421 35. Smith, N. J. & Kutas, M. Regression-based
1422 estimation of ERP waveforms: I. The rERP
1423 framework. *Psychophysiology* **52**, 157–168.
1424 ISSN: 0048-5772 (February 2015).
- 1425 36. Ding, N. & Simon, J. Z. Emergence of neural
1426 encoding of auditory objects while listening to
1427 competing speakers. en. *Proceedings of the*
1428 *National Academy of Sciences* **109**. Publisher:
1429 National Academy of Sciences Section:
1430 Biological Sciences, 11854–11859. ISSN:
1431 0027-8424, 1091-6490 (July 2012).
- 1432 37. Abrams, D. A., Nicol, T., Zecker, S. & Kraus, N.
1433 Right-Hemisphere Auditory Cortex Is
1434 Dominant for Coding Syllable Patterns in
1435 Speech. *The Journal of Neuroscience* **28**,
1436 3958–3965. ISSN: 0270-6474 (April 2008).
- 1437 38. Binder, J. R., Desai, R. H., Graves, W. W. &
1438 Conant, L. L. Where Is the Semantic System? A
1439 Critical Review and Meta-Analysis of 120
1440 Functional Neuroimaging Studies. en. *Cerebral*
1441 *Cortex* **19**. Publisher: Oxford Academic,
1442 2767–2796. ISSN: 1047-3211 (December 2009).
- 1443 39. Huth, A. G., de Heer, W. A., Griffiths, T. L.,
1444 Theunissen, F. E. & Gallant, J. L. Natural speech
1445 reveals the semantic maps that tile human
1446 cerebral cortex. en. *Nature* **532**. Number: 7600
1447 Publisher: Nature Publishing Group, 453–458.
1448 ISSN: 1476-4687 (April 2016).
- 1449 40. Matchin, W. & Hickok, G. The Cortical
1450 Organization of Syntax. eng. *Cerebral Cortex*
1451 *(New York, N.Y.: 1991)* **30**, 1481–1498. ISSN:
1452 1460-2199 (2020).
- 1453 41. Matchin, W., Brodbeck, C., Hammerly, C. &
1454 Lau, E. The temporal dynamics of structure
1455 and content in sentence comprehension:
1456 Evidence from fMRI-constrained MEG. en.
1457 *Human Brain Mapping* **40**. eprint:
1458 <https://onlinelibrary.wiley.com/doi/pdf/10.1002/hbm.24403>,
1459 663–678. ISSN: 1097-0193 (2019).
- 1460 42. Nelson, M. J. *et al.* Neurophysiological
1461 dynamics of phrase-structure building during
1462 sentence processing. en. *Proceedings of the*
1463 *National Academy of Sciences* **114**. Publisher:
1464 National Academy of Sciences Section: PNAS
1465 Plus, E3669–E3678. ISSN: 0027-8424,
1466 1091-6490 (May 2017).

- 1467 43. Lopopolo, A., Frank, S. L., Bosch, A. v. d. & Willems, R. M. Using stochastic language models (SLM) to map lexical, syntactic, and phonological information processing in the brain. en. *PLOS ONE* **12**. Publisher: Public Library of Science, e0177794. ISSN: 1932-6203 (2017).
- 1468
1469
1470
1471
1472
1473
- 1474 44. Brennan, J. R. & Hale, J. T. Hierarchical structure guides rapid linguistic predictions during naturalistic listening. en. *PLOS ONE* **14**. Publisher: Public Library of Science, e0207741. ISSN: 1932-6203 (January 2019).
- 1475
1476
1477
1478
- 1479 45. Van Herten, M., Kolk, H. H. J. & Chwilla, D. J. An ERP study of P600 effects elicited by semantic anomalies. eng. *Brain Research. Cognitive Brain Research* **22**, 241–255. ISSN: 0926-6410 (February 2005).
- 1480
1481
1482
1483
- 1484 46. Gwilliams, L., King, J.-R., Marantz, A. & Poeppel, D. *Neural dynamics of phoneme sequencing in real speech jointly encode order and invariant content* en. preprint (Neuroscience, April 2020). doi:[10.1101/2020.04.04.025684](https://doi.org/10.1101/2020.04.04.025684). <<http://biorxiv.org/lookup/doi/10.1101/2020.04.04.025684>> (visited on 09/20/2020).
- 1485
1486
1487
1488
1489
1490
1491
1492
- 1493 47. Kiebel, S. J., Daunizeau, J. & Friston, K. J. A Hierarchy of Time-Scales and the Brain. en. *PLOS Computational Biology* **4**. Publisher: Public Library of Science, e1000209. ISSN: 1553-7358 (November 2008).
- 1494
1495
1496
1497
- 1498 48. Marslen-Wilson, W. in *Lexical representation and process* 3–24 (The MIT Press, Cambridge, MA, US, 1989). ISBN: 978-0-262-13240-4.
- 1499
1500
- 1501 49. Rao, R. P. N. & Ballard, D. H. Predictive coding in the visual cortex: a functional interpretation of some extra-classical receptive-field effects. en. *Nature Neuroscience* **2**. Number: 1 Publisher: Nature Publishing Group, 79–87. ISSN: 1546-1726 (January 1999).
- 1502
1503
1504
1505
1506
- 1507 50. Heilbron, M. & Chait, M. Great expectations: Is there evidence for predictive coding in auditory cortex? *Neuroscience*. ISSN: 0306-4522. doi:[10.1016/j.neuroscience.2017.07.061](https://doi.org/10.1016/j.neuroscience.2017.07.061) (August 2017).
- 1508
1509
1510
1511
- 1512 51. Mantegna, F., Hintz, F., Ostarek, M., Alday, P. M. & Huettig, F. Distinguishing integration and prediction accounts of ERP N400 modulations in language processing through experimental design. en. *Neuropsychologia* **134**, 107199. ISSN: 0028-3932 (November 2019).
- 1513
1514
1515
1516
1517
- 1518 52. Nieuwland, M. S. *et al.* Dissociable effects of prediction and integration during language comprehension: evidence from a large-scale study using brain potentials. *Philosophical Transactions of the Royal Society B: Biological Sciences* **375**. Publisher: Royal Society, 20180522 (February 2020).
- 1519
1520
1521
1522
1523
1524
- 1525 53. Van den Brink, D., Brown, C. M. & Hagoort, P. Electrophysiological Evidence for Early Contextual Influences during Spoken-Word Recognition: N200 Versus N400 Effects. *Journal of Cognitive Neuroscience* **13**. Publisher: MIT Press, 967–985. ISSN: 0898-929X (October 2001).
- 1526
1527
1528
1529
1530
1531
- 1532 54. Friederici, A. D. Towards a neural basis of auditory sentence processing. en. *Trends in Cognitive Sciences* **6**, 78–84. ISSN: 1364-6613 (February 2002).
- 1533
1534
1535
- 1536 55. Weissbart, H., Kandylaki, K. D. & Reichenbach, T. Cortical Tracking of Surprisal during Continuous Speech Comprehension. eng. *Journal of Cognitive Neuroscience* **32**, 155–166. ISSN: 1530-8898 (2020).
- 1537
1538
1539
1540
- 1541 56. Heilbron, M., Ehinger, B., Hagoort, P. & de Lange, F. P. Tracking Naturalistic Linguistic Predictions with Deep Neural Language Models. *2019 Conference on Cognitive Computational Neuroscience*. arXiv: 1909.04400. doi:[10.32470/CCN.2019.1096-0](https://doi.org/10.32470/CCN.2019.1096-0). (Visited on 09/25/2019) (2019).
- 1542
1543
1544
1545
1546
1547

- 1548 57. Shain, C., Blank, I. A., van Schijndel, M.,
1549 Schuler, W. & Fedorenko, E. fMRI reveals
1550 language-specific predictive coding during
1551 naturalistic sentence comprehension. en.
1552 *Neuropsychologia* **138**, 107307. ISSN:
1553 0028-3932 (February 2020).
- 1554 58. Goldstein, A. *et al.* Thinking ahead: prediction
1555 in context as a keystone of language in
1556 humans and machines. en. *bioRxiv*. Publisher:
1557 Cold Spring Harbor Laboratory Section: New
1558 Results, 2020.12.02.403477 (December 2020).
- 1559 59. Ali, A., Ahmad, N., Groot, E. d.,
1560 Gerven, M. A. J. v. & Kietzmann, T. C. Predictive
1561 coding is a consequence of energy efficiency
1562 in recurrent neural networks. en. *bioRxiv*.
1563 Publisher: Cold Spring Harbor Laboratory
1564 Section: New Results, 2021.02.16.430904
1565 (February 2021).
- 1566 60. Heilbron, M., Richter, D., Ekman, M., Hagoort, P.
1567 & de Lange, F. P. Word contexts enhance the
1568 neural representation of individual letters in
1569 early visual cortex. en. *Nature Communications*
1570 **11**, 321. ISSN: 2041-1723 (January 2020).
- 1571 61. Broderick, M. P., Anderson, A. J. & Lalor, E. C.
1572 Semantic Context Enhances the Early Auditory
1573 Encoding of Natural Speech. en. *Journal of*
1574 *Neuroscience*, 0584–19. ISSN: 0270-6474,
1575 1529-2401 (August 2019).
- 1576 62. Sohoglu, E. & Davis, M. H. Rapid computations
1577 of spectrotemporal prediction error support
1578 perception of degraded speech. *eLife* **9** (eds
1579 King, A. J., Kok, P., Kok, P., Press, C. &
1580 Lalor, E. C.) Publisher: eLife Sciences
1581 Publications, Ltd, e58077. ISSN: 2050-084X
1582 (November 2020).
- 1583 63. Blank, H. & Davis, M. H. Prediction Errors but
1584 Not Sharpened Signals Simulate Multivoxel
1585 fMRI Patterns during Speech Perception. eng.
1586 *PLoS biology* **14**, e1002577. ISSN: 1545-7885
1587 (November 2016).
- 1588 64. McClelland, J. L., Hill, F., Rudolph, M.,
1589 Baldridge, J. & Schutze, H. Placing language in
1590 an integrated understanding system: Next
1591 steps toward human-level performance in
1592 neural language models. en. *Proceedings of the*
1593 *National Academy of Sciences*. Publisher:
1594 National Academy of Sciences Section:
1595 Perspective. ISSN: 0027-8424, 1091-6490.
1596 doi:[10.1073/pnas.1910416117](https://doi.org/10.1073/pnas.1910416117).
1597 <[https://www.pnas.org/content/early/
1598 2020/09/25/1910416117](https://www.pnas.org/content/early/2020/09/25/1910416117)> (visited on
1599 10/07/2020) (September 2020).
- 1600 65. Manning, C. D., Clark, K., Hewitt, J.,
1601 Khandelwal, U. & Levy, O. Emergent linguistic
1602 structure in artificial neural networks trained
1603 by self-supervision. en. *Proceedings of the*
1604 *National Academy of Sciences*. Publisher:
1605 National Academy of Sciences Section:
1606 Physical Sciences. ISSN: 0027-8424, 1091-6490.
1607 doi:[10.1073/pnas.1907367117](https://doi.org/10.1073/pnas.1907367117).
1608 <[https://www.pnas.org/content/early/
1609 2020/06/02/1907367117](https://www.pnas.org/content/early/2020/06/02/1907367117)> (visited on
1610 11/03/2020) (June 2020).
- 1611 66. Caucheteux, C. & King, J.-R. Language
1612 processing in brains and deep neural
1613 networks: computational convergence and its
1614 limits. en. *bioRxiv*, 2020.07.03.186288 (July
1615 2020).
- 1616 67. Schrimpf, M. *et al.* The neural architecture of
1617 language: Integrative reverse-engineering
1618 converges on a model for predictive
1619 processing. en. *bioRxiv*. Publisher: Cold Spring
1620 Harbor Laboratory Section: New Results,
1621 2020.06.26.174482 (October 2020).
- 1622 68. Toneva, M. & Wehbe, L. Interpreting and
1623 improving natural-language processing (in
1624 machines) with natural language-processing
1625 (in the brain). en. *Advances in Neural*
1626 *Information Processing Systems* **32**,
1627 14954–14964 (2019).
- 1628 69. Jain, S. & Huth, A. G. Incorporating Context
1629 into Language Encoding Models for fMRI. en.
1630 *bioRxiv*. Publisher: Cold Spring Harbor

- 1631 Laboratory Section: New Results, 327601
1632 (November 2018).
- 1633 70. Whittington, J. C. R. & Bogacz, R. An
1634 Approximation of the Error Backpropagation
1635 Algorithm in a Predictive Coding Network with
1636 Local Hebbian Synaptic Plasticity. eng. *Neural
1637 Computation* **29**, 1229–1262. ISSN: 1530-888X
1638 (May 2017).
- 1639 71. Millidge, B., Tschantz, A. & Buckley, C. L.
1640 Predictive Coding Approximates Backprop
1641 along Arbitrary Computation Graphs.
1642 *arXiv:2006.04182 [cs]*. arXiv: 2006.04182.
1643 <http://arxiv.org/abs/2006.04182>
1644 (visited on 02/07/2021) (October 2020).
- 1645 72. Issa, E. B., Cadieu, C. F. & DiCarlo, J. J. Neural
1646 dynamics at successive stages of the ventral
1647 visual stream are consistent with hierarchical
1648 error signals. *eLife* **7** (eds Connor, E., Marder, E.
1649 & Connor, E.) Publisher: eLife Sciences
1650 Publications, Ltd, e42870. ISSN: 2050-084X
1651 (November 2018).
- 1652 73. Schwiedrzik, C. M. & Freiwald, W. A. High-Level
1653 Prediction Signals in a Low-Level Area of the
1654 Macaque Face-Processing Hierarchy. eng.
1655 *Neuron* **96**, 89–97.e4. ISSN: 1097-4199
1656 (September 2017).
- 1657 74. Wacongne, C. *et al.* Evidence for a hierarchy of
1658 predictions and prediction errors in human
1659 cortex. eng. *Proceedings of the National
1660 Academy of Sciences of the United States of
1661 America* **108**, 20754–20759. ISSN: 1091-6490
1662 (December 2011).
- 1663 75. Oostenveld, R., Fries, P., Maris, E. &
1664 Schoffelen, J.-M. FieldTrip: Open source
1665 software for advanced analysis of MEG, EEG,
1666 and invasive electrophysiological data. eng.
1667 *Computational Intelligence and Neuroscience*
1668 **2011**, 156869. ISSN: 1687-5273 (2011).
- 1669 76. Vaswani, A. *et al.* *Attention is all you need in
1670 Advances in neural information processing
1671 systems* (2017), 5998–6008.
- 1672 77. Wolf, T. *et al.* HuggingFace's Transformers:
1673 State-of-the-art Natural Language Processing.
1674 *arXiv:1910.03771 [cs]*.
1675 <http://arxiv.org/abs/1910.03771>
1676 (visited on 07/30/2020) (July 2020).
- 1677 78. Daube, C., Ince, R. A. A. & Gross, J. Simple
1678 Acoustic Features Can Explain
1679 Phoneme-Based Predictions of Cortical
1680 Responses to Speech. en. *Current Biology* **29**,
1681 1924–1937.e9. ISSN: 0960-9822 (June 2019).
- 1682 79. Gramfort, A. *et al.* MNE software for
1683 processing MEG and EEG data. en. *NeuroImage*
1684 **86**, 446–460. ISSN: 1053-8119 (February 2014).
- 1685 80. Maris, E. & Oostenveld, R. Nonparametric
1686 statistical testing of EEG- and MEG-data. en.
1687 *Journal of Neuroscience Methods* **164**, 177–190.
1688 ISSN: 0165-0270 (August 2007).

Supplementary materials

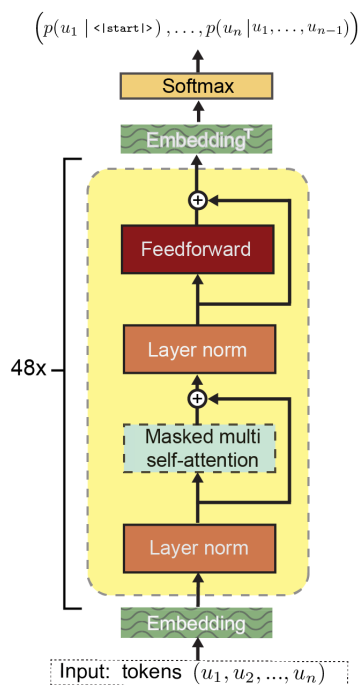


Figure S1 – GPT-2 ARCHITECTURE. Note that this panel is a re-rendered version of the original GPT schematic, slightly modified and re-arranged to match the architecture of GPT-2. For more details on the overall architecture and on the critical operation of self-attention, see *Methods*. In this graphic, Layer Norm refers to layer normalisation as described by Ba et al. Not visualised here is the initial tokenisation, mapping a sequence of characters into tokens.

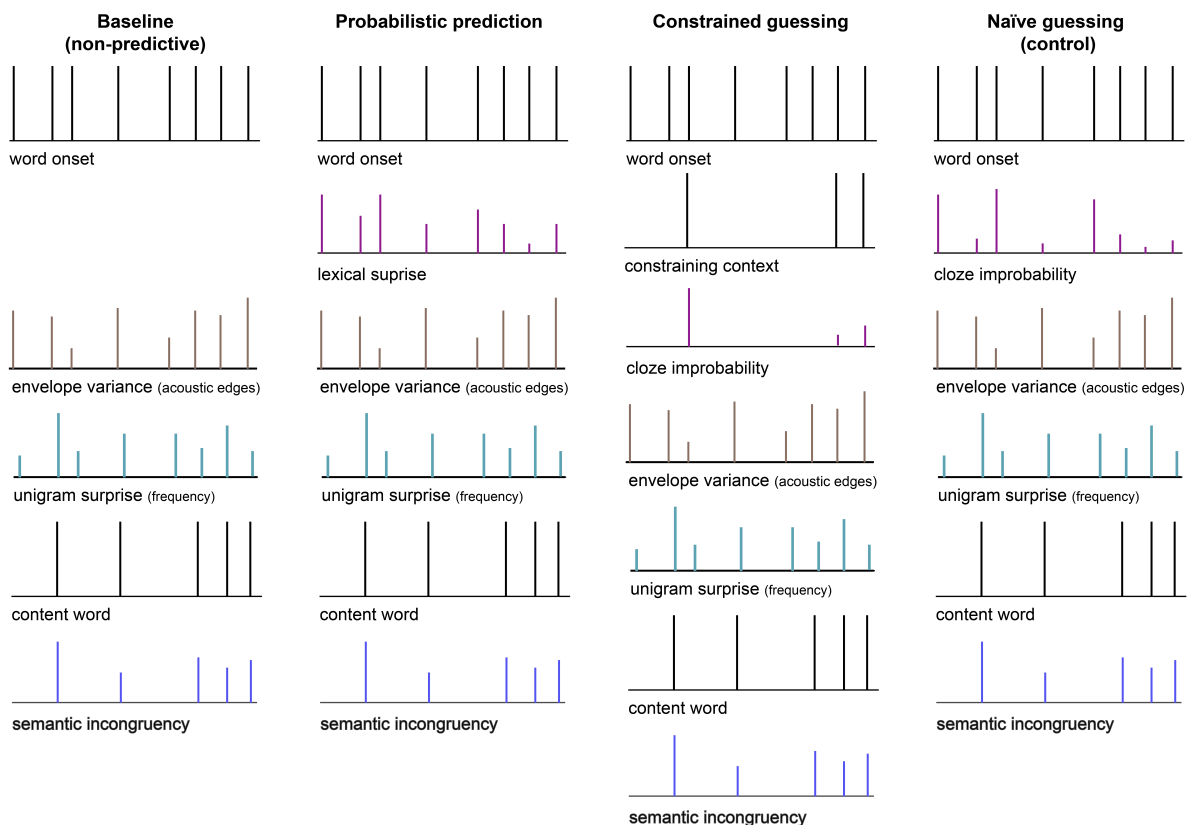


Figure S2 – WORD-LEVEL REGRESSION MODELS. Schematic of the main models plus the control model of the initial model comparison to test for predictive processing at the word level. Because we use a regression ERP/ERF scheme [35], aimed at capturing (modulations of) the evoked response to discrete events like words or phonemes, all regressors are modelled as impulses (see *Methods*).

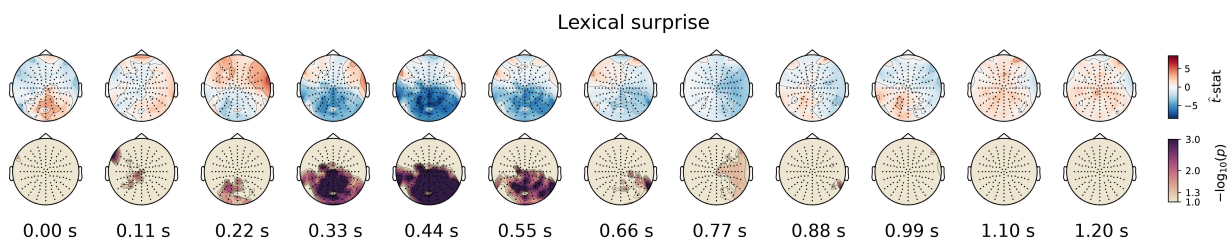


Figure S3 – FULL EEG TOPOGRAPHIES OF THE EFFECTS OF LEXICAL SURPRISE These topographies show the average t-statistics of the coefficients (upper row) and respective FWE-corrected significance (lower row) of the lexical surprise regressor from the *probabilistic prediction* model (Figure S2). As such, while Figure 2b shows the coefficients averaged over channels participating in the cluster (thereby only visualising *the effect*) these topographies visualise the results comprehensively across all channels over time.

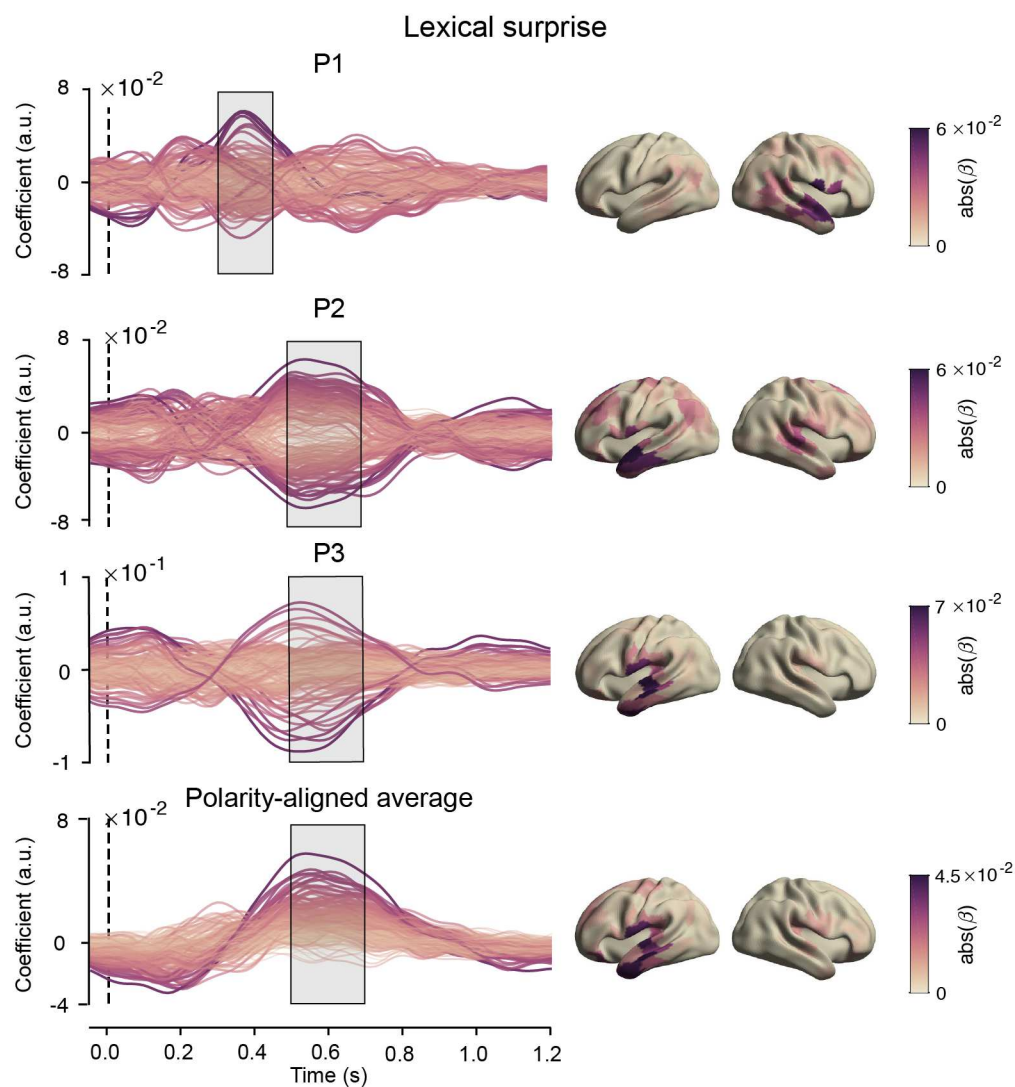


Figure S4 – COEFFICIENTS FOR LEXICAL SURPRISE FROM THE LEXICAL MODEL (FIGURE S2) Left column: timecourses of the coefficients at each MEG source-localised parcel for lexical surprise for all MEG participants, and the polarity-aligned average across them. Right column: Absolute value of the coefficients averaged across the highlighted period plotted across the brain.

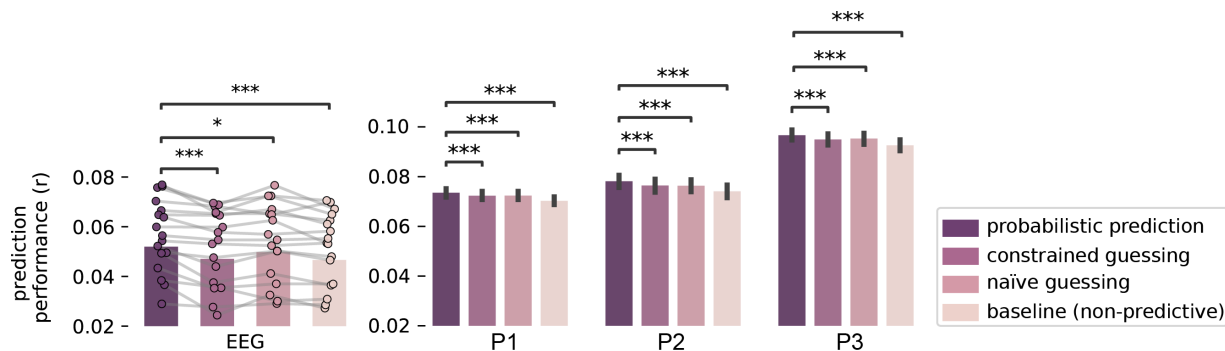


Figure S5 – MODEL COMPARISON RESULTS ACROSS ALL CHANNELS (EEG) AND THE FULL LANGUAGE NETWORK (MEG). Same as in Figure 2a, but now including the ‘naïve guessing’ control model. Like the constrained guessing model, this model included a linear estimate of word probability, but defined for every word rather than only for constraining contexts. This model was introduced to identify which of the two differences between the *probabilistic prediction* and *constrained guessing* model – i.e. assuming that predictions are (i) categorical vs. probabilistic and (ii) occasional vs. continuous – made the largest difference in model performance. As can be seen, the *naïve guessing* model performed considerably better than the *constrained guessing* model, but consistently worse than the *probabilistic prediction* model. This clearly shows that the modulatory effect of unexpectedness is not limited to only highly constraining contexts, but that that it applies much more generally – in line with the notion of continuous prediction.

Strictly speaking, the naïve guessing model formalises the hypothesis that the brain ‘naïvely’ makes *all-or-none* guesses about *every* upcoming word. Given that this hypothesis is a-priori so implausible, it may seem surprising that the model still performs comparably well. However, we should note that the probabilistic prediction regressor (*surprise*) and the categorical prediction regressor (linear (im)probability) are highly correlated (~ 0.7) because one is a monotonic function of the other. Therefore, we suggest the results are better interpreted the other way around: the fact that – despite being so correlated – the log-probability is consistently a better linear predictor of neural responses than the linear probability clearly supports predictive processing theories, which postulate that the neural response to a stimulus should be proportional to negative log-probability of that stimulus.

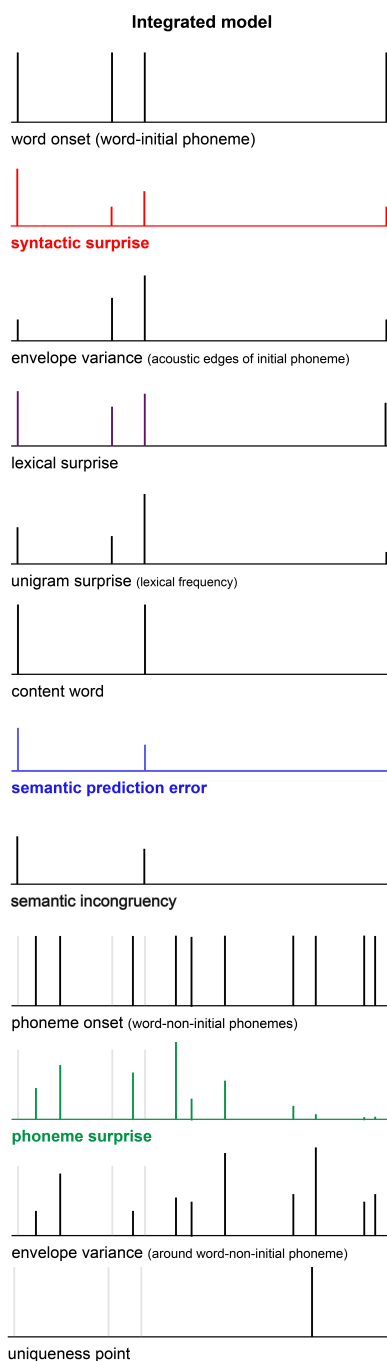


Figure S6 – REGRESSORS OF THE INTEGRATED FEATURE-SPECIFIC MODEL. Same as Figure [S5](#), but for the integrated feature-specific regression model. The three regressors of interest – syntactic surprise, semantic prediction error and phonemic surprise – are coloured, all control regressors are in black. Following the regression ERP/ERF scheme [35](#), aimed at capturing (modulations of) the evoked response to discrete events like words or phonemes, all regressors are modelled as impulses (see *Methods*). To avoid collinearity between word and phoneme regressors, phoneme regressors (both events and covariates) are restricted to all non-initial phonemes.

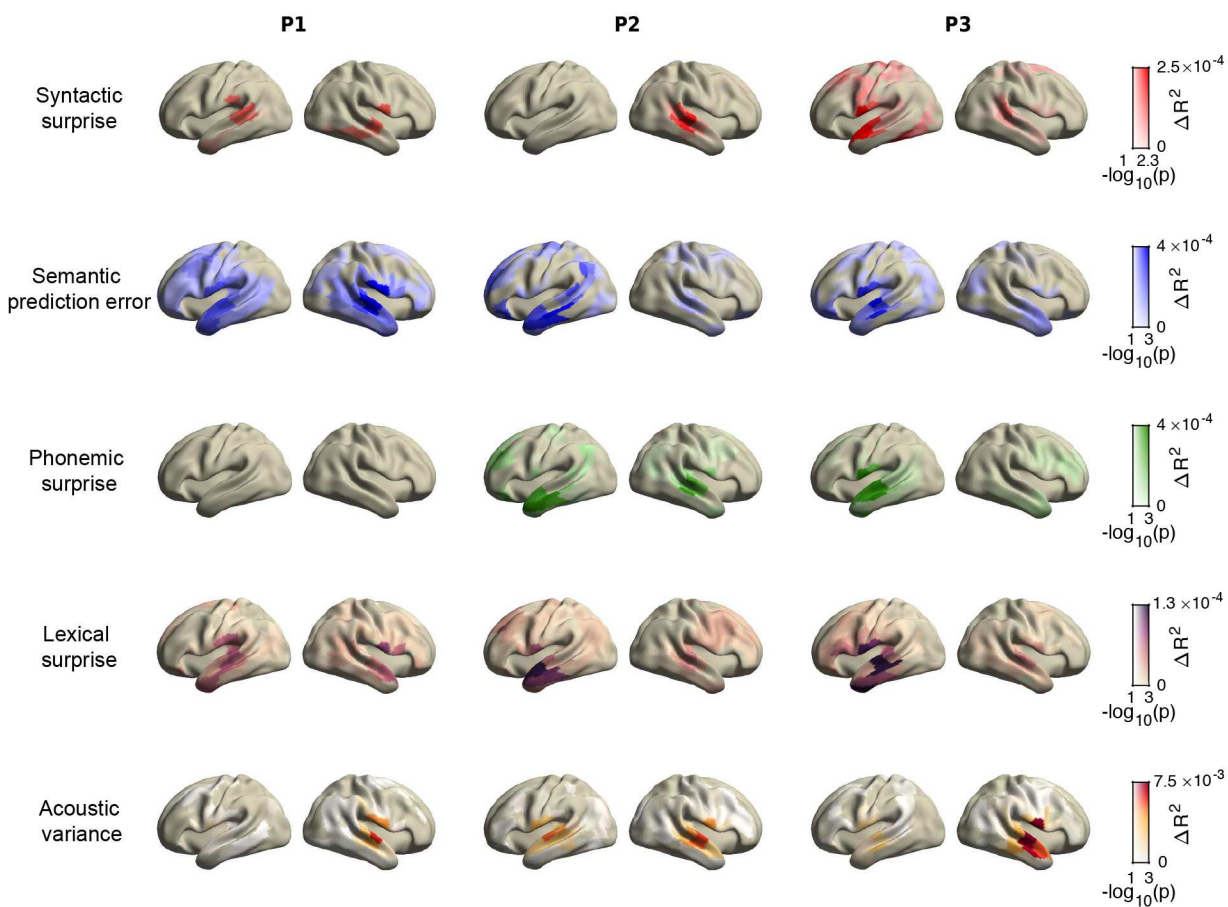


Figure S7 – UNIQUE EXPLAINED VARIANCE FOR FIVE REGRESSORS ACROSS THE BRAIN.

Same as Figure 4, but including 2 control regressors (lexical surprise and acoustic variance) for comparison. Colours indicate amount of additional variance explained by each regressor; opacity indicates the FWE-corrected statistical significance (across cross-validation folds). Note that $p < 0.05$ is equivalent to $-\log_{10}(p) > 1.3$.

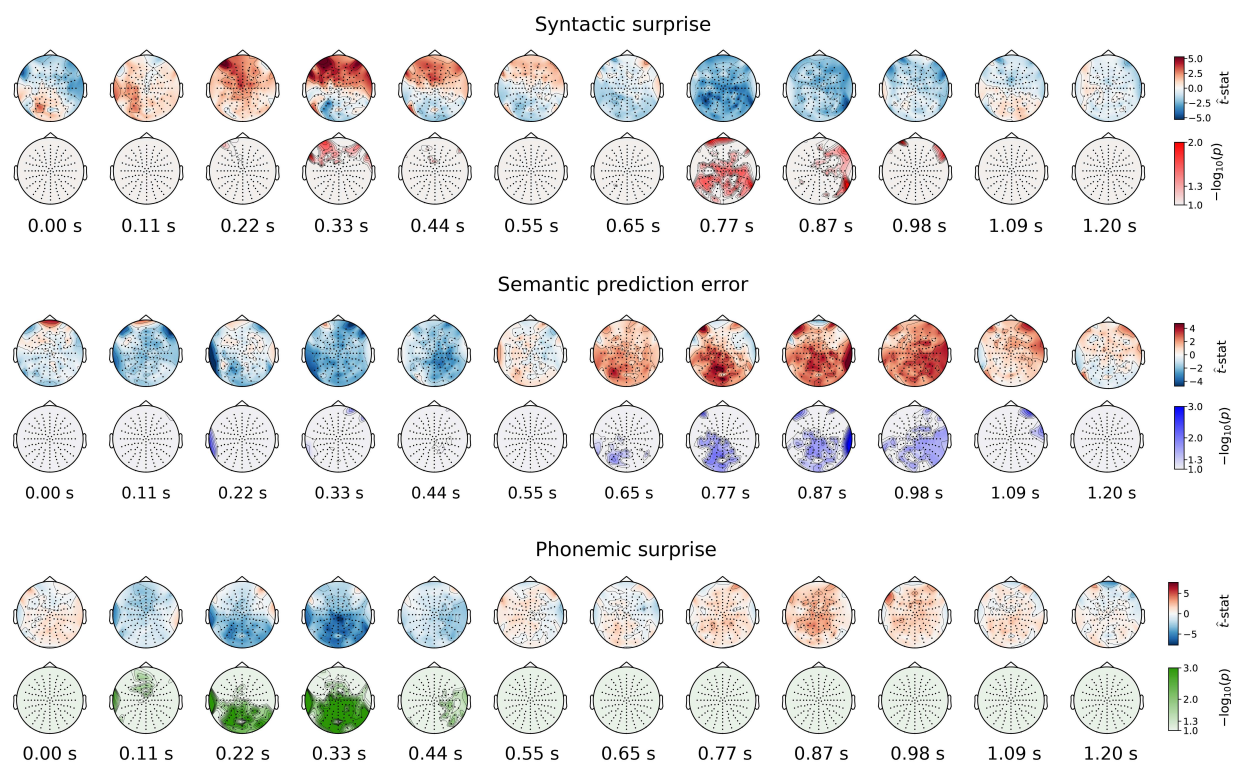


Figure S8 – FULL TOPOGRAPHIES OF THE COEFFICIENTS AND SIGNIFICANCE OF FEATURE-SPECIFIC PREDICTION ERRORS
For each feature-specific prediction error regressor, the topographies show the t-statistics of the coefficients (upper row) and the respective TFCE-corrected significance (lower row). So while Figure 5 only shows the coefficients averaged over channels participating in the cluster (thereby only visualising *the effect*) these topographies visualise the results comprehensively across all channels, over time.

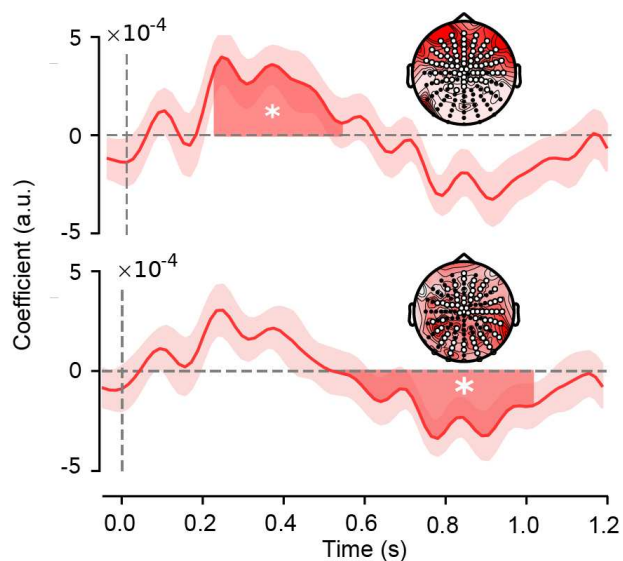


Figure S9 – SIGNIFICANT EFFECTS OF SYNTACTIC SURPRISE IN THE EEG DATA. Two significant effects were observed in the modulation functions for syntactic surprise: an early positive effect with a frontal topography (upper panel) and a later negative effect based on a distributed cluster (lower panel). The early effect tightly replicates recent model-based studies on EEG effects of syntactic surprise, and was also found in the MEG data. By contrast, the late effect of syntactic surprise is not in line with any earlier study (note that it is negative unlike the syntactic P600) and importantly was not replicated in the MEG data. Therefore we only consider the early effect a ‘main’ effect of syntactic surprise (visualised in the main Figure 5) and we advice to refrain from interpreting the late effect before it is independently replicated.

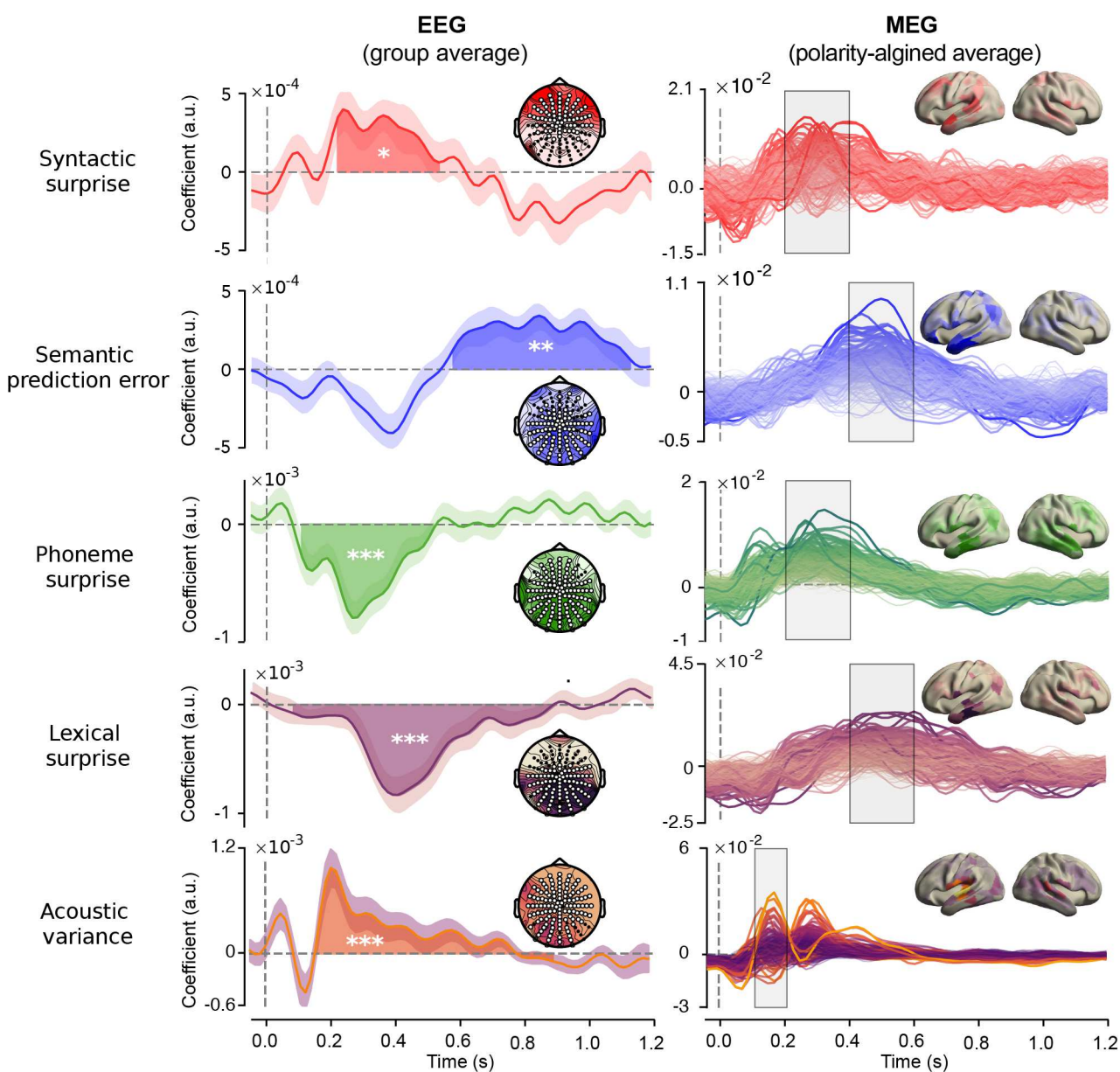


Figure S10 – COEFFICIENTS FOR EACH PREDICTION ERROR, PLUS TWO CONTROL VARIABLES.

EEG (left column): coefficient modulation function averaged across the channels participating for at least one sample in the significant clusters. Highlighted area indicates temporal extent of the cluster. Shaded area around waveform indicates bootstrapped standard errors. Stars indicate cluster-level significance; $p < 0.05$ (*), $p < 0.01$ (**), $p < 0.001$ (***). Insets represent channels assigned to the cluster (white dots) and the distribution of absolute values of t-statistics. MEG (right column): polarity aligned responses averaged across participants for all sources (same as in Figure 5) but without averaging over sources, and including two control variables). Insets represent topography of absolute value of coefficients averaged across the highlighted period. Note that due to polarity alignment, sign information is to be ignored for the MEG plots.

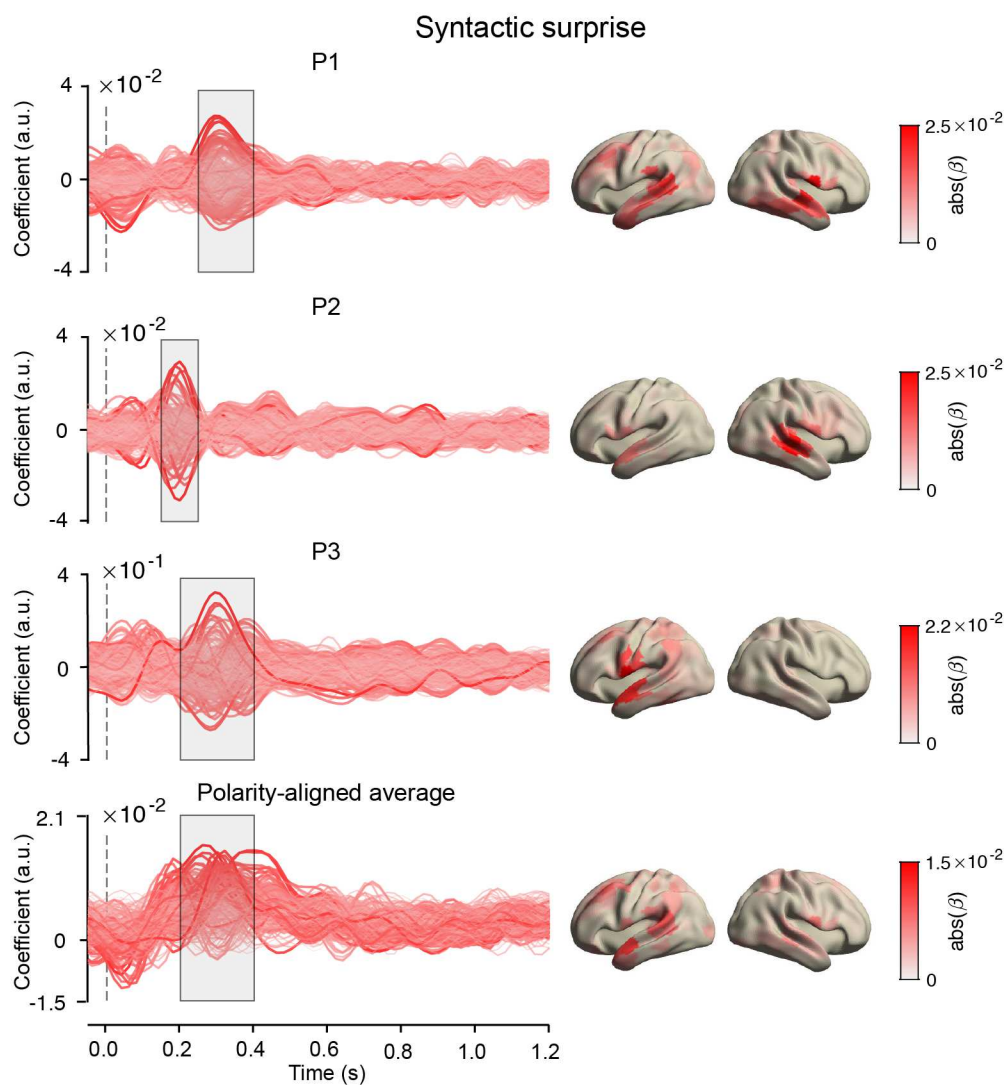


Figure S11 – COEFFICIENTS FOR SYNTACTIC SURPRISE FROM THE INTEGRATED MODEL (FIGURE S6)

Left column: coefficients for each source for each individual in the MEG experiment, and the polarity-aligned average across participants. Right column: absolute value of the coefficients across the brain, averaged across the highlighted time-period.

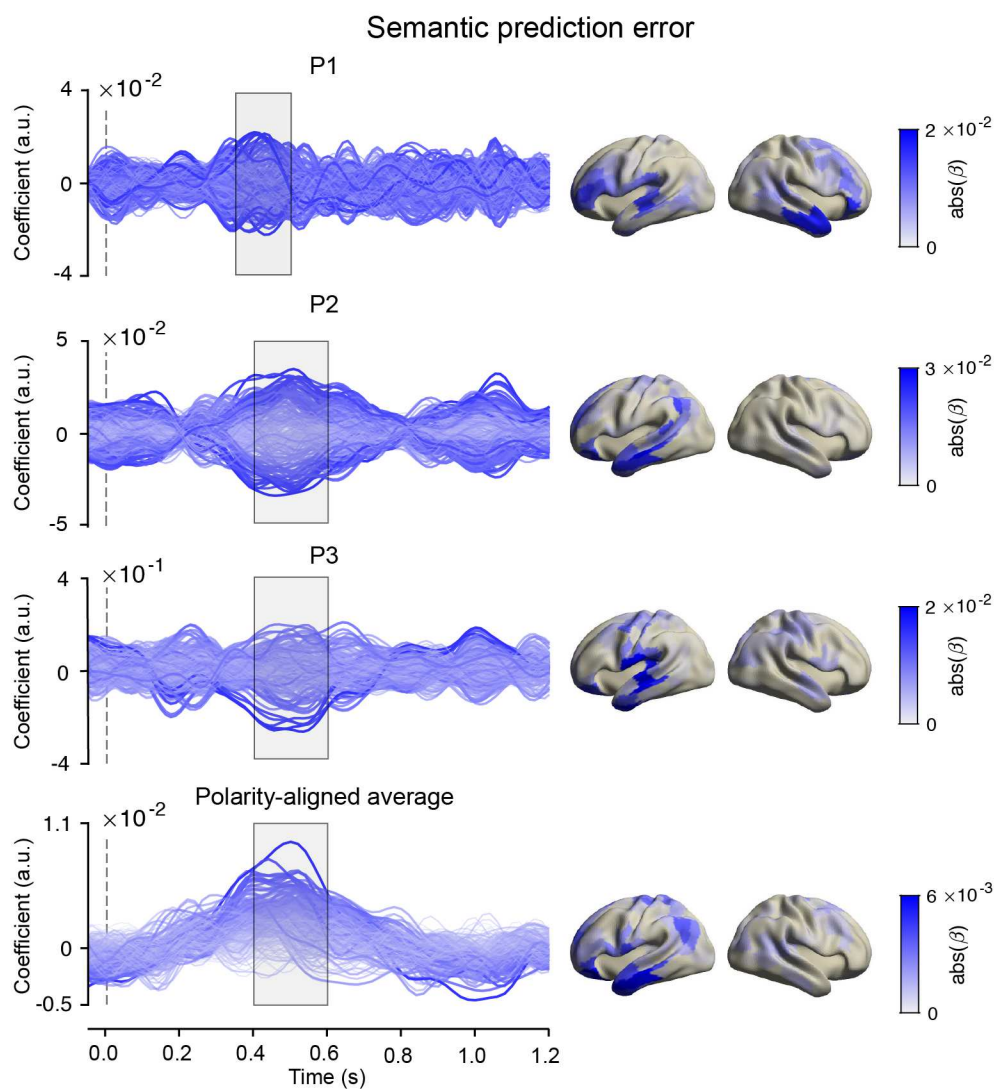


Figure S12 – COEFFICIENTS FOR SEMANTIC PREDICTION ERROR FROM THE INTEGRATED MODEL (FIGURE S6)

Left column: coefficients for each source for each individual in the MEG experiment, and the polarity-aligned average across participants. Right column: absolute value of the coefficients across the brain, averaged across the highlighted time-period.

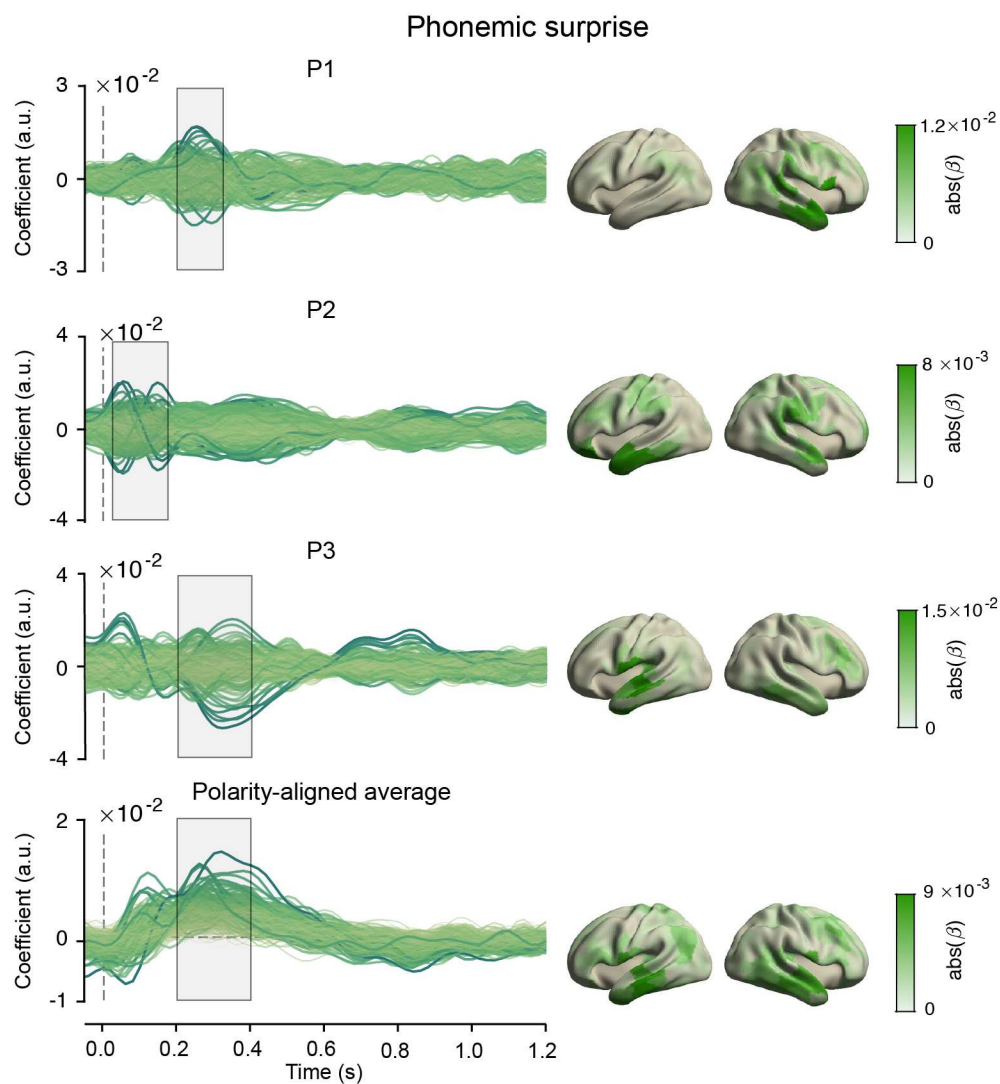


Figure S13 – COEFFICIENTS FOR PHONEMIC SURPRISE FROM THE INTEGRATED MODEL (FIGURE S6)

Left column: coefficients for each source for each individual in the MEG experiment, and the polarity-aligned average across participants. Right column: absolute value of the coefficients across the brain, averaged across the highlighted time-period.

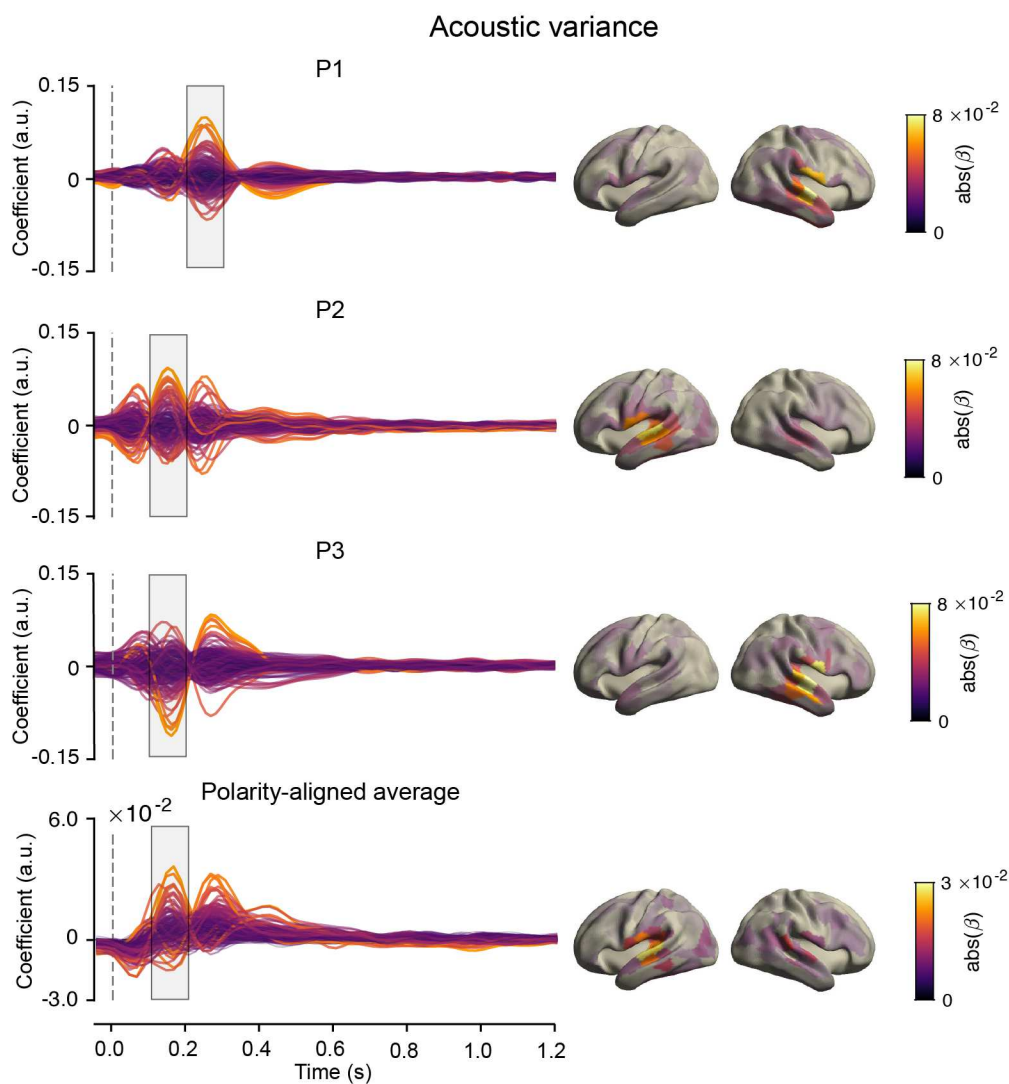


Figure S14 – COEFFICIENTS FOR ENVELOPE VARIABILITY FROM THE INTEGRATED MODEL (FIGURE S6)

Left column: coefficients for each source for each individual in the MEG experiment, and the polarity-aligned average across participants. Right column: absolute value of the coefficients across the brain, averaged across the highlighted time-period.

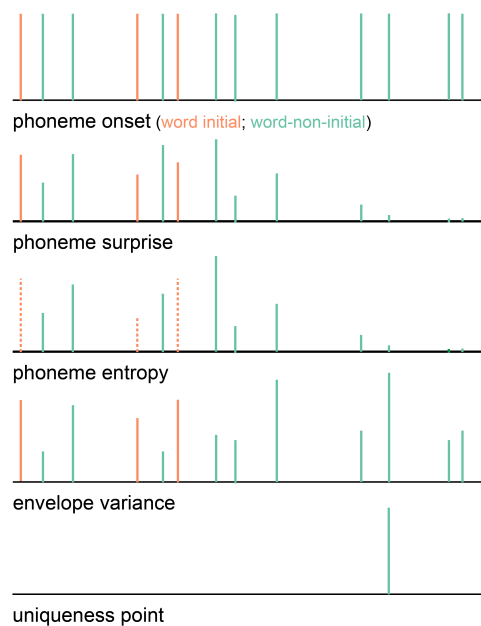


Figure S15 – REGRESSORS OF THE PHONEME MODEL. As indicated by the different colours, both the constants and covariates were modelled separately for word-initial and word-non-initial phonemes.

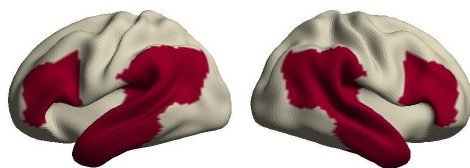


Figure S16 – LANGUAGE NETWORK DEFINITION The language network was defined as temporal cortex plus temporo-parietal junction, and IFG and dorsolateral prefrontal cortex; all bilaterally. In terms of Brodmann areas this corresponded to 20, 21, 22, 38, 39, 40, 41, 42, 44, 45, 46 and 47, amounting to a total of 100 out of 370 cortical parcels.