

1 **Predicting speech from a cortical hierarchy of**
2 **event-based timescales**

3 Lea-Maria Schmitt^{1,2}, Julia Erb^{1,2}, Sarah Tune^{1,2}, Anna Rysop³,
4 Gesa Hartwigsen³, Jonas Obleser^{1,2}

5 ¹ Department of Psychology, University of Lübeck,
6 Ratzeburger Allee 160, 23562 Lübeck, Germany

7 ² Center of Brain, Behavior and Metabolism, University of Lübeck,
8 Ratzeburger Allee 160, 23562 Lübeck, Germany

9 ³ Lise Meitner Research Group Cognition and Plasticity,
10 Max Planck Institute for Human Cognitive and Brain Sciences, Stephanstraße 1 A, 04103 Leipzig,
11 Germany

12 **Author correspondence:**

13 Lea-Maria Schmitt & Jonas Obleser

14 Department of Psychology

15 University of Lübeck

16 Maria-Goeppert-Straße 9 A

17 23562 Lübeck, Germany

18 *l.schmitt@uni-luebeck.de; jonas.obleser@uni-luebeck.de*

19 **Abstract**

20 How can anticipatory neural processes structure the temporal unfolding of context in our natural
21 environment? We here provide evidence for a neural coding scheme that sparsely updates
22 contextual representations at the boundary of events and gives rise to a hierarchical, multi-layered
23 organization of predictive language comprehension. Training artificial neural networks to predict
24 the next word in a story at five stacked timescales and then using model-based functional MRI, we
25 observe a sparse, event-based “surprisal hierarchy”. The hierarchy evolved along a temporo-parietal
26 pathway, with model-based surprisal at longest timescales represented in inferior parietal regions.
27 Along this hierarchy, surprisal at any given timescale gated bottom-up and top-down connectivity
28 to neighbouring timescales. In contrast, surprisal derived from a continuously updated context
29 influenced temporo-parietal activity only at short timescales. Representing context in the form of
30 increasingly coarse events constitutes a network architecture for making predictions that is both
31 computationally efficient and semantically rich.

32 **Keywords:** timescale hierarchy, speech prediction, natural language processing, surprisal, model-
33 based fMRI, artificial neural networks

34 Introduction

35 While the past predicts the future, not all context the past provides is equally informative: it might
36 be outdated, contradictory, or even irrelevant. Nevertheless, the brain as a “prediction machine”¹ is
37 seemingly equipped with a versatile repertoire of computations to overcome these contextual
38 ambiguities. A prominent example is speech, where a slip of the tongue may render the most recent
39 context uninformative, but we can still predict the next word from its remaining context. At much
40 longer time scales, we can re-use context that suddenly proves informative, as a speaker returns to
41 a topic discussed earlier.

42 Using natural language comprehension as a working model, we here ask: How does the brain
43 dynamically organize, evaluate, and update these complex contextual dependencies over time to
44 make accurate predictions?

45 A robust principle in cerebral cortex is the decomposition of temporal context into its
46 constituent timescales along a hierarchy from lower to higher-order areas, which is evident across
47 species^{2,3}, recording modalities^{4,5}, sensory modalities^{6,7}, and cognitive functions^{8,9}. For instance,
48 sensory cortices closely track rapid fluctuations of stimulus features and operate on short timescales
49 (e.g., Ref. ¹⁰). By contrast, association cortices integrate stimuli over an extended period and operate
50 on longer timescales (e.g., Ref. ¹¹).

51 Conceptually, such hierarchies of “temporal receptive windows” are often subsumed under
52 the framework of predictive coding¹²: A nested set of timescale-specific generative models informs
53 predictions on upcoming sensory input and is updated based on the actual input¹³. In particular,
54 context is thought to shape the prediction of incoming stimuli via feedback connections. These
55 connections would link each timescale to its immediate shorter timescale, while the prediction error
56 is propagated forward through the hierarchy^{1,14}. Indeed, hierarchical specialization has been shown
57 empirically to emerge from structural and functional large-scale connectivity across cortex^{15,16}. More
58 precisely, feedforward and feedback connections^{17,18} shown to carry prediction errors and
59 predictions^{19,20}, respectively, are a hallmark of hierarchical predictive coding.

60 However, studies on the neural underpinnings of predictive coding have primarily used
61 artificial stimuli of short temporal context (but see Ref. ²¹) and employed local vs. global violations of
62 expectations, effectively manifesting a two-level cortical hierarchy (but see Ref. ²²). We thus lack
63 understanding whether the hierarchical organization of prediction processes extends to natural
64 environments unfolding their temporal dependencies over a multitude of interrelated timescales.

65 With respect to functional organization in cortex, temporo-parietal areas are sensitive to a
66 rich set of hierarchies and timescales in speech²³⁻²⁷. Most relevant to the present work, semantic
67 context in a spoken story has been shown to map onto a gradient extending from early auditory
68 cortex representative of words up to intraparietal sulcus, representative of paragraphs²⁸. This

69 timescale-specific representation of context is reminiscent of the multi-layered generative models
70 proposed to underlie predictive coding^{29,30}. Compatible with this notion, previous studies on speech
71 comprehension found evidence for neural coding of prediction errors at the level of syllables³¹,
72 words³², or discourse³³.

73 Yet the interactions between multiple representational levels of speech in predicting
74 upcoming words remain unclear. Here, we ask whether the processing hierarchy enabling natural
75 speech comprehension is also implicated in evaluating the predictiveness of timescale-specific
76 semantic context and integrating informative context into predictions.

77 We do not know how context that unfolds at a particular timescale would be updated
78 cortically when the listener receives new bottom-up input. One attractively simple candidate
79 architecture is the *continuously updating processing hierarchy*. In a recent study, Chien and Honey³⁴
80 showed that neural responses to a story rapidly aligned across participants in areas with shorter, but
81 only later in areas with longer receptive windows. This response pattern was best explained by a
82 computational model which immediately integrates upcoming input with context representations
83 at all timescales. An important implication of such continuous updates is that all context
84 representations are continuously tuned to current processing demands.

85 A competing candidate architecture, however, is the *sparsely updating processing hierarchy*.
86 For example, it is known that scenes in a movie are encoded as event-specific neural responses³⁵ and
87 that more parietal receptive windows represent increasingly coarse events in movies³⁶. Such an
88 event hierarchy is effectively based on the boundaries of events: It calls for stable working memory
89 representations that are *sparsely* recombined with preceding events at higher processing stages only
90 at the end of an event. The simultaneous representation of distinct events in working memory allows
91 to draw on diverse context when making predictions. We here hypothesize that such a sparsely
92 updating network architecture is a more appropriate model for prediction processes in the brain.

93 In the present study, we recorded blood-oxygen-level-dependent (BOLD) responses while
94 participants listened to a narrated story, which provides rich semantic context and captures the full
95 dynamic range of speech predictability³⁷. Following the rationale that neural computations can be
96 inferred by comparing the fit of neural data to outputs from artificial neural networks with different
97 architectures^{38,39}, we derived context-specific surprisal associated with each word in the story from
98 single layers of long short-term memory (LSTM)-based language models with either a continuous or
99 a sparse updating rule.

100 Here, we show that the event-based organization of semantic context provides a valid model
101 of predictive processing in the brain. We show that a “surprisal hierarchy” of increasingly coarse
102 event timescales evolves along the temporo-parietal pathway, with stronger connectivity to
103 neighbouring timescales in states of higher word surprisal. Surprisal derived from continuously
104 updated context had a (non-hierarchical) effect on temporo-parietal activity only at short timescales.

105 Our results suggest that representing context in the form of increasingly coarse events constitutes a
106 network architecture that is both, computationally efficient and semantically rich for making
107 predictions.

108 **Results**

109 Thirty-four participants listened to a one-hour narrated story while their hemodynamic brain
110 responses were recorded using functional magnetic resonance imaging (fMRI). To emulate a
111 challenging listening scenario, the story audio was presented against a competing stream of
112 resynthesized natural sounds (for an analysis focussing on this acoustic representation see Ref. ⁴⁰).

113 The surprisal associated with each word in the story was modelled at multiple timescales of
114 semantic context by two artificial neural networks, one with a continuous updating rule (LSTM) and
115 another one with a sparse updating rule (HM-LSTM; Figure 1A).

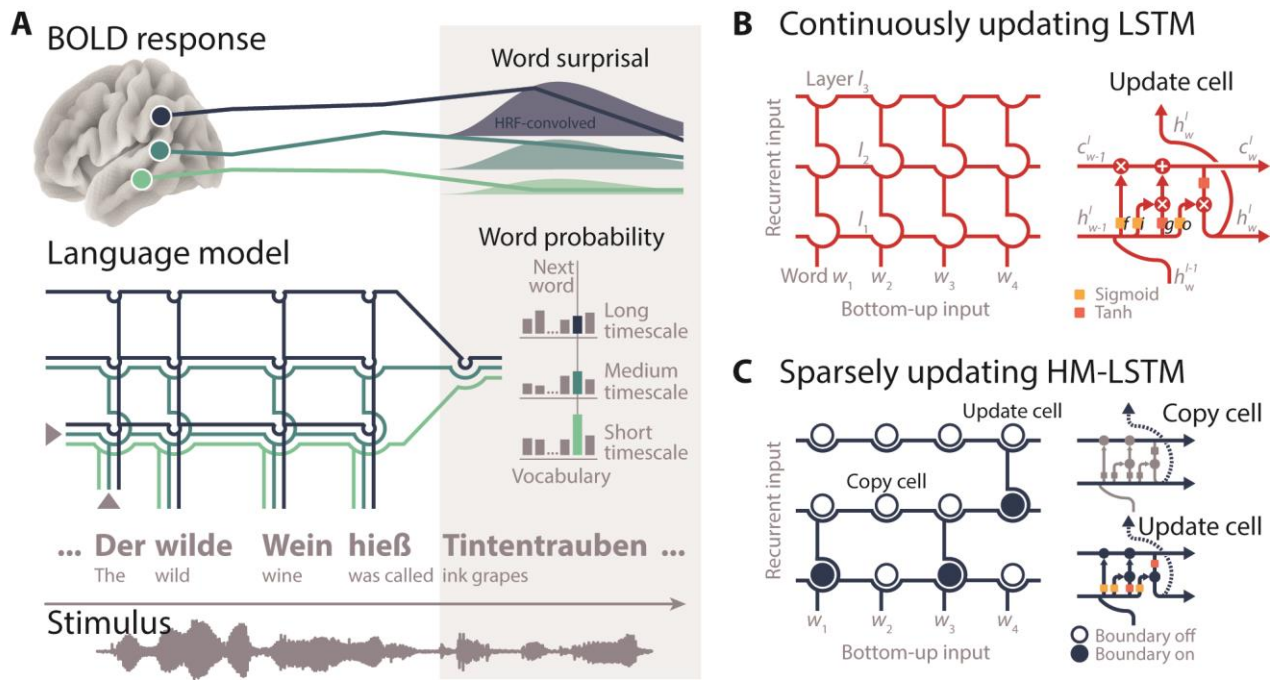
116 First, we encoded surprisal at multiple timescales into univariate neural responses and fit a
117 gradient to temporo-parietal peak locations of timescales. Second, we decoded timescale surprisal
118 from patterns of neural responses in single parcels and compared decoding accuracies between
119 language models. Finally, we investigated how surprisal gates the information flow between brain
120 regions sensitive to different timescales.

121 All encoding and decoding models were estimated separately per each language model,
122 using ridge regression with four-fold cross-validation.

123 **Two competing language models of hierarchical speech prediction**

124 We trained two artificial neural networks on more than 130 million words of running text to predict
125 an upcoming word by its preceding semantic context. More specifically, language models consisted
126 of long short-term memory cells (LSTM)⁴¹, which incorporate context that might become relevant at
127 *some* time (*cell state*) or that is relevant already to the prediction of the *next* word (*hidden state*). By
128 stacking five LSTM layers, our models operated on different timescales of context, with higher layers
129 coding for long-term dependencies between words.

130 In the continuously updating (or “vanilla”) LSTM, recurrent memory states are updated at
131 each layer with every new bottom-up word input (Figure 1B). A second model, the *hierarchical*
132 *multiscale* LSTM⁴², referred to as “sparsely updating HM-LSTM”, employs a revised updating rule
133 where information from a lower layer is only fed forward at the end of its representing timescale
134 (Figure 1C). This allows for less frequent updates between layers and stronger separation between
135 contextual information represented at different layers.



136 **Figure 1. Modelling neural speech prediction with artificial neural networks.** (A) Participants listened to
 137 a story (grey waveform) during fMRI. Based on its semantic context (“The wild wine was called”), a language
 138 model predicted each word in the story (“ink grapes”). The probability of the next word was read out from each
 139 layer of the model separately, with higher layers accumulating information across longer semantic timescales.
 140 Word probabilities were transformed to surprisal, convolved with the hemodynamic response function and
 141 mapped onto temporo-parietal BOLD time series. (B) Two language models were trained. With each new
 142 word-level input, the “continuously updating” long short-term memory (LSTM)⁴¹ combines “old” recurrent
 143 long-term (c_{w-1}^l) and short-term memory states (h_{w-1}^l) with “new” bottom-up semantic input (h_w^{l-1}) at each
 144 layer l . This allowed semantic information to continuously flow to higher layers with each incoming word. f:
 145 forget gate, i: input gate, g: candidate state, o: output gate. (C) The “sparsely updating” hierarchical multiscale
 146 LSTM (HM-LSTM)⁴² was designed to learn the hierarchical structure of text. An upper layer keeps its
 147 representation of context unchanged (copy mechanism) until a boundary indicates the end of a timescale on
 148 the lower layer and information is passed to the upper layer (update mechanism). Networks were unrolled for
 149 illustration only.

150 Three model-derived metrics of predictiveness at multiple timescales

151 For each word in the entire presented story (> 9,000 words), we determined its predictability given
 152 the semantic context of the preceding 500 words. Hidden states were combined across layers and
 153 mapped to an output module, which denotes the probability of occurring next for every word in a
 154 large vocabulary of candidate words. The word with the highest probability was selected as the
 155 predicted next word. Overall, the LSTM (proportion correct across words: 0.13) and the HM-LSTM
 156 (0.12; Supplementary Figure 1) were on par in accurately predicting the next word in the story.

157 To derive the predictability of words based on layer-specific context (or, for our purpose,
 158 timescale), we allowed information to freely flow through pre-trained networks, yet only mapped
 159 the hidden state of one layer to the output module by setting all other network weights to zero.
 160 Outputs from these “lesioned” language models represented the five timescales.

161 As the primary metric of predictiveness, we calculated the degree of surprisal associated with
162 the occurrence of a word given its context (i.e., negative logarithm of the probability assigned to the
163 *actual* next word). The surprisal evoked by an incoming word indexes the amount of information
164 that was not predictable from the context represented at a specific timescale^{43,44}. Of note, surprisal
165 was considerably higher for longer timescales in the LSTM ($p < 0.001$, Cohen's $d = 2.43$; compared to
166 slopes drawn from surprisal shuffled across timescales) but remained stable across timescales in the
167 HM-LSTM ($p = 0.955$, $d = 0.05$; direct comparison LSTM vs. HM-LSTM: $p < 0.001$, $d = 2.7$; Figure 2).

168 To determine the temporal integration window of each timescale, we scrambled input to the
169 network at different granularities corresponding to a binary logarithmic increase in the length of
170 intact context (i.e., 1–256 words). The LSTM showed no such increase of temporal integration
171 windows at higher layers (LSTM: $p = 0.219$, $d = 0.11$). In contrast, in the HM-LSTM, surprisal decreased
172 more strongly at longer compared to shorter timescales as more intact context became available
173 (HM-LSTM: $p = 0.027$, $d = 0.73$; LSTM vs. HM-LSTM: $p < 0.001$, $d = 0.76$; Figure 2).

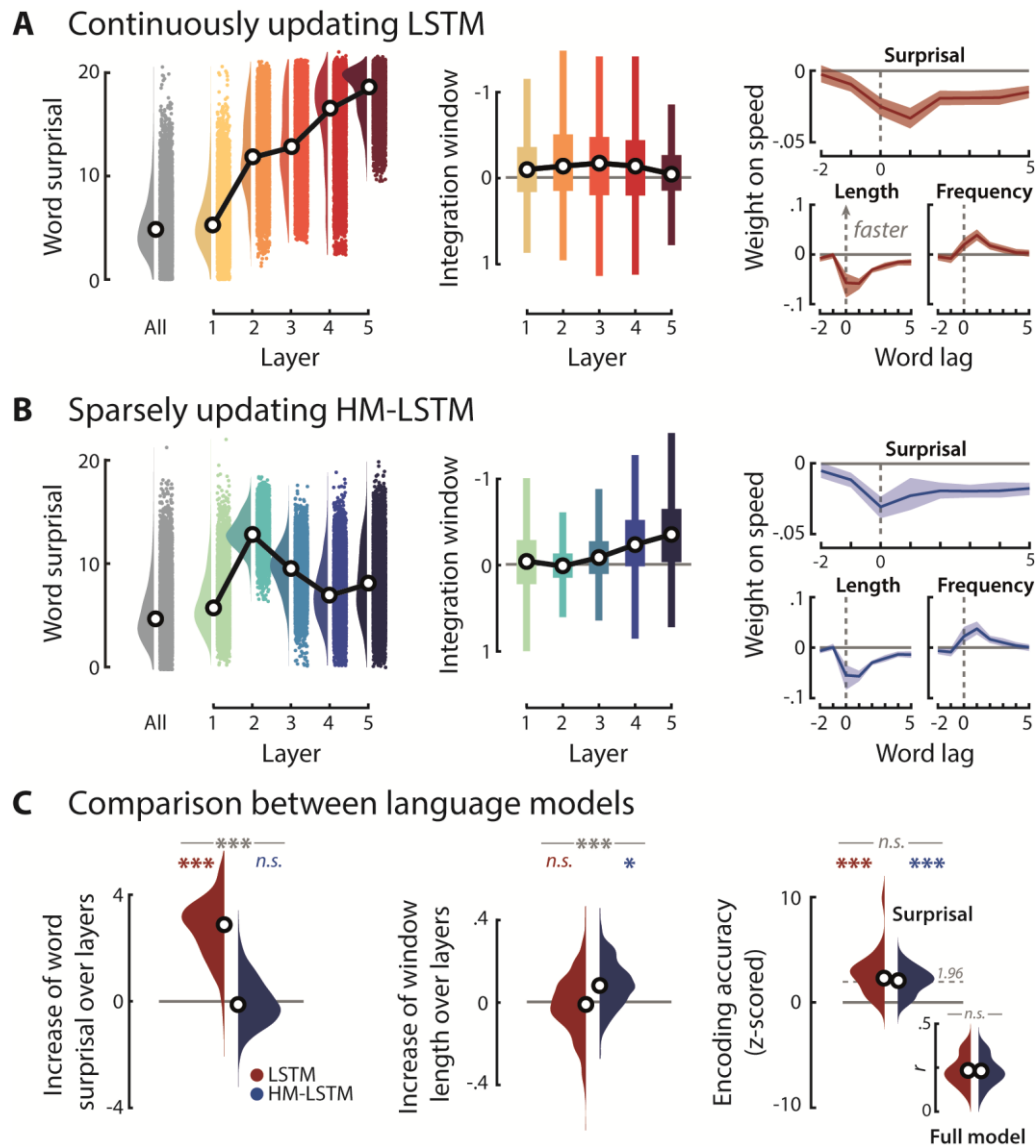
174 Our secondary metrics expressed the predictability of a word in relation to other words, that
175 is, (1) the entropy of the probability distribution predicted for individual words (indicative of the
176 difficulty to make a definite prediction) and (2) the dissimilarity of vector representations (or
177 embeddings) coding for the constituent linguistic features of the predicted and actual next word
178 (Product-moment correlation; indicative of conceptual (un-)relatedness).

179 We derived surprisal, entropy and dissimilarity associated with single words from “lesioned”
180 models at each of five timescale and from “full” models across all timescales, separately for each
181 language model. All features were convolved with the hemodynamic response function, and we will
182 collectively refer to them as “features of predictiveness” from here on.

183 **Higher model-derived surprisal of words slows down reading**

184 To test the behavioural relevance of model-based predictiveness, another 26 participants performed
185 a self-paced reading task where they read the transcribed story word-by-word on a noncumulative
186 display and pressed a button as soon as they had finished reading.

187 When regressing response speed onto time-lagged features of predictiveness and a set of
188 nuisance regressors (e.g., word length and frequency), we found that—as expected—reading speed
189 slowed down for words determined as more surprising by language models given the full context
190 across all timescales (Figure 2A).



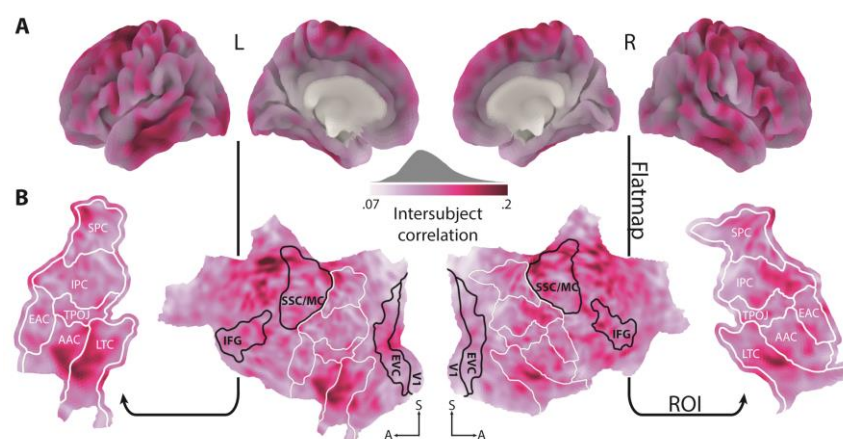
191 **Figure 2. Evaluating model-derived surprisal.** (A) Left: Word surprisal derived from the “full” LSTM model
 192 including all layers (grey distribution) and from single layers of “lesioned” LSTM models (coloured
 193 distributions); black circles represent grand-median surprisal. Middle: Input to the LSTM was scrambled at
 194 different granularities corresponding to an increase in the length of intact context (i.e., 1–256 words). For each
 195 layer of the LSTM, linear functions were fit to word surprisal across these context windows. A negative slope
 196 parameter indicates a stronger benefit (or lower surprisal) from longer context (i.e., larger integration window).
 197 Right: Speed in a self-paced reading task was modelled as a function of time-lagged predictiveness and a set
 198 of nuisance regressors. Weight profiles illustrate the temporal dynamics of the surprisal effect in the full model
 199 (top) in comparison to word length (bottom, left) and word frequency (bottom, right); positive weights indicate
 200 an increase in response speed; error bands represent \pm SEM. (B) Same as in A, but for the sparsely updating
 201 HM-LSTM. (C) Left: We fit linear functions to word surprisal across layers and compared resulting slope
 202 parameters to null distributions drawn from shuffled layers and between LSTM (red) and HM-LSTM (blue).
 203 Middle: Linear fit to integration windows across layers, indicating the benefit of higher layers from longer
 204 context. Right: Encoding accuracy in the self-paced reading task uniquely explained by the predictiveness of
 205 context (standardized to scrambled features of predictiveness); dotted grey line indicates critical significance
 206 level for single participants. Inset shows non-standardized encoding accuracies. *** $p < 0.001$, * $p < 0.05$, *n.s.*:
 207 not significant.

208 Further, we predicted response speed on held-out testing data and z-scored the resulting
209 encoding accuracy (i.e., Product-moment correlation of predicted and actual response speed) to a
210 null distribution drawn from scrambled features of predictiveness while only keeping nuisance
211 regressors intact. This yielded the unique contribution of the predictiveness of words (i.e., surprisal,
212 entropy and dissimilarity) to reading speed, which was significant for both language models (LSTM:
213 $p < 0.001$, $d = 1.51$; HM-LSTM: $p < 0.001$, $d = 1.64$, LSTM vs. HM-LSTM: $p = 0.975$, $d = 0.35$). Together,
214 these findings suggest that both language models picked up on processes of speech prediction that
215 shape behaviour.

216 **Selecting temporo-parietal regions of interest involved in speech processing**

217 We hypothesized that the speech prediction hierarchy is represented as a gradient along the
218 temporo-parietal pathway. This rather coarse region of interest was further refined to only include
219 regions implicated in processing of the listening task.

220 To this aim, we calculated pairwise intersubject correlations⁴⁵, which revealed consistent
221 cortical activity across participants in a broad bilateral language network. Responses were most
222 prominently shared in auditory association cortex and lateral temporal cortex as well as premotor
223 cortex, paracentral lobule and mid cingulate cortex (Figure 3A). Crucially, as sound textures
224 presented in the competing stream were randomly ordered across participants, this approach
225 allowed us to extract shared responses specific to the speech stream.



226 **Figure 3. Selection of regions of interest. (A)** When listening to a story against background noise, pairwise
227 intersubject correlations showed stronger synchronization of BOLD activity in cortical areas implicated in the
228 language network. **(B)** The cortical surface was flattened. All temporal and parietal parcels⁴⁶ highlighted by
229 white outlines were included as regions of interest (ROI) in the following analyses. Black outlined parcels serve
230 as reference point, only. EAC: early auditory cortex, AAC: auditory association cortex, LTC: lateral temporal
231 cortex, TPOJ: temporo-parieto-occipital junction, IPC: inferior parietal cortex, SPC: superior parietal cortex, V1:
232 primary auditory cortex, EVC: early visual cortex, IFG: inferior frontal gyrus, SSC/MC: somatosensory and motor
233 cortex. Maps were smoothed with an 8 mm FWHM Gaussian kernel for illustration only.

234 All further analyses were limited to those parcels in temporal and parietal cortex⁴⁶ that
235 showed significant median intersubject correlations in more than 80 % of vertices ($p_{\text{FDR}} < 0.01$; ranked
236 against a bootstrapped null distribution). The cortical sheet of the six parcels determined as regions
237 of interest (ROI) was flattened, resulting in a two-dimensional plane spanned by an anterior-posterior
238 and inferior-superior axis (Figure 3B).

239 We expected gradients of speech prediction to unfold along the inferior-superior axis, that
240 is, from temporal to parietal areas.

241 **Differential tuning to continuously and sparsely updated timescales of surprisal in temporo-** 242 **parietal cortex**

243 After the predictiveness of speech was encoded into neural activity of single vertices, we extracted
244 temporo-parietal weight maps of word surprisal at each timescale for both language models.

245 When performing spatial clustering on these weight maps (p_{vertex} and $p_{\text{cluster}} < 0.05$; compared
246 to scrambled surprisal by means of a cluster-based permutation test), we found large positive
247 clusters in both hemispheres for shorter timescales of the LSTM (Figure 4A, yellow outlines) but, if at
248 all, only focal clusters for longer timescales (Figure 4A, red outlines). Hence, temporo-parietal activity
249 primarily increased in response to words that were less predictable by the context provided at
250 shorter, continuously updated timescales. In contrast, clusters of distinct polarity, location and
251 extent were observed for all timescales of the HM-LSTM (Figure 4B), suggesting that even longer
252 timescales had the potency to modulate temporo-parietal activity when they were sparsely updated.

253 **Sparsely updated timescales of surprisal evolve along a temporo-parietal processing** 254 **hierarchy**

255 To probe the organization of timescales along a temporo-parietal gradient, we collapsed across the
256 anterior-posterior axis of weight maps and selected the local maximum with the largest positive
257 value on the inferior-superior axis. Fitting a linear function to those peak coordinates of timescales,
258 we found flat slope parameters indicating random ordering of LSTM timescales in both hemispheres
259 (left: $p = 0.458$, $d = -0.15$; right: $p = 0.716$, $d = -0.07$; compared to slopes drawn from coordinates
260 scrambled across timescales, Figure 4C). Conversely, we found steep positive slopes for the HM-LSTM
261 in both hemispheres (left: $p < 0.001$, $d = 0.72$; right: $p < 0.001$, $d = 0.75$; Figure 4C), reflecting the
262 representation of longer timescales in more parietal regions. On average, the left hemisphere
263 represented sparsely updated timescales 12 mm superior (along the unfolded cortical surface) to
264 their directly preceding timescale. Most relevant, this finding was underpinned by a significant
265 difference of slope parameters between the LSTM and HM-LSTM (left: $p = 0.005$, $d = 0.9$; right: $p <$
266 0.001 , $d = 0.89$; Figure 4C), demonstrating a temporo-parietal processing hierarchy of word surprisal
267 that preferably operates on sparsely updated timescales.

268 The absence of a gradient for continuously updated timescales was corroborated when
269 specifically targeting the dorsal processing stream. To this aim, we confined the first timescale to
270 peak in temporal regions and all other timescales to peak superior to the first timescale in more
271 parietal regions. Slope effects along the dorsal stream largely complied with those found along the
272 inferior-superior axis (LSTM: all $p \geq 0.134$; HM-LSTM: all $p \leq 0.006$; LSTM vs. HM-LSTM: all $p \geq 0.103$),
273 thereby ruling out the possibility that the presence of a competing ventral stream obscured the
274 consistent ordering of timescales.

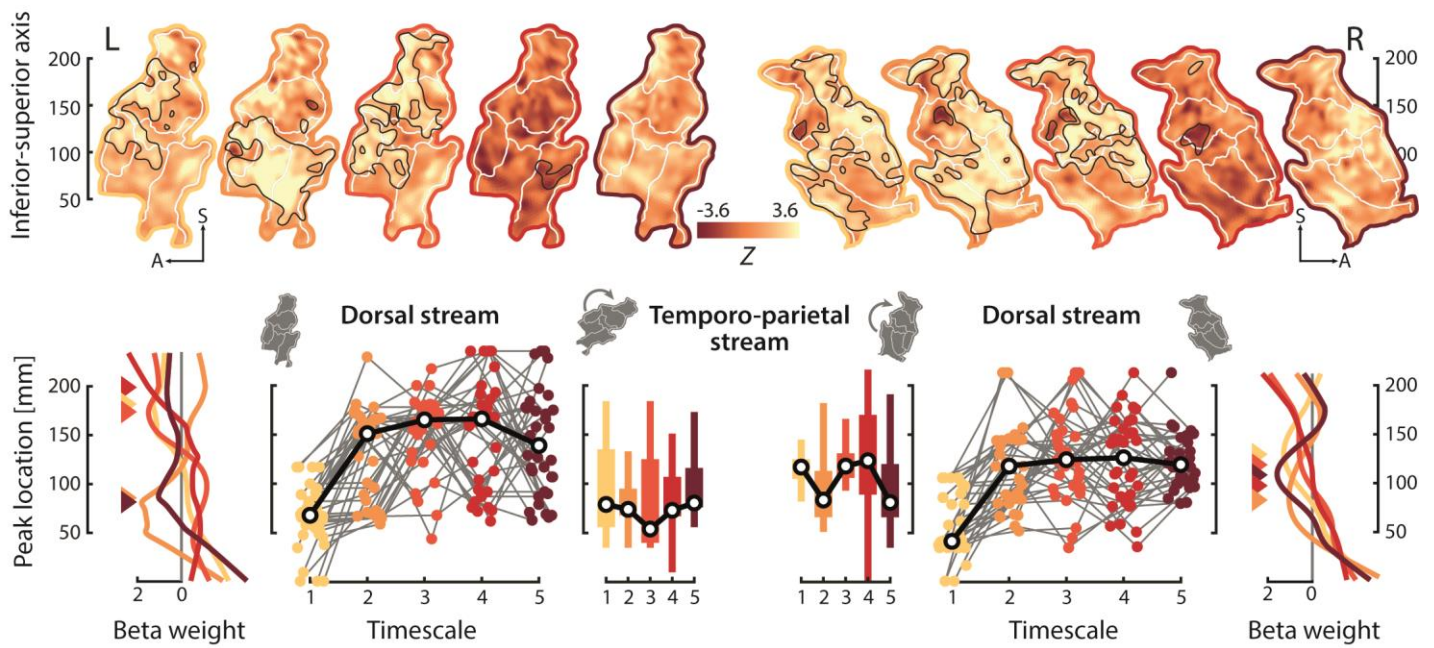
275 Additionally, rotating weight maps by -45° before collapsing across the first dimension
276 showed that sparsely updated timescales were not only processed along an inferior-superior but
277 also an anterior-posterior gradient in the left hemisphere (LSTM: $p = 0.921$, $d = 0.02$; HM-LSTM: $p =$
278 0.001 , $d = 0.65$; LSTM vs. HM-LSTM: $p = 0.011$, $d = 0.55$). As right-hemispheric maps already had a
279 strong initial rotation off the inferior-superior axis, rotating these maps merely confirmed that longer
280 timescales are processed in more superior regions (LSTM: $p = 0.355$, $d = -0.18$; HM-LSTM: $p < 0.001$, d
281 $= 1.39$; LSTM vs. HM-LSTM: $p = 0.039$, $d = 1.45$).

282 Unlike for the sparsely updated timescales of surprisal, neither the timescales of entropy (all
283 $p \geq 0.583$) nor dissimilarity (all $p \geq 0.623$) organized along a dorsal gradient (Supplementary Figure
284 2). Further, effects of HM-LSTM timescale surprisal were dissociable from a simple measure of
285 semantic incongruence between words in the story and their context at five timescales
286 logarithmically increasing in length (all $p \geq 0.5$; Product-moment correlation of target and average
287 context embedding). This highlights the specificity of the observed gradient to prediction processes
288 in general and word surprisal in particular.

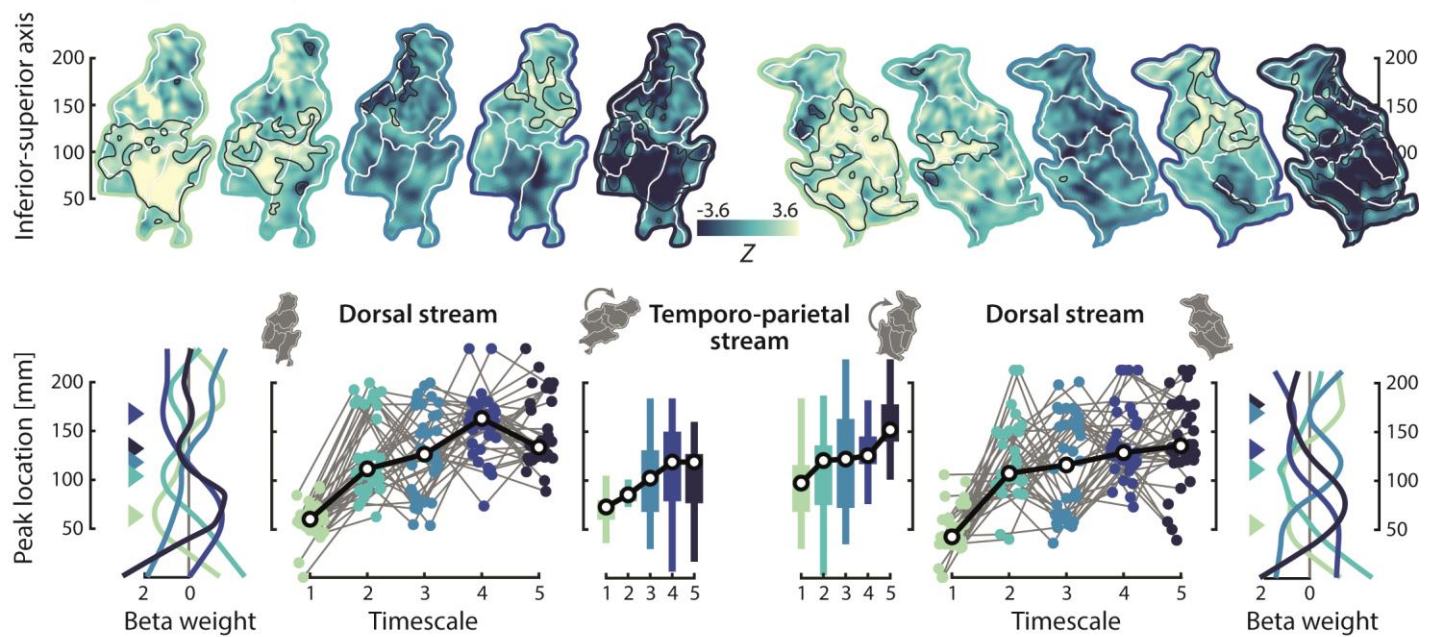
289 **A segregated stream of continuously updated timescales of surprisal?**

290 To determine the contribution of predictiveness to overall encoding accuracy on held-out data, we
291 z-scored accuracies relative to null distributions drawn from scrambled features of predictiveness
292 while keeping additional (spectro-temporal) acoustic and linguistic nuisance regressors intact.
293 Interestingly, the LSTM produced—in comparison to the HM-LSTM—better predictions in early
294 auditory cortex and supramarginal gyrus (p_{vertex} and $p_{\text{cluster}} < 0.05$; cluster-based permutation paired-
295 sample t -test; Figure 4C). On the other hand, predictions of the HM-LSTM seemed slightly more
296 accurate than for the LSTM along middle temporal gyrus, temporo-parieto-occipital junction and
297 angular gyrus, even though not statistically significant. Taking into account the broad clusters found
298 earlier specifically for shorter (but not longer) LSTM timescales, this poses the question whether
299 continuously updated timescales take full effect only in earlier medial temporal and anterior parietal
300 processing stages, whereas the sparsely updating processing hierarchy evolves along a separate
301 lateral temporal and posterior parietal route.

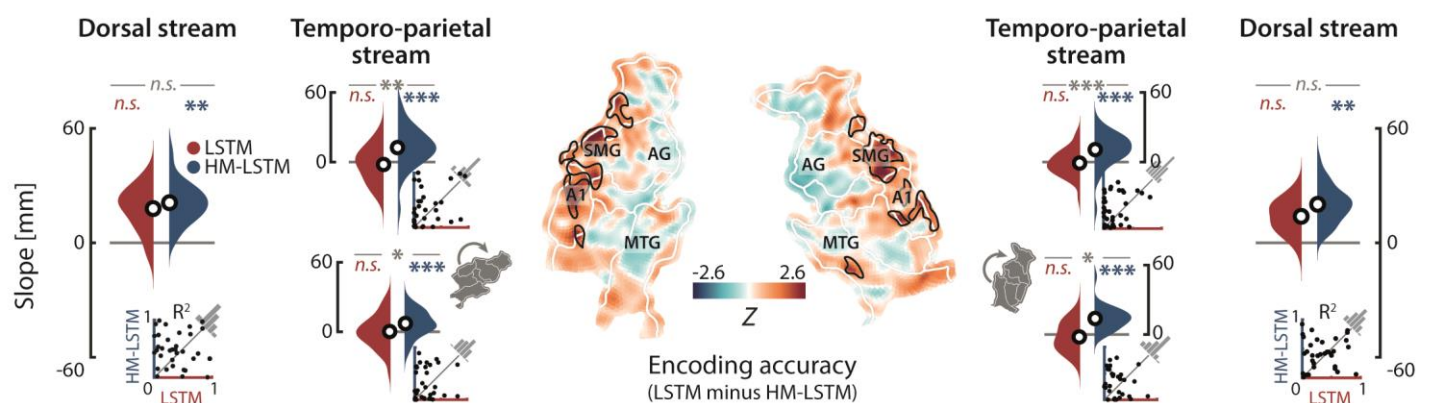
A Continuously updating LSTM



B Sparsely updating HM-LSTM



C Comparison between language models



303 **Figure 4. Encoding the timescales of surprisal. (A)** Temporo-parietal BOLD time series were mapped onto
304 the predictiveness of speech derived at five timescales from the continuously updating LSTM. Top row: Maps
305 show z-values from timescale-specific weights of surprisal tested against a null distribution drawn from
306 scrambled surprisal; black outlines indicate significant clusters; white outlines indicate parcels; coloured
307 outlines indicate short (light yellow) to long (dark red) timescales, separately for the left and right hemisphere.
308 Bottom row: For each timescale, we determined its peak coordinate along the inferior-superior axis (coloured
309 triangles), here shown for grand-average weight profiles for illustration only. Testing for a processing hierarchy
310 along the dorsal stream, timescales were restricted to peak superior to the first timescale; coloured dots
311 connected by grey lines represent peak coordinates of single participants; black circles represent grand-
312 median peak coordinates. Testing for a processing hierarchy along the temporo-parietal stream, timescales
313 were allowed to peak at any location. Maps were rotated by -45° to test simultaneous effects on the inferior-
314 superior and anterior-posterior axis in the left hemisphere. Note that right-hemispheric maps already had an
315 initial rotation off the inferior-superior axis, so rotating these maps resulted in testing for effects on the inferior-
316 superior axis only. **(B)** Encoding maps and timescale-specific peak locations for the sparsely updating HM-
317 LSTM. **(C)** Linear functions were fit to peak coordinates across timescales and resulting slope parameters were
318 compared to null distributions drawn from scrambled coordinates and between language models, separately
319 for the dorsal (not rotated) and temporo-parietal stream (top: not rotated; bottom: rotated) in both
320 hemispheres. Black circles represent grand-average slope parameters; insets depict coefficients of
321 determination for linear fits of single participants. Temporo-parietal encoding accuracies displayed on ROI
322 maps were z-scored to null distributions drawn from scrambled features of predictiveness and compared
323 between language models; black outlines indicate significant clusters; maps were smoothed with an 8 mm
324 FWHM Gaussian kernel for illustration only; SMG: supramarginal gyrus, AG: angular gyrus, A1: primary auditory
325 cortex; MTG: middle temporal gyrus. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, *n.s.*: not significant.

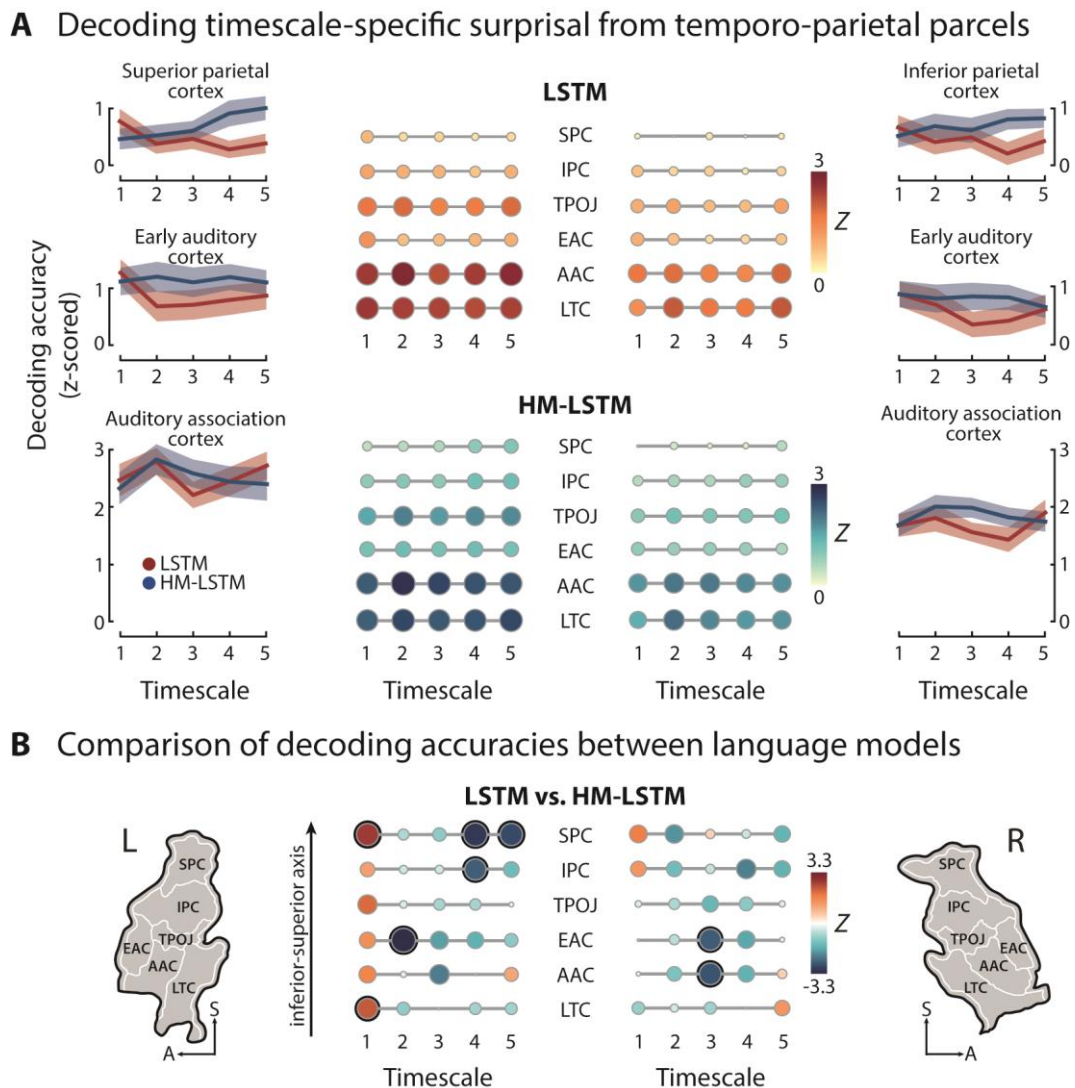
326 **Parietal regions preferentially represent sparsely updated, long timescales**

327 In a complementary decoding approach, we reconstructed the timescales of surprisal from patterns
328 of neural activity in single regions of interest (i.e., temporo-parietal parcels). Reconstructed
329 timescales were z-scored to scrambled features of predictiveness. Overall, bilateral auditory
330 association cortex and lateral temporal cortex yielded highest decoding accuracies on held-out
331 testing data for both language models (Figure 5A). More parietal regions showed comparably lower
332 decoding accuracies. Nevertheless, these accuracies can be deemed meaningful, as the pattern
333 mirrors the lower intersubject correlations in parietal compared to temporal regions (Figure 3),
334 which are commonly found during natural listening (e.g., Ref. ⁴⁷⁻⁴⁹). This indicates an overall greater
335 variability of neural responses in parietal regions irrespective of the timescales of surprisal.

336 Contrasting decoding accuracies between language models, temporo-parietal regions of
337 interest showed an overall preference for the shortest LSTM but longer HM-LSTM timescales (Figure
338 5B). This suggests that the predominance of the LSTM in early auditory cortex and supramarginal
339 gyrus observed for encoding accuracies of the encoding model is specific to the shortest timescale,
340 while lateral temporal and posterior parietal regions reflect longer timescales of the HM-LSTM,
341 lending further support to the functional dissociation of two routes of predictive processing in
342 speech prediction.

343 In particular, left-hemispheric early auditory cortex contained more information on medium,
344 sparsely updated than continuously updated timescales, whereas inferior and superior parietal
345 cortex preferentially represented long, sparsely updated timescales ($p_{\text{FDR}} < 0.05$). This finding

346 converges with the organization of the gradient described earlier for sparsely updated but not
 347 continuously updated timescales.



348 **Figure 5. Decoding the timescales of surprisal.** (A) Surprisal at different timescales was decoded from
 349 regions of interest. Matrices depict decoding accuracies determined on held-out testing data and z-scored to
 350 null distributions drawn from scrambled surprisal, separately for the LSTM (top) and HM-LSTM (bottom) in the
 351 left and right hemisphere; colour and size of circles scale to decoding accuracy. Of note, some z-scored
 352 decoding accuracies in more superior parcels fell below an average value of 1.96. However, z-scores were
 353 indicative of significance only on the level of single participants. Line plots illustrate patterns of decoding
 354 accuracies across timescales in select regions of interest; error bands represent \pm SEM. (B) Decoding accuracies
 355 were contrasted between language models by means of a permutation test on the mean of differences; black
 356 circles indicate $p_{\text{FDR}} < 0.05$; maps indicate location of parcels. EAC: early auditory cortex, AAC: auditory
 357 association cortex, LTC: lateral temporal cortex, TPOJ: temporo-parieto-occipital junction, IPC: inferior parietal
 358 cortex, SPC: superior parietal cortex.

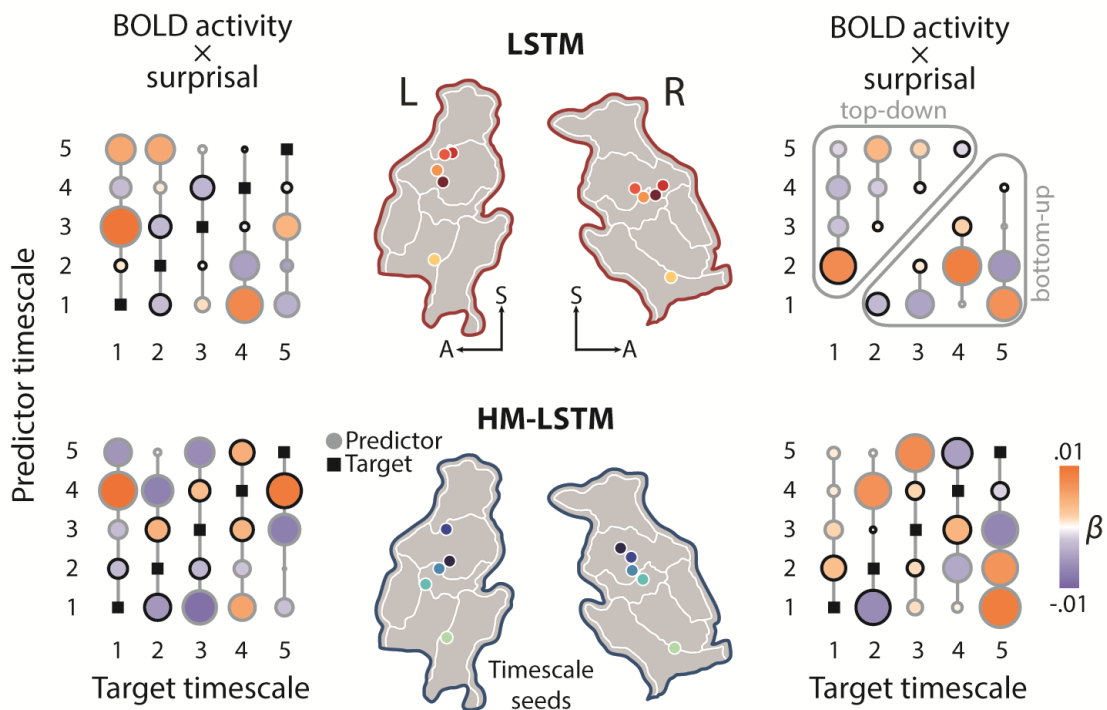
359 **Surprisal at sparsely updated timescales gates connectivity along the processing hierarchy**

360 After establishing the temporo-parietal processing hierarchy, we examined the modulatory effect of
361 surprisal on connectivity between peak locations of timescales taken from the encoding analysis. To
362 this aim, we created psychophysiological interactions between the BOLD response at the peak
363 location of one timescale and word surprisal at the same timescale. The BOLD response of each
364 (target) timescale was mapped onto psychophysiological interactions of all other (predictor)
365 timescales (Figure 6A).

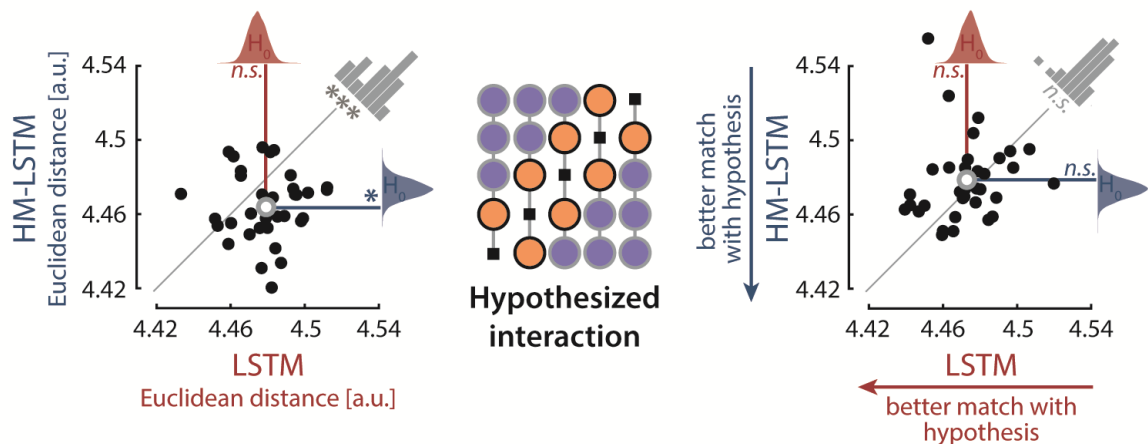
366 We hypothesized that coupling between brain regions representing two neighbouring
367 timescales increases when one timescale becomes unpredictable. Numerically, this can be expressed
368 by setting the weights of neighbouring timescales to 1 and all other predictor weights to -1 (Figure
369 6B). This hypothesized pattern of weights was not matched by the weights observed for the LSTM
370 (left: $p = 0.83$, $d = 0.27$; right: $p = 0.348$, $d = 0.1$; Euclidean distance compared to null distributions
371 drawn from target BOLD activity shifted in time), which was expected given that the continuously
372 updated timescales were not organized along a gradient in the first place. Critically, for sparsely
373 updated timescales of the HM-LSTM, surprisal-modulated connectivity in the left hemisphere not
374 only matched our hypothesis (left: $p = 0.032$, $d = 0.5$; right: $p = 0.853$, $d = 0.23$) but also matched our
375 hypothesis better than the LSTM (left: $p = 0.001$, $d = 0.79$; right: $p = 0.902$, $d = 0.3$).

376 To specify the directionality of information flow, we averaged weights of top-down
377 modulations (i.e., predictor weights of timescales longer than respective target timescales) and
378 bottom-up modulations. We found no difference between the modulatory strength of these top-
379 down and bottom-up connections (LSTM left: $p = 0.184$, $d = 0.23$; right: $p = 0.839$, $d = 0.4$; HM-LSTM
380 left: $p = 0.408$, $d = 0.4$; right: $p = 0.367$, $d = 0.16$).

A Functional connectivity modulated by timescale-specific surprisal



B Comparison of modulated connectivity between language models



381 **Figure 6. Surprisal-dependent modulation of effective connectivity.** (A) A sphere of 6 mm was centred on
 382 median peak locations of timescales as defined in the encoding analysis (coloured circles on temporo-parietal
 383 maps) and BOLD responses were averaged within these timescale seeds. BOLD time series at one (target)
 384 timescale were regressed onto psychophysiological interactions of all other (predictor) timescales (i.e.,
 385 pointwise product of timescale-specific BOLD and surprisal time series). For each target seed, we added a
 386 column vector of timescale-specific predictor weights to a 5-by-5 matrix with an empty main diagonal.
 387 Matrices were created separately for each language model (top: LSTM, bottom: HM-LSTM) and hemisphere.
 388 The upper triangle of a matrix indicates top-down, the lower triangle bottom-up information flow; black
 389 outlines of circles indicate timescale pairing for which we expected an increase of connectivity, a decrease of
 390 connectivity was expected for pairings with grey outlines (see hypothesized interaction pattern in 6B). (B) A
 391 hypothesized matrix of psychophysiological interactions was created, with positive weights on diagonals
 392 below and above the main diagonal, indicating increased connectivity between neighbouring timescales
 393 when surprisal is high. The Euclidean distance between observed and hypothesized matrices was compared
 394 to null distributions of distances drawn from target timescales shifted in time (coloured density plots),
 395 separately for the left and right hemisphere; black dots indicate distances of single participants; grey circles
 396 indicate mean distances. * p < 0.05; *** p < 0.001; n.s.: not significant.

397 **Discussion**

398 How are the complex temporal dependencies underlying natural speech processed in the brain to
399 inform predictions of upcoming speech? In the present study, we emulated these prediction
400 processes in two language models (artificial neural networks; LSTM vs. HM-LSTM), which critically
401 differed in how often semantic context representations are updated at multiple, hierarchically
402 organized timescales.

403 Surprisal as derived from both language models modulated reading times in the behavioural
404 reading task to a similar degree, while hemodynamic brain responses to surprisal during the listening
405 task differed between models: In early auditory cortex and supramarginal gyrus, the continuously
406 updating LSTM predicted activity better than the sparsely updating HM-LSTM. In general, surprisal
407 at the shortest timescale was decoded more precisely from temporo-parietal regions when derived
408 from the continuously updating LSTM than from the HM-LSTM.

409 In contrast, and in line with our initial hypothesis, temporo-parietal regions hierarchically
410 encoded the (sparsely updated) event-based surprisal provided by the layers of the HM-LSTM, with
411 longer timescales represented in inferior parietal regions. Moreover, higher timescale-specific
412 surprisal based on the HM-LSTM increased connectivity from receptive windows of a given timescale
413 to their immediately neighbouring (shorter or longer) timescales.

414 Together, these results provide evidence for the neurobiological parsimony of an event-
415 based processing hierarchy. In the present data, this was expressed in the simultaneous neural
416 representation of surprisal at multiple timescales, and in surprisal dynamically gating the
417 connectivity between these timescale-specific receptive windows.

418 **The event-based organization of context as a foundation for language prediction**

419 The spatial organization of timescale-specific receptive windows observed in the present study
420 converges with previous results, where bilateral primary auditory cortex coded for relatively shorter
421 timescales (e.g., words) and inferior parietal cortex coded for longer timescales (e.g., paragraphs; Ref.
422 ²⁸). Critically, this spatial overlap was found despite targeting different aspects of speech processing.

423 In our study, neural responses were expressed as a function of timescale-specific word
424 surprisal, a proxy tapping into prediction processes. In other studies, receptive windows were based
425 on the (in-)consistency of neural activity across participants in response to speech input at varying
426 timescales⁵⁰, which is typically linked to working memory formation. This implies that the same
427 neural system most likely fulfils distinct functions: Temporal receptive windows have been
428 suggested to store timescale-specific context in working memory, and, in parallel, exploit this
429 context to process information in the present⁹. In line with this theoretical account, our results

430 suggest that timescale-specific memory representations serve as the basis for the generative models
431 shaping predictions of upcoming speech.

432 The here observed temporo-parietal gradient of surprisal at sparsely updated
433 representations of context is specifically well in line with accounts of neural event segmentation^{51,52},
434 and with the notion of hierarchical multiscale network architectures more generally (here, HM-
435 LSTM). Taking a sentence from our listening task as an example, “The wild wine was called ink grapes”
436 is embedded in a brief event where the narrator describes how the bluish black of the grapes in a
437 backyard reminded her of the color of the night. At the same time, the sentence is embedded in a
438 larger event of the author wandering around the single rooms and the garden of her parents’ house,
439 and an even larger event of the author reliving the memory of walking the streets in her Romanian
440 hometown.

441 The HM-LSTM architecture resembles neural event segmentation in two decisive points.
442 First, the boundary detector allows revealing the event structure of context, similar to an increase of
443 neural activity indexing prediction errors at event boundaries⁵³⁻⁵⁵. Second, the sparse updates to
444 higher processing stages at event boundaries allow retaining multiple, stable context
445 representations in memory, similar to temporo-parieto-occipital receptive windows reflecting the
446 hierarchical event structure during movie watching³⁶. We directly tie in with this result by showing
447 that such hierarchical, event-based context enables neural prediction processes.

448 What are the computational and mechanistic implications of this contextual architecture for
449 the prediction of speech? Somewhat paradoxically, event models have been referred to as “an added
450 burden for an organism”⁵⁶. This argument is certainly plausible with regard to the size of the
451 parameter space, which increases in an artificial or, likewise, biological neural network by
452 introducing an additional boundary detector. At the expense of model parsimony, however, such an
453 event-based network allows for less updates in comparison to continuously updating networks like
454 the LSTM, where each new input to the model elicits computationally complex updates to all
455 timescales.

456 The trade-off between computational costs of boundary detector versus update frequency
457 is nicely illustrated by the fact that the sparsely updating HM-LSTM is considerably faster in making
458 predictions than the continuously updating LSTM⁴². Thus, from a functional perspective, keeping
459 layered representations of multiple events in memory allows to efficiently draw on diverse
460 information to make predictions on upcoming speech.

461 **Context-dependent surprisal as a gating mechanism for predictions and prediction errors**

462 The hallmark of prediction processes in our data is the increase in reading times and neural activity
463 observed in response to more surprising input^{19,57,58}. There are different computational ways in which

464 this “expectation suppression”⁵⁹ can be realized, namely integration difficulty, neural sharpening,
465 and predictive coding.

466 One take on expectation suppression is that surprising sensory input is more difficult to
467 integrate into already existing representations of context because it conveys a relatively larger
468 amount of new information^{60,61}. As the architecture of the sparsely updating HM-LSTM dictates that
469 new information is integrated into timescale-specific representations only at event boundaries,
470 integration difficulty should arise primarily at an event-by-event basis. However, surprisal indeed
471 varies on a word-by-word basis. This conceptual mismatch renders it highly unlikely that integration
472 difficulty accounts for our effects of event-based surprisal. The other two accounts both assume that
473 expectation suppression is indicative of prediction processes but differ in how these processes are
474 thought to be implemented in the brain.

475 Sharpening accounts argue that *unexpected* components of sensory input are suppressed
476 via feedback predictions^{62,63}, resulting in an overall decrease of neural activity in response to
477 expected input. Under the predictive coding account^{12,64}, the brain filters out (or “dampens”)
478 expected components of sensory input, so remaining neural activity (mostly related to prediction
479 errors) is overall smaller for more expected input. The similarity between hypothesized response
480 patterns makes it notoriously hard to disentangle those accounts^{32,63,65}.

481 Notably, however, a distinguishing feature of predictive coding is the specificity of
482 feedforward prediction error signals, which can be captured by modelling effective connectivity
483 between receptive windows of timescales. In agreement with the hierarchical information flow laid
484 out in predictive coding¹², surprisal in our study modulated connectivity via bidirectional links
485 between neighbouring receptive windows of longer and shorter event-based timescales in the left
486 hemisphere (Figure 6).

487 Surprisal in the event-based artificial neural network was modelled as the amount of
488 information an input word conveys that cannot be explained away by the context (or generative
489 model) represented at a specific timescale. Therefore, the increase of feedforward connectivity in
490 response to higher surprisal precisely aligns with the concept of prediction errors in predictive
491 coding⁶⁴.

492 In addition, the increase of feedback connectivity in response to higher surprisal accords
493 with an electrocorticography study in macaques by Chao and colleagues²². The study showed that
494 prediction errors evoked in tone sequences trigger feedback signals from prefrontal to anterior
495 temporal and early auditory cortex in alpha and beta frequency bands. Extending these previous
496 results, our findings suggest that surprisal initiates bottom-up prediction errors, indicative of
497 imprecise predictions, and top-down updates to predictions at processing stages of shorter events
498 to facilitate perception of new words.

499 As an interim conclusion, our findings have two important implications for frameworks of
500 prediction and prediction error: First, we show that a multi-layered hierarchy of predictive coding
501 (e.g., Ref. ¹⁴) applies well to higher-order semantic language processing. Second, predictive coding
502 remains a viable account of neural processing, also when put to test using complex temporal
503 dependencies underlying real-life stimuli.

504 **Implications for a larger network perspective on the event-based prediction hierarchy**

505 Dual stream models of language propose that speech processing is organized along a ventral and a
506 dorsal stream^{66,67}. In the present study, we found a hierarchy of speech prediction along the dorsal
507 stream, which emanated from early auditory cortex and extended well into parietal cortex (Figure
508 4B).

509 This result may seem at odds with other studies showing an additional mirror-symmetric
510 ventral gradient, in which more complex speech features are represented in more anterior temporal
511 regions⁶⁸. The ventral stream has been proposed to chunk speech features into increasingly abstract
512 concepts irrespective of their temporal presentation order⁶⁹. In contrast, we here modelled context
513 representations by respecting the temporal order of words, that is, the HM-LSTM integrates
514 incoming words into an event until words become too dissimilar to previous words and a new event
515 is created. Hence, the ventral stream may contribute to hierarchical speech prediction by exploiting
516 another, more nested facet of context.

517 The inferior frontal gyrus (IFG), alongside premotor cortex, is deemed the apex of the dorsal
518 stream⁶⁷, yet we here considered only the role of temporo-parietal cortex in speech prediction.
519 Previous studies showed that activity in IFG relies on longer timescales of speech being intact^{28,70},
520 that connectivity between IFG and superior temporal gyrus is driven by expectations^{71,72}, and that
521 right IFG is sensitive to the violation of non-local regularities^{73,74}. While this suggests an interplay
522 between frontal and temporo-parietal regions in hierarchical speech prediction, the precise
523 anticipatory mechanisms IFG exerts cognitive control over are just as unclear as how top-down
524 cognitive control and bottom-up sensory input are balanced along the hierarchy.

525 Beyond short-term semantic context, also long-term knowledge facilitates speech
526 prediction. In theory, both memory systems can be couched into the larger framework of the dual
527 reference frame system⁷⁵, where flexible sensory knowledge in parietal cortex interacts with stable
528 conceptual knowledge in hippocampus. Consistent with the key characteristics of the speech
529 prediction hierarchy, hippocampus codes for boundaries in the environment^{76,77}, hierarchically
530 organizes memories⁷⁸ and engages in predictive coding^{79,80}. As parietal cortex has been shown to
531 interface with hippocampus at event boundaries of longer timescales during movie watching³⁶, we
532 speculate that the hierarchy of speech prediction might extend from receptive windows in parietal
533 cortex to hippocampus.

534 Importantly, the event-based prediction hierarchy relies on a set of neural computations—
535 i.e., event segmentation, temporal receptive windows, predictive coding—available beyond the
536 domain of language. Our results thereby encourage future studies to probe its generalizability to
537 other species, sensory modalities, and cognitive functions.

538 **Alternative mechanisms of predictive processing in lieu of event-based timescales**

539 Although we only found a processing hierarchy for surprisal of sparsely updated timescales,
540 temporo-parietal regions were nevertheless sensitive to continuously updated timescales. In
541 particular, decoding accuracies suggested a predominance of the LSTM over the HM-LSTM at the
542 shortest timescale and encoding accuracies suggested a predominance in medial temporal and
543 anterior parietal regions (Figure 4C and 5B). This finding agrees with previous studies showing that
544 participants track changes to situational dimensions of narratives both “globally” at the end of an
545 event and “incrementally” within events⁸¹ and that computational models with continuous updates
546 to all hierarchical levels can explain the construction of context representations in temporo-parietal
547 regions³⁴. Could continuously updated context representations, after all, play an integral role in
548 successful speech prediction?

549 One potential explanation for the negligible neural effects at longer timescales is that the
550 continuously updating language model relies primarily on shorter timescales in predicting the next
551 word. This is supported by the considerably worse model performance observed for longer LSTM
552 timescales (i.e., higher average word surprisal) compared to both shorter LSTM timescales and,
553 despite comparable overall model performance, all HM-LSTM timescales (Supplementary Figure 1).
554 Interestingly, the accuracy in predicting reading speed from word surprisal was the same between
555 LSTM and HM-LSTM (Figure 2), suggesting that continuous updates make for an efficient mechanism
556 to generate equally accurate predictions while relying on less timescales. The strength of such
557 continuously updating models is that context representations are more integrated with what is
558 currently relevant for prediction.

559 A unifying account might be that continuous and sparse updating mechanisms form one
560 instead of two distinct processing streams. For example, Sainburg and colleagues⁸² showed that
561 short-range dependencies of acoustic speech features follow sequential Markovian processes,
562 whereas long-range dependencies follow hierarchical processes. This poses the question whether
563 such interactions of different sequencing mechanisms also better match the semantic structure of
564 speech. Future studies could test this by setting up hybrid language models with continuous
565 updates on shorter and sparse updates on longer timescales.

566 **Conclusion**

567 The present study bridges the gap between the hierarchical, temporally structured organization of
568 context in language comprehension on the one hand and the more general principles of hierarchical
569 predictive processing in cerebral cortex on the other hand.

570 Combining continuously narrated speech, artificial neural networks, and functional MRI
571 building on these networks' output allowed us, first, to sample the natural dynamic range of word-
572 to-word changes in predictiveness over a multi-level hierarchy. Second, we were able to
573 systematically compare the neural effects of different contextual updating mechanisms.

574 Our data demonstrate that the prediction processes in language comprehension build on
575 an event-based organization of semantic context along the temporo-parietal pathway. Not least, we
576 posit that such an event-based organization provides a blueprint for a semantically rich, yet
577 computationally efficient network architecture of anticipatory processing in complex naturalistic
578 environments.

579 **Data availability:** All functional data are publicly available on the Open Science Framework (OSF;
580 <https://osf.io/zbuah>). Custom code will be made available upon publication.

581 **Funding:** This research was supported by German Research Foundation (DFG) grants to JO (OB
582 352/2-1) and GH (HA 6314/4-1), and a European Research Council (ERC) consolidator grant to JO
583 (ERC-CoG-2014-646696).

584 **Acknowledgements:** We thank Anne Herrmann, Clara Mergner, Malte Naujokat, Anna Ruhe, and
585 Svenja Meyn for their help with data acquisition; Christine Sickert for her help in preparing the text
586 corpus; Martin Göttlich for setting up the MR sequences; Mattias Heinrich for discussions on natural
587 language processing; and Malte Wöstmann for suggestions to the reading-task design.

588 **Methods**

589 **Participants**

590 Thirty-seven healthy, young students took part in the fMRI listening study. The final sample included
591 $N = 34$ participants (18–32 years; $M = 24.65$; 18 female), as data from one participant was excluded
592 from all analyses due to strong head movement throughout the recording (mean framewise
593 displacement > 2 SD above group average⁸³) and two experimental sessions were aborted because
594 participants reported to not understand speech against noise. Another 26 students (19–32 years; M
595 $= 23.54$; 17 female) took part in the behavioural self-paced reading study.

596 All participants were right-handed German native speakers who reported no neurological,
597 psychiatric or hearing disorders. Participants gave written informed consent and received an
598 expense allowance of €10 per hour of testing. The study was conducted in accordance with the
599 Declaration of Helsinki and was approved by the local ethics committee of the University of Lübeck.

600 **Stimulus materials**

601 As a speech stimulus in the fMRI listening task, we used the first 64 minutes of an audio recording
602 featuring Herta Müller, a Nobel laureate in Literature, reminiscing about her childhood as part of the
603 German-speaking minority in the Romanian Banat (“Die Nacht ist aus Tinte gemacht”, 2009).

604 To emulate an acoustically challenging scenario in which listeners are likely to make use of
605 the semantic predictability of speech⁸⁴, this recording was energetically masked by a stream of
606 concatenated five-second sound textures at a signal-to-noise ratio of 0 dB. Sound textures were
607 synthesized from the spectro-temporal modulation content of 192 natural sounds (i.e., human and
608 animal vocalizations, music, tools, nature scenes⁸⁵), so that the noise stream did not provide any
609 semantic content potentially interfering with the prediction of upcoming speech. The order in which
610 sound textures were arranged was randomized across participants. For more details on how sound
611 textures in the present experiment were generated and how they were processed in auditory cortex,
612 see Ref. ⁴⁰.

613 The monaural speech and noise streams were sampled to 44.1 kHz and custom filters specific
614 to the left and right channel of the earphones used in the fMRI experiment were applied for
615 frequency response equalization. Finally, speech-in-noise stimuli were divided into eight excerpts à
616 8 minutes, which served as independent runs in the experiment.

617 A trained human external agent literally transcribed the speech stream. The text transcript
618 comprised 9,446 words, which were used as stimuli in the self-paced reading task and as input to our
619 language models. To automatically determine onset and offset times of all spoken words and
620 phonemes, we used the web service of the Bavarian Archive for Speech Signals (BAS)⁸⁶: First, the text
621 transcript was transformed to a canonical phonetic transcript encoded in SAM-PA by the G2P

622 module. Second, the most likely pronunciation for the phonetic transcript was determined within a
623 Markov model and aligned to the speech recording by the MAUS module. Fourteen part-of-speech
624 tags were assigned to the words in the text transcript using the pre-trained German language model
625 `de_core_news_sm` (2.2.5) from spaCy (<https://spacy.io/>). Based on these tags, words were classified
626 as content or function words. Word frequencies were derived from the subtitle-based SUBTLEX-DE
627 corpus⁸⁷ and transformed to standardized Zipf values⁸⁸ operating on a logarithmic scale from about
628 1 (word with a frequency of 1 per 100 million words) to 7 (1 per 1,000 words). The Zipf value of a
629 word not observed in the corpus was 1.59 (i.e., smallest possible value).

630 **Experimental procedures**

631 **Behavioural self-paced reading task.** While the transcribed story was presented word-by-word on
632 a noncumulative display, participants had the task to read each word once at a comfortable pace
633 and quickly press a button to reveal the next word as soon as they had finished reading. A timeout
634 of 6 seconds was implemented. The time interval between word appearance and button press was
635 logged as the reading time. After each run, participants answered three four-option multiple-choice
636 questions on the plot of the story (performance: $Ra = 58.33$ –100 % correct, $M = 79.17$, $SD = 10.87$)
637 and took a self-paced break. In total, each participant completed four out of eight runs, which were
638 randomly selected and presented in chronological order. Throughout the reading task, we recorded
639 movement and pupil dilation of participants' right eye at a sampling rate of 250 Hz in one continuous
640 shot with an eye tracker (EyeLink 1000, SR Research).

641 The experiment was controlled via the Psychophysics Toolbox⁸⁹ in MATLAB (R2017b,
642 MathWorks). All words were presented 20 % off from the left edge of the screen in white Arial font
643 on a grey background with a visual angle of approximately 18°. Participants used a response pad
644 (URP48, The Black Box Toolkit) to navigate the experiment with their right index finger. The
645 experimental session took approximately 40 minutes.

646 **Functional MRI listening task.** We instructed participants to carefully listen to the story while
647 ignoring the competing stream of sound textures as well as the MRI scanner noise in the background.
648 Each of the eight runs was initialized by 10 baseline MRI volumes after which a white fixation cross
649 appeared in the middle of a grey screen and playback of the 8-minute audio recording started. MRI
650 recording stopped with the end of playback and participants successively answered the same
651 questions used in the self-paced reading task via a response pad with four buttons (HHSC-2x4-C,
652 Current Designs). On average, participants answered 65.5 % of the questions correctly ($Ra = 38$ –100
653 %, $SD = 15.9$ %). There was a 20-second break between consecutive runs.

654 The experiment was run in MATLAB (R2016b) using the Psychophysics Toolbox. Stimuli were
655 presented at a subjectively comfortable sound pressure level via insert earphones (S14,

656 SENSIMETRICS) covered with circumaural air cushions. The experimenters monitored whether
657 participants kept their eyes open throughout the experiment via an eye tracker.

658 **MRI data acquisition.** MRI data were collected on a 3 Tesla Siemens MAGNETOM Skyra scanner
659 using a 64-channel head coil. During the listening task, continuous whole-brain fMRI data were
660 acquired in eight separate runs using an echo-planar imaging (EPI) sequence (repetition time (TR) =
661 947 ms, echo time (TE) = 28 ms, flip angle = 60°, voxel size = 2.5 × 2.5 × 2.5 mm, slice thickness = 2.5
662 mm, matrix size = 80 × 80, field of view = 200 × 200 mm, simultaneous multi-slice factor = 4). Fifty-
663 two axial slices were scanned in interleaved order. For each run, 519 volumes were recorded.

664 Before each second run, field maps were acquired with a gradient echo (GRE) sequence (TR
665 = 610 ms, TE₁ = 4.92 ms, TE₂ = 7.38 ms, flip angle = 60°, voxel size = 2.5 × 2.5 × 2.75 mm, matrix size
666 = 80 × 80, axial slice number = 62, slice thickness = 2.5 mm, slice gap = 10 %).

667 In the end of an experimental session, anatomical images were acquired using a T1-weighted
668 (T1w) MP-RAGE sequence (TR = 2,400 ms, TE = 3.16 ms, flip angle = 8°, voxel size = 1 × 1 × 1 mm,
669 matrix size = 256 × 256, sagittal slice number = 176) and a T2-weighted (T2w) SPACE sequence (TR =
670 3,200 ms, TE = 449 ms, flip angle = 120°, voxel size = 1 × 1 × 1 mm, matrix size = 256 × 256, sagittal
671 slice number = 176).

672 **Modelling the predictiveness of context at multiple timescales**

673 We trained two versions of a long short-term memory network (LSTM) with five layers to predict the
674 next word in a story given a sequence of semantic context: a “continuously updating LSTM” where
675 information is fed to a higher layer with each upcoming word, and a competing “sparsely updating
676 HM-LSTM” where information is fed to a higher layer only at the end of a timescale. The
677 predictiveness of context at multiple timescales was read out from single layers of both language
678 models for each word in the story presented to participants in experiments. Ultimately, we tested
679 how closely these derivatives of different network architectures match signatures of behavioural and
680 neural prediction processes.

681 **Representing words in vector space.** In natural language processing, it is common to represent a
682 word by its linguistic features in the form of high-dimensional vectors (or embeddings). As the
683 German language is morphologically rich and flexibly combines words into new compounds, there
684 are many rare words for which language models cannot learn good (if any) vector representations
685 on the word level. Therefore, we mapped all texts used for training, validating and testing our
686 language models to pre-trained *subword* vectors publicly available in the BPEmb collection⁹⁰. These
687 embeddings allow for the representation of *any* word by a combination of 100-dimensional
688 subwords from a finite vocabulary of 100,000 subwords. We further reduced this vocabulary to those

689 subwords that appeared at least once in any of the texts used for training, validating or testing our
690 language models (i.e., number of subwords in vocabulary $v = 91,645$). See Supplementary Text 1 for
691 a detailed description of the BPEmb vocabulary.

692 Matching our texts to subwords and their respective embeddings in the BPEmb vocabulary,
693 yielded the embedded text $t \in R^{w \times e}$ where w is the number of words and $e = 100$ is the number of
694 vector dimensions. On average, a word in the story was represented by 1.07 subwords ($Ra = 1-6$, SD
695 $= 0.33$). As single words were encoded by only one subword in 94.25 % of cases, we will refer to
696 subwords as words from here on.

697 **Architecture of language models.** When listening to a story, a fused representation of all spoken
698 words $\{w_1, w_2, \dots, w_p\}$ is maintained in memory and used as context information to make a
699 prediction about the upcoming word w_{p+1} . In natural language processing, this memory formation
700 is implemented via recurrent connections between the states of adjacent neural network cells. The
701 hidden state h_{p-1} stores all relevant context and is sequentially passed to the next cell where it is
702 updated with information from word w_p .

703 As such a simple recurrent neural network (RNN) tends to memorize only the most recent
704 past, the more complex LSTM⁴¹ became a standard model in time series forecasting. In an LSTM cell,
705 the state is split in two vectors: The cell state c_p acts as long-term memory, whereas the hidden state
706 h_p incorporates information relevant to the cell output (i.e., the prediction of the next word). The
707 integration of new information and the information flow between the two memory systems is
708 controlled by three gating mechanisms.

709 When stacking multiple LSTM cells on top of each other, semantic context gets hierarchically
710 organized in the model, with lower layers coding for short-term dependencies and higher layers
711 coding for long-term dependencies between words. The bottom-up input to the first layer remains
712 to be the embedded word w_p . However, the lower layer's hidden state h_p^{l-1} becomes the input to a
713 cell from the second layer on. Importantly, the hidden state and cell state are updated at each layer
714 with every new bottom-up input to the model.

715 A competing model that has been shown to slightly outperform the continuously updating
716 ("vanilla") LSTM in character-level language modelling is the hierarchical multiscale LSTM (HM-
717 LSTM)⁴². This model, referred to as "sparsely-updating HM-LSTM", employs a revised updating rule
718 where information from the lower layer is only fed forward at the end of a timescale (i.e., a sequence
719 of words closely related to each other).

720 Importantly, The HM-LSTM allows for a sparse updating rate, with lower layers operating on
721 short timescales and higher layers operating on longer timescales. Here, we used the simplified

722 version of the HM-LSTM⁹¹ with no top-down connections. See Supplementary Text 2 for a detailed
723 description of the model architecture including all relevant formulas.

724 **Prediction of the next word.** LSTM and HM-LSTM cells form the representations of semantic
725 information relevant to speech prediction, whereas the actual prediction of the next word takes
726 place in the output module. Here, hidden states at word position p are combined across the different
727 layers of the language model. The combined hidden state h_p^r is mapped to a fully connected dense
728 layer of as many neurons as there are words in the vocabulary and squashed to values in the interval
729 $[0,1]$, which sum to 1 (i.e., softmax function). Each neuron in resulting vector d_p indexes one
730 particular word in vocabulary v and denotes its probability of being the next word. Finally, the word
731 referring to the highest probability in the distribution is as the predicted next word s_p in a story. See
732 Supplementary Text 3 for a detailed description of word prediction including all relevant formulas.

733 **Training and evaluation of language models.** The objective of our language models was to
734 minimize the difference between the “predicted” probability distribution d_p (i.e., a vector of
735 probabilities ranging from 0 to 1) and the “actual” probability distribution corresponding to the next
736 word in a text (i.e., a vector of zeros with a one-hot encoded target word). To this end, we trained
737 models on mini-batches of 16 independent text sequences à 500 words and monitored model
738 performance by means of categorical cross-entropy between the “predicted” and “actual”
739 probability distribution of each word in a sequence. Based on model performance, trainable
740 parameters were updated after each mini batch using the Adam algorithm for stochastic gradient
741 optimization⁹².

742 Our text corpus comprised more than 130 million words including 4,400 political speeches⁹³
743 as well as 1,200 fictional and popular scientific books. All texts had at least 500 words; metadata,
744 page numbers, references and punctuations (except for hyphenated compound words) were
745 removed from documents. A held-out set of 10 randomly selected documents was used for
746 validation after each epoch of training (i.e., going through the complete training set once) and
747 allowed us to detect overfitting on the training set. Training automatically stopped after model
748 performance did not increase over two epochs for the validation data set.

749 Using a context window of 500 words, we aimed at roughly modelling timescales of the
750 length of common linguistic units in written language (i.e., words, phrases, sentences, and
751 paragraphs). Therefore, we only used a small range of values from three to seven to find the number
752 of layers—intended to represent distinct timescales—best suited to make good predictions.
753 Additionally, we tuned the number of units in LSTM and HM-LSTM cells of language models, using
754 values from 50 to 500 in steps of 50. Hyperparameters were evaluated on a single epoch using grid
755 search and the best combination of hyperparameters was chosen based on performance on the

756 validation set. Our final language models had five LSTM or HM-LSTM layers à 300 units and an output
757 module. The LSTM model had 31,428,745 and the HM-LSTM model had 31,431,570 trainable
758 parameters. Models were trained and evaluated with custom scripts in TensorFlow 2.1⁹⁴. See
759 Supplementary Text 4 for a detailed description of architectural choices.

760 **Deriving the predictiveness of timescales by “lesioning” the language models.** We used each
761 trained language model to determine the predictiveness of semantic context in the story presented
762 to participants in the behavioural and fMRI experiment. First, predictiveness was read out from “full”
763 models: We iteratively selected each word in the story as a target word and fed all 500 context words
764 preceding the target word to our language models. Note that the context for target words in the
765 very beginning of the story comprised less than 500 words. The “predicted” probability of each word
766 in the vocabulary was extracted from distribution d_p in the output module.

767 Second, predictiveness was read out from “lesioned” models, where we allowed information
768 to freely flow through networks, yet only considered semantic context represented at single layers
769 to generate the “predicted” probability distribution. These timescale-resolved probabilities were
770 created by setting weight matrix W_r^l of pre-trained models to zero for all layers of no interest, so that
771 the hidden state of only one layer is passed to the softmax function and all other layers have no
772 bearing on the final prediction. We iteratively set all but one layer to zero with each layer being the
773 only one influencing predictions once, resulting in five lesioned outputs for each language model.

774 We derived three measures of predictiveness from probability distributions. Our primary
775 measure was the degree of surprisal associated with the occurrence of a word given its context. Word
776 surprisal is the negative logarithm of the probability assigned to the actual next word in a story.

777 Secondary measures of predictiveness were used to explore the specificity of the processing
778 hierarchy to only some aspects of prediction processes. Word entropy reflects the amount of
779 uncertainty across the whole probability distribution, which is the negative sum of probabilities
780 multiplied by their natural logarithm. When high probabilities are assigned to only one or few words
781 in the vocabulary, entropy is low. On the other hand, entropy is high when semantic context is not
782 informative enough to narrow predictions down to a limited set of words, resulting in similar
783 probabilities for all candidate words. As all information necessary to determine the entropy of a word
784 is already available to participants before word presentation, entropy of word w_p was ascribed to
785 the previous word w_{p-1} . Whereas word surprisal quantifies the availability of information on the
786 actual next word, word entropy quantifies the overall difficulty of making any definite prediction.
787 Another secondary measure of predictiveness was the relatedness of the predicted next word to the
788 actual next word. This word similarity is expressed as the correlation of respective word embeddings.
789 A high positive Product-moment correlation indicates that the prediction is semantically close to the
790 target word, even though the model prediction might have been incorrect.

791 All three measures were calculated for each word in the story, separately for full models and
792 five lesioned models. This yielded an 18-dimensional feature space of predictiveness for the LSTM as
793 well as the HM-LSTM model, which was linked to BOLD activity and reading times in our analysis.

794 Additionally, we created a metric to dissociate neural effects of predictiveness from more
795 low-level effects of semantic dissimilarity between target words and their preceding context. To this
796 end, we correlated the embedding of each function word in the story with the average embedding
797 of a context window, and subtracted resulting Product-moment correlation coefficients from 1⁹⁵.
798 This measure of contextual dissimilarity was calculated at five timescales corresponding to a
799 logarithmic increase in context length (i.e., 2, 4, 8, 16, and 32 words).

800 To determine temporal integration windows of layers, we scrambled input to language
801 models at nine levels of granularity corresponding to a binary logarithmic increase in the length of
802 intact context (i.e., context windows of 1–256 words). For each layer, we fit linear functions to word
803 surprisal across context windows and extracted slope parameters indicating how much a layer
804 benefits from longer context being available when predicting the next word. On the second level,
805 we fit linear functions to these layer-specific integration windows to determine the context benefit
806 of higher layers over shorter layers. Resulting model-specific slopes were compared to a null
807 distribution of slopes computed by shuffling integration windows across layers (n = 10,000).
808 Additionally, slopes were compared between language models by means of a Monte Carlo
809 approximated permutation test (n = 10,000) on the difference of means.

810 **Convolving features with the hemodynamic response function.** We used three classes of
811 features to model brain responses: 18 features of the predictiveness of timescales (per language
812 model), 3 linguistic features, and 9 acoustic features. While we were primarily interested in modelling
813 effects of predictiveness, linguistic and acoustic features were used as nuisance regressors
814 potentially covarying with predictiveness. Linguistic features included information on when words
815 were presented (coded as 1), whether they were content or function words (coded as 1 and -1), and
816 which frequency they had. In Ref. ⁴⁰, we decomposed the speech-in-noise stimuli into a 288-
817 dimensional acoustic space of spectral, temporal and spectro-temporal modulations, which was
818 derived from a filter bank modelling auditory processing⁹⁶. Here, we reduced the number of acoustic
819 features to the first 9 principal components, which explained more than 80 % of variance in the
820 original acoustic space. All features were z-scored per run.

821 A set of 500 scrambled features of predictiveness was generated, which was used to estimate
822 null distributions of predictive processing. We applied the fast Fourier transform to single features,
823 randomly shifted the phase of frequency components, and inverted the transform to project the
824 data back into the time domain. This preserved power spectra of features but disrupted the temporal

825 alignment of frequencies. See Supplementary Text 5 for a detailed description of convolving features
826 with the hemodynamic response function (HRF).

827 **Data analysis**

828 See Supplementary Text 6 for a detailed description of structural and functional MRI data
829 preprocessing.

830 **Selection of regions of interest.** We hypothesized that the speech prediction hierarchy is
831 represented as a gradient along a temporo-parietal pathway. This rather coarse region of interest
832 was further refined to only include regions implicated in speech processing. To this end, we used
833 intersubject correlation⁴⁵ as a measure of neural activity consistently evoked across participants
834 listening to speech in noise. As we were primarily interested in shared responses to the speech
835 stream, this approach allowed us to leverage the inconsistency of the noise stream across
836 participants. The presentation of sound textures in different order likely evoked more
837 heterogeneous neural responses, leading to a diminished shared representation of the noise stream.
838 Therefore, we inferred that the shared neural responses we observed were largely driven by the
839 speech stream, which was the same for all participants.

840 At the first level, hyperaligned functional time series of each participant (see Supplementary
841 Text 6) were concatenated across experimental runs and correlated with every other participant on
842 a vertex-by-vertex basis, resulting in pairwise maps of intersubject Product-moment correlations. A
843 group map was created by calculating the median correlation coefficient across pairs of participants
844 for each vertex. At the second level, we applied a bootstrap hypothesis test with 10,000 iterations.
845 To create the null distribution, we iteratively resampled participants with replacement and derived
846 median group maps from their pairwise correlation maps. When the same participant was sampled
847 more than once in a bootstrap iteration, the pairwise correlation map of that participant with herself
848 was not included in the computation of the group map. The actual median intersubject correlation
849 was ranked against the normalized null distribution to obtain a p -value for each vertex. Intersubject
850 correlations were computed with the Python package BrainIAK⁹⁷.

851 Finally, we used a multi-modal parcellation⁴⁶ to select those lateral temporal and parietal
852 parcels of which at least 80 % of the vertices had a significant intersubject correlation ($p < 0.01$,
853 adjusted for false discovery rate; FDR)⁹⁸ in one hemisphere. The following parcels were included in
854 the region of interest (ROI): early auditory cortex (EAC), auditory association cortex (AAC), lateral
855 temporal cortex (LTC), temporo-parietal-occipital junction (TPOJ), inferior parietal cortex (IPC), and
856 superior parietal cortex (SPC). As the temporal MT+ complex is thought to be mainly involved in
857 visual processing, this region was not considered an appropriate candidate parcel. All further
858 analyses including MRI data were limited to the temporo-parietal ROI, which was organized along

859 the anterior-posterior (left: 124 mm, right: 167 mm) and inferior-superior axis (left: 234 mm, right:
860 212 mm).

861 **Functional data analysis.** The starting point of our analyses was the question whether the
862 timescales of speech prediction organize along a temporo-parietal processing hierarchy. In a forward
863 model, we encoded the predictiveness of timescales into univariate neural activity and fit a gradient
864 along the peak locations sensitive to specific timescales of surprisal. Next, we compared the
865 explanatory power of both language models in a backward model, which decoded surprisal at
866 different timescales from multivariate patterns of neural activity in temporo-parietal parcels. Finally,
867 we modelled functional connectivity between peak locations to test whether the timescales of
868 surprisal gate the information flow along the gradient.

869 **Encoding model.** The encoding approach (similar to e.g., Ref. ⁹⁹) allowed us to quantify for
870 each temporo-parietal vertex, which features of predictiveness it preferentially represents. Two
871 separate encoding models were estimated for each vertex in the ROI of single participants, one for
872 each language model. Besides the features of predictiveness specific to language models, both
873 models included the same linguistic and acoustic features as nuisance regressors. We modelled
874 neural activity as a function of 30 HRF-convolved features characterizing speech and noise stimuli
875 by:

$$876 \quad a = Sw + \epsilon,$$

877 where $a^{samples \times 1}$ is the activity vector (or BOLD time course) corresponding to a vertex,
878 $S^{samples \times features}$ is the stimulus matrix of features, $w^{features \times 1}$ is a vector of estimated model
879 weights, and $\epsilon^{samples \times 1}$ is a vector of random noise.

880 All models were estimated using ridge regression with four-fold cross validation. We paired
881 odd-numbered functional runs with their subsequent even-numbered run, resulting in four data
882 splits per participant. Each of the four data splits was selected as a testing set once; all other data
883 splits were used as a training set. Within each fold, generalized cross-validation¹⁰⁰ was carried out on
884 the training set to find an optimal estimate of regularization parameter λ from the data, searching
885 100 values evenly spaced on a logarithmic scale from 10^{-5} to 10^8 . Weights of predictiveness were
886 extracted from the model fit with the optimal regularization parameter and averaged across cross-
887 validation folds to obtain stable weights.

888 To evaluate the performance of encoding models and their ability to generalize to new data,
889 we applied the weights estimated on the training set to the features of the held-out testing set in
890 each cross-validation fold. The predicted BOLD time series was correlated with the actual BOLD time
891 series. The resulting Product-moment correlation coefficient is the encoding accuracy, which was
892 averaged across cross-validation folds and Fisher z-transformed.

893 Additionally, we created null distributions of weights and encoding accuracies by estimating
894 forward models on scrambled features of predictiveness (similar to e.g., Ref. ¹⁰¹). We set up 500
895 separate models, which included scrambled features of predictiveness but intact linguistic and
896 acoustic features. Models were estimated largely following the cross-validation scheme outlined for
897 observed data. However, we re-used optimal regularization parameters from non-scrambled models
898 of corresponding folds. All ridge regression models were implemented using the RidgeCV function
899 in the Python package scikit-learn¹⁰².

900 **Peak selection.** For both language models, we derived five temporo-parietal maps in the left
901 and right hemisphere of single participants: one weight map for each timescale of word surprisal.
902 Maps represented the sensitivity of brain regions to timescale surprisal; positive weights indicate
903 increasing BOLD activity to more surprising words.

904 To illustrate the location and extent of brain regions modulated by timescale surprisal, we
905 performed an analysis similar to cluster-based permutation tests in Fieldtrip¹⁰³. For each timescale,
906 vertex-wise weights observed across participants were tested against zero by means of a one-sample
907 *t*-test. We combined a vertex into a cluster with its adjacent vertices if it was significant at an alpha
908 level of 0.05 and had at least two significant neighbours. We clustered vertices with negative *t*-values
909 separately from vertices with positive *t*-values. The summed *t*-value of an observed cluster served as
910 the cluster-level statistic and was compared with a Monte Carlo approximated null distribution of
911 summed *t*-values. This null distribution was created by performing clustering on scrambled
912 partitions of timescale-specific weight maps and selecting the largest summed *t*-value for each
913 partition. An observed cluster was considered significant if its summed *t*-value was exceeded by no
914 more than 2.5% of summed *t*-values from scrambled partitions.

915 Beyond this rather coarse mapping of temporo-parietal brain regions onto the timescales of
916 surprisal, our main analysis focused on how timescale-specific peak locations distribute along the
917 inferior-superior axis only. Of note, we hypothesized that a hierarchy of speech prediction evolves
918 from temporal to parietal areas, which corresponds to the inferior-superior axis of our ROI. A window
919 with a height of 2 mm was shifted along the inferior-superior axis of the temporo-parietal ROI in
920 steps of 1 mm. All weights of a timescale falling into the window were averaged, thereby collapsing
921 across the anterior-posterior axis. The resulting one-dimensional weight profile of a timescale
922 spanned inferior to superior locations and was smoothed using robust linear regression over a
923 window of 70 mm. For each unilateral weight profile of single participants, local maxima (i.e., a
924 sample larger than its two neighbouring samples) were determined.

925 We applied two different approaches to select one peak location for each timescale from
926 these local maxima. In the naïve approach of peak selection, the local maximum with the highest
927 positive value was defined as a peak. As this approach makes it hard to find a consistent order of
928 timescales when surprisal is not just processed along the dorsal but also the ventral processing

929 stream, we also applied a pre-informed peak selection approach explicitly targeting the dorsal
930 stream. Here, the peak of the first timescale had to be in the inferior half of the axis (i.e., temporal
931 regions) and peaks of longer timescales had to be superior to the peak of the first timescale.
932 Whenever no timescale peak could be defined, the largest positive value was selected. Both peak
933 selection approaches yielded five timescale-specific coordinates on the inferior-superior axis for
934 each participant, hemisphere and language model.

935 Additionally, we applied naïve peak selection to weight maps, which were rotated by -45°
936 before collapsing across the first dimension. In the left hemisphere, an increase on that new axis
937 indicated a shift to more superior and posterior regions, thereby simultaneously modelling effects
938 on the inferior-superior and anterior-posterior axis. However, original right-hemispheric maps
939 already had a strong rotation off the inferior-superior axis, thus rotating these maps rather brought
940 them into alignment with the inferior-superior axis of non-rotated left-hemispheric maps.

941 **Gradient fitting.** We fit linear functions to coordinates of single participants across the
942 timescales of surprisal. Models included an intercept term and the slope parameter was extracted
943 from each fit. A positive slope indicates a gradient of timescale surprisal, where shorter timescales
944 are represented in more inferior (anterior) temporo-parietal regions than longer timescale, which are
945 represented in more superior (posterior) regions. We tested grand-average slope parameters against
946 a null distribution of slopes with 10,000 partitions, which was created by randomly shuffling the
947 coordinates of single participants across the timescales of surprisal and recalculating their slopes. As
948 the first timescale was pre-set to have the most inferior coordinate in the pre-informed peak
949 selection approach, this specific coordinate was not shuffled when calculating the null distribution
950 for this approach. To compare slope parameters between language models, we performed a Monte
951 Carlo approximated permutation test ($n = 10,000$), using the difference of means as a test statistic.
952 As secondary analyses, gradients of predictive processing were also calculated for the timescales of
953 word entropy and similarity. In a control analysis, a gradient was fit to timescale peaks following the
954 same procedure described above but replacing features of predictiveness by contextual dissimilarity
955 when estimating forward models.

956 To round off the encoding analysis, we compared temporo-parietal encoding accuracies
957 between both language models. As we were interested in effects specific to the predictiveness of
958 speech, encoding accuracies were z-scored to the null distribution of accuracies from scrambled
959 features of predictiveness. A cluster-based permutation paired-sample t -test was calculated ($n =$
960 $1,000$, vertex-specific alpha level: 0.05, cluster-specific alpha level: 0.05). In comparison to the cluster
961 test described above for the weight maps, we here created a null distribution of summed t -values by
962 contrasting accuracies of language models whose labels had been randomly shuffled in single
963 participants.

964 **Decoding model.** In our decoding approach (similar to e.g., Ref. ⁹⁹), we quantified how much
965 information multiple vertices jointly contain about a feature of predictiveness. For each language
966 model, five separate backward models were estimated in each of six temporo-parietal parcels of
967 single participants, one for each timescale of word surprisal. We modelled timescale-specific word
968 surprisal as a function of neural activity in all vertices forming a parcel by:

$$969 \quad s = Aw + \epsilon,$$

970 where $s^{samples \times 1}$ is the stimulus vector of a feature, $A^{samples \times vertices}$ is the activity matrix of BOLD
971 time courses corresponding to the vertices of a parcels, $w^{vertices \times 1}$ is a vector of model weights,
972 $\epsilon^{samples \times 1}$ is a vector of random noise.

973 The same cross-validation scheme as described for the encoding model was applied.
974 However, instead of predicting BOLD activity, we here reconstructed surprisal at different timescales.
975 By correlating the actual stimulus time series with the one predicted on the held-out testing set, we
976 obtained the decoding accuracy of a parcel. Decoding accuracies were z-scored to the null
977 distribution of accuracies determined for scrambled stimulus time series. We compared decoding
978 accuracies between language models in each hemisphere, parcel and timescale by means of a Monte
979 Carlo approximated permutation test ($n = 10,000$) on the difference of means. Resulting p -values
980 were corrected for multiple comparisons using FDR correction.

981 **Functional connectivity.** To model the information flow between brain regions sensitive to
982 the different timescales of word surprisal, we identified five unique seeds for both language models
983 in each temporo-parietal hemisphere. On the inferior-superior axis, we re-used the grand-median
984 coordinate of each timescale as localized in the pre-informed peak selection. The corresponding
985 coordinate on the anterior-posterior axis was localized by shifting a moving average with a window
986 centred on the inferior-superior coordinate along the anterior-posterior axis (width: 2 mm, height: 5
987 mm), and determining peak locations on smoothed weight profiles of single participants. Following,
988 we placed a sphere with a radius of 5 mm on peak coordinates from both axes and averaged BOLD
989 time courses of vertices falling within this sphere, yielding the timescale-specific neural activity of
990 seeds.

991 We expected increased information flow between seeds of adjacent timescales when one
992 timescale becomes uninformative for the prediction of upcoming speech. This modulatory influence
993 of surprisal on connectivity was modelled along the lines of a psychophysiological interaction
994 (PPI)¹⁰⁴. In a standard PPI analysis, the neural time series of one brain region is regressed onto the
995 pointwise product of an experimental stimulus and the neural time series of another brain region.
996 Here, we extended this approach by creating timescale-specific interactions: BOLD time series of
997 seeds were multiplied by their corresponding HRF-convolved surprisal time series but not any of the
998 surprisal time series at another timescale.

999 Functional connectivity was calculated for both language models in each participant and
1000 hemisphere. We set up five regression models, with every seed being selected as a target once. The
1001 physiological (BOLD) time series of the target seed was mapped onto the physiological,
1002 psychological (timescale-specific surprisal), and psychophysiological time series from all other
1003 (predictor) seeds. Models were estimated within the same cross-validation scheme outlined for the
1004 encoding model. We extracted all four weights from psychophysiological interaction terms of each
1005 target seed and arranged weights in a 5-by-5 matrix, with target seeds on the main diagonal and
1006 predictor seeds off the diagonal. This matrix of observed psychophysiological interactions was
1007 compared to a matrix with hypothesized interaction weights: The diagonals below and above the
1008 main diagonal were set to 1 (indicating increased coupling when surprisal at a neighbouring
1009 timescale is high), all other items were set to -1. We calculated the Euclidean distance of single-
1010 participant matrices to this hypothesized matrix. The mean of observed Euclidean distances was
1011 compared to a null distribution of 10,000 mean Euclidean distances calculated on BOLD time series
1012 of target seeds randomly shifted in time by the number of samples in 1 to 7 functional runs. Euclidean
1013 distances were compared between language models in each hemisphere by means of a Monte Carlo
1014 approximated permutation test ($n = 10,000$) on the difference of means.

1015 **Behavioural data analysis.** Reading times were used to test the behavioural relevance of the
1016 predictiveness determined by our language models. Trials with reading times shorter than 0.001
1017 seconds or longer than 6 seconds were considered invalid and excluded. Further, we converted
1018 reading times to speed (number of words per 100 seconds) and excluded trials exceeding 3 standard
1019 deviations within a run and participant from all further analyses. On average, 1.31 % of trials ($SD =$
1020 $1.12, Ra = 0.32-6.15$) were removed. Finally, reading speed was z-scored within runs.

1021 For each participant, we predicted reading speed in a forward model, adopting the same
1022 cross-validated ridge regression scheme used for the analysis of fMRI data. Our feature space
1023 included the predictiveness of words as well as word frequency, word length (number of letters),
1024 content vs. function words and trial number as nuisance regressors. As this was a high-pace task,
1025 some features might have unfolded their effect on reading speed only over the course of a few
1026 words. Therefore, we added time-lagged versions of features to the model, that is, shifting features
1027 by -2 to 5 word positions. There were no lagged versions of the predictor coding for trial number
1028 added to the model.

1029 To investigate whether predictiveness had an effect on reading speed beyond the effect of
1030 nuisance regressors, we compared the predictive accuracy of forward models in single participants
1031 to a null distribution of accuracies from models with scrambled features of predictiveness. The
1032 performance of language models was compared by z-scoring observed encoding accuracies to the

1033 null distribution and running a Monte Carlo approximated permutation test ($n = 10,000$) on the
1034 difference of means. This analysis was also carried out for the timescales of contextual dissimilarity.

1035 **References**

- 1036 1. Clark, A. Whatever next? Predictive brains, situated agents, and the future of cognitive
1037 science. *Behav. Brain Sci.* **36**, 181–204 (2013).
- 1038 2. Murray, J. D. *et al.* A hierarchy of intrinsic timescales across primate cortex. *Nat. Neurosci.* **17**,
1039 1661–1663 (2014).
- 1040 3. Zhang, W. & Yartsev, M. M. Correlated Neural Activity across the Brains of Socially Interacting
1041 Bats. *Cell* **178**, 413–428 (2019).
- 1042 4. La Camera, G. *et al.* Multiple Time Scales of Temporal Response in Pyramidal and Fast Spiking
1043 Cortical Neurons. *J. Neurophysiol.* **96**, 3448–3464 (2006).
- 1044 5. Burt, J. B. *et al.* Hierarchy of transcriptomic specialization across human cortex captured by
1045 structural neuroimaging topography. *Nat. Neurosci.* **21**, 1251–1259 (2018).
- 1046 6. Lakatos, P. *et al.* An Oscillatory Hierarchy Controlling Neuronal Excitability and Stimulus
1047 Processing in the Auditory Cortex. *J. Neurophysiol.* **94**, 1904–1911 (2005).
- 1048 7. Mattar, M. G., Kahn, D. A., Thompson-Schill, S. L. & Aguirre, G. K. Varying Timescales of
1049 Stimulus Integration Unite Neural Adaptation and Prototype Formation. *Curr. Biol.* **26**, 1669–
1050 1676 (2016).
- 1051 8. Lamme, V. A. F. & Roelfsema, P. R. The distinct modes of vision offered by feedforward and
1052 recurrent processing. *Trends Neurosci.* **23**, 571–579 (2000).
- 1053 9. Hasson, U., Chen, J. & Honey, C. J. Hierarchical process memory: memory as an integral
1054 component of information processing. *Trends Cogn. Sci.* **19**, 304–313 (2015).
- 1055 10. Buračas, G. T., Zador, A. M., DeWeese, M. R. & Albright, T. D. Efficient Discrimination of
1056 Temporal Patterns by Motion-Sensitive Neurons in Primate Visual Cortex. *Neuron* **20**, 959–969
1057 (1998).
- 1058 11. Runyan, C. A., Piasini, E., Panzeri, S. & Harvey, C. D. Distinct timescales of population coding
1059 across cortex. *Nature* **548**, 92–96 (2017).
- 1060 12. Friston, K. A theory of cortical responses. *Philos. Trans. R. Soc. B Biol. Sci.* **360**, 815–836 (2005).

- 1061 13. Keller, G. B. & Mrsic-Flogel, T. D. Predictive Processing: A Canonical Cortical Computation.
1062 *Neuron* **100**, 424–435 (2018).
- 1063 14. Kiebel, S. J., Daunizeau, J. & Friston, K. J. A Hierarchy of Time-Scales and the Brain. *PLoS*
1064 *Comput. Biol.* **4**, e1000209 (2008).
- 1065 15. Chaudhuri, R., Knoblauch, K., Gariel, M.-A., Kennedy, H. & Wang, X.-J. A Large-Scale Circuit
1066 Mechanism for Hierarchical Dynamical Processing in the Primate Cortex. *Neuron* **88**, 419–431
1067 (2015).
- 1068 16. Demirtaş, M. *et al.* Hierarchical Heterogeneity across Human Cortex Shapes Large-Scale
1069 Neural Dynamics. *Neuron* **101**, 1181–1194.e13 (2019).
- 1070 17. Bastos, A. M. *et al.* Visual Areas Exert Feedforward and Feedback Influences through Distinct
1071 Frequency Channels. *Neuron* **85**, 390–401 (2015).
- 1072 18. Cocchi, L. *et al.* A hierarchy of timescales explains distinct effects of local inhibition of primary
1073 visual cortex and frontal eye fields. *eLife* **5**, e15252 (2016).
- 1074 19. Wacongne, C. *et al.* Evidence for a hierarchy of predictions and prediction errors in human
1075 cortex. *Proc. Natl. Acad. Sci.* **108**, 20754–20759 (2011).
- 1076 20. Schwiedrzik, C. M. & Freiwald, W. A. High-Level Prediction Signals in a Low-Level Area of the
1077 Macaque Face-Processing Hierarchy. *Neuron* **96**, 89–97.e4 (2017).
- 1078 21. Donhauser, P. W. & Baillet, S. Two Distinct Neural Timescales for Predictive Speech Processing.
1079 *Neuron* **105**, 385–393.e9 (2020).
- 1080 22. Chao, Z. C., Takaura, K., Wang, L., Fujii, N. & Dehaene, S. Large-Scale Cortical Networks for
1081 Hierarchical Prediction and Prediction Error in the Primate Brain. *Neuron* **100**, 1252–1266.e3
1082 (2018).
- 1083 23. Honey, C. J. *et al.* Slow Cortical Dynamics and the Accumulation of Information over Long
1084 Timescales. *Neuron* **76**, 423–434 (2012).
- 1085 24. Stephens, G. J., Honey, C. J. & Hasson, U. A place for time: the spatiotemporal structure of
1086 neural dynamics during natural audition. *J. Neurophysiol.* **110**, 2019–2026 (2013).

- 1087 25. Mesgarani, N., Cheung, C., Johnson, K. & Chang, E. F. Phonetic Feature Encoding in Human
1088 Superior Temporal Gyrus. *Science* **343**, 1006–1010 (2014).
- 1089 26. Huth, A. G., de Heer, W. A., Griffiths, T. L., Theunissen, F. E. & Gallant, J. L. Natural speech
1090 reveals the semantic maps that tile human cerebral cortex. *Nature* **532**, 453–458 (2016).
- 1091 27. de Heer, W. A., Huth, A. G., Griffiths, T. L., Gallant, J. L. & Theunissen, F. E. The Hierarchical
1092 Cortical Organization of Human Speech Processing. *J. Neurosci.* **37**, 6539–6557 (2017).
- 1093 28. Lerner, Y., Honey, C. J., Silbert, L. J. & Hasson, U. Topographic Mapping of a Hierarchy of
1094 Temporal Receptive Windows Using a Narrated Story. *J. Neurosci.* **31**, 2906–2915 (2011).
- 1095 29. Bornkessel-Schlesewsky, I., Schlesewsky, M., Small, S. L. & Rauschecker, J. P. Neurobiological
1096 roots of language in primate audition: common computational properties. *Trends Cogn. Sci.*
1097 **19**, 142–150 (2015).
- 1098 30. Kuperberg, G. R. & Jaeger, T. F. What do we mean by prediction in language comprehension?
1099 *Lang. Cogn. Neurosci.* **31**, 32–59 (2016).
- 1100 31. Arnal, L. H., Wyart, V. & Giraud, A.-L. Transitions in neural oscillations reflect prediction errors
1101 generated in audiovisual speech. *Nat. Neurosci.* **14**, 797–801 (2011).
- 1102 32. Blank, H. & Davis, M. H. Prediction Errors but Not Sharpened Signals Simulate Multivoxel fMRI
1103 Patterns during Speech Perception. *PLOS Biol.* **14**, e1002577 (2016).
- 1104 33. Kandylaki, K. D. *et al.* Predicting ‘When’ in Discourse Engages the Human Dorsal Auditory
1105 Stream: An fMRI Study Using Naturalistic Stories. *J. Neurosci.* **36**, 12180–12191 (2016).
- 1106 34. Chien, H.-Y. S. & Honey, C. J. Constructing and Forgetting Temporal Context in the Human
1107 Cerebral Cortex. *Neuron* **106**, 675–686.e11 (2020).
- 1108 35. Zadbood, A., Chen, J., Leong, Y. C., Norman, K. A. & Hasson, U. How We Transmit Memories to
1109 Other Brains: Constructing Shared Neural Representations Via Communication. *Cereb. Cortex*
1110 **27**, 4988–5000 (2017).
- 1111 36. Baldassano, C. *et al.* Discovering Event Structure in Continuous Narrative Perception and
1112 Memory. *Neuron* **95**, 709–721.e5 (2017).

- 1113 37. Hamilton, L. S. & Huth, A. G. The revolution will not be controlled: natural stimuli in speech
1114 neuroscience. *Lang. Cogn. Neurosci.* 1–10 (2018).
- 1115 38. Cohen, J. D. *et al.* Computational approaches to fMRI analysis. *Nat. Neurosci.* **20**, 304–313
1116 (2017).
- 1117 39. Cichy, R. M. & Kaiser, D. Deep Neural Networks as Scientific Models. *Trends Cogn. Sci.* **23**, 305–
1118 317 (2019).
- 1119 40. Erb, J., Schmitt, L.-M. & Obleser, J. Temporal selectivity declines in the aging human auditory
1120 cortex. *eLife* **9**, e55300 (2020).
- 1121 41. Hochreiter, S. & Schmidhuber, J. Long short-term memory. (1997).
- 1122 42. Chung, J., Ahn, S. & Bengio, Y. Hierarchical Multiscale Recurrent Neural Networks.
1123 *ArXiv160901704 Cs* (2016).
- 1124 43. Hale, J. A probabilistic earley parser as a psycholinguistic model. in *Proceedings of the North*
1125 *American association of computational linguistics* 159–166 (Association for Computational
1126 Linguistics, 2001).
- 1127 44. Levy, R. Expectation-based syntactic comprehension. *Cognition* **106**, 1126–1177 (2008).
- 1128 45. Nastase, S. A., Gazzola, V., Hasson, U. & Keysers, C. Measuring shared responses across subjects
1129 using intersubject correlation. *Soc. Cogn. Affect. Neurosci.* nsz037 (2019).
- 1130 46. Glasser, M. F. *et al.* A multi-modal parcellation of human cerebral cortex. *Nature* **536**, 171–178
1131 (2016).
- 1132 47. Boldt, R. *et al.* Listening to an Audio Drama Activates Two Processing Networks, One for All
1133 Sounds, Another Exclusively for Speech. *PLoS ONE* **8**, e64489 (2013).
- 1134 48. Schmäzle, R., Häcker, F. E. K., Honey, C. J. & Hasson, U. Engaged listeners: shared neural
1135 processing of powerful political speeches. *Soc. Cogn. Affect. Neurosci.* **10**, 1137–1143 (2015).
- 1136 49. Regev, M. *et al.* Propagation of Information Along the Cortical Hierarchy as a Function of
1137 Attention While Reading and Listening to Stories. *Cereb. Cortex* **29**, 4017–4034 (2018).
- 1138 50. Hasson, U., Yang, E., Vallines, I., Heeger, D. J. & Rubin, N. A Hierarchy of Temporal Receptive
1139 Windows in Human Cortex. *J. Neurosci.* **28**, 2539–2550 (2008).

- 1140 51. Zacks, J. M., Speer, N. K., Swallow, K. M., Braver, T. S. & Reynolds, J. R. Event perception: A mind-
1141 brain perspective. *Psychol. Bull.* **133**, 273–293 (2007).
- 1142 52. Radvansky, G. A. Across the Event Horizon. *Curr. Dir. Psychol. Sci.* **21**, 269–272 (2012).
- 1143 53. Zacks, J. M. *et al.* Human brain activity time-locked to perceptual event boundaries. *Nat.*
1144 *Neurosci.* **4**, 651–655 (2001).
- 1145 54. Ditman, T., Holcomb, P. J. & Kuperberg, G. R. Time travel through language: Temporal shifts
1146 rapidly decrease information accessibility during reading. *Psychon. Bull. Rev.* **15**, 750–756
1147 (2008).
- 1148 55. Whitney, C. *et al.* Neural correlates of narrative shifts during auditory story comprehension.
1149 *NeuroImage* **47**, 360–366 (2009).
- 1150 56. Richmond, L. L. & Zacks, J. M. Constructing Experience: Event Models from Perception to
1151 Action. *Trends Cogn. Sci.* **21**, 962–980 (2017).
- 1152 57. Lieder, F., Stephan, K. E., Daunizeau, J., Garrido, M. I. & Friston, K. J. A Neurocomputational
1153 Model of the Mismatch Negativity. *PLoS Comput. Biol.* **9**, e1003288 (2013).
- 1154 58. Kumar, S., Kaposvari, P. & Vogels, R. Encoding of Predictable and Unpredictable Stimuli by
1155 Inferior Temporal Cortical Neurons. *J. Cogn. Neurosci.* **29**, 1445–1454 (2017).
- 1156 59. Todorovic, A. & de Lange, F. P. Repetition Suppression and Expectation Suppression Are
1157 Dissociable in Time in Early Auditory Evoked Fields. *J. Neurosci.* **32**, 13389–13395 (2012).
- 1158 60. Brown, C. & Hagoort, P. The Processing Nature of the N400: Evidence from Masked Priming. *J.*
1159 *Cogn. Neurosci.* **5**, 34–44 (1993).
- 1160 61. Hagoort, P., Baggio, G. & Willems, R. M. Semantic unification. in *The cognitive neurosciences*
1161 819–836 (MIT Press, 2009).
- 1162 62. Lee, T. S. & Mumford, D. Hierarchical Bayesian inference in the visual cortex. *J. Opt. Soc. Am. A*
1163 **20**, 1434 (2003).
- 1164 63. Kok, P., Jehee, J. F. M. & de Lange, F. P. Less Is More: Expectation Sharpens Representations in
1165 the Primary Visual Cortex. *Neuron* **75**, 265–270 (2012).

- 1166 64. Rao, R. P. N. & Ballard, D. H. Predictive coding in the visual cortex: a functional interpretation
1167 of some extra-classical receptive-field effects. *Nat. Neurosci.* **2**, 79–87 (1999).
- 1168 65. Bell, A. H., Summerfield, C., Morin, E. L., Malecek, N. J. & Ungerleider, L. G. Encoding of Stimulus
1169 Probability in Macaque Inferior Temporal Cortex. *Curr. Biol.* **26**, 2280–2290 (2016).
- 1170 66. Hickok, G. & Poeppel, D. Dorsal and ventral streams: a framework for understanding aspects
1171 of the functional anatomy of language. *Cognition* **92**, 67–99 (2004).
- 1172 67. Hickok, G. & Poeppel, D. The cortical organization of speech processing. *Nat. Rev. Neurosci.* **8**,
1173 393–402 (2007).
- 1174 68. DeWitt, I. & Rauschecker, J. P. Phoneme and word recognition in the auditory ventral stream.
1175 *Proc. Natl. Acad. Sci.* **109**, E505–E514 (2012).
- 1176 69. Bornkessel-Schlesewsky, I. & Schlewsky, M. Reconciling time, space and function: A new
1177 dorsal–ventral stream model of sentence comprehension. *Brain Lang.* **125**, 60–76 (2013).
- 1178 70. Wilson, S. M., Molnar-Szakacs, I. & Iacoboni, M. Beyond Superior Temporal Cortex: Intersubject
1179 Correlations in Narrative Speech Comprehension. *Cereb. Cortex* **18**, 230–242 (2008).
- 1180 71. Garrido, M. I., Rowe, E. G., Halász, V. & Mattingley, J. B. Bayesian Mapping Reveals That
1181 Attention Boosts Neural Responses to Predicted and Unpredicted Stimuli. *Cereb. Cortex* **28**,
1182 1771–1782 (2018).
- 1183 72. Phillips, H. N. *et al.* Convergent evidence for hierarchical prediction networks from human
1184 electrocorticography and magnetoencephalography. *Cortex* **82**, 192–205 (2016).
- 1185 73. Meyniel, F. & Dehaene, S. Brain networks for confidence weighting and hierarchical inference
1186 during probabilistic learning. *Proc. Natl. Acad. Sci.* **114**, E3859–E3868 (2017).
- 1187 74. Cheung, V. K. M., Meyer, L., Friederici, A. D. & Koelsch, S. The right inferior frontal gyrus
1188 processes nested non-local dependencies in music. *Sci. Rep.* **8**, (2018).
- 1189 75. Bottini, R. & Doeller, C. F. Knowledge Across Reference Frames: Cognitive Maps and Image
1190 Spaces. *Trends Cogn. Sci.* **24**, 606–619 (2020).
- 1191 76. Spiers, H. J., Hayman, R. M. A., Jovalekic, A., Marozzi, E. & Jeffery, K. J. Place Field Repetition and
1192 Purely Local Remapping in a Multicompartment Environment. *Cereb. Cortex* **25**, 10–25 (2015).

- 1193 77. Brunec, I. K., Moscovitch, M. & Barense, M. D. Boundaries Shape Cognitive Representations of
1194 Spaces and Events. *Trends Cogn. Sci.* **22**, 637–650 (2018).
- 1195 78. Alexander, A. S. & Nitz, D. A. Spatially Periodic Activation Patterns of Retrosplenial Cortex
1196 Encode Route Sub-spaces and Distance Traveled. *Curr. Biol.* **27**, 1551-1560.e4 (2017).
- 1197 79. Johnson, A. & Redish, A. D. Neural Ensembles in CA3 Transiently Encode Paths Forward of the
1198 Animal at a Decision Point. *J. Neurosci.* **27**, 12176–12189 (2007).
- 1199 80. Stachenfeld, K. L., Botvinick, M. M. & Gershman, S. J. The hippocampus as a predictive map.
1200 *Nat. Neurosci.* **20**, 1643–1653 (2017).
- 1201 81. Kurby, C. A. & Zacks, J. M. Starting from scratch and building brick by brick in comprehension.
1202 *Mem. Cognit.* **40**, 812–826 (2012).
- 1203 82. Sainburg, T., Theilman, B., Thielk, M. & Gentner, T. Q. Parallels in the sequential organization of
1204 birdsong and human speech. *Nat. Commun.* **10**, 3636 (2019).
- 1205 83. Power, J. D., Barnes, K. A., Snyder, A. Z., Schlaggar, B. L. & Petersen, S. E. Spurious but
1206 systematic correlations in functional connectivity MRI networks arise from subject motion.
1207 *NeuroImage* **59**, 2142–2154 (2012).
- 1208 84. Rysop, A. U., Schmitt, L.-M., Obleser, J. & Hartwigsen, G. Neural modelling of the semantic
1209 predictability gain under challenging listening conditions. *Hum. Brain Mapp.* **42**, 110–127
1210 (2021).
- 1211 85. McDermott, J. H. & Simoncelli, E. P. Sound Texture Perception via Statistics of the Auditory
1212 Periphery: Evidence from Sound Synthesis. *Neuron* **71**, 926–940 (2011).
- 1213 86. Kisler, T., Reichel, U. & Schiel, F. Multilingual processing of speech via web services. *Comput.*
1214 *Speech Lang.* **45**, 326–347 (2017).
- 1215 87. Brysbaert, M. *et al.* The Word Frequency Effect: A Review of Recent Developments and
1216 Implications for the Choice of Frequency Estimates in German. *Exp. Psychol.* **58**, 412–424
1217 (2011).

- 1218 88. van Heuven, W. J. B., Mandera, P., Keuleers, E. & Brysbaert, M. Subtlex-UK: A New and
1219 Improved Word Frequency Database for British English. *Q. J. Exp. Psychol.* **67**, 1176–1190
1220 (2014).
- 1221 89. Brainard, D. H. The Psychophysics Toolbox. *Spat. Vis.* **10**, 433–436 (1997).
- 1222 90. Heinzerling, B. & Strube, M. BPEmb: Tokenization-free Pre-trained Subword Embeddings in
1223 275 Languages. *ArXiv171002187 Cs* (2017).
- 1224 91. Kádár, Á., Côté, M.-A., Chrupała, G. & Alishahi, A. Revisiting the Hierarchical Multiscale LSTM.
1225 *ArXiv180703595 Cs* (2018).
- 1226 92. Kingma, D. P. & Ba, J. Adam: A Method for Stochastic Optimization. *ArXiv14126980 Cs* (2017).
- 1227 93. Barbaresi, A. A corpus of German political speeches from the 21st century. *11th Lang. Resour.*
1228 *Eval. Conf.* 792–797 (2018).
- 1229 94. Abadi, M. *et al.* TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed
1230 Systems. (2015).
- 1231 95. Broderick, M. P., Anderson, A. J., Di Liberto, G. M., Crosse, M. J. & Lalor, E. C.
1232 Electrophysiological Correlates of Semantic Dissimilarity Reflect the Comprehension of
1233 Natural, Narrative Speech. *Curr. Biol.* **28**, 803–809.e3 (2018).
- 1234 96. Chi, T., Ru, P. & Shamma, S. A. Multiresolution spectrotemporal analysis of complex sounds. *J.*
1235 *Acoust. Soc. Am.* **118**, 887–906 (2005).
- 1236 97. Kumar, M. *et al.* BrainIAK tutorials: User-friendly learning materials for advanced fMRI analysis.
1237 *PLOS Comput. Biol.* **16**, e1007549 (2020).
- 1238 98. Benjamini, Y. & Hochberg, Y. Controlling the False Discovery Rate: A Practical and Powerful
1239 Approach to Multiple Testing. *J. R. Stat. Soc. Ser. B Methodol.* **57**, 289–300 (1995).
- 1240 99. Santoro, R. *et al.* Reconstructing the spectrotemporal modulations of real-life sounds from
1241 fMRI response patterns. *Proc. Natl. Acad. Sci.* **114**, 4799–4804 (2017).
- 1242 100. Golub, G. H., Heath, M. & Wahba, G. Generalized Cross-Validation as a Method for Choosing a
1243 Good Ridge Parameter. *Technometrics* **21**, 215–223 (1979).

- 1244 101. Musall, S., Kaufman, M. T., Juavinett, A. L., Gluf, S. & Churchland, A. K. Single-trial neural
1245 dynamics are dominated by richly varied movements. *Nat. Neurosci.* **22**, 1677–1686 (2019).
- 1246 102. Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. *J. Mach. Learn. Res.* **12**, 2825–
1247 2830 (2011).
- 1248 103. Maris, E. & Oostenveld, R. Nonparametric statistical testing of EEG- and MEG-data. *J. Neurosci.*
1249 *Methods* **164**, 177–190 (2007).
- 1250 104. Friston, K. J. *et al.* Psychophysiological and Modulatory Interactions in Neuroimaging.
1251 *NeuroImage* **6**, 218–229 (1997).