



HAL
open science

A deep hierarchy of predictions enables assignment of semantic roles in online speech comprehension

Yaqing Su, Lucy J Macgregor, Itsaso Olasagasti, Anne-Lise Giraud Mamessier

► To cite this version:

Yaqing Su, Lucy J Macgregor, Itsaso Olasagasti, Anne-Lise Giraud Mamessier. A deep hierarchy of predictions enables assignment of semantic roles in online speech comprehension. 2022. hal-03995162v1

HAL Id: hal-03995162

<https://hal.science/hal-03995162v1>

Preprint submitted on 17 Feb 2023 (v1), last revised 8 Jul 2024 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0 International License

1 A deep hierarchy of predictions enables assignment of
2 semantic roles in online speech comprehension

3

4

5 Yaqing Su^{1,2,*}, Lucy J. MacGregor³, Itsaso Olasagasti^{1,2}, Anne-Lise Giraud^{1,2,4}

6

7

8 ¹Department of Fundamental Neuroscience, Faculty of Medicine, University of Geneva, Geneva,
9 Switzerland

10 ²Swiss National Centre of Competence in Research “Evolving Language” (NCCR EvolvingLanguage)

11 ³MRC Cognition and Brain Sciences Unit, University of Cambridge, UK

12 ⁴Institut Pasteur, Université Paris Cité, Inserm, Institut de l’Audition, F-75012 Paris, France

13

14

15

16 *Corresponding author

17 Email: yaqing.su@unige.ch (YS)

18 IO and ALG are joint senior authors

19 **Abstract (177 words)**

20 Understanding speech requires mapping fleeting and often ambiguous soundwaves to meaning.
21 While humans are known to exploit their capacity to contextualize to facilitate this process, how
22 internal knowledge is deployed on-line remains an open question. Here, we present a model that
23 extracts multiple levels of information from continuous speech online. The model applies linguistic
24 and nonlinguistic knowledge to speech processing, by periodically generating top-down predictions
25 and incorporating bottom-up incoming evidence in a nested temporal hierarchy. We show that a
26 nonlinguistic context level provides semantic predictions informed by sensory inputs, which are
27 crucial for disambiguating among multiple meanings of the same word. The explicit knowledge
28 hierarchy of the model enables a more holistic account of the neurophysiological responses to
29 speech compared to using lexical predictions generated by a neural-network language model (GPT-
30 2). We also show that hierarchical predictions reduce peripheral processing via minimizing
31 uncertainty and prediction error. With this proof-of-concept model we demonstrate that the
32 deployment of hierarchical predictions is a possible strategy for the brain to dynamically utilize
33 structured knowledge and make sense of the speech input.

34 **Introduction**

35 Understanding speech is a non-trivial feat. To extract information from ever-changing acoustic
36 signals, our brains must simultaneously “compress and recode linguistic input as rapidly as possible”
37 for multiple representation levels (1), while also keeping information in memory as we incrementally
38 build up the meaning of an utterance (2). No computational framework to date has captured the
39 transformation from continuous acoustic signal to abstract meaning: most speech processing models
40 focus on either the lower-level recognition from acoustic to lexicon (3-7), or the higher-level
41 linguistic manipulations without taking into account the constraint of elapsing time (8-13).

42 In addition to the challenge of fleeting time, speech signals are often ambiguous. However, humans
43 exhibit extraordinary flexibility in making sense of ambiguous speech. We constantly make
44 inferences based on our internal linguistic (e.g. syllabic composition of a word) and nonlinguistic
45 prior knowledge (e.g. speaker identity, semantic context) that are learned from our personal
46 experience. The influence of internal (prior) knowledge on speech perception takes place at all
47 processing levels, e.g. filling the gap of possibly obscured acoustic details (14-16), or interpreting a
48 sentence containing semantically ambiguous words (17, 18). Understanding how internal knowledge
49 is integrated with external input on the fly is key to deciphering speech processing in the brain, and
50 explaining the flexibility in human speech comprehension.

51 With the development of powerful neural networks (19-21), it is now possible for a model to
52 implicitly learn structured linguistic knowledge from an immense amount of written text, and apply
53 such knowledge in language tasks such as coherent text generation. Despite their remarkable
54 achievements in specific language tasks, these models are very resource-demanding and often make
55 egregious errors showing that their performance is not rooted in human-like understanding of the
56 language content (22, 23). Especially if trained and evaluated on tasks involving predicting the next
57 input (20, 21), e.g. a word, it is virtually impossible for such models to capture the abstract
58 processing necessary for human language comprehension extending beyond linguistic forms and
59 across cognitive domains (24, 25). A key aspect of speech understanding consists of applying
60 structured internal knowledge to extract relevant information from the input signal. How and what
61 internal knowledge is deployed depends on the listener's behavioral goal, which can range from
62 "understanding the message intended by the speaker" during a conversation to simply "predicting
63 the next word" during an experimental task. A language model exploiting built-in linguistic as well as
64 nonlinguistic knowledge, and driven by a behavioral goal, may hence be more powerful and
65 polyvalent than one based on recognition and short-range prediction.

66 Here, we propose a computational framework in which the use of linguistic and nonlinguistic
67 contextual knowledge allows the incremental extraction of multi-level information from the
68 continuous speech signal. The model achieves single-sentence understanding by assigning
69 appropriate values to semantic roles and making reasonable judgements about the nonlinguistic
70 context in which the sentence takes place. Such a process relies on a probabilistic generative model
71 that uses its linguistic and nonlinguistic knowledge to incrementally compose sentences. The
72 generative model has a top context level that determines 2nd-level semantic roles, which are
73 translated into a 3rd-level lemma sequences via linearized syntax rules. Each lemma produces a
74 sequence of continuous, bottom-level spectro-temporal patterns via two intermediate hierarchies,
75 integrating a syllable model (26) that was adapted from a biophysically plausible model of birdsong
76 recognition (27, 28). Importantly, context and semantic states are maintained throughout the
77 sentence but interact at the lemma rate, allowing the inverse model to modify previous estimates of
78 these states with incoming evidence. During model inversion, top-down and bottom-up messages
79 alternate at timescales of corresponding hierarchies, providing a possible solution to the "now-or-
80 never" bottleneck (1) that is also consistent with the predictive coding hypothesis of perception (29-
81 31).

82 With a small scope of knowledge adapted from stimuli in MacGregor et al. (32), the model can
83 extract contexts and semantic roles from ongoing speech signals and resolve semantic ambiguity

84 using new information; its beliefs about context and semantic roles, in turn, dynamically influence
85 message passing in lower levels. The linguistically informed model structure allows for hierarchy-
86 specific computational metrics that provide a more interpretable and holistic explanation of neural
87 speech responses than using next-word prediction statistics generated by GPT-2 (20), a large-scale
88 natural language model. In addition, we show that the prediction-update mechanism offers the
89 flexibility to balance between amount of processing and inference accuracy through the control of
90 weighting for bottom-up sensory cues versus top-down predictions.

91 This proof-of-concept model demonstrates a possible computational scheme of speech processing in
92 the brain in which top-down prediction serves as a key computational mechanism for information
93 exchange between hierarchies, driven by the goal of comprehension. Furthermore, correlations
94 between model-derived metrics and neural responses may provide insights into the functional roles
95 of various neuronal signals during speech perception.

96

97 Results

98 A deep hierarchical model of speech comprehension

99 We developed a model of speech processing based on the idea that the goal of the listener is to
100 understand the message conveyed by an utterance. Appropriate understanding entails retrieving
101 useful information from the utterance and optimally mapping it to the listener's knowledge of the
102 world, not restricted to linguistic representations (Fig 1A). Our model of the listener's internal
103 knowledge therefore consists of two parts that are both implemented as probabilistic generative
104 models. The first part exemplifies knowledge about the world by defining events and properties
105 constrained by specific nonlinguistic, situational contexts. For example, under the context of a tennis
106 game, the listener knows (that the speaker knows) about special winning serves, about runs to
107 return a ball etc. The serve or the run may be the central role in an event of winning a game, or
108 described as having a certain property (e.g.: being surprising). Under the different context of a poker
109 game, the listener knows some cards in the deck that can also be part of an event or entail some
110 property. The second part of the model converts these events or properties into linguistic forms by
111 choosing between a number of possible lemmas in an appropriate order, e.g. the special winning
112 serve can be expressed as a single word "ace" early in the sentence, and finally into spoken
113 utterances in the form of spectro-temporal sound patterns via a deep temporal hierarchy (Fig 1B).
114 These two parts are hierarchically linked via semantics and syntax. The inversion of this generative
115 "world knowledge" model fulfills the mapping from the sound patterns to abstract semantic roles
116 and contexts by estimating the probability of every possible value (*state*) of each element (*factor*) in
117 the knowledge hierarchy (Fig 1A), thus providing the listener with the means to understand the
118 utterance produced by the speaker.

119 In all, the model includes five levels, each consisting of several factors (represented in rectangles in
120 Fig 1A) which have multiple possible values (states) listed in Table 1 except for the *acoustic* factor,
121 which is a real-valued vector representing the signal amplitude of six acoustic channels. Probabilistic
122 mappings and transition probabilities between the values of the discrete factors in Table 1 are
123 defined in Methods and Appendix. The final output of the generative model (i.e. the input to the
124 perception model) is the continuous spectro-temporal pattern of the speech signal sampled at 1000
125 Hz and divided into six frequency channels (see Methods). Lengths of stimuli are fixed: each
126 sentence consists of 4 lemmas, each lemma of 3 syllables, and each syllable of 8 spectral vectors.
127 Every spectral vector is deployed into 25ms of time-varying continuous signal, thus each syllable
128 effectively has a duration of 200ms (33).

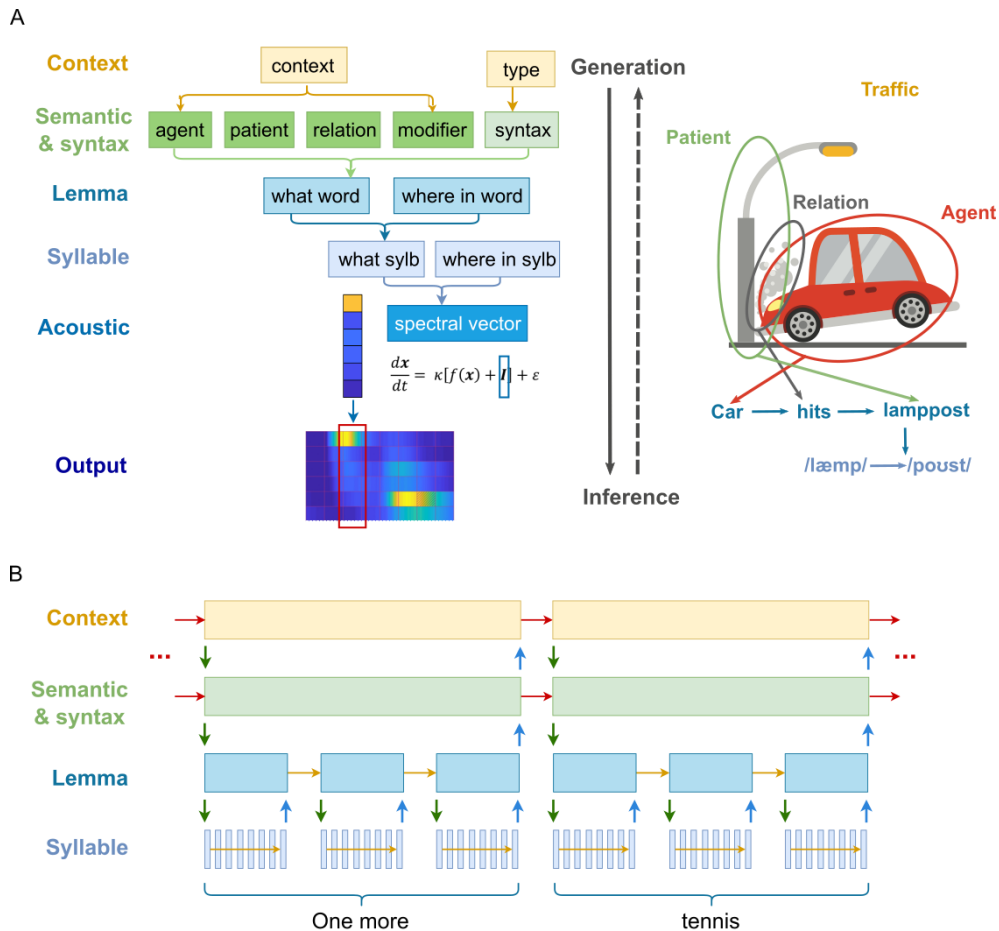


Fig 1. A generative model of speech and its inversion. A. Schematic of the generative model. Left: information conveyed in a speech signal is roughly separated into six hierarchies. To generate speech (solid downward arrow), the model first assigns values to semantic roles according to the contextual knowledge and determines a (linear) syntactic structure from the type of the message it's expressing. Together, semantics and syntax generate an ordered sequence of lemma units. Each lemma unit generates a sequence of syllables, which in turn generates a sequence of spectral vectors. Each spectral vector unit is then deployed as a continuous acoustic signal of 25 ms. Inference corresponds to the inversion of the generative process (dashed upward arrow). The model is divided into three parts that were implemented with different algorithms (see Methods). **Right:** cartoon (www.publicdomainpictures.net) illustrating how a sequence of syllables '/læmp-poust/' (lamppost) is generated from a traffic scene context. In describing a traffic accident, the speaker tries to convey its mental image of the scene consisting of an agent (the car), a patient (the lamppost) and the relation (the action of hitting) from the agent to the patient. With English vocabulary and grammar, it chooses one lemma corresponding to each element in the accident, and outputs (speaks) these lemmas in a specific order according to the syntactic rules. Each lemma is then expressed as a specific sequence of syllables. Importantly, the same lemma can be the result of different combinations of abstract information and syntactic rules. For example, in the sentence "The ball hits the floor", the word "hits" implies a different action than a car hitting a cyclist, whereas in "His songs are top hits" the relative position of the word implies an entity, not an action. **B. Temporal scheduling of hierarchical message passing during speech perception.** The generative model is inverted by alternating

top-down prediction (prior, green downward arrows) and bottom-up update (blue upward arrows). A supraordinate level initiates a sequence of evidence accumulation in its subordinate level and receives a state update at the end of such sequence. It then makes a transition and sends an updated prediction to the subordinate level and initiates another sequence of evidence accumulation. Such a process is repeatedly performed until the end of the sentence. Note that for the lemma and lower levels, states are generated anew each time when the supraordinate level makes a transition, i.e. no horizontal arrows between sending up an update and receiving a new prior. For the top two levels, however, states are maintained throughout the sequence (red horizontal arrows) or make transitions according to a set of rules (syntax).

129

Table 1. Factors and their possible values (states) in the model hierarchy

Hierarchy	Factor	Value (State)
Context	Context	tennis game, poker game, night out, car racing game
	Sentence type	event, property
Semantic & Syntax	Agent (semantic)	card a, winning serve, run, card j, neckband, score, buzz, null
	Patient (semantic)	tennis game, poker game, racing game, evening, null
	Relation (semantic)	win, ruin, be
	Modifier (semantic)	sufficient, unexpected, not pretty, not fair, high volume, high frequency
	Syntax	attribute, subject, verb, object, adjective
Lemma	Lemma	one more, that, ace, sprint, joker, tie, noise, wins, ruined, is, the tennis, the poker, the game, the evening, enough, surprising, ugly, unfair, loud, sharp
	Where in lemma	1-3
	Syllable	/eis/, /te/, /nis/, ... total of 32 including the silence syllable
	Where in syllable	1-8

130

*Note that these symbols are illustrative and not following IPA.

131

Next, we show how this model understands simple sentences and deals with semantic ambiguity,

132

and we demonstrate the role of top-down predictions in these processes. We assessed its

133

performance with different sentence stimuli and parameter settings, namely by varying the

134

perceptual bias among different contexts and the precision of the continuous module (see

135

Methods), focusing on: 1) the probability distributions that describe the model's beliefs (or

136

predictions) about possible states over time, and 2) divergence and entropy measures, which

137

summarize informational changes underlying the evolution of beliefs (see Methods). These

138

measures do not depend on the precise fine tuning of the model parameters, and are qualitatively

139

evaluated by whether the timing (when certain states are updated) and the outcome (what the

140 current beliefs are about different states) of the hierarchical inference concurs with human behavior
 141 in the language domain.

142 Stimuli are adapted from MacGregor et al. (2020) (32) and illustrate the use of internal knowledge to
 143 disambiguate speech. All sentence stimuli in the following sections share the same structure (see
 144 Table 2 for a complete list of possible sentences):

145 One more [MIDDLE WORD] wins [END WORD].

146 The MIDDLE WORD can have either one or multiple possible meanings, each meaning pointing to
 147 one context of the sentence. The END WORD either resolves the semantic ambiguity raised by the
 148 middle word or not. A disambiguating end word can also follow an unambiguous middle word
 149 without affecting its interpretation.

150 **Table 2. All possible sentences in the model**

Attribute	Subject	Verb	Object/Adjective	Context
One more/that	ace	wins	the game/ the poker	poker game
			the game/ the tennis	tennis game
		is	surprising/enough	poker or tennis game
	sprint	wins	the game/the tennis	tennis game
		is	surprising/enough	tennis game
	joker	wins	the game/the poker	poker game
		is	surprising/enough	poker game
	tie	ruined	the evening	night party/racing game
			the game	racing game
		is	ugly	night party
	unfair		racing game	
	noise	ruined	the evening/the game	night party/racing game
		is	loud/sharp	night party/racing game

151 **The use of knowledge about the world to interpret speech**

152 We first test how the model processes speech stimuli, with a focus on the timing of the incremental
 153 estimation process at the context and semantic levels, where “meaning” is extracted by assigning
 154 values to semantic roles.

155 Consider the following two sentences, A: “One more ace wins the tennis.” and B: “One more ace
 156 wins the game.” Both sentences contain the ambiguous word “ace”, which can be associated with a
 157 special serve in tennis or a special card in a poker game. The final word in the first sentence

158 disambiguates “ace” to mean a special serve because “the tennis” can only be generated from a
159 tennis game context, which applies to the whole sentence including the preceding “ace”. In the
160 second sentence, however, the ambiguity remains unresolved; the game can still refer to a tennis
161 or a poker game. In the latter case, the interpretation of the word “ace” will depend on the listener’s
162 preference. Unless specified otherwise, we introduce a prior preference for the poker context to
163 reflect the preference of the general population (32).

164 The word “ace” introduces ambiguity because it points to two possible states for *agent* (“tennis
165 serve” or “card A”), each of which points to a separate state for *context* (“tennis game” or “poker
166 game”, Table 2, ambiguous and disambiguating words in red). Figs 2A and 2B show the evolution of
167 the model’s beliefs about context and semantic factors for the two sentences. The ambiguity is
168 reflected in the posterior estimates of *agent* and *context* between the offset of “ace” and the
169 sentence ending word, where the model assigned nonzero probabilities to “card A” and “serve” as
170 the *agent*, and “poker game” and “tennis game” as the *context*, and near-zero probabilities for other
171 states (Fig 2A). Probabilities for poker-relevant states were higher (darker colors) due to the
172 contextual preference. The verb “wins” did not change the model’s estimation for the *agent* or the
173 *context*, but clarified the sentence *type* to be “event” and the *patient* to be nonempty, again with a
174 preference towards poker. After the model heard “the tennis” (Fig 2A), it immediately resolved its
175 beliefs of the *agent*, the *patient* and the *context* to the opposite of its prior preference. When the
176 sentence ended with “the game”, (Fig 2B), the model followed its preference with enhanced beliefs
177 as a result of the entropy reduction entailed by belief updating, but not as clearly resolved as with
178 “the tennis” (see next section).

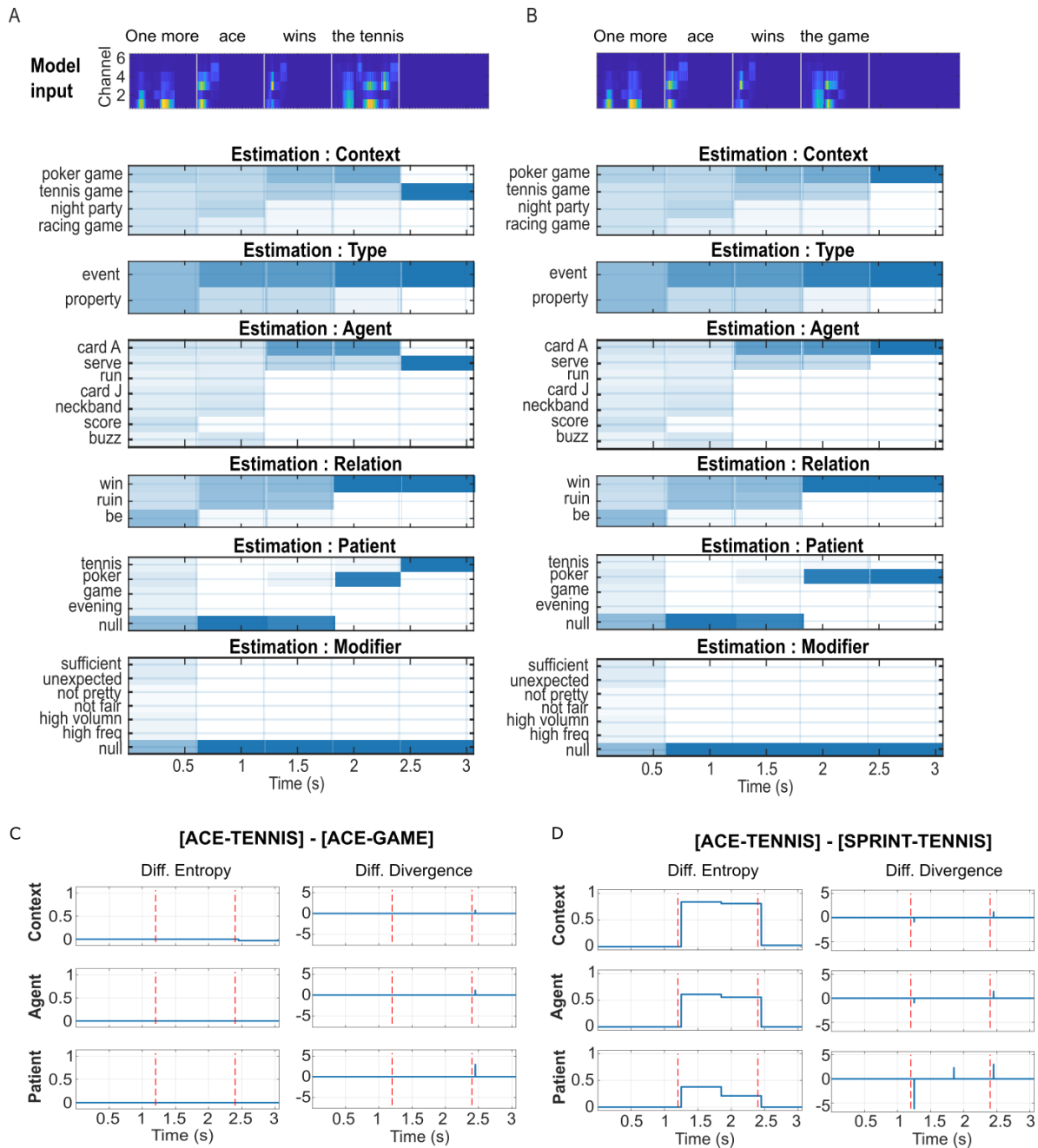


Fig 2. Semantic- and context-level model response to different sentence inputs. For all simulations, relative prior for context was set at the default of 1.5:1:1:1 for the four possibilities {‘poker game’, ‘tennis game’, ‘night party’, ‘racing game’}. **A. Top panel:** acoustic spectrogram of input sentence A: “one more ace wins the tennis”. Vertical grey lines mark the offset of each lemma, at which point updates were sent from the lemma level to semantic and context. **Lower panels:** estimation of posterior probabilities for the semantic (*agent*, *patient*, *relation*, *modifier*) and context states as the sentence unfolds. Possible values of each factor are labelled on the y axis. Blue scale blocks indicate the probability distribution for each factor, dark blue— $p=1$, white— $p=0$. The updating process is nearly instantaneous, and the main body of the n^{th} block (epoch corresponding to one lemma) is filled with the estimates after the $(n-1)^{\text{th}}$ update. The first input “one more” was not informative. The estimated distributions were slightly changed before and after

the offset of “one more” because the model still performed gradient descent to minimize free energy. After hearing “ace”, distributions for the *context* and the *agent* converged to either “poker game” or “tennis game” for *context*, and ‘card A’ or ‘serve’ for *agent*. Within these possibilities, probabilities for the poker *context* and the ‘card A’ *agent* were higher, reflecting the prior preference. Probabilities of “tennis” or “poker” as *patient* also increased. *Type*, *relation*, and *modifier* remain the same as in the previous epoch. After hearing ‘wins’, possibilities for *type* converged to ‘event’, and those for *relation* converged to ‘win’. Probabilities for ‘tennis’ and ‘poker’ for *patient* further increased, with a strong bias for “poker”, while the probability of a ‘null’ *patient* decreased to zero. In the last epoch, the model received a disambiguating phrase ‘the tennis’, and all factors are resolved to the correct state with a probability close to 1. **B. Acoustic input and probability estimation for the sentence “one more ace wins the game”.** The distributions are the same as in A before the last update. In the last epoch, the model receives an input, ‘the game’, that does not resolve the semantic and contextual ambiguity. As a result, distributions were further biased towards values corresponding to the ‘poker game’ context. **C. Entropy and Divergence derived from the sentence “one more ace wins the tennis” relative to the sentence “one more ace wins the game”.** The two vertical dashed lines mark the offset of the sentence middle word “ace” and the ending word, respectively. As the two sentences only differ in the ending word, both metrics differ only at sentence offset. Compared to “the game”, which does not completely resolve the ambiguity introduced by ‘ace’, ‘the tennis’ results in lower entropy in “context” (top left panel), indicating greater certainty about the estimate. The zero differences in entropy for agent and patient indicate that the model tends to believe in its bias for these two factors. “The tennis” also gives rise to higher divergence (right panels) at sentence offset. **D. Results from the sentence “one more ace wins the tennis” relative to “one more sprint wins the tennis”.** At its offset, the ambiguous word “ace” introduces higher entropy for all three factors compared to “sprint”, reflecting greater uncertainty about the hidden states. Uncertainty dominates divergence, which is indexed by a corresponding negative difference here. At sentence offset, entropy differences between the two sentences became minimal because the model has resolved hidden states of all hierarchies. The positive difference in divergence at the offset reflects the higher surprisal for “the tennis” when it follows “ace” compared to “sprint”.

179 The results in Fig 2A and 2B demonstrate how prior knowledge and preferences can dynamically
180 influence the extraction of semantic roles and contexts from the speech signal. This influence is not
181 only reflected in the perception of semantically ambiguous words, but also in the details of message
182 passing that give rise to its estimates. Fig 2C contrasts the inference processes between sentence
183 [ACE-TENNIS] and sentence [ACE-GAME] in Fig 2A and 2B using their derived information metrics
184 ([ACE-TENNIS] relative to [ACE-GAME]), focusing on the context, the agent, and the patient factors
185 that were most relevant for the set conditions. Fig 2D compares the same metrics between
186 sentences [ACE-TENNIS] and [SPRINT-TENNIS]. These contrasts were based on similar comparisons in
187 the M/EEG study of MacGregor et al. (32), where the authors identified two relevant findings. First,
188 they showed an effect of ambiguity on the magnitude of MEG sensor-space response activations

189 shortly after the word offset (increased activation for “ace” compared to “sprint”), which could be
190 interpreted as reflecting increased uncertainty. Secondly, they showed a (marginally significant)
191 effect of resolving ambiguity (increase in the difference of activation between “the tennis” after
192 “ace” vs. after “sprint” compared to between “the game” after “ace” vs. after “sprint”), which could
193 be interpreted as reflecting increased surprisal. Respectively, these two effects were qualitatively
194 captured by a difference in model-derived entropy (Fig 2D, left) and Kullback-Leibler (KL) divergence
195 (Fig 2C and 2D, right) in response to the sentence contrasts. However, a difference in entropy
196 between two conditions is often associated with a difference in divergence but in the opposite
197 direction, with magnitudes varying across hierarchies and across factors within the same hierarchy.
198 Thus, both semantic ambiguity and its resolution likely involve a complex combination of
199 computational processes of different types and hierarchies. Such a complexity is in line with the
200 finding of MacGregor et al. (32) that the two sensor-space phenomena were localized to different
201 but overlapping sources. Further dissociation between different computation processes should
202 involve correlating model-derived information metrics, importantly at different hierarchical levels
203 and factors, with source-, time- and frequency-specific responses (see Discussion).

204 While the direction of prior preference (e.g. poker over tennis) influences both the information
205 passing and the perceptual outcome (the state with highest posterior probability) as shown in Fig 2,
206 the degree of prior preference also has a subtle influence on message passing during the inference
207 process. With the same perceptual outcome, (S1 Fig A and D, either side of bias=1), the amount of
208 information maintained between belief updates as quantified by entropy, and the magnitude
209 information change induced by an update as quantified by the KL divergence, both vary
210 quantitatively with the model’s prior preference (S1 Fig B-C, E-F). Thus, model-derived information
211 metrics provide a means to relate the variability of neurophysiological responses to the perceptual
212 preferences of individual subjects.

213 **Semantic prediction influences low-level message passing**

214 The deployment of hierarchical prediction implies that high-level (*semantic, syntax* and *context*)
215 state estimates also dynamically influence the top-down predictions (priors) as well as the bottom-
216 up updates at lower (*lemma, syllable* and *acoustic*) levels. Figs 3A and 3B respectively show top-
217 down priors and posterior estimates at lemma and syllable levels with the same parameters as Fig
218 2A. The predictions reflect both prior knowledge and the updated estimates of superordinate levels,
219 in agreement with recent neurophysiological evidence that high-level (word) predictions constrain
220 low-level (phoneme) predictions (34). Posterior estimations of both levels immediately converged

221 onto the correct states after receiving the disambiguating input, for example the second syllable in
 222 the last lemma.

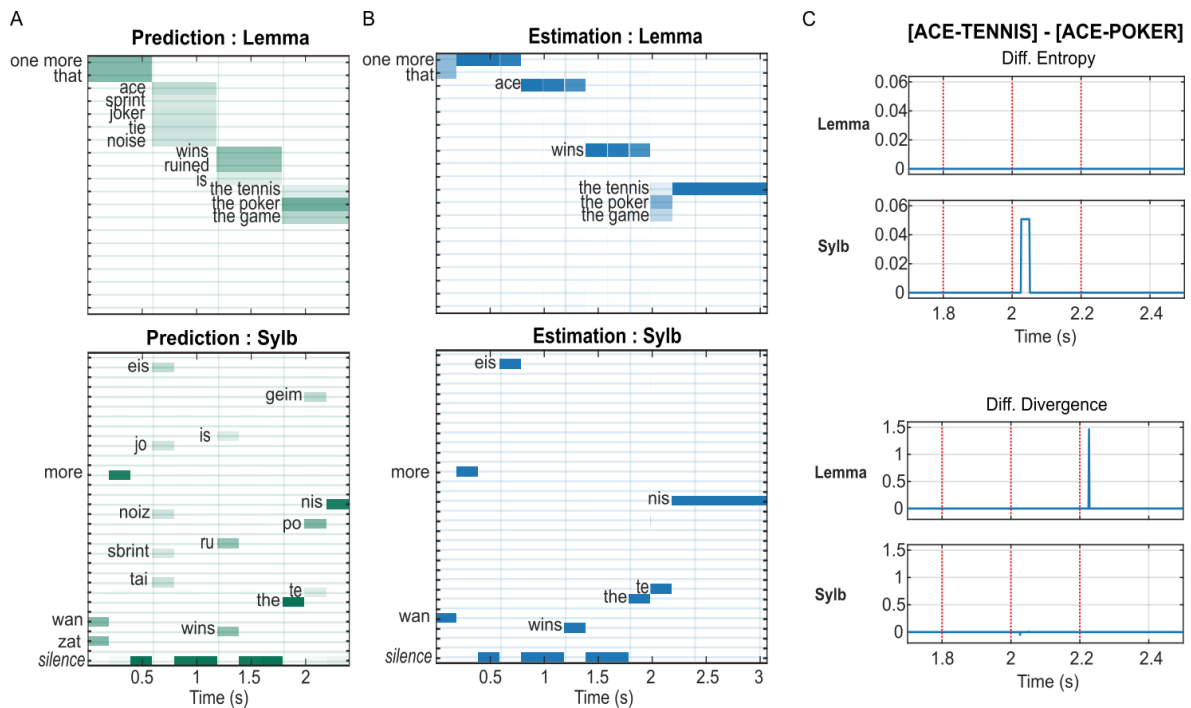


Fig 3. Influence of semantic state estimates on the prediction and updating of lemma and syllable states.

A. Semantic-to-lemma and lemma-to-syllable predictions (prior expectations) for the simulation in Fig 2A.

Vertical lines indicate offsets of each lemma input. In lemmas 1-3, syllable predictions (lower panel) are nearly certain after the first syllable because there was a one-to-one correspondence between the lemma and the first syllable. In lemma 4 (“the tennis”), the opposite is true because all possible lemmas start with the syllable “the”, diverging at the second syllable. Lemma predictions (top panel) depend on the current estimates at the superordinate level and the contextual bias, e.g. the prediction for the last lemma is highest for “the poker”, and lowest for “the tennis”.

B. Estimation of posterior probabilities for lemma and syllable states for the simulation in Fig 2A.

The model quickly recognizes each syllable (lower panel). The estimation for lemma states (upper panel) appears to lag for the duration of one syllable, because the lemma level receives a nearly instantaneous update at the **offset** of every syllable, and the grid between the i^{th} and $(i+1)^{\text{th}}$ updates is filled with the estimated distribution of the i^{th} update. For example, the estimation for the first lemma started with a 1:1 prior expectation between “one more” and “that”, then converged to “one more” after hearing the first syllable “one”. The estimate was not changed until the offset of “ace”, the first syllable of the second lemma. This is only due to our graphical representation and does not affect the update from lemma to semantics. **C. Upper panels: entropy derived from sentence [ACE-TENNIS] relative to sentence [ACE-POKER] for the lemma and the syllable levels in the proximity of the final lemma.**

Vertical dotted lines mark the onset of each syllable of the final lemma, either /the-te-nis/ or /the-po-ker/. A slightly higher syllable entropy after the onset of the second syllable for “the tennis” indicates the model took longer, i.e. more gradient descent steps, to converge to the less expected input /te/. **Lower panels: the**

difference between the divergence in response to the two sentences. A higher lemma divergence at the onset of the third syllable (the offset of the second syllable) for the lemma “the tennis” reflects that “tennis” is less expected than “poker” due to the preference at the context level.

223

224 For the sentence input “one more ace wins the poker”, the model makes the identical semantic-to-
225 lemma predictions as in Fig 3A (top panel), and nearly identical lemma-to-syllable predictions except
226 for the final syllable, which was informed by the preceding syllable /po/ in “the poker” (not shown).
227 Fig 3C shows the entropy and divergence derived from the posterior estimates of sentence [ACE-
228 TENNIS] relative to [ACE-POKER] for the lemma and syllable level, focusing on the final lemma.
229 Although the amplitudes of the differences are smaller than those at the semantic and the context
230 level (Fig 2), their presence indicates that lower-level processes likely also contribute to the
231 observed differences in neurophysiological response to semantically expected vs. unexpected
232 speech inputs, corroborating the finding that the neural encoding of phonological and acoustic
233 information of a word input is modulated by its semantic similarity to its preceding sentential
234 context (34). The influence of semantic prediction on lower-level message passing can also be
235 reflected in the processing of the same word embedded in different sentences, e.g. “the tennis” in
236 the sentence [ACE-TENNIS] vs. [SPRINT-TENNIS] (S2 Fig). Unlike the semantic and context levels,
237 however, the difference between “ace” and “sprint” at the acoustic and phonological levels was not
238 reflected in the low-level message passing (S2 Fig C).

239 **Interpreting neural speech response requires lexical prediction and** 240 **beyond**

241 Information metrics derived from our model suggest that the sensor-space effects observed in
242 MacGregor et al. (32) mainly reflect the message passing in semantic- and context-level processing
243 (Fig 2), rather than in the lemma (word) level (Fig 3, S2 Fig). Meanwhile, several recent studies have
244 successfully used word or phoneme prediction statistics derived from the output of natural language
245 models to explain variabilities in neural response to the semantic aspects of linguistic stimuli (35-38).
246 In doing so, the surprisal evoked by the received input given the preceding sentential context, and
247 less often the entropy of the prediction for the upcoming input, are used directly or indirectly (in
248 conjunction with additional regressors and regression models) as proxies of semantic knowledge to
249 identify the neuronal dynamics underlying semantic processing. To understand the extent to which
250 the output of a language model trained on next-word prediction can directly explain semantic- and
251 context-level effects on neurophysiological speech responses, we reanalyzed the neurophysiological

252 data of MacGregor et al. (32) using both explicit semantic properties as in the original study and
253 next-word prediction statistics from GPT-2 (20) (see Methods).

254 We first explored whether GPT-2 predictions captured the semantic ambiguity and disambiguation
255 in the stimuli. We adopt the terminology of MacGregor et al. (32), referring to the sentence-middle
256 word as “Target” and the sentence-ending word as “Resolution” (Table 3). Fig 4A shows the
257 distributions of prediction entropy after the ambiguous (blue) vs. unambiguous (orange) target
258 word. A one-way ANOVA indicates no significant difference between entropy in the two Target word
259 types (mean entropy: ambiguous = 4.734, unambiguous = 4.658; $p=0.59$). Fig 4B shows the
260 distributions of surprisal after receiving the resolving (blue) vs. unresolving (orange) Resolution
261 word, either following an ambiguous (left panel) or unambiguous (right panel) Target. A two-way
262 ANOVA showed that, although the surprisal values of resolving words are significantly higher than
263 those of unresolving words regardless of Target ambiguity (mean surprisal: resolving = 7.741,
264 unresolving = 5.937; $p<0.001$), there was no difference of surprisal depending on the preceding
265 ambiguity of the Target word (mean surprisal: prior ambiguity = 6.955, no prior ambiguity = 6.724; p
266 = 0.74), nor on the interaction between resolution and ambiguity ($p = 0.68$). Thus, similar to the
267 model’s lemma-level prediction metrics (S2 Fig C), GPT-2 entropy does not reflect the semantic
268 ambiguity of Target words, neither does the evoked surprisal capture the long-distance interaction
269 between Target and Resolution.

270 **Table 3. Example sentence input to the MEG subject and GPT-2**

Lead in	Target	Bridge	Resolve	Unresolve
The man knew that one more	ace	might be enough to win the	tennis	game
The woman hoped that one more	ace	might be enough to win the	tennis	game
The man knew that one more	sprint	might be enough to win the	tennis	game
The woman hoped that one more	sprint	might be enough to win the	tennis	game

271 We next compared how variabilities of semantic information and GPT-2 predictions correlate with
272 neurophysiological responses. In particular, we tested the effects of semantic properties (conceptual
273 replication of MacGregor et al. (32)) and GPT-2 prediction statistics on the MEG response during two
274 time windows around the Target offset and the Resolution offset. As in the original M/EEG study, we
275 focused on combined gradiometer pairs, which demonstrated the most robust effects, and two
276 analysis time windows around the Target offset and Resolution offset respectively.

277

denote sensor clusters that showed a prevalent positive effect of ambiguity within the time window. **D.**

Statistical test results for the effect of semantic ambiguity in the preceding context (left column) and GPT-2 prediction surprisal (right column) on MEG combined gradiometer data around the time of Resolution

offset. Top row: sensor-time maps for significance level of sensor clusters. Bottom rows: Topological distributions of the corresponding effects averaged over six 250ms time windows, spanning from -0.5 to 1s relative to the Resolution offset.

278 For the Target time window, we split the MEG response into two groups according to the property of
279 the Target word pair: 1. The GPT-2 entropy of the ambiguous Target is larger than that of its
280 unambiguous counterpart, and 2. The GPT-2 entropy of the ambiguous Target is smaller than that of
281 its unambiguous counterpart. S3 Fig A shows the distribution of entropy differences between
282 ambiguous and unambiguous Target word pairs (ambiguous minus unambiguous). Ambiguous Target
283 words with difference > 0 (i.e. in group 1, 29 pairs in total) and unambiguous Targets with difference
284 < 0 (i.e. in group 2, 29 pairs in total) contribute to the high-entropy group, and the rest contribute to
285 the low-entropy group. Such splitting ensures that the pair of Target words in the same sentence set
286 is always separated into two conditions, thus controlling possible confounds of the preceding
287 sentential context. Using a data-driven algorithm (see Methods), we identified sensor-time clusters
288 that showed a significant effect (two-tailed paired student's t-test, $p < 0.05$, same in the following
289 results) of semantic ambiguity by contrasting responses to ambiguous Target vs. unambiguous
290 Target words, (Fig. 4C, left column). We also identified clusters showing an effect of GPT-2 entropy
291 by contrasting responses to Target words with high vs. low entropies (Fig. 4C, right column). Sensor-
292 time statistical maps (Fig 4D, top row) as well as topographic plots over time (Fig 4D, bottom row)
293 indicate that these two effects are likely distributed differently both in space and time. The absence
294 of a significant correlation (Pearson's correlation $r = -0.04$, $p = 0.66$) between the sensor-wise effect
295 sizes of the two contrasts (S4 Fig A) also suggests that semantic ambiguity and GPT-2 prediction
296 entropy may account for different spatial aspects of the MEG responses. Interestingly, the positive
297 effect of GPT-2 entropy arose before the word offset, whereas the positive effect of semantic
298 ambiguity was only apparent after the word offset (Fig 4C, top row).

299 For the Resolution timepoint, responses to only the Resolve sentence ending were split into two
300 groups in a similar fashion as for Target: 1. The GPT-2 surprisal following the ambiguous Target was
301 larger than the same word following the unambiguous Target, and 2. The GPT-2 surprisal following
302 the ambiguous Target was smaller than the same word following the unambiguous Target. Thirty-six
303 out of the 58 sentences were labeled as being in group 1, 22 in group 2 (S3 Fig B). The contrast
304 between Resolution words following ambiguous vs. unambiguous Target words revealed an effect of
305 ambiguity of the previous context distributed among right temporal-parietal and midfrontal areas

306 spanning several time windows before and after the word offset (Fig. 4D, left column). The contrast
307 between Resolution words with high vs. low surprisal revealed an effect of GPT-2 prediction surprisal
308 with a different spatial distribution, and restricted to -250 to 250ms (Fig. 4D, right column). Similar
309 to the Target effects, the effect sizes of ambiguity and surprisal at Resolution offset were not
310 correlated ($r=0.001$, $p=0.99$, S4 Fig B) across sensor locations.

311 These results demonstrate that both GPT-2 word-prediction statistics and high-level semantic
312 properties contribute to the variability in neural speech responses, but their effects exhibit different
313 spatio-temporal distributions. Given that predictions from the GPT-2 output cannot directly capture
314 the semantic properties we investigate here (Fig 4A, B), the approach of interpreting the neural
315 response to speech (and more generally language) solely based on such predictions learned from
316 word sequence statistics overlooks important aspects of the dynamics underlying higher-level
317 language processing. Our model, on the other hand, explicitly depicts multiple levels of linguistic and
318 nonlinguistic processes under the same computational principles. Thus, it points to a more
319 interpretable and holistic approach to characterizing the functional network underlying speech
320 comprehension. A quantitative mapping between model and neural responses requires a nontrivial
321 expansion of the model and is beyond the scope of the current study (see Discussion).

322 **Top-down prediction reduces processing effort**

323 The model works by iteratively calculating the discrepancy between top-down predictions
324 (expectation of the input) and bottom-up input at each hierarchical level, and using such a
325 discrepancy to modify the state estimates of superordinate levels. This does not mean the model
326 needs to make the best prediction for the next input as in Fig 3A: hierarchical predictions are a
327 necessary computational mechanism in relaying information for making better inferences, even if
328 the actual input deviates from the predicted one. To examine how the prediction content may
329 influence the inference process, we ran the model using the same input as Fig 2A and 3B, “one more
330 ace wins the tennis”, but simulating the extreme case of uninformative (uniform distribution across
331 all possibilities) top-down predictions. We found that the predictive content influenced both the
332 model time course and final estimate.

333 Compared to the condition of informative top-down predictions (Fig 3B), when top-down predictions
334 were uninformative, the model still made correct inferences about every input, but with a slight
335 delay for syllables (Fig 5A). Fig 5B contrasts the entropy and cumulative divergence during the
336 inference process between the two conditions. Unsurprisingly, informative predictions lead to
337 reduced entropy (maintenance of possible items) and divergence (magnitude of updates after the

338 integration of new evidence), both contributing to fewer steps of gradient descent at each point of
 339 belief updating, hence less computation effort in terms of processing time and energy cost (39).

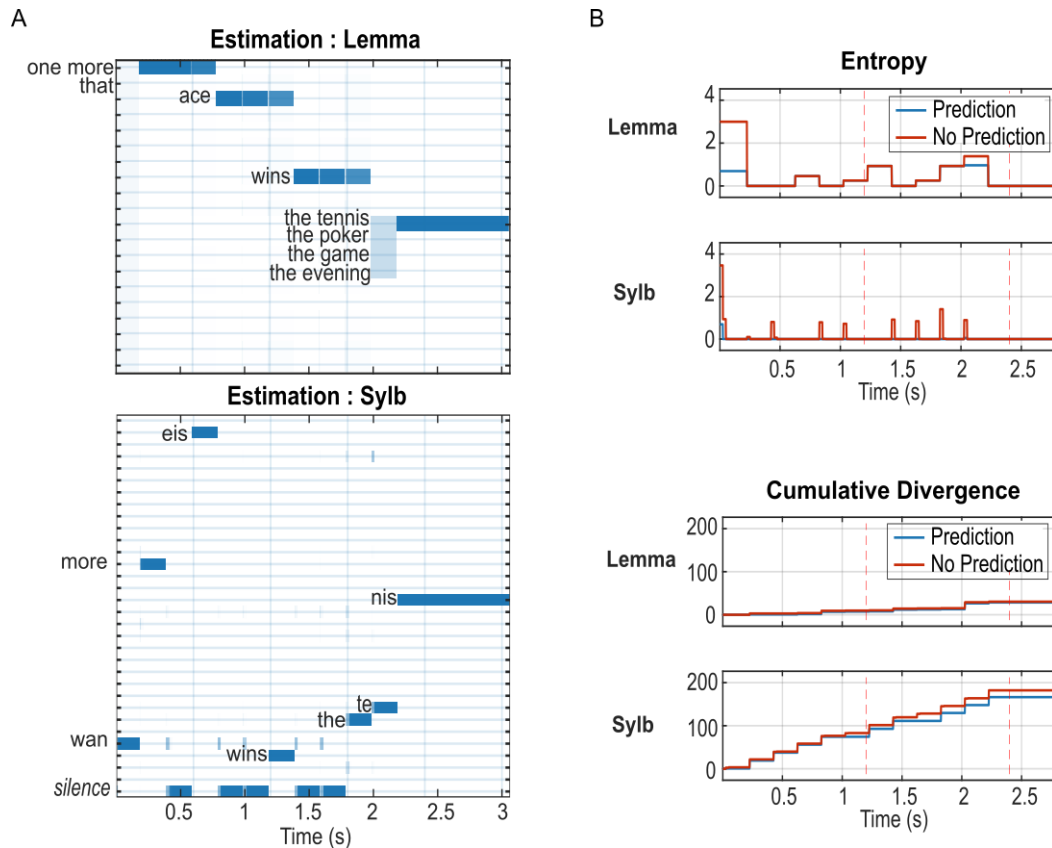


Fig 5. Influence of top-down predictions on syllable and lemma inference under high peripheral precision.

All results are simulated with the sentence “One more ace wins the tennis”. With uninformative predictions, model responses at the semantic and context levels are nearly identical to Fig 2A because the model reached the same, almost-certain lemma estimates at the time of semantic updating (at each lemma offset). Therefore we omit the higher-level results here and in Fig 6. **A. Estimation of posterior probabilities when top-down predictions are set to uniform distributions for all possible states.** Compared to Fig 4B, there is a slight delay for the convergence of every syllable indicated by the small vertical bars, each corresponding to one spectral vector, in more than one possible state. The inference for lemma states is not significantly changed: once the model is certain about the first (or the second in the case of the last lemma) syllable, it can quickly converge to the correct lemma using its internal knowledge. **B. Upper panels: entropy calculated from lemma and syllable states.** With uninformative top-down prediction (red), the entropy of syllable states was raised for a short duration (~1-2 spectral vectors) more often than with informative (blue) prediction (eight times throughout the sentence versus once at the sentence onset). The difference is less obvious at the lemma level except during the very first syllable and the /the/ syllable in the last lemma. **Lower panels: cumulative KL divergence for the two factors.** Overall, the cumulative divergence is smaller when informative prediction is available (blue).

340 So far, we have simulated the model with the ideal scenario of arbitrarily high precisions (see
341 Methods) at the continuous level. In general, a high precision implies that fine details of the input
342 are utilized to evaluate the mapping between the input signal and the generative model, analogous
343 to a perfect periphery that preserves the best possible spectro-temporal information from the
344 acoustic input. It has been suggested that top-down predictions may be especially important under
345 challenging situations, e.g. impaired auditory periphery (40). We tested the model with a broad
346 range of precisions to assess how precision affects online speech processing. In particular, we
347 lowered both the precision for the continuous state as well as for comparing the input with
348 predicted activity in the six frequency channels (see Methods), which is analogous to lesioning the
349 local computation supported by lateral connections and the cross-level information carried by
350 bottom-up connections, respectively (28, 41). Within a considerable degree of degradation, the
351 model performance is qualitatively the same as the intact model, in that it correctly infers the states
352 of all factors, but a strong difference arises in the time it takes to converge, especially in the case of
353 uninformative top-down predictions (S5 Fig A and B, precision= $\exp(6)$ vs. $\exp(16)$ in the intact
354 condition). Fig 6 shows the comparison of informative vs. uninformative predictions similar to Fig 5,
355 but with much lower peripheral precisions ($\exp(0)$). Syllable identification was delayed in both cases
356 when compared to their intact-periphery counterparts (Fig 6A vs. 4B, 6B vs. 5A), and the delay was
357 more pronounced with uninformative predictions. This dramatic delay with uninformative prediction
358 is accompanied by higher entropy (Fig 6C, upper panels) as well as divergence (Fig 6C, lower panels).
359 However, an increase in effort during syllable recognition may be important to avoid inaccurate
360 recognition: in Fig 6A, although the model saved processing time by relying on its prior knowledge, it
361 did so at the cost of incorrectly identifying the final lemma as “the poker”. The tradeoff between
362 processing and accuracy has been well-documented in the decision-making literature (42) and
363 neuroeconomics (43), which reveals that humans flexibly adapt their strategy in challenging
364 scenarios where high accuracy and low effort cannot be achieved simultaneously. Our results
365 suggest that such tradeoff can be manipulated via adjusting one’s reliance on top-down prediction
366 vs. bottom-up sensory information, an ability widely involved in perceptual processes including
367 inferencing others’ intention (44) and likely lacking in certain neuropsychological disorders such as
368 those inducing hallucinations (low sensory precision but high prediction precision) and autism
369 spectral disorder (extraordinarily high sensory precision) (45). Nevertheless, the effort-accuracy
370 tradeoff is also limited by the capacity of the sensory periphery: at extremely low precisions, the
371 model’s syllable recognition breaks down without the guidance of informative top-down prediction
372 (S5 Fig C and D, precision = $\exp(-4)$).

379 **Discussion**

380 The idea that our brains adaptively entertain internal models and that this facilitates language
381 comprehension underlies much current research in speech (language) perception. Nevertheless,
382 how internal knowledge is deployed in time, in relation to the timing of continuous speech
383 unfolding, is an open question, and may be key to achieve the form-meaning distinction in neural-
384 network language models (23, 24). Here, we attempt to establish a foundational framework that
385 dynamically exploits general knowledge in speech comprehension to bridge this gap. We implement
386 the listener's internal knowledge as a probabilistic generative model that consists of a non-linguistic
387 general knowledge (cognitive) model and multiple temporally organized hierarchies encoding
388 linguistic and acoustic knowledge. Speech perception, modeled as the inversion of this generative
389 model, involves interleaved top-down and bottom-up message passing in solving the computational
390 challenge of extracting meaning from ongoing, continuous speech. We show that the model makes
391 plausible inference of hierarchical information from semantically ambiguous speech stimuli and
392 demonstrate the influence of prior knowledge on the inference process, which is reflected in the
393 neural response to speech stimuli but not in next-word prediction statistics of a deep neural-
394 network language model (GPT-2) (20). We also show that hierarchical predictions can be exploited to
395 reduce processing effort. The model tries to mimic human language comprehension by jointly
396 implementing incrementality and prediction (46), and could potentially be expanded towards a
397 comprehensive model of natural language *understanding*, and guide the interpretation of
398 neurophysiological phenomena in realistic listening scenarios.

399 **Language comprehension as semantic role assignment**

400 Although we emphasize that speech (language) comprehension is driven by high-level behavioral
401 goals, to achieve comprehension the appropriate assignment of semantic roles is crucial for
402 (re)constructing the message conveyed in the utterance, e.g. the "mental image" in Fig 1A. Semantic
403 roles can be viewed as an interface between linguistic and nonlinguistic representations, the latter
404 being a fundamental, domain-general format of our internal abstraction of the world (24, 25) that is
405 shown to both behaviorally and neurophysiologically influence language comprehension (47, 48).
406 The process of semantic role assignment is central in psycholinguistic process theories (46, 49-51),
407 yet seldom reflected explicitly in existing computational models of language. A major challenge for
408 modeling semantic role assignment during language processing is in combining meaning extraction
409 with compositionality: words that carry semantic contents are presented in an order dictated by
410 compositional rules, thus the extraction of persisting meanings must take place dynamically

411 alongside the decomposition. These two aspects have only been addressed separately in some
412 existing models, e.g. topic models (9, 52) fulfill (lexical) semantic processing but ignore the word
413 order. On the other hand, the Discovery of Relation by Analogy model (11, 53) learns the time-based
414 binding rules that decompose words and phrases into hierarchical structures, but does not have
415 explicit representations of semantic knowledge.

416 A recent model of linguistic communication (12) did incorporate abstract nonlinguistic (geometric)
417 knowledge and compositionality, but lacked the incremental nature of the meaning-building process
418 in humans (2). The generative model encoded several templates of complete sentences and a set of
419 geometric properties. By applying nonlinguistic knowledge under the goal of resolving object
420 properties, the model generated sentences by picking the most probable sentence format and filling
421 specific positions with the most helpful descriptive words. The inverse model thus comprehends a
422 word sequence by inferring the sentence format and capturing keywords at the corresponding
423 positions. This template-matching strategy realized a form of meaning-structure conjunction.
424 However, it constrains the model comprehender to update its estimate of the sentence at the
425 sentence offset instead of on the fly during the sentence.

426 Our model achieves human-like speech (language) comprehension in that it applies syntactic rules to
427 dynamically update values assigned to semantic roles with each incoming lemma. It does not rely on
428 a direct representation of sentences, but incrementally builds up its understanding of an utterance
429 through incorporating new evidence into current beliefs of semantic roles. We share this notion with
430 the Sentence Gestalt (SG) model of language comprehension, which achieves dynamic thematic role
431 assignment from lexical inputs using a neural network trained on linguistic stimuli produced by a
432 probabilistic generative model (13, 54). The function of situation and thematic roles in this
433 generative model are homologous to that of the context (situation) and semantic (thematic) factors
434 of our model. However, while the SG model extracts thematic information from lexical input, a
435 central feature of our model is to deploy all the hierarchies from the *online* processing of continuous
436 speech to language comprehension. The variational Bayesian approach and the gradient-based
437 algorithms we used here have two particular advantages. First, it allow us to explicitly model the
438 interactions within and between meaningful computational hierarchies, and second, they can
439 account for dynamics of neuronal activities such as local field potentials (39, 55). We therefore
440 believe our model is better suited to our goal of explaining language processing within a potentially
441 unifying account of neuronal message passing, rather than in terms of neural-like network
442 activations (see next section).

443 The behavioral (nonlinguistic) goal of language comprehension is implemented minimally in the
444 current model as the task of inferring a simple context (situation) level, which represents the basic
445 “world knowledge” necessary for resolving semantic ambiguity. To implement cognitively more
446 elaborate language tasks, the context level in the model would need to include additional elements
447 that likely involve multiple decision hierarchies (56). Yet, while a model can include an arbitrary
448 number of hierarchies, there is not an infinity of corresponding specialized brain regions.
449 Computational hierarchies, especially those of higher cognitive functions that can expand to an
450 infinite depth, are therefore likely embodied by information exchanges among a limited number of
451 functionally specialized regions, through reciprocal interactions that can theoretically implement
452 unlimited hierarchical structures using only two abstract chunking levels (57, 58). These information
453 exchanges reflect the probabilistic mappings in the comprehender’s internal model, as shown in Figs
454 2 and 3, and play an important role for linking the model’s computational principles to
455 neurophysiological data of speech information processing in the human brain.

456 **Understanding neural information transfer through divergence and** 457 **entropy**

458 Brains process internal and external information with high efficiency. Two types of information
459 theoretic metrics have been of particular interest in establishing the connection between abstract
460 information and biophysical signals to probe the brain’s information processing capacity: surprisal
461 (related to, but distinct from divergence) and entropy. Efforts in associating neurophysiological
462 responses to surprisal for next-word expectation, either based on cloze probability tests (32, 59-61)
463 or the probabilistic distribution estimated by computational models (35-37, 62-65), largely credit
464 Levy’s influential work on expectation-based comprehension (10). Levy proposed a formal
465 relationship between incremental comprehension effort and the Kullback-Leibler divergence (KLD)
466 of syntactic structure inference before and after receiving a word input W , and proved that the KLD
467 reduced to the surprisal of W given the previous word string when conditioned on a constant extra-
468 sentential context that constrains comprehension. Although these studies robustly found
469 neurophysiological correlates of word surprisal, focusing on this aggregated measure without
470 explicitly modeling probabilistic representations above the word level may not be enough to tease
471 out the influence of high-level factors on language processing as was shown in Figs 3, 4 and S2 (62).
472 High-level processes presumably explain conflicting findings across studies on evoked response (66)
473 and underlying neuronal circuits (32, 36, 61) of word surprisal, because different experimental
474 paradigms likely tap into different language processing modes, making word surprisal too coarse a
475 measure. Here, we demonstrate the possibility to explicitly model information transduction above

476 lexical processing and use KLD as a universal metric to quantify information transfer, in line with
477 some predictive coding hypothesis that propose KLD to be driving the prediction error signal
478 transmitted between cortical hierarchies (55, 67).

479 Regarding entropy, the measure of information in a system (68) that represents the uncertainty in
480 linguistic stimuli, it has drawn less interest compared to surprisal metrics (32, 36, 69, 70). There is no
481 consensus on how information is maintained between two instantaneous belief updates, and
482 entropy may be valuable in investigating the *representation* of information in the brain. Intuitively,
483 higher entropy implies greater effort (more possibilities to be maintained), and less precise
484 estimates thus weaker top-down prediction influence, but it is unclear what neural activities can
485 underpin such effects. Noninvasive whole-brain imaging may inform us when and where the effort
486 takes place given that entropy and divergence can be properly dissociated (36), whereas the
487 biophysical implementation, e.g. neuronal firing patterns, may only be revealed by invasive
488 methods.

489 By showing that information passing across different processing levels contribute in a
490 complementary manner to the variability of the neurophysiological response to speech (Fig 4), our
491 model supports the neural processing of language as hierarchically organized information passing
492 among brain areas. Both KLD and entropy, as well as bottom-up prediction errors and top-down
493 priors that can be decomposed from KLD (71), are suitable metrics for such an investigation.
494 Although no definitive conclusion has been drawn on the anatomical circuits involved in high-level
495 (semantic and beyond) message passing during speech perception, a converging view is that the
496 extraction of different hierarchical representations is distributed in networks that perform multiple
497 subprocesses in parallel (72-75). Recent temporally and spatially resolved neuroimaging studies
498 suggest that neural oscillations are a good candidate mechanism for timed information transmission
499 in these subprocesses (67, 76-78). The discrete portion of our model, or in theory any model with
500 explicit structural and timing information (11, 53), can provide a template for organizing distributed
501 oscillatory activities into functional hierarchies through correlating latency- and frequency-specific
502 neuronal dynamics with model-derived information metrics. In general, sensory inputs sampled by
503 fast (gamma) oscillation are parsed into higher-level information as phase alignments of slow (theta,
504 delta) oscillations (26, 76, 79-82), which are found to be modulated by level-specific speech
505 information (32, 36, 61) and top-down coordination of mid-range (alpha, beta) oscillations (78, 79,
506 83-87). One promising avenue that exploits both model-derived computational metrics and neural
507 oscillations to disentangle neural information transfer is via a forward model that explains the
508 neurophysiological signal as a result of input-modulated changes in direction-specific connection

509 strengths between specific neural sources (brain areas), i.e. effective connectivity (88, 89). Through
510 hypothesis testing of specific brain areas and their connectivity patterns relevant for language
511 processing, direction (top-down or bottom-up) of information transfer can be distinguished by
512 frequency band-specific induced activities (90), and the functional hierarchy as well as the
513 computational roles of different connections may be mapped by regressing their modulation gain
514 with model-derived information metrics.

515 The proposed approach is fundamentally different from a purely data-driven one that identifies
516 neural response patterns correlated with pooled activities from hidden layers of a neural network
517 trained on specific tasks of next-input predictions such as in (62, 64, 65). The brain interacts with the
518 external stimuli, whether linguistic or not, in a structured fashion that is likely reused across
519 different domains (44, 58). Thus, a clear computational interpretation of brain activity patterns
520 requires an explicit representation of such structures that is lacking in most neural network models.

521 **Future development towards natural language understanding**

522 In this work, we provide a basic model that integrates linguistic and nonlinguistic world knowledge in
523 speech perception. Though the current work focuses on resolving ambiguity in semantic role
524 assignment within a reduced language and world model, the framework of a hierarchical generative
525 model is suitable for capturing various features of human language processing. For example,
526 additional branches can be “plugged-in” onto specific levels of the current generative model to
527 enable multi-modal speech processing. One possible case is to generate continuous lip movement
528 from each syllable (91, 92), in parallel with the syllable-to-acoustic generation. The inverse
529 (comprehension) model is then equipped to deal with audiovisual speech input, and can thus
530 potentially simulate known effects including using one modality to disambiguate the other (e.g. a
531 high-precision visual processing to mitigate noisy auditory input), or processing conflicting bimodal
532 inputs (e.g. relying more on the modality that has higher precision)(92). The additional branch can
533 also be attached to the context level to generate a sequence of events, such as a car speeds up and
534 hits a streetlight, to allow the inverse model to make inference about the shared context from both
535 linguistic (speech) and nonlinguistic inputs.

536 Another important feature of language processing is learning, which is also necessary for upscaling
537 the model to reflect the wealth of linguistic and nonlinguistic knowledge mastered by a real listener.
538 Language learning can be conceptualized as consisting of two complementary components: 1)
539 learning the structure of the generative model, including the possible states of different factors and
540 syntactic rules; 2) learning the parameters of the generative model, including priors, likelihoods, and

541 precisions, which are fixed in the current model. Although it is nontrivial to extend the current
542 model to include either type of learning, they could be achieved within the framework of
543 probabilistic generative models. For the first type, a plausible algorithm of statistical parameter
544 learning of structured contextual and semantic knowledge is the one proposed for the “topic” model
545 of semantic representation (9, 52). Griffiths et al. (9) also pointed to a possible way to integrate
546 complex syntax and semantic generative models by replacing one component in a syntax model (93)
547 with such a topic model. This would allow the syntax model to determine an appropriate semantic
548 component for the current timepoint and the semantic model to generate a corresponding word,
549 which is consistent with the way semantic and syntax factors interact in our current model. More
550 recently, Beck and colleagues (94) showed that a formal equivalence of the topic model can be
551 implemented via a probabilistic (neural) population code, providing a plausible path to a neural
552 implementation of the model. The second type of learning can be viewed by updating the relevant
553 parameters within a fixed structure learned from a structure-learning model. Such an updating
554 algorithm has been implemented within the dynamic expectation maximization (DEM) framework
555 that we currently use (95). To exploit the algorithm, the current generative model needs to be
556 modified to include a relevant task and associated rewards (both external and internal), so that the
557 model can actively adjust its parameters to optimize rewards. This way, top-down predictions can
558 evolve from naïve (e.g. uniform prior as we simulated in Results) to specific.

559 Overall, this model adopts a different and complementary perspective from the rapidly developing
560 world of large-scale natural language models (19-21) in that it puts upfront the gross biological
561 factors that motivate language in the first place (96-99), rather than those that seek to match human
562 performance via selected measurements in specific tasks. Recent interesting endeavors in merging
563 these two perspectives focus on adding more “neural features”, such as longer memory span and
564 domain-general knowledge beyond language, to improve natural language models (24, 25). While
565 this strategy is useful from the viewpoint of artificial language processing, it stays relatively removed
566 from the specific biological substrates of language and hence sheds little light on how human
567 language emerged and evolved under evolutionary pressure. Here, we propose a computational
568 framework to address more directly these fundamental questions by explicitly including
569 nonlinguistic components in the model architecture and using hierarchical (as opposed to
570 aggregated) prediction as a general computational strategy. Although here we focus on a passive
571 listener, a comprehensive model of human language understanding should also consider interactive
572 aspects of language, i.e. language production and multi-person communication (12) where language
573 serves as a medium to achieve shared goals (24, 100-103).

574 **Methods**

575 **Model for speech comprehension**

576 We model speech perception by inverting a generative model of speech that is able to generate
577 semantically meaningful sentences to express possible facts about the world. Since our main goal is
578 to illustrate the cognitive aspect of speech comprehension, we use the model to simulate a semantic
579 disambiguation task similar to MacGregor et al. (32). The task assesses the semantic ambiguity early
580 in a sentence, which is disambiguated later in the sentence on half of the trials. Speech inputs to the
581 model were synthesized short sentences adapted from MacGregor et al. (32).

582 In the next section we describe the speech stimuli, present the generative model, and briefly
583 describe the approximate inversion of the generative model as well as the two information theoretic
584 measures that could be related to measurable brain activity.

585 1. Speech stimuli

586 In the original design of MacGregor and colleagues, eighty sentence sets were constructed to
587 test the subjects' neural response to semantic ambiguity and disambiguation. Each set consists
588 of four sentences in which two sentence MIDDLE WORDS crossed with two sentence final words.
589 From the two sentence middle words, one was semantically ambiguous and from the two
590 sentence final words one disambiguated the ambiguous middle word, and the other did not
591 resolve the ambiguity. For example:

592 *The man knew that one more ACE might be enough to win the tennis.*

593 *The woman hoped that one more SPRINT might be enough to win the game.*

594 The middle word was either semantically ambiguous ("ace" can be a special serve in a tennis
595 game, or a poker card) or not ("sprint" only has one meaning of fast running); the two ending
596 words either resolved the ambiguity of the middle word ("tennis" resolves "ace" to mean the
597 special serve, not the poker card) or not ("game" can refer to either poker or tennis game). We
598 chose this set as part of input stimuli to the model, but reduced the sentences to essential
599 components for simplicity:

600 *One more ACE/SPRINT wins the tennis/game.*

601 The four sentences point to a minimum of two possible contexts, i.e. the nonlinguistic
602 backgrounds where they might be generated: all combinations can result from a "tennis game"

603 context, and the ACE-game combination can additionally result from a “poker game” context.
604 Importantly, in our model the context is directly related to the interpretation of the word “ace”.
605 To balance the number of plausible sentences for each context, we added another possible mid-
606 sentence word “joker”, which unambiguously refers to a poker card in the model’s knowledge.
607 We also introduced another possible sentence structure to add syntactic variability within the
608 same contexts:

609 *One more ACE/SPRINT is surprising/enough.*

610 The two syntactic structures correspond to two different types of a sentence: the “win”
611 sentences describe an event, whereas the “is” sentences describe a property of the subject.

612 We chose a total of two sentence sets from the original design. The other set (shortened
613 version) is:

614 *That TIE/NOISE ruined the game/evening.*

615 In these sentences, the subject “tie” can either mean a piece of cloth to wear around the neck
616 (“neckband” in the model) or equal scores in a game. The ending word “game” resolves it to the
617 latter meaning, whereas “evening” does not disambiguate between the two meanings. Similar to
618 set 1, we added the possibility of property-type sentences. Table 2 lists all possible sentences
619 and their corresponding contexts within the model’s knowledge (ambiguous and resolving words
620 are highlighted).

621 The input to the model consisted of acoustic spectrograms that were created using the Praat
622 (104) speech synthesizer with British accent, male speaker 1.

623 In this work we are not focusing on timing or parsing aspects, rather on how information is
624 incorporated into the inference process in an incremental manner and how the model’s
625 estimates about a preceding word can be revised upon new evidence during speech processing.
626 Therefore, we chose the syllable as the interface unit between continuous and symbolic
627 representations, and fixed the length of the input to simplify the model construction (see details
628 in Generative model). Each sentence consists of four lemma items (single words or two-word
629 phrases), and each lemma consists of three syllables. All syllables were normalized in length by
630 reducing the acoustic signal to 200 samples.

631 Specifically, in Praat, we first synthesized full words, then separated out syllables using the
632 TextGrid function. A 6-by-200 time-frequency (TF) matrix was created for each unique syllable by
633 averaging its spectro-temporal pattern into 6 log-spaced frequency channels (roughly spanning
634 from 150 Hz to 5 kHz) and 200 time bins in the same fashion as in Hovsepyan et al. (26). Each

635 sentence input to the model was then assembled by concatenating these TF matrices in the
636 appropriate order. Since we fixed the number of syllables in each word ($N_s = 3$), words
637 consisting of fewer syllables were padded with “silence” syllables, i.e. all-zero matrices. During
638 simulation, input was provided online in that 6-by-1 vectors from the padded TF matrix
639 representing the full sentence were presented to the model one after another, at the rate of
640 1000 Hz. In effect, all syllables were normalized to the same duration of 200ms. The same TF
641 matrices were used for the construction of the generative model as speech templates (see
642 section 2c for details).

643 2. Generative model

644 The generative model goes from a nonlinguistic, abstract representation of a message defined in
645 terms of semantic roles to a linearized linguistic sentence and its corresponding sound
646 spectrogram. The main idea of the model is that listeners have knowledge about the world that
647 explains how an utterance may be generated to express a message from a speaker.

648 In this miniature world, the modeled listener knows about a number of *contexts*, the scenarios
649 under which a message is generated (to distinguish them from names given to representation
650 levels in the model, we will use *italic* to refer to factors at each level; see below). Each message
651 can either be of an “event” *type* that describes an action within the context, or of a “property”
652 *type* that expresses a characteristic of an entity that exists in the context. *Context* and *type* are
653 nonlinguistic representations maintained throughout the message but make contact with
654 linguistic entities via semantics and syntax, which jointly determine an ordered sequence of
655 lemma that then generates the acoustic signal of an utterance that evolves over time.

656 As in the real world, connections from context to semantics and semantics to lemma are not
657 one-to-one, and ambiguity arises, for example, when two semantic items can be expressed as
658 the same lemma. In this case the model can output exactly the same utterance for two different
659 messages. When the model encounters such an ambiguous sentence during inference, it will
660 make its best guess based on its knowledge when ambiguity is present (see Model inversion).
661 For illustrative purposes, we only consider a minimum number of alternatives, sufficient to
662 create ambiguity, e.g. the word “ace” only has two possible meanings in the model. Also, while
663 the model generates a finite set of possible sentences, they are obtained in a compositional
664 fashion; they are not spelled out explicitly anywhere in the model, and must be incrementally
665 constructed according to the listener’s knowledge.

666 Specifically, the generative model (Figure 1A) is organized in three hierarchically related
667 submodels that differ in their temporal organization, with each submodel providing empirical

668 priors to the subordinate submodel, which then evolves in time according to its discrete or
 669 continuous dynamics for a fixed duration (as detailed below). Overall, this organization results in
 670 six hierarchically related levels of information carried by a speech utterance, from high to low
 671 (L_1 - L_6) we refer to them as: context, semantics and syntax, lemma, syllable, acoustic, and the
 672 continuous signal represented by time-frequency (TF) patterns that stands for the speech output
 673 signal.

674 Each level in the model consists of one or more factors representing the quantities of interest
 675 (e.g., *context*, *lemma*, *syllable* ...), illustrated as rectangles in Fig 1A. We use the term “states” or
 676 hidden states to refer to the values that a factor can take (e.g. in the model the factor *context*
 677 can be in one of four states {‘poker game’, ‘tennis game’, ‘night party’, ‘racing game’}. For a
 678 complete list of factors and their possible states of context to lemma levels see Table 1).

679 As an example, to generate a sentence to describe an event under a “tennis game” *context*, the
 680 model picks “tennis serve” as the agent, “tennis game” as the patient, and “win” as their
 681 relationship. When the syntactic rule indicates that the current semantic role to be expressed
 682 should be the agent, the model selects the lemma “ace”, which is then sequentially decomposed
 683 into three syllables /eis/, /silence/, /silence/. Each syllable corresponds to eight 6-by-1 spectral
 684 vectors that are deployed in time over a period of 25 ms each. The generative model therefore
 685 generates the output of continuous TF patterns as a sequence of “chunks” of 25 ms.

686 We next describe in detail the three submodels:

687 a. Discrete non-nested: context to lemma via semantic (dependency) and syntax (linearization)

688 The context level consists of two independent factors: the *context* c and the sentence *type*
 689 Ty . Together, they determine the probability distribution of four semantic roles: the *agent*
 690 s^A , the *relation* s^R , the *patient* s^P , and the *modifier* s^M . An important assumption of the model
 691 is that states of *context*, *type* and semantic roles are maintained throughout the sentence as
 692 if they had memory. These semantic roles generate a sequence of lemmas in the
 693 subordinate level, whose order is determined by the *syntax*, itself determined by the
 694 sentence *type*. This generative model for the first to the n^{th} lemma is (\vec{s} denotes the
 695 collection of all semantic factors $\vec{s} = \{s^A, s^R, s^P, s^M\}$):

$$\begin{aligned}
 & p(w^1, \dots, w^n, syn^1, \dots, syn^n, \vec{s}, c, Ty) = \\
 & p(w^1 | syn^1, \vec{s}) \dots p(w^n | syn^n, \vec{s}) p(\vec{s} | c, Ty) p(c) p(syn^1, \dots, syn^n | Ty) p(Ty) \quad (1)
 \end{aligned}$$

698 Here, $p(c)$ is the prior distribution for the *context*. The prior probability for the sentence type
 699 $p(Ty)$ was fixed to be equal between “property” and “event”.

700 The terms $p(\vec{s}|c, Ty)$ and $p(syn^1, \dots, syn^n|Ty)$ can be further expanded as:

701
$$p(\vec{s}|c, Ty) = p(s^A|c)p(s^R|c, Ty)p(s^P|c, Ty)p(s^M|c, Ty) \quad (2)$$

702
$$p(syn^1, \dots, syn^n|Ty) = p(syn^1|Ty) \dots p(syn^n|Ty) \quad (3)$$

703 When $Ty='event'$, the sentence consists of an *agent*, a *patient*, a *relation* between the *agent*
704 and the *patient*, and a null (empty) *modifier*. When $Ty='property'$, the sentence consists of
705 an *agent*, a *modifier* that describes the *agent*, a *relation* that links the *agent* and the
706 *modifier*, and a null *patient*.

707 To translate the static context, type and semantic states into ordered lemma sequences, we
708 constructed a minimal (linear) syntax model consistent with English grammar. We constrain
709 all possible sentences to have four syntactic elements syn^1 - syn^4 , values are {'attribute',
710 'subject', 'verb', 'object', 'adjective'}. The probability of syn^n is dependent solely on Ty .

711 The syntactic element syn^i is active during the i^{th} epoch, and each possible value of the
712 syntax (except 'attribute' that directly translates to a lemma item randomly determined
713 within {'one more', 'that'}) corresponds to one semantic factor (semantic factors in the
714 model include subject, verb, object and adjective):

715 Subject—*agent* ; Verb—*relation* ; Object—*patient* ; Adjective—*modifier*

716 Thus, sentences of the "event" type are always expressed in the form of subject-verb-object
717 (SVO), and those of the "property" type in the form of subject-verb-adjective (SVadj). In the
718 i^{th} lemma epoch, the model picks the current semantic factor via the value of syn_i and finds a
719 lemma to express the value (state) of this semantic factor, using its internal knowledge of
720 mapping between abstract, nonlinguistic concepts to lexical items (summarized in the form
721 of a dictionary in Appendix I). Note that the same meaning can be expressed by more than
722 one possible lemma, and several different meanings can result in the same lemma, causing
723 ambiguity. The mapping from L_2 to L_3 can be defined separately for each lemma as follows:

- 724
- 725 ● The first lemma (w^1 the attribute) does not depend on semantics or syntax and the
726 model would generate "one more" or "that" with equal probability ($p=0.5$).
 - 727 ● w^2 and w^3 are selected according to *agent* and *patient* values, respectively, which
728 are themselves constrained by context.
 - w^4 can be either a patient or a modifier depending on Ty .

729 Prior probabilities of context and type, as well as probabilistic mappings between levels
730 (eq.2-4), are all defined in the form of multidimensional arrays. Detailed expressions and
731 default values can be found in Appendix II.

732 b. Discrete nested: lemma to spectral

733 Over time, factors periodically make probabilistic transitions between states (not necessarily
734 different). Different model levels are connected in that during the generative process,
735 discrete hidden (true) states of factors in a superordinate level (L_n) determine the initial
736 state of one or more factors in the subordinate level (L_{n+1}). The L_{n+1} factors then make a fixed
737 number of state transitions. When the L_{n+1} sequence is finished, L_n makes one state
738 transition and initiates a new sequence at L_{n+1} . State transitioning of different factors within
739 the same level occurs at the same rate. We refer to the time between two transitions within
740 each level as one **epoch** of the level. Thus, model hierarchies are temporally organized in
741 that lower levels evolve at higher rates and are nested within their superordinate levels.

742 The formal definition of the discrete generative model is shown in eq.1, where the joint
743 probability distribution of the m^{th} outcome modality (here generally denoted by o^m , specified
744 in following sections) and hidden states (generally denoted by s^n) of the n^{th} factor up to a
745 time point τ , is determined by the priors over hidden states at the initial epoch $P(s^{n,1})$, the
746 likelihood mapping from states to outcome $P(o|s)$ over time $1:\tau$, and the transition
747 probabilities between hidden states of two consecutive time points $P(s^{n,t}|s^{n,t-1})$ up to $t=\tau$:

$$748 \quad P(o^{m,1:\tau}, s^{n,1:\tau}) = P(s^{n,1}) \prod_{\tau} P(o^{m,\tau}|s^{n,\tau})P(s^{n,\tau}|s^{n,\tau-1}) \quad (4)$$

749 For lower discrete levels, representational units unfold linearly in time, and a sequence of
750 subordinate units can be entirely embedded within the duration of one superordinate
751 epoch. Therefore, the corresponding models are implemented in a uniform way: the hidden
752 state consists of a “what” factor that indicates the value of the representation unit (e.g. the
753 lemma ‘the tennis’), and a “where” factor that points to the location of the outcome
754 (syllable) within the “what” state (e.g. the 2nd location of ‘tennis’ generates syllable ‘/nis/’).
755 During one epoch at each level (e.g. the entire duration of the lemma “the tennis”), the
756 value of the “what” factor remains unchanged with its transition probabilities set to the unit
757 matrix. The “where” factor transitions from 1 to the length of the “what” factor, which is the
758 number of its subordinate units during one epoch (three syllables per lemma). Together, the
759 “what” and “where” states at the lemma level generate a sequence of syllables by
760 determining the prior for “what” and “where” states in each syllable. In the same fashion,

761 each syllable determines the prior for each spectral vector. Thus, the syllable level goes
762 through 8 epochs, and for each epoch the output of the syllable level corresponds to a
763 spectral vector of dimension (1 x 6, number of frequency channels). This single vector
764 determines the prior for the continuous submodel.

765 Such temporal hierarchy is roughly represented in Figure 1B (downward arrows).

766 Unlike L_1 and L_2 states that are maintained throughout the sentence, states of the lemma
767 level and below are “memoryless”, in that they are generated anew by superordinate states
768 at the beginning of each epoch. This allows us to simplify the model inversion (see next
769 section) using a well-established framework that exploits the variational Bayes algorithm for
770 model inversion (71). The dynamic expectation maximization (DEM) framework of Friston et
771 al. (71) consists of two parts: hidden state estimation and action selection. In our model, the
772 listener does not perform any overt action (the state estimates do not affect state
773 transitioning), therefore the action selection part is omitted.

774 Using the notation of Eq.1, parameters of the generative model are defined in the form of
775 multidimensional arrays:

776 Probabilistic mapping from hidden states to outcomes:

$$777 \quad P(o^{m,\tau} | s^{1,\tau}, \dots, s^{N,\tau}) = \text{Cat}(A^m) \quad (5)$$

778 Probabilistic transition among hidden states:

$$779 \quad P(s^{n,\tau+1} | s^{n,\tau}) = \text{Cat}(B^{n,\tau}) \quad (6)$$

780 Prior beliefs about the initial hidden states:

$$781 \quad P(s^{n,1}) = \text{Cat}(D^n) \quad (7)$$

782 For each level we define **A**, **B**, **D** matrices according to the above description of hierarchical
783 “what” and “where” factors:

- 784 ● Probability mappings (matrix **A**) from a superordinate “what” to a subordinate
785 “what” states are deterministic, e.g. $p(\text{syllb}='one' | \text{lemma}='one\ more', \text{where}=1)=1$,
786 and no mapping is needed for “where” states;
- 787 ● Transition matrices (**B**) for “what” factors are all identity matrices, indicating that
788 the hidden state does not change within single epochs of the superordinate level;
- 789 ● Transition matrices for “where” factors are off-diagonal identity matrices, allowing
790 transition from one position to the next;

- 791 • Initial states (**D**) for “what” factors are set by the superordinate level, and always
792 start at position 1 for “where” factors.

793 c. Continuous: acoustic to output

794 The addition of an acoustic level between the syllable and the continuous levels is based on
795 a recent biophysically plausible model of syllable recognition, Precoss (26). In that model
796 syllables were encoded with continuous variables and represented, as is the case here, by an
797 ordered sequence of 8 spectral vectors (each vector having six components corresponding
798 to six frequency channels). In the current model we only implemented the bottom level of
799 the Precoss model (see also (28)), which deploys spectral vectors into continuous temporal
800 patterns. Specifically, the outcome of the syllable level sets the prior over the hidden cause,
801 a spectral vector **I** that drives the continuous model. It represents a chunk of the time-
802 frequency pattern determined by the “what” and “where” states of the syllable level s^ω and
803 s^γ respectively:

$$804 \quad I_f = \sum_{\omega=1}^{N_{syl}} \sum_{\gamma=1}^8 s^\omega s^\gamma V_{f\omega\gamma} + \epsilon^I \quad (8)$$

$$805 \quad V_{f\omega\gamma} = G_f(TF_{\omega\gamma}) - W_f \tanh(TF_{\omega\gamma}) \quad (9)$$

806 The noise terms ϵ^I is random Gaussian fluctuation. $TF_{\omega\gamma}$ stands for the average of the 6x200
807 TF matrix of syllable ω in the γ^{th} window of 25 ms. **G** and **W** are 6x6 connectivity matrices
808 that ensure the spectral vector **I** determines a global attractor of the Hopfield network that
809 sets the dynamics of the 6 frequency channels. Values of **G**, **W** and a scalar rate constant κ in
810 eq. 9-10 are the same as in Precoss:

$$811 \quad \frac{dx}{dt} = \kappa[-Gx + W \tanh x + I] + \epsilon^x \quad (10)$$

812 The continuous state of **x** determines the final output of the generative model **v**, which is
813 compared to the speech input during model inversion. As **x**, **v** is a 6x1 vector:

$$814 \quad v = x + \epsilon^v \quad (11)$$

815 The precision of the output signal depends on the magnitude of the random fluctuations in
816 the model (ϵ in eq. 8, 10, 11). During model inversion, the discrepancy between the input
817 and the prediction of the generative model, i.e. the prediction error, are weighted by the
818 corresponding precisions and used to update model estimates in generalized coordinates
819 (41). We manipulated the precisions for continuous state **x** and activities of frequency
820 channels **v** to simulate from intact (HP) to impaired (LP) periphery. The precision for top-

821 down priors from the syllable level, P_s , was kept high for all simulations (see Table 1 for
822 values used in different conditions).

823 The continuous generative model and its inversion were implemented using the ADEM
824 routine in the SPM12 software package (105), which integrates a generative process of
825 action. Because we focus on passive listening rather than interacting with the external
826 world, this generative process was set to identical to the generative model and without an
827 action variable. Precisions for the generative process were the same for all simulations
828 (Table 4).

829 **Table 4. Precisions**

Precision	Generative model: HP	Generative model: LP	Generative process
P^x	$\exp(16)$	$\exp(6), \exp(0), \exp(-4)$	$\exp(16)$
P^v	$\exp(16)$	$\exp(6), \exp(0), \exp(-4)$	$\exp(16)$
P^l	$\exp(8)$	$\exp(8)$	$\exp(8)$

830 3. Model inversion

831 The goal of the modeled listener is to estimate posterior probabilities of all hidden states given
832 observed evidence $p(s|o)$, which is the speech input to the model, here represented by TF
833 patterns sampled at 1000 Hz. This is achieved by the inversion of the above generative model
834 using the variational Bayesian approximation under the principle of minimizing free energy
835 (106). Although this same computational principle is applied throughout all model hierarchies,
836 the implementation is divided into three parts corresponding to the division of the generative
837 model. Because the three “submodels” are hierarchically related we follow and adapt the
838 approach proposed in (71), which shows how to invert models with hierarchically related
839 components through Bayesian model averaging. The variational Bayes approximation for each of
840 the three submodels is detailed below.

841 Overall, the scheme results in a nested estimation process (Figure 1B). For a discrete-state level
842 L_n , probability distributions over possible states within each factor are estimated at discrete
843 times over multiple inference epochs. Each epoch at level L_n starts as the estimated L_n states
844 generate predictions for initial states in the subordinate level L_{n+1} , and ends after a fixed number
845 of state transitions (epochs) at L_{n+1} . State estimations for L_n are then updated using the
846 discrepancy between the predicted and observed L_{n+1} states. The L_n factors make transitions into
847 the next epoch immediately following the update, and the same process is repeated with the
848 updated estimation. Different model hierarchies (from L_2 on) are nested in that the observed L_{n+1}
849 states are state estimations integrating information from L_{n+2} with the same alternating

850 prediction-update paradigm, but in a faster timescale. A schematic of such a hierarchical
 851 prediction-update process is illustrated in Figure 1B.

852 Since levels “lemma” to the continuous acoustic output conform to the class of generative
 853 models considered in (71), we use their derived gradient descent equations and
 854 implementation. Levels “context” and “semantic and syntax” do not conform to the same class
 855 of discrete models (due to their memory component and non-nested temporal characteristics);
 856 we therefore derived the corresponding gradient descent equations based on free energy
 857 minimization for our specific model of the top two levels Equations 2-4 (see Appendix III for the
 858 derivation) and incorporated them into the general framework of DEM (71).

859 The variational Bayes approximation for each of the three submodels is detailed below.

860 a. Lemma to context

861 For all discrete-state levels, the free energy F is generally defined as (106):

$$862 \quad Q(s) = \arg \min_{Q(s)} F \approx P(s|o) \quad (12)$$

$$863 \quad F = E_Q[\ln Q(s) - \ln P(o|s) - \ln P(s)] \quad (13)$$

864 In eq. 12 and 13, $Q(s)$ denotes the estimated posterior probability of hidden state s , $P(o|s)$
 865 the likelihood mapping defined in the generative model, and $P(s)$ the prior probability of s .
 866 The variational equations to find the $Q(s)$ that minimizes Free energy can be solved via
 867 gradient descent. We limit the number of gradient descent iterations to 16 in each update to
 868 reflect the time constraint in neuronal processes.

869 Although context/type and semantic/syntax are modeled as two hierarchies, we assign them
 870 the same temporal scheme for the prediction-update process at the rate of lemma units, i.e.
 871 they both generate top-down predictions prior to each new lemma input, and fulfill bottom-
 872 up updates at each lemma offset. Therefore, it is convenient to define their inference
 873 process in conjunction.

874 The posterior distribution $p(\text{syn}^1, \dots, \text{syn}^n, \vec{s}, c, ST | w^1, \dots, w^n)$ is approximated by a
 875 factorized one, $Q(\text{syn}^1) \dots Q(\text{syn}^n) Q(s^1) \dots Q(s^{n_s}) Q(c) Q(ST)$, and is parameterized as
 876 follows:

$$877 \quad Q(\text{syn}^\tau) : \text{syn}_k^{(\tau)}, \text{ or } \text{Cat}(\text{syn}^{(\tau)}), k = 1, \dots, \# \text{ of possible syntactic elements}, \tau = 1, \dots, n$$

$$878 \quad Q(s^\alpha) : s_j^{(\alpha)}, \text{ or } \text{Cat}(s^{(\alpha)}), j = 1, \dots, \# \text{ of possible states for semantic factor},$$

$$879 \quad \alpha = \{A, R, P, M\}$$

880 $Q(c) : c_m$, or $Cat(c)$, $m = 1, \dots$, # of possible states for context factor

881 $Q(Ty) : Ty_a$, or $Cat(Ty)$, $a = 1, \dots$, # of possible states for sentence type

882 Here, the model observation is the probability of the word being w^τ given the observed
 883 outcome o^τ , $p(w^\tau | o^\tau)$, which is gathered from lower-level models described in next sections.
 884 We denote $p(w^\tau | o^\tau)$ by a vector W_i^τ , where τ stands for the epoch, and i indexes the word in
 885 the dictionary. At the beginning of the sentence, the model predicts the first lemma input,
 886 which is, by definition, just one of the two possible attributes, ‘one more’ or ‘that’.

$$\begin{aligned}
 887 \quad p(w^1) &= \sum_{syn^1, \vec{s}, c, Ty} p(w^1 | syn^1, \vec{s}, c, Ty) p(syn^1, \vec{s}, c, Ty) \\
 888 \quad &= \sum_{syn^1} p(w^1 | syn^1) p(syn^1) = p(w^1 | syn^1 = \text{attribute}) \quad (14)
 \end{aligned}$$

889 The lower levels then calculate $p(w^1 | o^1)$ and provide an updated W_i^1 that incorporates the
 890 observation made from the first lemma. This is passed to the top levels to update L_1 and L_2
 891 states. Following this update, the next epoch is initiated with the prediction for w^2 . Because
 892 w^2 does not directly depend on lemma inputs before and after itself, we can derive the
 893 following informed prediction of w^2 from eq.2, where prior for L_1 and L_2 factors are replaced
 894 by their updated posterior estimates:

$$\begin{aligned}
 895 \quad p(w^2) &= \sum_{syn^2, \vec{s}, c, ST} p(w^2 | syn^2, \vec{s}, c, Ty) p(syn^2, \vec{s}, c, Ty | o^1) \\
 896 \quad &\approx \sum_{syn^2, \vec{s}, Ty} p(w^2 | syn^2, \vec{s}) p(syn^2 | Ty) Q^{(1)}(\vec{s}) Q^{(1)}(c) Q^{(1)}(Ty) \quad (15)
 \end{aligned}$$

897 Where we used:

$$\begin{aligned}
 898 \quad p(syn^2, \vec{s}, c, Ty | o^1) &\approx p(syn^2 | Ty) Q(\vec{s}, c, Ty | o^1) \\
 899 \quad &= p(syn^2 | Ty) Q^{(1)}(\vec{s}) Q^{(1)}(c) Q^{(1)}(Ty)
 \end{aligned}$$

900 During the second epoch, the model receives input of the second lemma and updates the
 901 estimation of W_i^2 . The updated W_i^2 is then exploited to update L_1 and L_2 states, which in turn
 902 provides the prediction for w^3 . The process is repeated until the end of the sentence.

903 The updating of L_1 and L_2 states, i.e. the estimation of their posterior probabilities after
 904 receiving the n^{th} lemma input relies on the minimization of the total free energy $F_{1,2}$ of the
 905 two levels (L_1, L_2)

$$\begin{aligned}
 906 \quad F_{1,2} \equiv & \sum_{syn^1:syn^n, \vec{s}, c, Ty} Q(syn^1, \dots, syn^n, \vec{s}, c, Ty) \left[\ln Q(syn^1, \dots, syn^n, \vec{s}, c, Ty) \right. \\
 907 & \left. - \sum_{w^1:w^n} Q(w^1, \dots, w^n) \ln p(w^1, \dots, w^n, syn^1, \dots, syn^n, \vec{s}, c, Ty) \right] \quad (16)
 \end{aligned}$$

908 The expanded expression of $F_{1,2}$ and derivation of the gradient descent equations can be
 909 found in Appendix III.

910 b. Spectral to lemma

911 The memoryless property of lower-level (lemma and below) states implies that the
 912 observation from the previous epoch does not directly affect the prediction for the new
 913 epoch, only indirectly through the evidence accumulated at superordinate levels. The
 914 framework from Friston et al. (71) is suitable for such construction. It uses the same
 915 algorithm of free-energy (inserting eq. 5-7 to eq. 12-13) minimization for posterior
 916 estimation, but this time there is conditional independence between factors in the same
 917 level. We implemented this part of the model by adapting the variational Bayesian routine in
 918 the DEM toolbox from the SPM12 software package.

919 c. Continuous to spectral

920 To enable the information exchange between the continuous and higher discrete levels that
 921 were not accounted for in (26), we implemented the inversion of the spectral-to-continuous
 922 generative model using the “mixed model” framework in (71). Essentially, the dynamics of
 923 spectral fluctuation determined by each spectral vector \mathbf{I} (eq.8) is treated as a separate
 924 model of continuous trajectories, and the posterior estimation of \mathbf{I} constitutes post-hoc
 925 model comparison that minimizes free energy in the continuous format. For a specific model
 926 m represented by spectral vector I_m , the free energy $F(t)_m$ can be computed as (adapted
 927 from (71)):

$$928 \quad F(t)_m = -\ln P(o_m) - \int_0^T L(t)_m dt \quad (17)$$

$$929 \quad L(t)_m = \ln P(o(t)|I_m) - \ln P(o(t)|I) \quad (18)$$

930 $P(o_m)$ indicates the likelihood for the m^{th} spectral vector (discrete). $P(o(t)|I_m)$ is the likelihood
 931 of observing the continuous input $o(t)$ given the m^{th} \mathbf{I} vector, and $P(o(t)|I)$ is the averaged
 932 likelihood over all possible \mathbf{I} vectors. In this way, the model compares the top-down
 933 prediction of \mathbf{I} and the estimate derived from the bottom-up evidence of integrated acoustic
 934 input over 25ms. Detailed explanation of the algorithm can be found in previous studies (71,

935 107). The software implementation was also adapted from existing routines in the DEM
936 toolbox of SPM12 (105).

937 **Information theoretic metrics**

938 Two metrics were derived from the belief updating process just described: the Kullback-Leibler (KL)
939 divergence (Div), which characterizes the discrepancy between the current and previous state
940 estimates of a factor, and entropy H that characterizes the uncertainty of the current state estimates
941 of the factor. We denote the posterior probability of the i^{th} possible state of an arbitrary factor at
942 time point τ as q_i^τ . The divergence and entropy are defined as:

$$943 \quad Div^\tau = - \sum_i q_i^\tau \ln q_i^{\tau-1} + \sum_i q_i^\tau \ln q_i^\tau \quad (19)$$

$$944 \quad H^\tau = - \sum_i q_i^\tau \ln q_i^\tau \quad (20)$$

945 These two (non-orthogonal) metrics provide a qualitative summary of the model response that can
946 be linked to neurophysiological signals (see Result and Discussion).

947 **Model guided MEG data analysis**

948 **Next-word prediction statistics from GPT-2 model**

949 We implemented a transformer pre-trained language model, GPT-2 (20) in Google Colab (108), to
950 obtain word prediction statistics of the sentence stimuli. The model is trained on ~40 GB text data
951 and generates next-word predictions given arbitrary sentence contexts. Inputs to the model were
952 sentences taken from (32), each sentence consisting of four parts (see Table 3 for an example set): a
953 lead-in phrase, a target word, a bridge phrase, and a resolution word. For every lead-in phrase, four
954 variations were played by crossing two different Target words and two different Resolution words.

955 **Target:** either with or without semantic ambiguity (Ambiguous vs. Unambiguous).

956 **Resolution:** either resolves the semantic ambiguity of the Ambiguous Target, or not (Resolve vs.
957 Unresolve).

958 For each set of (Target \times Resolution) combination, two versions of the lead-in phrase were available.
959 However, only one of the two lead-ins in each set was used for each subject in the MEG experiment,
960 i.e. each set of (Target \times Resolution) combination was played only once. Therefore, we averaged the
961 GPT-2 prediction metrics for the two versions. The bridge phrase was the same within each set,
962 regardless of other parts of the sentence.

963 The original speech stimuli in (32) contained sentence sets where the Target words were ambiguous
964 between two phonetically identical but morphologically different words. These sets were removed
965 for the GPT-2 analysis as well as for the MEG data analysis, resulting in 58 out of 80 sets.

966 Probability distributions of the next-word prediction of GPT-2 were obtained for two time points to
967 calculate the prediction entropy and surprisal, respectively:

968 1. After Target, i.e. the input to GPT-2 is [lead in] + [target]

969 We use the entropy H of this prediction as a proxy for the (semantic) ambiguity of the target
970 word, with the hypothesis that if a word has multiple meanings, different meanings will
971 predict different next words with similar probabilities, resulting in a flatter distribution
972 compared to the prediction from its unambiguous counterpart. H is calculated as follows,
973 where i indexes all words in the dictionary:

$$974 \quad H = - \sum_i p_i \ln p_i$$

975 2. Before Resolution, i.e. the input to GPT-2 is [lead in] + [target] + [bridge]

976 We calculate the surprisal S for each resolution word from the prediction probability as
977 follows, where r is the index for the resolution word in the dictionary:

$$978 \quad S = - \ln p_r$$

979 This surprisal is equivalent to the KL divergence of the posterior distribution after the
980 resolution word, because the distribution has collapsed to $p=1$ for the received word and 0
981 elsewhere.

982 **MEG sensor space analysis**

983 The MEEG module in SPM12 (105) was used for the MEG data preprocessing. Statistical analysis and
984 plotting of the preprocessed results were performed with the Fieldtrip Toolbox (109). We first
985 performed the identical preprocessing as MacGregor et al. (32) on head-adjusted raw MEG
986 responses to the 58 selected sentence sets for all 16 subjects. Briefly, raw recordings were first
987 bandpass filtered between 0.1 and 40 Hz, then epoched at the offsets of each keyword (Target or
988 Resolution). After baseline correction and the rejection of bad trials, combined gradiometer (RMS of
989 each of the 102 gradiometer pairs) responses were cropped into shorter time windows (-0.2~0.8s for
990 the Target offset, -0.5~1s for the Resolution offset) and averaged across trials for each subject. For
991 averaging, trials were split in the following way that allow for statistical tests for both the GPT-2
992 prediction metrics and the linguistic metrics of interest, i.e. semantic ambiguity at the Target offset
993 and resolution at the Resolution offset:

994 1. Target

995 Sentences were split into two groups: 1. The GPT-2 entropy for the Ambiguous word was
996 larger than the entropy for the Unambiguous word (Amb1, Uam1), and 2. The GPT-2 entropy
997 for the Ambiguous word was smaller than for the Unambiguous word (Amb2, Uam2).

998 2. Resolution

999 Sentences containing the Resolve words were split into two groups: 1. The GPT surprisal of
1000 the Resolve word following the Ambiguous target was larger than the Resolve word
1001 following the Unambiguous target (Res_Amb1, Res_Uam1), and 2. The GPT surprisal of the
1002 Resolve word following the Ambiguous target was smaller than following the Unambiguous
1003 target (Res_Amb2, Res_Uam2).

1004 To assess the effects of linguistic and GPT-2 metrics on the combined gradiometer data, we
1005 constructed the following four contrasts:

- 1006 1. [Amb1 + Amb2] vs. [Uam1 + Uam2]: effect of semantic ambiguity.
- 1007 2. [Amb1 + Uam2] vs. [Amb2 + Uam1]: effect of GPT-2 prediction entropy.
- 1008 3. [Res_Amb1 + Res_Amb2] vs. [Res_Uam1 + Res_Uam2]: effect of preceding ambiguity.
- 1009 4. [Res_Amb1 + Res_Uam2] vs. [Res_Uam1 + Res_Amb2]: effect of GPT-2 prediction surprisal.

1010 To test for differences between the two conditions within each contrast, we first took the average of
1011 the two averages in each condition within individual subjects, e.g. (Amb1 + Amb2)/2 for the
1012 ambiguous condition in contrast 1. This yields one sensor \times time response per condition and per
1013 subject. We then performed a paired t-test across subjects for each sensor and time point, resulting
1014 in a 2D parametric map of the test statistic. Clusters of sensors with $p_s < 0.05$ were identified on this
1015 map, each including at least 2 neighboring sensors. The statistical significance of each cluster was
1016 evaluated by comparing the maximum t-statistic of the cluster to a null distribution generated by
1017 randomly permuting the condition labels within each subject (5000 times across all 16 subjects). The
1018 cluster-level p-value (p_c) was the proportion of the t statistic in the permutation distribution larger
1019 than the maximum t statistic of the selected cluster. None of the clusters identified by the t-test
1020 survived the permutation test, therefore we report the five clusters with the highest t-statistics for
1021 the positive effect in each contrast. We also computed Cohen's d (110) from the grand average over
1022 time and across subjects of all the 102 combined gradiometer channel to evaluate the effect size of
1023 each contrast at single gradiometer pairs.

1024

1025 **Acknowledgements**

1026 We thank B. Bickel, S. van Ommen, D. Poeppel for critical feedback, NCCR TTF Data Science for
1027 support on the GPT-2 model, and E. Holmes for advice on the SPM software. This work was funded
1028 by Swiss National Science Foundation (grant number 320030B_182855) and NCCR Evolving
1029 Language, Swiss National Science Foundation Agreement #51NF40_180888.

1030

1031 **Data and code availability**

1032 Custom MATLAB code and simulation data will be made available upon request (mailto:
1033 yaqing.su@unige.ch).

1034 Reference

- 1035 1. Christiansen MH, Chater N. The Now-or-Never bottleneck: A fundamental constraint on
1036 language. *Behavioral and Brain Sciences*. 2016;39.
- 1037 2. Tanenhaus MK, Spiveyknowlton MJ, Eberhard KM, Sedivy JC. Integration of Visual and
1038 Linguistic Information in Spoken Language Comprehension. *Science*. 1995;268(5217):1632-4.
- 1039 3. Levinson SE. Continuously variable duration hidden Markov models for automatic speech
1040 recognition. *Computer Speech & Language*. 1986;1(1):29-45.
- 1041 4. McClelland JL, Elman JL. The Trace Model of Speech-Perception. *Cognitive Psychol.*
1042 1986;18(1):1-86.
- 1043 5. Norris D. Shortlist - a Connectionist Model of Continuous Speech Recognition. *Cognition*.
1044 1994;52(3):189-234.
- 1045 6. LeCun Y, Bengio Y. Convolutional networks for images, speech, and time series. *The*
1046 *handbook of brain theory and neural networks*. 1995;3361(10):1995.
- 1047 7. Friston KJ, Sajid N, Quiroga-Martinez DR, Parr T, Price CJ, Holmes E. Active listening. *Hearing*
1048 *Res*. 2021;399.
- 1049 8. Elman JL. Finding Structure in Time. *Cognitive Sci*. 1990;14(2):179-211.
- 1050 9. Griffiths TL, Steyvers M, Tenenbaum JB. Topics in semantic representation. *Psychol Rev*.
1051 2007;114(2):211-44.
- 1052 10. Levy R. Expectation-based syntactic comprehension. *Cognition*. 2008;106(3):1126-77.
- 1053 11. Martin AE, Doumas LA. A mechanism for the cortical computation of hierarchical linguistic
1054 structure. *PLoS Biol*. 2017;15(3):e2000663.
- 1055 12. Friston KJ, Parr T, Yufik Y, Sajid N, Price CJ, Holmes E. Generative models, linguistic
1056 communication and active inference. *Neurosci Biobehav R*. 2020;118:42-64.
- 1057 13. Stjohn MF, McClelland JL. Learning and Applying Contextual Constraints in Sentence
1058 Comprehension. *Artif Intell*. 1990;46(1-2):217-57.
- 1059 14. Warren RM. Perceptual restoration of missing speech sounds. *Science*. 1970;167(3917):392-
1060 3.
- 1061 15. Sohoglu E, Peelle JE, Carlyon RP, Davis MH. Predictive top-down integration of prior
1062 knowledge during speech perception. *J Neurosci*. 2012;32(25):8443-53.
- 1063 16. Leonard MK, Baud MO, Sjerps MJ, Chang EF. Perceptual restoration of masked speech in
1064 human cortex. *Nat Commun*. 2016;7.
- 1065 17. Swinney DA. Lexical Access during Sentence Comprehension - (Re)Consideration of Context
1066 Effects. *J Verb Learn Verb Be*. 1979;18(6):645-59.
- 1067 18. Rodd JM, Davis MH, Johnsrude IS. The neural mechanisms of speech comprehension: fMRI
1068 studies of semantic ambiguity. *Cereb Cortex*. 2005;15(8):1261-9.
- 1069 19. Devlin J, Chang M-W, Lee K, Toutanova K. Bert: Pre-training of deep bidirectional
1070 transformers for language understanding. *arXiv preprint arXiv:181004805*. 2018.
- 1071 20. Radford A, Wu J, Child R, Luan D, Amodei D, Sutskever I. Language models are unsupervised
1072 multitask learners. *OpenAI blog*. 2019;1(8):9.
- 1073 21. Brown TB, Mann B, Ryder N, Subbiah M, Kaplan J, Dhariwal P, et al. Language models are
1074 few-shot learners. *arXiv preprint arXiv:200514165*. 2020.
- 1075 22. Floridi L, Chiriatti M. GPT-3: Its Nature, Scope, Limits, and Consequences. *Mind Mach*.
1076 2020;30(4):681-94.
- 1077 23. Lake BM, Murphy GL. Word Meaning in Minds and Machines. *Psychological Review*. 2021.
- 1078 24. Bender EM, Koller A, editors. Climbing towards NLU: On meaning, form, and understanding
1079 in the age of data. *Proceedings of the 58th Annual Meeting of the Association for Computational*
1080 *Linguistics*; 2020.
- 1081 25. McClelland JL, Hill F, Rudolph M, Baldridge J, Schutze H. Placing language in an integrated
1082 understanding system: Next steps toward human-level performance in neural language models. *P*
1083 *Natl Acad Sci USA*. 2020;117(42):25966-74.

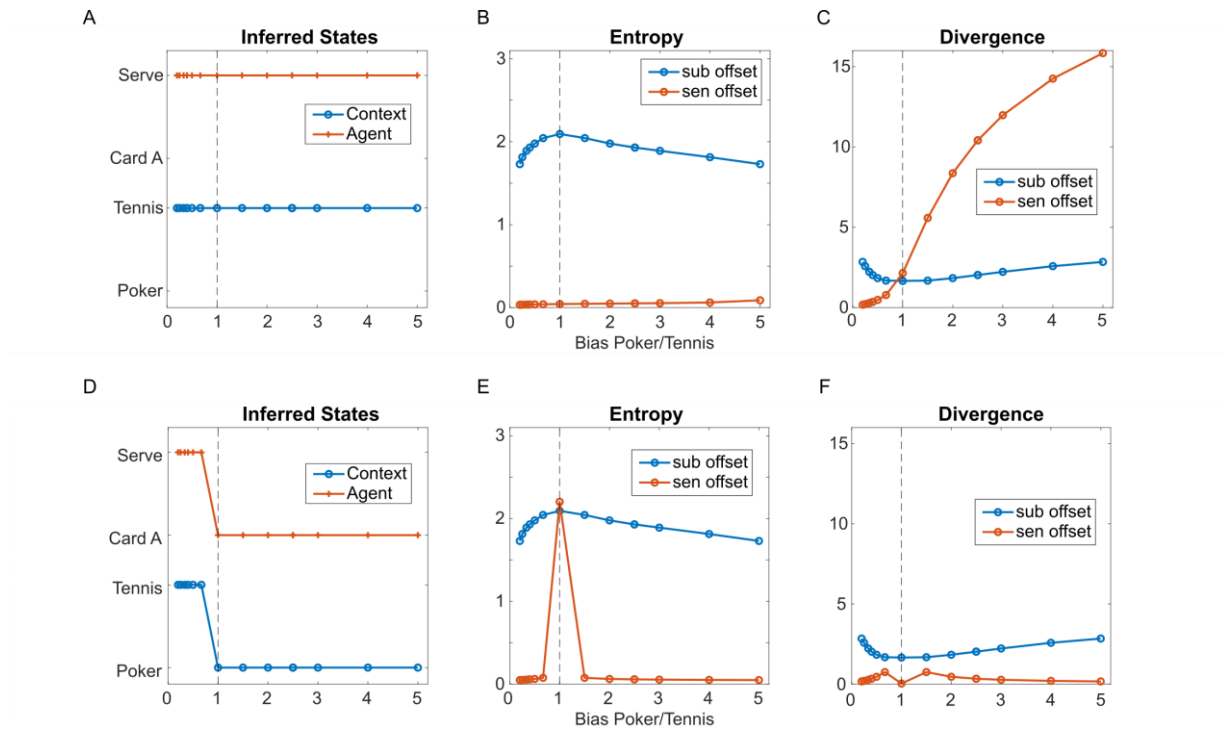
- 1084 26. Hovsepian S, Olasagasti I, Giraud AL. Combining predictive coding and neural oscillations
1085 enables online syllable recognition in natural speech. *Nat Commun.* 2020;11(1).
- 1086 27. Yildiz IB, Kiebel SJ. A Hierarchical Neuronal Model for Generation and Online Recognition of
1087 Birdsongs. *Plos Comput Biol.* 2011;7(12).
- 1088 28. Yildiz IB, von Kriegstein K, Kiebel SJ. From Birdsong to Human Speech Recognition: Bayesian
1089 Inference on a Hierarchy of Nonlinear Dynamical Systems. *Plos Comput Biol.* 2013;9(9).
- 1090 29. Rao RPN, Ballard DH. Predictive coding in the visual cortex: a functional interpretation of
1091 some extra-classical receptive-field effects. *Nat Neurosci.* 1999;2(1):79-87.
- 1092 30. Friston KJ. The free-energy principle: a rough guide to the brain? *Trends in Cognitive*
1093 *Sciences.* 2009;13(7):293-301.
- 1094 31. Clark A. Whatever next? Predictive brains, situated agents, and the future of cognitive
1095 science. *Behavioral and brain sciences.* 2013;36(3):181-204.
- 1096 32. MacGregor LJ, Rodd JM, Gilbert RA, Hauk O, Sohoglu E, Davis MH. The Neural Time Course of
1097 Semantic Ambiguity Resolution in Speech Comprehension. *J Cognitive Neurosci.* 2020;32(3):403-25.
- 1098 33. Greenberg S, Carvey H, Hitchcock L, Chang SY. Temporal properties of spontaneous speech -
1099 a syllable-centric perspective. *J Phonetics.* 2003;31(3-4):465-85.
- 1100 34. Broderick MP, Anderson AJ, Lalor EC. Semantic Context Enhances the Early Auditory
1101 Encoding of Natural Speech. *Journal of Neuroscience.* 2019;39(38):7564-75.
- 1102 35. Koskinen M, Kurimo M, Gross J, Hyvarinen A, Hari R. Brain activity reflects the predictability
1103 of word sequences in listened continuous speech. *Neuroimage.* 2020;219:116936.
- 1104 36. Donhauser PW, Baillet S. Two Distinct Neural Timescales for Predictive Speech Processing.
1105 *Neuron.* 2020;105(2):385-93 e9.
- 1106 37. Goldstein A, Zada Z, Buchnik E, Schain M, Price A, Aubrey B, et al. Thinking ahead: prediction
1107 in context as a keystone of language in humans and machines. *bioRxiv.* 2021:2020.12. 02.403477.
- 1108 38. Heilbron M, Armeni K, Schoffelen JM, Hagoort P, de Lange FP. A hierarchy of linguistic
1109 predictions during natural language comprehension. *Proc Natl Acad Sci U S A.*
1110 2022;119(32):e2201968119.
- 1111 39. Da Costa L, Parr T, Sengupta B, Friston K. Neural Dynamics under Active Inference:
1112 Plausibility and Efficiency of Information Processing. *Entropy-Switz.* 2021;23(4).
- 1113 40. Peelle JE. Listening Effort: How the Cognitive Consequences of Acoustic Challenge Are
1114 Reflected in Brain and Behavior. *Ear Hearing.* 2018;39(2):204-14.
- 1115 41. Friston KJ, Trujillo-Barreto N, Daunizeau J. DEM: A variational treatment of dynamic systems.
1116 *Neuroimage.* 2008;41(3):849-85.
- 1117 42. Payne JW, Bettman JR, Johnson EJ. Adaptive Strategy Selection in Decision-Making. *J Exp*
1118 *Psychol Learn.* 1988;14(3):534-52.
- 1119 43. Eckert MA, Teubner-Rhodes S, Vaden KI. Is Listening in Noise Worth It? The Neurobiology of
1120 Speech Recognition in Challenging Listening Conditions. *Ear Hearing.* 2016;37:101s-10s.
- 1121 44. Chambon V, Domenech P, Jacquet PO, Barbalat G, Bouton S, Pacherie E, et al. Neural coding
1122 of prior expectations in hierarchical intention inference. *Sci Rep-Uk.* 2017;7.
- 1123 45. Parr T, Rees G, Friston KJ. Computational Neuropsychology and Bayesian Inference. *Front*
1124 *Hum Neurosci.* 2018;12.
- 1125 46. Altmann GTM, Mirkovic J. Incrementality and Prediction in Human Sentence Processing.
1126 *Cognitive Sci.* 2009;33(4):583-609.
- 1127 47. Kutas M, Federmeier KD. Electrophysiology reveals semantic memory use in language
1128 comprehension. *Trends in Cognitive Sciences.* 2000;4(12):463-70.
- 1129 48. Unsworth N, McMillan BD. Mind Wandering and Reading Comprehension: Examining the
1130 Roles of Working Memory Capacity, Interest, Motivation, and Topic Experience. *J Exp Psychol Learn.*
1131 2013;39(3):832-42.
- 1132 49. Tanenhaus MK, Carlson G, Trueswell JC. The Role of Thematic Structures in Interpretation
1133 and Parsing. *Lang Cognitive Proc.* 1989;4(3-4):Si211-Si34.
- 1134 50. Altmann GTM. Thematic role assignment in context. *J Mem Lang.* 1999;41(1):124-45.

- 1135 51. McRae K, Ferretti TR, Amyote L. Thematic roles as verb-specific concepts. *Lang Cognitive*
1136 *Proc.* 1997;12(2-3):137-76.
- 1137 52. Blei DM, Griffiths TL, Jordan MI, Tenenbaum JB, editors. Hierarchical topic models and the
1138 nested Chinese restaurant process. NIPS; 2003.
- 1139 53. Martin AE. A Compositional Neural Architecture for Language. *J Cogn Neurosci.*
1140 2020;32(8):1407-27.
- 1141 54. Rabovsky M, Hansen SS, McClelland JL. Modelling the N400 brain potential as change in a
1142 probabilistic representation of meaning. *Nat Hum Behav.* 2018;2(9):693-705.
- 1143 55. Friston KJ, Kiebel S. Cortical circuits for perceptual inference. *Neural Networks.*
1144 2009;22(8):1093-104.
- 1145 56. Koechlin E, Summerfield C. An information theoretical approach to prefrontal executive
1146 function. *Trends in Cognitive Sciences.* 2007;11(6):229-35.
- 1147 57. Koechlin E, Jubault T. Broca's area and the hierarchical organization of human behavior.
1148 *Neuron.* 2006;50(6):963-74.
- 1149 58. Rouault M, Koechlin E. Prefrontal function and cognitive control: from action to language.
1150 *Curr Opin Behav Sci.* 2018;21:106-11.
- 1151 59. DeLong KA, Urbach TP, Kutas M. Probabilistic word pre-activation during language
1152 comprehension inferred from electrical brain activity. *Nat Neurosci.* 2005;8(8):1117-21.
- 1153 60. Wang L, Hagoort P, Jensen O. Gamma Oscillatory Activity Related to Language Prediction. *J*
1154 *Cognitive Neurosci.* 2018;30(8):1075-85.
- 1155 61. Mamashli F, Khan S, Obleser J, Friederici AD, Maess B. Oscillatory dynamics of cortical
1156 functional connections in semantic prediction. *Hum Brain Mapp.* 2019;40(6):1856-66.
- 1157 62. Caucheteux C, King JR. Brains and algorithms partially converge in natural language
1158 processing. *Commun Biol.* 2022;5(1).
- 1159 63. Heilbron M, Armeni K, Schoffelen J-M, Hagoort P, de Lange FP. A hierarchy of linguistic
1160 predictions during natural language comprehension. *bioRxiv.* 2021:2020.12. 03.410399.
- 1161 64. Schrimpf M, Blank IA, Tuckute G, Kauf C, Hosseini EA, Kanwisher N, et al. The neural
1162 architecture of language: Integrative modeling converges on predictive processing. *P Natl Acad Sci*
1163 *USA.* 2021;118(45).
- 1164 65. Caucheteux C, Gramfort A, King JR. Deep language algorithms predict semantic
1165 comprehension from brain activity. *Sci Rep-Uk.* 2022;12(1).
- 1166 66. Kuperberg GR. Neural mechanisms of language comprehension: Challenges to syntax. *Brain*
1167 *Res.* 2007;1146:23-49.
- 1168 67. Bastos AM, Usrey WM, Adams RA, Mangun GR, Fries P, Friston KJ. Canonical Microcircuits
1169 for Predictive Coding. *Neuron.* 2012;76(4):695-711.
- 1170 68. Shannon CE. A Mathematical Theory of Communication. *Bell Syst Tech J.* 1948;27(3):379-
1171 423.
- 1172 69. Willems RM, Frank SL, Nijhof AD, Hagoort P, van den Bosch A. Prediction During Natural
1173 Language Comprehension. *Cereb Cortex.* 2016;26(6):2506-16.
- 1174 70. Gwilliams L, King J-R, Marantz A, Poeppel D. Neural dynamics of phoneme sequencing in real
1175 speech jointly encode order and invariant content. *bioRxiv.* 2020:2020.04.04.025684.
- 1176 71. Friston KJ, Parr T, de Vries B. The graphical brain: Belief propagation and active inference.
1177 *Netw Neurosci.* 2017;1(4):381-414.
- 1178 72. Egorova N, Shtyrov Y, Pulvermuller F. Early and parallel processing of pragmatic and
1179 semantic information in speech acts: neurophysiological evidence. *Front Hum Neurosci.* 2013;7.
- 1180 73. Fedorenko E, Scott TL, Brunner P, Coon WG, Pritchett B, Schalk G, et al. Neural correlate of
1181 the construction of sentence meaning. *P Natl Acad Sci USA.* 2016;113(41):E6256-E62.
- 1182 74. Pulvermuller F. Neural reuse of action perception circuits for language, concepts and
1183 communication. *Prog Neurobiol.* 2018;160:1-44.
- 1184 75. Fairs A, Michelas A, Dufour S, Strijkers K. The Same Ultra-Rapid Parallel Brain Dynamics
1185 Underpin the Production and Perception of Speech. *Cerebral Cortex Communications.* 2021;2(3).

- 1186 76. Giraud AL, Poeppel D. Cortical oscillations and speech processing: emerging computational
1187 principles and operations. *Nat Neurosci.* 2012;15(4):511-7.
- 1188 77. Giraud AL, Arnal LH. Hierarchical Predictive Information Is Channeled by Asymmetric
1189 Oscillatory Activity. *Neuron.* 2018;100(5):1022-4.
- 1190 78. Bastos AM, Lundqvist M, Waite AS, Kopell N, Miller EK. Layer and rhythm specificity for
1191 predictive routing. *Proc Natl Acad Sci U S A.* 2020;117(49):31459-69.
- 1192 79. Arnal LH, Giraud AL. Cortical oscillations and sensory predictions. *Trends Cogn Sci.*
1193 2012;16(7):390-8.
- 1194 80. Ding N, Melloni L, Zhang H, Tian X, Poeppel D. Cortical tracking of hierarchical linguistic
1195 structures in connected speech. *Nat Neurosci.* 2016;19(1):158-64.
- 1196 81. Rimmele JM, Poeppel D, Ghitza O. Acoustically Driven Cortical δ Oscillations Underpin
1197 Prosodic Chunking. *Eneuro.* 2021;8(4).
- 1198 82. Lakatos P, Gross J, Thut G. A New Unifying Account of the Roles of Neuronal Entrainment.
1199 *Curr Biol.* 2019;29(18):R890-R905.
- 1200 83. Fontolan L, Morillon B, Liegeois-Chauvel C, Giraud AL. The contribution of frequency-specific
1201 activity to hierarchical information processing in the human auditory cortex. *Nat Commun.* 2014;5.
- 1202 84. Pefkou M, Arnal LH, Fontolan L, Giraud AL. theta-Band and beta-Band Neural Activity
1203 Reflects Independent Syllable Tracking and Comprehension of Time-Compressed Speech. *Journal of*
1204 *Neuroscience.* 2017;37(33):7930-8.
- 1205 85. Murphy E. Interfaces (travelling oscillations)+ recursion (delta-theta code)= language. *The*
1206 *Talking Species: Perspectives on the Evolutionary, Neuronal and Cultural Foundations of Language,*
1207 eds E Luef and M Manuela (Graz: Unipress Graz Verlag). 2018:251-69.
- 1208 86. Meyer L, Sun Y, Martin AE. Synchronous, but not entrained: exogenous and endogenous
1209 cortical rhythms of speech and language processing. *Lang Cogn Neurosci.* 2020;35(9):1089-99.
- 1210 87. Hovsepyan S, Olasagasti I, Giraud A-L. Rhythmic modulation of prediction errors: a possible
1211 role for the beta-range in speech processing. *bioRxiv.* 2022:2022.03.28.486037.
- 1212 88. Friston KJ. Functional and effective connectivity: a review. *Brain connectivity.* 2011;1(1):13-
1213 36.
- 1214 89. Kiebel SJ, Garrido MI, Moran R, Chen CC, Friston KJ. Dynamic Causal Modeling for EEG and
1215 MEG. *Human Brain Mapping.* 2009;30(6):1866-76.
- 1216 90. Chen CC, Kiebel SJ, Friston KJ. Dynamic causal modelling of induced responses. *Neuroimage.*
1217 2008;41(4):1293-312.
- 1218 91. Pelachaud C, Badler NI, Steedman M. Generating facial expressions for speech. *Cognitive Sci.*
1219 1996;20(1):1-46.
- 1220 92. Olasagasti I, Bouton S, Giraud AL. Prediction across sensory modalities: A
1221 neurocomputational model of the McGurk effect. *Cortex.* 2015;68:61-75.
- 1222 93. Griffiths T, Steyvers M, Blei D, Tenenbaum J. Integrating topics and syntax. *Advances in*
1223 *neural information processing systems.* 2004;17.
- 1224 94. Beck J, Heller K, Pouget A. Complex inference in neural circuits with probabilistic population
1225 codes and topic models. 2012.
- 1226 95. Friston KJ, Lin M, Frith CD, Pezzulo G, Hobson JA, Ondobaka S. Active Inference, Curiosity and
1227 Insight. *Neural Comput.* 2017;29(10):2633-83.
- 1228 96. Hauser MD, Chomsky N, Fitch WT. The faculty of language: What is it, who has it, and how
1229 did it evolve? *Science.* 2002;298(5598):1569-79.
- 1230 97. Corballis MC. The Evolution of Language. *Ann Ny Acad Sci.* 2009;1156:19-43.
- 1231 98. Greenfield PM. Language, Tools, and Brain - the Ontogeny and Phylogeny of Hierarchically
1232 Organized Sequential Behavior. *Behavioral and Brain Sciences.* 1991;14(4):531-50.
- 1233 99. Fitch WT. Evolutionary Developmental Biology and Human Language Evolution: Constraints
1234 on Adaptation. *Evol Biol.* 2012;39(4):613-37.
- 1235 100. Galantucci B, Fowler CA, Turvey MT. The motor theory of speech perception reviewed (vol
1236 13, pg 361, 2006). *Psychon B Rev.* 2006;13(4):742-.

- 1237 101. Hickok G, Poeppel D. Opinion - The cortical organization of speech processing. *Nat Rev*
1238 *Neurosci.* 2007;8(5):393-402.
- 1239 102. Pulvermuller F, Fadiga L. Active perception: sensorimotor circuits as a cortical basis for
1240 language. *Nat Rev Neurosci.* 2010;11(5):351-60.
- 1241 103. Castellucci GA, Kovach CK, Howard MA, Greenlee JDW, Long MA. A speech planning network
1242 for interactive language use. *Nature.* 2022.
- 1243 104. Boersma PW, David. Praat: doing phonetics by computer. 2021.
- 1244 105. Neuroimaging WTCf. SPM12. 2014.
- 1245 106. Friston KJ, Kilner J, Harrison L. A free energy principle for the brain. *J Physiol-Paris.*
1246 2006;100(1-3):70-87.
- 1247 107. Friston KJ, Penny W. Post hoc Bayesian model selection. *Neuroimage.* 2011;56(4):2089-99.
- 1248 108. Bisong E. Google Colaboratory. Building Machine Learning and Deep Learning Models on
1249 Google Cloud Platform: A Comprehensive Guide for Beginners. Berkeley, CA: Apress; 2019. p. 59-64.
- 1250 109. Oostenveld R, Fries P, Maris E, Schoffelen JM. FieldTrip: Open Source Software for Advanced
1251 Analysis of MEG, EEG, and Invasive Electrophysiological Data. *Comput Intel Neurosc.* 2011;2011.
- 1252 110. Cohen J. Statistical power analysis for the behavioral sciences: Routledge; 2013.
- 1253
- 1254

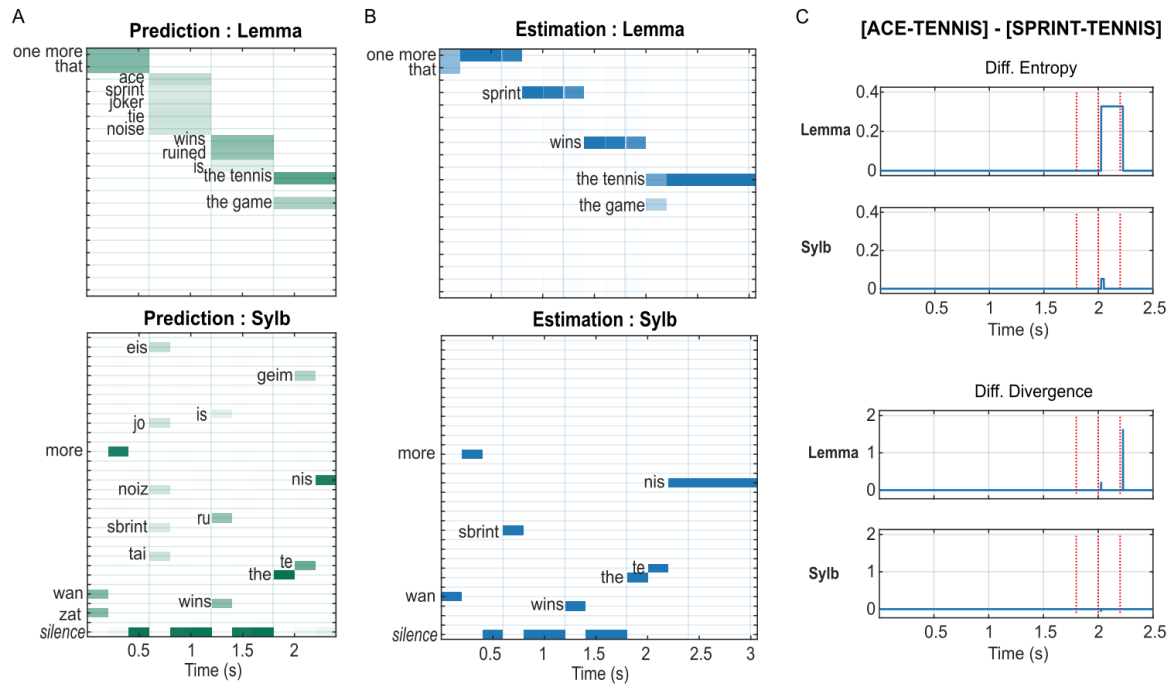
1255 **Supporting Figures**



S1 Fig. Effect of contextual bias ratio on the inference process. A-C: metrics derived from the sentence “One more ace wins the tennis” as function of contextual bias between “poker game” and “tennis game”. A bias of x implies that the prior probability ratio (the total probability is always normalized to 1) for context was set to $[x \ 1 \ 1 \ 1]$ for all 4 possible contexts {‘poker game’, ‘tennis game’, ‘night party’, ‘racing game’} for $x \geq 1$, and $[1 \ 1/x \ 1 \ 1]$ for $x < 1$ to balance the influence of the two irrelevant contexts. D-F: same metrics derived from sentence “One more ace wins the game”. A. Inferred states for the *context* (blue) and the *agent* (red) do not change with contextual bias, i.e. the model always resolved to the correct states. B. Sum of entropy across *context*, *agent* and *patient* at the subject word (“ace”) offset and the sentence offset. At the offset of “ace” (blue), the entropy is maximum at bias=1 and symmetric on both sides. At sentence offset (red), the entropy is overall lower than at the offset of “ace” and monotonically increases with a small slope, reflecting that the model was more certain about the state estimations at this point, but keeps a small possibility towards the poker game that increases with the bias towards the poker context. C. At the sentence offset, the divergence monotonically increases with bias towards poker reflecting the increasing difference between the expected context (poker) and the actual one (tennis). D. Inferred states for context and agent at the end of sentence B as a function of bias. For bias < 1 (preference for ‘tennis’ context), the inferred context is “tennis (game)” and inferred agent is “serve”. For bias ≥ 1 , the result corresponds to a preference for the “poker” context. E. Sum of entropy. For both time points, the entropy is at maximum when bias=1. Both curves are symmetrical by bias=1. The blue curve is the same as in B because the sentence input up to this point was the same. F. Sum of divergence across the same three factors at two critical time points. At the offset of “ace”, the divergence reached its minimum at bias=1 as a result of the uniform distribution over “poker” and “tennis” states, which is the least different from the previous time point. At the sentence offset, the stronger the bias (farther from 1), the smaller the difference between before and after hearing the final word. However, a notch is seen at bias=1 due to the uncertainty (S1 Fig E).

1256

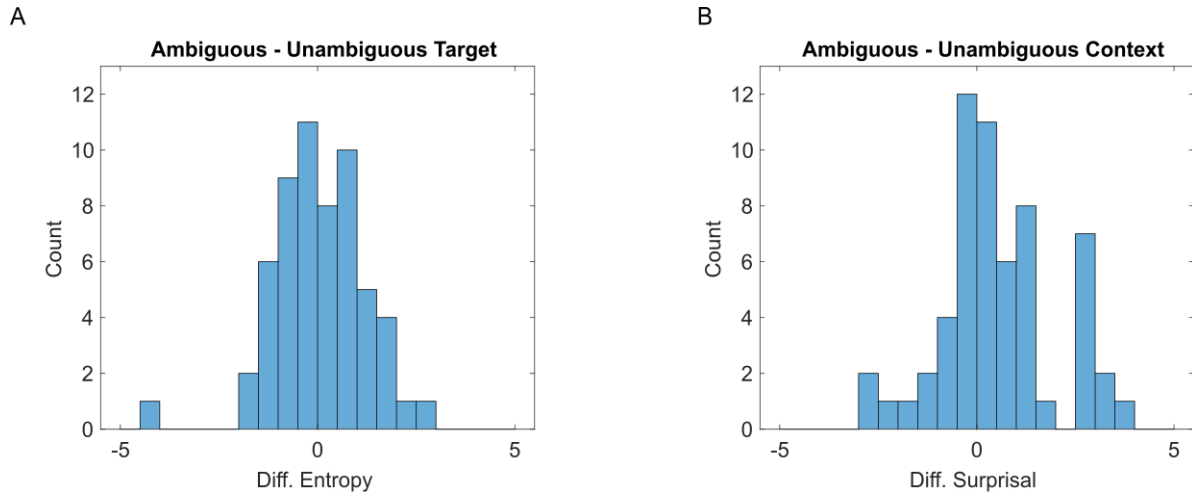
1257



S2 Fig. Message passing in the processing of the same word in different sentences. Figure specifications are the same as Fig 3. **A. Semantic-to-lemma and lemma-to-syllable predictions in response to sentence “one more sprint wins the tennis”.** The second lemma “sprint” influences the prediction for the final lemma as well as the corresponding syllables as compared to Fig 3A. **B. Estimation of posterior probabilities for lemma and syllable states for the sentence [SRPINT-tennis].** Similar to Fig 3B, the model instantly recognizes each syllable (lower panel). **C. Upper panels: entropy derived from sentence [ACE-TENNIS] minus sentence [SPRINT-TENNIS] for the lemma and the syllable levels for the entire sentence.** Vertical dotted lines mark the onset of each syllable of the final lemma. Entropies for both the lemma and the syllable level was higher for [ACE-TENNIS] after the onset of the second syllable, reflecting a greater complexity (three possible states compared to two in the sentence [SPRINT-TENNIS]) of the prediction of the final lemma. **Lower panels: the difference between the divergence in response to the two sentences.** A positive difference at the onset of the third syllable (the offset of the second syllable) indicates that the input “the tennis” is less expected in the sentence [ACE-TENNIS] due to the prior preference for the poker context, compared to in the sentence [SPRINT-TENNIS] where the context was already resolved to “poker game” after hearing “sprint”.

1258

1259

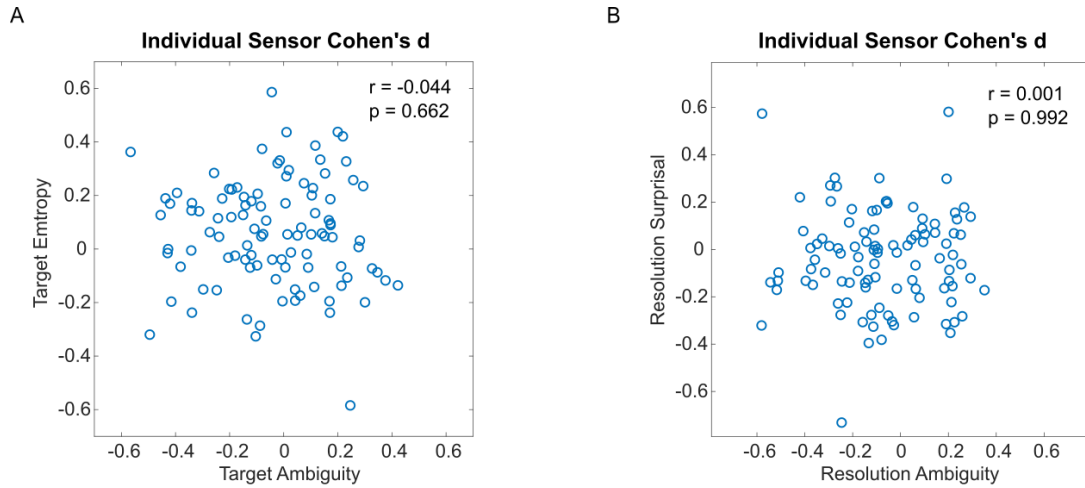


S3 Fig. A. Distribution for the difference of GPT-2 prediction entropy calculated from ambiguous vs. unambiguous Target words. Only the 58 selected sentences were included. **B.** Distribution for the difference of GPT-2 prediction surprisal calculated from the same Resolution words following ambiguous vs. unambiguous Target.

1260

1261

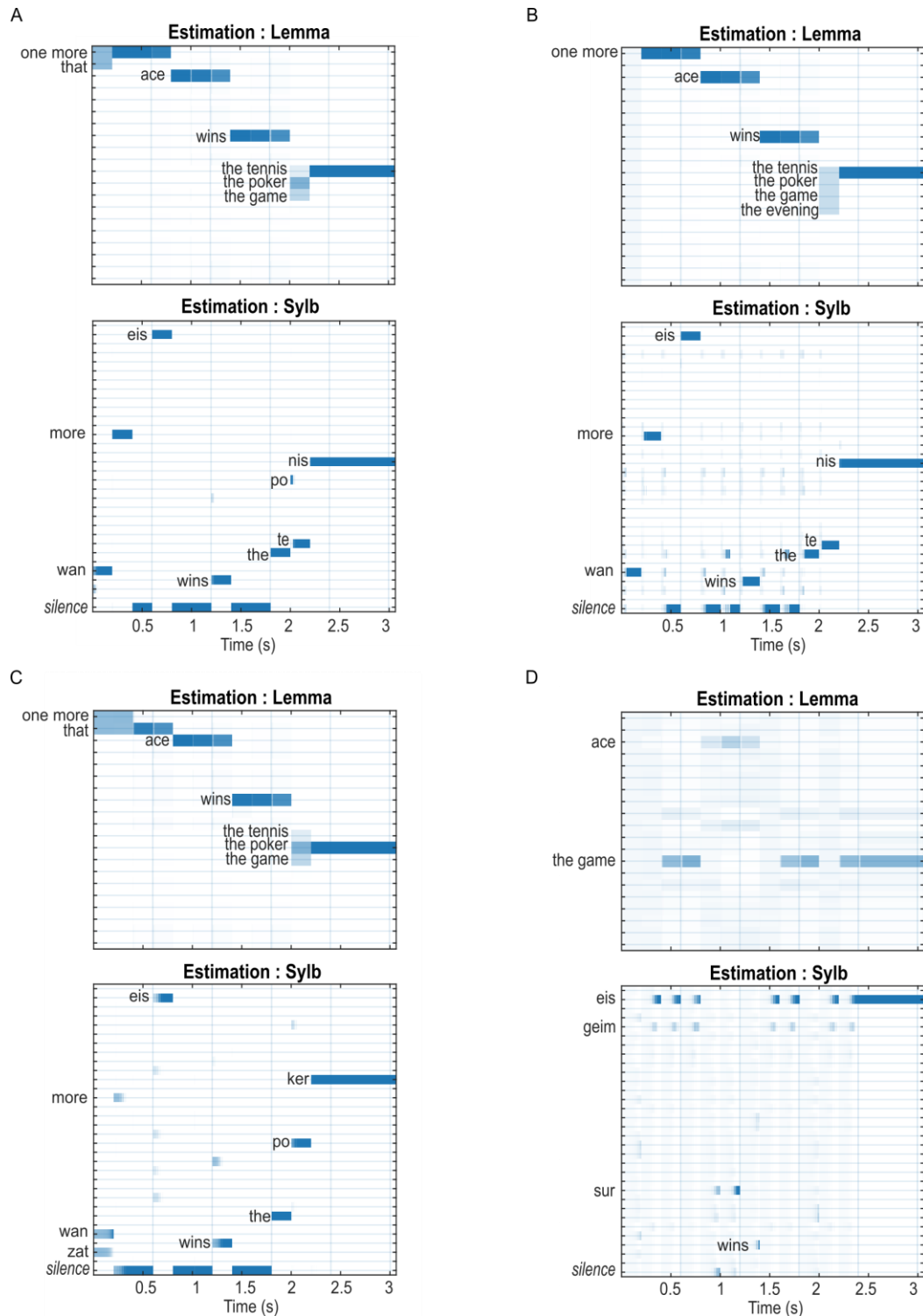
1262



S4 Fig. Comparison of effect sizes between semantic and GPT-2 prediction metrics. A. Cohen's d computed from the effect of semantic ambiguity (x-axis) and the effect of GPT-2 prediction entropy (y-axis) at Target offset for each of the 102 combined gradiometers. B. Cohen's d for the effect of preceding ambiguity (x-axis) vs. GPT-2 prediction surprisal (y-axis) at Resolution offset for each combined gradiometer.

1263

1264



S5 Fig. A, B: Inference of lemma and syntax states at moderately high precision ($\exp(6)$) with (A) or without (B) informative top-down predictions. The posterior estimates are very similar to the intact condition (Fig 4B and 5A, respectively) in that the model quickly converged onto the correct states after each update. However, longer delays to convergence can be observed at the syllable level with prediction, and both lemma and syllable levels without prediction, compared to their intact counterparts. **C, D: Inference of lemma and syntax states at extremely low precision ($\exp(-4)$) with (C) or without (D) informative top-down predictions.** The posterior estimates with informative prediction are qualitatively the same as the low-precision condition in Fig 6A but with longer delays before convergence. Without any top-down prediction, the model completely fails at the syllable level, hence cannot make accurate estimates for higher levels.