

To understand connected speech, listeners must construct a hierarchy of linguistic structures of different sizes, including syllables, words, phrases and sentences<sup>1-3</sup>. It remains puzzling how the brain simultaneously handles the distinct timescales of the different linguistic structures, for example, from a few hundred milliseconds for syllables to a few seconds for sentences<sup>4-14</sup>. Previous studies have suggested that cortical activity is synchronized to acoustic features of speech, approximately at the syllabic rate, providing an initial

corrected) and the response was highly consistent across listeners (Fig. 1c). Given that the phrasal- and sentential-rate rhythms were not conveyed by acoustic fluctuations at the corresponding frequencies (Fig. 1b), cortical responses at the phrasal and sentential rates must be a consequence of internal online structure building processes. Cortical activity at all the three peak frequencies was seen bilaterally (Fig. 1c). The response power averaged over sensors in each hemisphere was significantly stronger in the left hemisphere at the sentential rate ( $p = 0.014$ , paired two-sided  $t$  test), but not at the phrasal ( $p = 0.20$ , paired two-sided  $t$  test) or syllabic rates ( $p = 0.40$ , paired two-sided  $t$  test).

Are the responses at the phrasal and sentential rates indeed separate neural indices of processing at distinct linguistic levels or are they merely sub-harmonics of the syllabic rate response, generated by intrinsic cortical dynamical properties? We address this question by manipulating different levels of linguistic structure in the input. When the stimulus is a sequence of random syllables that preserves the acoustic properties of Chinese sentences (Fig. 1 and Supplementary Fig. 2), but eliminates the phrasal/sentential structure, only syllabic (acoustic) level tracking occurs ( $p = 1.1 \times 10^{-4}$  at 4 Hz, paired one-sided  $t$  test, FDR corrected; Fig. 2a). Furthermore, this manipulation preserves the position of each syllable in a sentence (Online Methods) and therefore further demonstrates that the phrasal- and sentential-rate responses are not a result of possible acoustic differences between the syllables in a sentence. When two adjacent syllables and morphemes combine into verb phrases, but there is no four-element sentential structure, phrasal-level tracking emerges at half of the syllabic rate ( $p = 8.6 \times 10^{-4}$  at 2 Hz and  $p = 2.7 \times 10^{-4}$  at 4 Hz, paired one-sided  $t$  test, FDR corrected; Fig. 2b). Similar responses are observed for noun phrases (Supplementary Fig. 3).

To test whether the phrase-level responses segregate from the sentence level, we constructed longer verb phrases that were unevenly divided into a monosyllabic verb followed by a three-syllable noun phrase (Fig. 2c). We expect that the neural responses to the long verb phrase to be tagged at 1 Hz, whereas the neural responses to the monosyllabic verb and the three-syllable noun phrase will present as harmonics of 1 Hz. Consistent with our hypothesis, cortical dynamics emerged at one-fourth of the syllabic rate, whereas the response at half of the syllabic rate is no longer detectable ( $p = 1.9 \times 10^{-4}$ ,  $1.7 \times 10^{-4}$  and  $9.3 \times 10^{-4}$  at 1, 3 and 4 Hz, respectively, paired one-sided  $t$  test, FDR corrected).

When listening to Chinese sentences (Fig. 1a), listeners who did not understand Chinese only showed responses to the syllabic (acoustic) rhythm ( $p = 3.0 \times 10^{-5}$  at 4 Hz, paired one-sided  $t$  test, FDR corrected; Fig. 2d), further supporting the argument that cortical responses to larger, abstract linguistic structures is a direct consequence of language comprehension.

If aligning cortical dynamics to the time course of linguistic constituent structure is a general mechanism required for comprehension, it must apply across languages. Indeed, when native English speakers were tested with English materials (Fig. 1a), their cortical activity also followed the time course of larger linguistic structures, that is, phrases and sentences ( $p = 4.1 \times 10^{-5}$ , syllabic rate; Fig. 2e;  $p = 3.9 \times 10^{-3}$

$4.3 \times 10^{-3}$  and  $6.8 \times 10^{-6}$  at the sentential, phrasal and syllabic rates, respectively; Fig. 2f; paired one-sided test, FDR corrected).

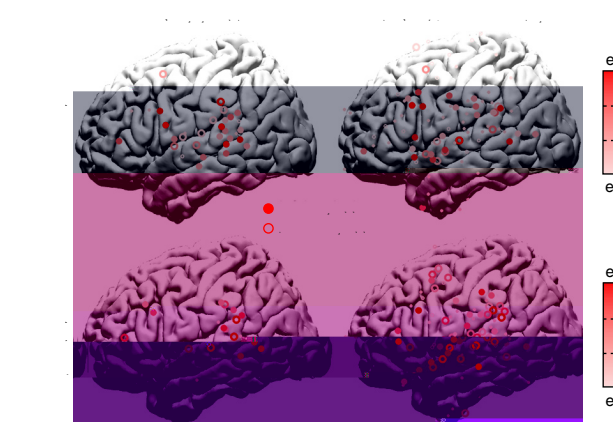
We found that concurrent neural tracking of multiple levels of linguistic structure was not confounded with the encoding of acoustic cues (Figs. 1 and 2). However, is this simply explained by the neural tracking of the predictability of smaller units? As a larger linguistic structure, such as a sentence, unfolds in time, its component units become more predictable. Thus, cortical networks solely tracking transitional probabilities across smaller units could show temporal dynamics matching the timescale of larger structures. To test this alternative hypothesis, we crafted a constant transitional probability Markovian Sentence Set (MSS) in which the transitional probability of lower level units was dissociated from the higher level structures (Fig. 3a and Supplementary Fig. 1e,f). The constant transitional probability MSS is contrasted with a varying transitional probability MSS, in which the transitional probability is low across sentential boundaries and high in a sentence (Fig. 3b,c). If cortical activity only encodes the transitional probability between lower level units (for example, acoustic chunks in the MSS) independent of the underlying syntactic structure, it can show tracking of the sentential structure for the varying probability MSS, but not for the constant probability MSS. In contrast with this prediction, indistinguishable neural responses to sentences were observed for both MSS (Fig. 3d), demonstrating that neural tracking of sentences is not confounded by transitional probability. Specifically, for the constant transitional probability MSS, the response was statistically significant at the sentential rate, twice the sentential rate and the syllable rate ( $p = 1.8 \times 10^{-4}$ ,  $2.3 \times 10^{-4}$  and

$2.7 \times 10^{-6}$ , respectively). For the varying transitional probability MSS, the response was statistically significant at the sentential rate and the syllable rate ( $p = 1.8 \times 10^{-4}$ ,  $2.3 \times 10^{-4}$  and

Localizing cortical sources of the sentential and phrasal rate responses using ECoG ( $N = 5$ ). Left, power envelope of high-gamma activity. Right, waveform of low-frequency activity. Electrodes in the right hemisphere were projected to the left hemisphere, and right hemisphere (left hemisphere) electrodes are shown by hollow (filled) circles. The figure only displays electrodes that showed statistically significant neural responses to sentences in Fig. 2 and no significant response to the acoustic control shown in Fig. 2. Significance was determined by bootstrap (FDR corrected) and the significance level is 0.05. The response strength, that is, the response at the target frequency relative to the mean response averaged over a 1-Hz wide neighboring region, is color coded. Electrodes with response strength less than 10 dB are shown by smaller symbols. The sentential and phrasal rate responses were seen in bilateral pSTG, TPJ and left IFG.

tracking of larger linguistic structures generalizes to sentences that are variable in duration (4–8 syllables) and syntactic structures. These sentences were again built on isochronous Chinese syllables, intermixed and sequentially presented without any acoustic gap at the sentence boundaries. Examples translated into English include “Don’t be nervous,” “The book is hard to read,” and “Over the street is a museum.”

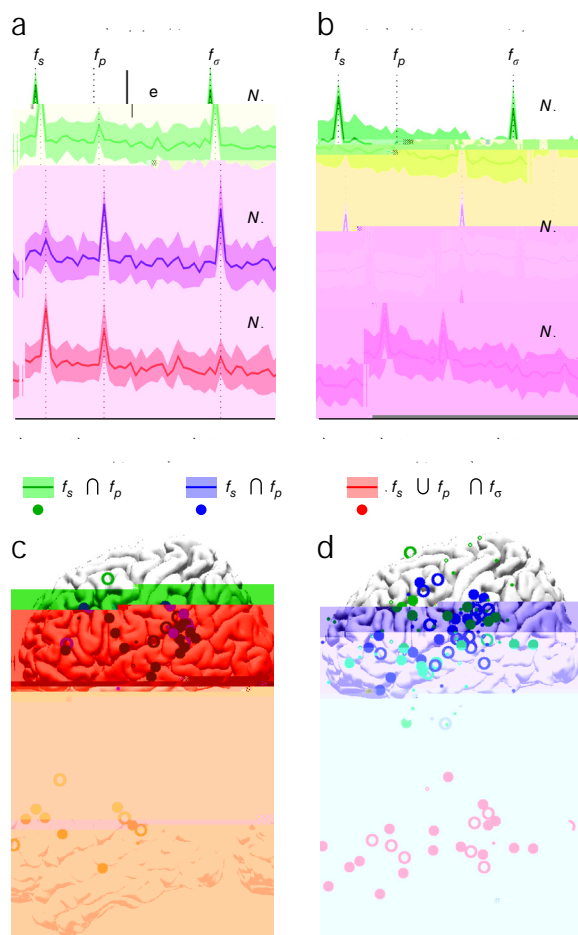
As these sentences have irregular durations that are not tagged by frequency, the MEG responses were analyzed in the time domain by averaging sentences of the same duration. To focus on sentential level processing, we low-pass filtered the response at 3.5 Hz. The MEG response (root mean square, r.m.s., over all sensors) rapidly increased after a sentence boundary and continuously changed throughout the duration of a sentence (Fig. 4a). To illustrate the detailed temporal



dynamics, we averaged the r.m.s. response over all sentences containing six or more syllables after aligning them to the sentence offset (Fig. 4b). During the last four syllables of a sentence, the r.m.s. response continuously and significantly decreased for every syllable, indicating that the neural response continuously changes during the course of a sentence rather than being a transient response only occurring at the sentence boundary.

A single-trial decoding analysis was performed to independently confirm that cortical activity tracks the duration of sentences (Fig. 4c). The decoder applied template matching for the response time course (leave-one-out cross-validation, Online Methods) and correctly determined the duration of  $34.9 \pm 0.6\%$  sentences (mean  $\pm$  s.e.m. over subjects, significantly above chance,  $p = 1.3 \times 10^{-7}$ , one-sided test).

After demonstrating cortical tracking of sentences, we further tested whether cortical activity also tracks the phrasal structure inside of a sentence. We constructed sentences that consist of a noun phrase followed by a verb phrase and manipulated the duration of the noun phrase (three syllable or four syllable). The cortical responses closely follow the duration of the noun phrase: the r.m.s. response gradually decreased in the noun phrase, then showed a transient increase after the onset of the verb phrase (Fig. 4d).



We found that large-scale neural activity measured by MEG concurrently follows the hierarchical linguistic structure of speech, but which neural networks generate such activity? To address this question, we recorded the ECoG responses to English sentences (Fig. 2e) and an acoustic control (Fig. 2f). ECoG signals are mesoscopic neurophysiological signals recorded by intracranial electrodes implanted in

Spatial dissociation between sentential-rate, phrasal-rate and syllabic-rate responses ( $N = 5$ ). ( ) The power spectrum of the power envelope of high-gamma activity. ( ) The power spectrum of low-frequency ECoG waveform. The top panels (green curves) show the response averaged over all electrodes that show a significant sentential-rate response but not a significant phrasal-rate response. Significance was determined by bootstrap (FDR corrected) and the significance level is 0.05. The shaded area is 1 s.d. over electrodes on each side. The blue curves show the response averaged over all electrodes that showed a significant phrasal-rate response, but not a significant sentential-rate response. The red curves show a significant sentential-rate or a significant phrasal-rate response, but not a significant syllabic response. ( , ) The topographic distribution of the three groups of electrodes analyzed in and . As in , electrodes showing a response greater than 10 dB are shown by larger symbols than electrodes showing a response weaker than 10 dB.

epilepsy patients for clinical evaluation (see **Supplementary Fig. 5** for the electrode coverage), and they possess better spatial resolution than MEG. We first analyzed the power of the ECoG signal in the high gamma band (70–200 Hz), as it highly correlates with multiunit firing<sup>23</sup>. The electrodes exhibiting significant sentential, phrasal and syllabic rate fluctuations in high gamma power are shown separately (**Fig. 5**). The sentential rate response clustered over the posterior and middle superior temporal gyrus (pSTG), bilaterally, with a second cluster over the left inferior frontal gyrus (IFG). Phrasal rate responses were also observed over the pSTG bilaterally. Notably, although the sentential and phrasal rate responses were observed in similar cortical areas, electrodes showing phrasal rate responses only partially overlapped with electrodes showing sentential rate responses in the pSTG (**Fig. 6**). For electrodes showing a significant response at either the sentential rate or the phrasal rate, the strength of the sentential rate response was negatively correlated with the strength of the phrasal rate response ( $r = -0.32$ ,  $P = 0.004$ , bootstrap). This phenomenon demonstrates spatially dissociable neural tracking of the sentential and phrasal structures.

Furthermore, some electrodes with a significant sentential or phrasal rate response showed no significant syllabic rate response ( $P < 0.05$ , FDR corrected, **Fig. 6**). In other words, there are cortical circuits specifically encoding larger, abstract linguistic structures without responding to syllabic-level acoustic features of speech. In addition, although the syllabic responses were not significantly different ( $P > 0.05$ , FDR corrected) for English sentences and the acoustic control in the MEG results, they were dissociable spatially in the ECoG results (**Fig. 7**). Electrodes showing significant syllabic responses ( $P < 0.05$ , FDR corrected) to sentences, but not the acoustic control, were seen in bilateral pSTG, bilateral anterior STG (aSTG), and left IFG.

We then analyzed neural tracking of the sentential, phrasal and syllabic rhythms in the low-frequency ECoG waveform (**Fig. 5**), which is a close neural correlate of MEG activity. Fourier analysis was directly applied to the ECoG waveform and the Fourier coefficients at 1, 2 and 4 Hz are extracted. Low-frequency ECoG activity is usually viewed as the dendritic input to a cortical area<sup>24</sup>. The low-frequency responses are more distributed than high-gamma activity, possibly reflecting the fact that the neural representations of different levels of linguistic structures serve as inputs to broad cortical areas. Sentential and phrasal rate responses are strong in STG, IFG and temporoparietal junction (TPJ). Compared with the acoustic control, the syllabic-rate response to sentences was stronger in broad cortical areas, including the temporal and frontal lobes (**Fig. 7**). Similar to the high-gamma activity, the low-frequency responses to the sentential and phrasal structures were not reflected in the same set of electrodes (**Fig. 6**).

For electrodes showing a significant response at either the sentential rate or the phrasal rate, the strength of the sentential rate response was also negatively correlated with the strength of the phrasal rate response ( $r = -0.21$ , significantly greater than 0,  $P = 0.023$ , bootstrap).

Our data show that the multiple timescales that are required for the processing of linguistic structures of different sizes emerge in cortical networks during speech comprehension. The neural sources for sentential, phrasal and syllabic rate responses are highly distributed and include cortical areas that have been found to be critical for prosodic (for example, right STG), syntactic and semantic (for example, left pSTG and left IFG) processing<sup>9,25–28</sup>. Neural integration on different timescales is likely to underlie the transformation from shorter lived neural representations of smaller linguistic units to longer lasting neural representations of larger linguistic structures<sup>11–14</sup>. These results underscore the undeniable existence of hierarchical structure building operations in language comprehension<sup>1,2</sup> and can be applied to objectively assess language processing in children and difficult-to-test populations, as well as animal preparations to allow for cross-species comparisons.

Concurrent neural tracking of hierarchical linguistic structures provides a plausible functional mechanism for temporally integrating smaller linguistic units into larger structures. In this form of concurrent neural tracking, the neural representation of smaller linguistic units is embedded at different phases of the neural activity tracking a higher level structure. Thus, it provides a possible mechanism to transform the hierarchical embedding of linguistic structures into hierarchical embedding of neural dynamics, which may facilitate information integration in time<sup>10,11</sup>. Low-frequency neural tracking of linguistic structures may further modulate higher frequency neural oscillations<sup>29–31</sup>, which have been proposed to provide additional roles for structure building<sup>7</sup>. In addition, multiple resources and computations are needed for syntactic analysis, for example, access to combinatorial syntactic subroutines, and such operations may correspond to neural computations on distinct frequency scales, which are coordinated by the low-frequency neural tracking of linguistic constituent structures. Furthermore, low-frequency neural activity and oscillations have been hypothesized as critical mechanisms to generate predictions about future events<sup>32</sup>. For language processing, it is likely that concurrent neural

Recent work has shown that cortex tracks the slow acoustic fluctuations of speech below 10 Hz (refs. 15–18,34,35), and this phenomenon is commonly described as ‘cortical entrainment’ to the syllabic rhythm of speech. It has been controversial whether such syllabic-level cortical entrainment is related to low-level auditory encoding or language processing<sup>6</sup>. Our findings demonstrate that processing goes well beyond stimulus-bound analysis: cortical activity is entrained to larger linguistic structures that are, by necessity, internally constructed, based on syntax. The emergence of slow cortical dynamics provides timescales suitable for the analysis of larger chunk sizes<sup>13,14</sup>.

A long-lasting controversy concerns how the neural responses to sensory stimuli are related to intrinsic, ongoing neural oscillations. This question is heavily debated for the neural response entrained to the syllabic rhythm of speech<sup>36</sup> and can also be asked for neural activity entrained to the time courses of larger linguistic structures. Our experiment was not designed to answer this question; however, we clearly found that cortical speech processing networks have the

9. Pallier, C., Devauchelle, A.-D. & Dehaene, S. Cortical representation of the constituent structure of sentences. *Proc. Natl. Acad. Sci. USA* **10**, 2522–2527 (2011).
10. Schroeder, C.E., Lakatos, P., Kajikawa, Y., Partan, S. & Puce, A. Neuronal oscillations and visual amplification of speech. *Trends Cogn. Sci.* **12**, 106–113 (2008).
11. Buzsáki, G. Neural syntax: cell assemblies, synapsemblies and readers. *Neuron* **67**, 362–385 (2010).
12. Bernacchia, A., Seo, H., Lee, D. & Wang, X.-J. A reservoir of time constants for memory traces in cortical neurons. *Nat. Neurosci.* **14**, 366–372 (2011).
13. Lerner, Y., Honey, C.J., Silbert, L.J. & Hasson, U. Topographic mapping of a hierarchy of temporal receptive windows using a narrated story. *J. Neurosci.* **31**, 2906–2915 (2011).
14. Kiebel, S.J., Daunizeau, J. & Friston, K.J. A hierarchy of time-scales and the brain. *PLoS Comput. Biol.* **4**, e1000209 (2008).
15. Luo, H. & Poeppel, D. Phase patterns of neuronal responses reliably discriminate speech in human auditory cortex. *Neuron* **59**, 1001–1010 (2007).
16. Ding, N. & Simon, J.Z. Emergence of neural encoding of auditory objects while listening to competing speakers. *Proc. Natl. Acad. Sci. USA* **109**, 11854–11859 (2012).
17. Zion Golumbic, E.M. *et al.* Mechanisms underlying selective neuronal tracking of attended speech at a “cocktail party”. *Neuron* **77**, 980–991 (2013).
18. Peelle, J.E., Gross, J. & Davis, M.H. Phase-locked responses to speech in human auditory cortex are enhanced during comprehension. *Cereb. Cortex* **23**, 1378–1387 (2013).
19. Pasley, B.N. *et al.* Reconstructing speech from human auditory cortex. *PLoS Biol.* **10**, e1001251 (2012).
20. Steinhauer, K., Alter, K. & Friederici, A.D. Brain potentials indicate immediate use of prosodic cues in natural speech processing. *Nat. Neurosci.* **2**, 191–196 (1999).
21. Peña, M., Bonatti, L.L., Nespor, M. & Mehler, J. Signal-driven computations in speech processing. *Science* **297**, 604–607 (2002).
22. Saffran, J.R., Aslin, R.N. & Newport, E.L. Statistical learning by 8-month-old infants. *Science* **270**, 1926–1928 (1996).

232

**Participants.** 34 native listeners of Mandarin Chinese (19–36 years old, mean 25 years old; 13 male) and 13 native listeners of American English (22–46 years old, mean 26 years old; 6 male) participated in the study. All Chinese listeners received high school education in China and 26 of them also received college education in China. None of the English listeners understood Chinese. All participants were right-handed<sup>51</sup>. Five experiments were run for Chinese listeners and two experiments for English listeners. Each experiment included eight listeners (except that the AMS experiment involved five listeners) and each listener participated in at most two experiments. The number of listeners per experiment was chosen based on previous MEG experiments on neural tracking of continuous speech. The sample size for previous experiments was typically between three and 12 (refs. 15,16), and the basic phenomenon reported here was replicated in all the seven experiments of the study ( $n = 47$  in total). The experimental procedures were approved by the New York University Institutional Review Board, and written informed consent was obtained from each participant before the experiment.

**Stimuli I: Chinese materials.** All Chinese materials were constructed based on an isochronous sequence of syllables. Even when the syllables were hierarchically grouped into linguistic constituents, no acoustic gaps were inserted between constituents. All syllables were synthesized independently using the Neospeech synthesizer (<http://www.neospeech.com/>, the male voice, Liang). The synthesized syllables were 75–354 ms in duration (mean duration 224 ms), and were adjusted to 250 ms by truncation or padding silence at the end. The last 25 ms of each syllable were smoothed by a cosine window.

**F - .** 50 four-syllable sentences were constructed, in which the first two syllables formed a noun phrase and the last two syllables formed a verb phrase (Supplementary Table 1). The noun phrase could be composed of either a single two-syllable noun or a one-syllable adjective followed by a one-syllable noun. The verb phrase could be composed of either a two-syllable verb or a one-syllable verb followed by a one-syllable noun object. In a normal trial, ten sentences were sequentially played and no acoustic gaps were inserted between sentences (Supplementary Fig. 1a). Due to the lack of phrasal and sentential level prosodic cues, the sound intensity of the stimulus, characterized by the sound envelope, only fluctuates at the syllabic rate but not at the phrasal or sentential rate (Supplementary Fig. 2). An outlier trial was the same as a normal trial except that the verb phrases in two sentences were exchanged, creating two nonsense sentences with incompatible subjects and predicates (an example in English would be “new plans rub skin”).

**F - .** Two types of four-syllable verb phrases were created. Type I verb phrase contained a one-syllable verb followed by a three-syllable noun phrase, which could be a compound noun or an adjective + noun phrase (Supplementary Fig. 1b and Supplementary Table 1). Type II verb phrase contained a two-syllable verb followed by a two-syllable noun (Supplementary Fig. 1c, all phrases listed in Supplementary Table 1). 50 instances were created for each type of verb phrases. In a normal trial, ten phrases of the same type were sequentially presented. An outlier trial was the same as a normal trial except that the verbs in two phrases were exchanged, creating two nonsense verb phrases with incompatible verbs and objects (an example in English would be “drink a long walk”).

**- .** The verb phrases (or the noun phrases) in the four-syllable sentences were presented in a sequence (Supplementary Fig. 1d). In a normal trial, 20 different phrases were played. In an outlier trial, one of the 20 phrases was replaced by two random syllables that did not constitute a sensible phrase.

**- .** The random syllabic sequences were generated based on the four-syllable sentences. Each four-syllable sentence was transformed into four random syllables using the following rule: the first syllable in the sentence was replaced by the first syllable of a randomly chosen sentence. The second syllable was replaced by the second syllable of another randomly chosen sentence and the same for the third and the fourth syllables. This way, if there were any consistent acoustic differences between the syllables at different positions in a sentence, those acoustic differences were preserved in the random syllabic sequences. Each normal trial contained 40 syllables. In outlier trials, four consecutive syllables were replaced by a Chinese idiom.

**B - .** In normal trials, ten four-syllable sentences were played but with all syllables being played backward in time. An outlier trial was the same as a normal trial except that four consecutive syllables at a random position were replaced by four random syllables that were not reversed in time.

**F - .** 50 common four-syllable idioms were selected (Supplementary Table 1), in which the first two syllables formed a noun phrase and the last two syllables formed a verb phrase. In a normal trial ten sentences were played. An outlier trial was the same as a normal trial except that the noun phrases in two idioms were exchanged, creating two nonexistent and semantically nonsensical idioms.

**- .** The sentence duration was varied between four and eight syllables. 40 sentences were constructed for each duration, resulting in a total of 200 sentences (listed in Supplementary Table 1). All 200 sentences were intermixed. In a normal trial, ten different sentences were sequentially played without inserting any acoustic gap in between sentences. In an outlier trial, one of the ten sentences was replaced by a syntactically correct but semantically anomalous sentence. Examples of nonsense sentences, translated into English, included “ancient history is drinking tea” and “take part in his portable hard drive”.

**- .** All sentences consisted of a noun phrase followed by a verb phrase (Supplementary Table 1). The noun phrase had three syllables for half of the sentences ( $n = 45$ ) and four syllables for the other half. A three-syllable noun phrase was followed by either a four-syllable verb phrase ( $n = 20$ ) or a five-syllable verb phrase ( $n = 25$ ). A four-syllable noun phrase was followed by a three-syllable verb phrase ( $n = 20$ ) or a four-syllable verb phrase ( $n = 25$ ). Sentences with different noun phrase durations and verb phrase durations were intermixed. In a normal trial 10 different sentences were played sequentially, without inserting any acoustic gap between phrases or sentences. In an outlier trial one sentence was replaced by a sentence with the same syntactic structure but that was semantically anomalous.

**A - .** Five sets of AMS were created. Each sentence consisted of three components, C1, C2 and C3. Each component (C1, C2 or C3) was independently chosen from three candidate syllables with equal probability. The grammar of the AMS is illustrated in Supplementary Figure 4a. In the experiments, sentences were played sequentially without any gap between sentences. Since all components were chosen independently and each component was chosen from three syllables with equal probability, all components were equally predictable regardless of its position in a sequence. In other words,  $P(C1) = P(C2) = P(C3) = P(C2|C1) = P(C3|C2) = P(C1|C3) = 1/3$ .

All Chinese syllables were synthesized independently and adjusted to 300 ms by truncation or padding silence at the end. In each trial, 60 sentences were played and no additional gap was inserted between sentences. Therefore, the syllables were played at a constant rate of 3.33 Hz and the sentences were played at a constant rate of 1.11 Hz. To make sure that neural encoding of the AMS was not confounded by acoustic properties of a particular set of syllables, five sets of AMS were created (Supplementary Table 1). No meaningful Chinese expressions are embedded in the AMS sequences.

**Stimuli II: English materials.** All English materials were synthesized using the MacinTalk Synthesizer (male voice Alex, in Mac OS X 10.7.5).

**F - E - .** 60 four-syllable English sentences were constructed (Supplementary Table 1), and each syllable was a monosyllabic word. All sentences had the same syntactic structure: adjective/pronoun + noun + verb + noun. Each syllable was synthesized independently, and all the synthesized syllables (250–347 ms in duration) were adjusted to 320 ms by padding silence at the end or truncation. The offset of each syllable was smoothed by a 25-ms cosine window. In each trial, 12 sentences were presented without any acoustic gap between them. In an outlier trial, 3 consecutive words from a random position were replaced by three random words so that the corresponding sentence(s) became ungrammatical.

**- .** Shuffled sequences were constructed as an unintelligible sound sequence that preserved the acoustic properties of the sentence sequences. All syllables in the four-syllable English sentences were segmented into five overlapping slices. Each slice was 72 ms in duration and overlapped with neighboring slices for 10 ms. The first 10 ms and the last 10 ms of each slice was smoothed by a linear ramp, except for the onset of the first slice and the offset of the last slice.





and neural tracking of those linguistic structures was analyzed in the frequency domain. For each trial, to avoid the transient response to the acoustic onset of each trial, the neural recordings were analyzed in a time window between the onset of the second sentence (or the fifth syllable if the stimulus contained no sentential structure) and the end of the trial. The single-trial responses were transformed into the frequency domain using the discrete Fourier transform (DFT). For all Chinese materials and the artificial Markovian language materials, nine sentences were analyzed in each trial, resulting in a frequency resolution of 1/9 of the sentential rate (~0.11 Hz). For the English sentences and the shuffled sequences, the trials were longer and the duration equivalent to 44 English syllables was analyzed, resulting in a frequency resolution of 1/44 of the syllabic rate, that is, 0.071 Hz.

The response topography (Fig. 1c) showed the power of the DFT coefficients at a given frequency and hemispheric lateralization was calculated by averaging the response power over the sensors in each hemisphere ( $n = 54$ ).

Given that the properties of the neural responses to linguistic structures and background neural activity might vary in different frequency bands, to treat each frequency band equally, a separate spatial filter was designed for every frequency bin in the DFT output to optimally estimate the response strength. The linear spatial filter was the DSS filter<sup>56</sup>. The output of the DSS filter was a weighted summation over all MEG sensors, and the weights were optimized to extract neural activity consistent over trials. In brief, if the DFT of the MEG response averaged over trials is  $X(\omega)$  and the autocorrelation matrix of single-trial MEG recordings is  $R(\omega)$ , the spatial filter is  $w = R^{-1}(\omega)X(\omega)$  (see the appendix of ref. 56). The spatial filter  $w$  is an  $157 \times 1$  vector (for the 157 sensors), the same size as  $X(\omega)$ , and  $R(\omega)$  is a  $157 \times 157$  matrix. The spatial filter could be viewed as a virtual sensor that was optimized to record phase-locked neural activity at each frequency.

Power of the scalar output of the spatial filter was calculated as  $|w^T X(\omega)|^2$ . The power spectrum was calculated by averaging the power over trials. The power spectrum was normalized by the power spectrum of the shuffled sequences. The normalized power spectrum was used to calculate the lateralization index (LI) as  $LI = (P_L - P_R) / (P_L + P_R)$ , where  $P_L$  and  $P_R$  are the power spectra of the left and right hemispheres, respectively. The LI ranges from -1 to 1, with -1 indicating left lateralization and 1 indicating right lateralization. The LI was used to calculate the lateralization index (LI) as  $LI = (P_L - P_R) / (P_L + P_R)$ , where  $P_L$  and  $P_R$  are the power spectra of the left and right hemispheres, respectively. The LI ranges from -1 to 1, with -1 indicating left lateralization and 1 indicating right lateralization.