# Task-dependent and automatic tracking of hierarchical linguistic structure

Sanne Ten Oever[1,2,3], Sara Carta[1,4,5], Greta Kaufeld[1], Andrea E. Martin[1,2*]

[1] Language and Computation in Neural Systems group, Max Planck Institute for Psycholinguistics, Nijmegen, The Netherlands

[2] Language and Computation in Neural Systems group, Donders Centre for Cognitive Neuroimaging, Nijmegen, The Netherlands

[3] Department of Cognitive Neuroscience, Faculty of Psychology and Neuroscience, Maastricht University, The Netherlands

[4] ADAPT Centre, School of Computer Science and Statistics, University of Dublin, Trinity College, Ireland

[5] CIMeC - Center for Mind/Brain Sciences, University of Trento, Italy

[*] Corresponding author: andrea.martin@mpi.nl

13 **Abstract**

14 Linguistic phrases are tracked in sentences even though there is no clear acoustic phrasal marker

15 in the physical signal. This phenomenon suggests an automatic tracking of abstract linguistic

16 structure that is endogenously generated by the brain. However, all studies investigating linguistic

17 tracking compare conditions where either relevant information at linguistic timescales is available,

18 or where this information is absent altogether (e.g., sentences versus word lists during passive

19 listening). It is therefore unclear whether tracking at these phrasal timescales is related to the

20 content of language, or rather, is a consequence of attending to the timescales that happen to match

21 behaviourally-relevant information. To investigate this question, we presented participants with

22 sentences and word lists while recording their brain activity with MEG. Participants performed

23 passive, syllable, word, and word-combination tasks corresponding to attending to rates they

24 would naturally attend to, syllable-rates, word-rates, and phrasal-rates, respectively. We replicated

25 overall findings of stronger phrasal-rate tracking measured with mutual information (MI) for

26 sentences compared to word lists across the classical language network. However, in the inferior

27 frontal gyrus (IFG) we found a task-effect suggesting stronger phrasal-rate tracking during the

28 word-combination task independent of the presence of linguistic structure, as well as stronger

29 delta-band connectivity during this task. These results suggest that extracting linguistic

30 information at phrasal-rates occurs automatically with or without the presence of an additional

31 task, but also that that IFG might be important for temporal integration across various perceptual

32 domains.

33

34 **Keywords:** sentence comprehension; mutual information; entrainment; temporal dynamics

## Introduction

Understanding spoken language requires a multitude of processes [1-3]. Acoustic patterns have to be segmented and mapped onto internally stored phonetic and syllabic representations [3-5]. These phonemes have to be combined and mapped onto words which then have to be mapped to abstract linguistic phrasal structures [2, 6]. Proficient speakers of a language seem to do this so naturally that one might almost forget the complex parallel and hierarchical processing which occurs during natural speech and language comprehension.

It has been shown that it is essential to track the temporal dynamics of the speech signal in order to understand its meaning [7, 8]. In natural speech, syllables follow up on each other in the theta range (3-8 Hz; [9-11]), while higher-level linguistic features such as words and phrases occur at lower rates (0.5-3 Hz; [9, 12, 13]). Tracking of syllabic features is stronger when one understands a language [14-16] and tracking of phrasal rates is more prominent when the signal contains phrasal information ([12, 13, 17]; e.g., word lists versus sentences). Importantly, phrasal tracking even occurs when there are no distinct acoustic modulations at the phrasal rate [12, 13, 17]. These results seem to suggest that tracking of relevant temporal timescales is critical for speech understanding.

An observation one could make regarding these findings is that tracking occurs only at the rates that are meaningful and thereby behaviourally relevant [12, 17]. For example, in word lists, word-rate is the slowest rate that is meaningful during natural listening. Modulations at slower phrasal rates might not be tracked as they do not contain behaviourally relevant information. In contrast, in sentences phrasal rates contain linguistic information and therefore these slower rates are also tracked. Thus, when listening to speech one automatically tries to extract the meaning, which requires extracting information at the highest linguistic level [3, 5]. However, it is unsure if the tracking at these slower rates is a unique feature of language processing or rather dependent on the level of attention to relevant temporal timescales.

As understanding language requires a multitude of processing, it is difficult to figure out what participants actually are doing when listening to natural speech. Moreover, designing a task in an experimental setting that does justice to this multitude of processing is difficult. This is probably why tasks in language studies vary vastly. Tasks include passively listening (e.g. [12], asking comprehension questions (e.g. [13], rating intelligibility (e.g. [14, 16], working memory

3

65  tasks (e.g. [18], or even syllable counting (e.g. [17]. It is unclear whether outcomes are dependent

66  on the specifics of the task. There has so far not been a study that investigates if task instructions

67  focusing on extracting information at different temporal rates or timescales have an influence on

68  the tracking that occurs on these timescales. It is therefore not clear whether tracking phrasal

69  timescales is unique for language stimuli which contain phrasal structures, or could also occur for

70  other acoustic materials where participants are instructed to pay attention to information happening

71  at these temporal rates or timescales.

72      To answer this question, we designed an experiment in which participants were instructed

73  to pay attention to different temporal modulation rates while listening to the same stimuli. We

74  presented participants with naturally spoken sentences and word lists and asked them to either

75  passively listen, or perform a task on the temporal scales corresponding to syllables, words, or

76  phrases. We recorded MEG while participants performed these tasks and investigated tracking as

77  well as power and connectivity at three nodes that are part of the language network: the superior

78  temporal gyrus (STG), the middle temporal gyrus (MTG), and the inferior frontal gyrus (IFG). We

79  hypothesized that if tracking is purely based on behavioural relevance, it should mostly depend on

80  the task instructions, rather than the nature of the stimuli. In contrast, if there is something

81  automatic and specific about language information, tracking should depend on the level of

82  linguistic information in the acoustic signal.

## Methods

*Participants.* In total twenty Dutch native speakers (16 females; age range: 18-59; mean age = 39.5) participated in the study. All were right-handed, reported normal hearing, had normal or corrected-to-normal vision, and did not have any history of dyslexia or other language related disorders. Participants performed a screening for their eligibility in the MEG and MRI and gave written informed consent. The study was approved by the ethical Commission for human research Arnhem/Nijmegen (project number CMO2014/288). Participants were reimbursed for their participation. One participant was excluded from the analysis as they did not finish the full session.

*Materials and design.* Materials were identical to the stimuli used in Kaufeld et al., [12]. They consisted of naturally spoken sentences or word lists which consisted of 10 words (see Table 1 for examples). The sentences contained two coordinate clauses with the following structure: [Adj N V N Conj Det Adj N V N]. All words were disyllabic except for the words "de" (*the*) and "en" (*and*). Word lists were word-scrambled versions of the original sentences which always followed the structure [V V Adj Adj Det Conj N N N N] or [N N N N Det Conj V V Adj Adj] to ensure that they were grammatically incorrect. In total sixty sentences were used. All sentences were presented at a comfortable sound level.

Participants were asked to perform four different tasks on these stimuli: a passive task, a syllable task, a word task, and a word combination task. For the passive task, participants did not need to perform any task other than comprehension – they only needed to press a button to go to the next trial. For the syllable task, participants heard after every sentence two part-of-speech sounds, each consisting of one syllable. The sound fragments were a randomly determined syllable from the previously presented sentence and a random syllable from all other sentences. Participants' task was to indicate via a button press which of the two sound fragments was part of the previous sentence. For the word task, two words were displayed

**Table 1. Stimuli and task examples**

| sentence | [bange helden] [plukken bloemen] en de [bruine vogels] [halen takken] |  |  |  |
|---|---|---|---|---|
| | [*timid heroes*] [*pluck flowers*] *and the* [*brown birds*] [*gather branches*] |  |  |  |
| word list | [helden bloemen] [vogels takken] de en [plukken halen] [bange bruine] |  |  |  |
| | [*heroes flowers*] [*birds branches*] *and the* [*pluck gather*] [*timid brown*] |  |  |  |

| | sentence | | word list | |
|---|---|---|---|---|
| | correct | incorrect | correct | incorrect |
| syllable | /bɑ/ | /lɑ/ | /bɑ/ | /lɑ/ |

| word | bloemen [*flowers*] | vaders [*fathers*] | bloemen [*flowers*] | vaders [*fathers*] |
|---|---|---|---|---|
| word combination | bange helden [*timid heroes*] | halen bloemen [*gather flowers*] | helden bloemen [*heroes flowers*] | vogels bloemen [*birds flowers*] |

**For each condition (sentence and word list) one example stimulus (top) and corresponding tasks are shown (bottom).**

107  on the screen after each trial (a random word from the just presented sentence and one random

108  word from all other sentences excluding "de" and "en"), and participants needed to indicate which

109  of the two words was part of the sentence before. For the word combination task, participants were

110  presented with two word pairs on the screen. Each of the four words was part of the just presented

111  sentence, but only one of the pairs was in the correct order. Participants needed to indicate which

112  of the two pairs was presented in the sentence before. Presented options for the sentence condition

113  were always a grammatically and semantically plausible combination of words. See Table 1 for an

114  example of the tasks for each condition (sentences and word lists). The three active tasks required

115  participants to focus on the syllabic (syllable task), word (word task), or phrasal (word combination

116  task) timescales.

117       *Procedure.* At the beginning of each trial, participants were instructed to look at a fixation

118  cross presented at the middle of the screen on a grey background. Audio recordings were presented

119  after a random interval between 1.5-3 seconds; 1 second after the end of the audio, the task was

120  presented. For the word and word combination task, this was the presentation of visual stimuli.

121  For the syllable task, this entailed presenting the sound fragments one after each other (with a

122  delay of 0.5 seconds in between). For the passive task, this was the instruction to press a button to

123  continue. In total there were eight blocks (two conditions * four tasks) each lasting about 8 minutes.

124  The order of the blocks was pseudo-randomized by independently randomizing the order of the

125  tasks and the conditions. We then always presented the same task twice in a row to avoid task-

126  switching costs. As a consequence, condition was always alternated (a possible order of blocks

127  would be: passive-sentence, passive-word list, word-sentence, word-word list, syllable-sentence,

128  syllable-word list, word combination-sentence, word combination-word list). After the main

129  experiment, an auditory localizer was collected which consisted of listening to 200ms sinewave

130  and broadband sounds (centred at 0.5, 1, and 2 kHz; for the broadband at a 10% frequency band)

131  at approximately equal loudness. Each sound had a 50ms linear on and off ramp and was presented

132  for 30 times (with random inter-stimulus interval between 1 and 2 seconds).

133   At arrival, participants filled out a screening. Electrodes to monitor eye movements and

134   heart beat were placed (left mastoid was used as ground electrode) at an impedance below 15

135   kiloOhm. Participants wore metal free clothes and fitted earmolds on which two of the three head

136   localizers were placed (together with a final head localizer placed at the nasion). They then

137   performed the experiment in the MEG. MEG was recorded using a 75-channel axial gradiometer

138   CTF MEG system at a sampling rate of 1.2 kHz. After every block participants had a break, during

139   which head position was corrected [19]. After the session, the headshape was collected using

140   Polhemus digitizer (using as fiducials the nasion and the entrance of the ear canals as positioned

141   with the earmolds). For each participant, an MRI was collected with a 3 T Siemens Skyra system

142   using the MPRAGE sequence (1mm isotropic). Also for the MRI acquisition participants wore the

143   earmolds with vitamin pills to optimize the alignment.

144   *Behavioural analysis.* We performed a linear mixed model analysis with fixed factors task

145   (syllable, word, and word combination) and condition (sentence and word list) as implemented by

146   lmer in R4.1.0. The dependent variable was accuracy. First, any outliers were removed (values

147   more extreme than median± 2.5 IQR). Then, we investigated what the best random model was,

148   including a random intercept or a random slope for one or two of the factors. The models with

149   varying random factors were compared with each other using an ANOVA. With no significant

150   difference, the model with the lowest number of factors was included (with minimally a random

151   intercept). Finally, lsmeans was used for follow-up tests using the kenward-roger method to

152   calculate the degrees of freedom from the linear mixed model. For significant interactions, we

153   investigated the effect of condition per task. For main effects, we investigated pairwise

154   comparisons. We corrected for multiple comparisons using adjusted Bonferroni corrections. For

155   all further reported statistical analyses for the MEG data, we followed the same procedure (except

156   that there was one more level of task, i.e. the passive task). To avoid exploding the amount of

157   comparisons, we a-priori decided for any task effects in the MEG analysis to only compare the

158   individual tasks with the phrase task.

159   *MEG pre-processing.* First source models from the MRI were made using a surface-based

160   approach in which grid points were defined on the cortical sheet using the automatic segmentation

161   of freesurfer6.0 [20] in combination with pre-processing tools from the HCP workbench1.3.2 [21]

162   to down-sample the mesh to 4k vertices per hemisphere. The MRI was co-registered to the MEG

163    using the previously defined fiducials as well as an automatic alignment of the MRI to the

164    Polhemus headshape using the Fieldtrip20211102 software [22].

165          Pre-processing involved epoching the data between -3 and +7.9 seconds (+3 relative to the

166    longest sentence of 4.9 sec) around sentence onset. We applied a dftfilter at 50, 100 and 150 Hz to

167    remove line noise, a Butterworth bandpass filter between 0.6 and 100 Hz, and performed baseline

168    correction (-0.2-0 sec baseline). Trials with excessive movements or squid jumps were removed

169    via visual inspection (20.1±18.5 trials removed; mean±standard deviation). Then data was

170    resampled to 300 Hz and we performed ICA decomposition to correct for eye blinks/movement

171    and heart beat artefacts (4.7±0.99 components removed; mean±standard deviation). Trials with

172    remaining artefacts were removed by visual inspection (11.3±12.4 trials removed; mean±standard

173    deviation). Then we applied a lcmv filter to transform the data to have single-trial source space

174    representations. A common filter across all trials was calculated using a fixed orientation and a

175    lambda of 5%. We only extracted time courses for our regions of interest (superior temporal gyrus

176    [1,29,32,33], medial temporal gyrus [6,8,14], and inferior frontal cortex [17,18,19]; numbers

177    correspond to label-coding from the aparc parcellations implemented in Freesurfer). These time

178    courses were baseline corrected (-0.2 to 0 seconds). To reduce computational load and to ensure

179    that we used relevant data within the ROI, we extracted the top 20 PCA components per ROI for

180    all following analyses based on a PCA using the time window of interest (0.5-3.7 seconds; 0.5 to

181    ensure that all initial evoked responses were not included and 3.7 as it corresponds to the shortest

182    trials).

183          *Mutual information analysis.* First, we extracted the speech envelopes by following

184    previous procedures [12, 13, 23]. The acoustic waveforms (third-order Butterworth filter) were

185    filtered in eight frequency bands (100-8000 Hz) equidistant on the cochlear frequency map [24].

186    The absolute of the Hilbert transform was computed, we low-passed the data at 100 Hz (third order

187    Butterworth) and then down-sampled to 300 Hz (matching the MEG sampling rate). Then, we

188    averaged across all bands.

189          Mutual information (MI) was calculated between the filtered speech envelopes and the

190    filtered MEG data at three different frequency bands corresponding to information content at

191    different linguistic hierarchical levels: phrase (0.8-1.1 Hz), word (1.9-2.8 Hz), and syllable (3.5-

192    5.0 Hz). Our main analysis focusses on the phrasal band, as that is where our previous study found

193    the strongest effects [12], but for completeness we also report on the other bands. Mutual

194    information was estimated after the evoked response (0.5 sec) until the end of the stimulus at five

195    different delays (60, 80, 100, 120, and 140 ms) and averaged across delays between the phase

196    estimations of the envelopes and MEG data. A single MI value was generated per condition per

197    ROI by concatenating all trials before calculating the MI (MEG and speech). Statistical analysis

198    was performed per ROI per frequency band.

199        *Power analysis.* Power analysis was performed to compare the MI results with absolute

200    power changes, as any MI differences could be a consequence of signal-to-noise differences in the

201    original data (which would be reflected in power effects). We first extracted the time-frequency

202    representation for all conditions and ROIs separately. To do so, we performed a wavelet analysis

203    with a width of 4, with a frequency of interest between 1 and 30 (step size of 1) and time of interest

204    between -0.2 and 3.7 sec (step size of 0.05 sec). We extracted the logarithm of the power and

205    baseline corrected the data in the frequency domain using a -0.3 and -0.1 sec window. For four

206    different frequency bands (delta: 0.5-3.0 Hz; theta: 3.0-8.0 Hz; alpha: 8.0-15.0 Hz; beta: 15.0-25.0

207    Hz) we extracted the mean power in the 0.5-3.7 sec time window per task, condition and ROI.

208    Again, our main analysis focusses on the delta band as that is where the main previous results were

209    found [12], but we also report on the other bands for completeness. For each ROI we performed

210    the statistical analysis on power as described in the behavioural analysis.

211        *Connectivity analysis.* For the coherence analysis we repeated all pre-processing as in the

212    power analysis, but separately for the left and right hemisphere (as we did not expect connections

213    for PCA across hemispheres), after which we averaged the connectivity measure (using the Fourier

214    spectrum and not the power spectrum). We used the debiased weighted phase lag index (WPLI)

215    for our connectivity measure, which ensures that no zero-lag phase differences are included in the

216    estimation (avoiding effects due to volume conduction). All connections between the three ROIs

217    were investigated for the mean WPLI for the four different frequency bands in the 0.5-3.7 sec time

218    window. Also in this case, the same statistical analysis was applied.
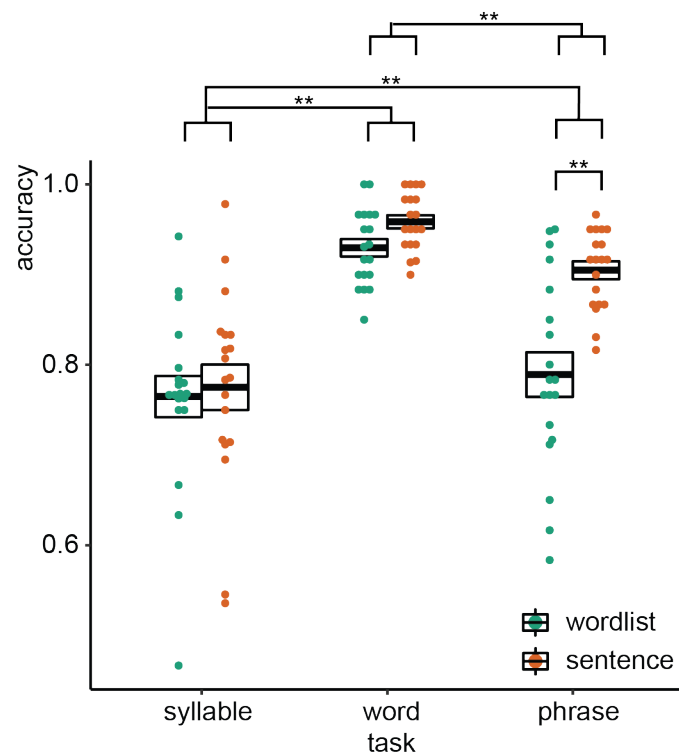
219        *Power control analysis.* The reliability of phase estimations is influenced by the signal-to-

220    noise ratio of the signal [25]. As a consequence, trials with generally high power have more reliable

221    phase estimations compared to low power trials. This could influence any measure relying on this

222    phase estimation, such as MI and connectivity [26, 27]. It is therefore possible that power

223    differences between conditions lead to differences between connectivity or MI. To ensure that our

224    reported effects are not due to signal-to-noise effects, we controlled any significant power

225    difference between conditions for the connectivity and MI analysis. To do this, we iteratively

226    removed the highest and lowest power trials between the mean highest and mean lowest of the two

227    relevant conditions (either collapsing trials across tasks/conditions or using individual conditions).

228    We repeated this until the original condition with the highest power had lower power than the other

229    condition. Then we repeated the analysis and statistics, investigating if the effect of interest was

230    still significant. The control analysis is reported along the main MI and connectivity sections.

231

## Results

*Behaviour.* Overall task performance was above chance and participants complied with task instructions (Figure 1). We found a significant interaction between condition and task $(F(2, 72.0) = 11.51, p < 0.001)$ as well as a main effect of task $(F(2, 19.7) = 44.19, p < 0.001)$ and condition $(F(2, 72.0) = 29.0, p < 0.001)$. We found that only for the word-combination (phrasal-level) task, the sentence condition had a significantly higher accuracy than the word list condition $(t(54.0) = 6.97, p < 0.001)$. For the other two tasks, no significant condition effect was found (syllable: $t(54.0) = 0.62, p = 1.000$; word list:



**Figure 1.** Behavioral results. Accuracy for the three different tasks. Double asterisks indicate significance at the 0.01 level.

$t(54.0) = 1.74, p = 0.176$). Investigating the main effect of task indicated a difference between all tasks (phrase-syllable: $t(18.0) = 3.71, p = 0.003$; phrase-word: $t(22.4) = -6.34, p < 0.001$; syllable-word: $t(19.2)=-8.67, p < 0.001$).

    *Mutual information.* The overall time-frequency response in the three different regions of interest using the top-20 PCA components was as expected, with an initial evoked response followed by a more sustained response to the ongoing speech (Figure 2). From these regions-of-interest, we extracted mutual information in three different frequency bands (phrasal, word, and syllable). Here, we focus on the phrasal band as this is the band that differentiates word lists from sentences and showed the strongest modulation for this contrast in our previous study [12]. Mutual Information results for all other bands are reported in the supplementary materials.

**Figure 2.** Anatomical regions of interests (ROIs). A) ROIs displayed one exemplar participant surface. B) Time-frequency response at each ROI. STG = superior temporal gyrus. MTG = medial temporal gyrus. IFG = inferior frontal gyrus.

259    For the phrasal timescale in STG, we found significantly higher MI in the sentence

260  compared to the word list condition ($F_{(3,126)} = 67.39$, $p < 0.001$; Figure 3). No other effects were

261  significant ($p > 0.1$). This finding paralleled the effect found in Kaufeld et al., [12]. For the MTG,

262  we saw a different picture: Besides the main effect of condition ($F_{(3,126)} = 50.24$, $p < 0.001$), an

263  interaction between task and condition was found ($F_{(3,126)} = 2.948$, $p = 0.035$). We next

264  investigated the effect of condition per task and found for all tasks except the passive task a

265  significant effect of condition, with stronger MI for the sentence condition (passive: $t_{(126)} = 1.07$,

266  $p = 0.865$; syllable: $t_{(126)} = 4.06$, $p = 0.003$; word: $t_{(126)} = 5.033$, $p < 0.001$; phrase: $t_{(126)} =$

267  $4.015$, $p = 0.003$). For the IFG, we found a main effect of condition ($F_{(3,108)} = 21.89$, $p < 0.001$)

268  as well as a main effect of task ($F_{(3,108)} = 2.74$, $p = 0.047$). The interaction was not significant



**Figure 3.** Mutual information (MI) analysis at the phrasal band (0.8-1.1 Hz) for the three different ROIs. Single and double asterisks indicate significance at the 0.05 and 0.01 level. T indicates trend level significance ($p < 0.1$). Inset at the top left of the graph indicate whether a main effect of condition was present (with higher MI for sentences versus wordlists).

12

269    $(F(3,108) = 1.49, p = 0.220)$. Comparing the phrasal task with the other tasks indicated higher MI

270    for the phrasal compared to the word task $(t(111) = 2.50, p = 0.028)$. We also found a trend for the

271    comparison between the phrasal and the syllable task $(t(111) = 2.17, p = 0.064)$, as well as the

272    phrasal and the passive task $(t(111) = 2.25, p = 0.052)$.

273          For the word and syllable frequency bands no interactions were found (all $p > 0.1$;

274    Supplementary Figure 1 and 2). For all six models there was a significant effect of condition, with

275    stronger MI for word lists compared to sentences (all $p < 0.001$). The main effect of task was not

276    significant in any of the models ($p > 0.1$; for the MTG syllable level there was a trend: $F(3,126) =$

277    $2.40, p = 0.071$).

278          When running the power control analysis, we did not find that significant effects in power

279    differences (see next section; mostly due to main effects of condition) influenced our tracking

280    results for any of the bands investigated.

281          *Power.* We repeated the linear mixed modelling using power instead of MI to investigate

282    if power changes paralleled the MI effects. For the delta band, we found for the STG a main effect

283    of condition $(F(1,18) = 6.11, p = 0.024)$ and task $(F(3,108) = 3.069, p = 0.031)$. For the interaction

284    we found a trend $(F(3,108) = 2.620, p = 0.054)$. Overall sentences had stronger delta power than

285    word lists. We found lower power for the phrase compared to the passive task $(t(111) = 2.31, p =$

286    $0.045)$ and lower power for the phrase compared to the syllable task $(t(111) = 2.43, p = 0.034)$.

287    There was no significant difference between the phrase and word task $(t(111) = 0.642, p = 1.00)$.

288          The MTG delta power effect overall paralleled the STG effects with a significant condition

289    $(F(1,124.94) = 12.339, p < 0.001)$ and task effect $(F(3,124.94) = 4.326, p = 0.006)$. The interaction

290    was trend significant $(F(3,124.94) = 2.58, p = 0.056)$. Pairwise comparisons of the task effect

291    showed significantly stronger power for the phrase compared to the passive task $(t(128) = 2.98, p$

292    $= 0.007)$ and lower power for the phrase compared to the syllable task $(t(128) = 3.10, p = 0.024)$.

293    The passive-word comparison was not significant $(t(128) = 2.577, p = 0.109)$. Finally, for the IFG
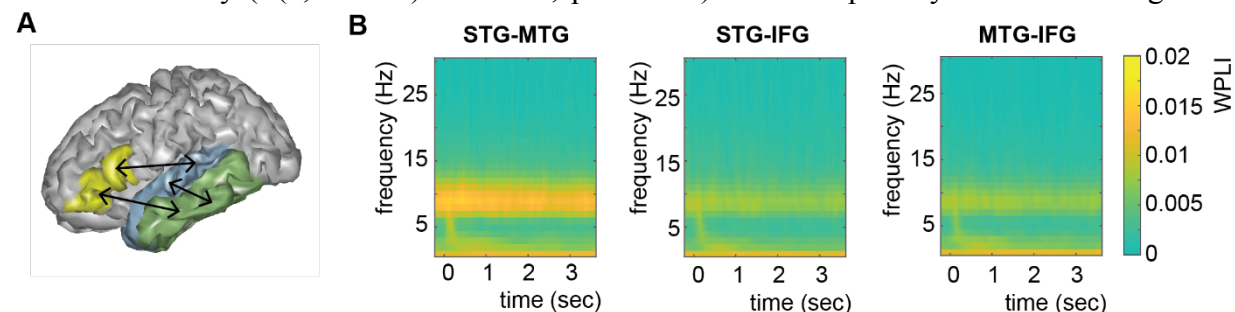
**Figure 6.** WPLI effects for the different ROIs. Single and double asterisks indicate significance at the 0.05 and 0.01 level after correcting for power differences between the two conditions (we plot the original data, not corrected for power, as we can only perform pairwise power and consequently data will be different for each control).

**Figure 4.** Power effects for the different ROIs. Single and double asterisks indicate significance at the 0.05 and 0.01 level. T indicates trend significance (p < 0.1) Inset at the left top of the graph indicate whether a main effect of condition was present (with higher activity for sentences versus wordlists).

294    we only found a trend effect for condition (F(1,123.27) = 4.15, p = 0.057), with stronger delta

295    power in the sentence condition.

296            The results for all other bands can be found in the supplementary materials (Supplementary

297    figure 3-5). In summary, no interaction effects were found for any of the models (all p > 0.1). In

298    all bands, power was generally higher for sentences than for word lists. Any task effect generally

299    showed stronger power for the lower hierarchical level (e.g. generally higher power for passive

300    versus phrasal tasks).

301            *Connectivity*. Overall connectivity patterns showed the strongest connectivity in the delta

302    and alpha frequency band (Figure 5). In the delta band, we found a main effect of task for the STG-

303    IFG connectivity (F(3, 122.06) = 4.1078, p = 0.008). Follow-up analysis showed a significant



**Figure 5.** Connectivity pattern between anatomical regions of interests (ROIs). A) ROI connections displayed one exemplar participant surface. B) Time-frequency weighted phase-lagged index (WPLI) response at each ROI.

14

304     difference between the phrasal and passive task ($t(125) = 3.254$, $p = 0.003$). The other comparisons

305     with the phrasal task were not significant. The effect of task remained significant even when

306     correcting for power differences between the passive and phrasal task ($F(1, 53.02) = 12.39$, $p <$

307     $0.001$; note the change in degrees of freedom as only the passive and phrasal task were included

308     in this mixed model as any power correction is done on pairs). Initially, we also found main effects

309     of condition for the delta and beta band for the MTG-IFG connectivity (stronger connectivity for

310     the sentence compared to the word list condition), however after controlling for power, these

311     effects did not remain significant.

**Discussion**

In the current study, we investigated the effects of 'additional' tasks on the neural tracking of sentences and word lists at temporal modulations that matched phrasal rates. Different nodes of the language network showed different tracking patterns. In STG, we found stronger tracking of phrase-timed dynamics in sentences compared to word lists, independent of task. However, in MTG we found this sentence-improved tracking only for active tasks. In IFG we also found an overall increase of tracking for sentences compared to word lists. Additionally, stronger phrasal tracking was found for the phrasal-level word-combination task compared to the other tasks (independent of stimulus type; note that for the syllable and passive comparison we found a trend), which was paralleled with increased IFG-STG connectivity in the delta band for the word combination task. This suggests that tracking at phrasal time-scales depends both on the linguistic information present in the signal, and on the specific task that is performed.

The findings reported in this study are in line with previous results, with overall stronger tracking of low frequency information in the sentences compared to the word list condition [12]. Crucially, for the stimuli used in our study it has been shown that the condition effects are not due to acoustic differences in the stimuli and also do not occur for reversed speech [12]. It is therefore most likely that our results reflect an automatic extraction of relevant phrase-level information in sentences, indicating the automatic processing of participants as they understand the meaning of the speech they hear using structural sentence information [2, 17, 28]. Overall, it did not seem that making participants pay attention to the temporal dynamics at the same hierarchical level through an additional task – instructing them to remember word combinations at the phrasal rate during word list presentation – could counter this main effect of condition.

Even though there was an overall main effect of condition, task did influence neural responses. Interestingly, the task effects differed for the three regions of interest. In the STG, we found no task effects, while in the MTG we found an interaction between task and condition. In the MTG increased phrasal-level tracking for sentences only occurred when participants were specifically instructed to perform an active task on the materials. It therefore seems that in MTG all levels of linguistic information are used to do an active language operation on the stimuli. This is in line with previous theoretical and empirical research suggesting a strong top-down modulatory response of speech processing in which predictions flow from the highest hierarchical

16

342 levels (e.g. syntax) down to lower levels (e.g. phonemes) to aid language understanding [5, 29,
343 30]. As in the word list condition no linguistic information is present at the phrasal-rate, this
344 information cannot be used to provide useful feedback for processing lower-level linguistic
345 information. Instead, it could have been expected that the same type of increased tracking should
346 have happened at the word-rate rather than the phrasal-rate for word lists (i.e., stronger word-rate
347 tracking for word lists for the active tasks versus passive task). This effect was not found; this
348 could either be attributed to different computational operations occurring at different hierarchical
349 levels or to signal-to-noise/signal detection issues.

350 It is interesting that MTG, but not STG, showed an interaction effect. Both MTG and STG
351 are strong hubs for language processing and have been involved in many studies which contrasted
352 pseudo-words and words [31-33]. It is likely that STG does the lower-level processing of the two
353 regions, as it is earlier in the cortical hierarchy, thereby being more involved in initial segmentation
354 and initial phonetic abstraction rather than a lexical interface [31]. This could also explain why
355 STG does not show task specific tracking effects; STG could be earlier in a workload bottleneck,
356 receiving feedback independent of task, while MTG-feedback is recruited only when active
357 linguistic operations are required. Alternatively, it is possible that either small differences in the
358 acoustics are detected by STG (even though this effect was not previously found with the same
359 stimuli [12]), or that our blocked designed put participants in a sentence or word list "mode" which
360 could have influenced the state of these early hierarchical regions.

361 The IFG was the only region that showed an increase in phrasal-rate tracking specifically
362 for the word-combination task. Note, however, that this was a weak effect, as the comparison
363 between the phrase task and the syllable and passive task only reached a trend towards significance.
364 Nonetheless, this effect is interesting for understanding the role of IFG in language. Traditionally,
365 IFG has been viewed as a hub for articulatory processing [31], but its role during speech
366 comprehension, specifically in syntactic processing, has also been acknowledged [1, 29, 34-36].
367 Integrating information across time and relative timing is essential for syntactic processing [2, 35,
368 37], and IFG feedback has been shown to occur in temporal dynamics at lower (delta) rates during
369 sentence processing [38, 39]. However, it has also been shown that syntactic-independent verbal
370 working memory chunking tasks recruit the IFG [35, 40-42]. This is in line with our findings that
371 show that IFG is involved when we need to integrate across temporal domains either in a language-

17

372   specific domain (sentences versus word lists) or for language-unspecific tasks (word combination

373   versus other tasks). We also show increased delta-connectivity with STG for the only temporal-

374   integration tasks in our study (i.e., the word combination task), independent of the linguistic

375   features in the signal. Our results therefore support a role of the IFG as a combinatorial hub
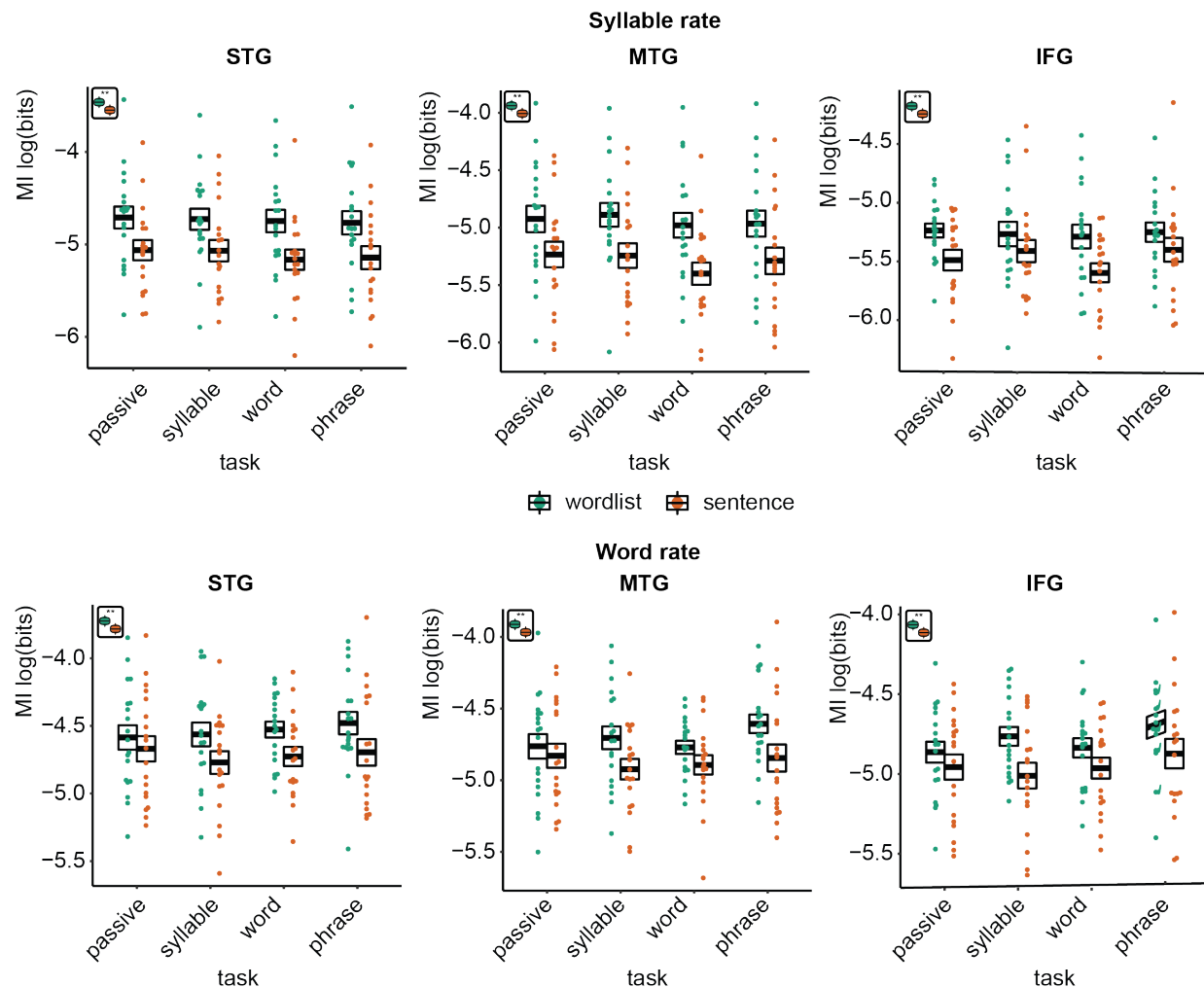
376   integrating information across time [43-45].

377   In the current study we investigated power as a neural readout during language

378   comprehension from speech. This was both to ensure that any tracking effects we found were not

379   due to overall signal-to-noise (SNR) differences, as well as to investigate task-and-condition

380   dependent computations. SNR is better for conditions with higher power, which therefore leads to

381   more reliable phase estimations, critical for computing MI as well as connectivity [25]. We will

382   therefore discuss the power differences as well as their consequences for the interpretation of the

383   MI and connectivity results. Generally, it seemed that there was stronger power in the sentence

384   compared to the word list condition in the delta band. However, the pattern was very different than

385   the MI pattern. For the power, the word list-sentence difference was the biggest in the passive

386   condition. In contrast, for the MI there was either no task difference (in STG) or even a stronger

387   effect for the active tasks (MTG; note that the power interaction was trend significant STG and

388   MTG). We therefore think it unlikely that our MI effects were purely driven by SNR differences,

389   and our power control analysis is consistent with this interpretation. Instead, power seems to reflect

390   a different computation than the tracking, where more complex tasks generally lead to lower power

391   across almost all tested frequency bands. As most of our frequency bands are on the low side of

392   the spectrum (up to beta), it is expected that more complex tasks reduce the low-frequency power

393   [46, 47]. It is interesting to observe that this did not reduce the connectivity for the delta band

394   between IFG and STG, but rather increased it. It has been suggested that low power can potentially

395   increase the available computational space, as it increases the entropy in the signal [48, 49].

396   Finally, in the power comparisons for the theta, alpha, and beta band we found stronger power for

397   the sentence compared to the word list condition, which could reflect that listening to a natural

398   sentence is generally less effortful than listening to a word list.

399   In the current manuscript we describe tracking of ongoing temporal dynamics. However,

400   the neural origin of this tracking is unknown. While we can be sure that modulations in the phrasal-

401   rate follow changes in the phrasal-rate of the acoustic input, it is unclear what the mechanism

18

402 behind this modulation is. It is possible that there is stronger alignment of neural oscillations with

403 the acoustic input at the phrasal rate [50, 51]. However, it could as well be that there is a phrasal

404 time-scale or slower operation happening while processing the incoming input (which de facto is

405 at the same time-scale as the phrasal structure occurring in the input). This operation, in response

406 to stimulus input, could just as well induce the patterns we observe [52, 53]. Finally, it is possible

407 that there are specific responses as a consequence of the syntactic structure, task, or statistical

408 regularities occurring as specific events at phrasal time-scales [51, 54, 55].

409     It is difficult to decide on the most natural task in an experimental setting, that best reflects

410 how we use language in a natural setting. This is probably why such a vast number of different

411 tasks have been used in the literature. Our study (and many before us) indicates that during passive

412 listening, we naturally attend to all levels of linguistic hierarchy. This is consistent with the widely

413 accepted notion that the meaning of a natural sentence requires understanding the compositionality

414 of words in a grammatical structure. For most research questions in language, it therefore is

415 understandable to use a task that mimics this automatic natural understanding of a sentence. Here,

416 we show that automatic understanding of linguistic information, and all the processing that this

417 entails, cannot be countered to substantially change the consequences for neural readout, even

418 when explicitly instructing participants to pay attention to particular time-scales.

419 **Supplementary figures**



**Supplementary Figure 1.** Mutual information (MI) analysis at the syllable (3.5-5.0 Hz) and word rate (1.9-2.8 Hz) for the three different ROIs. Double asterisks indicate significance at the 0.01 level. Inset at the top left of the graph indicate whether a main effect of condition was present (with higher MI for wordlists versus sentences).

420

421

**Supplementary figure 2.** Power effects for the different ROIs and different bands. Single and double asterisks indicate significance at the 0.05 and 0.01 level. T indicates trend significance (p < 0.1) Inset at the top left of the graph indicate whether a main effect of condition was present (with higher activity for sentences versus wordlists).

422

**Supplementary figure 3.** WPLI effects for the different ROIs and different bands. Connectivity is displayed before correcting for power differences. None of the effects survived correcting for power differences.

423

**Acknowledgments**

**References**

1.      Friederici AD. The brain basis of language processing: from structure to function. Physiol Rev. 2011;91(4):1357-92.

2.      Martin AE. A compositional neural architecture for language. J Cognit Neurosci. 2020:1-20.

3.      Halle M, Stevens K. Speech recognition: A model and a program for research. IRE transactions on information theory. 1962;8(2):155-9.

4.      Marslen-Wilson WD, Welsh A. Processing interactions and lexical access during word recognition in continuous speech. Cognitive psychology. 1978;10(1):29-63.

5.      Martin AE. Language processing as cue integration: Grounding the psychology of language in perception and neurophysiology. Frontiers in psychology. 2016;7:120.

6.      Pinker S, Jackendoff R. The faculty of language: what's special about it? Cognition. 2005;95(2):201-36.

7.      Giraud AL, Poeppel D. Cortical oscillations and speech processing: emerging computational principles and operations. Nat Neurosci. 2012;15(4):511-7.

8.      Peelle JE, Davis MH. Neural oscillations carry speech rhythm through to comprehension. Frontiers in Psychology. 2012;3.

9.      Rosen S. Temporal information in speech: acoustic, auditory and linguistic aspects. Philosophical Transactions of the Royal Society of London Series B: Biological Sciences. 1992;336(1278):367-73.

10.     Ding N, Patel AD, Chen L, Butler H, Luo C, Poeppel D. Temporal modulations in speech and music. Neurosci Biobehav Rev. 2017;81:181-7.

11.     Pellegrino F, Coupé C, Marsico E. A cross-language perspective on speech information rate. Language. 2011:539-58.

12.     Kaufeld G, Bosker HR, Ten Oever S, Alday PM, Meyer AS, Martin AE. Linguistic structure and meaning organize neural oscillations into a content-specific hierarchy. J Neurosci. 2020;40(49):9467-75.

13.     Keitel A, Gross J, Kayser C. Perceptually relevant speech tracking in auditory and motor cortex reflects distinct linguistic features. PLoS Biol. 2018;16(3):e2004473.

14.     Luo H, Poeppel D. Phase patterns of neuronal responses reliably discriminate speech in human auditory cortex. Neuron. 2007;54(6):1001-10.

15.     Zoefel B, Archer-Boyd A, Davis MH. Phase entrainment of brain oscillations causally modulates neural responses to intelligible speech. Curr Biol. 2018;28(3):401-8. e5.

16.     Doelling KB, Arnal LH, Ghitza O, Poeppel D. Acoustic landmarks drive delta–theta oscillations to enable speech comprehension by facilitating perceptual parsing. NeuroImage. 2014;85:761-8.

17.     Ding N, Melloni L, Zhang H, Tian X, Poeppel D. Cortical tracking of hierarchical linguistic structures in connected speech. Nat Neurosci. 2016;19(1):158-64.

18.     Kayser SJ, Ince RA, Gross J, Kayser C. Irregular speech rate dissociates auditory cortical entrainment, evoked responses, and frontal alpha. J Neurosci. 2015;35(44):14691-701.

19.     Stolk A, Todorovic A, Schoffelen J-M, Oostenveld R. Online and offline tools for head movement compensation in MEG. NeuroImage. 2013;68:39-48.

20.     Fischl B. FreeSurfer. NeuroImage. 2012;62(2):774-81.

21.     Glasser MF, Sotiropoulos SN, Wilson JA, Coalson TS, Fischl B, Andersson JL, et al. The minimal preprocessing pipelines for the Human Connectome Project. NeuroImage. 2013;80:105-24.
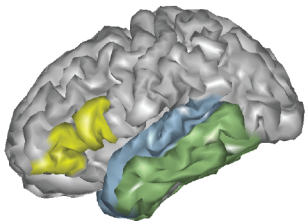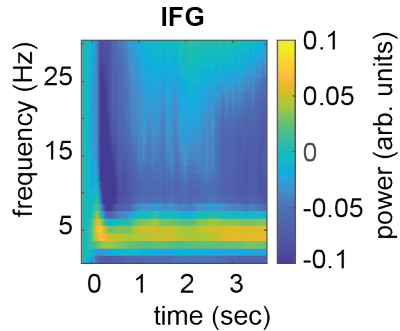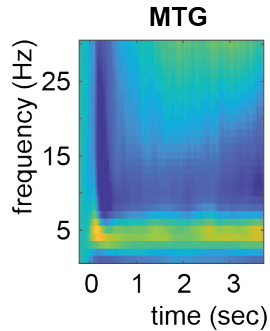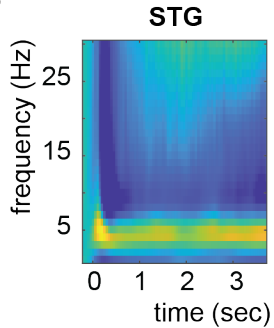
470   22.     Oostenveld R, Fries P, Maris E, Schoffelen J-M. FieldTrip: open source software for advanced
471   analysis of MEG, EEG, and invasive electrophysiological data. Computational intelligence and
472   neuroscience. 2011;2011:1.
473   23.     Gross J, Hoogenboom N, Thut G, Schyns P, Panzeri S, Belin P, et al. Speech rhythms and
474   multiplexed oscillatory sensory coding in the human brain. PLoS Biol. 2013;11(12):e1001752.
475   24.     Smith ZM, Delgutte B, Oxenham AJ. Chimaeric sounds reveal dichotomies in auditory perception.
476   Nature. 2002;416(6876):87-90.
477   25.     Zar JH. Biostatistical Analysis. 4 ed. Englewood Cliffs, New Jersey: Prentice Hall; 1998.
478   26.     Bastos AM, Schoffelen J-M. A tutorial review of functional connectivity analysis methods and their
479   interpretational pitfalls. Frontiers in systems neuroscience. 2016;9:175.
480   27.     Ince RA, Giordano BL, Kayser C, Rousselet GA, Gross J, Schyns PG. A statistical framework for
481   neuroimaging data analysis based on mutual information estimated via a gaussian copula. Hum Brain
482   Mapp. 2017;38(3):1541-73.
483   28.     Yahav PH-s, Golumbic EZ. Linguistic processing of task-irrelevant speech at a Cocktail Party. Elife.
484   2021;10:e65096.
485   29.     Hagoort P. The core and beyond in the language-ready brain. Neurosci Biobehav Rev.
486   2017;81:194-204.
487   30.     Federmeier KD. Thinking ahead: The role and roots of prediction in language comprehension.
488   Psychophysiology. 2007;44(4):491-505.
489   31.     Hickok G, Poeppel D. The cortical organization of speech processing. Nature Reviews
490   Neuroscience. 2007;8(5):393-402.
491   32.     Dronkers NF. The neural architecture of the language comprehension network: converging
492   evidence from lesion and connectivity analyses. Frontiers in systems neuroscience. 2011;5:1.
493   33.     Vouloumanos A, Kiehl KA, Werker JF, Liddle PF. Detection of sounds in the auditory stream: event-
494   related fMRI evidence for differential activation to speech and nonspeech. J Cognit Neurosci.
495   2001;13(7):994-1005.
496   34.     Nelson MJ, El Karoui I, Giber K, Yang X, Cohen L, Koopman H, et al. Neurophysiological dynamics
497   of phrase-structure building during sentence processing. Proc Natl Acad Sci. 2017;114(18):E3669-E78.
498   35.     Dehaene S, Meyniel F, Wacongne C, Wang L, Pallier C. The neural representation of sequences:
499   from transition probabilities to algebraic patterns and linguistic trees. Neuron. 2015;88(1):2-19.
500   36.     Zaccarella E, Meyer L, Makuuchi M, Friederici AD. Building by syntax: the neural basis of minimal
501   linguistic structures. Cereb Cortex. 2017;27(1):411-21.
502   37.     Martin AE, Doumas LA. Predicate learning in neural systems: using oscillations to discover latent
503   structure. Current Opinion in Behavioral Sciences. 2019;29:77-83.
504   38.     Park H, Ince RA, Schyns PG, Thut G, Gross J. Frontal top-down signals increase coupling of auditory
505   low-frequency oscillations to continuous speech in human listeners. Curr Biol. 2015;25(12):1649-53.
506   39.     Keitel A, Gross J. Individual human brain areas can be identified from their characteristic spectral
507   activation fingerprints. PLoS Biol. 2016;14(6):e1002498.
508   40.     Osaka N, Osaka M, Kondo H, Morishita M, Fukuyama H, Shibasaki H. The neural basis of executive
509   function in working memory: an fMRI study based on individual differences. NeuroImage. 2004;21(2):623-
510   31.
511   41.     Fegen D, Buchsbaum BR, D'Esposito M. The effect of rehearsal rate and memory load on verbal
512   working memory. NeuroImage. 2015;105:120-31.
513   42.     Koelsch S, Schulze K, Sammler D, Fritz T, Müller K, Gruber O. Functional architecture of verbal and
514   tonal working memory: an FMRI study. Hum Brain Mapp. 2009;30(3):859-73.
515   43.     Gelfand JR, Bookheimer SY. Dissociating neural mechanisms of temporal sequencing and
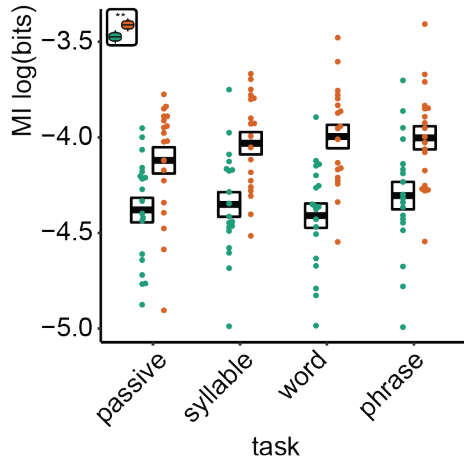516   processing phonemes. Neuron. 2003;38(5):831-42.

517    44.    Schapiro AC, Rogers TT, Cordova NI, Turk-Browne NB, Botvinick MM. Neural representations of
518    events arise from temporal community structure. Nat Neurosci. 2013;16(4):486-92.
519    45.    Skipper JI. The NOLB model: A model of the natural organization of language and the brain. 2015.
520    46.    Jensen O, Mazaheri A. Shaping functional architecture by oscillatory alpha activity: gating by
521    inhibition. Front Hum Neurosci. 2010;4.
522    47.    Klimesch W. EEG alpha and theta oscillations reflect cognitive and memory performance: a review
523    and analysis. Brain Res Rev. 1999;29(2):169-95.
524    48.    Hanslmayr S, Staudigl T, Fellner M-C. Oscillatory power decreases and long-term memory: the
525    information via desynchronization hypothesis. Front Hum Neurosci. 2012;6:74.
526    49.    Ten Oever S, Sack AT. Oscillatory phase shapes syllable perception. Proc Natl Acad Sci.
527    2015;112(52):15833-7.
528    50.    Lakatos P, Karmos G, Mehta AD, Ulbert I, Schroeder CE. Entrainment of neuronal oscillations as a
529    mechanism of attentional selection. science. 2008;320(5872):110-3.
530    51.    Obleser J, Kayser C. Neural entrainment and attentional selection in the listening brain. Trends
531    Cogn Sci. 2019;23(11):913-26.
532    52.    Meyer L, Sun Y, Martin AE. Synchronous, but not entrained: Exogenous and endogenous cortical
533    rhythms of speech and language processing. Language, Cognition and Neuroscience. 2019:1-11.
534    53.    Zoefel B, Ten Oever S, Sack AT. The involvement of endogenous neural oscillations in the
535    processing of rhythmic input: More than a regular repetition of evoked neural responses. Front Neurosci.
536    2018;12:95.
537    54.    Ten Oever S, Martin AE. An oscillating computational model can track pseudo-rhythmic speech by
538    using linguistic predictions. Elife. 2021;10:e68066.
539    55.    Frank SL, Yang J. Lexical representation explains cortical entrainment during speech
540    comprehension. PloS one. 2018;13(5):e0197304.

541

**A**

**B**

STG    MTG    IFG

STG / MTG / IFG

MI log(bits) vs task (passive, syllable, word, phrase)

Legend: wordlist (green), sentence (orange)

**STG** — **MTG** — **IFG**

power (arb. units) vs task (passive, syllable, word, phrase)

wordlist    sentence

Figure showing Syllable rate (top row) and Word rate (bottom row) results across STG, MTG, and IFG regions, plotting MI log(bits) by task (passive, syllable, word, phrase) for wordlist (green) and sentence (orange) conditions.

**Theta**

**STG-MTG beta** | **STG-IFG beta** | **MTG-IFG beta**

**Alpha**

**STG-MTG beta** | **STG-IFG beta** | **MTG-IFG beta**

**Beta**

**STG-MTG beta** | **STG-IFG beta** | **MTG-IFG beta**

wordlist | sentence