1    **Title:** Cortical tracking of continuous speech under bimodal divided attention

2    **Abbreviated title:** Continuous speech processing under bimodal divided attention

3    Zilong Xie[1†], Christian Brodbeck[2†], Bharath Chandrasekaran[3]

4    [1]School of Communication Science and Disorders, Florida State University, Tallahassee, FL

5    [2]Department of Psychological Sciences, University of Connecticut, Storrs, CT

6    [3]Department of Communication Science and Disorders, University of Pittsburgh, Pittsburgh, PA

7

8    **Corresponding authors**

9    Zilong Xie (zx22c@fsu.edu); Bharath Chandrasekaran (b.chandra@pitt.edu)

10

[†] Zilong Xie and Christian Brodbeck should be considered joint first author.

## Abstract

Speech processing often occurs amidst competing inputs from other modalities, e.g., listening to the radio while driving. We examined the extent to which *dividing* attention between auditory and visual modalities (bimodal divided attention) impacts neural processing of natural continuous speech from acoustic to linguistic levels of representation. We recorded electroencephalographic (EEG) responses when human participants performed a challenging primary visual task, imposing low or high cognitive load while listening to audiobook stories as a secondary task. The two dual-task conditions were contrasted with an auditory single-task condition in which participants attended to stories while ignoring visual stimuli. Behaviorally, the high load dual-task condition was associated with lower speech comprehension accuracy relative to the other two conditions. We fitted multivariate temporal response function encoding models to predict EEG responses from acoustic and linguistic speech features at different representation levels, including auditory spectrograms and information-theoretic models of sublexical-, word-form-, and sentence-level representations. Neural tracking of most acoustic and linguistic features remained unchanged with increasing dual-task load, despite unambiguous behavioral and neural evidence of the high load dual-task condition being more demanding. Compared to the auditory single-task condition, dual-task conditions selectively reduced neural tracking of only some acoustic and linguistic features, mainly at latencies >200 ms, while earlier latencies were surprisingly unaffected. These findings indicate that behavioral effects of bimodal divided attention on continuous speech processing occur not due to impaired early sensory representations but likely at later cognitive processing stages. Crossmodal attention-related mechanisms may not be uniform across different speech processing levels.

# Introduction

Speech processing often occurs amidst competing inputs from other sensory modalities, e.g., listening to the radio while driving. In such situations, listeners must allocate attention across modalities to effectively select the most relevant information within a modality. This raises the question of whether and how *dividing* attention between modalities (e.g., audition and vision; bimodal divided attention) affects the processing of natural continuous speech.

Resource-based theoretical frameworks have been invoked to scaffold the understanding of mechanisms governing crossmodal attention (Wahn & König, 2017). Two contrastive resource-based accounts (modality-specific versus supramodal) yield different hypotheses regarding the effects of bimodal divided attention on continuous speech processing. Per the *modality-specific* account, each sensory modality is allocated a limited pool of attentional resources, and these pools of attentional resources operate independently of each other (Alais et al., 2006; Arrighi et al., 2011; Duncan et al., 1997; Keitel et al., 2013; Parks et al., 2011; Porcu et al., 2014). In contrast, per the *supramodal* account, different sensory modalities share a central, limited pool of attentional resources. The availability of resources to one modality is inversely related to the amount of resources used by other modalities (Broadbent, 1958; Ciaramitaro et al., 2017; Klemen et al., 2009; Macdonald & Lavie, 2011; Molloy et al., 2015).

Empirical evidence regarding bimodal divided attention effects on speech processing primarily comes from experimenter-constrained tasks (e.g., Gennari et al., 2018; Kasper et al., 2014; Mattys et al., 2009, 2014; Mattys & Palmer, 2015; Mattys & Wiget, 2011). Many studies have shown the detrimental effects of bimodal divided attention on the acoustic processing of simplified, controlled speech stimuli (e.g., syllable or single words) (Gennari et al., 2018; Mattys et al., 2014; Mattys & Palmer, 2015; Mattys & Wiget, 2011), which is consistent with the

68   *supramodal* account of attention. Speech processing entails mapping acoustic features into

69   linguistic representations of increasing complexity (Brodbeck & Simon, 2020; Hickok &

70   Poeppel, 2007), raising the question of how bimodal divided attention affects linguistic

71   representations beyond acoustic processing. Behavioral studies with simple speech stimuli

72   indicate that reduced acoustic processing under bimodal divided attention may lead to

73   compensatory changes manifested by increased reliance on higher-order linguistic knowledge

74   during auditory lexical perception (Mattys et al., 2009). However, to date, there is a lack of a

75   systematic and holistic analysis of divided attention-related changes across different levels

76   (acoustic-to-linguistic) of natural continuous speech processing, which is distinctly different

77   from processing simple speech stimuli (Gaston et al., 2022; Hamilton & Huth, 2020).

78        Here, we assessed electroencephalography (EEG) to provide a systematic and holistic

79   analysis of the acoustic and linguistic processing of continuous speech (Brodbeck & Simon,

80   2020; Gillis et al., 2022). The continuous speech paradigm uses the multivariate temporal

81   response function approach (Crosse et al., 2016; Ding & Simon, 2012) to predict neural

82   responses from a combination of hypothesis-driven acoustic and linguistic properties of

83   continuous speech. The predictive power of each speech property is used to quantify the

84   corresponding processing levels (Brodbeck & Simon, 2020; Gillis et al., 2022). The spectro-

85   temporal acoustic properties included envelope-based spectrogram and acoustic onset

86   spectrogram. The linguistic properties included measures of informativeness (surprisal and

87   entropy) based on the information-theoretic framework (Brodbeck et al., 2018). Prior work

88   suggests that both acoustic and linguistic representations are strongly modulated by *selective*

89   attention, within the auditory modality and across modalities. Attentional effects are

90    disproportionality more robust on the linguistic representations than acoustic-based

91    representations (Brodbeck et al., 2018, 2020).

92        Here we integrated the continuous speech paradigm with an audiovisual dual-task

93    paradigm to examine the effects of bimodal divided attention on the acoustic and linguistic

94    processing of continuous speech. In the dual-task paradigm, participants performed a challenging

95    primary visuospatial task that imposed low or high cognitive load while listening to audiobook

96    stories as a secondary task. The two dual-task conditions were contrasted with an auditory

97    single-task condition in which participants attended to the story while ignoring visual stimuli.

98    We hypothesized that compared to the auditory single-task condition, dual-task conditions would

99    lead to reduced acoustic and linguistic representations of continuous speech, especially at high

100   cognitive load. However, we hypothesized that linguistic representations may be affected to a

101   relatively greater extent based on evidence from the literature on selective attention. These

102   hypotheses are aligned with the *supramodal* account of crossmodal attention.


## Materials and Methods

103

### Experimental design

104

105   Bimodal divided attention was manipulated via a dual-task paradigm. Specifically, participants

106   performed a primary visuospatial *n*-back task of varying (high or low) cognitive load (Jaeggi et

107   al., 2007) while listening to continuous speech as a secondary task. We designated the visual task

108   as the primary task to maximize the chance of observing the bimodal divided attention effects on

109   continuous speech processing. The cognitive load of the dual-task paradigm was manipulated via

110   3- and 0-back tasks on the visuospatial stimuli (blue squares; Figure 1A and 1B). The dual-task

111   conditions were contrasted with an auditory single-task condition (Figure 1C), in which

112    participants explicitly attended to the auditory stimuli while ignoring the visual stimuli. To

113    obtain a behavioral measure for the auditory task, participants were instructed to respond to two

114    multiple-choice comprehension questions on the story segments at the end of each trial. Detailed

115    task instructions are presented in the section on *Experimental procedure.*

116        Each task condition consisted of 15 trials of visual stimuli paired with 15 unique story

117    segments and were presented in separate blocks. The order of the story segments was fixed and

118    identical across participants in order to maintain the continuity of the storyline. The order of task

119    conditions was counterbalanced across participants. Each trial of visual stimuli ended later than

120    the corresponding story segment. Such offset gaps were not significantly different across task

121    conditions [$F(2, 42) = .01$, $p = .99$]. The experiment was controlled with E-Prime 2.0.10

122    (Schneider et al., 2002).

123

124                                     **[Fig 1 about here]**

125

126    **Figure 1.** Trial design illustrations for (A) high load dual-task (3-back visual ask), (B) low load

127    dual-task (0-back visual task), and (C) auditory single-task condition. In the two dual-task

128    conditions, the primary task was to respond to the visual stimuli and the secondary task was to

129    attend to auditory stimuli (story segments of about 60 seconds). In the high load condition (A),

130    participants responded only when the current blue square matched the one 3 positions back

131    (examples highlighted in red squares). In the low load condition (B), participants responded only

132    when the current blue square matched the first square in each trial (highlighted in the red square).

133    In the auditory single-task condition (C), participants were instructed to attend to the auditory

134      stimulus and ignore the visual stimuli. At the end of each trial, participants responded to two

135      multiple-choice comprehension questions for the story segments. ISI: interstimulus interval.

136

## Participants

138      Adult native American English speakers (N = 18) were recruited from the Austin, Texas,

139      community. Data from one participant were excluded due to technical problems. Data from

140      another participant were excluded because their story comprehension accuracy was lower for the

141      auditory single-task condition (66.67%) than the two dual-task conditions (73.37% for low load

142      and 76.67% for high load). We interpreted this result as that this participant did not understand or

143      follow the task instructions. The final sample consisted of sixteen participants (18 to 23 years

144      old; 11 females, five males; 14 right-handed and two left-handed). The sample size was selected

145      based on prior work examining the effects of bimodal attention on the neural processing of

146      speech stimuli (e.g., Gennari et al., 2018; Kasper et al., 2014). Previous studies have shown that

147      music training can influence speech processing (e.g., Bidelman & Alain, 2015). Therefore, we

148      recruited only participants without a history of or significant formal music training (<= four

149      years of continuous training, not currently practicing). All participants had normal air and bone-

150      conduction audiometric thresholds, defined as <= 20 dB hearing level for octave frequencies

151      from 0.25 to 8 kHz. The thresholds were measured via an Interacoustics Equinox 2.0 PC-Based

152      Audiometer. Additional inclusion criteria are as follows: no history of psychological or

153      neurological disorders, no use of neuropsychiatric medication, and having normal or corrected-

154      to-normal vision. Before the experiment, all participants provided written, informed consent.

155      Participants received monetary compensation for their participation. The Institutional Review

156      Board at the University of Texas at Austin approved the experimental protocols.

## Stimuli and apparatus

The stimuli were composed of visual and auditory materials. The visual stimuli (Figure 1) were blue squares at one of eight loci around a white fixation cross in the center of a black screen, adapted from Jaeggi et al. (2007). The duration for individual squares was 500 ms, and the interval between consecutive squares was 2500 ms. Twenty-three squares were included in a trial, lasting 69 seconds. The stimuli were displayed on a VIEWPixx/EEG LCD monitor with a scanning LED-backlight design [29.1 cm (height) $\times$ 52.2 cm (width); display resolution: 1920 $\times$ 1080; refresh rate: 120 Hz] at an ~110 cm distance from participants' eyes.

The auditory stimuli were English audiobook stories selected from a classic work of fiction, *Alice's Adventures in Wonderland* (Chapters 1-7, http://librivox.org/alices-adventures-in-wonderland-by-lewis-carroll-5). The audiobook was narrated by an adult male American English speaker at a sampling rate of 22.05 kHz. The chapters were divided into 45 segments (each ~60 seconds long). Each segment began where the story ended in the previous segment. In each segment, silent periods of more than 500 ms were shortened to 500 ms. The story stimuli were presented diotically via insert earphones (ER-3; Etymotic Research, Elk Grove Village, IL) to the participants at a 70 dB sound pressure level. A trial of visual stimuli (23 blue squares) was presented concurrently with each story segment, with the segment beginning later (3 seconds after the onset of the visual trial) and ending earlier relative to the visual trial.

## Experimental procedure

### High and low load dual-task

The cognitive load of the dual-task conditions was manipulated via the visual task. For the high load condition, the visual task required participants to respond when the current blue square matched the one three-position back in the sequence (i.e., 3-back task, Figure 1A). For the low

180    load condition, the visual task required participants to respond when the current blue square

181    matched the first square in the sequence (i.e., 0-back task, Figure 1B). We randomized the

182    location of the first square across trials. Matched squares were treated as targets, and unmatched

183    ones were non-targets. Note that targets could appear only starting from the fourth square in the

184    sequence for a given trial in the 3-back task. In other words, targets would be among the last 20

185    squares in the sequence on a given trial. We designed the 0-back task to match that. Six of the

186    last 20 squares were set as targets for both task conditions, and the remaining 14 were non-

187    targets. The target locations were randomized across trials.

188        Participants responded to the targets by pressing buttons on a game controller.

189    Participants were told that speed and accuracy were equally important. Participants were

190    required to rest their fixations on a white cross in the middle of the screen. To encourage

191    engagement, accuracy feedback on the visual task was displayed after their responses. The

192    number of button presses was not significantly different between 3- and 0-back visual tasks

193    [$t(15)$ = .96, $p$ = .36]. After the visual task, participants responded to two multiple-choice

194    comprehension questions for the auditory stories. Participants had unlimited time to respond to

195    the story questions. No feedback was provided after their responses.

196        Critically, to manipulate the priority of the auditory and visual tasks, participants were

197    instructed to focus primarily on the visual task and attend to the auditory stimulus as a secondary

198    task. They were explicitly told that their data could not be used if their performance on the visual

199    task was poor.

## Auditory single-task

201    In this condition, participants were instructed to focus on the story segments and ignore the

202    visual stimuli (Figure 1C). Participants were required to keep their eyes open and rest their

203  fixations on a white cross in the middle of the screen. At the end of each trial, participants

204  responded to two multiple-choice questions to assess their comprehension of the story segments.

205  Participants had unlimited time to respond to questions. Visual feedback about the accuracy of

206  the story question was displayed following their responses.

# Electrophysiological data acquisition and preprocessing

## Acquisition

209  The experiment was conducted in a dark, acoustically shielded booth. Participants were seated in

210  a comfortable chair during tasks. Electroencephalography (EEG) data were recorded using the

211  Easycap recording cap (www.easycap.de) with 64 actiCAP active electrodes (Brain Products,

212  Gilching, Munich, Germany) at a sampling rate of 5 kHz. The electrode locations were

213  determined according to the extended 10-20 system (Oostenveld & Praamstra, 2001), with a

214  common ground at the Fpz electrode site and TP9 as the reference. The electrode impedances

215  were below 20 kΩ.

216  The EEG data were acquired using BrainVision actiCHAmp amplifier (Brain Products, Gilching,

217  Munich, Germany) linked to BrainVision Pycorder software 1.0.7.

## Preprocessing

219  The EEG data were preprocessed offline in MNE-Python (Gramfort et al., 2013), and the

220  deconvolution analysis was implemented with the Eelbrain package (Brodbeck et al., 2021). The

221  data were re-referenced to the average of the electrodes TP9 and TP10, and then band-pass

222  filtered from 1 to 15 Hz using a zero-phase overlap-add finite impulse response filter (hamming

223  window) with default settings in MNE-Python. Independent component analysis was applied to

224  EEG data combined across the three task conditions in individual participants using the

225    extended-infomax algorithm. Artifact-related components (mainly ocular artifacts) were

226    identified according to the topographical distribution and time course via visual inspection and

227    then removed. After that, the EEG data were segmented into epochs that were time-locked to the

228    onsets of corresponding story segments, and then downsampled to 100 Hz. The maximum

229    possible duration of the epochs was set to 61 seconds.

230    ## Assessing neural tracking of visual and auditory stimuli

231          To assess the neural representation of speech, we used the multivariate temporal response

232    function (mTRF) approach (Crosse et al., 2016; Ding & Simon, 2012). In this approach, the EEG

233    signal is predicted using time-delayed multiple regression. We first generated several visual,

234    acoustic, and linguistic models (see below). Each model was used to define several predictor

235    variables, each implementing a specific linking hypothesis for predicting brain activity from the

236    corresponding model. We then tested the predictive power of different combinations of predictor

237    variables to evaluate which acoustic and linguistic models are associated with predictive power

238    for the EEG data. Each predictor variable thus operationalizes a hypothesis that EEG responses

239    are modulated by a given property of the speech signal, which would indicate neural

240    representations arising from a corresponding acoustic or linguistic model. Figure 2 displays an

241    example of each predictor variable. In the following paragraphs, we provide more detailed

242    descriptions.

243                                      **[Fig 2 about here]**

244

245    **Figure 2.** An excerpt of raw EEG responses from all 64 electrodes (top row) and the predictor

246    variables (subsequent rows) used to model the EEG responses. Note that visual predictors consist

247     of a separate one-dimensional array with impulses for onsets and offsets of the blue squares.

248     They are combined into a single predictor in this example for illustration purposes.

249

## Visual model

251     Because the visual stimuli were temporally sparse, visual responses were modeled analogously

252     to a visual evoked potential. The visual predictor was a one-dimensional time series with an

253     impulse at the onsets and offsets of the blue squares. We did not separate predictors for targets

254     and non-targets because this study was not intended to explore differences in neural processing

255     of visual targets and non-targets, and thus there were not enough targets to estimate stable

256     responses.

## Acoustic model

258     The acoustic model was designed to assess EEG responses related to representations of acoustic

259     spectro-temporal features. All acoustic predictors were derived from 256-band gammatone-based

260     spectrograms of the speech stimuli, with cut-off frequencies from 0.02 to 5 kHz. The 256-band

261     spectrograms were downsampled to 1 kHz and scaled with an exponent of 0.6. A *spectrogram*

262     predictor was then created by summing the 256-band spectrograms in eight logarithmically

263     spaced frequency bands. In addition, an *onset spectrogram* predictor was defined to detect and

264     control for representations of acoustic edges. These were generated using an auditory edge

265     detection model (Brodbeck et al., 2020; Fishbach et al., 2001) and applied to each frequency

266     band of the 256-band spectrograms. The onsets across these 256 bands were also summed into

267     eight logarithmically spaced frequency bands to generate 8-band onset spectrogram predictors.

## Linguistic models

269     Linguistic processing was assessed using information-theoretic models. These models assume

270     that listeners maintain predictive models of speech, which can be linked to brain activity through

271     surprisal and entropy measures (Brodbeck et al., 2018). Previous work suggests that listeners

272     maintain multiple such predictive models, differing in complexity, in parallel (Brodbeck et al.,

273     2022). The predictive models were all defined over phoneme sequences, determined for each

274     stimulus via forced alignment using the Montreal Forced Aligner (MFA) (McAuliffe et al.,

275     2017). The predictors based on the specific information-theoretic models (described in

276     subsequent sections) all consisted of time series with an impulse of variable size at each

277     phoneme onset. In order to provide a control for responses related to linguistic segmentation, two

278     additional predictors were included: A *word onset* predictor with a unit (value of 1) impulse at

279     the onset of each word-initial phoneme and a *phoneme onset* predictor with a unit impulse at the

280     onsets of all other phonemes.

281         **Sublexical model**. The sublexical model assumes that listeners predict upcoming

282     phonemes or speech sounds based on a local context, consisting of only a few preceding sounds.

283     To implement such a model, all sentences from the SUBTLEX-US corpus (Keuleers et al., 2010)

284     were transcribed into phoneme sequences without word boundary markers,  and a 5-gram model

285     (Heafield, 2011) was trained on these phoneme sequences. This model was then applied to the

286     experimental stimuli to compute a probability distribution over phonemes at each phoneme

287     position, conditional on the four preceding phonemes. This distribution was used to calculate a

288     *sublexical surprisal* predictor (the surprisal of encountering phoneme $ph_k$ at position $k$ in the

289     stimulus is $-log_2\left(p(ph_k|context)\right)$, and a *sublexical entropy* predictor (the entropy at

290     phoneme position $ph_k$ reflects the uncertainty about what the next phoneme, $ph_{k+1}$, will be

291     $-\sum_{ph}^{phonemes} p(ph_{k+1} = \text{ph}|context)\log_2 p(ph_{k+1} = \text{ph}|context))$. Surprisal is a measure of

292     the amount of new information provided by a stimulus; a response to sublexical surprisal is thus

293     evidence that listeners integrate information on a sublexical level. A response to entropy

294     additionally suggests that listeners create a probabilistic expectation about future input before

295     encountering the phoneme (Pickering & Gambi, 2018). A response to either of those predictors

296     would provide evidence that listeners maintain a sublexical language model.

297     **Word-form model.** The word-form model aims to predict the surface form of the word

298     that is currently being heard, but without access to any information preceding the word. To

299     implement this model, a phonological lexicon was generated by combining pronunciations from

300     the MFA English dictionary and the Carnegie Mellon University Pronouncing Dictionary

301     (http://www.speech.cs.cmu.edu/cgi-bin/cmudict). The word-form model was implemented based

302     on the cohort model of word recognition (Brodbeck et al., 2018; Marslen-Wilson, 1987). Each

303     word was assigned a prior probability based on its count frequency in the SUBTLEX corpus

304     (Keuleers et al., 2010), substituting a count of 1 for missing words. For each word in the speech

305     stimuli, the cohort model was then implemented by starting with the complete lexicon and, for

306     each subsequent phoneme of the word, incrementally removing words that were not compatible

307     with that phoneme in that position. The changing probability distribution over the lexicon was

308     then used to define two predictors, each with a value at each phoneme position: *phoneme*

309     *surprisal* (log inverse of the posterior probability of the phoneme given the preceding phonemes)

310     and *cohort entropy* (Shannon entropy over all words currently in the cohort,

311     $\sum_{word}^{lexicon} p(word|context) log_2 p(word|context)$). This model implements the hypothesis that

312     listeners recognize words using a probabilistic model that takes into account all the information

313     since the last word boundary (i.e., where the word started), but that does not further take into

314     account any context when considering possible word forms as candidates.

315    **Sentence model.** The sentence model is very similar to the word-form model, with the

316    only difference being that the prior expectation for each word is modulated by the sentence

317    context. To implement this, a lexical 5-gram model (Heafield, 2011) was trained on the whole

318    SUBTLEX-US database (Keuleers et al., 2010). This 5-gram model was used to set the prior

319    probability for each word in the lexicon based on the preceding four words before applying the

320    procedure described for the word-form model above. The same two linking hypotheses were

321    used to define predictor variables (*phoneme surprisal* and *cohort entropy*). The sentence model

322    implements the hypothesis that listeners use a wider context including multiple words, when

323    modulating their phoneme-by-phoneme perception and expectations.

## Estimation of neural tracking

325    We used forward encoding mTRF models to predict EEG responses from the predictors

326    described above. The mTRF models were fitted to the EEG responses at individual electrodes

327    using the boosting algorithm implemented in the Eelbrain toolbox (Brodbeck et al., 2021). The

328    predictive power of the mTRF models was evaluated by how accurately they could predict EEG

329    responses from novel trials of the same condition. This was quantified through the *z*-transformed

330    Pearson's correlation coefficient between predicted and actual EEG responses (i.e., prediction

331    accuracy). Higher prediction accuracy indicates better neural tracking of the corresponding

332    predictor.

333    The mTRFs were estimated separately for each subject and condition using a 5-fold

334    cross-validation strategy. First, the trials were divided into 5 test sets. For each test set, EEG

335    responses were predicted from the average of 4 mTRF models, estimated from the remaining 4

336    datasets, each with 3 of the remaining 4 sets serving as training data, and one as validation set.

337    The mTRFs were generated from a basis of 50 ms Hamming windows with stimulus-EEG lag

338     from -100 to 500 ms (window center). The mTRFs were estimated jointly for all predictors with

339     coordinate descent to minimize the $\ell_1$ error. After each step, the change in error was evaluated in

340     the validation set, and if there was an increase in error, the TRF for the predictor responsible for

341     the change was frozen (in its state before the change). This continued until the whole mTRF was

342     frozen. A single measure of prediction accuracy (fisher $z$-transformed correlation between

343     predicted and measured response, see above) was calculated after concatenating the predicted

344     responses from the 5 test sets. For analysis of the response functions, the mTRFs were averaged

345     across all the test sets. For the visual predictor, the TRFs to onsets and offsets were combined for

346     an effective response function with lags from -100 to 1000 ms relative to visual stimulus onset

347     (because the visual stimulus always lasted exactly 500 ms).

348        To estimate the neural tracking of a given predictor (or a combination of predictors), we

349     calculated the change in prediction accuracy (i.e., $\Delta z$) when the predictor(s) of interest was(ere)

350     removed from the full model that included all the visual, acoustic, and linguistic predictors. This

351     procedure tests for variability in the responses that can be *uniquely* attributed to the predictor(s)

352     under investigation and cannot be explained by any other predictors. Such a strong test is

353     warranted because different properties of natural, connected speech tend to be correlated in time.

354     Note that the analysis of the mTRFs themselves could not implement such strict control, and thus

355     we cannot exclude the possibility that response functions include components that are

356     confounded with other, correlated speech features. For this reason we focus our interpretation on

357     tests of predictive power more than mTRF comparisons.

## Statistical analysis

359     All statistical analyses, if unspecified, were implemented in the R software (version 4.2.1; Team,

360     2022).

361    First, we examined the effect of task condition (auditory single-task, or low or high load

362    dual-task) on behavioral performance, and neural visual, acoustic, and linguistic processing

363    separately. A paired T-test (two-sided), or one- or two-way repeated-measures analysis of

364    variance (ANOVA), whichever was appropriate, was performed with an alpha level of .05. The

365    reported $p$ values of those analyses were adjusted using the False Discover Rate (FDR) method

366    (Benjamini & Hochberg, 1995). We also calculated effect sizes [Cohen's $d$ for T-tests and partial

367    eta squared ($\eta^2_p$) for ANOVAs] and Bayes Factors (BF). The Bayes Factors were computed

368    using appropriate functions from the R package 'BayesFactor' (version 0.9.12.4.4; Morey et al.,

369    2022). Post hoc analysis, if necessary, was performed using paired T-tests (two-sided). FDR-

370    corrected $p$ values, effect sizes (Cohen's $d$), and Bayes Factors (BF) were reported. More

371    analysis details are provided in the following paragraphs.

372    Behavioral performance was quantified by three measures, including the accuracy and

373    reaction time (RT) for the visual task and the accuracy for the auditory task. Visual accuracy was

374    calculated as the difference in hit rates (i.e., correctly responding to a target) and false alarm

375    rates (i.e., identifying a non-target as being a target). Visual RT was calculated as the median RT

376    for hits only (Jaeggi et al., 2007; Snodgrass & Corwin, 1988). Auditory accuracy was calculated

377    as the percentage of correctly answered story questions.

378    The extent of neural visual processing was determined using a mass-univariate analysis,

379    comparing the predictive power ($z$) between the full model and a model missing the visual

380    predictor. For this, we averaged the prediction accuracy for visual predictors across task

381    conditions at individual electrodes and tested whether the averaged difference in prediction

382    accuracy ($\Delta z$) was greater than zero using a mass-univariate, one-sample T-test (one-sided). This

383    was implemented in the Eelbrain package. The mass-univariate test was a cluster-based

384    permutation test, using a *t*-value equivalent to uncorrected $p \leq 0.05$ as the cluster-forming

385    threshold. A corrected *p*-value was computed for each cluster based on the cluster-mass statistic

386    in a null distribution from 10,000 permutations (or a complete set of all possible permutations, in

387    cases where this was fewer than 10,000) (Maris & Oostenveld, 2007). We reported the largest *t*

388    value from the cluster, i.e., $t_{max}$, as an estimate of effect size (Brodbeck et al., 2018). Neural

389    acoustic and linguistic processing were analyzed in the same manner.

390        We followed each of these analyses by examining the extent to which task conditions

391    modulated neural tracking of individual predictors, or subsets of predictors. To this end we used

392    the significant cluster from the mass-univariate analysis as region of interest (ROI) to extract $\Delta z$

393    values averaged across the electrodes in the cluster, but for each condition separately. Regarding

394    neural acoustic processing, we examined the spectrogram and onset spectrogram predictors

395    separately. Regarding neural linguistic processing, we conducted three sets of analyses to

396    examine individual linguistic predictors. First, a two-way repeated-measures ANOVA was

397    performed to examine the effects of context levels (sublexical, word-form, and sentence) and

398    task condition on prediction accuracy. Second, a two-way repeated-measures ANOVA was

399    performed to examine the effects of predictor type (entropy and surprisal) and task condition on

400    prediction accuracy. Third, a one-way repeated-measures ANOVA was performed to examine

401    the effect of task condition on the prediction accuracy of word onsets. Further, if a significant

402    effect of task condition was observed from any of those analyses, we conducted follow-up

403    analyses to examine whether task conditions eliminated neural tracking of the corresponding

404    predictor(s) by testing whether the prediction accuracy at individual task conditions was greater

405    than zero using one-sample T-tests (one-sided).

406    Finally, we examined the effect of task conditions on the mTRFs from predictors which

407    showed significant task conditions effects on prediction accuracy. The predictors include visual

408    predictors, onset spectrogram, three context levels (sublexical, word-form, and sentence), and

409    two predictor types (entropy and surprisal). We calculated the global field power (GFP) of

410    mTRFs across the corresponding ROI from the above analyses of prediction accuracy. We then

411    compared the GFP of mTRFs between task conditions using mass-univariate paired T-tests (two-

412    sided). The mTRF analyses were implemented in the Eelbrain package with default parameters

413    except for the analysis time window. For visual predictors, we concatenated the mTRFs for

414    visual onsets and offsets to analyze the response to visual stimuli as a whole. For the onset

415    spectrogram, we averaged the mTRFs across the eight frequency bands. The analysis time

416    window was 0 to 1000 ms for visual predictors and 0 to 450 ms for auditory predictors.

417    # Results

418

419    Table 1 summarizes the key findings regarding the effect of task condition on behavioral

420    performance, and neural visual, acoustic, and linguistic processing.

421    **Table 1. Task Effects on Continuous Speech Processing**

| Type | Measure | | Key Result |
|---|---|---|---|
| Behavioral | Visual accuracy | | Low load > High load |
| | Visual RT | | Low load < High load |
| | Auditory Accuracy | | Auditory single-task = Low load > High load |
| | | | |
| Neural (Δz) | Visual | | Auditory single-task < Low load < High load |
| | Acoustic | Gammatone spectrogram | No significant task effect |
| | | Onset spectrogram | Auditory single-task > Low load = High load |
| | | | |

| | | | |
|---|---|---|---|
| | Linguistic | Sublexical, word-form, and sentence context | Auditory single-task > Low load = High load |
| | | Entropy and surprisal | *Entropy*: Auditory single-task > Low load = High load |
| | | | *Surprisal*: Auditory single-task > Low load = High load |
| | | Word onset | No significant task effect |

## Divided attention and visual load impair behavioral performance

Figure 3A and 3B display the accuracy and RT of the visual task for individual participants. Compared to the low load (0-back) condition, the high load (3-back) condition was associated with lower accuracy [low load: mean = 99.54% ($SD$ = 0.82) vs. high load: mean = 63.31% ($SD$ = 21.85), $t$ (15) = 6.60, $p$ < .001, Cohen's d = 1.65, BF = 2.59 × 10$^3$] and slower RT [low load: mean = 453.11 ms ($SD$ = 67.54) vs. high load: mean = 785.24 ms ($SD$ = 233.41), $t$ (15) = -5.33, $p$ < .001, Cohen's d = 1.33, BF = 330.30]. These results confirmed that the manipulation of cognitive load in the visual task was successful.


**[Fig 3 about here]**


**Figure 3.** Behavioral performance on visual and auditory tasks. (A) Accuracy on the low load (0-back) and high load (3-back) visual tasks, which was calculated as the difference in hit rates (i.e., correctly responding to a target) and false alarm rates (i.e., identifying a non-target as being a target). (B) Reaction time (RT) on the low load (0-back) and high load (3-back) visual tasks, which was calculated for hits only. (C) Accuracy on the auditory task, which was calculated as the percentage of correctly answered story questions. Individual lines in (A) to (C) denote individual participants (n = 16). (D) Correlation between the change in auditory accuracy [i.e.,

440    (low load – high load)/low load] and the change in visual RT [i.e., (high load – low load)/low

441    load]. The gray line is the linear regression line. N.s. $p > .05$, *** $p < .001$.

442

443

444         Figure 3C displays the auditory task accuracy for individual participants. The mean

445    accuracy was 88.96% ($SD = 5.93$) in the auditory single-task condition, 84.58% ($SD = 11.86$) in

446    the low load dual-task condition, and 65.63% ($SD = 12.75$) in the high load dual-task condition.

447    The effect of task condition was significant [$F(2,30) = 36.59$, $p < .001$; $\eta^2_p = .71$, BF = 6.75 ×

448    $10^6$]. Post hoc analysis revealed that auditory task accuracy was significantly lower in the high

449    load dual-task condition compared to the other two conditions: vs. auditory single-task, $t(15) =$

450    7.38, $p < .001$, Cohen's d = 1.84, BF = 8.31 × $10^3$; vs. low load dual-task, $t(15) = 6.34$, $p < .001$,

451    Cohen's d = 1.58, BF = 1.70 × $10^3$. The auditory task accuracy was not significantly different

452    between auditory single-task and low load dual-task conditions [$t(15) = 1.75$, $p = .10$, Cohen's d

453    = 0.44, BF = 0.88].

454         Further, we examined the relationship between visual and auditory task performance

455    during the dual-task conditions. The change in auditory accuracy [i.e., (low load – high load)/low

456    load] was negatively correlated with the change in visual RT [i.e., (high load – low load)/low

457    load] (Spearman's ρ = - .46, uncorrected $p = .038$, one-sided; Figure 3D), such that listeners who

458    slowed down more on the visual task from low to high load conditions tended to have a smaller

459    drop in auditory accuracy. The change in auditory accuracy was not significantly correlated with

460    the change in visual accuracy (Spearman's ρ = - .29, uncorrected $p = .28$, one-sided).

461         These results demonstrate that divided (vs. selective) attention and increasing visual load

462    impair behavioral visual and auditory performance.

## Neural tracking of visual stimuli is strongly modulated by divided attention and visual load

To assess neural tracking of visual stimuli, we focused on the predictive power of visual predictors while controlling for all speech-related predictors (acoustic and linguistic). Adding predictors for visual stimuli to a model including only speech predictors significantly improved its predictive power (prediction accuracy averaged across task conditions; $t_{max} = 12.93$, $p < .001$), providing evidence for neural tracking of visual stimuli. The cluster-based test resulted in a single significant cluster that spread across all electrodes, with the largest effects on parietal and occipital electrodes (Figure 4A).

**[Fig 4 about here]**

**Figure 4. Neural tracking of visual stimuli across task conditions**. Visual stimuli were associated with a robust response, which further increased with task-relevance and -load. (A) Topography showing the increase in prediction accuracy ($\Delta z$) due to visual predictors, which was significantly above zero in a single cluster encompassing the highlighted (yellow) electrodes. (B) Prediction accuracy across task conditions. Blue lines denote individual participants: Thicker lines indicate higher prediction accuracy for the high vs. low load condition, and thinner lines indicate lower accuracy for the high vs low load condition. Red asterisks denote $p$ values for comparison between conditions. Error bars denote the 95% within-subject confidence interval (Loftus & Masson, 1994). (C) Global field power (GFP) of the visual temporal response functions (TRFs). Visual stimuli lasted from 0 to 500 ms. Shaded areas denote within-subject standard errors around the mean (for color labels see panel B). Horizontal lines denote time windows in which TRFs were significantly different between conditions. (D)

487     Topographies of selected times in panel C (grey vertical lines). A-ST: auditory single-task, Lo-

488     DT: low load dual-task, Hi-DT: high load dual-task. ** $p < .01$, *** $p < .001$.

489

490

491         Importantly, the predictive power of the visual predictors was modulated by task

492     condition [$F(2,30) = 46.10$, $p < .001$, $\eta^2_p = .76$, BH $= 6.09 \times 10^7$]. As shown in Figure 4B, the

493     high load dual-task condition was associated with the highest predictive power (mean $= 0.075$,

494     $SD = 0.029$), followed by the low load dual-task condition (mean $= 0.053$, $SD = 0.022$), and

495     lowest for the auditory single-task condition (mean $= 0.020$, $SD = 0.012$): high load dual-task vs.

496     auditory single-task, $t(15) = 9.52$, $p < .001$, Cohen's d $= 2.38$, BF $= 1.50 \times 10^5$; high load dual-

497     task vs. low load dual-task, $t(15) = 3.34$, $p = .005$, Cohen's d $= .84$, BF $= 10.7$; low load dual-

498     task vs. auditory single-task, $t(15) = 6.64$, $p < .001$, Cohen's d $= 1.66$, BF $= 2.72 \times 10^3$. Together,

499     these results suggest that neural tracking of visual stimuli was successively enhanced with

500     increasing load of the visual task.

501         We analyzed mTRFs to further clarify how the difference in model predictive power was

502     reflected in brain responses. Visual mTRFs can be conceptualized as evoked responses to the

503     visual stimuli. Consistent with results for prediction accuracy, the mTRFs were also modulated

504     by the task condition (Figure 4C). The high load dual-task condition showed larger mTRF

505     amplitudes than the auditory single-task condition from 0 to 680 ms ($p < .001$) and the low load

506     dual-task condition from 190 to 380 ms ($p < .001$). The mTRF amplitudes for the low load dual-

507     task condition were larger than the auditory single-task condition from 70 to 210 ms ($p = .009$)

508     and from 260 to 600 ms ($p < .001$).

## Divided attention, but not visual load, reduces late neural tracking of acoustic features

The acoustic predictors significantly contributed to model prediction beyond the visual and linguistic predictors in a cluster that spread across almost all electrodes, with maxima at temporal sites ($t_{max}$ = 12.00, $p$ < .001; Figure 5A). As expected, these results provide evidence for robust neural tracking of acoustic properties of speech.


**[Fig 5 about here]**


**Figure 5. Neural tracking of acoustic information across task conditions**. (A) Increase in prediction accuracy (Δz) due to acoustic predictors of speech (gammatone and onset spectrogram), which was significantly above zero in a cluster encompassing the highlighted (yellow) electrodes. Blue dots denote individual participants. (B) Prediction accuracy across task conditions for acoustic predictors, i.e., combined gammatone and onset spectrogram. (C) and (D) Prediction accuracy across task conditions for gammatone spectrogram and onset spectrogram separately. Topographies highlight electrodes (yellow) with prediction accuracy that was significantly above zero. Black asterisks denote $p$ values for testing against (above) zero at individual conditions. (B) to (D) Blue lines denote individual participants: Thicker lines indicate lower accuracy for high vs. low load condition, and thinner lines indicate higher accuracy for high vs. low load condition. Red asterisks denote $p$ values for comparison between conditions. Error bars denote 95% confidence interval. (E) and (F) Global field power (GFP; top) of mTRFs and related topographies (bottom) for gammatone and onset spectrogram. The mTRFs were averaged across the eight frequency bands. Shaded areas denote within-subject standard errors around the mean. Horizontal lines denote time windows in which mTRFs were significantly

533 different between conditions. Topographies are shown for selected times indicated in GFPs (grey

534 vertical lines). A-ST: auditory single-task, Lo-DT: low load dual-task, Hi-DT: high load dual-

535 task. * $p < .05$, ** $p < .01$, *** $p < .001$.

536

537

538

539        The prediction accuracy for acoustic predictors was modulated by task condition [$F(2,30)$

540 $= 14.83$, $p < .001$, $\eta^2_p = .50$, BF $= 581.38$; Figure 5B]. Post hoc analysis showed that the

541 prediction accuracy significantly dropped in the two dual-task conditions compared to the

542 auditory single-task condition [vs. low load dual-task, $t(15) = 3.84$, $p = .002$, Cohen's d = 0.96,

543 BF = 25.60; vs. high load dual-task, $t(15) = 4.78$, $p < .001$, Cohen's d = 1.20, BF = 130.60]. The

544 prediction accuracy was not significantly different between the dual-task conditions [$t(15) = .77$,

545 $p = .45$, Cohen's d = 0.19, BF = 0.33]. These results suggest that neural tracking of acoustic

546 information was reduced when directing attention from one task (auditory) to two tasks (visual

547 and auditory).

548        Then, we assessed whether the effect of task condition could be attributed to specific

549 acoustic predictors. The two acoustic predictors both independently contributed to overall model

550 prediction (gammatone spectrogram: $t_{max} = 6.08$, $p < .001$, Figure 5C; onset spectrogram: $t_{max} =$

551 9.91, $p < .001$, Figure 5D). The effect of task condition on the prediction accuracy was

552 significant for onset spectrogram [$F(2,30) = 4.93$, $p = .033$, $\eta^2_p = 0.25$, BF = 4.17] but not for

553 gammatone spectrogram [$F(2,30) = .70$, $p = .59$, $\eta^2_p = 0.04$, BF = 0.26]. Post hoc analysis

554 revealed that the prediction accuracy for onset spectrogram significantly dropped in the two

555 dual-task conditions compared to the auditory single-task condition [vs. low load dual-task, $t(15)$

556   = 2.61, $p$ = .030, Cohen's d = 0.65, BF = 3.14; vs. high load dual-task, $t(15)$ = 2.89, $p$ = .030,

557   Cohen's d = 0.72, BF = 4.94]. The prediction accuracy was not significantly different between

558   the dual-task conditions [$t(15)$ = - .79, $p$ = .44, Cohen's d = 0.20, BF = 0.34].

559       Considering the modulation by task condition, we further examined whether divided

560   attention eliminated neural tracking of onset spectrogram. The prediction accuracy at individual

561   task conditions was significantly above zero (all uncorrected $p$s < .001, Cohen's d > 1.20, BF >

562   256.40; Figure 5D), suggesting that directing attention from one task to two tasks did not

563   eliminate the neural tracking of acoustic onsets.

564       Finally, we examined the effect of task condition on the mTRFs for the onset

565   spectrogram (Figure 5F). mTRFs to a continuous stimulus like the auditory spectrogram can be

566   conceived of as evoked responses to an elementary event in the stimulus (i.e., the impulse

567   response). The mTRF amplitudes in the auditory single-task condition were larger compared to

568   the high load dual-task condition from 200 to 260 ms ($p$ = .003). Further, a visual inspection of

569   the mTRFs from individual subjects revealed two relatively reliable peaks at about 56 (P1) and

570   152 (P2) ms. Latencies of these peaks were not significantly different across conditions [56 ms:

571   $F(2,30)$ = .65, uncorrected $p$ = .94, $\eta^2_p$ = 0.004; 152 ms: $F(2,30)$ = .62, uncorrected $p$ = .54, $\eta^2_p$ =

572   0.04].

573       In sum, acoustic tracking was very similar across conditions, with only a slight reduction

574   in the tracking of acoustic onsets in the divided attention tasks, compared to the single task. This

575   difference was likely explained by a reduction in a relatively late response component, starting at

576   200 ms.

## Divided attention, but not visual load, reduces late tracking of linguistic information

The linguistic predictors significantly contributed to model prediction beyond the visual and acoustic predictors ($t_{max} = 4.95$, $p < .001$; Figure 6A). The cluster-based test indicated that the effect of linguistic predictors was primarily observed for temporal-central electrodes. These results provide evidence for neural tracking of linguistic properties of speech.


**[Fig 6 about here]**


**Figure 6. Neural tracking of linguistic information across task conditions**. (A) Increase in prediction accuracy (Δz) due to linguistic predictors of speech (word onsets, phoneme onsets, sublexical surprisal and entropy, word-form surprisal and entropy, and sentence surprisal and entropy), which was significantly above zero across highlighted (yellow) electrodes in the topography. Blue dots denote individual participants. (B) Prediction accuracy for combined linguistic predictors across conditions. Blue lines denote individual participants: Thicker lines indicate lower accuracy for high vs. low load condition, and thinner lines indicate higher accuracy for high vs. low load condition. Red asterisks denote *p* values for comparison between conditions. (C) Prediction accuracy for three context levels (sublexical, word-form, and sentence) across conditions. Each level includes entropy and surprisal predictors. (D) Global field power (GFP; top) of mTRFs and related topographies (bottom) for each context level. The mTRF GFPs were averaged across entropy and surprisal. (E) Prediction accuracy for entropy and surprisal. Each predictor includes the three context levels. (F) GFP of mTRFs and related topographies for entropy and surprisal. (B), (C), and (E) Error bars denote 95% confidence interval. (D) and (F) Shaded areas denote standard errors around the mean. Horizontal lines

601    denote time windows in which the mTRFs were significantly different between task conditions.

602    Topographies are shown for selected times indicated in GFPs (grey vertical lines). A-ST:

603    auditory single-task, Lo-DT: low load dual-task, Hi-DT: high load dual-task. * $p < .05$, *** $p <$

604    .001.

605

606

607

608         The prediction accuracy for linguistic predictors was modulated by task condition

609    $[F(1.41, 21.15) = 6.66, p = .029, \eta^2_p = 0.31, BF = 10.82;$ Figure 6B]. The prediction accuracy

610    significantly dropped in the two dual-task conditions compared to the auditory single-task

611    condition [vs. low load dual-task, $t(15) = 2.83, p = .029$, Cohen's d = 0.71, BF = 4.49; vs. high

612    load dual-task, $t(15) = 2.61, p = .029$, Cohen's d = 0.65, BF = 3.16]. The prediction accuracy was

613    not significantly different between the two dual-task conditions [$t(15) = -1.80, p = .091$, Cohen's

614    d = 0.45, BF = 0.95]. These results suggest that neural tracking of linguistic information was

615    reduced when directing attention from one task to two tasks.

616         Next, we conducted three sets of analyses to assess whether the effect of task condition

617    could be attributed to specific linguistic properties.

618    **Task effects appear to be similar across different context levels**

619    The first analysis focused on the three context levels (sublexical, word-form, and sentence). Each

620    level independently contributed significantly to model prediction (sublexical: $t_{max} = 5.22, p <$

621    .001; word-form: $t_{max} = 3.92, p < .001$; sentence: $t_{max} = 4.98, p < .001$; Figure 6C). A two-way

622    repeated-measures ANOVA showed that the interaction between context level and task condition

623    was not significant [$F(2.55, 38.18) = 1.19, p = .40, \eta^2_p = 0.073, BF = 0.19$]. The main effect of

624    context level was not significant [$F(2,30) = .32$, $p = .77$, $\eta^2_p = 0.021$, BF = 0.08]. But the main

625    effect of task condition was significant [$F(1.2,18.01) = 8.46$, $p = .021$, $\eta^2_p = 0.36$, BF = $1.40 \times$

626    $10^3$]. Post hoc analysis showed that the prediction accuracy was significantly reduced from the

627    auditory single-task condition to the low load [$t(15) = 2.90$, $p = .016$, Cohen's d = 0.73, BF =

628    5.08] and high load dual-task conditions [$t(15) = 4.27$, $p = .002$, Cohen's d = 1.07, BF = 54.63].

629    But the prediction accuracy was not significantly different between the low and high load dual-

630    task condition[$t(15) = 1.02$, $p = .32$, Cohen's d = 0.26, BF = 0.40]. Further, we found similar

631    patterns of results when restricting the two-way repeated-measures ANOVA analysis to the dual-

632    task conditions. In sum, patterns of task condition effects observed for linguistic predictors

633    appeared to be similar across the different linguistic models.

634         Considering the modulation by context level and task condition, we further examined

635    whether divided attention eliminated neural tracking of any of these predictors. The prediction

636    accuracies for all predictors at individual task conditions were significantly above zero (all

637    uncorrected $p$s < .03, Cohen's d > 0.53, BF > 1.51).

638         Regarding mTRFs, the effect of task condition was not significant for sublexical or word-

639    form context but was for sentence context (Figure 6D). The mTRF amplitude of sentence context

640    in the auditory single-task condition was larger compared to the low load dual-task condition

641    from 400 to 430 ms ($p = .036$). Topographies suggest that this is due to a broadly distributed

642    more negative component in the single task condition.

643         Initial response peaks to linguistic features appear relatively early. This is consistent with

644    previous results (Brodbeck et al., 2022) and might be partly because forced alignment, which

645    was used to determine phoneme timing, does not take into account coarticulation effects. Some

646　　information about upcoming phonetic features might thus have systematically precede the

647　　estimates of phoneme onset times we used.

## Neural tracking of surprisal might increase with visual load

649　　　　The second analysis focused on entropy and surprisal. The two predictors independently

650　　contributed significantly to model prediction (entropy: $t_{max}$ = 5.51, $p$ < .001; surprisal: $t_{max}$ =

651　　3.91, $p$ = .001; Figure 6E). A two-way repeated-measures ANOVA showed that the interaction

652　　between predictor type (entropy vs. surprisal) and task condition was not significant [$F$(1.32,

653　　19.86) = 1.29, $p$ = .40, $\eta^2_p$ = 0.079, BF = 0.31]. The main effect of predictor type was not

654　　significant [$F$(1,15) = .31, $p$ = .65, $\eta^2_p$ = 0.02, BF = 0.23]. But the main effect of task condition

655　　was significant [$F$(1.35,20.2) = 9.85, $p$ = .011, $\eta^2_p$ = 0.40, BF = 890.10]. Post hoc analysis

656　　showed that, when averaging across surprisal and entropy, the prediction accuracy was

657　　significantly reduced from the auditory single-task condition to the low load [$t$(15) = 3.28, $p$ =

658　　.008, Cohen's d = 0.82, BF = 9.56] and high load dual-task conditions [$t$(15) = 4.11, $p$ = .003,

659　　Cohen's d = 1.03, BF = 40.91]. Numerically, the prediction accuracy was improved from the low

660　　load to high load dual-task condition, but this difference was not significant [$t$(15) = 1.27, $p$ =

661　　.22, Cohen's d = 0.32, BF = 0.51].

662　　　　Because of theoretical predictions of enhanced reliance on linguistic representations

663　　during higher visual task load (see *Introduction* and *Discussion*), we further restricted the two-

664　　way repeated-measures ANOVA analysis to the dual-task conditions. The interaction between

665　　predictor type and task condition was significant [$F$(1,15) = 5.75, uncorrected $p$ = .03 (FDR-

666　　corrected $p$ = .063), $\eta^2_p$ = 0.28, BF = 1.23]. There was no significant main effect of predictor type

667　　[$F$(1,15) = 1.92, $p$ = .33, $\eta^2_p$ = 0.11, BF = 0.40] or task condition [$F$(1,15) = 1.62, $p$ = .36, $\eta^2_p$ =

668　　0.097, BF = 0.79]. Post hoc analysis showed that for entropy, the prediction accuracy was not

669     different between the dual-task conditions [$t(15) = .10$, $p = .92$, Cohen's d = 0.03, BF = 0.26].

670     But for surprisal, the prediction accuracy was significantly improved from the low load to high

671     load dual-task condition [$t(15) = 2.20$, uncorrected $p = .044$, Cohen's d = 0.55, BH = 1.66].

672          Considering the modulation by predictor type and task condition, we further examined

673     whether divided attention eliminated neural tracking of entropy or surprisal. The prediction

674     accuracies for both predictors at individual task conditions were significantly above zero (all

675     uncorrected $p$s < .01, Cohen's d > 0.66, BF > 3.36), except for the surprisal predictors at the low

676     load dual-task condition (uncorrected $p = .059$, Cohen's d = 0.41, BF = 0.78).

677          Regarding mTRFs, the mTRF amplitude of entropy in the low load dual-task condition

678     was smaller than the high load dual-task condition from 160 to 200 ms (uncorrected $p = .037$).

679     The mTRF amplitude of surprisal in the low load dual-task condition was smaller compared to

680     the auditory single-task condition from 380 to 430 ms (uncorrected $p = .012$). We did not

681     observe a significant effect of task load, although the mTRF to surprisal during high visual load

682     was numerically stronger than low load from 200 ms onwards.

683     **Divided attention or visual load does not affect neural tracking of word onsets**

684     The third set of analysis focused on word onset. This predictor independently contributed

685     significantly to model prediction ($t_{max} = 4.51$, $p < .001$; Figure 7A). But the effect of task

686     condition on prediction accuracy was not significant [$F(2,30) = .07$, $p = .93$, $\eta^2_p = 0.005$, BF =

687     0.17; Figure 7B].

688

689                              **[Fig 7 about here]**

690

691 **Figure 7. Neural tracking of word onsets across task conditions**. (A) Topography showing

692 the increase in prediction accuracy (Δz) due to word onsets, which was significantly above zero

693 across highlighted (yellow) electrodes. (B) Prediction accuracy across conditions. Blue lines

694 denote individual participants: Thicker lines indicate higher accuracy for the high vs. low load

695 condition, and thinner lines indicate lower accuracy for the high vs low load condition. Error

696 bars denote 95% within-subject confidence interval (Loftus & Masson, 1994). (C) Global field

697 power (GFP) of mTRFs. Shaded areas denote within-subject standard error around the mean. (D)

698 Topographies of selected times in panel C (grey vertical lines). * $p < .05$, ** $p < .01$.

699

700 Taken together, results suggest that directing attention from one task to two tasks may

701 reduce but does not eliminate the neural tracking of linguistic features of speech. However,

702 increasing visual load does not lead to a further reduction. On the contrary, an increasing load of

703 the dual task might even be associated with enhanced neural tracking of phoneme surprisal.

704 However, this effect should be interpreted with care because the effect was not significant when

705 analyzing all linguistic predictors as a group or after correction for multiple comparisons.


# Discussion

706

707 We examined the extent to which bimodal divided attention influences acoustic and linguistic

708 representations of natural continuous speech. Compared to unimodal auditory speech processing,

709 the visual tasks affected acoustic onsets (but not the acoustic spectrogram, Figure 5D) and

710 linguistic representations related to predictive processing (but not to lexical segmentation, Figure

711 6E). Surprisingly, we did not find evidence of further reduction (at any processing level) with

712 increased visual (dual) task load, despite unambiguous behavioral and neural evidence of the

713 high load task as being more demanding (Figures 3 and 4).

## Locus of effects of bimodal divided attention on continuous speech processing

We noted a striking dissociation between the impact of the dual-task on behavioral performance in the speech comprehension task, and a relative lack of impact on neural speech processing. Behaviorally, the dual-task load clearly impacted listeners' ability to answer auditory comprehension questions. However, neural tracking of acoustic and linguistic speech features was affected only at late response components, and remained largely unchanged with varying dual-task load. This neural and behavioral dissociation suggests that bimodal divided attention largely only impacts late, post-perceptual processes of continuous speech processing. The significant and unchanged responses related to predictive processing using the sentence context suggest that listeners could track multi-word sequences regardless of dual-task load. We posit that the decreased behavioral performance originates from higher-order cognitive processes that are not adequately described by probabilistic word-sequence models, such as semantic integration and memory formation.

Previous behavioral research has suggested that increased dual-task load is associated with reduced acoustic sensitivity during speech recognition (Mattys et al., 2014; Mattys & Palmer, 2015; Mattys & Wiget, 2011). In our data, the dual-task did not alter early acoustic responses and only had subtle effects on later (> 200 ms) responses (see Figures 5E and 5F). This suggests that the effect of bimodal divided attention may not be on basic acoustic representations per se, but on secondary acoustic analysis stages or on how these representations are accessed and used by higher-order processes.

## Implications for resource-based accounts

A common framework for understanding effects under dual-task paradigms is resource-based (Wahn & König, 2017). When two tasks draw from a limited pool of shared resources, increased load in one task is associated with poorer performance in another task. Such a hypothesis is often referred to as the *supramodal* account of crossmodal attention (Broadbent, 1958; Ciaramitaro et al., 2017; Klemen et al., 2009; Macdonald & Lavie, 2011; Molloy et al., 2015). In contrast, if the increased load in one task does not affect a corresponding decrease in another, the two modalities can draw on separate resource pools. Such a hypothesis is consistent with a *modality-specific* account of crossmodal attention (Alais et al., 2006; Arrighi et al., 2011; Duncan et al., 1997; Keitel et al., 2013; Parks et al., 2011; Porcu et al., 2014).

Here, we observed a reduction in neural tracking of speech acoustic and linguistic features under bimodal divided attention, consistent with previous studies demonstrating detrimental effects of bimodal divided attention for simplified speech stimuli such as syllables (Gennari et al., 2018), words (Kasper et al., 2014), and sentences (Salo et al., 2015). In conjunction with the co-occurring improved neural tracking of visual stimuli, this finding appears to suggest a tradeoff between attending to the auditory versus visual modalities. Hence, these results appear to align with the *supramodal* hypothesis of the dual-task effects that the auditory and visual tasks of our study draw on a limited pool of shared resources (Broadbent, 1958; Ciaramitaro et al., 2017; Klemen et al., 2009; Macdonald & Lavie, 2011; Molloy et al., 2015).

However, a *supramodal* hypothesis of the dual-task effects does not seem to fit other key results from our study. First, the impact of bimodal divided attention is specific to certain features of the speech signals: we found bimodal attention effects for acoustic onsets but not acoustic envelope representations, and for predictive linguistic processing, indexed through

759 information-theoretic variables, but not lexical segmentation, indexed through the word-onset

760 predictors. In each case, the impact of divided attention is not a generally reduced representation

761 but is restricted to only specific response components in the response time courses (the mTRFs).

762       Furthermore, a resource-based account would suggest that when the visual load is further

763 increased, available resources for speech representations should further decrease, which is not

764 what we observed. Instead, adding a visual task exacted a cost on neural speech representations,

765 but this cost did not scale with the task load. In contrast to these neural effects, task load did

766 affect behavioral performance on the auditory task. These divergent results may require an

767 explanation involving different resource pools (Wahn & König, 2017). For example, there may

768 be a resource pool for sensory processing, which is sensitive to divided attention but not task

769 load, thus, is relatively modality-specific. There may be a second resource pool, which is

770 sensitive to task load and affects higher-order story comprehension, thus, is relatively

771 supramodal.

## Selective *versus* divided attention on speech processing

773 Previous studies on continuous speech processing have shown that selective attention within and

774 between modalities strongly modulates neural processing of both acoustic and linguistic features

775 of continuous speech, and the attentional effects seem to be even stronger for linguistic

776 processing (Brodbeck et al., 2018; Broderick et al., 2018; Ding et al., 2018; Kiremitçi et al.,

777 2021; Vanthornhout et al., 2019; Yahav & Golumbic, 2021). Specifically, neural tracking of

778 acoustic features is reduced and delayed but not eliminated for unattended speech, but the

779 tracking of linguistic features is virtually abolished. A parsimonious null hypothesis, consistent

780 with the notion of a shared resource pool, is that speech representations during divided attention

781    ought to be halfway between attended and ignored speech. Our results suggest that this is not the

782    case.

783         First, certain speech features (e.g., acoustic spectrogram and word onsets) that have been

784    shown in prior work to be modulated by selective attention are insensitive to bimodal divided

785    attention. Second, unlike prior work demonstrating differential selective attention effects on the

786    relative balance of acoustic vs. linguistic processing, we did not observe a greater reduction in

787    linguistic processing than acoustic processing with the manipulation of divided attention. The

788    neural tracking of both feature classes is reduced but not eliminated. Third, for those features

789    showing modulation by divided attention, we did not observe any delay in the neural responses

790    as reflected in the mTRFs (Figures 5 and 6). Fourth, the effect of divided attention emerged

791    largely at later stages (after ~200 ms) with the earlier latencies relatively unaffected. Thus, the

792    effect of bimodal divided attention on neural continuous speech processing appears to be feature-

793    specific and occurs relatively late in processing.

794         These differences indicate that selective and divided attention are subserved by distinct

795    mechanisms. Relative to selective attention, bimodal divided attention tasks may be associated

796    with additional recruitment of frontal regions that interact with sensory cortices (Gennari et al.,

797    2018; Johnson & Zatorre, 2006; Loose et al., 2003). A stronger engagement of frontal regions

798    has been associated with poorer task performance (Gennari et al., 2018; Johnson & Zatorre,

799    2006). These neural findings appear to align with the argument that the costs of bimodal divided

800    attention may come from limitations of executive control to coordinate processes related to two

801    tasks rather than a competition for shared sensory resources (Katus & Eimer, 2019; Loose et al.,

802    2003). The differential effects of selective and divided attention on continuous speech processing

803    suggest that the costs of selective attention are more likely to originate from 'filter' mechanisms

804    (Broadbent, 1958; Lachter et al., 2004) that pass task-relevant signals but block task-irrelevant

805    others, instead of the re-allocation of shared resources. Nevertheless, future studies are needed to

806    elucidate mechanisms underlying differences in continuous speech processing between selective

807    and divided attention.

## Increased responses to surprisal with dual-task load

809    We found that increasing visual load increased responses to phoneme surprisal, but not entropy.

810    This effect was statistically only seen after excluding the auditory single-task condition and

811    should thus be interpreted with care, but it is consistent with several extant findings. The

812    dissociation between entropy and surprisal is consistent with recent evidence that these two

813    processes may reflect different neural processes (Gaston et al., 2022). Neural responses

814    associated with surprisal may reflect prediction errors that signal the difference between

815    predicted and observed phonemes. Such prediction error signals may be boosted when attention

816    is directed to the speech stimuli (e.g., auditory single-task; Auksztulewicz & Friston, 2015;

817    Smout et al., 2019) or when attention to the speech stimuli is directed away to demanding

818    crossmodal tasks (e.g., high-load visual tasks; Xie et al., 2018). The increased response to

819    surprisal might also reflect a shift toward more reliance on linguistic representations during

820    speech processing when resources for auditory processing were constrained under divided

821    attention of higher load (Mattys et al., 2009; Mattys & Wiget, 2011).

## Neural tracking of word onsets was not affected by divided attention

823    Tracking of word onsets might reflect lexical segmentation (Sanders et al., 2002; Sanders &

824    Neville, 2003) and, along with other linguistic features, is strongly affected by selective attention

825    (Brodbeck et al., 2018). It has been suggested that neural responses to word onsets reflect the

826  dynamic allocation of attention to time windows that contain word onsets (Astheimer & Sanders,

827  2009). However, our results indicate that tracking of word onsets is robust to manipulations of

828  attentional load by adding a visual task and increasing dual-task load. This suggests that the

829  word-onset attention effect may draw on a relatively unshared resource pool, or that the word-

830  onset responses reflect a more mechanistic aspect of lexical segmentation.

## Conclusion

832  This study demonstrates a striking dissociation between the impact of dual-task load on

833  behavioral speech comprehension performance and a relative lack of impact on time-locked

834  neural representations of continuous speech. The behavioral effects of bimodal divided attention

835  on continuous speech processing occur not due to impaired early sensory representations but

836  likely at later cognitive processing stages.

837

## References

839  Alais, D., Morrone, C., & Burr, D. (2006). Separate attentional resources for vision and audition.

840      *Proceedings of the Royal Society B: Biological Sciences*, *273*(1592), 1339–1345.

841  Arrighi, R., Lunardi, R., & Burr, D. (2011). Vision and audition do not share attentional

842      resources in sustained tasks. *Frontiers in Psychology*, *2*, 56.

843  Astheimer, L. B., & Sanders, L. D. (2009). Listeners modulate temporally selective attention

844      during natural speech processing. *Biological Psychology*, *80*(1), 23–34.

845  Auksztulewicz, R., & Friston, K. (2015). Attentional enhancement of auditory mismatch

846      responses: A DCM/MEG study. *Cerebral Cortex*, *25*(11), 4273–4283.

847  Benjamini, Y., & Hochberg, Y. (1995). Controlling the false discovery rate: A practical and

848      powerful approach to multiple testing. *Journal of the Royal Statistical Society: Series B*

849      *(Methodological)*, *57*(1), 289–300.

850  Bidelman, G. M., & Alain, C. (2015). Musical training orchestrates coordinated neuroplasticity

851      in auditory brainstem and cortex to counteract age-related declines in categorical vowel

852      perception. *Journal of Neuroscience*, *35*(3), 1240–1249.

853  Broadbent, D. E. (1958). *Perception and communication*. Pergamon Press.

854  Brodbeck, C., Bhattasali, S., Heredia, A. A. C., Resnik, P., Simon, J. Z., & Lau, E. (2022).

855      Parallel processing in speech perception with local and global representations of

856      linguistic context. *Elife*, *11*, e72056.

857  Brodbeck, C., Das, P., Kulasingham, J. P., Bhattasali, S., Gaston, P., Resnik, P., & Simon, J. Z.

858      (2021). Eelbrain: A Python toolkit for time-continuous analysis with temporal response

859      functions. *BioRxiv*.

860  Brodbeck, C., Hong, L. E., & Simon, J. Z. (2018). Rapid transformation from auditory to

861      linguistic representations of continuous speech. *Current Biology*, *28*(24), 3976–3983.

862  Brodbeck, C., Jiao, A., Hong, L. E., & Simon, J. Z. (2020). Neural speech restoration at the

863      cocktail party: Auditory cortex recovers masked speech of both attended and ignored

864      speakers. *PLoS Biology*, *18*(10), e3000883.

865  Brodbeck, C., & Simon, J. Z. (2020). Continuous speech processing. *Current Opinion in*

866      *Physiology*, *18*, 25–31.

867  Broderick, M. P., Anderson, A. J., Di Liberto, G. M., Crosse, M. J., & Lalor, E. C. (2018).

868      Electrophysiological correlates of semantic dissimilarity reflect the comprehension of

869      natural, narrative speech. *Current Biology*, *28*(5), 803–809.

870    Ciaramitaro, V. M., Chow, H. M., & Eglington, L. G. (2017). Cross-modal attention influences

871        auditory contrast sensitivity: Decreasing visual load improves auditory thresholds for

872        amplitude-and frequency-modulated sounds. *Journal of Vision*, *17*(3), 20–20.

873    Crosse, M. J., Di Liberto, G. M., Bednar, A., & Lalor, E. C. (2016). The multivariate temporal

874        response function (mTRF) toolbox: A MATLAB toolbox for relating neural signals to

875        continuous stimuli. *Frontiers in Human Neuroscience*, *10*, 604.

876    Ding, N., Pan, X., Luo, C., Su, N., Zhang, W., & Zhang, J. (2018). Attention is required for

877        knowledge-based sequential grouping: Insights from the integration of syllables into

878        words. *Journal of Neuroscience*, *38*(5), 1178–1188.

879    Ding, N., & Simon, J. Z. (2012). Neural coding of continuous speech in auditory cortex during

880        monaural and dichotic listening. *Journal of Neurophysiology*, *107*(1), 78–89.

881    Duncan, J., Martens, S., & Ward, R. (1997). Restricted attentional capacity within but not

882        between sensory modalities. *Nature*, *387*(6635), 808–810.

883    Fishbach, A., Nelken, I., & Yeshurun, Y. (2001). Auditory edge detection: A neural model for

884        physiological and psychoacoustical responses to amplitude transients. *Journal of*

885        *Neurophysiology*, *85*(6), 2303–2323.

886    Gaston, P., Brodbeck, C., Phillips, C., & Lau, E. (2022). Auditory Word Comprehension is Less

887        Incremental in Isolated Words. *Neurobiology of Language*, 1–50.

888        https://doi.org/10.1162/nol_a_00084

889    Gennari, S. P., Millman, R. E., Hymers, M., & Mattys, S. L. (2018). Anterior paracingulate and

890        cingulate cortex mediates the effects of cognitive load on speech sound discrimination.

891        *NeuroImage*, *178*, 735–743.

892   Gillis, M., Van Canneyt, J., Francart, T., & Vanthornhout, J. (2022). Neural tracking as a

893        diagnostic tool to assess the auditory pathway. *Hearing Research*, 108607.

894   Gramfort, A., Luessi, M., Larson, E., Engemann, D. A., Strohmeier, D., Brodbeck, C., Goj, R.,

895        Jas, M., Brooks, T., & Parkkonen, L. (2013). MEG and EEG data analysis with MNE-

896        Python. *Frontiers in Neuroscience*, 267.

897   Hamilton, L. S., & Huth, A. G. (2020). The revolution will not be controlled: Natural stimuli in

898        speech neuroscience. *Language, Cognition and Neuroscience*, *35*(5), 573–582.

899   Heafield, K. (2011). *KenLM: Faster and smaller language model queries*. 187–197.

900   Hickok, G., & Poeppel, D. (2007). The cortical organization of speech processing. *Nature*

901        *Reviews Neuroscience*, *8*(5), 393–402.

902   Jaeggi, S. M., Buschkuehl, M., Etienne, A., Ozdoba, C., Perrig, W. J., & Nirkko, A. C. (2007).

903        On how high performers keep cool brains in situations of cognitive overload. *Cognitive,*

904        *Affective, & Behavioral Neuroscience*, *7*(2), 75–89.

905   Johnson, J. A., & Zatorre, R. J. (2006). Neural substrates for dividing and focusing attention

906        between simultaneous auditory and visual events. *Neuroimage*, *31*(4), 1673–1681.

907   Kasper, R. W., Cecotti, H., Touryan, J., Eckstein, M. P., & Giesbrecht, B. (2014). Isolating the

908        neural mechanisms of interference during continuous multisensory dual-task

909        performance. *Journal of Cognitive Neuroscience*, *26*(3), 476–489.

910   Katus, T., & Eimer, M. (2019). The sources of dual-task costs in multisensory working memory

911        tasks. *Journal of Cognitive Neuroscience*, *31*(2), 175–185.

912   Keitel, C., Maess, B., Schröger, E., & Müller, M. M. (2013). Early visual and auditory

913        processing rely on modality-specific attentional resources. *Neuroimage*, *70*, 240–249.

914    Keuleers, E., Brysbaert, M., & New, B. (2010). SUBTLEX-NL: A new measure for Dutch word

915        frequency based on film subtitles. *Behavior Research Methods*, *42*(3), 643–650.

916    Kiremitçi, I., Yilmaz, Ö., Çelik, E., Shahdloo, M., Huth, A. G., & Çukur, T. (2021). Attentional

917        modulation of hierarchical speech representations in a multitalker environment. *Cerebral*

918        *Cortex*, *31*(11), 4986–5005.

919    Klemen, J., Büchel, C., & Rose, M. (2009). Perceptual load interacts with stimulus processing

920        across sensory modalities. *European Journal of Neuroscience*, *29*(12), 2426–2434.

921    Lachter, J., Forster, K. I., & Ruthruff, E. (2004). Forty-five years after Broadbent (1958): Still no

922        identification without attention. *Psychological Review*, *111*(4), 880.

923    Loftus, G. R., & Masson, M. E. (1994). Using confidence intervals in within-subject designs.

924        *Psychonomic Bulletin & Review*, *1*(4), 476–490.

925    Loose, R., Kaufmann, C., Auer, D. P., & Lange, K. W. (2003). Human prefrontal and sensory

926        cortical activity during divided attention tasks. *Human Brain Mapping*, *18*(4), 249–259.

927    Macdonald, J. S., & Lavie, N. (2011). Visual perceptual load induces inattentional deafness.

928        *Attention, Perception, & Psychophysics*, *73*(6), 1780–1789.

929    Maris, E., & Oostenveld, R. (2007). Nonparametric statistical testing of EEG-and MEG-data.

930        *Journal of Neuroscience Methods*, *164*(1), 177–190.

931    Marslen-Wilson, W. D. (1987). Functional parallelism in spoken word-recognition. *Cognition*,

932        *25*(1–2), 71–102.

933    Mattys, S. L., Barden, K., & Samuel, A. G. (2014). Extrinsic cognitive load impairs low-level

934        speech perception. *Psychonomic Bulletin & Review*, *21*(3), 748–754.

935    Mattys, S. L., Brooks, J., & Cooke, M. (2009). Recognizing speech under a processing load:

936        Dissociating energetic from informational factors. *Cognitive Psychology*, *59*(3), 203–243.

937    Mattys, S. L., & Palmer, S. D. (2015). Divided attention disrupts perceptual encoding during

938         speech recognition. *The Journal of the Acoustical Society of America*, *137*(3), 1464–

939         1472.

940    Mattys, S. L., & Wiget, L. (2011). Effects of cognitive load on speech recognition. *Journal of*

941         *Memory and Language*, *65*(2), 145–160.

942    McAuliffe, M., Socolof, M., Mihuc, S., Wagner, M., & Sonderegger, M. (2017). *Montreal*

943         *Forced Aligner: Trainable Text-Speech Alignment Using Kaldi. 2017*, 498–502.

944    Molloy, K., Griffiths, T. D., Chait, M., & Lavie, N. (2015). Inattentional deafness: Visual load

945         leads to time-specific suppression of auditory evoked responses. *Journal of*

946         *Neuroscience*, *35*(49), 16046–16054.

947    Morey, R. D., Rouder, J. N., Jamil, T., & Morey, M. R. D. (2022). *BayesFactor: Computation of*

948         *Bayes Factors for Common Designs*. https://CRAN.R-project.org/package=BayesFactor

949    Oostenveld, R., & Praamstra, P. (2001). The five percent electrode system for high-resolution

950         EEG and ERP measurements. *Clinical Neurophysiology*, *112*(4), 713–719.

951    Parks, N. A., Hilimire, M. R., & Corballis, P. M. (2011). Steady-state signatures of visual

952         perceptual load, multimodal distractor filtering, and neural competition. *Journal of*

953         *Cognitive Neuroscience*, *23*(5), 1113–1124.

954    Pickering, M. J., & Gambi, C. (2018). Predicting while comprehending language: A theory and

955         review. *Psychological Bulletin*, *144*(10), 1002.

956    Porcu, E., Keitel, C., & Müller, M. M. (2014). Visual, auditory and tactile stimuli compete for

957         early sensory processing capacities within but not between senses. *Neuroimage*, *97*, 224–

958         235.

959  Salo, E., Rinne, T., Salonen, O., & Alho, K. (2015). Brain activations during bimodal dual tasks

960      depend on the nature and combination of component tasks. *Frontiers in Human*

961      *Neuroscience*, *9*, 102.

962  Sanders, L. D., & Neville, H. J. (2003). An ERP study of continuous speech processing: I.

963      Segmentation, semantics, and syntax in native speakers. *Cognitive Brain Research*, *15*(3),

964      228–240.

965  Sanders, L. D., Newport, E. L., & Neville, H. J. (2002). Segmenting nonsense: An event-related

966      potential index of perceived onsets in continuous speech. *Nature Neuroscience*, *5*(7),

967      700–703.

968  Schneider, W., Eschman, A., & Zuccolotto, A. (2002). *E-Prime: User's guide. Reference guide.*

969      *Getting started guide*. Psychology Software Tools, Incorporated.

970  Smout, C. A., Tang, M. F., Garrido, M. I., & Mattingley, J. B. (2019). Attention promotes the

971      neural encoding of prediction errors. *PLoS Biology*, *17*(2), e2006812.

972  Snodgrass, J. G., & Corwin, J. (1988). Pragmatics of measuring recognition memory:

973      Applications to dementia and amnesia. *Journal of Experimental Psychology: General*,

974      *117*(1), 34.

975  Team, R. C. (2022). *R: A language and environment for statistical computing*. https://www.R-

976      project.org/.

977  Vanthornhout, J., Decruy, L., & Francart, T. (2019). Effect of task and attention on neural

978      tracking of speech. *Frontiers in Neuroscience*, *13*, 977.

979  Wahn, B., & König, P. (2017). Is attentional resource allocation across sensory modalities task-

980      dependent? *Advances in Cognitive Psychology*.

981    Xie, Z., Reetzke, R., & Chandrasekaran, B. (2018). Taking attention away from the auditory

982         modality: Context-dependent effects on early sensory encoding of speech. *Neuroscience*,

983         *384*, 64–75.

984    Yahav, P. H., & Golumbic, E. Z. (2021). Linguistic processing of task-irrelevant speech at a

985         cocktail party. *Elife*, *10*, e65096.

986

**A) High Load Dual-Task (3-back visual task)**

Target

Target

**B) Low Load Dual-Task (0-back visual task)**

Target

Target

**C) Auditory single-task**

ISI = 2500 ms
500 ms

69 s

**A)**

Pred. acc. (Δz)
0.08
0

**B)**

Pred. acc. (Δz)

*** 
*** 
**

0.15
0.10
0.05
0.00

Auditory single-task | Low load dual-task | High load dual-task

**C)**

RMS (normalized; ×10⁻⁴)

30
20
10
0

Lo- vs Hi-DT
A-ST vs. Lo-DT
A-ST vs. Hi-DT

0  200  400  600  800  1000
Time (ms)

**D)**

150 ms | 240 ms | 360 ms | 480 ms | 680 ms

Auditory single-task

Low load dual-task

High load dual-task

-50   0   50
Voltage(Normalized)

**A) (Vis+Aco+Lin) vs. (Vis+Lin)**

Pred. acc. (Δz; ×10⁻³)

**B) Acoustic**

Auditory single-task | Low load dual-task | High load dual-task

**C) Gammatone Spectrogram**

Pred. acc. (Δz; ×10⁻³)

**D) Onset Spectrogram**

**E)**

Gammatone Spectrogram

RMS (normalized; ×10⁻⁴)

Time (ms)

60 ms | 120 ms | 160 ms | 420 ms

Auditory single-task

Low load dual-task

High load dual-task

**F)**

Onset Spectrogram

A-ST vs. Hi-DT

56 ms | 152 ms | 250 ms | 340 ms

Voltage(Normalized)

-15 0 15

**A) (Vis+Aco+Lin) vs. (Vis+Aco)**
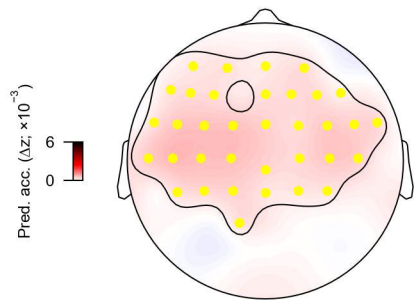
**B) Lingustic**

**C) Context Level**

**E) Entropy vs. Surprisal**

**D)**

**F)**