# Neural speech tracking benefit of lip movements predicts

# behavioral deterioration when the speaker's mouth is

# occluded

Patrick Reisinger[1]*, Marlies Gillis[2], Nina Suess[1], Jonas Vanthornhout[2], Chandra Leon Haider[1],

Thomas Hartmann[1], Anne Hauswald[1], Konrad Schwarz[3], Tom Francart[2†], Nathan Weisz[1,4†]

[1]Paris-Lodron-University of Salzburg, Department of Psychology, Centre for Cognitive

Neuroscience, Salzburg, Austria

[2]Experimental Oto-Rhino-Laryngology, Department of Neurosciences, Leuven Brain Institute,

KU Leuven, Leuven, Belgium

[3]MED-EL GmbH, Innsbruck, Austria

[4]Neuroscience Institute, Christian Doppler University Hospital, Paracelsus Medical University

Salzburg, Salzburg, Austria

*Corresponding author: patrick.reisinger@plus.ac.at

†Shared last authorship

## Abstract

Observing lip movements of a speaker is known to facilitate speech understanding, especially in challenging listening situations. Converging evidence from neuroscientific studies shows enhanced processing of audiovisual stimuli. However, the interindividual variability of this visual benefit and its consequences on behavior are unknown. Here, we analyzed source-localized magnetoencephalographic (MEG) responses from normal-hearing participants listening to audiovisual speech with or without an additional distractor speaker. Using temporal response functions (TRFs), we show that neural responses to lip movements are, in general, enhanced when speech is challenging. After conducting a crucial control for speech acoustics, we show that lip movements effectively contribute to higher neural speech tracking, particularly when a distractor speaker is present. However, the extent of this visual benefit varied greatly among participants. Probing the behavioral relevance, we show that individuals who benefit more from lip movement information in terms of neural speech tracking, show a stronger drop in performance and an increase in perceived difficulty when the mouth is occluded by a surgical face mask. By contrast, no effect was found when the mouth was not occluded. We provide novel insights on how the benefit of lip movements in terms of neural speech tracking varies among individuals. Furthermore, we reveal its behavioral relevance by demonstrating negative consequences for behavior when visual speech is absent. Our results also offer potential implications for future objective assessments of audiovisual speech perception.

## Introduction

Face masks are an important tool in preventing the spread of contagious diseases such as COVID-19 (e.g. Chu et al., 2020; Suñer et al., 2022). However, as many have subjectively experienced first hand, the use of face masks also impairs speech perception, and not only by attenuating sound. More importantly, they occlude facial expressions, such as lip movements (e.g. Brown et al., 2021; Rahne et al., 2021), that provide visual information for a relevant speech stream. This is particularly critical when speech is challenging, such as in the classic cocktail party situation, where multiple conversations are happening simultaneously (Cherry, 1953). In such situations, the brain separates auditory information of interest from competing input (McDermott, 2009). Ideally, visual information is available to support this process, with numerous studies demonstrating that visual speech features enhance the understanding of degraded auditory input (e.g. Grant & Seitz, 2000; Remez, 2012; Ross et al., 2007; Sumby & Pollack, 1954). This concept is known as inverse effectiveness (Meredith & Stein, 1983; van de Rijt et al., 2019). Among visual speech features, lip movements are the most important, playing a crucial role in the perception of challenging speech (Erber, 1975; Peelle & Sommers, 2015). This is especially intriguing given the substantial interindividual differences in lip-reading performance among normal, as well as hearing-impaired, populations (Suess, Hauswald, Zehentner, et al., 2022; Summerfield et al., 1992). Despite our imperfect lip-reading abilities, the human brain effectively uses lip movements to facilitate the perception of challenging speech, with the neural mechanisms and regions involved still under debate (Ross et al., 2022; Zhang & Du, 2022).

Previous studies have shown beneficial effects of visual speech on the representation of speech in the brain. An MEG study by Park et al. (2016) showed enhanced entrainment between lip movements and speech-related brain areas when congruent audiovisual speech was presented. Other studies have shown that the incorporation of visual speech enhances the ability of the brain to track acoustic speech (Crosse et al., 2015; Crosse, Liberto, et al., 2016; Golumbic et al., 2013). Interestingly, when silent lip movements are presented, the brain also tracks the unheard acoustic speech envelope (e.g. Hauswald et al., 2018) or spectral fine details (Suess, Hauswald, Reisinger, et al., 2022). Despite these findings, two questions remain unanswered: First, it is unknown how individuals vary in their benefit of lip movements at the neural level. Given the aforementioned interindividual differences in lip-reading performance, a high degree of variability could also be expected here. Importantly, lip movements are correlated with acoustic speech features (Chandrasekaran et al., 2009), so it is essential to control for acoustic-related brain

66    activity. Second, it is unknown if the individual benefit of lip movements is of behavioral relevance,
67    as, for example, when the lips are occluded with a face mask, as has been common during the
68    COVID-19 pandemic. Given the negative impact of face masks on behavioral measures (e.g.
69    Rahne et al., 2021; Toscano & Toscano, 2021; Truong et al., 2021), a relationship is plausible:
70    Individuals who benefit more should, in principle, also show poorer behavioral outcomes when no
71    lip movements are available, as they are deprived of critical visual information.

72    A suitable method to obtain the individual benefit of lip movements is neural tracking (Obleser &
73    Kayser, 2019). Besides frequency-based coherence and mutual information, temporal response
74    functions (TRFs) have gained widespread popularity (Brodbeck & Simon, 2020; Crosse et al.,
75    2021). TRFs typically aim to predict the M/EEG-recorded neural response to one or more stimulus
76    features, and the prediction is correlated with the original signal to quantify neural tracking. This
77    approach has so far extended our understanding of speech processing from acoustic features
78    (Lalor et al., 2009) to higher-level linguistic features (Brodbeck, Hong, et al., 2018; Broderick et
79    al., 2018; Gillis et al., 2021). Crucially, neural tracking can be used to disentangle the
80    aforementioned intercorrelation of audiovisual speech by controlling for acoustic speech features
81    (Gillis et al., 2022). This could reveal the "pure" individual benefit of lip movements to neural
82    speech tracking in audiovisual settings, which has not yet been shown.

83    Neural speech tracking has been proposed as an objective measure for speech intelligibility
84    (Schmitt et al., 2022; Vanthornhout et al., 2018) along with a whole range of auditory and linguistic
85    processes (Gillis et al., 2022). Previously, acoustic neural speech tracking has been related to
86    behavioral measures such as speech intelligibility (Chen et al., 2023; Ding & Simon, 2013).
87    Studies that involve visual speech features have established a relationship between the neural
88    tracking of visual speech cues, so-called visemes, and lip-reading performance (Nidiffer et al.,
89    2021) or lip movements and speech comprehension (Park et al., 2016). In sum, these findings
90    strongly suggest a meaningful relationship between neural speech tracking and behavioral
91    measures. Regarding the aspect of interindividual differences, Schubert et al. (2023) showed that
92    the MEG-derived tendency of individuals to predict upcoming tones facilitates neural speech
93    tracking, and this relationship generalizes to various audio-only listening situations. Here, we aim
94    to combine both aspects by evaluating the relationship between interindividual differences and
95    behavioral measures. In particular, we probe the behavioral relevance of the individual benefit of
96    lip movements, especially when critical visual information is not available. Addressing this could
97    further strengthen the case for the behavioral relevance of neural speech tracking as an objective
98    measure of speech processing.

99    Here, we used MEG and an audiovisual speech paradigm with one or two speakers to investigate
100   the benefit of lip movements and its behavioral relevance. Utilizing a state-of-the-art neural
101   tracking framework with source-localized TRFs (see Figure 1), we show that lip movements are
102   processed more strongly when speech is challenging. Additionally, we show that the neural
103   tracking of lip movements is enhanced in multi speaker settings. When controlled for acoustic
104   speech features, the obtained benefit of lip movements is, in general, more enhanced in the multi
105   speaker condition, with substantial interindividual variability. Using Bayesian modeling, we show
106   that acoustic speech tracking is related to behavioral measures. Crucially, we demonstrate that
107   individuals who benefit more from lip movements show a stronger drop in performance and report
108   a higher subjective difficulty when the mouth is occluded by a surgical face mask. In terms of
109   neural tracking, our results suggest that individuals benefit from lip movements in a highly variable
110   manner. We also establish a novel link between the neural benefit of visual speech and behavior
111   when no visual speech information is available.

## Material and Methods

*Participants*

114   The data was collected as part of a recent study (Haider et al., 2022), in which 30 native speakers
115   of German participated. One participant was excluded because signal source separation could
116   not be applied to the MEG dataset. This led to a final sample size of 29 participants (12 females,
117   $M_{age}$ = 26.79, $SD_{age}$ = 4.87 years). All participants reported normal vision and hearing (thresholds
118   did not exceed 25 dB HL at any frequency from 125 to 8000 Hz), the latter verified by a standard
119   clinical audiometer (AS608 Basic; Interacoustics A/S, Middelfart, Denmark). Additional exclusion
120   criteria included non-removable magnetic objects and any psychiatric or neurologic history. All
121   participants signed an informed consent and were reimbursed at a rate of 10 € per hour. The
122   experimental protocol was approved by the ethics committee of the Paris-Lodron-University of
123   Salzburg and was conducted in accordance with the Declaration of Helsinki.

*Stimuli and experimental design*

125   The experimental procedure was implemented in MATLAB 9.10 (The MathWorks Inc., Natick,
126   Massachusetts, USA) using custom scripts. Presentation of stimuli and response collection was
127   achieved with the Objective Psychophysics Toolbox (o_ptb; Hartmann & Weisz, 2020), which
128   adds a class-based abstraction layer onto the Psychophysics Toolbox version 3.0.16 (Brainard,

129    1997; Kleiner et al., 2007; Pelli, 1997). Stimuli and triggers were generated and emitted via the

130    VPixx system (DATAPixx2 display driver, PROPixx DLP LED projector, RESPONSEPixx

131    response box by VPixx Technologies Inc., Saint-Bruno, Canada). Videos were back-projected

132    onto a translucent screen with a screen diagonal of 74 cm (~110 cm in front of the participants),

133    with a refresh rate of 120 Hz and a resolution of 1920×1080 pixels. Timings were measured with

134    the Black Box ToolKit v2 (The Black Box ToolKit Ltd., Sheffield, UK) to ensure accurate stimulus
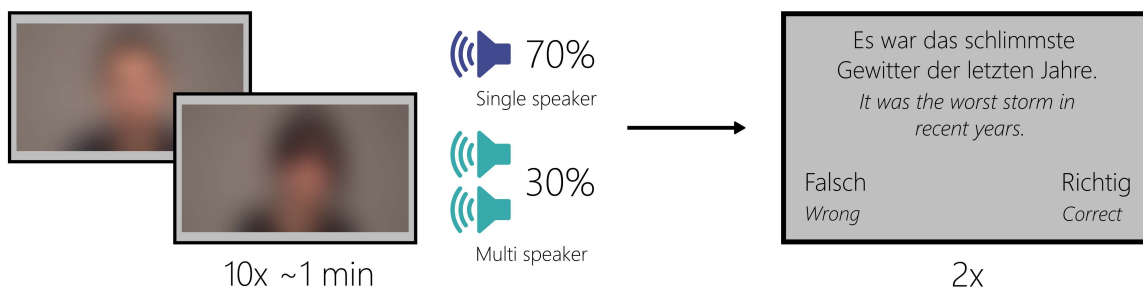
135    presentation and triggering.

136    The audiovisual stimuli were excerpts from four German stories, two of each read out loud by a

137    female or male speaker (female: "Die Schokoladenvilla - Zeit des Schicksals. Die Vorgeschichte

138    zu Band 3" by Maria Nikolai, "Die Federn des Windes" by Manuel Timm; male: "Das Gestüt am

139    See. Charlottes großer Traum" by Paula Mattis and "Gegen den Willen der Väter" by Klaus

140    Tiberius Schmidt). A Sony NEX-FS100 (Sony Group Corporation, Tokyo, Japan) camera with a

141    sampling rate of 25 Hz and a RØDE NTG2 microphone (RØDE Microphones Pty. Ltd., Sydney,

142    Australia) with a sampling rate of 48 kHz were used to record the stimuli. Each of the four stories

143    was recorded twice, once with and once without a surgical face mask (type IIR three-layer

144    disposable medical mask). These eight videos were cut into 10 segments of about one minute

145    each (M = 64.29 s, SD = 4.87 s), resulting in 80 videos. In order to rule out sex-specific effects,

146    40 videos (20 with a female speaker and 20 with a male speaker) were presented to each

147    participant. The speakers' syllable rates were analyzed using Praat (Boersma, 2001; de Jong &

148    Wempe, 2009) and varied between 3.7 Hz and 4.6 Hz (M = 4.1 Hz). The audio-only distractor

149    speech consisted of pre-recorded audiobooks (see Schubert et al., 2023), read by either a female

150    or a male speaker.

151    Before the experiment, a standard clinical audiometry was performed (for details, see

152    *Participants*). The MEG measurement started with a 5-minute resting-state recording (not

153    analyzed in this manuscript). Next, the participant's individual hearing threshold was determined

154    in order to adjust the stimulation volume. If the participant reported that the stimulation was not

155    loud enough or comfortable, the volume was manually adjusted to the participant's requirements.
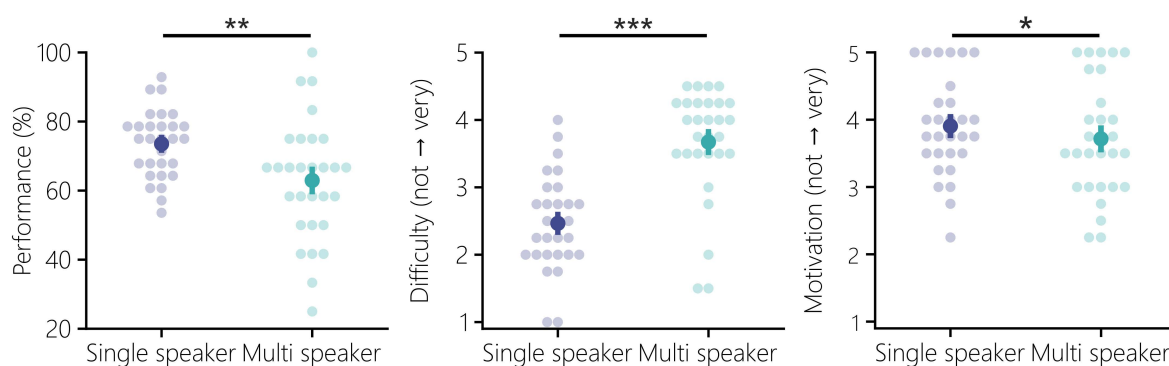
156    The actual experiment consisted of four stimulation blocks, one for each of the four stories, with

157    two featuring each sex. Each story was presented as a block of 10 ~one-minute trials in

158    chronological order to preserve the story content (Figure 1A). In every block, a same-sex audio-

159    only distractor speaker was added to three randomly selected trials, with a 5-second delay and

160    volume equal to the target speaker. The resulting ratio of 30% multi speaker trials and 70% single

161 speaker trials per block was chosen because of a different data analysis method in Haider et al.

162 (2022). The distractor speech started with a delay of 5 seconds to give participants time to attend

163 the target speaker. In two randomly selected blocks, the target speaker wore a face mask (only

164 the corresponding behavioral data was used here, see *Statistical analysis and Bayesian*

165 *modeling*). Two unstandardized correct or wrong statements about semantic content were

166 presented after each trial to assess comprehension performance and to maintain attention (Figure

167 1A). On four occasions in each block, participants also rated subjective difficulty and motivation

168 on a five-point Likert scale (not depicted in Figure 1A). The participants responded by pressing

169 buttons. The total duration of the experiment was ~2 h, including preparation.
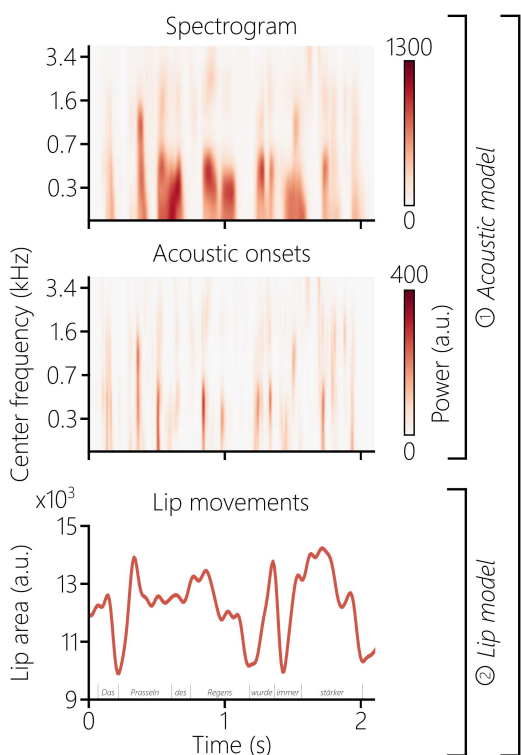
**Figure 1.** *Experimental design, behavioral results and analysis framework.* (A) Each block consisted of 10 ~1-min trials of continuous audiovisual speech by either a female or male speaker (single speaker condition). In 30% of these 10 trials, a same-sex audio-only distractor speaker was added (multi speaker condition). After every block, two comprehension statements had to be rated as correct or wrong. (B) Performance on the comprehension statements in the multi speaker condition was lower than in the single speaker condition ($p = .003$, $r_C = 0.64$). Subjective difficulty ratings, reported on a five-point Likert scale, were higher in the multi speaker condition ($p = 9.00e^{-06}$, $r_C = -0.95$). The reported motivation was lower in the multi speaker condition ($p = .024$, $r_C = 0.62$). The middle dots represent the mean, and the bars, the standard error of the mean. (C) Three stimulus features (spectrogram, acoustic onsets and lip movements) extracted from the audiovisual stimuli are shown for an example sentence. Higher values in the lip area unit represent a wider opening of the mouth and vice versa. Three forward models were calculated: (1) one using only acoustic features, (2) one using only lip movements, and (3) one combining all features. Together with the corresponding source-localized MEG data, the boosting algorithm was used to calculate the models. Exemplary minimum-norm source estimates are shown for a representative participant. The resulting TRFs (a.u.) and neural tracking (expressed as Pearson's *r*) were analyzed in functional regions of interest (fROIs), obtained either via the acoustic or lip model of the multi speaker condition. The TRFs and prediction accuracies shown are from a representative participant reflecting the group-level results. To obtain the benefit of lip movements, acoustic features were controlled by subtracting the prediction accuracies in an acoustic+lip fROI of the acoustic model from the combined model. The benefit of lip movements was expressed as a percentage change. *$p < .05$, **$p < .01$, ***$p < .001$; *Speakers have been blurred due to a bioRxiv policy on the inclusion of faces.*

*MEG data acquisition and preprocessing*

Before entering the magnetically shielded room, five head position indicator (HPI) coils were applied on the scalp. Electrodes for electrooculography (EOG; vertical and horizontal eye movements) and electrocardiography (ECG) were also applied (recorded data not used here). Fiducial landmarks (nasion and left/right pre-auricular points), the HPI locations and ~300 head shape points were sampled with a Polhemus FASTRAK digitizer (Polhemus, Colchester, Vermont, USA).

Magnetic brain activity was recorded with a Neuromag Triux whole-head MEG system (MEGIN Oy, Espoo, Finland) using a sampling rate of 1000 Hz (hardware filters: 0.1-330 Hz). The signals were acquired from 102 magnetometers and 204 orthogonally placed planar gradiometers at 102 different positions. The system is placed in a standard passive magnetically shielded room (AK3b; Vacuumschmelze GmbH & Co. KG, Hanau, Germany).

203   A signal space separation (SSS; Taulu & Kajola, 2005; Taulu & Simola, 2006) algorithm

204   implemented in MaxFilter version 2.2.15 provided by the MEG manufacturer was used. The

205   algorithm removes external noise from the MEG signal (mainly 16.6 Hz, and 50 Hz, plus

206   harmonics) and realigns the data to a common standard head position (to [0 0 40] mm, -*trans*

207   *default* MaxFilter parameter) across different blocks, based on the measured head position at the

208   beginning of each block.

209   Preprocessing of the raw data was done in MATLAB 9.8 using the FieldTrip toolbox (revision

210   f7adf3ab0; Oostenveld et al., 2011). A low-pass filter of 10 Hz (hamming-windowed sinc FIR filter,

211   onepass-zerophase, order: 1320, transition width: 2.5 Hz) was applied, and the data was

212   downsampled to 100 Hz. Afterwards, a high-pass filter of 1 Hz (hamming-windowed sinc FIR filter,

213   onepass-zerophase, order: 166, transition width: 2.0 Hz) was applied.

214   Independent component analysis (ICA) was used to remove eye and cardiac artifacts (data was

215   filtered between 1-100 Hz, sampling rate: 1000 Hz) via the infomax algorithm ("runica"

216   implementation in EEGLAB; Bell & Sejnowski, 1995; Delorme & Makeig, 2004) applied to a

217   random block of the main experiment. Prior to the ICA computation, we performed a principal

218   component analysis (PCA) with 50 components in order to ease the convergence of the ICA

219   algorithm. After visual identification of artifact-related components, an average of 2.38

220   components per participant were removed (SD = 0.68).

221   The cleaned data was epoched into trials that matched the length of the audiovisual stimuli. To

222   account for an auditory stimulus delay introduced by the tubes of the sound system, the data were

223   shifted by 16.5 ms. In the multi speaker condition, the first 5 seconds of data were removed to

224   match the onset of the distractor speech. The last eight trials were removed to equalize the data

225   length between the single speaker and multi speaker conditions. To prepare the data for the

226   following steps, the trials in each condition were concatenated. This resulted in a data length of

227   ~6 min per condition.

228   *Source localization*

229   Source projection of the data was done with MNE-Python 1.1.0 running on Python 3.9.7 (Gramfort

230   et al., 2013, 2014). A semi-automatic coregistration pipeline was used to coregister the FreeSurfer

231   "fsaverage" template brain (Fischl, 2012) to each participant's head shape. After an initial fit using

232   the three fiducial landmarks, the coregistration was refined with the Iterative Closest Point (ICP)

233   algorithm (Besl & McKay, 1992). Head shape points that were more than 5 mm away from the

234    scalp were automatically omitted. The subsequent final fit was visually inspected to confirm its

235    accuracy. This semi-automatic approach performs comparably to manual coregistration pipelines

236    (Houck & Claus, 2020).

237    A single-layer boundary element model (BEM; Akalin-Acar & Gençer, 2004) was computed to

238    create a BEM solution for the "fsaverage" template brain. Next, a volumetric source space with a

239    grid of 7 mm was defined, containing a total of 5222 sources (Kulasingham et al., 2020). In order

240    to remove non-relevant regions and shorten computation times, subcortical structures along the

241    midline were removed, reducing the source space to 3053 sources (similar to Das et al., 2020).

242    Subsequently, the forward operator (i.e. lead field matrix) was computed using the individual

243    coregistrations, the BEM and the volume source space.

244    Afterwards, the data were projected to the defined sources using the Minimum Norm Estimate

245    method (MNE; Hämäläinen & Ilmoniemi, 1994). MNE is known to be biased towards superficial

246    sources, which can be reduced by applying depth weighting with a coefficient between 0.6 and

247    0.8 (Lin et al., 2006). For creating the MNE inverse operator, depth weighting with a coefficient of

248    0.8 was used (e.g. Brodbeck et al., 2018). The required noise covariance matrix was estimated

249    with an empty-room MEG recording relative to the participant's measurement date with the same

250    preprocessing settings as the MEG data of the actual experiment (see *MEG data acquisition and*

251    *preprocessing*). The MNE inverse operator was then applied to the concatenated MEG data with

252    $\ell$2 regularization (signal-to-noise ratio (SNR) = 3 dB, $\lambda^2 = \frac{1}{SNR^2}$) and three free-orientation dipoles

253    orthogonally at each source.

254    *Extraction of stimulus features*

255    Since the focus of this study is on audiovisual speech, we extracted acoustic (spectrograms and

256    acoustic onsets) and visual (lip movements) speech features from the stimuli (for examples see

257    Figure 1C). The spectrograms of the auditory stimuli were obtained using the Gammatone

258    Filterbank Toolkit 1.0 (Heeris, 2013), with frequency cutoffs at 20 and 5000 Hz, 256 filter channels

259    and a window time of 0.01 s. This toolkit computes a spectrogram representation on the basis of

260    a set of Gammatone filters which are inspired by the human auditory system (Slaney, 1998). The

261    resulting filter outputs with logarithmic center frequencies were averaged into eight frequency

262    bands (frequencies <100 Hz were omitted; Gillis et al., 2021). Each frequency band was scaled

263    with exponent 0.6 (Biesmans et al., 2017) and downsampled to 100 Hz, which is the same

264    sampling frequency as the preprocessed MEG data.

265 Acoustic onset representations were calculated for each frequency band of the spectrograms
266 using an auditory edge detection model (Fishbach et al., 2001). The resulting spectrograms of
267 the acoustic onsets are valuable predictors of MEG responses to speech stimuli (Brodbeck et al.,
268 2020; Daube et al., 2019). A delay layer with 10 delays from 3 to 5 ms, a saturation scaling factor
269 of 30 and a receptive field based on the derivative of a Gaussian window (SD = 2 ms) were used
270 (Gillis et al., 2021). Each frequency band was downsampled to 100 Hz.

271 The lip movements of every speaker were extracted from the videos with a MATLAB script
272 adapted from Suess et al. (2022; originally by Park et al., 2016). Within the lip contour, the area,
273 and the horizontal and vertical axis were calculated. Only the area was used for the analysis,
274 which leads to results comparable to using the vertical axis (Park et al., 2016). The lip area signal
275 was upsampled from 25 Hz to 100 Hz using FFT-based interpolation.

276 *Forward models*

277 A linear forward modeling approach was used to predict the MEG response to the aforementioned
278 stimulus features (see Figure 1C). These approaches are based on the idea that the brain's
279 response to a stimulus is a continuous function in time (Lalor et al., 2006). The boosting algorithm
280 (David et al., 2007), implemented in eelbrain 0.38 (running on Python 3.9.7; Brodbeck et al.,
281 2022), was used to predict MNE source-localized MEG responses to stimulus features ("MNE-
282 boosting"; Brodbeck, Presacco, et al., 2018). For multiple stimulus features, the linear forward
283 model can be formulated as:

284
$$\hat{y}_t = \sum_{i=0}^{n} \sum_{\tau=\tau_{min}}^{\tau_{max}} h_{i,\tau}\, x_{i,t-\tau}$$

285 For every $n$ stimulus feature, the algorithm finds an optimal filter kernel $h$, which is also known as
286 a temporal response function (TRF). When $n$ stimulus features is > 1, $h$ is referred to as
287 multivariate TRF (mTRF). The term $\tau$ denotes the delays between the predicted brain response
288 $\hat{y}_t$ and stimulus feature $x$ (for further details see Brodbeck et al., 2022). TRFs reflect responses
289 to continuous data instead of averaged responses to discrete events (Crosse et al., 2021). For
290 the estimation of the TRFs, the stimulus features and MEG data were normalized (z-scored), and
291 an integration window from -100 to 600 ms with a kernel basis of 50 ms Hamming windows was
292 defined. To prevent overfitting, early stopping based on the ℓ2 norm was used. By using four-fold
293 nested crossvalidation (two training folds, one validation fold, and one test fold), each partition

294   served as a test set once (Brodbeck et al., 2022). TRFs were estimated for each of the three free-
295   orientation dipoles independently at all 3053 sources (see *Source localization*). The spectrogram
296   and acoustic onset mTRFs were averaged over the frequency dimension. To account for
297   interindividual anatomical differences, TRFs were spatially smoothed with a Gaussian kernel (SD
298   = 5 mm; Kulasingham et al., 2020). The vector norm of the smoothed TRFs was taken, resulting
299   in one TRF per source.

300   To obtain a measure of neural tracking, the predicted brain response $\hat{y}_t$ is correlated with the
301   original response to calculate the prediction accuracy and computed as the average dot product
302   over time (expressed as Pearson correlation coefficient *r*). This correlation can be interpreted as
303   follows: The higher the prediction accuracy, the higher the neural tracking (Gillis et al., 2022).

304   In order to investigate the neural processing of the audiovisual speech features, we calculated
305   three different forward models per condition and participant (see Figure 1C for the analysis
306   framework). The acoustic model consisted of the two acoustic stimulus features (spectrogram
307   and acoustic onsets) and – also applicable to all other models – the corresponding MNE source-
308   localized MEG data. The lip model contained only the lip movements as a stimulus feature.
309   Additionally, a combined acoustic+lip model was calculated to control for acoustic features in a
310   subsequent analysis.

311   We defined functional regions of interest (fROIs; Nieto-Castanon et al., 2003) by creating labels
312   based on the 90th percentile of the whole-brain prediction accuracies in the multi speaker
313   condition (similar to Suess, Hauswald, Reisinger, et al., 2022). The multi speaker condition was
314   chosen for extracting the fROIs because it potentially incorporates all included stimulus features,
315   due to its higher demand (Golumbic et al., 2013). This was done separately for the acoustic and
316   lip models to map their unique neural sources (see Figure 1C). According to the "aparc"
317   FreeSurfer parcellation (Desikan et al., 2006), the acoustic fROI mainly involved sources in the
318   temporal, lateral parietal and posterior frontal lobes. The superior parietal and lateral occipital
319   lobes made up the majority of the lip fROI. To obtain an audiovisual fROI for the acoustic+lip
320   model, we combined the labels of the acoustic and lip fROIs.

321   For every model, the TRFs in their respective fROI were averaged and, exclusively for Figure 2A,
322   smoothed over time with a 50 ms Hamming window. Grand-average TRF magnitude peaks were
323   detected with scipy version 1.8.0 (running on Python 3.9.7; Virtanen et al., 2020) and visualized
324   as a difference between the multi and single speaker conditions. To suppress regression artifacts
325   that typically occur (Crosse, Di Liberto, et al., 2016), TRFs were visualized between -50 and 550

13

326    ms. Prediction accuracies in the fROIs were Fisher z-transformed, then averaged, and then the

327    z-values were back-transformed to Pearson correlation coefficients (Corey et al., 1998). For the

328    lower panels of each model in Figure 2B, the prediction accuracies of the acoustic and lip models

329    were averaged in their respective fROIs. Figures were created with the built-in plotting functions

330    of eelbrain and seaborn version 0.12.0 (running on Python 3.9.7; Waskom, 2021).

331    In order to answer the question whether or not lip movements enhance neural tracking, a control

332    for acoustic features is needed. This is particularly important due to the intercorrelation of speech

333    features (Chandrasekaran et al., 2009; Daube et al., 2019). To investigate the individual benefit

334    of lip movements, we used the averaged prediction accuracies in the audiovisual fROI and

335    subtracted the acoustic model from the acoustic+lip model (for a general overview on control

336    approaches see Gillis et al., 2022). The resulting individual benefit of lip movements was

337    expressed as percentage change (see Figure 2C).

338    *Statistical analysis and Bayesian modeling*

339    All frequentist statistical tests were conducted with built-in functions from eelbrain and the

340    statistical package pingouin version 0.5.2 (running on Python 3.9.7; Vallat, 2018). The three

341    behavioral measures (performance, difficulty, and motivation; Figure 1B) were statistically

342    compared between the two conditions (single speaker and multi speaker) using a Wilcoxon

343    signed-rank test and the matched-pairs rank-biserial correlation $r_C$ was reported as effect size

344    (King et al., 2018).

345    The TRFs corresponding to the three stimulus features (spectrogram, acoustic onsets and lip

346    movements; Figure 2A), were tested for statistical difference between the two conditions using a

347    cluster-based permutation test with threshold-free cluster enhancement (TFCE; dependent

348    samples t-test, 10000 randomizations, Maris & Oostenveld, 2007; Smith & Nichols, 2009). Due to

349    the previously mentioned TRF regression artifacts, the time window for the test was limited to -50

350    to 550 ms. Depending on the direction of the cluster, the maximum or minimum *t*-value was

351    reported and Cohen's *d* of the averaged temporal extent of the cluster was calculated.

352    We tested the non-averaged prediction accuracies in the acoustic and lip fROIs (Figure 2B) with

353    a cluster-based permutation test with TFCE (dependent samples t-test, 10000 randomizations).

354    According to the cluster's direction, the maximum or minimum *t*-value was reported, and Cohen's

355    *d* of the cluster's averaged spatial extent was calculated. Additionally, averaged prediction

356    accuracies in the acoustic and lip fROIs were statistically tested with a dependent-samples t-test,

357    and Cohen's *d* was reported as effect size. In the audiovisual fROI, the prediction accuracies and

358    benefit of lip movements (Figure 2C) were tested with a dependent-samples t-test, and Cohen's

359    *d* was reported as effect size. If the data were not normally distributed according to a Shapiro-

360    Wilk test, the Wilcoxon signed-rank test was used, and the matched-pairs rank-biserial correlation

361    $r_C$ was reported as effect size. The distribution of the benefit of lip movements was assessed

362    using the bimodality coefficient (Freeman & Dale, 2013).

363    To investigate if neural tracking is predictive for behavior, we calculated Bayesian multilevel

364    models in R version 4.2.2 (R Core Team, 2022) with the Stan-based package brms version 2.18.4

365    (Bürkner, 2017; Carpenter et al., 2017). Neural tracking (i.e. the averaged prediction accuracies

366    within the respective fROI) was used to separately predict the three behavioral measures. A

367    random intercept was added for each participant to account for repeated measures (single

368    speaker and multi speaker). The models were fitted independently for the acoustic and lip models

369    (Figure 3). According to the Wilkinson notation (Wilkinson & Rogers, 1973), the general formula

370    was:

371    $$behavioral\ measure \sim 1\ +\ neural\ tracking\ +\ (1\ |\ participant)$$

372    We wanted to test whether the individual benefit of lip movements to neural speech tracking (see

373    *Forward models*) yields any behavioral relevance. For this, we also used the behavioral data of

374    the otherwise unanalyzed conditions with a face mask (see *Stimuli and experimental design*). We

375    fitted Bayesian multilevel models with the individual benefit of lip movements to separately predict

376    the behavioral measures when the speaker wore a face mask or not (Figure 4). The general

377    formula was:

378    $$behavioral\ measure \sim 1\ +\ benefit\ of\ lip\ movements\ +\ (1\ |\ participant)$$

379    Before doing so, we fitted control models to show the effect of the conditions on the behavioral

380    measures when the lips are occluded (see Supplementary Table 1). Additional control models to

381    test the effect of the benefit of lip movements on the behavioral data without a face mask were

382    also fitted (see Supplementary Table 2 for model fits). In all described models, a random intercept

383    was included for each participant to account for repeated measures (single speaker and multi

384    speaker).

385    Weakly or non-informative default priors of brms were used, whose influence on the results is

386    negligible (Bürkner, 2017, 2018). For model calculation, all numerical variables were z-scored,

15

387    and standardized regression coefficients (*b*) were reported with 89% credible intervals (CIs; i.e.

388    Bayesian uncertainty intervals, McElreath, 2020). In addition, we report posterior probabilities

389    ($PP_{b>0}$) with values closer to 100%, providing evidence that the effect is greater than zero, and

390    closer to 0% that the effect was reversed (i.e. smaller than zero). If the 89% CIs for an estimate

391    did not include zero and $PP_{b>0}$ was below 5.5% or above 94.5%, the effects were considered

392    statistically significant.

393    All models were fitted with a Student-t distribution, as indicated by graphical posterior predictive

394    checks, Pareto $\hat{k}$ diagnostics (Vehtari, Simpson, et al., 2022) and leave-one-out crossvalidation

395    via loo version 2.5.1 (Vehtari et al., 2017; Vehtari, Gabry, et al., 2022). Common algorithm-

396    agnostic (Vehtari et al., 2021) and algorithm-specific diagnostics (Betancourt, 2018) showed that

397    all Bayesian multilevel models converged. For all relevant parameters, the convergence

398    diagnostic $\hat{R} < 1.01$ and effective sample size (ESS) > 400 indicated that there were no divergent

399    transitions. Figures were created with ggplot2 version 3.4.0 (Wickham, 2016) and ggdist version

400    3.2.0 (Kay, 2022). Unstandardized *b*'s were used for the fitted values of the models in Figures 3

401    and 4.

402    *Data and Code Availability*

403    Preprocessed data and code are publicly available at GitHub (https://github.com/reispat/

404    av_speech_mask).

405    **Results**

406    Twenty-nine participants listened to audiobooks with a corresponding video of the speaker and a

407    randomly occurring audio-only distractor. Source-localized MEG responses to acoustic features

408    (spectrogram and acoustic onsets) and lip movements were predicted using forward models

409    (TRFs). We compared the TRFs between the two conditions and evaluated neural tracking of the

410    acoustic features and lip movements. The individual benefit of lip movements was obtained by

411    controlling for acoustic features and was compared between conditions. Using Bayesian

412    multilevel modeling, we predicted the behavioral measures with neural tracking. We also probed

413    the individual benefit of lip movements for their behavioral relevance by predicting the behavioral

414    measures when the lips were occluded with a surgical face mask or not.

415    *Listening situations with multiple speakers are behaviorally more demanding*

416    Participants performed worse in the multi speaker condition (M = 62.93%, SD = 17.34%),

417    compared to the single speaker condition (M = 73.52%, SD = 9.71%; $W$ = 73.00, $p$ = .003, $r_C$ =

418    0.64). In the multi speaker condition, subjective difficulty ratings were higher (M = 3.67, SD = 0.82)

419    than in the single speaker condition (M = 2.47, SD = 0.71; $W$ = 11.50, $p$ = 9.00e$^{-06}$, $r_C$ = -0.95).

420    Motivation was rated higher in the single speaker condition (M = 3.91, SD = 0.74) compared to

421    the multi speaker condition (M = 3.72, SD = 0.85; $W$ = 29.00, $p$ = .024, $r_C$ = 0.62). Overall,

422    behavioral data showed that in the multi speaker condition, participants performed worse,

423    reported the task to be more difficult and were less motivated (Figure 1B).

424    *Neural responses to lip movements are enhanced in a multi speaker setting*
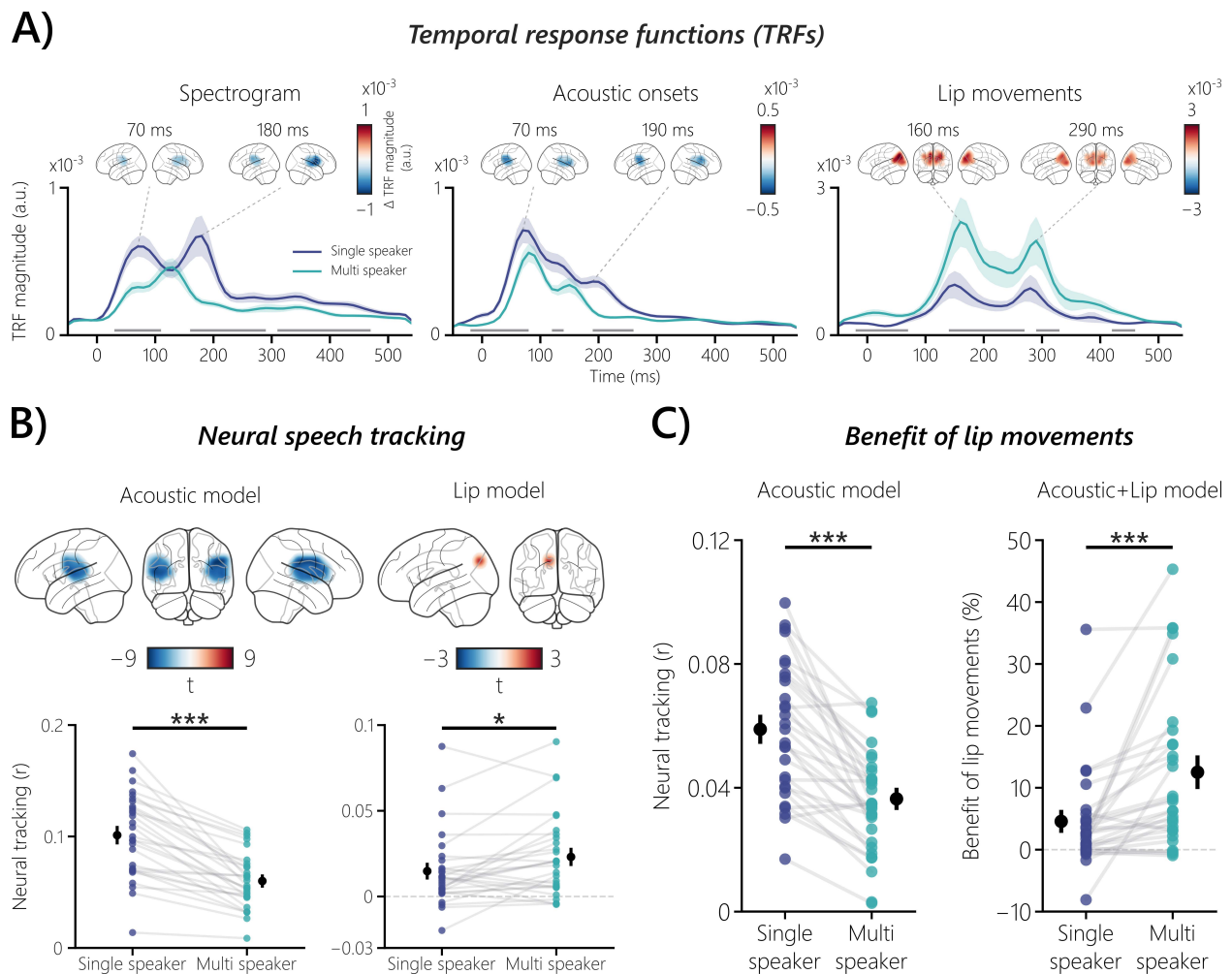
425    First, we analyzed the neural responses to acoustic and visual speech features by statistically

426    comparing the corresponding TRFs between the single- and multi speaker conditions within their

427    respective fROIs (Figure 2A). The spectrogram TRFs showed a significant difference between

428    conditions, with three clusters extending from early (30 to 110 ms; $t$ = -5.26, $p$ = .0001, $d$ = -0.81),

429    middle (160 to 290 ms; $t$ = -3.78, $p$ = .003, $d$ = -1.00) and late (310 to 470 ms; $t$ = -5.58, $p$ = .0001,

430    $d$ = -1.02) time ranges. Grand-average TRF peaks are more pronounced in the single speaker

431    condition, with two peaks at 70 and 180 ms. While the first peak is also present in the multi speaker

432    condition, the second peak appeared 50 ms earlier than the single speaker setting. The latter

433    peak caused the largest differences in the magnitudes of the TRFs, which are most prominent in

434    the right hemisphere of the fROI.

435    The TRFs to acoustic onsets showed a significant difference between single- and multi speaker

436    speech, with three clusters extending from early (-20 to 80 ms; $t$ = -5.39, $p$ < .001, $d$ = -1.10;

437    Figure 2A), mid (120 to 140 ms; $t$ = -4.54, $p$ = .004, $d$ = -1.43) and mid-late (190 to 260 ms; $t$ = -

438    6.11, $p$ < .001, $d$ = -1.13) time windows. The TRFs showed two peaks at 70 and 190 ms in the

439    single speaker condition. Similar to the spectrogram TRFs, the first peak in the multi speaker

440    condition is at the same time point as in the single speaker condition and the second peak is 50

441    ms earlier. The magnitude differences across peaks and hemispheres are not substantially

442    different.

443    TRFs to lip movements show an opposite pattern to the TRFs to acoustic features, with stronger

444    processing in the multi speaker condition. Significant condition differences in the TRFs to lip

445    movements between single- and multi speaker speech were found, with four clusters extending

17

446    from early (-20 to 70 ms; $t$ = 4.41, $p$ = .0005, $d$ = 0.86; Figure 2A), mid (140 to 270 ms; $t$ = 3.97,

447    $p$ = .001, $d$ = 0.88), mid-late (290 to 330 ms; $t$ = 3.34, $p$ = .01, $d$ = 0.91) and late (420 to 460 ms;

448    $t$ = 3.90, $p$ = .002, $d$ = 0.90) time windows. The latencies of the peaks were later in general (160

449    and 290 ms), as compared to the acoustic TRFs, which is also in line with the longer duration for

450    a stimulus to reach the visual system (Thorpe et al., 1996; VanRullen & Thorpe, 2001). In the

451    single speaker condition, they are delayed by 10 ms, and magnitude differences are most

452    prominent in the first peak and left hemisphere.

453    Our initial analysis showed that neural responses to acoustic features are stronger when speech

454    is clear. In contrast, neural responses to lip movements were enhanced in a multi speaker

455    environment. The stronger processing of lip movements suggests a greater reliance on the lips

456    of a speaker when speech is harder to understand.

**Figure 2.** *Neural responses to audiovisual speech features, neural speech tracking, and the benefit of lip movements.* (A) The three plots show grand-averaged TRFs for the stimulus features in their respective fROIs and the peak magnitude contrasts (multi speaker vs. single speaker) between the two conditions in the involved sources. For the acoustic features, TRF magnitudes were generally enhanced when speech was clear, with significant differences ranging from $p = .004$ to $p < .001$ ($d = -0.81$ to $-1.43$). In contrast, the TRF to lip movements showed an enhanced magnitude in the multi speaker condition ($p = .01$ to $p = .0005$ and effect sizes from $d = 0.86$ to $0.91$). The shaded areas of the respective conditions represent the standard error of the mean (SEM). Gray bars indicate the temporal extent of significant differences ($p < .05$) between the two conditions. (B) Neural speech tracking is shown for the non-averaged and averaged fROIs of the acoustic and lip models. Acoustic neural tracking was higher in the single speaker condition, with significant left- and right-hemispheric differences (both $p < .001$ with $d$ from $-1.30$ to $-1.47$; averaged: $p = 8.76e^{-09}$, $d = -1.30$). Lip movements were tracked higher in the multi speaker condition ($p = .037$, $d = 0.51$; averaged: $p = .026$, $r_C = 0.48$). In the averaged plots, the black dots represent the mean, and the corresponding bars the SEM, of the respective condition. (C) In a combined acoustic and lip fROI, the acoustic model showed higher neural tracking in the single speaker condition ($p = 7.68e^{-08}$, $d = 1.18$). The benefit of lip movements was obtained by subtracting the acoustic model from the acoustic+lip model and expressed as percentage change. Lip movements especially enhanced neural tracking in the multi speaker condition ($p = .00003$, $r_C = 0.89$). Participants showed high interindividual variability with a visual benefit of up to 45.37%, but also only a small benefit or no benefit at all. The black dots represent the mean, and the corresponding bars the SEM, of the respective condition. $^*p < .05$, $^{**}p < .01$, $^{***}p < .001$

*The cocktail party diametrically affects acoustic and visual neural speech tracking*

So far, the TRF results indicate a stronger neural response to lip movements and a weaker one to acoustic features when there is more than one simultaneous speaker. We also wanted to answer the question whether neural tracking of audiovisual speech features differs between the single speaker and multi speaker conditions in their respective fROIs (Figure 2B; see Figure S1 for whole-brain neural tracking of the audiovisual speech features). Acoustic neural tracking in the non-averaged acoustic fROI showed a significant condition difference in the left ($t = -8.04$, $p < .001$, $d = -1.47$) and right ($t = -9.26$, $p < .001$, $d = -1.30$) hemispheres. Averaged acoustic neural tracking was higher in the single speaker condition than in the multi speaker condition ($t(28) = -8.07$, $p = 8.76e^{-09}$, $d = -1.30$). Neural tracking of lip movements showed a significant condition difference in the left hemisphere ($t = 3.83$, $p = .037$, $d = 0.51$; Figure 2B), with a focal inferior parietal area involved. When averaging over sources, neural tracking was higher in the multi speaker condition than in the single speaker condition ($W = 114.00$, $p = .026$, $r_C = 0.48$).

490    Overall, the results showed that neural tracking was enhanced for acoustic features when speech

491    is clear, and higher for lip movements when there are multiple speakers. This is in line with the

492    observed neural responses.

493    *Lip movements enhance neural speech tracking more in multi speaker situations*

494    When there are two speakers, we have so far demonstrated that lip movements are processed

495    more strongly and lead to higher neural tracking compared to one speaker. However, their unique

496    contribution to neural tracking is still unknown, due to the intercorrelation of speech features

497    (Chandrasekaran et al., 2009; Daube et al., 2019). To address this, we controlled for the acoustic

498    features so as to obtain the unique benefit of lip movements over and above acoustic speech

499    features. First, the acoustic model was evaluated in the audiovisual fROI (Figure 2C). Acoustic

500    neural tracking was higher in the single speaker condition than in the multi speaker condition

501    ($t(28) = -7.20$, $p = 7.68e^{-08}$, $d = 1.18$). The acoustic model served as a baseline and was subtracted

502    from a combined acoustic+lip model and expressed as percentage change. The obtained benefit

503    of lip movements was higher in the multi speaker condition than in the single speaker condition

504    ($W = 24.00$, $p = .00003$, $r_C = 0.89$). The benefit of lip movements showed high interindividual

505    variability and seemed to follow a bimodal distribution (Figure 2C), which was confirmed by a

506    bimodality coefficient of 0.68 (values > 0.555 indicate bimodality; Pfister et al., 2013).

507    These results strongly indicate that lip movements enhance neural tracking, especially in multi-

508    talker speech. However, substantial interindividual variability was observed, with participants

509    showing an individual benefit of lip movements of up to 45.37% in the multi speaker condition,

510    while others showed only a small benefit or no benefit at all. In the next steps, we will probe the

511    behavioral relevance of the benefit that lip movements provide to neural speech tracking by

512    depriving individuals of this source of information.
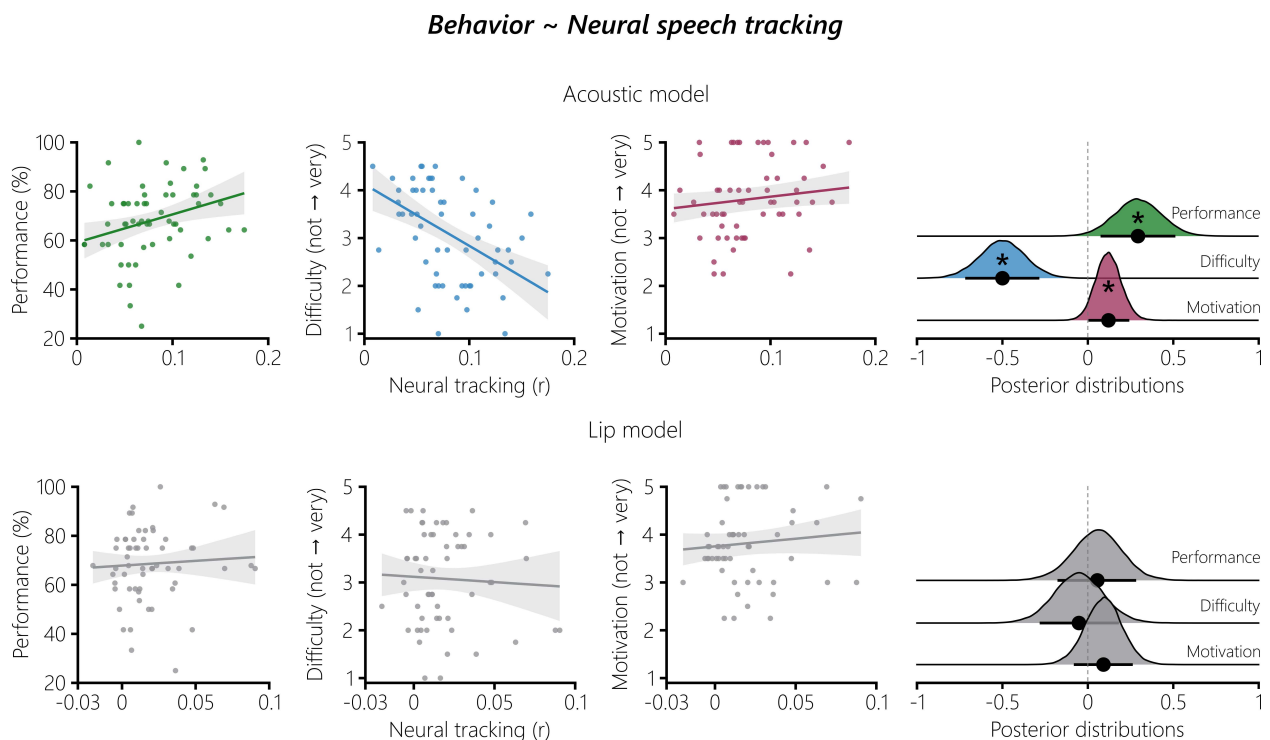
513    *Only acoustic neural speech tracking predicts behavior*

514    Having established that listening situations with two speakers affect neural tracking of acoustic

515    and visual speech features in a diametrical way, we were further interested if neural tracking is

516    able to predict the behavioral measures. We calculated Bayesian multilevel models to predict the

517    three behavioral measures (performance, difficulty and motivation; see Figure 1B) with the

518    averaged neural tracking of the acoustic and lip models (Figure 3). In the acoustic model, higher

519    neural tracking was linked to higher performance ($b = 0.29$, 89% CI = [0.07, 0.51], $PP_{b>0} =$

520    98.37%). Lower neural tracking predicted higher difficulty ratings ($b = -0.50$, 89% CI = [-0.72, -

521    0.29], $PP_{b>0}$ = 0.01%). When neural tracking was high, the motivation ratings were also higher (*b*

522    = 0.12, 89% CI = [0.004, 0.24], $PP_{b>0}$ = 95.05%).

523    Neural tracking of lip movements was not related to performance (*b* = 0.06, 89% CI = [-0.18, 0.28],

524    $PP_{b>0}$ = 65.61%; Figure 3). We also observed no evidence for an effect of the difficulty (*b* = -0.05,

525    89% CI = [-0.28, 0.18], $PP_{b>0}$ = 35.63%) or motivation (*b* = 0.09, 89% CI = [-0.08, 0.26], $PP_{b>0}$ =

526    80.40%) ratings.

527    These results indicate that acoustic neural speech tracking predicts behavior: The higher the

528    neural speech tracking, the higher the performance and motivation ratings. Lower acoustic neural

529    speech tracking was linked to higher difficulty ratings. In contrast, neural speech tracking of lip

530    movements did not predict behavior.

### *Behavior ~ Neural speech tracking*



531    **Figure 3.** *Relating behavior to neural speech tracking.* Bayesian multilevel models were fitted to predict the

532    behavioral measures with neural speech tracking. Higher acoustic neural speech tracking was linked to

533    higher performance, lower difficulty ratings and higher motivation ratings. No evidence for an effect was

534    observed for the neural tracking of lip movements. The shaded areas show the 89% CIs of the respective

535    model. The distributions on the right show the posterior draws of the three models. The black dots represent

536    the mean standardized regression coefficient *b* of the corresponding model. The corresponding bars show

537    the 89% CI. If zero was not part of the 89% CI, the effect was considered significant (*).
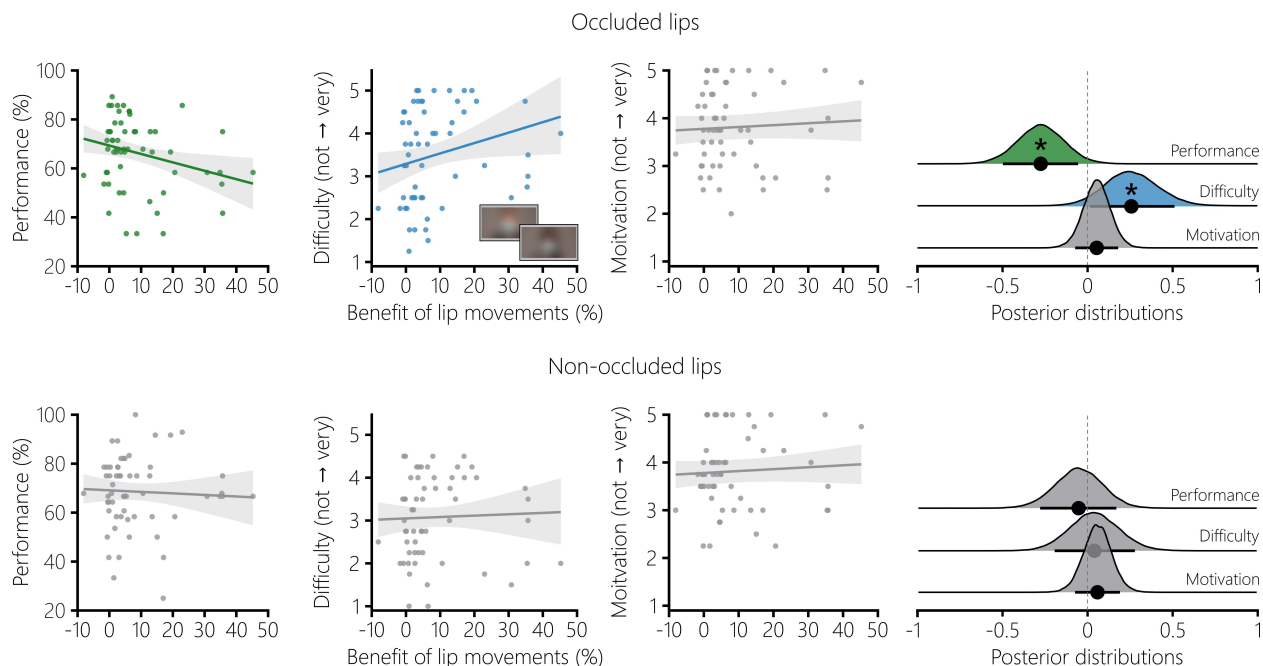
538   *Stronger benefit of lip movements predicts behavioral deterioration when lips are occluded*

539   Given the finding that lip movements enhance neural speech tracking (Figure 2C), we were
540   interested in whether this visual benefit is behaviorally relevant. To do so, we also used the
541   behavioral data from the otherwise unanalyzed conditions in which the mouth was occluded by a
542   surgical face mask (see Figure 4 for an example). Given that critical visual information is missing
543   in these conditions, individuals who show a strong benefit of lip movements on a neural level
544   should show poorer behavioral outcomes. An initial analysis showed that the effect of the
545   conditions with a surgical face mask on behavior followed the same pattern as those with non-
546   occluded lips (see Figure 1B), although with no effect on the motivation ratings. These control
547   models are reported in Supplementary Table 1.

548   While the effects on a solely behavioral level seem not to differ when the lips are occluded or not,
549   predicting the behavioral measures with the lip benefit showed the expected outcome (Figure 4):
550   Participants that had a higher benefit of lip movements in terms of neural tracking showed a
551   decline in performance ($b$ = -0.27, 89% CI = [-0.49, -0.06], $PP_{b>0}$ = 2.21%) and reported the task
552   to be more difficult ($b$ = 0.25, 89% CI = [0.01, 0.51], $PP_{b>0}$ = 95.41%). The motivation ratings did
553   not yield an effect ($b$ = 0.05, 89% CI = [-0.07, 0.18], $PP_{b>0}$ = 76.14%).

554   Interestingly, we were not able to establish a link between the benefit of lip movements to the
555   behavioral data when the lips were not occluded (Figure 4; see Supplementary Table 2 for model
556   fits). Taken together, these findings support a behavioral relevance of the benefit of lip
557   movements. Individuals that benefit more from lip movements on a neural level performed worse
558   and reported the task to be more difficult when the mouth of the speaker was covered by a surgical
559   face mask.

**Figure 4.** *Relating the benefit of lip movements to behavior.* The benefit of lip movements was used to predict the behavioral measures when the lips are occluded or not. The values of the fitted Bayesian multilevel models are shown with a depiction of the conditions in which the speakers wore a surgical face mask. When the benefit of lip movements was high, performance was lower, and difficulty was reported higher. No evidence for an effect was observed for the motivation rating. The behavioral measures when the lips were not occluded were not linked to the benefit of lip movements. The shaded areas show the 89% CIs of the respective model. The distributions on the right show the posterior draws of the three models. The black dots represent the mean standardized regression coefficient *b* of the corresponding model. The corresponding bars show the 89% CI. If zero was not part of the 89% CI, the effect was considered significant (*). *Speakers have been blurred due to a bioRxiv policy on the inclusion of faces.*

## Discussion

Neural speech tracking is widely used to study the neural processing of continuous speech, though primarily with audio-only stimuli (Brodbeck, Hong, et al., 2018; Chalas et al., 2022; Di Liberto et al., 2015; Keitel et al., 2018). Recent studies have used audiovisual speech settings, but without directly modeling the visual speech features (Crosse, Liberto, et al., 2016; Golumbic et al., 2013) or not incorporating their temporal dynamics due to the use of frequency-based methods (Aller et al., 2022; Bröhl et al., 2022; Park et al., 2016). Here, we show, for the first time, the temporal dynamics and cortical origins of TRFs obtained from lip movements in an audiovisual

578    setting with one or two speakers. Using these neural responses, we demonstrate that the neural

579    tracking of lip movements is enhanced in a multi speaker situation compared to a single speaker.

580    When controlling for acoustic speech features, we show that the obtained benefit of lip

581    movements is enhanced in the multi speaker condition, although with high interindividual

582    variability. Using Bayesian modeling, we demonstrate that acoustic neural speech tracking

583    predicts the behavioral measures. Furthermore, individuals who displayed a higher benefit of lip

584    movements showed a stronger behavioral decline when the mouth was occluded with a surgical

585    face mask. Our findings show that individuals vary highly in their visual speech benefit and provide

586    new insights into the behavioral relevance of neural speech tracking.

587    *Neural responses to audiovisual speech*

588    Similar to Brodbeck, Hong, et al. (2018), neural responses to acoustic features in the two-speaker

589    paradigm were generally weaker. The TRFs to lip movements showed an opposite pattern, with

590    an enhanced magnitude in the multi speaker condition (Figure 2A), and with substantially later

591    peaks compared to the TRF to acoustic features. This is in line with Bourguignon et al. (2020),

592    where initial TRF peaks at 115 and 159 ms were shown from two significant sources, overlapping

593    with our involved parietal and occipital sources (see Figure 1C). However, the TRFs in their work

594    were modeled to lip movements from silent videos, which precludes a comparison between

595    different listening situations. Our findings also strengthen the argument that TRFs to visual speech

596    are qualitatively different from TRFs to acoustic speech (for coherence, see Park et al., 2016),

597    despite the high intercorrelation of speech features (Chandrasekaran et al., 2009).

598    *Neural tracking of audiovisual speech*

599    Based on the source-localized neural tracking, we determined fROIs via a data-driven approach

600    – separately for the acoustic features and lip movements (see Figure 1C). The fROIs for the

601    acoustic speech features involved sources along temporal, parietal and posterior frontal regions,

602    covering regions that are related to speech perception (Franken et al., 2022). Previous studies

603    source-localized TRFs in audio-only settings, though commonly restricting the analysis to

604    temporal regions (e.g. Brodbeck, Hong, et al., 2018; Kulasingham et al., 2020). The fROIs for the

605    lip movements involved parietal and occipital regions, in line with previous studies that source-

606    localized the neural tracking of lip movements (Aller et al., 2022; Bourguignon et al., 2020;

607    Hauswald et al., 2018). Similar to Park et al. (2016), we also observed neural tracking of lip

608    movements in temporal regions (see Figure S1), but with less involvement of the primary visual

609    cortex and prominent only in the single speaker condition. Due to our approach of defining our

24

610    fROIs based on the multi speaker condition, we removed any involvement of auditory regions in

611    the lip fROIs. In contrast to Park et al. (2016), we did not observe neural tracking of lip movements

612    in motor regions, resulting in no involvement of related sources in the lip fROIs.

613    When analyzing neural speech tracking in the acoustic fROIs, we showed a large effect with

614    enhanced tracking in the single speaker condition compared to the multi speaker condition (Figure

615    2B). We did not find a previous study that showed such a statistical contrast, which could be due

616    to the general focus on neural tracking of attended versus unattended speech, especially to

617    decode auditory attention (e.g. Ciccarelli et al., 2019; Geirnaert et al., 2021; Mirkovic et al., 2015;

618    J. A. O'Sullivan et al., 2015; Schäfer et al., 2018). On a group level, the neural tracking of lip

619    movements showed an enhancement in the multi speaker condition (Figure 2B). When comparing

620    the involved sources of the corresponding lip fROI, we found a medium effect in the left superior

621    parietal cortex. This is well in line with Park et al. (2016), showing an effect in left occipital and

622    parietal cortex when comparing two similar conditions to our design ("AV congruent vs. All

623    congruent"), although after partializing out auditory-related coherence. When we averaged the

624    neural tracking of lip movements, we observed interesting patterns, with participants showing no

625    meaningful neural tracking (i.e. close to zero or negative correlations) when there was one

626    speaker, but when speech became challenging, their neural tracking reached positive values.

627    Notably, this pattern was reversed for some participants, suggesting that not all of them used the

628    lip movements in the same manner. To investigate this further, eye tracking should be used to

629    identify which face regions participants fixated when attending audiovisual speech (e.g. Rennig &

630    Beauchamp, 2018) or to additionally incorporate a recently proposed phenomenon termed "ocular

631    speech tracking" (Gehmacher et al., 2023). Altogether, this is the first time that neural tracking of

632    lip movements has been quantified in the context of TRFs, although with substantially smaller

633    correlations as compared to acoustic speech tracking. Other algorithms, such as ridge regression,

634    could, in principle, yield higher values due to their optimization towards maximizing neural tracking

635    values (for a comparison of algorithms, see Kulasingham & Simon, 2022).

636    *Benefit of lip movements*

637    We first compared the neural tracking of audiovisual speech between single speaker and multi

638    speaker conditions in an isolated manner. Due to the aforementioned inter-correlation of speech

639    features (Chandrasekaran et al., 2009; Daube et al., 2019), this approach could not rule out any

640    acoustic contributions to the neural tracking of lip movements or vice versa. To reveal the unique

641    benefit of lip movements and to incorporate regions that are part of models of audiovisual speech

25

642   perception (Bernstein & Liebenthal, 2014) and multisensory integration (Peelle & Sommers,
643   2015), we combined both fROIs and controlled for acoustic speech features. Within the TRF
644   framework, we provide first evidence that lip movements enhance acoustic-controlled neural
645   speech tracking (Figure 2C). A general enhancement was observed for both single- and multi
646   speaker speech, which is in line with behavioral findings that visual speech features enhance
647   intelligibility under clear speech conditions as well (Blackburn et al., 2019; Stacey et al., 2016).
648   When comparing the two conditions, we observed a large effect, showing a higher benefit of lip
649   movements in the multi speaker condition. Our findings are also well in line with a previous study
650   (Park et al., 2016) that used partial coherence to remove auditory-related contributions, showing
651   higher coherence in a challenging audiovisual speech situation compared to a condition where
652   the audiovisual input was congruent.

653   Analogous to behavioral findings in Aller et al. (2022), the benefit of lip movements showed high
654   interindividual variability (see Figure 2C) and followed a bimodal distribution. Some individuals
655   benefited massively from lip movements, while others showed only a small benefit or none at all.
656   Interestingly, one individual even showed a negative influence when adding lip movements to the
657   acoustic model when there was only one speaker. As soon as speech became challenging, that
658   individual benefited from the lip information. Overall, these findings are in line with the beneficial
659   effects of visual speech when listening is challenging (e.g. Grant & Seitz, 2000; Remez, 2012;
660   Ross et al., 2007; Sumby & Pollack, 1954). Given our moderate sample size, we refrained from
661   conducting further analysis by defining groups of individuals who showed a higher or lower benefit
662   of lip movements. Future studies should include more participants, as well as hearing-impaired
663   populations. A recent study that used neural tracking showed an increased audiovisual speech
664   benefit when speech was noisy (Puschmann et al., 2019). This could also provide a clearer picture
665   of how individuals benefit from lip movements in terms of neural tracking. Previous studies used
666   only the acoustic envelope to investigate the benefit of visual speech features on neural speech
667   tracking (Crosse, Liberto, et al., 2016; Golumbic et al., 2013). Here, we also incorporated lip
668   movements to provide a more complete picture of the unique benefit of visual speech features in
669   audiovisual settings with naturalistic stimuli (Hamilton & Huth, 2020; A. E. O'Sullivan et al., 2019).

670   *Predicting behavior with neural tracking*

671   Our initial analysis of the behavioral measures suggests a higher cognitive demand when speech
672   was challenging (Figure 1B). Participants displayed lower task performance, higher difficulty
673   ratings and lower motivation ratings when more than one speaker was involved (Figure 1B). The

674 influence of challenging speech is also reflected in the findings of neural speech tracking (Figure
675 2B). Building on these results, we used Bayesian multilevel modeling to establish a link between
676 neural speech tracking and behavior (Figure 3). Higher acoustic neural tracking is related to
677 higher task performance, a finding also reported in a study that used vocoded speech (Chen et
678 al., 2023). We also show that higher acoustic neural tracking is related to lower difficulty ratings.
679 This is in line with a study that showed a positive relationship between speech intelligibility ratings
680 and acoustic neural tracking, though using speech-in-noise (Ding & Simon, 2013). Higher
681 motivation ratings were associated with higher acoustic neural tracking – in contrast to Schubert
682 et al. (2023) – showing no relationship between the two measures. We were not able to establish
683 any link between the neural tracking of lip movements and the behavioral measures. It is important
684 to note here that the analyzed neural tracking of lip movements was not yet controlled for speech
685 acoustics (Gillis et al., 2022), which could confound any relationship with behavior. A recent MEG
686 study impressively showed that the neural tracking of acoustic speech features can explain
687 cortical responses to higher-order linguistic features, such as phoneme onsets (Daube et al.,
688 2019), emphasizing the importance of controlling acoustics (see also Gillis et al., 2021).

689 The COVID-19 pandemic established the use of face masks on a global scale (Feng et al., 2020).
690 However, it has been demonstrated that covering the mouth has adverse effects on behavioral
691 measures, such as speech perception (e.g. Rahne et al., 2021). On a neural level, Haider et al.
692 (2022) showed that surgical face masks impair the neural tracking of acoustic and higher-order
693 segmentational speech features. However, the consequences of an absence of visual speech
694 were not analyzed in this study. Here, we establish a relationship between behavioral measures
695 and the individual benefit of visual speech on neural tracking. When the speaker wore a surgical
696 face mask, individuals that benefit more from lip movements displayed lower task performance
697 and higher difficulty ratings. Strikingly, no effect was found when the speaker did not wear a
698 surgical face mask. Overall, our results suggest that individuals who use lip movements more
699 effectively show behavioral deterioration when visual speech is absent. However, further studies
700 with larger sample sizes are needed to disentangle the potential influence of experimental
701 conditions on this relationship, e.g. using Bayesian mediation analysis (Nuijten et al., 2015; Yuan
702 & MacKinnon, 2009).

703 *Conclusion*

704 The current study provides first evidence for the substantial interindividual variability in the neural
705 tracking of lip movements and its relationship to behavior. First, we show that neural responses

27

706 to lip movements are more pronounced when speech is challenging, compared to when speech

707 is clear. We show that lip movements effectively enhance neural speech tracking in brain regions

708 related to audiovisual speech, with high interindividual variability. Furthermore, we demonstrate

709 that this individual visual benefit is behaviorally relevant. Individuals that benefit more from lip

710 movements have a lower task performance and rate the task to be more difficult when the speaker

711 wears a surgical face mask. Remarkably, this relationship is completely absent when the speaker

712 did not wear a mask. Our results provide new insights into the individual differences in the neural

713 tracking of lip movements and offer potential implications for future clinical and audiological

714 settings to objectively assess audiovisual speech perception.

## Acknowledgments

## Author Contributions

725 P.R. analyzed the data, created the figures and wrote the manuscript. M.G. and J.V. analyzed the

726 data and edited the manuscript. N.S. and T.H. provided input on data analysis and edited the

727 manuscript. C.H. designed the study, collected the original dataset and edited the manuscript.

728 A.H. designed the study and edited the manuscript. K.S. edited the manuscript. T.F. supervised

729 the project and edited the manuscript. N.W. acquired the funding, supervised the project and

730 edited the manuscript.

## Conflict of interest statement

732 K.S. is an employee of MED-EL GmbH. All other authors declare no competing interests.

## References

Akalin-Acar, Z., & Gençer, N. G. (2004). An advanced boundary element method (BEM) implementation for the forward problem of electromagnetic source imaging. *Physics in Medicine & Biology*, *49*(21), 5011. https://doi.org/10.1088/0031-9155/49/21/012

Aller, M., Økland, H. S., MacGregor, L. J., Blank, H., & Davis, M. H. (2022). Differential auditory and visual phase-locking are observed during audio-visual benefit and silent lip-reading for speech perception. *Journal of Neuroscience*, *42*(31), 6108–6120. https://doi.org/10.1523/JNEUROSCI.2476-21.2022

Bell, A. J., & Sejnowski, T. J. (1995). An Information-Maximization Approach to Blind Separation and Blind Deconvolution. *Neural Computation*, *7*(6), 1129–1159. https://doi.org/10.1162/neco.1995.7.6.1129

Bernstein, L. E., & Liebenthal, E. (2014). Neural pathways for visual speech perception. *Frontiers in Neuroscience*, *8*. https://www.frontiersin.org/article/10.3389/fnins.2014.00386

Besl, P. J., & McKay, N. D. (1992). A method for registration of 3-D shapes. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, *14*(2), 239–256. https://doi.org/10.1109/34.121791

Betancourt, M. (2018). *A Conceptual Introduction to Hamiltonian Monte Carlo* (arXiv:1701.02434). arXiv. https://doi.org/10.48550/arXiv.1701.02434

Biesmans, W., Das, N., Francart, T., & Bertrand, A. (2017). Auditory-Inspired Speech Envelope Extraction Methods for Improved EEG-Based Auditory Attention Detection in a Cocktail Party Scenario. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, *25*(5), 402–412. https://doi.org/10.1109/TNSRE.2016.2571900

756   Blackburn, C. L., Kitterick, P. T., Jones, G., Sumner, C. J., & Stacey, P. C. (2019). Visual

757         Speech Benefit in Clear and Degraded Speech Depends on the Auditory Intelligibility of

758         the Talker and the Number of Background Talkers. *Trends in Hearing*, *23*.

759         https://doi.org/10.1177/2331216519837866

760   Boersma, P. (2001). Praat, a system for doing phonetics by computer. *Glot. Int.*, *5*(9), 341–345.

761   Bourguignon, M., Baart, M., Kapnoula, E. C., & Molinaro, N. (2020). Lip-Reading Enables the

762         Brain to Synthesize Auditory Features of Unknown Silent Speech. *Journal of*

763         *Neuroscience*, *40*(5), 1053–1065. https://doi.org/10.1523/JNEUROSCI.1101-19.2019

764   Brainard, D. H. (1997). The Psychophysics Toolbox. *Spatial Vision*, *10*(4), 433–436.

765         https://doi.org/10.1163/156856897X00357

766   Brodbeck, C., Das, P., Gillis, M., Kulasingham, J. P., Bhattasali, S., Gaston, P., Resnik, P., &

767         Simon, J. Z. (2022). *Eelbrain: A Python toolkit for time-continuous analysis with temporal*

768         *response functions*. bioRxiv. https://doi.org/10.1101/2021.08.01.454687

769   Brodbeck, C., Hong, L. E., & Simon, J. Z. (2018). Rapid Transformation from Auditory to

770         Linguistic Representations of Continuous Speech. *Current Biology*, *28*(24), 3976-

771         3983.e5. https://doi.org/10.1016/j.cub.2018.10.042

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

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

774         speakers. *PLOS Biology*, *18*(10), e3000883.

775         https://doi.org/10.1371/journal.pbio.3000883

776   Brodbeck, C., Presacco, A., & Simon, J. Z. (2018). Neural source dynamics of brain responses

777         to continuous stimuli: Speech processing from acoustics to comprehension.

778         *NeuroImage*, *172*, 162–174. https://doi.org/10.1016/j.neuroimage.2018.01.042

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

780         *Physiology*, *18*, 25–31. https://doi.org/10.1016/j.cophys.2020.07.014

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

782   Electrophysiological Correlates of Semantic Dissimilarity Reflect the Comprehension of

783   Natural, Narrative Speech. *Current Biology*, *28*(5), 803-809.e3.

784   https://doi.org/10.1016/j.cub.2018.01.080

785 Bröhl, F., Keitel, A., & Kayser, C. (2022). MEG Activity in Visual and Auditory Cortices

786   Represents Acoustic Speech-Related Information during Silent Lip Reading. *ENeuro*,

787   *9*(3). https://doi.org/10.1523/ENEURO.0209-22.2022

788 Brown, V. A., Van Engen, K. J., & Peelle, J. E. (2021). Face mask type affects audiovisual

789   speech intelligibility and subjective listening effort in young and older adults. *Cognitive*

790   *Research: Principles and Implications*, *6*(1), 49. https://doi.org/10.1186/s41235-021-

791   00314-0

792 Bürkner, P.-C. (2017). brms: An R Package for Bayesian Multilevel Models Using Stan. *Journal*

793   *of Statistical Software*, *80*, 1–28. https://doi.org/10.18637/jss.v080.i01

794 Bürkner, P.-C. (2018). Advanced Bayesian Multilevel Modeling with the R Package brms. *The R*

795   *Journal*, *10*(1), 395–411.

796 Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M., Brubaker, M.,

797   Guo, J., Li, P., & Riddell, A. (2017). Stan: A Probabilistic Programming Language.

798   *Journal of Statistical Software*, *76*, 1–32. https://doi.org/10.18637/jss.v076.i01

799 Chalas, N., Daube, C., Kluger, D. S., Abbasi, O., Nitsch, R., & Gross, J. (2022). Multivariate

800   analysis of speech envelope tracking reveals coupling beyond auditory cortex.

801   *NeuroImage*, *258*, 119395. https://doi.org/10.1016/j.neuroimage.2022.119395

802 Chandrasekaran, C., Trubanova, A., Stillittano, S., Caplier, A., & Ghazanfar, A. A. (2009). The

803   Natural Statistics of Audiovisual Speech. *PLOS Computational Biology*, *5*(7), e1000436.

804   https://doi.org/10.1371/journal.pcbi.1000436

805   Chen, Y.-P., Schmidt, F., Keitel, A., Rösch, S., Hauswald, A., & Weisz, N. (2023). Speech

806        intelligibility changes the temporal evolution of neural speech tracking. *NeuroImage*, *268*,

807        119894. https://doi.org/10.1016/j.neuroimage.2023.119894

808   Cherry, E. C. (1953). Some Experiments on the Recognition of Speech, with One and with Two

809        Ears. *The Journal of the Acoustical Society of America*, *25*(5), 975–979.

810        https://doi.org/10.1121/1.1907229

811   Chu, D. K., Akl, E. A., Duda, S., Solo, K., Yaacoub, S., Schünemann, H. J., Chu, D. K., Akl, E.

812        A., El-harakeh, A., Bognanni, A., Lotfi, T., Loeb, M., Hajizadeh, A., Bak, A., Izcovich, A.,

813        Cuello-Garcia, C. A., Chen, C., Harris, D. J., Borowiack, E., … Schünemann, H. J.

814        (2020). Physical distancing, face masks, and eye protection to prevent person-to-person

815        transmission of SARS-CoV-2 and COVID-19: A systematic review and meta-analysis.

816        *The Lancet*, *395*(10242), 1973–1987. https://doi.org/10.1016/S0140-6736(20)31142-9

817   Ciccarelli, G., Nolan, M., Perricone, J., Calamia, P. T., Haro, S., O'Sullivan, J., Mesgarani, N.,

818        Quatieri, T. F., & Smalt, C. J. (2019). Comparison of Two-Talker Attention Decoding

819        from EEG with Nonlinear Neural Networks and Linear Methods. *Scientific Reports*, *9*(1),

820        11538. https://doi.org/10.1038/s41598-019-47795-0

821   Corey, D. M., Dunlap, W. P., & Burke, M. J. (1998). Averaging Correlations: Expected Values

822        and Bias in Combined Pearson rs and Fisher's z Transformations. *The Journal of*

823        *General Psychology*, *125*(3), 245–261. https://doi.org/10.1080/00221309809595548

824   Crosse, M. J., Butler, J. S., & Lalor, E. C. (2015). Congruent Visual Speech Enhances Cortical

825        Entrainment to Continuous Auditory Speech in Noise-Free Conditions. *Journal of*

826        *Neuroscience*, *35*(42), 14195–14204. https://doi.org/10.1523/JNEUROSCI.1829-15.2015

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

828        response function (mTRF) toolbox: A Matlab toolbox for relating neural signals to

829        continuous stimuli. *Frontiers in Human Neuroscience*, *10*.

830        https://doi.org/10.3389/fnhum.2016.00604

831    Crosse, M. J., Liberto, G. M. D., & Lalor, E. C. (2016). Eye Can Hear Clearly Now: Inverse

832        Effectiveness in Natural Audiovisual Speech Processing Relies on Long-Term

833        Crossmodal Temporal Integration. *Journal of Neuroscience*, *36*(38), 9888–9895.

834        https://doi.org/10.1523/JNEUROSCI.1396-16.2016

835    Crosse, M. J., Zuk, N. J., Di Liberto, G. M., Nidiffer, A. R., Molholm, S., & Lalor, E. C. (2021).

836        Linear Modeling of Neurophysiological Responses to Speech and Other Continuous

837        Stimuli: Methodological Considerations for Applied Research. *Frontiers in Neuroscience*,

838        *15*, 1350. https://doi.org/10.3389/fnins.2021.705621

839    Das, P., Brodbeck, C., Simon, J. Z., & Babadi, B. (2020). Neuro-current response functions: A

840        unified approach to MEG source analysis under the continuous stimuli paradigm.

841        *NeuroImage*, *211*, 116528. https://doi.org/10.1016/j.neuroimage.2020.116528

842    Daube, C., Ince, R. A. A., & Gross, J. (2019). Simple Acoustic Features Can Explain Phoneme-

843        Based Predictions of Cortical Responses to Speech. *Current Biology*, *29*(12), 1924-

844        1937.e9. https://doi.org/10.1016/j.cub.2019.04.067

845    David, S. V., Mesgarani, N., & Shamma, S. A. (2007). Estimating sparse spectro-temporal

846        receptive fields with natural stimuli. *Network: Computation in Neural Systems*, *18*(3),

847        191–212. https://doi.org/10.1080/09548980701609235

848    de Jong, N. H., & Wempe, T. (2009). Praat script to detect syllable nuclei and measure speech

849        rate automatically. *Behavior Research Methods*, *41*(2), 385–390.

850        https://doi.org/10.3758/BRM.41.2.385

851    Delorme, A., & Makeig, S. (2004). EEGLAB: An open source toolbox for analysis of single-trial

852        EEG dynamics including independent component analysis. *Journal of Neuroscience*

853        *Methods*, *134*(1), 9–21. https://doi.org/10.1016/j.jneumeth.2003.10.009

854    Desikan, R. S., Ségonne, F., Fischl, B., Quinn, B. T., Dickerson, B. C., Blacker, D., Buckner, R.

855        L., Dale, A. M., Maguire, R. P., Hyman, B. T., Albert, M. S., & Killiany, R. J. (2006). An

856        automated labeling system for subdividing the human cerebral cortex on MRI scans into

857        gyral based regions of interest. *NeuroImage*, *31*(3), 968–980.

858        https://doi.org/10.1016/j.neuroimage.2006.01.021

859    Di Liberto, G. M., O'Sullivan, J. A., & Lalor, E. C. (2015). Low-Frequency Cortical Entrainment to

860        Speech Reflects Phoneme-Level Processing. *Current Biology*, *25*(19), 2457–2465.

861        https://doi.org/10.1016/j.cub.2015.08.030

862    Ding, N., & Simon, J. Z. (2013). Adaptive Temporal Encoding Leads to a Background-

863        Insensitive Cortical Representation of Speech. *Journal of Neuroscience*, *33*(13), 5728–

864        5735. https://doi.org/10.1523/JNEUROSCI.5297-12.2013

865    Erber, N. P. (1975). Auditory-Visual Perception of Speech. *Journal of Speech and Hearing*

866        *Disorders*, *40*(4), 481–492. https://doi.org/10.1044/jshd.4004.481

867    Feng, S., Shen, C., Xia, N., Song, W., Fan, M., & Cowling, B. J. (2020). Rational use of face

868        masks in the COVID-19 pandemic. *The Lancet Respiratory Medicine*, *8*(5), 434–436.

869        https://doi.org/10.1016/S2213-2600(20)30134-X

870    Fischl, B. (2012). FreeSurfer. *NeuroImage*, *62*(2), 774–781.

871        https://doi.org/10.1016/j.neuroimage.2012.01.021

872    Fishbach, A., Nelken, I., & Yeshurun, Y. (2001). Auditory Edge Detection: A Neural Model for

873        Physiological and Psychoacoustical Responses to Amplitude Transients. *Journal of*

874        *Neurophysiology*, *85*(6), 2303–2323. https://doi.org/10.1152/jn.2001.85.6.2303

875    Franken, M. K., Liu, B. C., & Ostry, D. J. (2022). Towards a somatosensory theory of speech

876        perception. *Journal of Neurophysiology*, *128*(6), 1683–1695.

877        https://doi.org/10.1152/jn.00381.2022

878    Freeman, J. B., & Dale, R. (2013). Assessing bimodality to detect the presence of a dual

879         cognitive process. *Behavior Research Methods*, *45*(1), 83–97.

880         https://doi.org/10.3758/s13428-012-0225-x

881    Gehmacher, Q., Schubert, J., Schmidt, F., Hartmann, T., Reisinger, P., Rösch, S., Schwarz, K.,

882         Popov, T., Chait, M., & Weisz, N. (2023). *Eye movements track prioritized auditory*

883         *features in selective attention to natural speech*. bioRxiv.

884         https://doi.org/10.1101/2023.01.23.525171

885    Geirnaert, S., Vandecappelle, S., Alickovic, E., de Cheveigne, A., Lalor, E., Meyer, B. T., Miran,

886         S., Francart, T., & Bertrand, A. (2021). Electroencephalography-Based Auditory

887         Attention Decoding: Toward Neurosteered Hearing Devices. *IEEE Signal Processing*

888         *Magazine*, *38*(4), 89–102. https://doi.org/10.1109/MSP.2021.3075932

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

890         diagnostic tool to assess the auditory pathway. *Hearing Research*, *426*, 108607.

891         https://doi.org/10.1016/j.heares.2022.108607

892    Gillis, M., Vanthornhout, J., Simon, J. Z., Francart, T., & Brodbeck, C. (2021). Neural Markers of

893         Speech Comprehension: Measuring EEG Tracking of Linguistic Speech

894         Representations, Controlling the Speech Acoustics. *Journal of Neuroscience*, *41*(50),

895         10316–10329. https://doi.org/10.1523/JNEUROSCI.0812-21.2021

896    Golumbic, E. Z., Cogan, G. B., Schroeder, C. E., & Poeppel, D. (2013). Visual Input Enhances

897         Selective Speech Envelope Tracking in Auditory Cortex at a "Cocktail Party." *Journal of*

898         *Neuroscience*, *33*(4), 1417–1426. https://doi.org/10.1523/JNEUROSCI.3675-12.2013

899    Gramfort, A., Luessi, M., Larson, E., Engemann, D. A., Strohmeier, D., Brodbeck, C.,

900         Parkkonen, L., & Hämäläinen, M. S. (2014). MNE software for processing MEG and

901         EEG data. *NeuroImage*, *86*, 446–460. https://doi.org/10.1016/j.neuroimage.2013.10.027

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

903          M., Brooks, T., Parkkonen, L., & Hämäläinen, M. (2013). MEG and EEG data analysis

904          with MNE-Python. *Frontiers in Neuroscience*, *7*.

905          https://www.frontiersin.org/articles/10.3389/fnins.2013.00267

906    Grant, K. W., & Seitz, P.-F. (2000). The use of visible speech cues for improving auditory

907          detection of spoken sentences. *The Journal of the Acoustical Society of America*,

908          *108*(3), 1197–1208. https://doi.org/10.1121/1.1288668

909    Haider, C. L., Suess, N., Hauswald, A., Park, H., & Weisz, N. (2022). Masking of the mouth area

910          impairs reconstruction of acoustic speech features and higher-level segmentational

911          features in the presence of a distractor speaker. *NeuroImage*, *252*, 119044.

912          https://doi.org/10.1016/j.neuroimage.2022.119044

913    Hämäläinen, M. S., & Ilmoniemi, R. J. (1994). Interpreting magnetic fields of the brain: Minimum

914          norm estimates. *Medical & Biological Engineering & Computing*, *32*(1), 35–42.

915          https://doi.org/10.1007/BF02512476

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

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

918          https://doi.org/10.1080/23273798.2018.1499946

919    Hartmann, T., & Weisz, N. (2020). An introduction to the Objective Psychophysics Toolbox

920          (o_ptb). *Frontiers in Psychology*, *11*. https://doi.org/10.3389/fpsyg.2020.585437

921    Hauswald, A., Lithari, C., Collignon, O., Leonardelli, E., & Weisz, N. (2018). A Visual Cortical

922          Network for Deriving Phonological Information from Intelligible Lip Movements. *Current*

923          *Biology*, *28*(9), 1453-1459.e3. https://doi.org/10.1016/j.cub.2018.03.044

924    Heeris, J. (2013). *Gammatone Filterbank Toolkit*. https://github.com/detly/gammatone

925    Houck, J. M., & Claus, E. D. (2020). A comparison of automated and manual co-registration for

926          magnetoencephalography. *PLOS ONE*, *15*(4), e0232100.

927     https://doi.org/10.1371/journal.pone.0232100

928     Kay, M. (2022). *ggdist: Visualizations of distributions and uncertainty*. Zenodo.

929     https://doi.org/10.5281/zenodo.6862765

930     Keitel, A., Gross, J., & Kayser, C. (2018). Perceptually relevant speech tracking in auditory and

931     motor cortex reflects distinct linguistic features. *PLOS Biology*, *16*(3), e2004473.

932     https://doi.org/10.1371/journal.pbio.2004473

933     King, B. M., Rosopa, P. J., & Minium, E. W. (2018). *Statistical Reasoning in the Behavioral*

934     *Sciences* (7th Edition). John Wiley & Sons.

935     Kleiner, M., Brainard, D., & Pelli, D. (2007). What's new in Psychtoolbox-3? *Perception*, *36*, 1–

936     16.

937     Kulasingham, J. P., Brodbeck, C., Presacco, A., Kuchinsky, S. E., Anderson, S., & Simon, J. Z.

938     (2020). High gamma cortical processing of continuous speech in younger and older

939     listeners. *NeuroImage*, *222*, 117291. https://doi.org/10.1016/j.neuroimage.2020.117291

940     Kulasingham, J. P., & Simon, J. Z. (2022). Algorithms for Estimating Time-Locked Neural

941     Response Components in Cortical Processing of Continuous Speech. *IEEE*

942     *Transactions on Biomedical Engineering*, 1–9.

943     https://doi.org/10.1109/TBME.2022.3185005

944     Lalor, E. C., Pearlmutter, B. A., Reilly, R. B., McDarby, G., & Foxe, J. J. (2006). The VESPA: A

945     method for the rapid estimation of a visual evoked potential. *NeuroImage*, *32*(4), 1549–

946     1561. https://doi.org/10.1016/j.neuroimage.2006.05.054

947     Lalor, E. C., Power, A. J., Reilly, R. B., & Foxe, J. J. (2009). Resolving Precise Temporal

948     Processing Properties of the Auditory System Using Continuous Stimuli. *Journal of*

949     *Neurophysiology*, *102*(1), 349–359. https://doi.org/10.1152/jn.90896.2008

950  Lin, F.-H., Witzel, T., Ahlfors, S. P., Stufflebeam, S. M., Belliveau, J. W., & Hämäläinen, M. S.

951      (2006). Assessing and improving the spatial accuracy in MEG source localization by

952      depth-weighted minimum-norm estimates. *NeuroImage*, *31*(1), 160–171.

953      https://doi.org/10.1016/j.neuroimage.2005.11.054

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

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

956      https://doi.org/10.1016/j.jneumeth.2007.03.024

957  McDermott, J. H. (2009). The cocktail party problem. *Current Biology*, *19*(22), R1024–R1027.

958      https://doi.org/10.1016/j.cub.2009.09.005

959  McElreath, R. (2020). *Statistical Rethinking: A Bayesian Course with Examples in R and STAN*

960      (2nd ed.). Chapman and Hall/CRC. https://doi.org/10.1201/9780429029608

961  Meredith, M. A., & Stein, B. E. (1983). Interactions Among Converging Sensory Inputs in the

962      Superior Colliculus. *Science*, *221*(4608), 389–391.

963      https://doi.org/10.1126/science.6867718

964  Mirkovic, B., Debener, S., Jaeger, M., & Vos, M. D. (2015). Decoding the attended speech

965      stream with multi-channel EEG: Implications for online, daily-life applications. *Journal of*

966      *Neural Engineering*, *12*(4), 046007. https://doi.org/10.1088/1741-2560/12/4/046007

967  Nidiffer, A. R., Cao, C. Z., O'Sullivan, A., & Lalor, E. C. (2021). *A linguistic representation in the*

968      *visual system underlies successful lipreading*. bioRxiv.

969      https://www.biorxiv.org/content/10.1101/2021.02.09.430299v1

970  Nieto-Castanon, A., Ghosh, S. S., Tourville, J. A., & Guenther, F. H. (2003). Region of interest

971      based analysis of functional imaging data. *NeuroImage*, *19*(4), 1303–1316.

972      https://doi.org/10.1016/S1053-8119(03)00188-5

973  Nuijten, M. B., Wetzels, R., Matzke, D., Dolan, C. V., & Wagenmakers, E.-J. (2015). A default

974      Bayesian hypothesis test for mediation. *Behavior Research Methods*, *47*(1), 85–97.

975          https://doi.org/10.3758/s13428-014-0470-2

976    Obleser, J., & Kayser, C. (2019). Neural Entrainment and Attentional Selection in the Listening

977          Brain. *Trends in Cognitive Sciences*, *23*(11), 913–926.

978          https://doi.org/10.1016/j.tics.2019.08.004

979    Oostenveld, R., Fries, P., Maris, E., & Schoffelen, J.-M. (2011). *FieldTrip: Open Source*

980          *Software for Advanced Analysis of MEG, EEG, and Invasive Electrophysiological Data*

981          [Research article]. Computational Intelligence and Neuroscience.

982          https://doi.org/10.1155/2011/156869

983    O'Sullivan, A. E., Lim, C. Y., & Lalor, E. C. (2019). Look at me when I'm talking to you: Selective

984          attention at a multisensory cocktail party can be decoded using stimulus reconstruction

985          and alpha power modulations. *European Journal of Neuroscience*, *50*(8), 3282–3295.

986          https://doi.org/10.1111/ejn.14425

987    O'Sullivan, J. A., Power, A. J., Mesgarani, N., Rajaram, S., Foxe, J. J., Shinn-Cunningham, B.

988          G., Slaney, M., Shamma, S. A., & Lalor, E. C. (2015). Attentional Selection in a Cocktail

989          Party Environment Can Be Decoded from Single-Trial EEG. *Cerebral Cortex*, *25*(7),

990          1697–1706. https://doi.org/10.1093/cercor/bht355

991    Park, H., Kayser, C., Thut, G., & Gross, J. (2016). Lip movements entrain the observers' low-

992          frequency brain oscillations to facilitate speech intelligibility. *ELife*, *5*, e14521.

993          https://doi.org/10.7554/eLife.14521

994    Peelle, J. E., & Sommers, M. S. (2015). Prediction and constraint in audiovisual speech

995          perception. *Cortex*, *68*, 169–181. https://doi.org/10.1016/j.cortex.2015.03.006

996    Pelli, D. G. (1997). The VideoToolbox software for visual psychophysics: Transforming numbers

997          into movies. *Spatial Vision*, *10*(4), 437–442. https://doi.org/10.1163/156856897X00366

998    Pfister, R., Schwarz, K., Janczyk, M., Dale, R., & Freeman, J. (2013). Good things peak in pairs:

999          A note on the bimodality coefficient. *Frontiers in Psychology*, *4*.

1000        https://www.frontiersin.org/articles/10.3389/fpsyg.2013.00700

1001   Puschmann, S., Daeglau, M., Stropahl, M., Mirkovic, B., Rosemann, S., Thiel, C. M., &

1002        Debener, S. (2019). Hearing-impaired listeners show increased audiovisual benefit when

1003        listening to speech in noise. *NeuroImage*, *196*, 261–268.

1004        https://doi.org/10.1016/j.neuroimage.2019.04.017

1005   R Core Team. (2022). *R: A Language and Environment for Statistical Computing*. R Foundation

1006        for Statistical Computing. https://www.R-project.org/

1007   Rahne, T., Fröhlich, L., Plontke, S., & Wagner, L. (2021). Influence of surgical and N95 face

1008        masks on speech perception and listening effort in noise. *PLOS ONE*, *16*(7), e0253874.

1009        https://doi.org/10.1371/journal.pone.0253874

1010   Remez, R. E. (2012). Three puzzles of multimodal speech perception. In E. Vatikiotis-Bateson,

1011        G. Bailly, & P. Perrier (Eds.), *Audiovisual Speech Processing* (pp. 4–20). Cambridge

1012        University Press. https://doi.org/10.1017/CBO9780511843891.003

1013   Rennig, J., & Beauchamp, M. S. (2018). Free viewing of talking faces reveals mouth and eye

1014        preferring regions of the human superior temporal sulcus. *NeuroImage*, *183*, 25–36.

1015        https://doi.org/10.1016/j.neuroimage.2018.08.008

1016   Ross, L. A., Molholm, S., Butler, J. S., Bene, V. A. D., & Foxe, J. J. (2022). Neural correlates of

1017        multisensory enhancement in audiovisual narrative speech perception: A fMRI

1018        investigation. *NeuroImage*, *263*, 119598.

1019        https://doi.org/10.1016/j.neuroimage.2022.119598

1020   Ross, L. A., Saint-Amour, D., Leavitt, V. M., Javitt, D. C., & Foxe, J. J. (2007). Do You See What

1021        I Am Saying? Exploring Visual Enhancement of Speech Comprehension in Noisy

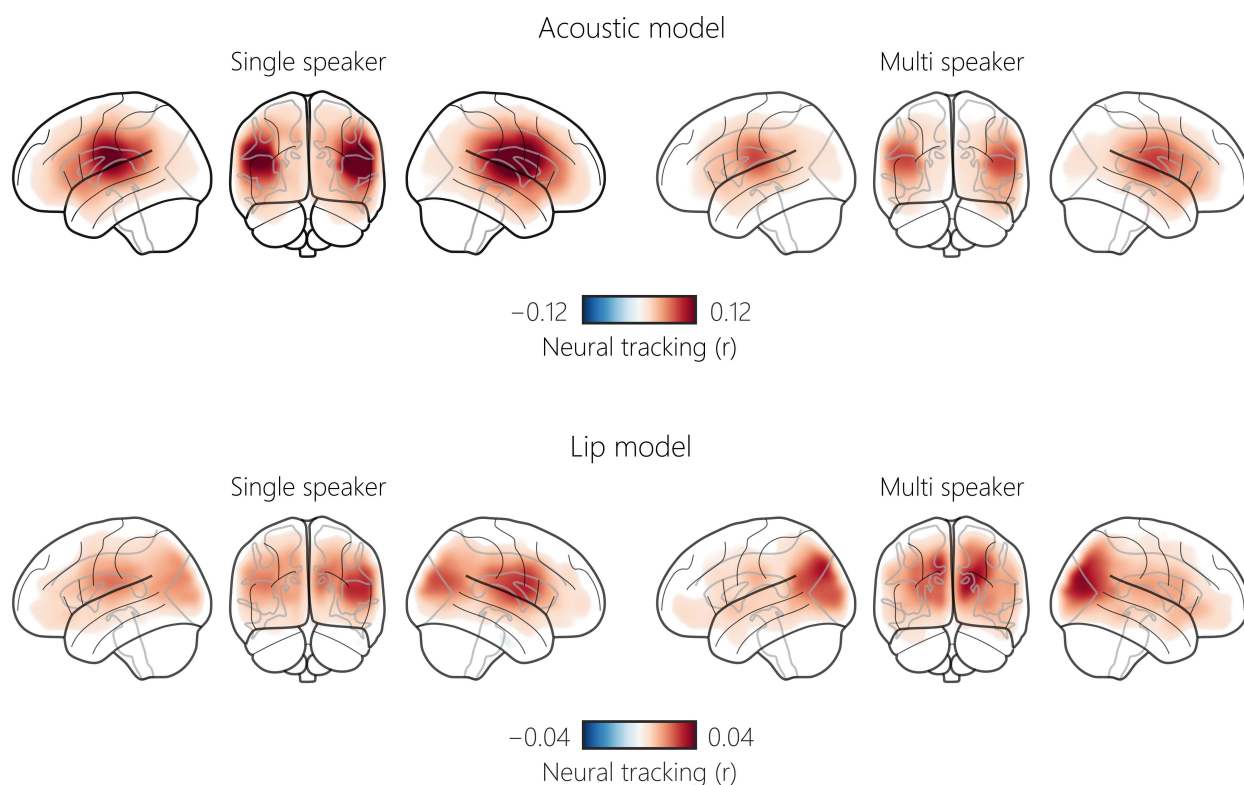1022        Environments. *Cerebral Cortex*, *17*(5), 1147–1153. https://doi.org/10.1093/cercor/bhl024

1023    Schäfer, P. J., Corona-Strauss, F. I., Hannemann, R., Hillyard, S. A., & Strauss, D. J. (2018).

1024    Testing the Limits of the Stimulus Reconstruction Approach: Auditory Attention Decoding

1025    in a Four-Speaker Free Field Environment. *Trends in Hearing*, *22*.

1026    https://doi.org/10.1177/2331216518816600

1027    Schmitt, R., Meyer, M., & Giroud, N. (2022). Better speech-in-noise comprehension is

1028    associated with enhanced neural speech tracking in older adults with hearing

1029    impairment. *Cortex*, *151*, 133–146. https://doi.org/10.1016/j.cortex.2022.02.017

1030    Schubert, J., Schmidt, F., Gehmacher, Q., Bresgen, A., & Weisz, N. (2023). Cortical speech

1031    tracking is related to individual prediction tendencies. *Cerebral Cortex*, bhac528.

1032    https://doi.org/10.1093/cercor/bhac528

1033    Slaney, M. (1998). Auditory toolbox. *Interval Research Corporation*, *10*(1998), 1194.

1034    Smith, S. M., & Nichols, T. E. (2009). Threshold-free cluster enhancement: Addressing

1035    problems of smoothing, threshold dependence and localisation in cluster inference.

1036    *NeuroImage*, *44*(1), 83–98. https://doi.org/10.1016/j.neuroimage.2008.03.061

1037    Stacey, P. C., Kitterick, P. T., Morris, S. D., & Sumner, C. J. (2016). The contribution of visual

1038    information to the perception of speech in noise with and without informative temporal

1039    fine structure. *Hearing Research*, *336*, 17–28.

1040    https://doi.org/10.1016/j.heares.2016.04.002

1041    Suess, N., Hauswald, A., Reisinger, P., Rösch, S., Keitel, A., & Weisz, N. (2022). Cortical

1042    Tracking of Formant Modulations Derived from Silently Presented Lip Movements and

1043    Its Decline with Age. *Cerebral Cortex*, bhab518. https://doi.org/10.1093/cercor/bhab518

1044    Suess, N., Hauswald, A., Zehentner, V., Depireux, J., Herzog, G., Rösch, S., & Weisz, N.

1045    (2022). Influence of linguistic properties and hearing impairment on visual speech

1046    perception skills in the German language. *PLOS ONE*, *17*(9), e0275585.

1047    https://doi.org/10.1371/journal.pone.0275585

1048 Sumby, W. H., & Pollack, I. (1954). Visual Contribution to Speech Intelligibility in Noise. *The*

1049 *Journal of the Acoustical Society of America*, *26*(2), 212–215.

1050 https://doi.org/10.1121/1.1907309

1051 Summerfield, Q., Bruce, V., Cowey, A., Ellis, A. W., & Perrett, D. I. (1992). Lipreading and

1052 audio-visual speech perception. *Philosophical Transactions of the Royal Society of*

1053 *London. Series B: Biological Sciences*, *335*(1273), 71–78.

1054 https://doi.org/10.1098/rstb.1992.0009

1055 Suñer, C., Coma, E., Ouchi, D., Hermosilla, E., Baro, B., Rodríguez-Arias, M. À., Puig, J.,

1056 Clotet, B., Medina, M., & Mitjà, O. (2022). Association between two mass-gathering

1057 outdoor events and incidence of SARS-CoV-2 infections during the fifth wave of COVID-

1058 19 in north-east Spain: A population-based control-matched analysis. *The Lancet*

1059 *Regional Health - Europe*, *15*, 100337. https://doi.org/10.1016/j.lanepe.2022.100337

1060 Taulu, S., & Kajola, M. (2005). Presentation of electromagnetic multichannel data: The signal

1061 space separation method. *Journal of Applied Physics*, *97*(12), 124905.

1062 https://doi.org/10.1063/1.1935742

1063 Taulu, S., & Simola, J. (2006). Spatiotemporal signal space separation method for rejecting

1064 nearby interference in MEG measurements. *Physics in Medicine & Biology*, *51*(7), 1759.

1065 https://doi.org/10.1088/0031-9155/51/7/008

1066 Thorpe, S., Fize, D., & Marlot, C. (1996). Speed of processing in the human visual system.

1067 *Nature*, *381*(6582), 520–522. https://doi.org/10.1038/381520a0

1068 Vallat, R. (2018). Pingouin: Statistics in Python. *Journal of Open Source Software*, *3*(31), 1026.

1069 https://doi.org/10.21105/joss.01026

1070 van de Rijt, L. P. H., Roye, A., Mylanus, E. A. M., van Opstal, A. J., & van Wanrooij, M. M.

1071 (2019). The Principle of Inverse Effectiveness in Audiovisual Speech Perception.

1072 *Frontiers in Human Neuroscience*, *13*.

1073    https://www.frontiersin.org/articles/10.3389/fnhum.2019.00335

1074  VanRullen, R., & Thorpe, S. J. (2001). The Time Course of Visual Processing: From Early

1075    Perception to Decision-Making. *Journal of Cognitive Neuroscience*, *13*(4), 454–461.

1076    https://doi.org/10.1162/08989290152001880

1077  Vanthornhout, J., Decruy, L., Wouters, J., Simon, J. Z., & Francart, T. (2018). Speech

1078    Intelligibility Predicted from Neural Entrainment of the Speech Envelope. *Journal of the*

1079    *Association for Research in Otolaryngology*, *19*(2), 181–191.

1080    https://doi.org/10.1007/s10162-018-0654-z

1081  Vehtari, A., Gabry, J., Magnusson, M., Yao, Y., Bürkner, P.-C., Paananen, T., & Gelman, A.

1082    (2022). *loo: Efficient leave-one-out cross-validation and WAIC for Bayesian models*.

1083    https://mc-stan.org/loo/

1084  Vehtari, A., Gelman, A., & Gabry, J. (2017). Practical Bayesian model evaluation using leave-

1085    one-out cross-validation and WAIC. *Statistics and Computing*, *27*(5), 1413–1432.

1086    https://doi.org/10.1007/s11222-016-9696-4

1087  Vehtari, A., Gelman, A., Simpson, D., Carpenter, B., & Bürkner, P.-C. (2021). Rank-

1088    Normalization, Folding, and Localization: An Improved $\hat{R}$ for Assessing Convergence of

1089    MCMC (with Discussion). *Bayesian Analysis*, *16*(2), 667–718. https://doi.org/10.1214/20-

1090    BA1221

1091  Vehtari, A., Simpson, D., Gelman, A., Yao, Y., & Gabry, J. (2022). *Pareto Smoothed Importance*

1092    *Sampling* (arXiv:1507.02646). arXiv. https://doi.org/10.48550/arXiv.1507.02646

1093  Virtanen, P., Gommers, R., Oliphant, T. E., Haberland, M., Reddy, T., Cournapeau, D.,

1094    Burovski, E., Peterson, P., Weckesser, W., Bright, J., van der Walt, S. J., Brett, M.,

1095    Wilson, J., Millman, K. J., Mayorov, N., Nelson, A. R. J., Jones, E., Kern, R., Larson, E.,

1096    … van Mulbregt, P. (2020). SciPy 1.0: Fundamental algorithms for scientific computing

1097    in Python. *Nature Methods*, *17*(3), Article 3. https://doi.org/10.1038/s41592-019-0686-2

1098 Waskom, M. L. (2021). seaborn: Statistical data visualization. *Journal of Open Source Software*,

1099 *6*(60), 3021. https://doi.org/10.21105/joss.03021

1100 Wickham, H. (2016). *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York.

1101 https://ggplot2.tidyverse.org

1102 Wilkinson, G. N., & Rogers, C. E. (1973). Symbolic Description of Factorial Models for Analysis

1103 of Variance. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, *22*(3),

1104 392–399. https://doi.org/10.2307/2346786

1105 Yuan, Y., & MacKinnon, D. P. (2009). Bayesian mediation analysis. *Psychological Methods*, *14*,

1106 301–322. https://doi.org/10.1037/a0016972

1107 Zhang, L., & Du, Y. (2022). Lip movements enhance speech representations and effective

1108 connectivity in auditory dorsal stream. *NeuroImage*, *257*, 119311.

1109 https://doi.org/10.1016/j.neuroimage.2022.119311

1110    **Supplementary Materials**



**Figure S1.** *Whole-brain neural tracking of the audiovisual speech features.* Neural tracking (*r*) of all sources is shown for the acoustic model (spectrogram and acoustic onsets) and the lip model (lip movements).

|  | b | 89% CI | Conditions $PP_{b>0}$ |
|---|---|---|---|
| Performance (occluded lips) | -0.77 | [-1.13, -0.41]* | 0.07%* |
| Difficulty (occluded lips) | 1.26 | [1.04, 1.49]* | 100%* |
| Motivation (occluded lips) | -0.11 | [-0.27, 0.04] | 11.11% |

1113    **Supplementary Table 1.** *Effects of conditions on behavior when the lips are occluded.* The formula was:

1114    *behavioral measure ~ 1 + conditions + (1 | participant).* *89% CI not including zero and $PP_{b>0}$ below 5.5%

1115    or above 94.5% (i.e. significant effect).

|  | $b$ | 89% CI | $PP_{b>0}$ |
|---|---|---|---|
| Performance | -0.05 | [-0.28, 0.17] | 36.09% |
| Difficulty | 0.04 | [-0.19, 0.28] | 60.86% |
| Motivation | 0.06 | [-0.08, 0.19] | 76.64% |

Benefit of lip movements spans the three value columns above the header row.

1116 **Supplementary Table 2.** *Effects of the benefit of lip movements on behavior when the lips are*
1117 *not occluded.*