



University of Dundee

MEG Activity in Visual and Auditory Cortices Represents Acoustic Speech-Related Information during Silent Lip Reading

Bröhl, Felix; Keitel, Anne; Kayser, Christoph

Published in: eNeuro

DOI: 10.1523/ENEURO.0209-22.2022

Publication date: 2022

Licence: CC BY

Document Version Peer reviewed version

Link to publication in Discovery Research Portal

Citation for published version (APA): Bröhl, F., Keitel, A., & Kayser, C. (2022). MEG Activity in Visual and Auditory Cortices Represents Acoustic Speech-Related Information during Silent Lip Reading. *eNeuro*, *9*(3). https://doi.org/10.1523/ENEURO.0209-22.2022

General rights

Copyright and moral rights for the publications made accessible in Discovery Research Portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

• Users may download and print one copy of any publication from Discovery Research Portal for the purpose of private study or research.

You may not further distribute the material or use it for any profit-making activity or commercial gain.
You may freely distribute the URL identifying the publication in the public portal.

Take down policy If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.



Research Article: New Research | Cognition and Behavior

MEG activity in visual and auditory cortices represents acoustic speech-related information during silent lip reading

https://doi.org/10.1523/ENEURO.0209-22.2022

Cite as: eNeuro 2022; 10.1523/ENEURO.0209-22.2022

Received: 30 May 2022 Accepted: 6 June 2022

This Early Release article has been peer-reviewed and accepted, but has not been through the composition and copyediting processes. The final version may differ slightly in style or formatting and will contain links to any extended data.

Alerts: Sign up at www.eneuro.org/alerts to receive customized email alerts when the fully formatted version of this article is published.

Copyright © 2022 Bröhl et al.

This is an open-access article distributed under the terms of the Creative Commons Attribution 4.0 International license, which permits unrestricted use, distribution and reproduction in any medium provided that the original work is properly attributed.

1 MEG activity in visual and auditory cortices represents acoustic speech-

- 2 related information during silent lip reading
- 3 Short title: Processing of visual speech in the brain
- 4 Felix Bröhl^{1*}, Anne Keitel², Christoph Kayser¹
- 5 ¹Department for Cognitive Neuroscience, Faculty of Biology, Bielefeld University, Universitätsstr. 25, 33615, Bielefeld, Germany
- 6 ²Psychology, University of Dundee, Scrymgeour Building, Dundee DD1 4HN, UK
- 7 * Corresponding author: Felix Bröhl (felix.broehl@uni-bielefeld.de)

8 Highlights

9

- Visual and auditory cortices represent unheard acoustic features during lip reading
- 10 Auditory cortex emphasizes the acoustic envelope
- 11 Visual cortex emphasizes a pitch signature
- 12 Tracking of unheard features in auditory cortex is associated with behavior
- 13 Number of pages: 30
- 14 Number of figures: 5
- 15 Number of tables: 2
- 16 Number of words: Abstract 209, Introduction 749, Discussion 1843

17 Declaration of competing interest

18 We declare no conflict of interest.

19 Acknowledgement

- 20 This work was supported by the UK Biotechnology and Biological Sciences Research Council (BBSRC,
- 21 BB/L027534/1) and the European Research Council (ERC-2014-CoG; grant No 646657)

22 Abstract

23 Speech is an intrinsically multisensory signal and seeing the speaker's lips forms a cornerstone of 24 communication in acoustically impoverished environments. Still, it remains unclear how the brain exploits 25 visual speech for comprehension. Previous work debated whether lip signals are mainly processed along 26 the auditory pathways or whether the visual system directly implements speech-related processes. To 27 probe this, we systematically characterized dynamic representations of multiple acoustic and visual speech-28 derived features in source localized MEG recordings that were obtained while participants listened to 29 speech or viewed silent speech. Using a mutual-information framework we provide a comprehensive 30 assessment of how well temporal and occipital cortices reflect the physically presented signals and unique 31 aspects of acoustic features that were physically absent but may be critical for comprehension. Our results 32 demonstrate that both cortices feature a functionally specific form of multisensory restoration: during lip 33 reading they reflect unheard acoustic features, independent of co-existing representations of the visible lip 34 movements. This restoration emphasizes the unheard pitch signature in occipital cortex and the speech 35 envelope in temporal cortex and is predictive of lip reading performance. These findings suggest that when 36 seeing the speaker's lips, the brain engages both visual and auditory pathways to support comprehension 37 by exploiting multisensory correspondences between lip movements and spectro-temporal acoustic cues.

38 Significance statement

Lip reading is central for speech comprehension in acoustically impoverished environments. Recent studies show that the auditory and visual cortex can represent acoustic speech features from purely visual speech. It is still unclear, however, what information is represented in these cortices and if this phenomenon is related to lip reading comprehension. Using a comprehensive conditional mutual information analysis applied to magnetoencephalographic data, we demonstrate that signatures of acoustic speech arise in both cortices in parallel, even when discounting for the physically presented stimulus. In addition, the auditory but not the visual cortex activity was related to successful lip reading across participants.

46 Keywords

47 Speech entrainment, lip reading, audio-visual, speech tracking, language, MEG

48 **1. Introduction**

49 Speech is an intrinsically multisensory stimulus that can be conveyed via acoustic and visual signals. It 50 remains debated how the brain exploits the information derived from visual speech (Besle et al., 2008; 51 Calvert et al., 1997; Calvert and Campbell, 2003; Grant and Seitz, 2000). One view is that the visual system 52 directly contributes to establishing speech representations (Bernstein et al., 2011; O'Sullivan et al., 2017; 53 Ozker et al., 2018), as oro-facial movements provide temporal information that can be predictive of 54 concurrent acoustic signals and allow mapping visual cues onto phonological representations (Campbell, 2008; Lazard and Giraud, 2017). The visual cortex tracks dynamic lip signals (Park et al., 2016) and, as 55 56 suggested recently, may also directly 'restore' the acoustic envelope of the visually presented speech 57 (Hauswald et al., 2018; Suess et al., 2022). Another view is that visual speech is mainly represented in 58 regions of the auditory pathways, possibly exploiting speech-specific processes of this system. Along this 59 line, a recent study suggested that the early auditory cortex may also be capable of reflecting the unheard 60 acoustic envelope of a spoken narrative (Bourguignon et al., 2020). Importantly, the evidence that visual 61 speech is reflected along both auditory and the visual pathways may not be mutually exclusive, as both may 62 contribute to a supramodal frame of reference for speech (Arnal et al., 2009; Rauschecker, 2012).

63 To probe the respective involvement of visual and auditory cortices in representing visual speech, many 64 previous studies presented syllables or isolated words as stimuli (Calvert et al., 1997; Calvert and Campbell, 65 2003; Pekkola et al., 2005). However, these results come short of how they translate to continuous or 66 natural speech. Furthermore, many studies did not probe a direct link to behavioral performance, leaving it 67 unclear whether potential cerebral representations derived from visual speech are behaviorally relevant 68 (Bourguignon et al., 2020; Ludman et al., 2000; Mégevand et al., 2020). The latter can be particularly 69 challenging given that pure lip reading performance for every-day speech is often low (Altieri et al., 2011; 70 Grant and Seitz, 2000).

71 The present study rests on the assumption that probing the roles of visual and auditory cortices in 72 representing visual speech requires data from a paradigm based on continuous speech with carefully 73 controlled levels of lip reading performance. In previous work we established such a paradigm and 74 collected MEG data from participants during a word recognition task based on syntactically similar 75 sentences that were presented either purely acoustically or purely visually. In the auditory condition 76 participants were presented with the acoustic signal embedded in background noise, while in the visual 77 condition they watched the muted speaker. The individual sentences were constructed from a closed-set of 78 linguistic items with a common syntactic structure, similar to matrix-sentences used in standardized 79 hearing assessment (Hagerman, 1982; Kollmeier et al., 2015). With this we achieved a comparable level of 80 word recognition performance during auditory- and visual-only conditions and verified that this dataset 81 allows linking neural representations of lexical information and speech dynamics to behavior (Keitel et al., 82 2020, 2018).

83 We leverage this paradigm to probe the roles of visual and auditory pathways in representing visual speech 84 and facilitating lip reading performance. This led us to formulate the two following questions: First, 85 whether representations of restored (e.g. unheard acoustic) features are independent of those physically 86 present (e.g. lip movements). Second, we asked whether these representations of restored features are 87 tied to word recognition performance. To be able to compare cerebral signatures of visual and acoustic 88 speech-derived features, we rooted this analysis on the following systematic assessment: we quantified 89 how well source-localized MEG signals track multiple acoustic and visual speech-derived features 90 independently of each other, both when these features are physically present (e.g. the lip contour when 91 watching the speaker) or absent (e.g. the pitch contour when watching the speaker). We focused this 92 analysis on two regions of interest centered on early auditory and visual cortices, which previous studies 93 have implied in supporting lip reading (Bourguignon et al., 2020; Hauswald et al., 2018; Park et al., 2016). 94 The analysis was based on a mutual information approach that has been used in previous work to probe 95 dynamic speech representations and which is well suited to address the statistical dependency between 96 multiple variables (Daube et al., 2019; Keitel et al., 2018). Our results show that both occipital and temporal 97 regions reflect unheard acoustic speech-derived features independently of the physically present lip 98 movements. This 'restoration' of acoustic information in the temporal, but not the occipital, cortex is 99 predictive of word recognition performance across participants.

100 **2. Materials and Methods**

The data analyzed in this study has been collected and analyzed in previous studies (Keitel et al., 2020, 2018). The analyses conducted here pose new questions and provide novel results beyond the previous work.

104 2.1 Participants and data acquisition

105 Data was collected from 20 native British-English speaking participants (9 female, age 23.6 ± 5.8 years mean 106 \pm SD). Due to prominent environmental artefacts in the MEG recordings, data from two participants were 107 excluded from further analysis. Thus, the analyzed data was from 18 participants (7 female, age 24 ± 6.0 108 years mean ± SD). All participants were screened to exclude hearing impairment prior to data collection 109 using the quick hearing check questionnaire (Koike et al., 1994), had normal or corrected-to-normal vision and were all right-handed (Oldfield, 1971). All participants provided written informed consent and received 110 monetary compensation of 10 £/h. The experiment was approved by the College of Science and 111 Engineering, University of Glasgow (approval number 300140078) and conducted in compliance with the 112 113 Declaration of Helsinki.

MEG data was collected using a 248-magnetometer whole-head MEG system (MAGNES 3600 WH, 4-D Neuroimaging) with a sample rate of 1 kHz. Head positions were measured at the beginning and end of each run, using five coils placed on the participants' heads. Coil positions were co-digitized with the participant's head-shape (FASTRAK[®], Polhemus Inc., VT, USA). Participants were seated in an upright position in front of a screen. Visual stimuli were displayed with a DLP projector at 25 frames per second, a resolution of 1280 × 720 pixels, and covered a visual field of 25 × 19 degrees. Acoustic stimuli were transmitted binaurally through plastic earpieces and 370-cm long plastic tubes connected to a sound pressure transducer and were presented in stereo at a sampling rate of 22,050 Hz.

122 2.2 Stimulus material

123 The stimulus material comprised two structurally equivalent sets of 90 unique closed-set English sentences. Specifically, along the idea of matrix-style sentences using in standardized hearing assessment (Hagerman, 124 125 1982; Kollmeier et al., 2015), each sentence was constructed with the same sequence of linguistic 126 elements, the order of which can be described with the following pattern [filler phrase, time phrase, name, 127 verb, numeral, adjective, noun]. One such sentence for example was 'I forgot to mention (filler phrase), last 128 Thursday morning (time phrase) Mary (name) obtained (verb) four (numeral) beautiful (adjective) journals 129 (noun)'. For each element, a list of 18 different options was created and sentences were constructed so 130 that each single element was repeated ten times. Sentence elements were randomly combined within each set of 90 sentences. This procedure yielded 180 structurally similar but distinct sentences. To measure 131 132 word recognition performance for each sentence, a target word was defined in each sentence: either the 133 adjective (first set of sentences) or the numeral (second set). Sentences lasted on average 5.4 ± 0.4 s (mean 134 \pm SD, ranging from 4.6 s to 6.5 s) and lasted a total of approximately 22 minutes. The speech material was 135 spoken by a male British actor, who was tasked to speak clearly and naturally and to move as little as 136 possible while speaking to assure that the lips center stayed at the same place in each video frame. 137 Audiovisual recordings were gathered with a high-performance camcorder (Sony PMW-EX1) and an 138 external microphone in a sound attenuating booth.

139 Participants were presented with audio-only (A-only), audiovisual or visual-only (V-only) speech material in 140 three conditions (Keitel et al., 2020). However, for the present analysis we only focus on the A-only and V-141 only conditions, as in these one can best dissociate visual- and auditory-related speech representations 142 given that only one physical stimulus was present. Furthermore, during the AV condition word recognition 143 performance was near-ceiling (Keitel et al., 2020), making it difficult to link cerebral and behavioral data. 144 Because performance would have been at ceiling with clear speech in the A-only condition, the acoustic 145 speech was embedded in environmental noise. This noise for each trial was generated by randomly 146 selecting 50 individual sounds from a set of sounds recorded from natural, everyday sources or scenes (e.g. 147 car horns, talking people, traffic). These sounds were then added together to create a distracting noise 148 scene for the duration of each trial. For each participant the individual noise level was further adjusted, as 149 described previously (Keitel et al., 2020). This resulted in an average performance of approximately 70% 150 correct for both A-only and V-only conditions and allowed us to dissociate between correct and incorrect151 word recognition.

152 2.3 Experimental Design

153 Each participant was presented with each of the 180 sentences in three conditions (A-only, V-only and AV). 154 The order of the conditions was fixed for all participants as A-only, AV and then V-only. This order exposed 155 the participants to the stimuli twice before the lip reading task, which helped to increase performance and 156 render it comparable to the A-only task. Each condition was divided into 4 blocks of 45 sentences each, 157 with two blocks being 'adjective' and two 'number' blocks. For each participant, the order of sentences within each block was randomized. The first sentence of each block was a 'dummy' trial that was 158 159 subsequently excluded from analysis. During each trial, participants either fixated a dot (in A condition) or a small cross overlaid onto the mouth of the speaker's face (in V condition). In the A condition, each sentence 160 161 was presented as the respective audio recording, i.e. the spoken sentence, together with the background 162 noise. In the V condition, only the video of the speaker's face was presented clearly and no sound was 163 present. After each trial, four words were presented as response options (either four adjectives or four 164 written numbers) on the screen and participants had to indicate using a button press which word they had 165 perceived. Inter-trial intervals were set to last about two seconds.

166 2.4 Preprocessing of stimulus material

167 From the stimulus material we extracted the following auditory and visual features. Based on previous 168 literature that demonstrated robust encoding of the amplitude envelope, it's temporal derivative and the fundamental frequency of speech, we derived these features from the acoustic speech recordings (Bröhl 169 170 and Kayser, 2021; Oganian and Chang, 2019; Teoh et al., 2019). To derive the broadband envelope we 171 filtered the acoustic waveform into twelve logarithmically spaced bands between 0.1 and 10 kHz (zerophase 3rd order Butterworth filter with boundaries: 0.1, 0.22, 0.4, 0.68, 1.1, 1.7, 2.7, 4.2, 6.5, 10 kHz) and 172 173 subsequently took the absolute value of the Hilbert transform for each band. The broadband amplitude 174 envelope (hereon termed aud env) was then derived by taking the average across all twelve band-limited 175 envelopes and was subsequently down-sampled to 50 Hz. We computed the slope of this broadband 176 envelope (hereon termed aud slope) by taking its first derivative. To characterize the pitch contour we 177 extracted the fundamental frequency (hereon termed aud pitch) over time using the Praat software ('to 178 Pitch' method with predefined parameters) (Boersma and van Heuven, 2001). This was done using the 179 original acoustic waveform at a sampling rate of 22,050 Hz. The resulting pitch contour was again down 180 sampled to 50 Hz. All three acoustic features together are labelled AudFeat in the following.

181 In a similar fashion we derived the horizontal opening of the lips, the area covered by the lip opening, and 182 its derivative from the video recordings. The lips were detected based on the color of the lips in the video 183 material using a custom-made algorithm. From these we determined the contour of the lip opening based

185 were visually inspected to ensure accurate tracking of the lips. From this segmentation of the lip opening 186 we derived the total opening (in pixels) (hereon termed *lip area*) and estimates of the respective diameters along the horizontal axes (hereon termed lip width): these were defined between the outermost points 187 188 along the horizontal axis. These signals were initially sampled at the video rate of 25 fps. As for the auditory 189 features, we computed the slope of the lip area (hereon termed lip slope). The time series of these visual 190 features were then linearly interpolated to a sample rate of 50 Hz. Because the horizontal and vertical 191 mouth openings are partially correlated with each other and with the total mouth opening, we selected the 192 total area and the horizontal width as signals of interest, as the latter is specifically informative about the 193 acoustic formant structure (Plass et al., 2020). We grouped the total lip area, it's temporal derivative and 194 the lip-width as signatures of lip features (LipFeat), which are of the same dimensionality as the acoustic 195 features (AudFeat) described above. 196 For comparison with previous studies (Chandrasekaran et al., 2009; Giordano et al., 2017; Hauswald et al., 197 2018; Park et al., 2016) we quantified the power spectra of these features and their cross-coherences using 198 MATLAB's 'pwelch' and 'mscoher' functions using a window length of 1 s with 50% overlap and otherwise

predefined parameters. The resulting spectra were log transformed and averaged across sentences. To
 visualize the cross-coherences we first obtained key frequency ranges of interest from our main results (c.f.
 Fig. 3) and averaged the coherences within two ranges of interest (0.5 - 1 Hz and 1 - 3 Hz).

on luminance values and deriving connected components from these (Giordano et al., 2017). The results

202 2.5 MEG preprocessing

184

203 Preprocessing of MEG data was carried out using custom MATLAB scripts and the FieldTrip toolbox 204 (Oostenveld et al., 2011). Each experimental block was processed separately. Individual trials were 205 extracted from continuous data starting 2 s before sound onset and until 10 s after sound onset. The MEG 206 data were denoised using a reference signal. Known faulty channels (N = 7) were removed. Trials with 207 SQUID jumps (3.5% of trials) were detected and removed using FieldTrip procedures with a cut-off z-value 208 of 30. Data were band-pass filtered between 0.2 and 150 Hz using a zero-phase 4th order Butterworth filter 209 and subsequently down sampled to 300 Hz before further artefact rejection. Data were visually inspected 210 to find noisy channels (4.37 \pm 3.38 on average across blocks and participants) and trials (0.66 \pm 1.03 on 211 average across blocks and participants). Noise cleaning was performed using independent component 212 analysis with 30 principal components (2.5 components removed on average). Data were further down 213 sampled to 50 Hz and bandpass filtered between 0.8 and 30 Hz using a zero-phase 3rd order Butterworth 214 filter for subsequent analysis.

215 2.6 MEG source reconstruction

Source reconstruction was performed using Fieldtrip, SPM8, and the Freesurfer toolbox based on T1weighted structural magnetic resonance images (MRIs) for each participant. These were co-registered to 218 the MEG coordinate system using a semi-automatic procedure (Gross et al., 2013; Keitel et al., 2017). MRIs 219 were then segmented and linearly normalized to a template brain (MNI space). We projected sensor-level 220 time series into source space using a frequency-specific linear constraint minimum variance (LCMV) 221 beamformer (Van Veen et al., 1997) with a regularization parameter of 7% and optimal dipole orientation 222 (singular value decomposition method). The grid points had a spacing of 6 mm, thus resulting in 12,337 223 points. For whole-brain analyses, a subset of grid points corresponding to cortical gray matter regions only 224 was selected (using the AAL atlas, Tzourio-Mazoyer et al., 2002), yielding 6,490 points in total. Within these 225 we defined temporal and occipital regions of interest (ROI) based on the brainnetome atlas (Yu et al., 226 2011). The individual ROIs were chosen based on previous studies that demonstrate the encoding of acoustic and visual speech features in occipital and superior temporal regions (Di Liberto et al., 2018; 227 228 Giordano et al., 2017; Keitel et al., 2020; Teng et al., 2018). As temporal ROI we included Brodmann area 229 41/42, caudal area 22 (A22c), rostral area 22 (A22r) and TE1.0 and TE1.2. As occipital ROI we defined the 230 middle occipital gyrus (mOccG), occipital polar gyrus (OPC), inferior occipital gyrus (iOccG) and the medial 231 superior occipital gyrus (msOccG).

232 2.7 MEG analysis

233 The questions outlined in the introduction require quantifying how well the source reconstructed MEG data 234 reflect the visual and or acoustic features. For this we relied on a previously established and validated 235 mutual information (MI) framework (Ince et al., 2017). The analysis relies on the notion that a significant 236 temporal relation between a cerebral signal and sensory features is indicating the cerebral encoding (or 237 tracking) of the respective features in temporally entrained brain activity (Bröhl and Kayser, 2021; Keitel et 238 al., 2018; Park et al., 2016). In the following we use the term 'tracking' when referring to such putative 239 cerebral representations characterized using MI (Obleser and Kayser, 2019). To quantify the tracking of a given stimulus feature, or of a feature group, we concatenated the trial-wise MEG data and features along 240 241 the time dimension and filtered these (using 3rd order Butterworth IIR filters) into typical frequency bands 242 used to study dynamic speech encoding: 0.5 - 1 Hz, 1 - 3 Hz, 2 - 4 Hz, 3 - 6 Hz and 4 - 8 Hz (and 0.5 - 8 Hz). 243 These were chosen to cover the typical modulation spectra of these features (Fig. 1C,D) and similar to 244 previous work (Bröhl and Kayser, 2021; Etard and Reichenbach, 2019; van Bree et al., 2020; Zuk et al., 245 2021). The first 500 ms of each sentence were discarded to remove the influence of the transient sound-246 onset response. To compute the MI between filtered MEG and stimulus features, we relied on a complex-247 valued representation of each signal, which allowed us to include both the amplitude and phase 248 information in the analysis: we first derived the analytic signal of both the MEG and stimulus feature(s) 249 using the Hilbert transform and then calculated the MI using the Gaussian copula approach including the 250 real and imaginary part of the Hilbert signals (Daube et al., 2019; Ince et al., 2017).

In a first step, we used this framework to visualize the tracking of AudFeat and LipFeat within the entire source space (Fig. 2A,B). This was mainly done to assert that the predefined ROIs used for the subsequent

253 analysis indeed covered the relevant tracking of these features. This analysis relied on a frequency range 254 from 0.5 to 8 Hz and a range of stimulus-to-brain lags from 60 to 140 ms after stimulus onset. As a second 255 step, we then quantified the tracking of auditory or visual features and their dependencies specifically 256 within these ROIs and individual frequency bands (Fig. 3,4,5). To facilitate these analyses, we first 257 determined the optimal lags for each feature, ROI and frequency band, given that the encoding latencies 258 may differ between features and regions (Giordano et al., 2017). For this we determined at the group-level 259 and for each set of features (i.e. AudFeat and LipFeat) and for each ROI and frequency band the respective 260 lag yielding the largest group-level MI value (across participants and both A-only and V-only trials). This was 261 done by computing the MI between each set of features and the MEG in a range of lags between 0 and 500 262 ms in 20 ms steps. For the subsequent analyses, we used these optimal lags and computed averaged MI 263 values in a time window of -60 to 60 ms around these lags (computed in 20 ms steps).

264 The first question of this study as outlined in the introduction concerns the tracking (MI) of individual feature groups in temporal and occipital ROIs and the two experimental conditions, respectively, and more 265 so if a given ROI reflects a given feature (e.g. the unheard acoustic envelope) independently of the 266 267 physically present other feature (e.g. the visible lip movements in the visual-only condition). To quantify 268 whether the tracking of each feature group (in a given ROI and frequency band) is statistically redundant 269 with (or possibly complementary to) the other group, we calculated the conditional mutual information 270 (CMI) between MEG and one feature group, partialling out the respective other group (Fig. 3, CMI values) 271 (Giordano et al., 2017; Ince et al., 2017). Specifically, the CMI measure allows us to quantify the unique 272 information shared between a variable and the MEG while controlling for the information provided by the 273 conditional variable. Mathematically, it can be described as

274

281

I(X;Y|Z) = H(X,Z) + H(Y,Z) - H(X,Y,Z) - H(Z)

where *I* denotes the conditional mutual information and *H* the joint entropies between combinations of variables *X*, *Y* and the conditional variable *Z*. Similarly, we also determined the CMI between the MEG and each individual feature, obtained by partialling out all other visual and auditory features (Fig. 4). To be able to compare the MI and CMI estimates directly, we ensured that both estimates had comparable statistical biases. To achieve this, we effectively derived the MI as a conditional estimate, in which we partialled out a statistically-unrelated variable. That is, we defined

MI(feature ; MEG) I *MI(feature ; MEG* | *time_shifted_feature)*

Here, time_shifted_feature is a representation of the respective feature(s) with a random time lag and hence no expected causal relation to the MEG. Each MI estimate was obtained by averaging this estimate over 2,000 repetitions of a randomly generated time-shifted feature vector. To render the (conditional) MI estimates meaningful relative to the expectation of zero MI between MEG and stimulus features, we furthermore subtracted an estimate of the null-baseline of no systematic relation between signals. This was obtained by computing (conditional) MI values after randomly time-shifting the stimulus feature(s) and averaging the resulting surrogate MI estimates over 100 randomizations.

289 2.8 Relating MI to word recognition performance

290 The behavioral performance for each participant and condition was obtained as the percent correctly (PC) 291 reported target words (obtained in a 4-choice task). To probe the second question of whether the tracking 292 of restored features relates to word recognition performance we relied on partial regression. Specifically, 293 we probed the linear relation of word recognition performance and feature tracking across participants 294 while accounting for potential spurious correlations between these due to variations in the individual 295 signal-to-noise ratio in each participants' MEG data. We predicted the PC in the visual-only trials based on i) 296 the individual MI for aud env in the temporal ROI and the MI for aud pitch in the and occipital ROI as the 297 primary variables of interest, and ii) the tracking of LipFeat (MI) in the occipital ROI in visual trials and iii) 298 the tracking of AudFeat in the temporal ROI in auditory trials. The last two serve as potentially confounding 299 variables, as they provide a proxy to the overall SNR of the speech and lip tracking in the respective dataset. 300 By focusing on aud env / aud pitch in the temporal/occipital ROIs respectively, we predicted task 301 performance based on the individual features that were most associated with the tracking of AudFeat (c.f. Fig. 4C,D). To establish these regression models, we z-scored the MI values of interest (variables i - iii) and 302 the PC across participants. For the confounding variables, we applied the z-scoring for each frequency band 303 304 and subsequently averaged the z-scored values across bands. For each frequency band, we created a single 305 model containing all target and confounding variables. From the respective models we obtained the 306 significance of each predictor of interest. Furthermore, we compared the predictive power of this full 307 model with that of a reduced model not featuring the predictors of interest (variable i). From the 308 likelihoods of each model we derived the relative Bayes factor (BF) between these based on the respective 309 BIC values obtained from each model. For visualization we used partial residual plots using the procedure 310 described by Velleman and Welsch (Velleman and Welsch, 1981). This procedure was applied to each 311 individual feature of interest (i.e. aud env and aud pitch).

312 2.9 Statistical analysis

313 Statistical testing of mutual information data was based on a non-parametric randomization approach 314 incorporating corrections for multiple comparisons (Nichols and Holmes, 2003). To test whether the group-315 level median MI (or CMI) values were significantly higher than expected based on the null hypothesis of no 316 systematic temporal relation between sensory features and MEG, we proceeded in a similar fashion as in 317 previous work (Bröhl and Kayser, 2021; Giordano et al., 2017): we obtained a distribution of 2,000 MI 318 values between randomly time-shifted MEG and the stimulus vectors, while keeping the temporal relation 319 of individual features to each other constant. This distribution was obtained for each participant, frequency 320 band, feature group (AudFeat and LipFeat), ROI (temporal, occipital) and condition (A-only, V-only) 321 separately. To correct for multiple comparisons, we generated a single random distribution by pooling the 322 randomly generated MI values across all dimensions except frequency bands, given that the MI values 323 decreased considerably across bands (c.f. Fig. 3), and selecting the maximum 2000 values, thereby creating 324 a random maximum null distribution (Nichols and Holmes, 2003). We then tested the group-level median 325 against the 99th percentile of this maximum distribution as a significance threshold, which effectively implements a one-sided randomization test at p < 0.01 corrected for all dimensions except frequency 326 bands. To test for differences between MI and CMI values for a given condition, band and ROI, we also used 327 328 a permutation approach combined with a Wilcoxon signed-rank test: first, we established the respective 329 true Wilcoxon z-statistic between MI and CMI values; then we created a distribution of surrogate z-330 statistics under the null hypothesis of no systematic group-level effect, obtained by randomly permuting 331 the labels of MI and CMI values 5,000 times. From this we obtained the maximum across features, bands, 332 ROIs and conditions to correct for multiple comparisons and used the 99th percentile of this randomization 333 distribution to determine the significance of individual tests.

334 The CMI values for individual features in Figure 4 were compared using a one-way repeated measure 335 Kruskal-Wallis rank test, followed by a post-hoc Tukey Kramer multiple comparison. We used the same procedure to test for differences between CMI values in the sub-areas composing each ROI (Table 1). To 336 337 test CMI values between hemispheres, we used a Wilcoxon signed rank test (Table 2). The resulting p-338 values were corrected for false discovery rate using the Benjamini-Hochberg procedure within each set of 339 comparisons (Benjamini and Hochberg, 1995). In all tests an alpha level of α < 0.01 was deemed significant. 340 For all statistical tests we provide exact p-values, except for randomization tests where the approximate pvalues were smaller than the inverse of the number of randomizations. 341



342

Fig. 1. Stimulus material and experimental methodology. Acoustic and visual features were extracted from audiovisual speech material and were used to quantify their cerebral tracking during audio-only and visual-only presentations. (A) The stimulus material consisted of 180 audiovisual recordings of a trained actor speaking individual 11

346 English sentences. For visualization here only the mouth is shown, but participants were presented with the entire face. 347 From the video recordings we extracted three features describing the dynamics of the lip aperture: the area of lip 348 opening (lip area), its slope (lip slope), and the width of lip opening (lip width); collectively termed 'LipFeat'. From the 349 audio waveform we extracted three acoustic features: the broadband envelope (aud env), its slope (aud slope), and a 350 measure of dominant pitch (aud pitch); collectively termed 'AudFeat'. (B) Trial-averaged percent correctly (PC) reported 351 target words in auditory (A-only) and visual-only (V-only) conditions, with dots representing individual participants. (C) 352 Logarithmic power spectra for individual stimulus features. For reference, a 1/f spectrum is shown as a dashed grey 353 line. (D) Coherence between pairs of features averaged within two predefined frequency bands (0.5 - 1 Hz left; 1 - 3 Hz 354 right, see Methods for details).

355 **3. Results**

356 3.1 Acoustic and visual features are tracked in temporal and occipital cortices

Participants were presented with either spoken speech (in A-only trials) or a silent video of the speaking 357 face (in V-only trials) and were asked to report a target word for each sentence in a 4-choice word 358 359 recognition task. The behavioral data show that participants were well able to detect the correct word both 360 during acoustic speech embedded in noise and during lip reading and achieved overall similar levels of 361 performance in both conditions (Fig. 1B, median fraction correct responses for A-only = 0.7, V-only = 0.71; n = 18). To quantify the tracking of relevant features, we defined three auditory (AudFeat) and three visual 362 363 (LipFeat) features respectively based on the acoustic waveform and the lip trajectory (Fig. 1A). An analysis of their temporal coherences revealed that they were coherent in the frequency bands of interest (e.g. 1 - 3 364 Hz, envelope-lip area coherence of ~0.2) (Fig. 1D). The overall pattern of coherence and the degree of 365 temporal relation between acoustic features and lip movements in the present material is comparable with 366 those in other datasets (Chandrasekaran et al., 2009; Giordano et al., 2017; Hauswald et al., 2018; Park et 367 368 al., 2016).

Previous work has shown that in the dataset analyzed here temporal and occipital brain regions reflect auditory and visual speech signals respectively (Keitel et al., 2020). We extend this observation to the entire group of acoustic (AudFeat; Fig. 2A) or lip features (LipFeat; Fig. 2B) using a mutual information (MI) approach. The whole-brain maps show the expected prevalence of acoustic (visual) tracking in temporal (occipital) regions. Given that our main questions concerned the tracking of features specifically in occipital and temporal brain regions, we focused the subsequent work on atlas-based regions of interest (Fig. 2C; the temporal ROI shaded in mint and the occipital ROI shaded in purple, see methods for details).

12



376

Fig. 2. Tracking of auditory and visual features in MEG source space. The figure shows group-level median MI values
for auditory (AudFeat; panel A) and lip features (LipFeat; panel B) in the frequency range from 0.5 - 8 Hz (n = 18
participants). (C) Colored shading indicates regions of interest: temporal region in mint includes Brodmann area 41/42,
caudal area 22 (A22c), rostral area 22 (A22r) and TE1.0 and TE1.2; occipital region in purple includes middle occipital
gyrus (mOccG), occipital polar gyrus (OPC), inferior occipital gyrus (iOccG) and medial superior occipital gyrus
(msOccG).

383 3.2 Temporal and occipital cortex represent acoustic speech features during silent lip reading

To address the question of whether temporal and occipital cortices represent auditory and visual speech 384 features during lip reading, we performed a comprehensive analysis of the tracking of both sets of features 385 386 across a range of frequency bands during auditory (A-only) and visual (V-only) conditions (MI values; Figure 387 3). To further quantify whether the tracking of each feature group is possibly redundant with the tracking 388 of the respective other feature group, we derived CMI values for each feature group, obtained by partialling 389 out the respective other group (CMI values). By comparing MI and CMI values we can test, for example, 390 whether the temporal ROI tracks the unheard speech envelope during silent lip reading also when 391 discounting for the actually presented lip trajectory. In the following we discuss the results per sensory 392 modality and region of interest.

As expected, when listening to speech (A-only), the temporal ROI significantly tracks auditory features (AudFeat) in all frequency bands tested (Fig. 3, top row, red MI data; non-parametric randomization test, all bands: $p < 5 \times 10^{-5}$). This tracking persists when discounting potential contributions of the not-seen visual features (red CMI data all individually significant: $p < 5 \times 10^{-5}$), though in some bands the CMI values were significantly lower than the unconditional MI (Wilcoxon signed rank test comparing MI vs. CMI, 2 - 4 Hz: z = 398 3.59, 3 - 6 Hz: z = 3.68, 4 - 8 Hz: z = 3.42, all comparisons: $p < 2 \times 10^{-5}$). During the same auditory trials, lip 399 features are only marginally reflected in the temporal ROI, as shown by low but significant MI and CMI 400 values above 1Hz (Fig. 3, top row, cyan MI and CMI data; all bands above 1 Hz: $p < 5 \times 10^{-5}$). This tracking of 401 visual features was significantly reduced when partialling out the physically presented auditory features (2 -4 Hz: z = 3.59, 3 - 6 Hz: z = 3.68, 4 - 8 Hz: z = 3.42, all comparisons: $p < 2 \times 10^{-5}$).

403 During lip reading (V-only), the temporal ROI tracks the unheard auditory features, particularly below 1 Hz 404 (Fig. 3, 2nd row, red MI data; all bands: $p < 5 \times 10^{-5}$). Except in the 2 - 4 Hz range, the temporal ROI tracks 405 the unheard AudFeat to a similar degree as when discounting the actually presented visual signal (significant red CMI values, all bands: $p < 5 \times 10^{-5}$) as there were no significant differences between MI and 406 CMI values except one band (2 - 4 Hz: z = 3.42, $p < 1 \times 10^{-4}$, see asterisks). The physically presented lip 407 movements during these V-only trials were also tracked significantly in the temporal ROI (Fig. 3, 2nd row; 408 cyan MI and CMI data, 1 - 6 Hz: $p < 5 \times 10^{-5}$) but the CMI values were only marginally above chance level, 409 410 suggesting that genuine visual representations in the temporal region is weak.

As expected, during lip reading (V-only) the occipital ROI tracks lip features (LipFeat) across frequency 411 bands (Fig. 3, bottom row, cyan MI values; all bands: $p < 5 \times 10^{-5}$). Again, this tracking persists after 412 413 partialling out the non-presented acoustic features (cyan CMI values; all bands: $p < 5 \times 10^{-5}$), although the CMI values were significantly lower than the MI (all bands above 1 Hz: $z \ge 3.72$, $p < 2 \times 10^{-5}$). This indicates 414 415 some redundancy between the tracking of the physically present lip trajectory and that of the unheard 416 auditory features. Confirming this, occipital tracking of the physically presented lip signals emerges in 417 parallel with that of the non-presented auditory features (Fig. 3, bottom panel, red MI data; all bands: p < 5 imes 10⁻⁵). This occipital tracking of unheard auditory features was significantly reduced when partialling out 418 the lip signal (MI vs. CMI data; all bands above 1 Hz: $z \ge 3.72$, $p < 2 \times 10^{-5}$) but remained statistically 419 significant (red CMI data; all bands: $p < 5 \times 10^{-5}$). 420

Finally, when listening to speech (A-only), the occipital ROI shows significant but weak tracking of auditory (Fig. 3, 3rd row, red MI data; 1 - 6 Hz: $p < 5 \times 10^{-5}$) and visual features (cyan MI data; only 3 - 6 Hz: $p < 5 \times 10^{-5}$), suggesting that purely acoustic signals have a weak influence on the occipital brain region.

Collectively, these results show the expected representations of auditory features in temporal cortex during listening to speech and of lip features in occipital cortex during lip reading. In addition, they reveal that during lip reading, both temporal and occipital regions represent unheard auditory features and do so independently of co-existing representations of the physically presented lip movements. In the auditory cortex this 'restoration' of auditory signals prevails in the low delta band (0.5 - 1 Hz), in the visual cortex this emerges in multiple bands.



AudFeat LipFeat



431 Fig. 3. Feature tracking across regions of interest and conditions. For both conditions (A-only and V-only) and ROIs 432 (temporal and occipital) the figure illustrates the strength of feature tracking for presented and physically not-present 433 features (MI values) and the respective strength of tracking after partialling out the respective other feature group 434 (CMI values). Each panel depicts (from left to right) the MI for AudFeat, the CMI for AudFeat partialling out LipFeat, the 435 MI for LipFeat, and the CMI for LipFeat partialling out AudFeat. Dots represent individual participants (n = 18). Bars 436 indicate the median, 25th and 75th percentile. The grey dashed line indicates the 99th percentile of the frequency-437 specific randomized maximum distribution correcting for all other dimensions. Conditions below a group-level 438 significance threshold of 0.01 are greyed out. Brackets with asterisks indicate significant differences between MI and 439 CMI values, based on a Wilcoxon signed-rank test (* p < 0.01, ** p < 0.005, *** p < 0.001). Units for MI and CMI are in 440 bits.

To obtain an estimate of the effect size of the restoration of the unheard AudFeat during lip reading we expressed the respective CMI values relative to those of the tracking of the respectively modality-preferred inputs of each ROI (Fig. 4A,B): for the temporal region the tracking of AudFeat during A-only trials and for the occipital region the tracking of LipFeat during V-only trials. In the temporal ROI, the restoration effect size, i.e. the tracking of AudFeat during lip reading, was about a third as strong as this feature's tracking while directly listening to speech (Fig. 4A; AudFeat_{V-only} /AudFeat_{A-only}: 0.5 - 1 Hz: median = 0.37, 1 - 3 Hz:

eNeuro Accepted Manuscript

447 median = 0.24). In the occipital ROI, the tracking of AudFeat was about half as strong or stronger compared 448 to the tracking of lip features when seeing the speaker (Fig. 4B; AudFeat_{V-only} /LipFeat_{V-only}; 0.5 - 1 Hz: 449 median = 0.84, 1 - 3 Hz: median = 0.4). Albeit smaller than the tracking of the respective modality-preferred 450 sensory inputs, the restoration of unheard auditory features still results in a prominent signature in 451 temporally aligned brain activity in both cortices.

452 3.3 Feature tracking is bilateral and prevails across anatomical brain areas

453 Having established the tracking of auditory and lip features in both temporal and occipital ROIs, we probed whether this tracking is possibly lateralized in a statistical sense and whether it potentially differs among 454 455 the individual anatomical areas grouped into temporal and occipital ROIs respectively. While these analyses 456 do not directly concern our main hypotheses outlined in the introduction, the issue of lateralization is pervasive in the literature on speech, and hence is addressed here for the sake of completeness. For this 457 458 analysis we focused on the conditional tracking of each feature group. Comparing CMI values among 459 anatomical areas (averaged across hemispheres) for each ROI (occipital, temporal), frequency band (0.5 - 1 and 1 - 3 Hz), condition and feature group revealed a significant effect of area for AudFeat tracking in the 460 temporal ROI during A-only trials (Table 1; 0.5 - 1 Hz: $\chi^{2}(3) = 27.02$, p = 4.7 x 10⁻⁶, $\epsilon^{2} = 0.35$; 1 - 3 Hz: $\chi^{2}(3) = 27.02$, p = 4.7 x 10⁻⁶, $\epsilon^{2} = 0.35$; 1 - 3 Hz: $\chi^{2}(3) = 27.02$, p = 4.7 x 10⁻⁶, $\epsilon^{2} = 0.35$; 1 - 3 Hz: $\chi^{2}(3) = 27.02$, p = 4.7 x 10⁻⁶, $\epsilon^{2} = 0.35$; 1 - 3 Hz: $\chi^{2}(3) = 27.02$, p = 4.7 x 10⁻⁶, $\epsilon^{2} = 0.35$; 1 - 3 Hz: $\chi^{2}(3) = 27.02$, p = 4.7 x 10⁻⁶, $\epsilon^{2} = 0.35$; 1 - 3 Hz: $\chi^{2}(3) = 27.02$, p = 4.7 x 10⁻⁶, $\epsilon^{2} = 0.35$; 1 - 3 Hz: $\chi^{2}(3) = 27.02$, p = 4.7 x 10⁻⁶, $\epsilon^{2} = 0.35$; 1 - 3 Hz: $\chi^{2}(3) = 27.02$, p = 4.7 x 10⁻⁶, $\epsilon^{2} = 0.35$; 1 - 3 Hz: $\chi^{2}(3) = 27.02$, p = 4.7 x 10⁻⁶, $\epsilon^{2} = 0.35$; 1 - 3 Hz: $\chi^{2}(3) = 27.02$, p = 4.7 x 10⁻⁶, $\epsilon^{2} = 0.35$; 1 - 3 Hz: $\chi^{2}(3) = 27.02$, p = 4.7 x 10⁻⁶, $\epsilon^{2} = 0.35$; 1 - 3 Hz: $\chi^{2}(3) = 27.02$, p = 4.7 x 10⁻⁶, $\epsilon^{2} = 0.35$; 1 - 3 Hz: $\chi^{2}(3) = 27.02$, p = 4.7 x 10⁻⁶, $\epsilon^{2} = 0.35$; 1 - 3 Hz: $\chi^{2}(3) = 27.02$, p = 4.7 x 10⁻⁶, $\epsilon^{2} = 0.35$; 1 - 3 Hz: $\chi^{2}(3) = 27.02$, p = 4.7 x 10⁻⁶, $\epsilon^{2} = 0.35$; 1 - 3 Hz: $\chi^{2}(3) = 27.02$, p = 4.7 x 10⁻⁶, $\epsilon^{2} = 0.35$; 1 - 3 Hz: $\chi^{2}(3) = 27.02$, p = 4.7 x 10⁻⁶, $\epsilon^{2} = 0.35$; 1 - 3 Hz: $\chi^{2}(3) = 27.02$, p = 4.7 x 10⁻⁶, $\epsilon^{2} = 0.35$; 1 - 3 Hz: $\chi^{2}(3) = 27.02$, p = 4.7 x 10⁻⁶, $\epsilon^{2} = 0.35$; 1 - 3 Hz: $\chi^{2}(3) = 27.02$, p = 4.7 x 10⁻⁶, $\epsilon^{2} = 0.35$; 1 - 3 Hz: $\chi^{2}(3) = 27.02$, p = 4.7 x 10⁻⁶, $\epsilon^{2} = 0.35$; 1 - 3 Hz; $\chi^{2}(3) = 27.02$, p = 4.7 x 10⁻⁶, $\epsilon^{2} = 0.35$; 1 - 3 Hz; $\chi^{2}(3) = 27.02$, p = 4.7 x 10⁻⁶, $\epsilon^{2} = 0.35$; 1 - 3 Hz; $\chi^{2}(3) = 0.35$; 1 - 3 Hz; 1 - 461 29.62, p = 2.7 x 10^{-6} , ε^2 = 0.39; p-values FDR-corrected). Post hoc comparisons revealed that in both bands, 462 tracking of AudFeat was higher in A41/42 and A22c compared to TE1.0/1.2 and A22r (Tukey-Kramer test, all 463 tests $p < 10^{-5}$). The effect of Area was close to but not significant for LipFeat tracking in the occipital ROI 464 during V-only trials (0.5 - 1 Hz: $\chi^2(3)$ = 12.3, p = 0.026, ε^2 = 0.14; 1 - 3 Hz: $\chi^2(3)$ = 14.57, p = 0.012, ε^2 = 0.17). 465 466 Importantly, these results suggest that while the tracking of auditory features was stronger in the the early 467 auditory region during A-only trials, the restoration of unheard auditory features during lip reading 468 emerges to a similar degree among the individual temporal and occipital areas.

We performed a similar analysis comparing the CMI values within temporal or occipital ROIs between hemispheres. This revealed no significant effects of hemispheres (Table 2), hence offering no evidence for a

471 statistical lateralization of feature tracking in the present data.

		0.5 - 1 Hz			1 - 3 Hz				
ROI	Anatomical area	AudCMI	Chisq; pval	LipCMI	Chisq; pval	AudCMI	Chisq; pval	LipCMI	Chisq; pval
A-only trials	A-only trials								
temporal	A41/42	0.97	27.02;	0.096	2.47 ; 0.59	0.19	29.62 ; 2 7a 05	0.032	5.14 ; 0.32
	TE1.0/1.2	0.56	4.7e-05	0.099		0.11	2.70-05	0.028	
	A22c	0.86		0.098		0.18		0.033	
	A22r	0.5		0.093		0.095		0.029	
occipital	mOccG	0.1	3.50 ; 0.47	0.073	0.66 ; 0.88	0.025	2.71 ; 0.58	0.02	1.97 ; 0.66

	OPC	0.09		0.069		0.025		0.02	
	iOccG	0.11		0.068		0.027		0.021	
	msOccG	0.11		0.074		0.029		0.021	
V-only trials									
temporal	A41/42	0.33	4.59 ; 0.36	0.1	5.14 ; 0.32	0.034	5.24 ; 0.32	0.027	1.00 ; 0.85
	TE1.0/1.2	0.25		0.088		0.031		0.028	
	A22c	0.34		0.1		0.035		0.027	
	A22r	0.22		0.085		0.028		0.029	
occipital	mOccG	0.13	3.90 ; 0.44	0.17	12.30;	0.045	8.20 ; 0.13	0.15	14.57;
	OPC	0.14		0.19	0.020	0.06		0.2	0.012
	iOccG	0.14		0.2		0.048		0.17	
	msOccG	0.11		0.11		0.039		0.082	

472

473 Table 1. Feature tracking in individual anatomical areas within temporal and occipital ROIs. The table lists

474 CMI values of either set of features (AudCMI, LipCMI) and a statistical comparison between the individual

475 atlas-defined areas of the temporal and occipital ROIs (Kruskal-Wallis tests, reporting chi-squares (Chisq)

476 and p-values (pval)). Bold numbers indicate statistically significant results. P-values are FDR-corrected within

477 this table.

			0.5 -	1 Hz			1 - 3	3 Hz	
ROI	Hemisphere	AudCMI	z; pval	LipCMI	z; pval	AudCMI	z; pval	LipCMI	z; pval
A-only trials									
temporal	left	0.8	1.20;	0.094	-0.33;	0.13	-0.81;	0.028	-1.11;
	right	0.64	0.39	0.099	0.74	0.15	0.39	0.031	0.39
occipital	left	0.1	-0.37;	0.07	-0.33;	0.027	0.81;	0.021	0.63;
	right	0.1	0.74	0.072	0.74	0.025	0.59	0.02	0.05
V-only trials	V-only trials								
temporal	left	0.31	0.89;	0.098	0.76;	0.035	0.85;	0.025	-1.85;
	right	0.26	0.59	0.091	0.59	0.03	0.59	0.03	0.20
occipital	left	0.11	-2.24;	0.14	-2.98;	0.043	-1.68;	0.13	-2.07;
	right	0.15	0.2	0.2	0.040	0.054	0.5	0.18	0.21

478

479Tab. 2. Feature tracking in each hemisphere. The table lists CMI values of either set of features (AudCMI,480LipCMI) and a statistical comparison between hemispheres of each ROI (Wilcoxon signed rank tests,

481 reporting z values (z) and p-values (pval)). P-values are FDR-corrected within this table.

482 3.4 Occipital cortex reflects pitch more than other acoustic features during lip reading

Having established that occipital and temporal regions track unheard auditory features, we then asked how individual features contribute to these representations. For this we focused on the following condition: the tracking of AudFeat in the delta range in V-only trials (Fig. 4C,D). We quantified the CMI for each individual feature, while discounting the evidence about all other left-out visual and auditory features, hence focusing on the unique tracking of each individual acoustic feature.

For the temporal ROI this revealed the prominent tracking of aud env (Fig. 4C). In the in the 0.5 - 1 Hz band only the CMI for aud env was above chance ($p < 5 \times 10^{-5}$) and there was a significant effect of feature (Kruskal-Wallis rank test $\chi^2(2) = 9.27$, $p = 9.1 \times 10^{-4}$, $\varepsilon^2 = 0.14$). Post-hoc tests revealed that the CMI for aud env differed significantly from that of aud slope (Tukey-Kramer test, $p = 6.2 \times 10^{-4}$; the other comparisons were not significant; p = 0.35 for env vs. slope and p = 0.22 for slope vs. pitch). In the 1 - 3 Hz band, the tracking of all auditory features was significant (all features: $p < 5 \times 10^{-5}$) and there was no significant effect of features ($\chi^2(2) = 4.14$, p = 0.13, $\varepsilon^2 = 0.04$).

For the occipital ROI, this revealed a dominance of aud pitch (Fig. 4D). In the 0.5 - 1 Hz band, only the CMI 495 of aud pitch was above chance ($p < 5 \times 10^{-5}$), a direct comparison revealed a significant effect of features 496 (0.5 - 1 Hz: χ^2 (2) = 18.28, p = 1.07 x 10⁻⁴, ϵ^2 = 0.32) and post-hoc tests revealed a significant difference 497 between aud pitch and aud slope (p = 7.03 x 10^{-5}), while the other comparisons were not significant (p = 498 0.26 for pitch vs. env and p = 0.02 for env vs. slope). In the 1 - 3 Hz range, the tracking of all features was 499 significant (all features: $p < 5 \times 10^{-5}$), there was a significant effect of features $\chi^2(2) = 19.2$, $p = 6.77 \times 10^{-5}$, ϵ^2 500 = 0.34), and post-hoc tests revealed a significant difference between pitch and slope ($p = 3.61 \times 10^{-5}$), while 501 502 the other comparisons were not significant (p = 0.05 for pitch vs. env and p = 0.12 for env vs. slope). 503 Collectively these results suggest that the restoration of acoustic signals in the occipital region emphasizes 504 spectral pitch, while in the temporal region this emphasizes the temporal speech envelope.





506 Fig. 4. Modality dominance and tracking of individual auditory features during lip reading. (A,B) Comparison of the 507 tracking of unheard AudFeat over the tracking of the modality-preferred sensory input in each ROI (i.e. AudFeat during 508 A-only trials in the temporal ROI; LipFeat during V-only trials in the occipital ROI). (C,D) Tracking of individual auditory 509 features during V-only trials conditioned on all other auditory and lip features in temporal (C) and occipital (D) ROIs. 510 Brackets with asterisks indicate levels of significance from one-way Kruskal-Wallis rank test with post-hoc Tukey-Kramer testing (* p < 0.01, ** p < 0.005, *** p < 0.001). Dots represent individual data points. Bars indicate the 511 512 median, 25th and 75th percentile. The grey dashed line indicates the 99th percentile of the frequency-specific 513 randomized maximum distribution correction for all other features. Units in (A) and (B) are a ratio, in panels (C) and (D) 514 units are in bits.

515 **3.5 Tracking of auditory features is associated with lip reading performance**

Finally, we probed the second main question of whether the restoration of unheard auditory features during silent lip reading relates to word recognition performance. For this we probed the predictive power of the MI about specific auditory features in either ROI for word recognition performance during V-only trials (Fig. 5). We specifically focused on the tracking of aud env in the temporal ROI and of aud pitch in the occipital ROI as the dominant feature-specific representations (c.f. Fig. 4C,D). Using linear models we predicted word recognition scores across participants based on the tracking indices of interest and while discounting for potential confounds from differences in signal-to-noise ratio in the MEG data.

523	The results show that variations in word recognition scores are well predicted by the collective measures of
524	feature tracking (0.5 - 1 Hz: R^2 = 0.74, 1 - 3 Hz: R^2 = 0.8). Importantly, the tracking of aud env in the
525	temporal ROI was significantly predictive of lip reading performance (Fig. 5, 0.5 - 1 Hz: β = 0.6, p = 0.037; 1 -
526	3 Hz: aud env β = 0.6, p = 2.8 x 10 ⁻⁴), while tracking of pitch in the occipital ROI was not (0.5 - 1 Hz: β = -0.13,
527	p = 0.56; 1 - 3 Hz: β = -0.026, p = 0.91). This conclusion is also supported by Bayes factors for the added
528	predictive power of aud env and aud pitch to these models (aud env in the temporal ROI; 0.5 - 1 Hz: BF =
529	3.12; 1 - 3 Hz: BF = 26.34; aud pitch in the occipital ROI; 0.5 - 1 Hz BF = 0.3; 1 - 3 Hz BF = 0.24).



530

Fig. 5. Association between lip reading performance and tracking of auditory features. Across participants the tracking of aud env during V-only trials in the temporal ROI but not the tracking of aud pitch in the occipital ROI was significantly associated with word recognition performance (PC) across participants in visual trials. Graphs show partial residual plots, dots represent individual data points and the line indicates the linear fit to the target variable from the full regression model.

537 Natural face-to-face speech is intrinsically multidimensional and provides the auditory and visual pathways 538 with partly distinct acoustic and visual information. These pathways could in principle focus mainly on the 539 processing of their modality-specific signals, effectively keeping the two input modalities largely separated. 540 Yet, many studies highlight the intricate multisensory nature of speech-related representations in the brain, 541 including multisensory convergence at early stages of the hierarchy (Bernstein and Liebenthal, 2014; Crosse 542 et al., 2015; Schroeder et al., 2008; Schroeder and Lakatos, 2009) as well as in classically amodal speech 543 regions (Keitel et al., 2020; Mégevand et al., 2020; Scott, 2019). However, as the present results suggest, 544 the auditory and visual pathways are also capable of apparent 'restoring' information about an absent 545 modality-specific speech component: while seeing a silent speaker, both auditory and visual cortices track the temporal dynamics of the speech envelope and the pitch contour respectively, in a manner that is 546 547 independent on the physically visible lip movements. These 'restored' representations of acoustic features 548 relate to participants' word recognition, suggesting that they may form a central component of silent lip 549 reading.

550 4.1 Auditory and visual cortex reflect acoustic speech features during lip reading

551 We systematically quantified the tracking of auditory and visual speech features during unisensory auditory and visual (lip reading) conditions in dynamically entrained brain activity. As expected, this confirmed that 552 553 early auditory and visual regions reflect acoustic and visual features respectively at the time scales of delta 554 (< 4 Hz) and theta (4 - 8 Hz) band activity, in line with previous work (Aiken and Picton, 2008; Bauer et al., 2020; Doelling et al., 2014; Giraud and Poeppel, 2012; Haegens and Zion Golumbic, 2018; Obleser and 555 556 Kayser, 2019). In addition, we found that during lip reading both regions contained significant information 557 about unheard auditory features, also when discounting for the physically presented lip movements. This 558 representation of acoustic features prevailed in low delta in auditory and delta and theta bands in visual 559 cortex. Interestingly, this representation emphasized the temporal speech envelope in auditory cortex and 560 spectral pitch in visual cortex. These results not only support that both regions are active during lip reading 561 (Besle et al., 2008; Calvert et al., 1997; Calvert and Campbell, 2003; Ludman et al., 2000; Luo et al., 2010), 562 but directly show that they contain temporally and feature-specific representations derived from lip 563 movements that are relevant for comprehension.

564 These results advance our understanding of how the brain exploits lip movements in a number of ways. The 565 restoration of auditory features during silent lip reading has been suggested in previous studies, one 566 quantifying the coherence of temporal brain activity with the non-presented speech envelope 567 (Bourguignon et al., 2020) and others quantifying the coherence between occipital activity and the 568 envelope (Hauswald et al., 2018; Suess et al., 2022). Yet, these studies differed in their precise experimental 569 designs, their statistical procedures revealing the 'restoration' effect, and did not probe a direct link to 570 behavioral performance. The present data demonstrate that such tracking of auditory speech-derived 571 features indeed emerges in parallel and in the same participants. Our data reveal the restoration of

21

572 unheard acoustic features also when discounting the physically present lip signals (i.e. when using 573 conditional mutual information). This finding is important, as the mere coherence of brain activity with the 574 acoustic speech envelope may otherwise simply reflect amodal information contained in the physically-575 present visual speech that is directly redundant with the acoustic domain (Daube et al., 2019).

576 Furthermore, they show that this effect is largely bilateral and emerges across a number of anatomically-577 identified areas, suggesting that it forms a generic property of the respective pathways. Interestingly, the 578 unheard auditory features during V-only trials were restored dominantly in the lower frequencies (0.5 - 3 579 Hz), similarly to recent results (Bourguignon et al., 2020). In principle, activity at these slow timescales may 580 possibly reflect oro-facial cues such as head, eye or eyebrow movements (Munhall et al., 2004; Schroeder et al., 2008). We aimed to mitigate such confounds by instructing the speaker to move their head as little as 581 582 possible and to avoid gestures, and by instructing participants to focus their gaze on the speaker's lips. 583 Moreover, our results align with recent work showing the restoration of the unheard acoustic envelope even when controlling for the speaker's movement during visual presentation (Bourguignon et al., 2020). 584 One may speculate whether this restoration reflects the synthesis of speech-specific elements. However, 585 586 linguistic elements at this time scale mostly encompass phrasal structures, prosody or speech rhythm 587 (Gross et al., 2013; Keitel et al., 2018; Meyer et al., 2017) and few of these are probably restored during lip 588 reading in detail. Possibly, the restoration of the unheard envelope based on lip movements reflects processes for the temporal segmentation of speech-related information based on low-frequency activity 589 590 (Doelling et al., 2014; Ghitza, 2017; Nidiffer et al., 2021).

591 These results come with an important caveat: the capability to read from lips alone is generally low in naïve 592 listeners (Altieri et al., 2011; Grant and Seitz, 2000), which poses an intricate problem when studying the 593 cerebral basis of lip reading. To solicit a sufficient number of trials with successful lip reading and to balance 594 word recognition performance between visual-only and auditory-only trials, we relied on a specifically 595 designed experimental paradigm with two critical features. First, this paradigm relied on sentences 596 constructed based on a repeating set of linguistic elements and a forced-choice task with a closed set of 597 options. This limits the generalizability of the results towards naturally-produced every day speech, as 598 participants could in principle learn the mapping of only target words onto lip movements and choose the 599 most likely one during the course of the experiment. Although we did not strictly control for this, both the 600 chosen elements in each sentence as well as the target and distractor words were chosen randomly. 601 Second, to familiarize participants with the material, the A-only condition preceded the V-only condition 602 during the experiment. This may allow for memory-related processes to contribute to the observed 603 restoration effects. However, the use of 180 syntactically similar but unique sentences makes it in our view 604 highly unlikely that participants solely relied on the stimulus repetition and memory to solve the word 605 recognition task. Rather, we believe that the restoration in the occipital cortex reflects the active parsing of 606 the lip movement signal and engages specific visuo-phonetic transformations, as speculated previously 607 (Hauswald et al., 2018; Nidiffer et al., 2021). This poses a possible solution to how the brain finds the best 608 match between visually perceived speech and a word from a limited set of options. Nevertheless, the visual 609 system might be primed to perceive visual speech after being familiarized with the underlying acoustic 610 stimulus in a previous condition. Even in naturalistic listening situations one is likely to do so when 611 observing a moving face. Therefore, we do not expect any priming of the visual system to confound the 612 nature of lipreading in this paradigm compared to real-life situations. More so, this alludes to the origin of 613 speech-related information during lip reading in general, as bottom-up processes may be aided by 614 sentence-level predictions or expectations that contribute in a top-down manner and partially predict 615 acoustic and lexical information based on the immediately preceding material (Cope et al., 2017). Given 616 that lip reading performance was higher than in other studies or real-life circumstances (Altieri et al., 2011; 617 Grant and Seitz, 2000) it is possible that top-down processes exerted a stronger influence on early visual 618 and auditory cortices in this data compared to real-life circumstances.

619 4.2 Lip reading activates a network of occipital and temporal regions

620 Previous work has shown that lip movements activate a network of temporal, parietal and frontal regions 621 (Bourguignon et al., 2020; Calvert et al., 1997; Capek et al., 2008; O'Sullivan et al., 2017; Ozker et al., 2018; 622 Paulesu et al., 2003; Pekkola et al., 2005) and that both occipital and motor regions can align their neural 623 activity to the dynamics of lip movements (Park et al., 2018, 2016). The present data substantiate this, but 624 also show that the representation of the physically visible lip trajectory along visual pathways is 625 accompanied by the representation of spectral pitch, a type of selectivity not directly revealed previously 626 (Suess et al., 2022). Spectral features are vital for a variety of listening tasks (Albouy et al., 2020; Bröhl and 627 Kayser, 2021; Ding and Simon, 2013; Tivadar et al., 2020, 2018), and oro-facial movements provide concise 628 information about the spectral domain. Importantly, seeing the speaker's mouth allows discriminating 629 formant frequencies and provides a comprehension benefit particularly when spectral features are 630 degraded in the underlying acoustic speech (Plass et al., 2020). This suggests a direct and comprehension-631 relevant link between the dynamics of the lip contour and spectral speech features (Campbell, 2008). 632 Hence, a representation of acoustic features during silent lip reading may underlie the mapping of lip 633 movements onto phonological units such as visemes, a form of language-specific representation emerging 634 along visual pathways (Nidiffer et al., 2021; O'Sullivan et al., 2017). This emphasize the role of the visual 635 system as an active agent during audio-visual speech processing.

Our results corroborate the notion that multisensory speech reception is not contingent only on high-level and amodal representations. Rather, they suggest that the brain likely exploits cross-modal correspondences between auditory and visual speech along a number of dimensions, including basic temporal properties (Bizley et al., 2016; Chandrasekaran et al., 2009) as well as mid-level features, such as pitch or visual object features, whose representation is traditionally considered to be modality specific (Crosse et al., 2015; Plass et al., 2020; Schroeder et al., 2008; Zion Golumbic et al., 2013). Previous work has debated whether visual speech is mainly encoded along the auditory pathways or whether occipital regions 643 contribute genuine speech-specific representations (O'Sullivan et al., 2017; Ozker et al., 2018). Our results 644 speak in favor of occipital regions supporting speech reception by establishing multiple forms of speech-645 related information, including those aligned with the acoustic domain revealed here, and those establishing visemic categories based on complementary visual signals (Nidiffer et al., 2021; Suess et al., 2022). Which 646 647 precise occipital areas and by which patterns of connectivity they contribute to comprehension remains to 648 be investigated, but both kinds of representations may well emerge from distinct temporal-occipital networks (Bernstein and Liebenthal, 2014). While visemic information may be driven by object-related 649 650 lateral occipital regions, the more auditory-aligned representations such as the restoration of spectral 651 signatures may be directly driven by the connectivity between occipital areas and superior temporal 652 regions, which play a key role for audio-visual speech integration (Arnal et al., 2009; Lazard and Giraud, 653 2017). In the auditory cortex, the alignment of neural activity to the unheard speech envelope may reflect 654 the predictive influence of visual signals on guiding the excitability of auditory pathways via low frequency 655 oscillations (Schroeder et al., 2008). This alignment of auditory cortical activity to attended or expected 656 sounds is well documented and has been proposed as a cornerstone of multisensory speech integration in 657 general (Lakatos et al., 2008; Schroeder and Lakatos, 2009; Stefanics et al., 2010), and as shown here, 658 directly relates to participants comprehension performance.

659 Credit author statement

- 660 Conceptualization: A.K., C.K., Project administration: A.K., C.K., Funding acquisition: C.K., Methodology: F.B.,
- 661 A.K., C.K., Software: F.B., A.K., Formal Analysis: F.B., Investigation: F.B., Data Curation: F.B., Supervision:
- 662 C.K., Writing Original Draft: F.B., C.K., Writing Review & Editing: F.B., A.K., C.K.

663 Data and code availability

- Data and code used in this study are publicly available on the Data Server of the University of Bielefeld
- 665 (https://gitlab.ub.uni-bielefeld.de/felix.broehl/fb02).

666 References

667	Aiken, S.J., Picton, T.W., 2008. Human cortical responses to the speech envelope. Ear Hear. 29, 139–157.
668	https://doi.org/10.1097/AUD.0b013e31816453dc
669	Albouy, P., Benjamin, L., Morillon, B., Zatorre, R.J., 2020. Distinct sensitivity to spectrotemporal modulation
670	supports brain asymmetry for speech and melody. Science (80). 367, 1043–1047.
671	https://doi.org/10.1126/science.aaz3468
672	Altieri, N.A., Pisoni, D.B., Townsend, J.T., 2011. Some normative data on lip-reading skills (L). J. Acoust. Soc.
673	Am. 130, 1–4. https://doi.org/10.1121/1.3593376
674	Arnal, L.H., Morillon, B., Kell, C.A., Giraud, A.L., 2009. Dual neural routing of visual facilitation in speech
675	processing. J. Neurosci. 29, 13445–13453. https://doi.org/10.1523/JNEUROSCI.3194-09.2009
676	Bauer, A.K.R., Debener, S., Nobre, A.C., 2020. Synchronisation of Neural Oscillations and Cross-modal
677	Influences. Trends Cogn. Sci. https://doi.org/10.1016/j.tics.2020.03.003
678	Benjamini, Y., Hochberg, Y., 1995. Controlling the False Discovery Rate: A Practical and Powerful Approach
679	to Multiple Testing. J. R. Stat. Soc. Ser. B.
680	Bernstein, L.E., Jiang, J., Pantazis, D., Lu, ZL., Joshi, A., 2011. Visual phonetic processing localized using
681	speech and nonspeech face gestures in video and point-light displays. Hum. Brain Mapp. 32, 1660–
682	1676. https://doi.org/10.1002/hbm.21139
683	Bernstein, L.E., Liebenthal, E., 2014. Neural pathways for visual speech perception. Front. Neurosci.
684	https://doi.org/10.3389/fnins.2014.00386
685	Besle, J., Fischer, C., Bidet-Caulet, A., Lecaignard, F., Bertrand, O., Giard, M.H., 2008. Visual activation and
686	audiovisual interactions in the auditory cortex during speech perception: Intracranial recordings in
687	humans. J. Neurosci. 28, 14301–14310. https://doi.org/10.1523/JNEUROSCI.2875-08.2008
688	Bizley, J.K., Maddox, R.K., Lee, A.K.C., 2016. Defining Auditory-Visual Objects: Behavioral Tests and
689	Physiological Mechanisms. Trends Neurosci. 39, 74–85. https://doi.org/10.1016/j.tins.2015.12.007
690	Boersma, P., van Heuven, V., 2001. PRAAT, a system for doing phonetics by computer. Glot Int. 5, 341–347.
691	Bourguignon, M., Baart, M., Kapnoula, E.C., Molinaro, N., 2020. Lip-reading enables the brain to synthesize
692	auditory features of unknown silent speech. J. Neurosci. 40, 1053–1065.
693	https://doi.org/10.1523/JNEUROSCI.1101-19.2019
694	Bröhl, F., Kayser, C., 2021. Delta/theta band EEG differentially tracks low and high frequency speech-
695	derived envelopes. Neuroimage 233, 117958. https://doi.org/10.1016/j.neuroimage.2021.117958
696	Calvert, G.A., Bullmore, E.T., Brammer, M.J., Campbell, R., Williams, S.C.R., McGuire, P.K., Woodruff, P.W.R.,
697	Iversen, S.D., David, A.S., 1997. Activation of auditory cortex during silent lipreading. Science (80).
698	276, 593–596. https://doi.org/10.1126/science.276.5312.593
699	Calvert, G.A., Campbell, R., 2003. Reading speech from still and moving faces: The neural substrates of
700	visible speech. J. Cogn. Neurosci. 15, 57–70. https://doi.org/10.1162/089892903321107828
701	Campbell, R., 2008. The processing of audio-visual speech: Empirical and neural bases. Philos. Trans. R. Soc.

702	B Biol. Sci. 363, 1001–1010. https://doi.org/10.1098/rstb.2007.2155
703	Capek, C.M., MacSweeney, M., Woll, B., Waters, D., McGuire, P.K., David, A.S., Brammer, M.J., Campbell, R.,
704	2008. Cortical circuits for silent speechreading in deaf and hearing people. Neuropsychologia 46,
705	1233–1241. https://doi.org/10.1016/j.neuropsychologia.2007.11.026
706	Chandrasekaran, C., Trubanova, A., Stillittano, S., Caplier, A., Ghazanfar, A.A., 2009. The natural statistics of
707	audiovisual speech. PLoS Comput. Biol. 5, e1000436. https://doi.org/10.1371/journal.pcbi.1000436
708	Cope, T.E., Sohoglu, E., Sedley, W., Patterson, K., Jones, P.S., Wiggins, J., Dawson, C., Grube, M., Carlyon,
709	R.P., Griffiths, T.D., Davis, M.H., Rowe, J.B., 2017. Evidence for causal top-down frontal contributions
710	to predictive processes in speech perception. Nat. Commun. 8, 2154. https://doi.org/10.1038/s41467-
711	017-01958-7
712	Crosse, M.J., Butler, J.S., Lalor, E.C., 2015. Congruent Visual Speech Enhances Cortical Entrainment to
713	Continuous Auditory Speech in Noise-Free Conditions. J. Neurosci. 35, 14195–204.
714	https://doi.org/10.1523/JNEUROSCI.1829-15.2015
715	Daube, C., Ince, R.A.A., Gross, J., 2019. Simple Acoustic Features Can Explain Phoneme-Based Predictions of
716	Cortical Responses to Speech. Curr. Biol. 29, 1924-1937.e9.
717	https://doi.org/10.1016/J.CUB.2019.04.067
718	Di Liberto, G.M., Lalor, E.C., Millman, R.E., 2018. Causal cortical dynamics of a predictive enhancement of
719	speech intelligibility. Neuroimage 166, 247–258. https://doi.org/10.1016/j.neuroimage.2017.10.066
720	Ding, N., Simon, J.Z., 2013. Adaptive temporal encoding leads to a background-insensitive cortical
721	representation of speech. J. Neurosci. 33, 5728–35. https://doi.org/10.1523/JNEUROSCI.5297-12.2013
722	Doelling, K.B., Arnal, L.H., Ghitza, O., Poeppel, D., 2014. Acoustic landmarks drive delta-theta oscillations to
723	enable speech comprehension by facilitating perceptual parsing. Neuroimage 85, 761–768.
724	https://doi.org/10.1016/j.neuroimage.2013.06.035
725	Etard, O., Reichenbach, T., 2019. Neural Speech Tracking in the Theta and in the Delta Frequency Band
726	Differentially Encode Clarity and Comprehension of Speech in Noise. J. Neurosci. 39, 5750–5759.
727	https://doi.org/10.1523/JNEUROSCI.1828-18.2019
728	Ghitza, O., 2017. Acoustic-driven delta rhythms as prosodic markers. Lang. Cogn. Neurosci. 32, 545–561.
729	https://doi.org/10.1080/23273798.2016.1232419
730	Giordano, B.L., Ince, R.A.A., Gross, J., Schyns, P.G., Panzeri, S., Kayser, C., 2017. Contributions of local
731	speech encoding and functional connectivity to audio-visual speech perception. Elife 6.
732	https://doi.org/10.7554/eLife.24763
733	Giraud, A.L., Poeppel, D., 2012. Cortical oscillations and speech processing: Emerging computational
734	principles and operations. Nat. Neurosci. 15, 511–517. https://doi.org/10.1038/nn.3063
735	Grant, K.W., Seitz, PF., 2000. The use of visible speech cues for improving auditory detection of spoken
736	sentences. J. Acoust. Soc. Am. 108, 1197. https://doi.org/10.1121/1.1288668
737	Gross, J., Hoogenboom, N., Thut, G., Schyns, P., Panzeri, S., Belin, P., Garrod, S., 2013. Speech Rhythms and

738	Multiplexed Oscillatory Sensory Coding in the Human Brain. PLoS Biol. 11, e1001752.
739	https://doi.org/10.1371/journal.pbio.1001752
740	Haegens, S., Zion Golumbic, E., 2018. Rhythmic facilitation of sensory processing: A critical review.
741	Neurosci. Biobehav. Rev. 86, 150–165. https://doi.org/10.1016/j.neubiorev.2017.12.002
742	Hagerman, B., 1982. Sentences for Testing Speech Intelligibility in Noise. Scand. Audiol. 11, 79–87.
743	https://doi.org/10.3109/01050398209076203
744	Hauswald, A., Lithari, C., Collignon, O., Leonardelli, E., Weisz, N., 2018. A Visual Cortical Network for
745	Deriving Phonological Information from Intelligible Lip Movements. Curr. Biol. 28, 1453-1459.e3.
746	https://doi.org/10.1016/j.cub.2018.03.044
747	Ince, R.A.A., Giordano, B.L., Kayser, C., Rousselet, G.A., Gross, J., Schyns, P.G., 2017. A statistical framework
748	for neuroimaging data analysis based on mutual information estimated via a gaussian copula. Hum.
749	Brain Mapp. 38, 1541–1573. https://doi.org/10.1002/hbm.23471
750	Keitel, A., Gross, J., Kayser, C., 2020. Shared and modality-specific brain regions that mediate auditory and
751	visual word comprehension. Elife 9, 1–23. https://doi.org/10.7554/ELIFE.56972
752	Keitel, A., Gross, J., Kayser, C., 2018. Perceptually relevant speech tracking in auditory and motor cortex
753	reflects distinct linguistic features. PLoS Biol. 16, e2004473.
754	https://doi.org/10.1371/journal.pbio.2004473
755	Keitel, A., Ince, R.A.A., Gross, J., Kayser, C., 2017. Auditory cortical delta-entrainment interacts with
756	oscillatory power in multiple fronto-parietal networks. Neuroimage 147, 32–42.
757	https://doi.org/10.1016/j.neuroimage.2016.11.062
758	Koike, K.J., Hurst, M.K., Wetmore, S.J., 1994. Correlation between the American Academy of
759	Otolaryngology-Head and Neck Surgery Five-Minute Hearing Test and standard audiologic data.
760	Otolaryngol Head Neck Surg. 111, 625–632. https://doi.org/10.1016/S0194-5998(94)70531-3
761	Kollmeier, B., Warzybok, A., Hochmuth, S., Zokoll, M.A., Uslar, V., Brand, T., Wagener, K.C., 2015. The
762	multilingual matrix test: Principles, applications, and comparison across languages: A review. Int. J.
763	Audiol. 54, 3–16. https://doi.org/10.3109/14992027.2015.1020971
764	Lakatos, P., Karmos, G., Mehta, A.D., Ulbert, I., Schroeder, C.E., 2008. Entrainment of neuronal oscillations
765	as a mechanism of attentional selection. Science 320, 110–3.
766	https://doi.org/10.1126/science.1154735
767	Lazard, D.S., Giraud, A.L., 2017. Faster phonological processing and right occipito-temporal coupling in deaf
768	adults signal poor cochlear implant outcome. Nat. Commun. 8, 1–9.
769	https://doi.org/10.1038/ncomms14872
770	Ludman, C.N., Summerfield, A.Q., Hall, D., Elliott, M., Foster, J., Hykin, J.L., Bowtell, R., Morris, P.G., 2000.
771	Lip-reading ability and patterns of cortical activation studied using fMRI. Br. J. Audiol. 34, 225–230.
772	https://doi.org/10.3109/03005364000000132
773	Luo, H., Liu, Z., Poeppel, D., 2010. Auditory Cortex Tracks Both Auditory and Visual Stimulus Dynamics Using

774	Low-Frequency Neuronal Phase Modulation. PLoS Biol. 8, e1000445.
775	https://doi.org/10.1371/journal.pbio.1000445
776	Mégevand, P., Mercier, M.R., Groppe, D.M., Golumbic, E.Z., Mesgarani, N., Beauchamp, M.S., Schroeder,
777	C.E., Mehta, A.D., 2020. Crossmodal Phase Reset and Evoked Responses Provide Complementary
778	Mechanisms for the Influence of Visual Speech in Auditory Cortex. J. Neurosci. 405597.
779	https://doi.org/10.1101/405597
780	Meyer, L., Henry, M.J., Gaston, P., Schmuck, N., Friederici, A.D., 2017. Linguistic bias modulates
781	interpretation of speech via neural delta-band oscillations. Cereb. Cortex 27, 4293–4302.
782	https://doi.org/10.1093/cercor/bhw228
783	Munhall, K.G., Jones, J.A., Callan, D.E., Kuratate, T., Vatikiotis-Bateson, E., 2004. Visual Prosody and Speech
784	Intelligibility: Head Movement Improves Auditory Speech Perception. Psychol. Sci. 15, 133–137.
785	https://doi.org/10.1111/j.0963-7214.2004.01502010.x
786	Nichols, T., Holmes, A., 2003. Nonparametric Permutation Tests for Functional Neuroimaging. Hum. Brain
787	Funct. Second Ed. 25, 887–910. https://doi.org/10.1016/B978-012264841-0/50048-2
788	Nidiffer, A.R., Cao, C.Z., O'Sullivan, A.E., Lalor, E.C., 2021. A linguistic representation in the visual system
789	underlies successful lipreading. bioRxiv. https://doi.org/10.1101/2021.02.09.430299
790	O'Sullivan, A.E., Crosse, M.J., Di Liberto, G.M., Lalor, E.C., 2017. Visual cortical entrainment to motion and
791	categorical speech features during silent lipreading. Front. Hum. Neurosci. 10, 1–11.
792	https://doi.org/10.3389/fnhum.2016.00679
793	Obleser, J., Kayser, C., 2019. Neural Entrainment and Attentional Selection in the Listening Brain. Trends
794	Cogn. Sci. 23, 913–926. https://doi.org/10.1016/j.tics.2019.08.004
795	Oganian, Y., Chang, E.F., 2019. A speech envelope landmark for syllable encoding in human superior
796	temporal gyrus. Sci. Adv. 5, 1–14. https://doi.org/10.1126/sciadv.aay6279
797	Oldfield, R.C., 1971. The assessment and analysis of handedness: The Edinburgh inventory.
798	Neuropsychologia 9, 97–113. https://doi.org/10.1016/0028-3932(71)90067-4
799	Oostenveld, R., Fries, P., Maris, E., Schoffelen, JM., 2011. FieldTrip: Open source software for advanced
800	analysis of MEG, EEG, and invasive electrophysiological data. Comput. Intell. Neurosci. 2011, 156869.
801	https://doi.org/10.1155/2011/156869
802	Ozker, M., Yoshor, D., Beauchamp, M.S., 2018. Frontal cortex selects representations of the talker's mouth
803	to aid in speech perception. Elife 7, 1–14. https://doi.org/10.7554/eLife.30387
804	Park, H., Ince, R.A.A., Schyns, P.G., Thut, G., Gross, J., 2018. Representational interactions during
805	audiovisual speech entrainment: Redundancy in left posterior superior temporal gyrus and synergy in
806	left motor cortex. PLOS Biol. 16, e2006558. https://doi.org/10.1371/journal.pbio.2006558
807	Park, H., Kayser, C., Thut, G., Gross, J., 2016. Lip movements entrain the observers' low-frequency brain
808	oscillations to facilitate speech intelligibility. Elife 5. https://doi.org/10.7554/eLife.14521
809	Paulesu, E., Perani, D., Blasi, V., Silani, G., Borghese, N.A., De Giovanni, U., Sensolo, S., Fazio, F., 2003. A

810	functional-anatomical model for lipreading. J. Neurophysiol. 90, 2005–2013.
811	https://doi.org/10.1152/jn.00926.2002
812	Pekkola, J., Ojanen, V., Autti, T., Jääskeläinen, I.P., Möttönen, R., Tarkiainen, A., Sams, M., 2005. Primary
813	auditory cortex activation by visual speech: An fMRI study at 3 T. Neuroreport 16, 125–128.
814	https://doi.org/10.1097/00001756-200502080-00010
815	Plass, J., Brang, D., Suzuki, S., Grabowecky, M., 2020. Vision perceptually restores auditory spectral
816	dynamics in speech. Proc. Natl. Acad. Sci. U. S. A. 117, 16920–16927.
817	https://doi.org/10.1073/pnas.2002887117
818	Rauschecker, J.P., 2012. Ventral and dorsal streams in the evolution of speech and language. Front. Evol.
819	Neurosci. 4, 5–8. https://doi.org/10.3389/fnevo.2012.00007
820	Schroeder, C.E., Lakatos, P., 2009. Low-frequency neuronal oscillations as instruments of sensory selection.
821	Trends Neurosci. 32, 9–18. https://doi.org/10.1016/j.tins.2008.09.012
822	Schroeder, C.E., Lakatos, P., Kajikawa, Y., Partan, S., Puce, A., 2008. Neuronal oscillations and visual
823	amplification of speech. Trends Cogn. Sci. 12, 106–113. https://doi.org/10.1016/j.tics.2008.01.002
824	Scott, S.K., 2019. From speech and talkers to the social world: The neural processing of human spoken
825	language. Science (80). https://doi.org/10.1126/science.aax0288
826	Stefanics, G., Hangya, B., Hernadi, I., Winkler, I., Lakatos, P., Ulbert, I., 2010. Phase Entrainment of Human
827	Delta Oscillations Can Mediate the Effects of Expectation on Reaction Speed. J. Neurosci. 30, 13578-
828	13585. https://doi.org/10.1523/JNEUROSCI.0703-10.2010
829	Suess, N., Hauswald, A., Reisinger, P., Rösch, S., Keitel, A., Weisz, N., 2022. Cortical Tracking of Formant
830	Modulations Derived from Silently Presented Lip Movements and Its Decline with Age. Cereb. Cortex
831	1–16. https://doi.org/10.1093/cercor/bhab518
832	Teng, X., Tian, X., Doelling, K., Poeppel, D., 2018. Theta band oscillations reflect more than entrainment:
833	behavioral and neural evidence demonstrates an active chunking process. Eur. J. Neurosci. 48, 2770–
834	2782. https://doi.org/10.1111/ejn.13742
835	Teoh, E.S., Cappelloni, M.S., Lalor, E.C., 2019. Prosodic pitch processing is represented in delta-band EEG
836	and is dissociable from the cortical tracking of other acoustic and phonetic features. Eur. J. Neurosci.
837	50, 3831–3842. https://doi.org/10.1111/ejn.14510
838	Tivadar, R.I., Gaglianese, A., Murray, M.M., 2020. Auditory Enhancement of Illusory Contour Perception.
839	Multisens. Res. 34, 1–15. https://doi.org/10.1163/22134808-bja10018
840	Tivadar, R.I., Retsa, C., Turoman, N., Matusz, P.J., Murray, M.M., 2018. Sounds enhance visual completion
841	processes. Neuroimage 179, 480–488. https://doi.org/10.1016/j.neuroimage.2018.06.070
842	Tzourio-Mazoyer, N., Landeau, B., Papathanassiou, D., Crivello, F., Etard, O., Delcroix, N., Mazoyer, B., Joliot,
843	M., 2002. Automated anatomical labeling of activations in SPM using a macroscopic anatomical
844	parcellation of the MNI MRI single-subject brain. Neuroimage 15, 273–289.

846	van Bree, S., Sohoglu, E., Davis, M.H., Zoefel, B., 2020. Sustained neural rhythms reveal endogenous
847	oscillations supporting speech perception, PLoS Biology. https://doi.org/10.1101/2020.06.26.170761
848	Van Veen, B.D., Van Drongelen, W., Yuchtman, M., Suzuki, A., 1997. Localization of brain electrical activity
849	via linearly constrained minimum variance spatial filtering. IEEE Trans. Biomed. Eng. 44, 867–880.
850	https://doi.org/10.1109/10.623056
851	Velleman, P.F., Welsch, R.E., 1981. Efficient computing of regression diagnostics. Am. Stat. 35, 234–242.
852	https://doi.org/10.1080/00031305.1981.10479362
853	Yu, C., Zhou, Y., Liu, Y., Jiang, T., Dong, H., Zhang, Y., Walter, M., 2011. Functional segregation of the human
854	cingulate cortex is confirmed by functional connectivity based neuroanatomical parcellation.
855	Neuroimage 54, 2571–2581. https://doi.org/10.1016/j.neuroimage.2010.11.018
856	Zion Golumbic, E., Cogan, G.B., Schroeder, C.E., Poeppel, D., 2013. Visual input enhances selective speech
857	envelope tracking in auditory cortex at a "Cocktail Party." J. Neurosci. 33, 1417–1426.
858	https://doi.org/10.1523/JNEUROSCI.3675-12.2013
859	Zuk, N.J., Murphy, J.W., Reilly, R.B., Lalor, E.C., 2021. Envelope reconstruction of speech and music
860	highlights stronger tracking of speech at low frequencies, PLOS Computational Biology.
861	https://doi.org/10.1371/journal.pcbi.1009358



eNeuro Accepted Manuscript



eNeuro Accepted Manuscript



AudFeat LipFeat





