1    **Left frontal motor delta oscillations reflect the temporal integration of multimodal speech**

2

3    Emmanuel Biau [1,2], Benjamin G. Schultz [2], Thomas C. Gunter [3] and Sonja A. Kotz [2,3].

4

5    [1] School of Psychology, University of Birmingham, B15 2TT, Birmingham, UK.

6    [2]Basic and Applied NeuroDynamics Laboratory, Department of Neuropsychology and

7    Psychopharmacology, University of Maastricht, 6200 MD, Maastricht, Netherlands.

8    [3]Department of Neuropsychology, Max Planck Institute for Human Cognitive and Brain Sciences,

9    04103, Leipzig, Germany.

10

11   Corresponding authors:

12   Dr. Emmanuel Biau

13   Address: School of Psychology, University of Birmingham, B15 2TT, Birmingham, UK.

14   Email: e.biau@bham.ac.uk

15

16   Prof. Dr. Sonja A. Kotz

17   Address: Basic and Applied NeuroDynamics Laboratory, Department of Neuropsychology and

18   Psychopharmacology, University of Maastricht, 6200 MD, Maastricht, Netherlands.

19   Email: sonja.kotz@maastrichtuniversity.nl

20

21

22

23

24

25

26

27

28

29

1

30    **ABSTRACT**

31        During multimodal speech perception, slow delta oscillations (~1 - 3 Hz) in the listener's brain

32    synchronize with speech signal, likely reflecting signal decomposition at the service of

33    comprehension. In particular, fluctuations imposed onto the speech amplitude envelope by a

34    speaker's prosody seem to temporally align with articulatory and body gestures, thus providing

35    two complementary sensations to the speech signal's temporal structure. Further, endogenous

36    delta oscillations in the left motor cortex align with speech and music beat, suggesting a role in

37    the temporal integration of (quasi)-rhythmic stimulations. We propose that delta activity

38    facilitates the temporal alignment of a listener's oscillatory activity with the prosodic fluctuations

39    in a speaker's speech during multimodal speech perception. We recorded EEG responses in an

40    audiovisual synchrony detection task while participants watched videos of a speaker. To test the

41    temporal alignment of visual and auditory prosodic features, we filtered the speech signal to

42    remove verbal content. Results confirm (i) that participants accurately detected audiovisual

43    synchrony, and (ii) greater delta power in left frontal motor regions in response to audiovisual

44    asynchrony. The latter effect correlated with behavioural performance, and (iii) decreased delta-

45    beta coupling in the left frontal motor regions when listeners could not accurately integrate visual

46    and auditory prosodies. Together, these findings suggest that endogenous delta oscillations align

47    fluctuating prosodic information conveyed by distinct sensory modalities onto a common

48    temporal organisation in multimodal speech perception.

49    **KEYWORDS**

51

52

53

54

55

56

57

**INTRODUCTION**

Speaker prosody displays perceptible fluctuations in the speech amplitude envelope that allow a listener to segment and parse incoming speech (Ghitza, 2017). While not isochronous, prosody imposes a temporal structure with regular alterations of strong and weak accentuated cues occurring at ~1 - 3 Hz delta rate (Ding et al., 2016; Doelling et al., 2014; Ghitza, 2017; Pell & Davis, 2012). Delta oscillation responses track and align with these prosodic events in auditory cortex to extract the temporal structure of speech (i.e. "neural entrainment"; Giraud & Poeppel, 2012; Keitel et al., 2017; Kosem & van Wassenhove, 2017; Meyer, Sun & Martin, 2019). Beyond segmentation, prosody is present in the visual and auditory domains and may facilitate the listener's brain activity to synchronize with multimodal information in social interactions (Esteve-Gibert & Guellaï, 2018; Kotz, Ravignani & Fitch, 2018). The term "visual prosody" encompasses communicative gestures (i.e. hand, head, face, and body movements) whose prominent phase temporally coincides with acoustic prosodic anchors such as intonational phrases, pitch accents, and boundary tones (Biau et al., 2016; Chandrasekaran et al., 2009; Munhall et al., 2004; Wagner et al., 2014). For example, the famous cocktail party effect illustrates how listeners rely on temporal alignment between gestures and sounds to improve speech perception (Cherry, 1953; Obermeier, Dolk & Gunter, 2012; Sumby & Pollack, 1954). Together these facts raise the following question: How does the brain integrate multiple dynamic visual and auditory prosody streams to facilitate multimodal speech perception?

The present study investigated how delta oscillations mark the temporal integration of audiovisual prosody in speech. We refer to temporal integration as the mechanism that integrates visual and auditory prosody, leading to the improved perception of their temporal representation in speech. Delta activity in the motor cortex has been associated with the temporal integration of rhythmic stimuli including speech, as its phase aligns with the onsets of predictable events (Morillon et al., 2019; Morillon & Schroeder, 2015; Saleh et al., 2010). Keitel et al. (2018) showed that left motor delta activity tracked temporally predictable slow phrasal features in auditory sentences and predicted speech comprehension. This suggests that this region integrates perceptually relevant regularities in the signal to facilitate comprehension. The

86    authors also found delta-beta cross-frequency coupling in the left motor region, in line with
87    previous research showing that motor beta oscillations also respond to the temporal integration
88    of rhythmic auditory tones or visual cues stimulations (Fujioka, Ross & Trainor, 2015; Saleh et al.,
89    2010). These findings led to the hypothesis that motor delta oscillations are involved in the
90    temporal integration of speech by mediating top-down control through cross-frequency coupling
91    with beta activity (Arnal, 2012; Arnal, Doelling & Poeppel, 2015; Morillon and Baillet, 2017). In
92    other words, delta activity could reflect how the brain gathers multiple temporal representations
93    of input across modalities in the left motor cortex, and generates predictions to improve signal
94    processing. Finally, the left motor cortex including the left inferior frontal gyrus is involved in
95    gestures and speech integration, making it a critical candidate of the present study (Biau et al.,
96    2016; Park et al., 2016; Zhao et al., 2018).

97    We propose that visual and auditory prosodic cues encoded in the visual and auditory sensory
98    cortices respectively, provide two representations of the speech signal that are integrated in the
99    left motor in speech perception. Such crossmodal temporal integration is reflected by delta
100   responses in this region. To test this hypothesis, we manipulated the temporal alignment of
101   filtered speech with corresponding whole body or masked head movements. Participants
102   performed an audiovisual synchrony detection task by attending to videos of a single speaker
103   engaged in a conversation, while we recorded their electroencephalogram (EEG). First, we tested
104   whether listeners relate the respective temporal structures of visual and auditory prosodic
105   features in multimodal speech. Second, we explored delta oscillations in response to audiovisual
106   asynchrony in multimodal stimuli. Third, we tested whether delta-beta coupling in the left motor
107   cortex predicted multimodal synchrony detection in speech perception.

108   **METHODS**

109   **Participants**

110   Twenty-six native Dutch speakers (mean age = 22.24, SD = 4.24; 15 females) were recruited at
111   Maastricht University and received €10 for participating in the experiment after giving informed
112   consent. All participants were right-handed and had normal or corrected-to-normal vision and

4

113 hearing. The protocol of the study was approved by the Research Ethical Committee of
114 Maastricht University. Data from three participants were removed from the final analysis due to
115 technical problems.

116 **Stimuli**

117 Short videos were extracted from a longer video recording used in a previous study (Gunter &
118 Weinbrenner, 2017). The videos depicted a female actor and an experimenter (both German
119 native speakers) engaged in a question-answer conversation. The actor sat on a chair, moved
120 freely, and was visible from her knees up to the top of her head. Relevant segments containing
121 the actor' answers separate from the experimenter were selected to create the current stimulus
122 set (N = 54). Each of the 54 segments was 10 seconds long (600 frames at 60 frames per second;
123 FPS). The audio track was extracted to be low-pass filtered with Hann band windowing procedure
124 (from 0 Hz to 400 Hz; 20 Hz smoothing) using Praat (Boersma & Weenink, 2015). In doing so, we
125 altered speech intelligibility removing verbal content while keeping the prosodic contour of the
126 signal. Peak frequencies were extracted from the audio and video files through Fourier
127 transformations that calculated the frequency at which the peak amplitude occurred within a
128 range of 0.5Hz to 8Hz. For videos, the average magnitude of grayscale pixel changes between
129 consecutive frames was used to determine the frequency of movement and gesture (see Table
130 1).

131

| Row Labels | Mean Peak Freq. | SD Peak Freq. | Min. Peak Freq. | Max Peak Freq. |
|---|---|---|---|---|
| **Audio*** | 2.74 | 1.44 | 0.86 | 5.86 |
| **Body*** | 3.37 | 1.25 | 0.86 | 6.13 |
| **Head-mask Face** | 2.27 | 1.53 | 0.86 | 6.13 |
| **Head-mask Full** | 3.10 | 1.45 | 0.86 | 6.13 |
| **No-mask Face** | 3.59 | 0.91 | 1.00 | 3.99 |
| **No-mask Full** | 3.65 | 1.02 | 0.86 | 6.13 |

132

133 Table 1. Summary statistics of peak frequencies obtained for video and audio signals using a
134 Fourier transformation (*Features remaining identical across the no mask and head mask
135 conditions).

5

136 We applied two visual manipulations to each of the 54 speech segments: (1) The presence or

137 absence of a visual mask (no mask, head-mask), and (2) the original temporal alignment of the

138 audiovisual information or a temporal shift of the audio signal relative to the video onset

139 (synchronous, asynchronous). In the no-mask condition, the speaker's body and face were fully

140 visible. In the head-mask condition, the head of the speaker was blurred to degrade visual

141 prosody conveyed by the speaker's lips. The mask was created by applying a low-pass Gaussian

142 filter on the upper third of the original video containing the speaker's face, attenuating a high

143 frequency signal. This manipulation removed fine-grained facial expressions from the video while

144 slow gestures remained intact (see Figure 1). In the synchronous condition, the original temporal

145 alignment between visual and auditory onsets was intact. To create an asynchronous condition,

146 we inserted a delay between the visual and auditory onsets by shifting the sound onset by +400

147 ms relative to the video onset (i.e., 24 frames). This manipulation maintained the natural order

148 of visual information preceding auditory information. A 400ms lag was used to ensure that the

149 delay was long enough to detect audiovisual asynchrony; it was also based on the time-window

150 of multisensory integration established in previous studies (Biau et al., 2016; Biau & Soto-Faraco,

151 2013; Jessen & Kotz, 2015; Obermeier & Gunter, 2014). Further, a central white fixation cross

152 was displayed in each video to allow participants to focus their gaze on a central cue while

153 attending audiovisual stimuli. Altogether, this created four conditions: Head-Mask Synchronous

154 (HMS), Head-Mask Asynchronous (HMA), No-Mask Synchronous (NMS), and No-Mask

155 Asynchronous (NMA) (see Figure 1A). 18 additional video clips, in which the central white

156 fixation-cross turned red, were used as fillers, counterbalanced across conditions (colour change

157 onset jittered between 5 and 9 seconds after the video onset; ~ 8 % of total stimuli, not included

158 in the final data analysis). We used the fillers in a memory test to focus the participants' attention

159 on the videos during the experiment. Finally, audio files were recombined with their

160 corresponding video files for each condition. Videos were edited using Adobe Premiere Pro CS3

161 and exported using the following parameters: Pixel resolution 1920 × 1080, 60 FPS compressor

162 Indeo video 5.10, AVI format, audio sample rate 48 kHz, 16 bits Mono.

163 **Apparatus**

6

164    The audio files were presented through EEG-compatible air tubes (ER3C Tubal Insert Earphones,

165    Etymotic Research). Videos were presented on a 27 inch Iiyama G-MASTER (GB2760HSU-B1) TN

166    display with a 1ms response time, a refresh rate of 144Hz, and a native resolution of 1920 x 1080

167    pixels connected to the stimulus presentation computer (Intel i7-6700 CPU @ 3.40 GHz, 32 GB,

168    running 64-bit Windows 7, NVIDIA GeForce FTX 1080 GTX GPU). Stimuli were presented using a

169    custom MATLAB script (MATLAB and Statistics Toolbox Release 2015b, The MathWorks, Inc.,

170    Natick, Massachusetts, United States) that called VideoLAN Client (VLC; VideoLAN Client, 2017;

171    http://www.videolan.org/) to play the videos. EEG data were collected using BrainVision

172    Recorder (Brain Products, GmbH, 2017) software on an Intel Xeon E5-1650 PC (3.5 GHz, 32GB

173    RAM) running Windows 7. Video onsets were synchronized to EEG data using the Schultz

174    Cigarette Burn Toolbox (Schultz, Biau, & Kotz, 2020).

175    **Procedure**

176    Participants were seated approximately 60 cm apart from the monitor in a sound attenuated

177    booth while videos were displayed on a computer screen. Participants watched 234 videos

178    organised in nine blocks of 26 randomised trials (i.e., 6 stimuli per condition + 2 fillers). The task

179    was a two-alternative forced choice synchrony detection task (Figure 1B). Participants attended

180    both the audio and video stimuli. Each trial began with a central white fixation cross (jittered

181    duration 500 +/- 250 ms), followed by the stimulus. After the video ended, participants decided

182    whether the audio and the video signals were synchronous or asynchronous by pressing the "1"

183    or "2" key on the keyboard without time pressure (counterbalanced across participants).

184    Additionally, participants were asked to count internally the number of times they observed a

185    red cross in a video clip, and reported it at the end of the experiment. Filler trials were not

186    included in behavioural and EEG analyses but the total number of reported red crosses served to

187    check that attention was maintained throughout the experiment. Before the experiment,

188    participants received five practice trials where they were presented with one example of each

189    condition to ensure they understood the instructions. At the end of the experiment, participants

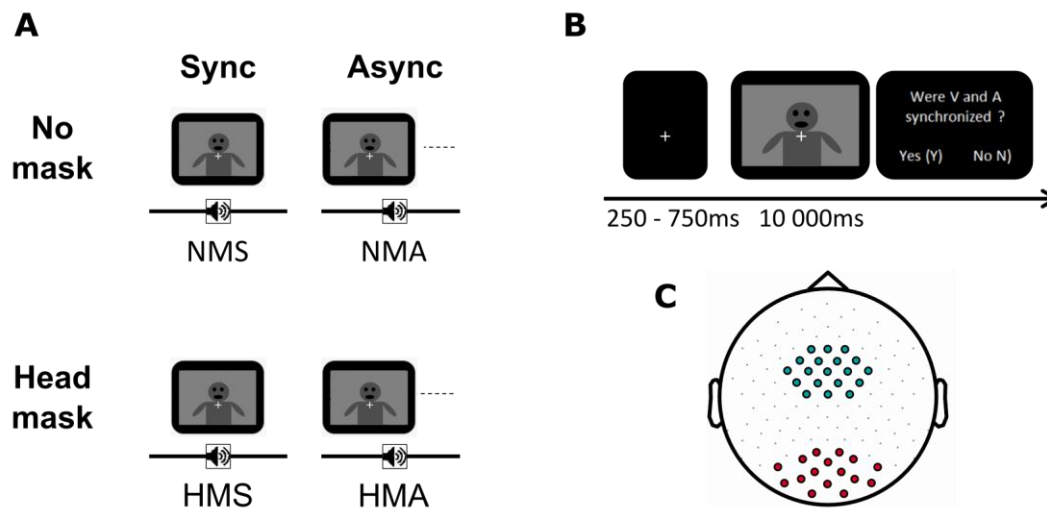190    were asked if they could identify the speaker's language and to report it.

191

*Figure 1.* Experimental procedure of the audiovisual synchrony detection task. (A) For each item, the audio signal was the same across all four versions. Visual information was manipulated for the factor mask (no-mask or head-mask); audiovisual stimuli were temporally aligned in the synchronous conditions (NMS, HMS), and the audio signal was temporally delayed (400ms) in the asynchronous conditions (NMA, HMA). (B) Example of one trial timeline. (C) Distribution of the electrodes covering the motor region of interest (ROI; blue circles) and the control region of non-interest in the visual area (RONI; red circles). N.B: The images in A and B have been modified for anonymity purpose here.

**EEG recording and preprocessing**

Electrophysiological data were recorded at 1000 Hz with 128 active electrodes (ActiCap, Brain Vision Recorder, Brain Products) according to the 10-20 international standard, and impedances were kept below 10 kΩ. The ground electrode was located at AFz, and the reference electrode was placed at the right mastoid (TP10).

Offline EEG preprocessing: EEG data were preprocessed offline using Fieldtrip (Oostenveld et al., 2011) and SPM8 toolboxes (Wellcome Trust Centre for Neuroimaging). Continuous EEG signals

8

208   were bandpass filtered between 1Hz and 100 Hz and bandstop filtered (48-52 Hz and 98-102 Hz)

209   to remove line noise at 50 and 100 Hz. Data were epoched from 1000 ms before stimulus onset

210   to 11000 ms after stimulus onset. Trials and channels with artefacts were excluded by visual

211   inspection before applying an independent component analysis (ICA) to remove components

212   related to ocular artefacts. Excluded channels were then interpolated using the method of

213   triangulation of nearest. After re-referencing the data to an average reference, the remaining

214   trials with artefacts were manually rejected by a final visual inspection (on average, 13.57 ± 8.32

215   trials across conditions per participant).

216   **EEG data analyses at the scalp level**

217   For each participant, time-frequency representations (TFRs) were computed using a Morlet

218   wavelet (width: 5 cycles) from 1 to 40 Hz (1 Hz step), with 20 ms time steps. Power in the hit trials

219   (i.e. correct synchrony detection in synchronous conditions and correct asynchrony detection in

220   the asynchronous condition) was calculated first and then averaged across trials in the four

221   conditions. The power was normalised relative to a pre-stimulus baseline (-700 to -200 ms with

222   respect to stimulus onset) to determine increases or decreases of power dependent on the

223   conditions. As entrainment necessitates several cycles from recurrent stimulations to build up

224   (Doelling et al. 2014; Thut et al., 2011; Zoefel et al. 2018) and the slower frequency in our band

225   of interest was 2Hz (corresponding to a period of 500 ms), we defined a time window of interest

226   from + 3 to + 9 seconds after stimulus onset. This time window ensured that neural activity

227   sufficiently entrained to the temporal structure of the stimuli, and that the responses evoked by

228   the stimulus onset-offsets did not influence the results. In the identified regions of interest and

229   non-interest (see Results section), normalised mean power across pool electrodes in the 2-3 Hz

230   frequency band was computed for the four conditions and exported for further statistical

231   analyses.

232   **EEG data analyses at the source level**

233   <u>Source localisation</u>: We used the Montreal Neurological Institute (MNI) MRI template and a

234   template volume conduction model from Fieldtrip. The 128 electrode positions on the

9

235    volunteer's head were defined by using a Polhemus FASTRAK device (Colchester), recorded with

236    the Brainstorm toolbox implemented in MATLAB (Tadel et al., 2011), and realigned to the

237    template head model using Fieldtrip. The template volume conduction model and the electrode

238    template were used to prepare the source models. Leadfields were computed based on scalp

239    potentials and source activity was reconstructed applying a linearly constrained minimum

240    variance (LCMV) beamforming approach implemented in Fieldtrip (van Veen et al., 1997; Wang

241    et al., 2018). Source analyses were run on potential data (i.e., average referenced) and time-

242    series data were reconstructed on 2020 virtual electrodes for each participant. Time-frequency

243    analysis was computed at each of 2020 virtual sources with the exact same approach to scalp

244    level analyses. The maximum voxel activation regions were defined by using the automated

245    anatomical labelling atlas (AAL).

246    <u>Phase-amplitude coupling (PAC) between delta and beta oscillations</u>: We applied a modulation

247    index (MI) analysis in the time-window of interest to quantify delta-beta PAC in the significant

248    cluster revealed by source localisation in the NMA-NMS contrast (Tort et al., 2010). First, the

249    power spectrum (1 - 30 Hz) was estimated across all grids of the significant cluster and trials by

250    applying a 1/f correction time-frequency decomposition method with wavelet for each

251    participant (Griffith et al., 2019). The most prominent power spectrum peaks in the delta and

252    beta bands were then extracted and saved as the individual delta and beta peaks. Across

253    participants, the mean delta peak was at 2.38 Hz and the mean beta peak was at 24.16 Hz. The

254    number of trials were balanced by identifying the condition with the smallest number of incorrect

255    trials, and taking 80 percent of the smaller sample for all the conditions (NMA$_{hit}$, NMA$_{miss}$, NMS$_{hit}$

256    and NMs$_{miss}$ ; HMA$_{hit}$, HMA$_{miss}$, HMS$_{hit}$ and HMs$_{miss}$). Subsampled trials were concatenated and

257    the operation was repeated for 50 iterations in each condition (Keitel et al., 2018). The grids of

258    interest resulted from the significant cluster identified in the left frontal motor cortex by the

259    source localisation analysis (contrast NMA-NMS; number of significant grids = 92). Phase and

260    power were derived from Hilbert-transformed time series and filtered around the delta peak (±

261    0.5 Hz) and beta peak (± 5 Hz) based on the frequency window. For each trial and grid source,

262    beta power was binned into 12 equidistant bins of 30° according to the delta phase. The MI was

263    computed by comparing the observed distribution to a uniform distribution for each trial and

264    grid. The PAC was then averaged across the left frontal motor grids and 50 iterations in each

265    condition. Finally, we investigated whether the delta-beta coupling was specifically localised in

266    the region of interest, identified by the source localisation analysis (i.e., left frontal motor area).

267    We compared the delta-beta PAC between masks (no-mask, head-mask) based on the results

268    from the cluster of interest. In contrast to the PAC analysis in the region of interest, the difference

269    of trial numbers between conditions was balanced by taking 80 percent of the smaller sample

270    between all the conditions. Subsampled trials were concatenated and the operation was

271    repeated for 40 iterations (to circumvent computational resource limits reached by concatenated

272    epoch lengths). The delta-beta PAC was then averaged across all iterations at each grid (n = 2020)

273    and conditions across participants.

274    **Experimental design and statistical analysis**

275    <u>Audiovisual synchrony detection task</u>: The experiment used a full within-subject design. The

276    effect of asynchrony and its interaction with the head-mask in audiovisual integration was

277    assessed by means of $d'$ sensitivity index and reaction times for correct trials. The correct

278    responses (Hits and correct rejections in the synchronous conditions) and errors (misses and false

279    alarms FA in the asynchronous conditions) as well as the reaction times of the hits (comprised

280    between mean reaction times ± two standard deviations range), were computed in each

281    condition for each participant. Then, the d' scores for synchrony detection in the no-mask and

282    head-mask contrasts were calculated for each participant as follows: $d' = Z (Hit_{rate}) - Z (FA_{rate})$.

283    The $d'$ index allows taking into account response bias by comparing hits and false alarms to assess

284    whether participants actually discriminated synchrony and asynchrony. Additionally, the decision

285    criterion $c$ was computed as follows: $c = 0.5 \times (Hit_{rate} - FA_{rate}) / 2$ to determine the decision shift

286    between no-mask and head-mask contrasts. The effects of masking the speaker's face and

287    audiovisual synchrony on reaction times were assessed using two-way repeated-measure

288    ANOVAs with the factors mask (no-mask, head-mask), synchrony (synchronous, asynchronous),

289    and the interaction between mask and synchrony, using SPSS (IBM Corp. Released 2015. IBM

290    SPSS Statistics for Windows, Version 23.0. Armonk, NY: IBM Corp.). In the case of significant

291    interactions, *post-hoc t*-tests were Bonferroni-corrected. To test for the modulation of sensitivity

11

292    to synchrony depending on speaker's face information, the *d'* and *c* criterion in no-mask and

293    head-mask contrasts were individually tested against zero by means of one-sample *t*-tests.

294    Further, the difference of *d'* between the no-mask and head-mask contrasts was assessed

295    applying a paired-samples *t*-test and the effect size was defined using Cohen's *d*.

296    EEG data at the scalp level: EEG data of correct trials at the scalp and source levels were

297    statistically analysed. The differences of mean power between two contrasts (NMA-NMS and

298    HMA-HMS) at the electrode level were statistically assessed by applying dependent *t*-tests using

299    Monte-Carlo cluster-based permutation tests (Maris & Oostenveld, 2007) with an alpha cluster-

300    forming threshold set at 0.05, three minimum neighbour channels, 5000 iterations, and cluster

301    selection based on maximum size. Cluster-based permutation statistics were applied for the time

302    window of interest in the delta 2-3 Hz band across all the electrodes. To address whether centro-

303    frontal delta oscillations responses reflect temporal speech analysis rather than pure signal

304    processing via neural entrainment, we performed the same tests on the theta band (4 - 8 Hz),

305    which plays a role in the integration of the syllabic structure of speech (Giraud & Poeppel, 2012).

306    We expected to find modulations of delta but not theta oscillations for audiovisual asynchrony

307    in the region of interest if motor delta responses reflect temporal integration. In the identified

308    regions of interest and non-interest (see Results section), normalised mean power across pool

309    electrodes in the 2-3 Hz delta and 4-8 Hz theta frequency bands was computed for the four

310    conditions and exported. Statistical differences of power in relevant contrasts were assessed by

311    means of two-way repeated-measure ANOVAs.

312    EEG data source localisation: Differences in delta power for the two contrasts (no-mask: NMA-

313    NMS and head-mask: HMA-HMS) were assessed by applying dependent *t*-tests using Monte-

314    Carlo cluster-based permutation tests as performed for the scalp level analyses. For visualisation

315    of the source localisation results, the power differences in the two contrasts were grand-

316    averaged across participants, and the grand average power differences were interpolated to the

317    MNI MRI template for visualization. Only voxels surpassing the statistical significance threshold

318    are depicted in both contrasts (significant *t*-values at alpha = 0.05, multiple comparison cluster-

319    corrected).

320     <u>Delta-beta PAC</u>: First, statistical differences of mean PAC across conditions in the region of

321     interest were assessed applying a three-way repeated-measure ANOVA with the factors mask

322     (no-mask and head-mask), synchrony (synchronous and asynchronous), and correctness (correct

323     and incorrect trials). Second, statistical differences of whole brain delta-beta PAC was assessed

324     by applying dependent *t*-tests using Monte-Carlo cluster-based permutation tests as described

325     above (whereas *t*-tests were one-tailed here as we had a strong hypothesis about delta-beta PAC

326     modulation directionality based on results at region of interest level).

327     <u>Correlations between synchrony detection performance and delta oscillations in the identified</u>

328     <u>left motor cluster</u>: We addressed whether delta responses in the left motor area correlated with

329     audiovisual synchrony perception. We performed Pearson correlations between the difference

330     of mean power in the 2-3 Hz band at source level, and the difference of correct response rates

331     within contrast. For each participant, we computed the 2-3 Hz power at the grids sources from

332     the significant cluster established in the NMA-NMS contrast source analysis (all significant grids

333     were situated in the left central and frontal gyrus areas; n = 92). Power was averaged across the

334     92 grids in the four conditions separately (NMS, NMA, HMS and HMA), and we calculated the

335     mean difference separately in the no-mask (NMA-NMS) and head-mask (HMA-HMS) contrasts to

336     obtain two delta power values per participant. Similarly, the difference of correct response rates

337     was calculated in the no-mask and head-mask contrasts, resulting in two behaviour values per

338     participant. The statistical relationship between behaviour and delta power was assessed

339     applying Pearson's correlation tests.

340     **RESULTS**

341     Participants reported 18.26 ± SD = 1.51 red crosses (out of 18) at the end of the experiment.

342     Additionally, they correctly identified the speaker's native language (they all responded

343     "German"), although they could not report any semantic content. These results confirmed that

344     participants correctly paid attention to both the audio and video signals.

345     **Listeners successfully temporally aligned visual and auditory prosodic features to achieve**
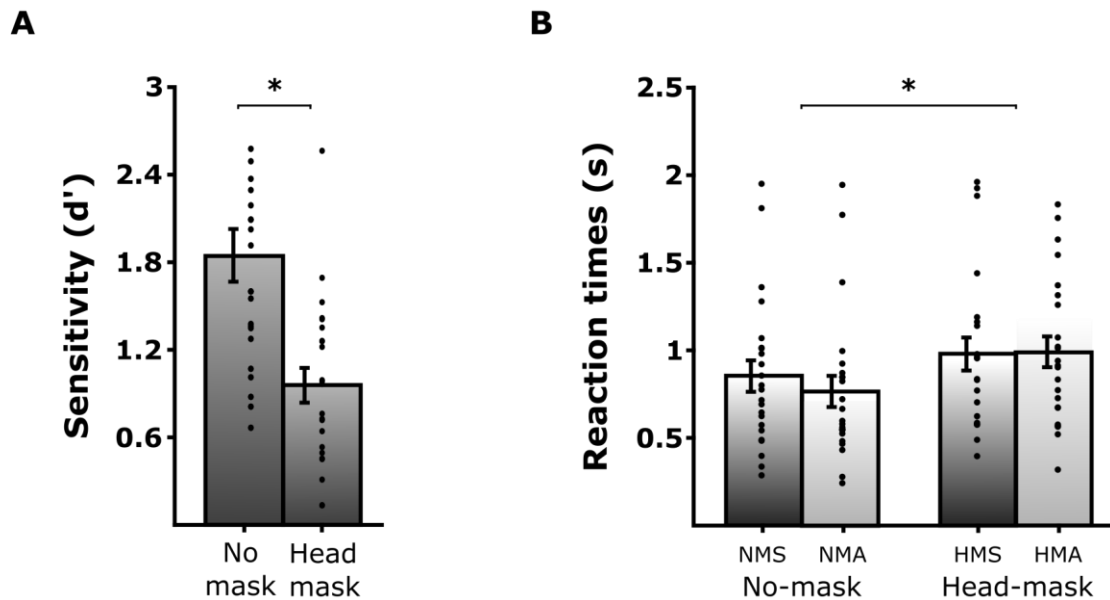
346     **multimodal speech integration.**

13

*Figure 2.* Behavioural performances in the synchrony detection task. (A) Average *d'* scores and (B) reaction times of hits across conditions (± standard error of the mean; grey dots represent individual averages; n = 23). Significant contrasts are marked by stars ($p < 0.05$).

D' scores are reported in Figure 2A. Two independent one-sample *t*-tests revealed that the mean *d'* in the no-mask (NMS and NMA) and head-mask contrasts (HMS and HMA) were significantly greater than $\mu$ = zero (no-mask: $t(1,22) = 10.25$; $p < 0.001$; *Cohen's d* = 3.04; head-mask: $t(1,22) = 8.07$; $p < 0.001$; *Cohen's d* = 2.38). A paired-samples *t*-test performed on *d'* for no-mask and head-mask contrasts confirmed that participants detected synchrony better in the no-mask contrast than in the head-mask contrast ($t(1,22) = 6.96$; $p = < 0.001$, two-tailed; *Cohen's d* = 1.46). Finally, two independent one-sample *t*-tests revealed that the mean *c* criterion in the no-mask contrast was not significantly different from zero (0.11 ± 0.40; $t(1,22) = 1.32$; $p = 0.1$; *Cohen's d* = 0.39), whereas it was significantly below zero in the head-mask contrast (-0.53 ± 0.28; $t(1,22) = -9.01$; $p < 0.001$; *Cohen's d* = 2.68). These results established that when the speaker's face was head-masked, participants were significantly biased to respond "synchrony" when the stimulus was asynchronous than in the no-mask contrast (i.e., liberal guessing). A two-way repeated-measure ANOVA revealed a significant main effect of mask on reaction times ($F(1, 22) = 16.50$, $p$

14

365    < 0.01; $\eta_p^2$ = 0.43). No significant effect of synchrony ($F$(1, 22) = 0.67, $p$ = 0.42; $\eta_p^2$ = 0.03) or

366    interaction between mask and synchrony was found ($F$(1, 22) = 2.32, $p$ = 0.14; $\eta_p^2$ = 0.1). These

367    results showed that accurate responses were faster when the face of the speaker was not masked

368    compared to head-masked, in line with greater difficulty in integrating video and audio signals

369    together when visual information was degraded (Figure 2B).

370    Together, the behavioural performance supports the hypothesis that participants integrated the

371    temporal structure of slow prosodic features in integrating visual and auditory information during

372    multimodal speech perception. Further, when visual information carried by the speaker's face

373    was degraded with a head-mask, the sensitivity to audiovisual synchrony decreased. This

374    suggests that successful multimodal integration of speech requires visual information of the head

375    and face.

376    **Delta oscillations in the left frontal-motor cortex reflect temporal integration of audio and**

377    **visual prosody and shape multimodal speech perception.**

378    We then addressed whether delta oscillations in the left motor cortex relate to multimodal

379    temporal integration, and whether responses depend on the amount of visual information

380    available. First, a cluster-based permutation tests revealed a significant increase in delta power

381    (2-3 Hz) in response to the audiovisual asynchrony when the speaker's face was visible (no-mask:

382    NMA-NMS) but not when it was masked (head-mask: HMA-HMS) (NMA-NMS: $p$ < 0.001, cluster

383    statistic = 117.23; HMA-HMS: zero positive cluster statistic; multiple comparisons are cluster-

384    corrected). No significant negative clusters were found in both contrasts. Importantly, the

385    topography of the significant delta cluster in the no-mask contrast showed a main fronto-central

386    response when video and audio signals were asynchronous, in line with the expected source

387    localization of delta in the motor region (Figure 3B; Puzzo et al., 2010; Stegemöller et al., 2017).

388    To assess the potential interaction of visual information and audiovisual synchrony perception in

389    this motor region of interest, we defined a set of electrodes as the region of interest (ROI)

390    representative of the delta response topography: F1, Fz, F2, FFC3h, FFC1h, FFC2h, FFC4h, FC3,

391    FC1, FCz, FC2, FC4, FCC3h, FCC1h, FCC2h, fCC4h, C1, Cz and C2 (Figure 1C). The mean delta power

392    across the electrodes of the ROI was computed in the four conditions separately, and confirmed

15

393    an increase of induced delta activity compared to the pre-stimulus baseline (NMS: 0.64 ± 0.17;

394    NMA: 0.74 ± 0.15; HMS: 0.70 ± 0.16 and HMA: 0.68 ± 0.20; see Figure 3A and 3C). A two-way

395    repeated-measure ANOVA revealed a significant interaction between the factors mask and

396    synchrony for delta power ($F(1, 22) = 5.78$, $p = 0.03$; $\eta_p^2 = 0.21$). Bonferroni-corrected pairwise

397    comparisons showed that in the no-mask contrast, delta power was significantly greater in the

398    asynchronous (NMA) than synchronous (NMS) condition ($p < 0.001$). In contrast, asynchrony did

399    not affect delta power responses in the head mask contrast ($p = 0.52$). Second, when video and

400    audio signals were presented in asynchrony, delta power was significantly greater when the

401    speaker's face was not masked as compared to head-masked (NMA vs. HMA; $p = 0.038$). When

402    the video and audio signals were synchronous, the presence of the mask on the speaker's face

403    did not significantly modulate the delta power (NMS vs. HMS; $p = 0.11$).
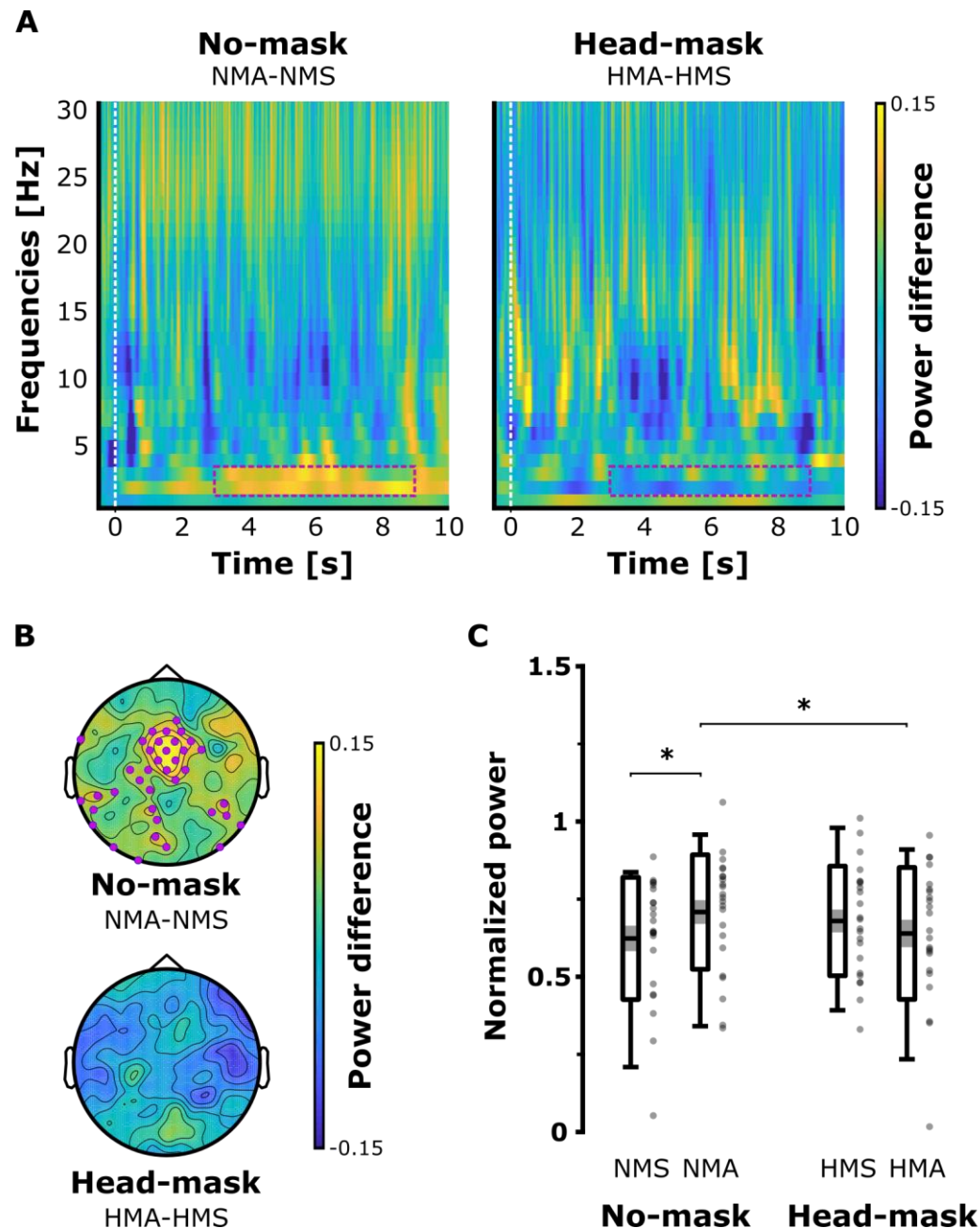
404

405

406
407

*Figure 3.* Delta responses to audiovisual asynchrony at the scalp level. (A) Time-frequency spectra of the mean power differences in the motor ROI between asynchronous and synchronous conditions in the no-mask (NMA-NMS; left) and head-mask (HMA-HMS; right) contrasts. The white dashed lines correspond to the onset of the video and the window of interest is marked by the pink dashed rectangles. (B) Topographical distribution of the difference of 2-3 Hz delta power in the time-window of interest, in the no-mask (NMA-NMS; top) and head-mask (HMA-HMS;

414    bottom) contrasts. The pink dots display electrodes with significant *t*-values (alpha threshold =

415    0.05). (C) Delta power across the electrodes of interest in the four conditions (2-3 Hz band).

416    Significant contrasts are marked by stars ($p < 0.05$).

417

418    Second, to separate the influence of audiovisual speech integration from sensory processing,

419    delta responses were also examined in a control region of non-interest (RONI; Figure 1C). The

420    region of non-interest was located in the occipital cortex where we did not expect higher

421    audiovisual speech analysis to take place as visual information was identical between

422    synchronous and asynchronous conditions within mask contrasts (RONI electrodes: PPO1h,

423    PPO2h, PO3, POz, PO4, POO1, POO2, POO9h, O1, Oz, O2, POO10h, Ol1h, Ol2h, O9 and O10). We

424    compared the effect of audiovisual asynchrony between the identified motor region (ROI) and

425    the visual sensory area (RONI) to confirm that delta response modulations did not reflect signal

426    processing only (Figure 4A). The mean differences of 2-3Hz delta power (NMA-NMS and HMA-

427    HMS) were computed in the regions of interest and non-interest in the same time-window as

428    previously (Figure 4B; ROI: NMA-NMS = 0.1 ± 0.09; HMA-HMS = -0.03± 0.19; RONI: NMA-NMS =

429    0.05 ± 0.10; HMA-HMS = 0.01 ± 0.24). A two-way repeated-measures ANOVA with the mean

430    factors region (ROI or RONI) and mask (no-mask or head-mask) was performed to assess whether

431    the responses of delta oscillations to asynchrony reflected multimodal speech analysis or purely

432    signal processing taking place in sensory areas (i.e. visual occipital areas). Results revealed a

433    significant interaction between region and mask ($F(1, 22) = 5.75$, $p = 0.025$; $\eta_p^2 = 0.21$). First,

434    Bonferroni-corrected pairwise comparisons showed that in the no-mask contrast the delta power

435    difference NMA-NMS (but not HMA-HMS) was significantly greater in the region of interest than

436    in the region of non-interest (respectively p = 0.025 and p = 0.572). Only in the region of interest

437    the difference of power NMA-NMS was significantly greater than HMA-HMS (respectively p =

438    0.019 and p = 0.113). No main effect of mask ($F(1, 22) = 0.25$, $p = 0.622$; $\eta_p^2 = 0.21$) or region ($F(1,$

439    $22) = 2.18$, $p = 0.154$; $\eta_p^2 = 0.09$) was found.
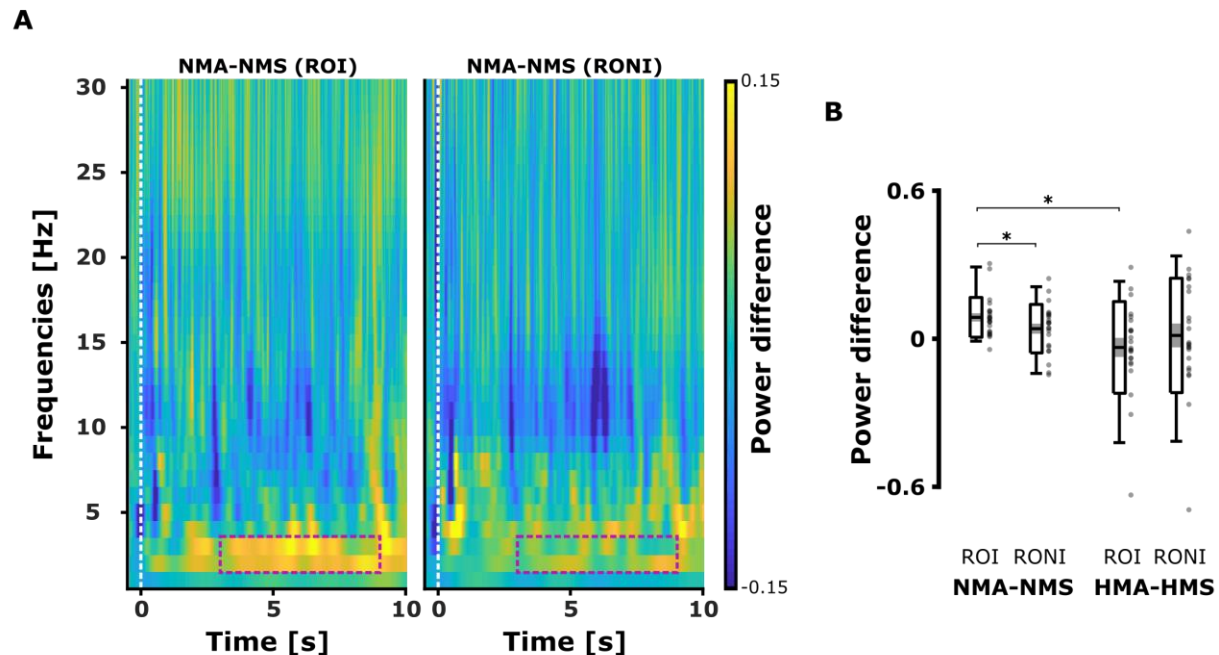
440

441

442

18

**A**



*Figure 4*. Comparisons between the motor region of interest (ROI) and the visual region of non-interest (RONI). (A) TFRs of the difference of spectrum in the no-mask contrast (NMA-NMS) in the ROI and RONI. (B) The mean differences of 2-3Hz delta power (NMA-NMS and HMA-HMS) were computed in the regions of interest and non-interest. Significant contrasts are marked by stars ($p < 0.05$).

Third, the mean power in the 4 - 8 Hz band was computed in the four conditions separately from the ROI electrodes and confirmed an increase of theta activity compared to the pre-stimulus onset baseline (NMS: 0.86 ± 0.25; NMA: 0.85 ± 0.18; HMS: 0.83 ± 0.16 and HMA: 0.81 ± 0.24). A two-way repeated-measure ANOVA revealed no significant main effect of mask ($F(1, 22) = 2.77$, $p = 0.11$; $\eta_p^2 = 0.11$), synchrony ($F(1, 22) = 0.27$, $p = 0.606$; $\eta_p^2 = 0.01$) or interaction between the factors mask and synchrony ($F(1, 22) = 0.05$, $p = 0.825$; $\eta_p^2 < 0.01$) on theta power in the region of interest. Further, the cluster-based permutation tests revealed no significant modulation of theta power by audiovisual asynchrony in any of the mask contrasts (NMA-NMS: no significant cluster; HMA-HMS: no significant cluster; multiple comparisons are cluster-corrected). These results confirmed that audiovisual asynchrony specifically modulated delta power over the expected fronto-central region. Further, delta responses were attenuated when listeners were

19

461    less able to integrate visual and auditory features, supporting the role of delta activity in the

462    temporal integration of multimodal speech.

463

464    Next, we analysed the source localisation of the delta power modulations observed when

465    video and audio signals were presented in asynchrony in both no-mask and head-mask contrasts.

466    Cluster-based permutation *t*-tests between synchronous and asynchronous conditions revealed

467    that asynchrony significantly increased delta oscillation responses when the head of the speaker

468    was visible (NMA-NMS: $p$ = 0.042; cluster statistic = 233.02) but not when it was head-masked

469    (HMA-HMS: $p$ = 0.27; cluster statistic = 38.27). The projections of the significant *t*-values on the

470    brain's surface showed an increase of delta power originating mainly in the left precentral region

471    and the left inferior frontal gyrus (Figure 5A). The source results support the topographies of the

472    delta power modulations observed at the scalp level, which revealed fronto-central differences

473    in the no-mask contrast only (Figure 3B). Further, we tested whether delta power responses from

474    the left motor areas correlated with the synchrony detection performance in the no-mask and

475    head-mask contrasts (Figure 5B). Pearson correlations revealed a positive correlation between

476    the hit rates and delta power differences in the no-mask contrast (NMA-NMS: r = 0.45; $p$ = 0.031,

477    two-tailed), but not in the head-mask contrast (HMA-HMS: r = 0.03; $p$ = 0.9, two-tailed). These

478    results confirmed that when participants were able to perceive synchrony between video and

479    audio signals (no-mask contrast), the amplitude of delta power modulations positively correlated

480    with accuracy. In contrast, when participants were less able to discriminate temporal alignment

481    between visual and auditory information (head-mask contrast), left motor delta oscillations did

482    not significantly correlate with behavioural performance.
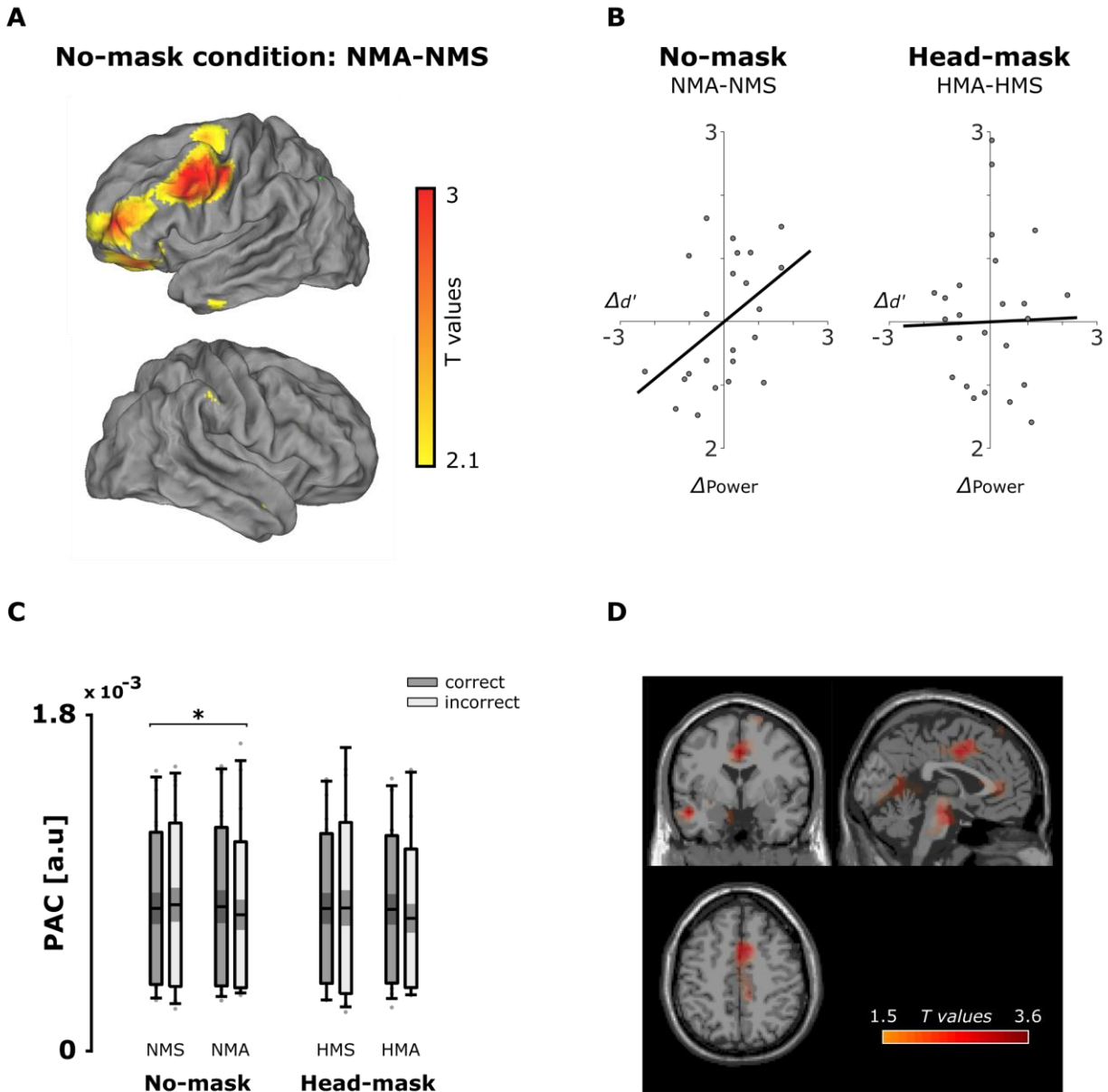
483

Figure 5. Delta oscillation responses to audiovisual asynchrony at the source level for no-mask and head-mask contrasts. (A) Contrast NMA$_{hit}$ - NMS$_{hit}$ projected onto the brain's surface (significance *t*-values; cluster-corrected at alpha threshold = 0.05). The maximum voxel MNI coordinates is located left precentrally [-50 19 40] but significant activations were also found in the left inferior frontal gyrus (pars triangularis; maximum voxel MNI coordinates [-30 31 0]). No significant difference was found when the head of the speaker was masked (HMA$_{hit}$ - HMS$_{hit}$ contrast; not represented). (B) Scatterplots of audiovisual synchrony detection performance and

492    delta power in the significant cluster region (left frontal-motor area). The difference of delta

493    power in the left motor cluster (x-axis) correlated with the difference of audiovisual synchrony

494    perception (y-axis) between synchronous and asynchronous conditions only when the face of the

495    speaker was visible and participants could integrate video and audio onsets (no-mask contrast).

496    (C) PAC analysis in the left frontal motor cluster. The figure represents the modulation of delta-

497    beta PAC in a significant cluster, depending on the mask and audiovisual synchrony. Significance

498    is indicated by an asterisk ($p < 0.05$, Bonferroni-corrected). Delta-beta PAC from the left frontal

499    motor area was greater in the no-mask as compared to the head-mask contrast in general, but

500    did not discriminate between hit and miss trials. Significant contrasts are marked by stars ($p <$

501    $0.05$). (D) Delta-beta PAC difference between no-mask ($NMA_{hit} + NMS_{hit}$) and head-mask ($HMA_{hit}$

502    $+ HMS_{hit}$) contrasts in the whole brain. Results revealed significant maximum differences located

503    in the superior motor area (MNI coordinates [0 11 50]) and in the left middle temporal lobe (MNI

504    coordinates [-50 -1 -20]).

505

506    **Delta-beta PAC reflects temporal integration in audiovisual speech perception, but is not**

507    **limited to the left motor region.**

508    Finally, we assessed whether delta-beta PAC modulations in the left frontal-motor area would

509    reflect sensitivity to audiovisual synchrony in speech. First, a three-way repeated-measure

510    ANOVA (main factors: mask, synchrony and correctness) revealed a main effect of mask on delta-

511    beta PAC with delta-beta phase-coupling being significantly greater in the no-mask than in the

512    head-mask contrast ($F(2,22) = 4.72$; $p = 0.041$; $\eta_p^2 = 0.18$; see Figure 5C). No further significant

513    main effect or interaction were found. These results suggest that left frontal-motor delta-beta

514    PAC relates to the temporal integration of audiovisual speech as it responded depending on

515    whether listeners were able to match visual and auditory prosodic features (no-mask contrast)

516    or not (head-mask contrast). Second, we investigated whether the delta-beta PAC difference

517    between no-mask and head-mask contrasts was restricted to the left motor areas. As accuracy

518    and synchrony did not affect delta-beta PAC in the cluster of interest, we selected only the hit

519    trials for the delta-beta PAC analysis at the whole brain level and put synchronous and

520    asynchronous trials together within contrasts (i.e. NMCs: $NMA_{hit} + NMS_{hit}$; HMCs: $HMA_{hit} +$

521    HMS$_{hit}$). The cluster-based permutation tests revealed one significant positive cluster peaking in

522    the superior motor area and in the left middle temporal lobe (although not exclusively; see Figure

523    5D), confirming that delta-beta PAC was significantly greater in the no-mask (NMCs) contrast as

524    compared to the head-mask (HMCs) contrast (NMCs - HMCs : $p$ = 0.043, cluster statistic = 216.69).

525    In summary, the EEG results mirrored the behavioural evidence as delta responses were distinctly

526    modulated by audiovisual synchrony only when participants could view the face and visible

527    articulators (no-mask contrast). Delta activity from the left motor region increased when visual

528    and auditory information were misaligned (NMA), reflecting greater difficulty to match visual and

529    auditory prosodic features as compared to the synchronous condition (NMS). Further, left

530    frontal-motor delta responses predicted synchrony detection performance only when

531    participants were able to properly integrate visual and auditory features (no-mask contrast), but

532    not when they guessed (head-mask contrast). Finally, the cross-frequency coupling analysis

533    showed that delta-beta PAC in the left frontal motor cluster of interest also increased when

534    listeners were able to match prosodic features between modalities (no-mask contrast) as

535    compared to when they guessed (head-mask contrast). These results suggest that delta-beta PAC

536    in expected motor areas (although not exclusive) are sensitive to temporal integration of

537    audiovisual speech information, and may predict whether listeners integrate visual and auditory

538    prosodic features in asynchrony detection.

**DISCUSSION**

540    The present study investigated the role of motor delta oscillations during the temporal

541    integration of multimodal prosodic features in speech perception. Behavioural results showed

542    that listeners processed both prosodic features in multimodal speech with sufficient visual

543    information. At the brain level, the perception of audiovisual asynchrony induced an increase in

544    delta activity in the expected left motor cortex (extending to the inferior frontal gyrus), which

545    correlated with the participants' sensitivity to audiovisual synchrony. In contrast, participants

546    were less able to discriminate audiovisual synchrony when the speaker's facial information was

547    masked, which was characterised by an absence of delta activity response in the EEG. Finally,

548    delta-beta PAC in the left frontal-motor areas decreased significantly when listeners could not

549 integrate efficiently visual and auditory prosodic features in speech perception. Altogether, our
550 results indicate that the delta time-scale provides a flexible framework to synchronise the
551 listener's brain activity with the temporal organization of external audiovisual speech. In this
552 framework, the oscillatory activity can gather and realign multiple temporal representations of
553 the visual and auditory speech features in the left motor cortex to improve dynamic signal
554 processing.

555 Synchrony detection performance confirmed our first hypothesis that listeners integrate
556 prosodic events in multimodal speech perception. This finding was expected as visual information
557 complements auditory information and often improves speech perception (Sumby & Pollack,
558 1954; van Wassenhove et al., 2005). Speaker's articulatory movements and gestures temporally
559 aligned with acoustic prosodic cues, providing listeners with a reliable delta temporal structure
560 of the speech signal (Biau et al., 2016; Esteve-Gibert & Guellaï, 2018; Wagner et al., 2014).
561 Therefore, participants likely use these salient prosodic events as landmarks present in two
562 different sensory streams to align and integrate them into a coherent multisensory speech
563 percept. These results suggest that the temporal structure focuses the listeners' attention within
564 brief time-windows containing common multimodal prosodic events to facilitate their
565 integration. This is in line with the theory of dynamic attending stating that non-random external
566 stimulation drives periodic attention allocation towards critical events (Large & Jones, 1999). In
567 contrast, when the speaker's face was masked, participants could not integrate the temporal
568 correspondence between visual and auditory prosodic anchors properly, and were less able to
569 perceive multimodal speech. Noteworthy, the differences of performance between the no-mask
570 and head-mask contrasts indicate that participants likely relied on complementary information
571 conveyed by the speaker's head, face, and fine articulatory gesture information to achieve the
572 integration of the visual prosodic signal (Cross, Butler, & Lalor, 2015).

573 Further, the EEG results revealed an increase in motor delta activity in response to audiovisual
574 asynchrony, confirming its role in temporal integration of multimodal prosodic features. Previous
575 literature associated motor delta oscillations with the perception of rhythmic auditory inputs
576 (Keitel et al., 2018; Morillon et al., 2019; Morillon & Schroeder, 2015). The present results extend

24

577    these findings to the temporal integration of non-isochronous events that act as punctual "snap

578    fasteners" integrating visual and auditory signals within relevant time-windows. As long as they

579    provide the brain with a dominant temporal structure aligning multiple sensory inputs, salient

580    prosodic features do not have to be perfectly regular to engage delta responses in the motor

581    cortex. The present EEG results corroborate this hypothesis in three ways: First, we did not

582    observe any different delta responses in auditory and visual cortices when audiovisual stimuli

583    were synchronous. This would have reflected low-level feature tracking during early sensory

584    processing (Cross, Butler & Lalor, 2015; Ghitza, 2017; Gross et al., 2013; Mai, Minett & Wang,

585    2016). Next, audiovisual asynchrony would likely decrease pure entrainment by making signal

586    tracking more difficult than when different channels of the same input are processed in

587    synchrony. Further, we found no theta activity in response to audiovisual asynchrony that would

588    have indicated an effect driven specifically by the prosodic features' rate (e.g., lip movements)

589    rather than temporal integration of sensory input. Second, the difference in delta power in the

590    left motor cortex correlated positively with performance between the synchronous and

591    asynchronous conditions in the no-mask contrast. Moreover, the fact that performance in the

592    synchronous and asynchronous conditions was similar when the face of the speaker was visible

593    suggests an increase of difficulty to integrate the two temporal representations of speech signal.

594    This extra cognitive load may be reflected by an increase of delta activity responses in the left

595    motor cortex. Third, participants did not perceive audiovisual synchrony when the speaker's

596    facial information was blurred, which was reflected by weaker responses in the left motor cortex

597    and no significant difference between synchronous and asynchronous conditions. Importantly,

598    the responses found in the left inferior frontal gyrus align well with previous research that

599    established a role in crossmodal information integration between gestures and speech (Park et

600    al., 2018; Willems, Ozyürek & Hagoort, 2009; Zhao et al., 2018). Here, participants perceived

601    information carried in the two modalities and likely integrated gestures' kinematics with auditory

602    envelope modulations to perform the synchrony detection task. Further investigations will need

603    to address whether the response modulations in the left IFG were specific to gesture-speech

604    temporal integration or could be reproduced using moving dots following gestures' dynamics. In

605    contrast, we found no difference of activation in additional regions associated with multimodal

25

606 speech integration such as the left posterior superior temporal sulcus (Marstaller & Burianová,
607 2014). However, it is possible that in the present context delta oscillations did not reflect
608 multisensory integration *per se* but temporal integration taking place in the left motor cortex and
609 IFG.

610 Finally, the cross-frequency coupling analysis revealed that delta-beta coupling in the left
611 frontal motor cortex increased when listeners perceived audiovisual (mis)alignment (no-mask
612 contrast). This finding indicates that delta-beta PAC contributes to temporal integration of
613 prosody as well. Potentially, delta-beta coupling may support the latter mechanisms taking place
614 after proper temporal integration of the visual and auditory prosodic features, e.g. auditory-
615 motor communication. Park et al. (2015) showed that the left frontal-motor areas modulated the
616 phase of delta oscillations in the left auditory cortex by means of top-down control in speech
617 perception. Reciprocally, delta-beta PAC in the auditory cortex respond to the modulations of
618 rhythmic regularity in auditory speech perception (Chang, Bosnyak and Trailor, 2019). Further,
619 Keitel et al. (2018) reported that delta-beta PAC in the left motor cortex predicted behavioural
620 performance in speech comprehension. Future research will need to unravel whether delta-beta
621 coupling provides a ubiquitous means of cross-regional communication to align temporally
622 different dynamic inputs in sensory cortices (Arnal, 2012; Fujioka, Ross & Trainor, 2015; Morillon
623 et al., 2019). For instance, Fontolan et al. (2014) reported that delta-beta coupling in the
624 associative auditory cortex modulated the phase of gamma activity related to phonological
625 processing in the primary auditory cortex in auditory sentence perception (Giraud & Poeppel,
626 2012). Alternatively, delta-beta PAC may drive the periodicity of attention to critical time-
627 windows containing relevant accentuated speech information, which fits with the dynamic
628 attention theory (Large & Jones, 1999).

629 To our knowledge, our results show for the first time how delta activity provides an interface
630 between external dynamic stimulation and inner brain oscillations to facilitate multimodal
631 speech perception. We propose that motor delta oscillations align together distinct
632 representations of non-verbal and auditory prosodic features encoded separately in their
633 respective sensory cortices. The slow time-scale of delta (1-3Hz) may also offer the brain some

634    flexibility to create a coherent multimodal percept despite the natural delay between visual and

635    auditory signal onsets in speech (Chandrasekaran et al., 2009). In social interactions where

636    conditions change quickly, such a delta framework would help listeners to align the related

637    speech streams in a bottleneck fashion to maintain a stable synchronization with the speaker's

638    flow (Kotz, Ravignani & Fitch, 2018). In contrast, when the temporal structure of events from two

639    contemporary streams cannot be integrated in critical delta time-windows, they are

640    discriminated against each other. When video and audio signal onsets were misaligned by 400ms,

641    the alignment of the visual and auditory neural representations in the delta-phase became likely

642    impossible, leading to the detection of asynchrony. Further investigations will need to address

643    whether this potential mechanism exists with other time-scales present in both speech signal and

644    endogenous oscillations. For instance, we cannot fully discard that the prosodic contour in our

645    stimuli still contained a syllable structure embedded in it (e.g. at onsets and stress peaks).

646    Further, lip movements and auditory envelope convey syllabic information occurring at a theta

647    rate (4-8 Hz) providing other robust temporal information in the speech signal during face-to-

648    face conversations (Chandrasekaran et al., 2009; Giraud & Poeppel, 2012). Therefore, delta and

649    theta activities may actually couple to strengthen speaker-listener synchronization in social

650    communicative interactions. Future research needs to investigate the potential role of a delta-

651    theta coupling in speech perception.

652    **CONCLUSION**

653    Our findings show that delta power and delta-beta phase-amplitude coupling in the left motor

654    cortex reflect the temporal integration of visual and auditory prosodic events, and shaped

655    multimodal integration in speech perception. We propose that the delta time-scale provides a

656    reliable framework allowing endogenous activity to align multiple prosodic features conveyed in

657    distinct sensory modalities in a common temporal organization during speech perception.

658    **ACKNOWLEDGMENT**

27

662     **RESOURCE SHARING**

663     Consent for sharing data at the level of the individual participant was received. Data for individual

664     participants and associated scripts will be made available after publication of the manuscript.

665     Further information or requests should be directed to the corresponding authors.

666     **REFERENCES**

667     Arnal, L. H. (2012). Predicting "When" Using the Motor System's Beta-Band Oscillations. *Frontiers*
668         *in Human Neuroscience*, *6*, 225

669     Arnal, L. H., Doelling, K. B., & Poeppel, D. (2015). Delta-Beta Coupled Oscillations Underlie
670         Temporal Prediction Accuracy. *Cerebral Cortex*, *25*(9), 3077–3085

671     Biau, E., & Soto-Faraco, S. (2013). Beat gestures modulate auditory integration in speech
672         perception. *Brain and Language*, *124*(2), 143–152

673     Biau, E., Fernandez, L. M., Holle, H., Avila, C., & Soto-Faraco, S. (2016). Hand gestures as visual
674         prosody: BOLD responses to audio-visual alignment are modulated by the communicative
675         nature of the stimuli. *NeuroImage,132*, 129–137

676     Boersma, P., and Weenink, D. (2015). Praat: Doing Phonetics by Computer. Version 5.4.17

677     Chandrasekaran, C., Trubanova, A., Stillittano, S., Caplier, A., & Ghazanfar, A. A. (2009). The
678         natural statistics of audiovisual speech. *PLoS Computational Biology*, *5*(7), e1000436.
679         https://doi.org/10.1371/journal.pcbi.1000436

680     Chang, A., Bosnyak, D. J., & Trainor, L. J. (2019). Rhythmicity facilitates pitch discrimination:
681         Differential roles of low and high frequency neural oscillations. *NeuroImage*, *198*, 31–43.
682         https://doi.org/10.1016/j.neuroimage.2019.05.007

683     Cherry, E. C. (1953). Some Experiments on the Recognition of Speech, with One and with Two
684         Ears. The Journal of the Acoustical Society of America. 25 (5): 975–79

685    Crosse M.J., Butler J.S., Lalor E.C. (2015). Congruent visual speech enhances cortical entrainment
686         to continuous auditory speech in noise-free conditions. The Journal of Neuroscience
687         35:14195–14204

688    Ding, N., Melloni, L., Zhang, H., Tian, X., & Poeppel, D. (2016). Cortical tracking of hierarchical
689         linguistic structures in connected speech. *Nature Neuroscience*, *19*(1), 158–164

690    Doelling, K. B., Arnal, L. H., Ghitza, O., & Poeppel, D. (2014). Acoustic landmarks drive delta-theta
691         oscillations to enable speech comprehension by facilitating perceptual parsing.
692         *NeuroImage*, *85 Pt 2*, 761–768

693    Esteve-Gibert N., & Guellaï B. (2018). Prosody in the Auditory and Visual Domains: A
694         Developmental Perspective. *Frontiers in Psychology*, 9:338. doi:10.3389/fpsyg.2018.00338

695    Fontolan, L., Morillon, B., Liegeois-Chauvel, C., and Giraud, A.-L. (2014). The contribution of
696         frequency-specific activity to hierarchical information processing in the human auditory
697         cortex. *Nat. Commun.* 5:4694. doi: 10.1038/ncomms5694

698    Fujioka, T., Ross, B., & Trainor, L. J. (2015). Beta-Band Oscillations Represent Auditory Beat and
699         Its Metrical Hierarchy in Perception and Imagery. *Journal of Neuroscience*, *35*(45), 15187–
700         15198

701    Ghitza, O. (2017). Acoustic-driven delta rhythms as prosodic markers. *Language, Cognition and*
702         *Neuroscience*, *32*(5), 545–561. https://doi.org/10.1080/23273798.2016.1232419

703    Giraud, A.-L., & Poeppel, D. (2012). Cortical oscillations and speech processing: emerging
704         computational principles and operations. *Nature Neuroscience*, *15*(4), 511–517

705    Griffiths, B. J., Parish, G., Roux, F., Michelmann, S., Plas, M. Van Der, Kolibius, D., … Hanslmayr, S.
706         (2019). Directional coupling of slow and fast hippocampal gamma with neocortical alpha /
707         beta oscillations in human episodic memory. *Proceedings of the National Academy of*
708         *Sciences*, 1–9. https://doi.org/10.1073/pnas.1914180116

709    Gross, J., Hoogenboom, N., Thut, G., Schyns, P., Panzeri, S., Belin, P., & Garrod, S. (2013). Speech
710         rhythms and multiplexed oscillatory sensory coding in the human brain. *PLoS Biology*,
711         *11*(12), e1001752

712    Gunter, T. C., & Douglas Weinbrenner, J. E. (2017). When to Take a Gesture Seriously: On How
713         We Use and Prioritize Communicative Cues. *Journal of Cognitive Neuroscience, 29*(8), 1355-
714         1367

715  Jessen, S., & Kotz, S. A. (2015). Affect differentially modulates brain activation in uni- and
716      multisensory body-voice perception. *Neuropsychologia*, *66*, 134–143

717  Keitel, A., Ince, R. A. A., Gross, J., & Kayser, C. (2017). Auditory cortical delta-entrainment
718      interacts with oscillatory power in multiple fronto-parietal networks. *NeuroImage*, *147*, 32–
719      42. https://doi.org/10.1016/j.neuroimage.2016.11.062

720  Keitel, A., Gross, J., & Kayser, C. (2018). Perceptually relevant speech tracking in auditory and
721      motor cortex reflects distinct linguistic features. *PLoS Biology*, *16*(3), e2004473

722  Kösem, A., & van Wassenhove, V. (2017). Distinct contributions of low- and high-frequency neural
723      oscillations to speech comprehension. *Language, Cognition and Neuroscience*, *32*(5), 536–
724      544

725  Kotz S.A., Ravignani A., Fitch W.T. The Evolution of Rhythm Processing. *Trends Cogn Sci*.
726      2018;22(10):896-910. doi:10.1016/j.tics.2018.08.002

727  Large, E.W., & Jones, M.R. (1999). The dynamics of attending: how people track time-varying
728      events. *Psychol. Rev.* 106, 119

729  Mai, G., Minett, J. W., & Wang, W. S. Y. (2016). Delta, theta, beta, and gamma brain oscillations
730      index levels of auditory sentence processing. *NeuroImage*, *133*, 516–528.
731      https://doi.org/10.1016/j.neuroimage.2016.02.064

732  Maris, E., & Oostenveld, R. (2007). Nonparametric statistical testing of EEG- and MEG-data.
733      *Journal of Neuroscience Methods*, *164*(1), 177–190

734  Marstaller L, Burianov_a H.(2014). The multisensory perception of co-speech gestures - a review
735      and meta-analysis of neuroimaging studies. *J Neurolinguistics, 30*, 69-77.

736  Meyer, L., Sun, Y., & Martin, A. E. (2019). Synchronous, but not entrained: exogenous and
737      endogenous cortical rhythms of speech and language processing. *Language, Cognition and
738      Neuroscience*, 1–11. https://doi.org/10.1080/23273798.2019.1693050

739  Morillon, B., & Schroeder, C. E. (2015). Neuronal oscillations as a mechanistic substrate of
740      auditory temporal prediction. *Annals of the New York Academy of Sciences*, *1337*(1), 26–31.
741      https://doi.org/10.1111/nyas.12629

742 Morillon, B., & Baillet, S. (2017). Motor origin of temporal predictions in auditory attention.
743      *Proceedings of the National Academy of Sciences of the United States of America*, *114*(42),
744      E8913–E8921

745 Morillon, B., Arnal, L. H., Schroeder, C. E., & Keitel, A. (2019). Prominence of delta oscillatory
746      rhythms in the motor cortex and their relevance for auditory and speech perception.
747      *Neuroscience and Biobehavioral Reviews*, Vol. 107, pp. 136–142.
748      https://doi.org/10.1016/j.neubiorev.2019.09.012

749 Munhall, K. G., Jones, J. A., Callan, D. E., Kuratate, T., & Vatikiotis-Bateson, E. (2004). Visual
750      prosody and speech intelligibility: head movement improves auditory speech perception.
751      *Psychological Science*, *15*(2), 133–137

752 Obermeier, C., Dolk, T., & Gunter, T. (2012). The benefit of gestures during communication:
753      Evidence from hearing and hearing-impaired individuals. *Cortex, 48*, 857-870

754 Obermeier, C., & Gunter, T. C. (2014). Multisensory Integration: The Case of a Time Window of
755      Gesture-Speech Integration. *Journal of Cognitive Neuroscience*, 1–16

756 Oostenveld R, Fries P, Maris E, Schoffelen JM (2011) FieldTrip: open source software for advanced
757      analysis of MEG, EEG, and invasive electrophysiological data. Comput Intell Neurosci
758      2011:156869

759 Park, H., Ince, R. A. A., Schyns, P. G., Thut, G., & Gross, J. (2015). Frontal Top-Down Signals
760      Increase Coupling of Auditory Low-Frequency Oscillations to Continuous Speech in Human
761      Listeners. *Current Biology*, *25*(12), 1649–1653

762 Park, H., Kayser, C., Thut, G., & Gross, J. (2016). Lip movements entrain the observers' low-
763      frequency brain oscillations to facilitate speech intelligibility. *ELife*, *5*

764 Park, H., Ince, R., Schyns, P. G., Thut, G., & Gross, J. (2018). Representational interactions during
765      audiovisual speech entrainment: Redundancy in left posterior superior temporal gyrus and
766      synergy in left motor cortex. *PLoS biology*, *16*(8), e2006558.
767      https://doi.org/10.1371/journal.pbio.2006558

768 Peelle, J. E., & Davis, M. H. (2012). Neural Oscillations Carry Speech Rhythm through to
769      Comprehension. *Frontiers in Psychology*, *3*, 320

770     Puzzo, I., Cooper, N. R., Vetter, P., & Russo, R. (2010). EEG activation differences in the pre-motor
771          cortex and supplementary motor area between normal individuals with high and low traits
772          of autism. *Brain Research*, *1342*, 104–110

773     Saleh, M., Reimer, J., Penn, R., Ojakangas, C. L., & Hatsopoulos, N. G. (2010). Fast and Slow
774          Oscillations in Human Primary Motor Cortex Predict Oncoming Behaviorally Relevant Cues.
775          *Neuron*, *65*(4), 461–471

776     Schultz, B. G., Biau, E., & Kotz, S. A. (2020). An open-source toolbox for measuring dynamic video
777          framerates and synchronizing video stimuli with neural and behavioral responses. *Journal*
778          *of Neuroscience Methods*, 108830

779     Stegemöller, E. L., Allen, D. P., Simuni, T., & MacKinnon, C. D. (2017). Altered premotor cortical
780          oscillations during repetitive movement in persons with Parkinson's disease. *Behavioural*
781          *Brain Research*, *317*, 141–146

782     Sumby, W. H., & Pollack, I. (1954). Visual Contribution to Speech Intelligibility in Noise. *Journal of*
783          *the Acoustical Society of America*, *26*(2), 212–215

784     Tadel F, Baillet S, Mosher JC, Pantazis D, Leahy RM (2011) Brainstorm: a user-friendly application
785          for MEG/EEG analysis. Comput Intell Neurosci 2011:879716

786     Thut, G., Veniero, D., Romei, V., Miniussi, C., Schyns, P., & Gross, J. (2011). Rhythmic TMS causes
787          local entrainment of natural oscillatory signatures. *Current Biology*, *21*(14), 1176–1185.
788          https://doi.org/10.1016/j.cub.2011.05.049

789     Tort, A. B. L., Komorowski, R., Eichenbaum, H., & Kopell, N. (2010). Measuring phase-amplitude
790          coupling between neuronal oscillations of different frequencies. *Journal of Neurophysiology*,
791          *104*(2), 1195–1210. https://doi.org/10.1152/jn.00106.2010

792     van Veen, B. D., van Drongelen, W., Yuchtman, M., & Suzuki, A. (1997). Localization of brain
793          electrical activity via linearly constrained minimum variance spatial filtering. *IEEE*
794          *Transactions on Bio-Medical Engineering*, *44*(9), 867–880

795     van Wassenhove, V., Grant, K. W., & Poeppel, D. (2005). Visual speech speeds up the neural
796          processing of auditory speech. *Proceedings of the National Academy of Sciences of the*
797          *United States of America*, *102*(4), 1181–1186

798     Wagner, P., Malisz, Z., & Kopp, S. (2014). Gesture and speech in interaction: An overview. *Speech*
799          *Communication*, *57*, 209–232

800 Wang, D., Clouter, A., Chen, Q., Shapiro, K. L., & Hanslmayr, S. (2018). Single-trial phase
801     entrainment of theta oscillations in sensory regions predicts human associative memory
802     performance. *Journal of Neuroscience*, *38*(28), 6299–6309

803 Willems, R. M., Ozyürek, A., & Hagoort, P. (2009). Differential roles for left inferior frontal and
804     superior temporal cortex in multimodal integration of action and
805     language. *NeuroImage*, *47*(4), 1992–2004.
806     https://doi.org/10.1016/j.neuroimage.2009.05.066

807 Zhao, W., Riggs, K., Schindler, I., & Holle, H. (2018). Transcranial Magnetic Stimulation over Left
808     Inferior Frontal and Posterior Temporal Cortex Disrupts Gesture-Speech Integration. *The*
809     *Journal of Neuroscience*, *38*(8), 1891–1900. https://doi.org/10.1523/JNEUROSCI.1748-
810     17.2017

811 Zoefel, B., Archer-Boyd, A., & Davis, M. H. (2018). Phase Entrainment of Brain Oscillations Causally
812     Modulates Neural Responses to Intelligible Speech. *Current Biology*, *28*(3), 401-408.e5.
813     https://doi.org/10.1016/j.cub.2017.11.071