# JNeurosci
## THE JOURNAL OF NEUROSCIENCE

*Research Articles: Behavioral/Cognitive*

# Left motor delta oscillations reflect asynchrony detection in multisensory speech perception

1    **Left motor delta oscillations reflect asynchrony detection in multisensory speech perception**

2

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

4

5    [1]Department of Psychology, University of Liverpool, L69 7ZA, Liverpool, United Kingdom.

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

9    Sciences, 04103, Leipzig, Germany.

10

11   Corresponding authors:

12   Dr. Emmanuel Biau

13   Address: Department of Psychology, University of Liverpool, L69 7ZA, Liverpool, United

14   Kingdom.

15   Email : e.biau@liverpool.ac.uk

16

17   Prof. Dr. Sonja A. Kotz

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

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

20   Email: sonja.kotz@maastrichtuniversity.nl

21

22

23

24

25

26

27

28

29

30

## ABSTRACT

32    During multisensory speech perception, slow delta oscillations (~1 - 3 Hz) in the listener's

33    brain synchronize with the speech signal, likely engaging in speech signal decomposition.

34    Notable fluctuations in the speech amplitude envelope, resounding speaker prosody,

35    temporally align with articulatory and body gestures and both provide complementary

36    sensations that temporally structure speech. Further, delta oscillations in the left motor cortex

37    seem to align with speech and musical beats, suggesting their possible role in the temporal

38    structuring of (quasi)-rhythmic stimulation. We extended the role of delta oscillations to audio-

39    visual asynchrony detection as a test case of the temporal analysis of multisensory prosody

40    fluctuations in speech. We recorded EEG responses in an audio-visual asynchrony detection

41    task while participants watched videos of a speaker. We filtered the speech signal to remove

42    verbal content and examined how visual and auditory prosodic features temporally (mis-)align.

43    Results confirm (i) that participants accurately detected audio-visual asynchrony, and (ii)

44    increased delta power in the left motor cortex in response to audio-visual asynchrony. The

45    difference of delta power between asynchronous and synchronous conditions predicted

46    behavioural performance, and (iii) decreased delta-beta coupling in the left motor cortex when

47    listeners could not accurately map visual and auditory prosodies. Finally, both behavioural and

48    neurophysiological evidence was altered when a speaker's face was degraded by a visual mask.

49    Together, these findings suggest that motor delta oscillations support asynchrony detection of

50    multisensory prosodic fluctuation in speech.

## KEYWORDS

52    Audio-visual asynchrony; multisensory speech; delta oscillations; prosody; motor cortex.

## SIGNIFICANCE STATEMENT

54    Speech perception is facilitated by regular prosodic fluctuations that temporally structure

55    the auditory signal. Auditory speech processing involves the left motor cortex and associated

56    delta oscillations. However, visual prosody (i.e., a speaker's body movements) complements

57  auditory prosody, and it is unclear how the brain temporally analyses different prosodic

58  features in multisensory speech perception. We combined an audio-visual asynchrony

59  detection task with electroencephalographic recordings to investigate how delta oscillations

60  support the temporal analysis of multisensory speech. Results confirmed that asynchrony

61  detection of visual and auditory prosodies leads to increased delta power in left motor cortex

62  and correlates with performance. We conclude that delta oscillations are invoked in an effort to

63  resolve denoted temporal asynchrony in multisensory speech perception.

64  **INTRODUCTION**

65  Speaker prosody displays perceptible fluctuations in the speech amplitude envelope,

66  allowing a listener to segment and parse incoming speech (Ghitza, 2017). While not

67  isochronous, prosody imposes a temporal structure with regular alterations of strong and weak

68  accentuated cues occurring at ~1 - 3 Hz delta rate (Ding et al., 2016; Doelling et al., 2014;

69  Ghitza, 2017; Pell & Davis, 2012). Populations of neurons in visual and auditory cortices

70  synchronize their firing responses with the onsets of predictable events, structuring sensory

71  signals at a delta rate. Such "neural entrainment" reflects the early stages in sensory processing

72  by which neural oscillations might track temporally relevant signal features. The neural

73  representation of such sensory features must temporally align, and there is evidence that delta

74  oscillations play a role in unisensory as well as multisensory integration (Giraud & Poeppel,

75  2012; Keitel et al., 2017; Kösem & van Wassenhove, 2017; Meyer, Sun & Martin, 2019). For

76  example, using a temporal order judgement task, Kösem, Gramfort and van Wassenhove (2014)

77  showed that the phase shifts of entrained delta oscillations in the auditory cortex linearly

78  mapped participants' perception of audio-visual simultaneity. Other studies described an

79  interaction of delta oscillations in visual and auditory cortices for audio-visual speech (Crosse et

80  al., 2016; Mercier et al., 2015). Crosse, Liberto, and Lalor (2016) reported that speech envelope

81  tracking in the auditory cortex improved through visual information, particularly in the delta

82  range. Beyond segmentation, prosody presents in visual and auditory information and facilitate

83  synchronization of multimodal information in social interaction (Esteve-Gibert & Guellaï, 2018;

84  Kotz, Ravignani & Fitch, 2018). The term "visual prosody" encompasses communicative

3

85  gestures (i.e., hand, head, face, and body movements) whose prominent phase temporally

86  coincides with acoustic prosodic features such as intonational phrases, pitch accents, and

87  boundary tones (Biau et al., 2016; Chandrasekaran et al., 2009; Munhall et al., 2004; Wagner et

88  al., 2014). For example, listeners rely on the successful temporal analysis of gestures and

89  sounds in speech perception (Cherry, 1953; Obermeier, Dolk & Gunter, 2012; Sumby & Pollack,

90  1954). Together this raises the following questions: How does the brain temporally align

91  multiple dynamic prosodies in multisensory speech perception?

92  The present study investigated whether delta oscillations respond to manipulation of

93  temporal alignment in multisensory speech (i.e., dynamic non-verbal visual and auditory

94  prosodies). We refer to temporal alignment as the mechanism by which the brain attempts to

95  integrate the quasi-rhythmic structure of visual and auditory prosodies in multisensory speech.

96  Delta activity in the motor cortex has been associated with the temporal analysis of rhythmic

97  stimuli as its phase aligns with the onsets of predictable events (Morillon et al., 2019; Morillon

98  & Schroeder, 2015; Saleh et al., 2010). In speech, Keitel and colleagues (2018) showed that left

99  motor delta activity tracked temporally predictable slow phrasal features in auditory sentences

100  and predicted successful speech comprehension. This suggests that this region responds to

101  perceptually relevant regularities in the signal to improve comprehension. Keitel et al. (2018)

102  also found delta-beta cross-frequency coupling in the left motor region, in line with previous

103  research, showing that motor beta oscillations respond to the temporal alignment of rhythmic

104  auditory tones or visual cues (Fujioka, Ross & Trainor, 2015; Saleh et al., 2010). These findings

105  led to the hypothesis that delta oscillations are involved in the temporal analysis of speech by

106  mediating top-down control through cross-frequency coupling with beta activity (Arnal, 2012;

107  Arnal, Doelling & Poeppel, 2015; Morillon and Baillet, 2017). In other words, delta activity could

108  reflect how the brain gathers and temporally analyses different sensory inputs in left motor

109  cortex and generates predictions to improve (multisensory) signal processing. Finally, the left

110  motor cortex, including the left inferior frontal gyrus, is involved in gestures and speech

111  integration (Biau et al., 2016; Park et al., 2016; Zhao et al., 2018).

4

112    We propose that visual and auditory prosodic features encoded in the visual and auditory

113    sensory cortices provide two representations of the speech signal, and their (un-)successful

114    temporal alignment may recruit the left motor cortex during speech perception. To test this

115    hypothesis, we manipulated the temporal structure of filtered multisensory speech, including

116    whole body or masked head movements. Participants performed an audio-visual synchrony

117    detection task and watched small video clips of a single speaker engaged in a conversation. We

118    also recorded their electroencephalogram (EEG). Firstly, we tested behaviourally how

119    successfully listeners temporally align visual and auditory prosodic features in multisensory

120    speech. We then analysed modulations of delta oscillations in response to audio-visual

121    asynchrony to find out whether and to which degree they index (un-)successful temporal

122    alignment in multisensory speech perception. Thirdly, we tested whether delta-beta coupling in

123    the left motor cortex predicts multisensory (a-)synchrony detection in speech perception.

124    **METHODS**

125    **Participants**

126    We recruited twenty-six native Dutch speakers (mean age = 22.24, SD = 4.24; 15 females) at

127    Maastricht University, who received €10 for participating in the experiment after giving

128    informed consent. All participants were right-handed and had normal or corrected-to-normal

129    vision and hearing. The protocol of the study was approved by the Research Ethical Committee

130    of Maastricht University. Data from three participants were removed from the final analysis due

131    to technical problems.

132    **Stimuli**

133    Short videos were extracted from a longer video recording used in a previous study (Gunter &

134    Weinbrenner, 2017). The videos depicted a female actor and an experimenter (both German

135    native speakers) engaged in a question-answer conversation. The actor sat on a chair, moved

136    freely, and was visible from her knees up to the top of her head. Relevant segments containing

137    the actor' answers separate from the experimenter were selected to create the current

138    stimulus set (N = 54). Each of the 54 segments was 10 seconds long (600 frames at 60 frames

5

139  per second; FPS). The audio track was extracted to be low pass filtered with Hann band

140  windowing procedure (from 0 Hz to 400 Hz; 20 Hz smoothing) using Praat (Boersma & Weenink,

141  2015). In doing so, we altered speech intelligibility removing verbal content while keeping the

142  prosodic contour of the signal. Peak frequencies were extracted from the audio and video files

143  through Fourier transformations that calculated the frequency at which the peak amplitude

144  occurred within a range of 0.5Hz to 8Hz. For videos, the average magnitude of grayscale pixel

145  changes between consecutive frames was used to determine the frequency of movement and

146  gesture (see Table 1).

147                                    **[Insert Table 1]**

148  *Table 1*. Summary statistics of peak frequencies obtained for video and audio signals using a

149  Fourier transformation. Video signals: Full (Head + Body), Head only (either original head

150  information or head-masked) and Body only (lower body part without head information). Audio

151  signals: Audio only. The Audio and Body only measures (*) are consistent across the no-mask

152  and head-mask conditions. Frequencies are shown for the no-mask and head-mask videos for

153  the masked area (Head only) and all pixels (Full video).

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

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

156  audio-visual information or a temporal shift of the audio signal relative to the video onset

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

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

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

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

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

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

163  temporal alignment between visual and auditory onsets was intact. To create an asynchronous

164  condition, we inserted a delay between the visual and auditory onsets by shifting the sound

165  onset by +400 ms relative to the video onset (i.e., 24 frames). This manipulation maintained the

166  natural order of visual information preceding auditory information in the synchronous

167    condition in ecologically valid contexts but with a longer duration. In the current experiment,

168    audio onsets did not precede corresponding video onsets. This lag duration was based on a

169    time-window of multisensory integration established in previous studies (Biau et al., 2016; Biau

170    & Soto-Faraco, 2013; Jessen & Kotz, 2015; Obermeier & Gunter, 2014). A 400 ms lag ensured

171    that the delay between video and audio onsets was long enough for participants to detect

172    audio-visual asynchrony at a success rate of approximatively 80% ± 5%. This delay was

173    established in a pilot experiment with different participants (n = 16). Results confirmed that

174    participants detected both synchrony and asynchrony between visual and auditory information

175    in the audio-visual stimuli similarly (correct response rates in the synchronous condition: 0.79 ±

176    0.11 and asynchronous condition: 0.78 ± 0.12; t(1,15) = -0.26; *p* = 0.797; two-tailed; *Cohen's d* =

177    0.07). This was done to ensure we retained enough correct response trials in both conditions

178    for further EEG analyses. Further, a central white fixation cross was displayed in each video to

179    allow participants to focus their gaze on a central cue while attending audio-visual stimuli.

180    Altogether, this created four conditions: Head-Mask Synchronous (HMS), Head-Mask

181    Asynchronous (HMA), No-Mask Synchronous (NMS), and No-Mask Asynchronous (NMA) (see

182    Figure 1A). 18 additional video clips, in which the central white fixation-cross turned red, were

183    used as fillers, counterbalanced across conditions (colour change onset jittered between 5 and

184    9 seconds after the video onset; ~ 8 % of total stimuli, not included in the final data analysis).

185    We used the fillers in a memory test to focus the participants' attention on the videos during

186    the experiment. Finally, audio files were recombined with corresponding video files in each

187    condition. Videos were edited using Adobe Premiere Pro CS3 and exported using the following

188    parameters: Pixel resolution 1920 × 1080, 60 FPS compressor Indeo video 5.10, AVI format,

189    audio sample rate 48 kHz, 16 bits Mono.

190    **Apparatus**

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

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

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

194    1080 pixels connected to the stimulus presentation computer (Intel i7-6700 CPU @ 3.40 GHz,

195   32 GB, running 64-bit Windows 7, NVIDIA GeForce FTX 1080 GTX GPU). Stimuli were presented

196   using a custom MATLAB script (MATLAB and Statistics Toolbox Release 2015b, The MathWorks,

197   Inc., Natick, Massachusetts, United States) that called VideoLAN Client (VLC; VideoLAN Client,

198   2017; http://www.videolan.org/) to play the videos. EEG data were collected using BrainVision

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

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

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

202   **Procedure**

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

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

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

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

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

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

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

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

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

212   red cross in a video clip and reported it at the end of the experiment. This secondary task

213   ensured that the participants carefully attended both visual and auditory information during

214   the experiment. Further, we chose a relatively easy task, ensuring that performance in the

215   audio-visual synchrony detection task was not affected. Filler trials were not included in

216   behavioural and EEG analyses but the total number of reported red crosses served to check

217   that attention was maintained throughout the experiment. Before the experiment, participants

218   received five practice trials where they were presented with one example of each condition to

219   ensure they understood the instructions. At the end of the experiment, participants were asked

220   if they could identify the speaker's language and to report it.

221                                  **[Insert Figure 1]**

8

222   *Figure 1.* Experimental procedure of the audio-visual asynchrony detection task. (A) The four

223   experimental conditions. For each item, the audio signal was the same across all four versions.

224   Visual information was manipulated by the presence or absence of a mask (no-mask or head-

225   mask). Video and sound were either temporally aligned in the synchronous conditions (NMS,

226   HMS), or temporally misaligned by 400ms in the asynchronous conditions (NMA, HMA). (B)

227   Example of one trial timeline. (C) Distribution of the electrodes covering the motor region of

228   interest (ROI; blue circles) and the control region of non-interest in the visual area (RONI; red

229   circles).

230   **EEG recording and pre-processing**

231   Electrophysiological data were recorded at 1000 Hz with 128 active electrodes (ActiCap, Brain

232   Vision Recorder, Brain Products) according to the 10-20 international standard, and impedances

233   were kept below 10 kΩ. The ground electrode was located at AFz, and the reference electrode

234   was placed at the right mastoid (TP10).

235   Offline EEG pre-processing: EEG data were pre-processed offline using Fieldtrip (Oostenveld et

236   al., 2011) and SPM8 toolboxes (Wellcome Trust Centre for Neuroimaging). Continuous EEG

237   signals were bandpass filtered (standard non-causal two-pass Butterworth filters) between 0.1

238   Hz and 100 Hz and bandstop filtered (48-52 Hz and 98-102 Hz) to remove line noise at 50 and

239   100 Hz. Data were epoched from 1000 ms before stimulus onset to 11000 ms after stimulus

240   onset. Trials and channels with artefacts were excluded by visual inspection before applying an

241   independent component analysis (ICA) to remove components related to ocular artefacts.

242   Excluded channels were then interpolated using the method of triangulation of nearest. After

243   re-referencing the data to an average reference, the remaining trials with artefacts were

244   manually rejected by a final visual inspection (on average, 13.57 ± 8.32 trials across conditions

245   per participant).

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

247   Time-frequency analysis was applied to each electrode using a Morlet wavelet (width: 5 cycles,

248   from 1 to 40 Hz with 1 Hz step and 20 ms time steps) and frequency analyses were performed

9

249    for each trial prior to averaging across trials in the four conditions. The power was normalised

250    relative to a pre-stimulus baseline (-700 to -200 ms with respect to stimulus onset) to

251    determine increases or decreases of power dependent on the conditions. The peak frequency

252    analysis applied to the video and audio signals of the audio-visual clips revealed that prosodic

253    features conveyed activity mainly between 2 to 3 Hz (see Table 1), which determined our

254    frequency band of interest. In the present study, oscillatory delta activity was assessed by

255    means of power information, that is by taking the average power across the 2-to-3 Hz

256    frequencies (power and peak frequency values) and investigating its modulations as a function

257    of the audio-visual speech analysis. Further, as the spontaneous speech signal is not

258    isochronous, phase information likely would be too noisy to extract meaningful information.

259    This is the reason why we did not investigate phase modulations here. As entrainment

260    necessitates several cycles from recurrent stimulations to build up (Doelling et al. 2014; Thut et

261    al., 2011; Zoefel et al. 2018) and the slower frequency in our band of interest was 2Hz

262    (corresponding to a period of 500 ms), we defined a time window of interest from + 3 to + 9

263    seconds after stimulus onset. This time window ensured that neural activity sufficiently

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

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

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

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

268    analyses.

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

270    Source localisation: We used the Montreal Neurological Institute (MNI) MRI template and a

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

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

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

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

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

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

10

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

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

279   series data were reconstructed in 2020 virtual electrodes for each participant. Time-frequency

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

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

282   anatomical labelling atlas (AAL).

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

284   index (MI; Tort et al., 2010) analysis in the time-window of interest to quantify delta-beta PAC

285   in the significant cluster revealed by source localisation in the contrast NMA- NMS (i.e.

286   difference of delta power in the NMA condition minus NMS condition). Firstly, the power

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

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

289   participant (Griffith et al., 2019). Fractal activity was attenuated by subtracting the linear fit of

290   $1/f$ characteristic from the data to isolate oscillatory components before extracting the power

291   peaks. For each epoch, the spectral power was first calculated by applying a constant time-

292   frequency decomposition method with 5-cycles wavelet across all frequencies (from 1 to 30

293   Hz). This ensured that a single slope was computed and subsequently subtracted from the

294   signal. This step generated two vectors: one vector contained the values of each wavelet

295   frequency A, while the other vector contained the power spectrum for each electrode-sample

296   pair B. Both vectors were then put into log-space to provide a linear line in order to get the

297   slope and intercept of the $1/f$ curve. The linear equation A$x$ = B was resolved using least-

298   squares regression, where $x$ is an unknown constant describing the curvature of the

299   $1/f$ characteristic. The 1/$f$ fit A$x$ was then subtracted from the log-transformed power spectrum

300   B. Peaks of 1/$f$–corrected absolute power were then identified in the delta (1-3 Hz) and beta

301   (20-30 Hz) bands of interest for each trial. The most prominent power spectrum peaks in the

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

303   Across participants, the mean delta peak was at 2.1 Hz and the mean beta peak was at 24.16

304   Hz. To obtain an equal number of correct and incorrect trials across conditions, the same

305   number of trials between all conditions was determined by taking 80% of the smallest number

11

306    of available trials across all the conditions ($NMS_{correct}$, $NMS_{incorrect}$, $NMA_{correct}$, $NMA_{incorrect}$,

307    $HMS_{correct}$, $HMS_{incorrect}$, $HMA_{correct}$ and $HMA_{incorrect}$; average minimum number of trials: 6.61 ±

308    3.37). The 80% subsampling was done to ensure that some participants were not

309    overrepresented in the resampling procedure due to using 100% of their available data, as well

310    as to vary the set of trials in the condition determining the minimum number of trials across

311    iterations (Keitel et al., 2018). Subsampled trials were concatenated, and the operation was

312    repeated for 50 iterations in each condition to provide enough random trials to compute the

313    PAC (i.e., 50 trials per condition per participant). The grids of interest were identified during

314    source localisation (see Results section) and correspond to the grids at which the difference of

315    2-3 Hz delta power between the condition NMA minus NMS was significant (i.e., contrast NMA-

316    NMS; number of significant grids = 92). Second, the time-series of each grid source of the left

317    motor cluster were duplicated and filtered separately: the first time-series was filtered around

318    the theta peak (± 0.5 Hz) and the second time-series was filtered around the beta peak (± 5 Hz).

319    Third, the Hilbert transform was applied to the delta and beta filtered time-series to extract the

320    phase of the former and the power of the latter. Fourth, beta power was binned into 12

321    equidistant bins of 30° according to the delta phase. The binning was computed for each trial

322    and grid source separately. The MI was computed by comparing the observed distribution to a

323    uniform distribution for each trial and grid. The PAC was then averaged across the left motor

324    grids and 50 iterations in each condition. Finally, we investigated whether the delta-beta

325    coupling was specifically localised in the region of interest, identified by the source localisation

326    analysis (i.e., left motor area), or extend to further regions in the brain. We compared the delta-

327    beta PAC between masks in a whole brain PAC analysis (no-mask and head-mask, correct trials

328    only as the ROI analysis did not established a relationship between PAC and behavioural

329    performance). The difference of trial numbers between conditions was balanced by taking 80%

330    of the smallest sample of available correct trials between all the four conditions ($NMS_{correct}$,

331    $NMA_{correct}$, $HMS_{correct}$ and $HMA_{correct}$; average minimum number of trials: 18.78 ± 5.77).

332    Subsampled trials were concatenated, and the operation was repeated for 40 iterations (to

333    circumvent computational resource limits reached by concatenated epoch lengths). The delta-

334   beta PAC was then averaged across all iterations at each grid (n = 2020) and conditions across

335   participants.

336   **Experimental design and statistical analysis**

337   Audio-visual asynchrony detection task: The experiment used a full within-subject design. The

338   effect of asynchrony and its interaction with the head-mask in audio-visual speech perception

339   was assessed by means of the d' sensitivity index (Macmillan & Kaplan, 1985). To calculate the

340   d' index, the hit trials (i.e., "yes" responses in synchronous conditions NMS and HMS) and false

341   alarm trials (FA, i.e., "yes" responses in the asynchronous conditions NMA and HMA) were

342   computed for each participant. The d' scores for asynchrony detection in the no-mask and

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

344   The $d'$ index allows considering response bias by comparing hits and false alarms to assess

345   whether participants actually discriminated synchrony and asynchrony. Additionally, the

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

347   decision shift between no-mask and head-mask conditions. Further, the mean correct response

348   rates were computed for each participant (hit and correct rejection trials, respectively from the

349   synchronous and asynchronous conditions). Finally, the mean reaction times of the correct

350   trials were computed for each participant (hit and correct rejection trials; comprised between

351   mean reaction times ± two standard deviations range). The effects of masking the speaker's

352   face and audio-visual asynchrony on correct response rates and reaction times were assessed

353   using two-way repeated-measure ANOVAs with the factors mask (no-mask, head-mask),

354   asynchrony (synchronous, asynchronous), and the interaction between mask and asynchrony,

355   using SPSS (IBM Corp. Released 2015. IBM SPSS Statistics for Windows, Version 23.0. Armonk,

356   NY: IBM Corp.). In the case of significant interactions, *post-hoc t*-tests were Bonferroni-

357   corrected. To test whether participants' sensitivity to asynchrony was dependent on

358   information conveyed by head and facial movements, the $d'$ and $c$ criterion in the no-mask and

359   head-mask conditions were individually tested against zero by means of one-sample t-tests.

360   Further, the difference of $d'$ between the no-mask and head-mask conditions was assessed

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

13

362  EEG data at the scalp level: EEG data of correct trials at the scalp level were statistically

363  analysed (NMA$_{correct}$, NMS$_{correct}$, HMA$_{correct}$ and HMS$_{correct}$). We first tested if delta power

364  responses were modulated and dependent on the participants' sensitivity to audio-visual

365  asynchrony in multisensory speech perception. The differences of mean power between the

366  two contrasts NMA-NMS and HMA-HMS (NMA-NMS: difference of power NMA minus NMS;

367  HMA-HMS: difference of power HMA minus HMS) at the electrode level were statistically

368  assessed by applying dependent *t*-tests using Monte-Carlo cluster-based permutation tests

369  (Maris & Oostenveld, 2007) with an alpha cluster-forming threshold set at 0.05, three minimum

370  neighbour channels, 5000 iterations, and cluster selection based on maximum size. Cluster-

371  based permutation statistics were applied for the time window of interest in the delta 2-3 Hz

372  band across all the electrodes. Further, to test whether changes in fronto-central delta

373  oscillations reflect the temporal analysis of multisensory speech rather than purely sensory-

374  driven response activity, we performed the same tests on the theta band (4 - 8 Hz), which

375  tracks the syllabic structure of speech (Giraud & Poeppel, 2012). We expected to find

376  modulations of delta but not theta oscillations for audio-visual asynchrony in the region of

377  interest if motor delta responses reflect temporal analysis. In the identified regions of interest

378  and non-interest (see Results section), normalised mean power across pooled electrodes in the

379  2-3 Hz delta and 4-8 Hz theta frequency bands was computed for the four conditions and

380  exported. This step allowed confirming that delta oscillations responded to audio-visual speech

381  perception independently from the conditions, with an increase of power as compared to the

382  pre-onset baseline (i.e., positive values meaning an increase of power, while negative values

383  meaning a decrease of power in audio-visual speech perception). Statistical differences of

384  power in relevant contrasts were assessed by means of two-way repeated-measure ANOVAs.

385  EEG data source localisation: We tested whether delta power responses at the source level

386  depend on the participants' sensitivity to audio-visual asynchrony in multisensory speech

387  perception. Differences in delta power for the two contrasts NMA-NMS (difference of power

388  NMA minus NMS) and HMA-HMS (difference of power HMA minus HMS) were assessed by

389  applying dependent *t*-tests using Monte-Carlo cluster-based permutation tests at source level

390  as performed for the scalp level analysis. For visualisation of the source localisation results, the

391 power differences in the two contrasts were grand averaged across participants, and the grand

392 average power differences were interpolated to the MNI MRI template for visualization. Only

393 voxels surpassing the statistical significance threshold are depicted in both contrasts (significant

394 *t*-values at alpha = 0.05, multiple comparison cluster-corrected).

395 Delta-beta PAC: Cross-frequency analyses were performed to investigate whether left motor

396 delta-beta PAC is associated with the successful detection of audio-visual speech asynchrony,

397 dependent on whether listeners were able to match visual and auditory prosodies (no-mask

398 conditions) or not (head-mask conditions). First, statistical differences of mean PAC across

399 conditions in the region of interest were assessed applying a three-way repeated-measure

400 ANOVA with the factors mask (no-mask and head-mask), asynchrony (synchronous and

401 asynchronous), and correctness (correct and incorrect trials). Second, statistical differences of

402 whole brain delta-beta PAC were assessed by applying dependent *t*-tests using Monte-Carlo

403 cluster-based permutation tests as described above (whereas *t*-tests were one-tailed here as

404 we had a strong hypothesis about delta-beta PAC modulation directionality based on results at

405 region of interest level).

406 Correlations between performance in synchrony detection and delta oscillations in the

407 identified left motor cluster: We examined the relationship between neural activity and

408 sensitivity to audio-visual asynchrony in multisensory stimuli. By means of Pearson correlations,

409 we tested whether the difference of delta power in the left motor cortex [$\Delta_{power}$ = delta

410 power$_{asynchronous}$ - delta power$_{synchronous}$] predicted differences in correct responses [$\Delta_{CR}$ =

411 CR$_{asynchronous}$ - CR$_{synchronous}$] in the no-mask and head-mask conditions. The purpose of this

412 analysis was to link the increase of left motor delta power in response to audio-visual

413 asynchrony and the participants' sensitivity to the temporal analysis of multisensory speech. A

414 positive correlation between the two variables would establish that an increase of delta power

415 predicts improved asynchrony detection when audio-visual stimuli are asynchronous. Increased

416 sensitivity to asynchrony corresponding with increased delta power would support our

417 hypothesis on the role of left motor cortex in the temporal analysis of audio-visual speech. For

418 each participant, we computed the 2-3 Hz power at the grids sources from the significant

15

419   cluster established in the NMA-NMS contrast source analysis (i.e., significant grids situated in

420   the left central and frontal gyrus areas of interest). Power was averaged across grids in the four

421   conditions separately (NMS$_{correct}$, NMA$_{correct}$, HMS$_{correct}$ and HMA$_{correct}$), and we calculated the

422   mean difference ($\Delta_{Power}$) separately in the no-mask (NMA-NMS) and head-mask (HMA-HMS)

423   contrasts to obtain two delta power values per participant. Similarly, the difference of correct

424   response rates ($\Delta_{CR}$) was calculated in the no-mask and head-mask contrasts, resulting in two

425   behaviour values per participant. The statistical relationship between behaviour ($\Delta_{CR}$) and delta

426   power ($\Delta_{Power}$) was assessed applying Pearson correlation tests.

427   <u>Difference of delta power between correct and incorrect trials across conditions:</u> We tested

428   whether correctness (correct vs incorrect trials) predicted delta power differences in the left

429   motor cortex across conditions (NMA, NMS, HMA and HMS). To circumvent the unbalanced

430   number of trials between correct and incorrect trials within conditions (which was expected

431   according to our experimental procedure targeting ~75-85% of accuracy), we performed

432   permutations tests on the difference of delta power trials$_{correct}$ - trials$_{incorrect}$ between the

433   original data and 5000 permuted data as follows: First, delta power (2-3Hz) in the time-window

434   of interest was computed at source level for all trials and conditions (NMS$_{correct}$, NMS$_{incorrect}$,

435   NMA$_{correct}$, NMA$_{incorrect}$, HMS$_{correct}$, HMS$_{incorrect}$ HMA$_{correct}$ and HMA$_{incorrect}$). Second, correct and

436   incorrect labels were randomly shuffled across trials in each condition. Third, for each iteration

437   two equal samples of shuffled correct and incorrect trials were generated by taking the smallest

438   number of available trials in each condition (i.e. between the original number of correct and

439   incorrect trials). Fourth, the mean delta power from the left motor cluster identified in the

440   source localisation step was computed separately for the shuffled correct and incorrect trials in

441   each condition. Then, the mean difference of delta power trials$_{correct}$ - trials$_{incorrect}$ was

442   computed for each iteration in the NMA, NMS, HMA and HMS conditions. Fifth, in each

443   condition a one-sample t-test against zero (two-tailed) was performed on the difference of

444   delta power trials$_{correct}$ - trials$_{incorrect}$ from the original data to determine the original effect size

445   (t-value$_{original}$), as well as from every permutated data set (i.e. 5000 t-values$_{permut}$). Finally, the

446   5000 t-values from the t-tests were ranked and the *p*-value in each condition was calculated as

16

447    $p$ = [(number of absolute t-values$_{permut}$ +1) > (absolute t-value$_{original}$ +1)]/(number of

448    permutations +1).

449    Distance between delta peak frequencies in the stimulus and delta peak frequencies induced in

450    the left motor cortex: We wanted to confirm that modulation of delta oscillations in left motor

451    cortex does not reflect mere stimulation frequencies, i.e., purely sensory entrainment, but

452    reflects process-driven temporal analysis of visual and auditory prosodies: The rationale was

453    that in the former case, one would assume that tracking the dominant stimulus oscillation

454    would entrain neural delta responses in the exact same frequency. In the latter case, an

455    increase in neural delta activity would reflect the temporal analysis of sensory-specific

456    oscillations independent of their respective frequencies. If true, there should be no direct

457    mapping of the frequency of the delta power maxima in left motor cortex, and the frequency of

458    dominant delta activity conveyed by the multisensory stimuli (see Table 1). To test this

459    assumption, we probed the absolute distance (i.e., absolute difference) between the

460    distribution of peak frequencies of the delta power induced in stimulus perception and the

461    delta peak signal frequencies in the corresponding video clips (Full, Head only, Body only and

462    Audio only signals;  see Table 1). Firstly, we determined the individual delta peak (1-3 Hz) of

463    each participant in every trial (correct trials only: NMS$_{correct}$, NMA$_{correct}$, HMS$_{correct}$ and

464    HMA$_{correct}$) as described previously (see methodology in the previous phase-amplitude coupling

465    section). Secondly, for each trial we calculated the absolute distance between the delta peak

466    frequency of the neural power and every signal of the stimulus presented in the corresponding

467    trial. This step resulted in four absolute distance scores per trial per participant. We averaged

468    the absolute distance scores across participants for each stimulus. Finally, absolute difference

469    scores were sorted by conditions (NMS, NMA, HMS and HMA), and stimulus signals (Full, Head

470    only, Body only and Audio only). To assess statistically the distance between the peak

471    frequencies of delta power and stimuli, we tested the mean of each score distribution against

472    zero with a one-sample t-test (one-tailed). P-values were corrected for multiple comparisons by

473    applying a Bonferroni correction ($\alpha$ = 0.05/total number of comparisons). A one-way repeated-

474    measures ANOVA assessed the statistical difference between the multiple cases of absolute

475    difference (16 in total = two masks x two synchronies x four signals). Similarly, we tested the

17

476    consistency of the delta power frequency maxima in left motor cortex across all trials. This

477    should confirm that any observed variations in delta activity reflect a difference in amplitude

478    modulation on the power of the same delta activity rather than different oscillations across

479    conditions. We computed the delta peak frequency of the neural power of every participant for

480    each stimulus in all four conditions ($NMS_{correct}$, $NMA_{correct}$, $HMS_{correct}$ and $HMA_{correct}$). To

481    statistically assess the consistency of activity across all trials and independent of all conditions,

482    we averaged the EEG delta peak frequency across participants for each stimulus in all four

483    conditions separately (i.e., 54 scores per condition). We then applied a two-way repeated-

484    measure ANOVA with the factors mask (no-mask and head-mask) and synchrony (synchronous

485    and asynchronous).

486    **RESULTS**

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

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

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

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

491    **Listeners successfully temporally analysed visual and auditory prosodic features to denote**

492    **audio-visual asynchrony in multisensory speech perception.**

493                                    **[Insert Figure 2]**

494    *Figure 2.* Behavioural performances in the asynchrony detection task. (A) Average *d'* scores and

495    correct response rates (± standard error of the mean; grey dots represent individual averages; n

496    = 23). (B) Reaction times of correct responses across conditions (± standard error of the mean;

497    grey dots represent individual averages). Significant contrasts are marked by stars ($p < 0.05$).

498    *D'* scores are reported in Figure 2A (left panel). To test whether participants perceived audio-

499    visual asynchrony in both the no-mask and head-mask conditions, we preformed two

500    independent one-sample *t*-tests. Results showed that the mean *d'* was significantly greater than

501    zero in the no-mask and head-mask conditions, confirming that participants were sensitive to

18

502 audio-visual asynchrony in both cases (no-mask: $t(1,22) = 10.25$; $p < 0.001$; *Cohen's d* = 3.04;

503 head-mask: $t(1,22) = 8.07$; $p < 0.001$; *Cohen's d* = 2.38). A paired-samples *t*-test comparing the

504 *d'* between the no-mask and the head-mask conditions tested the hypothesis that participants

505 detected asynchrony better in the no-mask conditions. Results confirmed that it was indeed the

506 case ($t(1,22) = 6.96$; $p$ = < 0.001, two-tailed; *Cohen's d* = 1.46). To assess whether participants

507 tended to respond "synchrony" more often (i.e., a liberal response bias) independently from

508 their actual sensitivity to audio-visual synchrony, and whether this response bias differed when

509 a head-mask was present, we performed two independent one-sample t-tests on the mean

510 criterion c in the no-mask and head-mask conditions. Results revealed that the mean *c* criterion

511 was significantly more negative in the head-mask conditions (-0.53 ± 0.28; $t(1,22) = -9.01$; $p <$

512 0.001; *Cohen's d* = 2.68) but not different from zero in the no-mask conditions (0.11 ± 0.40;

513 $t(1,22) = 1.32$; $p = 0.1$; *Cohen's d* = 0.39). This confirmed that when the speaker's face was head-

514 masked, participants were significantly more biased toward responding "synchrony" than in the

515 no-mask conditions (i.e., a liberal response bias).

516     The mean correct response rates across conditions are depicted in Figure 2A (right panel).

517 NMS: 0.78 ± 0.09; NMA: 0.80 ± 0.13; HMS: 0.82 ± 0.11; HMA: 0.48 ± 0.14. To test whether the

518 presence of the head-mask affected participants' perception of audio-visual asynchrony, we

519 performed a two-way repeated-measure ANOVA with the main factors mask and asynchrony

520 on accuracy. Results confirmed a significant interaction between the mask and asynchrony

521 ($F(1,22) = 82.04$; $p < 0.001$; $\eta_p^2 = 0.789$). Bonferroni-corrected pairwise comparisons showed

522 that performance decreased significantly only in the asynchronous condition of the head-mask

523 conditions (HMA) but not in the three other conditions (NMS, NMA and HMS; no significant

524 difference between them). These results show that the synchrony between visual and auditory

525 stimulus information predicted participants' performance differently, dependent on the

526 presence or absence of the head-mask. The test also revealed a significant main effect of mask

527 ($F(1, 22) = 115.22$, $p < 0.001$; $\eta_p^2 = 0.84$) and asynchrony ($F(1, 22) = 34.52$, $p < 0.001$; $\eta_p^2 = 0.61$)

528 for correct response rates.

529     Reaction times across conditions are reported in Figure 2B. Similarly, a two-way repeated-

530 measure ANOVA with the main factors mask and asynchrony was performed on the reaction

19

531  times. Results revealed a significant main effect of mask on reaction times ($F(1, 22) = 16.50$, $p <$

532  $0.01$; $\eta_p^2 = 0.43$). No significant effect of asynchrony ($F(1, 22) = 0.67$, $p = 0.42$; $\eta_p^2 = 0.03$) or an

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

534  results show that accurate responses were faster when the face of the speaker was not masked

535  compared to head masked.

536

537  Together, the behavioural results support our hypothesis that participants can successfully

538  temporally analyse slow auditory and visual prosodic features in an audio-visual asynchrony

539  detection task. This sensitivity to audio-visual (a)synchrony was altered by the amount of

540  available visual information: On the one hand, the temporal analysis of visual and auditory

541  prosodic information did not change the participants' sensitivity to audio-visual asynchrony in

542  the no-mask conditions. Therefore, the no-mask conditions represent a case of successful

543  temporal analysis in audio-visual speech perception. On the other hand, participants were both

544  slower and less accurate in detecting audio-visual asynchrony in the head-mask conditions,

545  which represents the case of less successful temporal analysis of audio-visual speech

546  perception. Response accuracy in HMS did not differ from the no-mask conditions (although

547  participants were slower in responding correctly), whereas it decreased to chance-level in the

548  asynchronous head-mask condition (HMA). Consequently, the visual mask affected participants'

549  sensitivity to audio-visual asynchrony in both HMS and HMA conditions to a different degree,

550  likely due to the delay between visual and auditory stimulus onsets.

551  **Delta oscillations in the left motor cortex denote asynchrony between the visual and auditory**

552  **prosodies in multisensory speech perception.**

553  We then addressed whether delta oscillations in the left motor cortex relate to the temporal

554  analysis of multisensory information, and whether responses depend on the amount of visual

555  information available. First, a cluster-based permutation tests revealed a significant increase in

556  delta power (2-3 Hz) in response to the audio-visual asynchrony when the speaker's face was

557  visible (no-mask: NMA-NMS) but not when it was masked (head-mask: HMA-HMS) (NMA-NMS:

558  $p < 0.001$, cluster statistic = 117.23; HMA-HMS: No positive cluster; multiple comparisons are

20

559    cluster-corrected). No significant negative clusters were found in both contrasts. Importantly,

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

561    central response when video and audio signals were asynchronous, in line with the expected

562    source localization of delta in the motor region (Figure 3B; Puzzo et al., 2010; Stegemöller et al.,

563    2017). To assess the potential interaction of visual information and audio-visual asynchrony

564    detection in this motor region of interest, we defined a set of electrodes as the region of

565    interest (ROI) representative of the delta response topography: F1, Fz, F2, FFC3h, FFC1h, FFC2h,

566    FFC4h, FC3, FC1, FCz, FC2, FC4, FCC3h, FCC1h, FCC2h, fCC4h, C1, Cz and C2 (Figure 1C). The

567    mean delta power across the electrodes of the ROI was computed separately in the four

568    conditions and confirmed an increase of induced delta activity compared to the pre-stimulus

569    baseline (NMS: 0.64 ± 0.17; NMA: 0.74 ± 0.15; HMS: 0.70 ± 0.16 and HMA: 0.68 ± 0.20; see

570    Figure 3A and 3C). A two-way repeated-measure ANOVA revealed a significant interaction

571    between the factors mask and asynchrony for delta power ($F$(1, 22) = 5.78, $p$ = 0.03; $\eta_p^2$ = 0.21).

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

573    was significantly greater in the asynchronous (NMA) than synchronous (NMS) condition ($p$ =

574    0.02), whereas asynchrony did not affect delta power responses in the head mask contrast ($p$ >

575    0.5). No further pairwise comparison was significant in the post hoc tests. The significant

576    interaction established that the detection of temporal (a)synchrony of visual and auditory

577    information modulated increases in delta power differently and dependent on the availability

578    of  visual information (i.e., no-mask versus head-mask).

579                                        **[Insert Figure 3]**

580    *Figure 3.* Delta responses to audio-visual asynchrony at the scalp level. (A) Time-frequency

581    spectra of the mean power differences in the motor ROI between asynchronous and

582    synchronous conditions in the no-mask (NMA-NMS; left) and head-mask (HMA-HMS; right)

583    contrasts. The white dashed lines correspond to the onset of the video and the window of

584    interest is marked by the pink dashed rectangles. (B) Topographical distribution of the

585    difference of 2-3 Hz delta power in the time-window of interest, in the no-mask (NMA-NMS;

586    top) and head-mask (HMA-HMS; bottom) contrasts. The pink dots display electrodes with

587   significant *t*-values (alpha threshold = 0.05). (C) Delta power across the electrodes of interest in

588   the four conditions (2-3 Hz band). Significant contrasts are marked by stars (*p* < 0.05).

589   Secondly, to separate the influence of audio-visual speech (a)synchrony perception from

590   sensory processing, delta responses were also examined in a control visual region of non-

591   interest (RONI; Figure 1C). The region of non-interest was located in the occipital cortex where

592   we did not expect higher audio-visual speech analysis to take place as visual information was

593   identical between synchronous and asynchronous conditions within mask contrasts (RONI

594   electrodes: PPO1h, PPO2h, PO3, POz, PO4, POO1, POO2, POO9h, O1, Oz, O2, POO10h, Ol1h,

595   Ol2h, O9 and O10). We compared the effect of audio-visual asynchrony between the identified

596   motor region (ROI) and the visual sensory area (RONI) to confirm that delta response

597   modulations did not reflect signal processing only (Figure 4A). The mean differences of 2-3Hz

598   delta power (NMA-NMS and HMA-HMS) were computed in the regions of interest and non-

599   interest at the same time-window (Figure 4B; ROI: NMA-NMS = 0.1 ± 0.09; HMA-HMS = -0.03±

600   0.19; RONI: NMA-NMS = 0.05 ± 0.10; HMA-HMS = 0.01 ± 0.24). A two-way repeated-measures

601   ANOVA with the mean factors region (ROI or RONI) and mask (no-mask or head-mask) was

602   performed to assess whether the responses of delta oscillations to asynchrony reflected

603   multisensory speech analysis or purely signal processing taking place in sensory areas (i.e.,

604   visual occipital areas). Results revealed a significant interaction between region and mask ($F$(1,

605   22) = 5.75, *p* = 0.025; $\eta_p^2$ = 0.21). First, Bonferroni-corrected pairwise comparisons showed that

606   in the no-mask contrast the delta power difference NMA-NMS (but not HMA-HMS) was

607   significantly greater in the region of interest than in the region of non-interest (respectively p =

608   0.025 and p = 0.572). Only in the region of interest the difference of power NMA-NMS was

609   significantly greater than HMA-HMS (respectively *p* = 0.019 and *p* = 0.113). No main effect of

610   mask ($F$(1, 22) = 0.25, *p* = 0.622; $\eta_p^2$ = 0.21) or region ($F$(1, 22) = 2.18, *p* = 0.154; $\eta_p^2$ = 0.09) was

611   found.

612                                                   **[Insert Figure 4]**

613   *Figure 4*. Comparisons between the motor region of interest (ROI) and the visual region of non-

614   interest (RONI). (A) TFRs of the difference of spectrum in the no-mask contrast (NMA-NMS) in

615    the ROI and RONI. (B) The mean differences of 2-3Hz delta power (NMA-NMS and HMA-HMS)

616    were computed in the regions of interest and non-interest. Significant contrasts are marked by

617    stars ($p < 0.05$).

618    Thirdly, the mean power in the 4 - 8 Hz band was computed in the four conditions separately

619    from the ROI electrodes and confirmed an increase of theta activity compared to the pre-

620    stimulus onset baseline (NMS: 0.86 ± 0.25; NMA: 0.85 ± 0.18; HMS: 0.83 ± 0.16 and HMA: 0.81

621    ± 0.24). A two-way repeated-measure ANOVA revealed no significant main effect of mask ($F$(1,

622    22) = 2.77, $p$ = 0.11; $\eta_p^2$ = 0.11), asynchrony ($F$(1, 22) = 0.27, $p$ = 0.606; $\eta_p^2$ = 0.01) or interaction

623    between the factors mask and asynchrony ($F$(1, 22) = 0.05, $p$ = 0.825; $\eta_p^2 < 0.01$) on theta power

624    in the region of interest. Further, the cluster-based permutation tests revealed no significant

625    modulation of theta power by audio-visual asynchrony in any of the mask contrasts (NMA-NMS:

626    no significant cluster; HMA-HMS: no significant cluster; multiple comparisons are cluster-

627    corrected). These results confirmed that audio-visual asynchrony detection modulated delta

628    power over the expected fronto-central region. Further, the delta power response was

629    attenuated when listeners were less able to integrate visual and auditory prosodies (i.e., in the

630    head-mask as compared to the no-mask conditions). This result suggests that increased delta

631    activity in left motor cortex plays a role in audio-visual asynchrony detection as it only increased

632    in asynchronous but not synchronous multisensory speech perception. Therefore, delta activity

633    might be associated with the brain's effort to resolve mismatches between visual and auditory

634    prosodies in the temporal analysis of multisensory speech.

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

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

637    Cluster-based permutation $t$-tests between synchronous and asynchronous conditions at the

638    source level revealed that asynchrony significantly increased delta oscillation responses when

639    the head of the speaker was visible (NMA-NMS: $p$ = 0.042; cluster statistic = 233.02) but not

640    when it was head-masked (HMA-HMS: $p$ = 0.27; cluster statistic = 38.27). The projections of the

641    significant $t$-values on the brain's surface showed an increase of delta power originating mainly

642    in the left precentral region and the left inferior frontal gyrus (Figure 5A). The source results

643 support the topographies of the delta power modulations observed at the scalp level, which

644 revealed fronto-central differences in the no-mask contrast only (Figure 3B). Similar to the scalp

645 level analysis, we computed the mean 2-3Hz power across the significant grids in all four

646 conditions in the time-window of interest. Power was normalised relative to the pre-stimulus

647 baseline to determine an increase of delta power during stimulus presentation in all four

648 conditions. Four one-sample *t*-tests against zero confirmed a significant increase of delta power

649 in response to audio-visual speech perception in all four conditions (respectively NMS: 0.69 ±

650 0.15; $t(1,22)$ = 22.25; *p* < 0.001; *Cohen's d* = 4.64; NMA: 0.76 ± 0.17; $t(1,22)$ = 21.81; *p* < 0.001;

651 *Cohen's d* = 4.55; HMS: 0.70 ± 0.17; $t(1,22)$ = 19.84; *p* < 0.001; *Cohen's d* = 4.14 and HMA: 0.69 ±

652 0.20; $t(1,22)$ = 16.74; *p* < 0.001; *Cohen's d* = 3.49). We then performed a two-way repeated-

653 measure ANOVA with the main factors mask and asynchrony on mean power as in the scalp

654 level analysis, but the test did not reveal any significant effects (mask: $F(1, 22)$ = 1.42, *p* = 0.25;

655 $\eta_p^2$ = 0.061; asynchrony: $F(1, 22)$ = 1.75, *p* = 0.20; $\eta_p^2$ = 0.074; mask*asynchrony: $F(1, 22)$ = 1.94,

656 *p* = 0.18; $\eta_p^2$ = 0.081). Further, we tested whether the modulation of delta responses in the left

657 motor areas by audio-visual asynchrony predicted detection performance in the no-mask and

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

659 the correct response rate differences ($\Delta_{CR}$) and delta power differences ($\Delta_{Power}$) in the no-mask

660 contrast (NMA-NMS: *r* = 0.36; *p* = 0.046, one-tailed) but not in the head-mask contrast HMA

661 minus HMS (HMA-HMS: *r* = 0.04; *p* = 0.43, one-tailed). These results confirmed that when

662 participants perceived asynchrony between video and audio signals (no-mask conditions), the

663 difference in delta power between asynchronous and synchronous conditions predicted

664 detection accuracy. This was not the case when participants were less able to detect temporal

665 alignment between visual and auditory information (head-mask conditions).

666 **[Insert Figure 5]**

667 *Figure 5.* Delta oscillation responses to audio-visual asynchrony at the source level for no-mask

668 and head-mask contrasts. (A) Contrast NMA – NMS projected onto the brain's surface

669 (significance *t*-values; cluster-corrected at alpha threshold = 0.05). The maximum voxel MNI

670 coordinates is located left precentrally [-50 19 40] but significant activation was also found in

24

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 – HMS

contrast; not represented). (B) Scatterplots of audio-visual asynchrony detection performance

and delta power in the significant cluster region (left motor cortex). The difference of delta

power in the left motor cluster ($\Delta_{Power}$; x-axis; z scores) correlated with the difference of audio-

visual asynchrony detection ($\Delta_{CR}$; y-axis; z scores) between asynchronous and synchronous

conditions only when the face of the speaker was visible, and participants could integrate video

and audio onsets (no-mask conditions). (C) Average delta power differences between correct

and incorrect trials from the significant left motor cluster in the four conditions NMS, NMA,

HMS and HMA (± standard error of the mean; grey dots represent individual averages; *n* = 23;

outliers not represented). Significant differences from zero are marked by stars ($p < 0.05$). (D)

Left panel: Peak frequency correspondence between delta activity carried in the video clips and

delta power responses induced in the left motor cluster. The bars represent the mean absolute

distance between the delta peak frequencies in the stimulus and the peak frequencies of neural

delta power induced during the corresponding trial (± standard error of the mean). Peak

frequency matching was assessed for the synchronous and asynchronous conditions in the two

mask conditions (no-mask and head-mask), and the different signal types of each stimulus: Full,

Head only, Body only and Audio only (see Table 1). The mean of the absolute difference scores

were significantly greater than zero in all conditions and for all the stimulus signals. (D) Right

panel: Consistency of delta peaks across the ordered stimuli in all four conditions. The upper

panel displays the mean delta peak frequencies in the left motor cortex across all participants

(± standard error of the mean) for each stimulus in the no-mask conditions (black squares:

NMS; orange squares: NMA). The lower panel displays the mean EEG delta peak frequencies

across all participants (± standard error of the mean) for each stimulus in the head-mask

conditions (black squares: HMS; orange squares: HMA). The variations in delta responses

observed across all conditions reflect a difference of power amplitude modulation on the same

oscillatory activity.

We tested whether correctness predicted delta power response modulations in the significant

cluster identified in the left motor cluster (Figure 5C). Permutation tests revealed that delta

25

700 power from the left motor cortex was significantly greater in the correct trials as compared to

701 incorrect trials in the asynchronous conditions no-mask (NMA: *Cohen's d* = 0.552; *p* = 0.002)

702 and head-mask (HMA: *Cohen's d* = 0.401; p = 0.011). In contrast, no significant difference of

703 delta power between correct and incorrect trials were found in the synchronous conditions

704 (NMS: *Cohen's d* = 0.28; *p* = 0.429; HMS: *Cohen's d* = 0.232; *p* = 0.332). These results showed

705 that increases of delta power in the left motor cortex predicted sensitivity to audio-visual

706 alignment in the asynchronous conditions but not in the synchronous conditions. Further, we

707 aimed to control that delta responses induced in the left motor cortex during audio-visual

708 speech perception did not reflect purely stimulus driven entrainment to the delta activity

709 carried in the visual or auditory signals of the video clips (Figure 5D left panel). The mean peak

710 of delta power in the left motor cortex in the four conditions was respectively for NMS: 2.06 ±

711 0.13 Hz; NMA: 2.07 ± 0.18 Hz; HMS: 2.03 ± 0.17 Hz and HMA: 2.06 ± 0.21 Hz. To statistically

712 assess the distance between the peak frequencies of left motor delta power and stimulus delta

713 activity, we tested the mean of each score distribution against zero with a one-sample t-test

714 (one-tailed). Results revealed that the mean absolute distance was significantly greater than

715 zero in all conditions ($p < p_{corrected}$). Further, a one-way repeated-measure ANOVA tested the

716 difference of absolute distance between all conditions. Results revealed a significant effect of

717 condition ($F(15,848) = 2.995$; $p < 0.001$; $\eta_p^2 = 0.05$). However, Bonferroni-corrected pairwise

718 comparisons revealed only a single marginal tendency for a difference between the Full$_{no-mask}$

719 $_{synchronous}$ and Head$_{head-mask\ synchrony}$ absolute distances ($p = 0.09$; 1$^{st}$ and 11$^{th}$ bars on Figure 5D).

720 These results confirm that the delta responses induced in the left motor cortex significantly

721 deviated from stimulus-related delta frequency, thus did not just reflect entrainment. Finally,

722 concerning the consistency of neural delta power across trials in all conditions, results revealed

723 no significant main effect of mask or an interaction between mask and asynchrony for motor

724 delta peak frequency (mask: $F(1,53) = 1.84$; $p = 0.181$; $\eta_p^2 = 0.034$; asynchrony: $F(1,53) = 1.11$;

725 $p = 0.741$; $\eta_p^2 = 0.002$; interaction between mask*asynchrony: $F(1,53) = 1.33$; $p = 0.717$; $\eta_p^2 =$

726 0.003; see Figure 5D right panel). These results confirm that any observed variation in motor

727 delta activity between the experimental conditions cannot be explained as mere stimulus

728 frequency.

26

729 **Delta-beta PAC reflects sensitivity to audio-visual temporal asynchrony in speech perception**

730 **but is not limited to the left motor cortex.**

731 Finally, we assessed whether delta-beta PAC modulations in the left motor area reflect

732 sensitivity to audio-visual asynchrony in speech perception. First, a three-way repeated-

733 measure ANOVA (main factors: mask, asynchrony and correctness) revealed a main effect of

734 mask on delta-beta PAC with delta-beta phase-coupling being significantly greater in the no-

735 mask than in the head-mask conditions ($F(2,22) = 4.72$; $p = 0.041$; $\eta_p^2 = 0.18$; see Figure 6A). No

736 further significant main effects or interactions were found. These results show greater left

737 motor cortex delta-beta PAC when participants were more sensitive to asynchronous audio-

738 visual speech in the no-mask conditions than when they were less able to match visual and

739 auditory prosodic features (head-mask conditions). Nevertheless, we cannot fully discard that

740 delta-beta PAC also increased in the head-mask condition as compared to a control baseline

741 condition (e.g., visual or auditory only condition). Secondly, we investigated whether the delta-

742 beta PAC difference between no-mask and head-mask conditions was restricted to the left

743 motor areas. As accuracy and asynchrony did not affect delta-beta PAC in the cluster of

744 interest, we selected only correct trials for the delta-beta PAC analysis at the whole brain level

745 and combined synchronous and asynchronous trials within no-mask and head-mask conditions

746 (i.e., NMCs: NMA+ NMS; HMCs: HMA + HMS). The cluster-based permutation tests revealed

747 one significant positive cluster peaking in the superior motor area and in the left middle

748 temporal lobe (although not exclusively; see Figure 6B), confirming that delta-beta PAC was

749 significantly larger in the no-mask (NMCs) compared to the head-mask (HMCs) case (NMCs -

750 HMCs : $p = 0.043$, cluster statistic = 216.69).

751                                             **[Insert Figure 6]**

752 *Figure 6*: Phase-amplitude coupling between delta and beta oscillations. (A) PAC analysis in the

753 left motor cluster. The figure represents the modulation of delta-beta PAC in a significant

754 cluster, dependent on the mask and audio-visual asynchrony. Significance is indicated by an

755 asterisk ($p < 0.05$, Bonferroni-corrected). Delta-beta PAC from the left motor cortex was greater

756 in the no-mask than the head-mask conditions but did not discriminate between correct and

27

757 incorrect trials. Significant contrasts are marked by stars (*p* < 0.05). (B) Delta-beta PAC

758 difference between no-mask (NMA+ NMS) and head-mask (HMA + HMS) case in the whole

759 brain. Results revealed significant maximum differences located in the superior motor area

760 (MNI coordinates [0 11 50]) and in the left middle temporal lobe (MNI coordinates [-50 -1 -20]).

761     In summary, the EEG results mirrored the behavioural results as modulations in left motor

762 delta power reflect the successful detection of audio-visual asynchrony when participants were

763 able to see face and visible articulators (no-mask conditions), but not in the head-mask

764 conditions. An increase in left motor delta power only predicted differential sensitivity to audio-

765 visual asynchrony in the no-mask conditions and related to correctly perceiving asynchronous

766 audio-visual speech. Importantly, a control analysis confirmed that variations in left motor delta

767 activity reflect an amplitude difference based on the same oscillatory activity across all stimuli

768 rather than oscillation differences of stimulation *per se*. Lastly, delta-beta PAC in the left motor

769 cortex was greater when listeners detected audio-visual asynchrony more accurately during

770 speech perception (i.e., no-mask as compared to head-mask conditions). Nevertheless, this

771 result did not exclude that delta-beta PAC also increased in the head-mask conditions, but to a

772 lesser extent.

773 **DISCUSSION**

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

775 analysis of multisensory prosodic features in speech perception. The behavioural results of the

776 audio-visual asynchrony detection task confirmed that listeners processed both prosodies in

777 multisensory speech perception when sufficient visual information was available. At the brain

778 level, the perception of audio-visual asynchrony induced an increase in left motor delta activity

779 (extending to the inferior frontal gyrus). Further, the difference of delta power between

780 asynchronous and synchronous conditions predicted participants' sensitivity of audio-visual

781 asynchrony. In contrast, participants were less able to discriminate audio-visual information

782 when a speaker's facial information was masked. This is evident in the absence of difference in

783 delta activity between asynchronous and synchronous conditions. Finally, delta-beta PAC in the

784 left motor cortex was significantly greater when listeners were more accurate in perceiving

785  asynchrony between visual and auditory information during multisensory speech perception
786  (no-mask vs. head-mask conditions). Altogether, the current results indicate that the delta
787  timescale provides a flexible framework to synchronise a listener's brain activity with
788  multisensory speech input. Thus, motor delta activity seems to play a role in detecting temporal
789  mismatches between visual and auditory prosodies and is as a sensitive measure of (un-
790  )successful temporal analysis in multisensory speech perception.

791  Behaviourally, the results of the asynchrony detection task confirm our first hypothesis, that
792  is, listeners temporal analyse prosodic events in multisensory speech perception. This finding
793  was expected as visual information complements auditory information and often improves
794  speech perception (Sumby & Pollack, 1954; van Wassenhove et al., 2005). Speaker's
795  articulatory movements and gestures temporally aligned with acoustic prosodic cues, providing
796  listeners with a reliable temporal structure of the speech signal in the delta range (Biau et al.,
797  2016; Esteve-Gibert & Guellaï, 2018; Wagner et al., 2014). Participants likely use these salient
798  prosodic events as landmarks to align them into a coherent multisensory speech percept. The
799  results suggest that successful temporal analysis can focus the listeners' attention within brief
800  time-windows containing complementary multisensory prosodic events. This is in line with the
801  theory of dynamic attending, stating that non-random external stimulation drives periodic
802  attention allocation towards critical events (Large & Jones, 1999). Noteworthy, the differences
803  of performance between the no-mask and head-mask conditions indicate that participants
804  likely relied on complementary information conveyed by the speaker's head, face, and fine
805  articulatory gestures to achieve the integration of the visual prosodic signal (Cross, Butler, &
806  Lalor, 2015). Of note, when the speaker's face was masked, participants' response accuracy
807  decreased significantly while they remained sensitive to the temporal alignment of the audio-
808  visual signals in the synchronous condition (HMS). Our results suggest that participants adopted
809  a liberal guessing strategy and tended to respond "synchronous" more often in the head-mask
810  conditions (i.e., negative $c$ criterion and decrease of d' as compared to the no-mask conditions).
811  Therefore, we assume that if participants were not sensitive at all to audio-visual temporal
812  alignment, such a bias would have only increased and led to responding "synchronous" even
813  more systematically. Consequently, correct response rates would have significantly increased in

29

814 the HMS as compared to the NMS and NMA conditions. Our results show that this is not the

815 case as HMS accuracy was equivalent to NMS and NMA accuracy. Rather, performance in the

816 HMA condition decreased to chance level. While somewhat speculative, a comparable delta

817 power increase in the HMS and NMA conditions at the scalp level (Figure 3C) may reflect a

818 similar increased effort to reach comparable accuracy levels when integrating a blurred visual

819 signal with an auditory signal. Such an increased effort was not observed in the HMA condition,

820 potentially due to the larger delay between visual and auditory signal onsets that prevented

821 participants to align them. In other words, increased delta power reflects an analytic effort and

822 explains comparable delta power patterns in the NMA and HMS conditions. Lastly, although we

823 applied a unique head-mask to obscure visual facial prosody, this is technically a gradual

824 masking approach because blurring the face did not prevent the participants from using other

825 available visual prosodic information (e.g., head nods, upper parts of the speaker's body, and

826 hand gestures). Nevertheless, future studies could adopt a more fine-grained gradual masking

827 approach by using different levels of visual degradation, masking different effectors (e.g.,

828 mouth, head, hands, and breathing) to examine which movements best carry information

829 needed for successful temporal analysis in multisensory speech perception.

830 The EEG results confirmed an increase in motor delta activity in response to audio-visual

831 asynchrony detection, extending the role of delta activity to the temporal analysis of

832 multisensory prosodies. Previous literature associated delta oscillations in the motor cortex

833 with the perception of auditory rhythmic stimulation (Keitel et al., 2018; Morillon et al., 2019;

834 Morillon & Schroeder, 2015). The present results extend these findings to the temporal analysis

835 of non-isochronous events that act as punctual "snap fasteners" streaming visual and auditory

836 signals within relevant time-windows. As long as they provide the brain with sufficient time for

837 the temporal analysis of multiple sensory inputs, salient prosodic features do not have to be

838 perfectly regular to trigger delta motor responses. The present EEG results corroborate this

839 hypothesis in three ways: First, we did not observe different delta responses in auditory and

840 visual cortices when audio-visual stimuli were synchronous. This would have reflected low-level

841 feature tracking in early sensory processing (Cross, Butler & Lalor, 2015; Ghitza, 2017; Gross et

842 al., 2013; Mai, Minett & Wang, 2016). Further, a control analysis confirmed that delta

843  responses in left motor cortex did not simply reflect stimulus entrainment as they significantly

844  differed from the frequency of the audio-visual stimuli. Next, audio-visual asynchrony would

845  likely decrease pure entrainment by making signal tracking more difficult than when different

846  channels of the same input are processed synchronously. Further, we found no theta activity in

847  response to audio-visual asynchrony at the scalp level that would have indicated an effect

848  driven specifically by the rate of the prosodic features (e.g., lip movements). Additionally,

849  differences in left motor cortex delta power only predicted accuracy in the no-mask contrast.

850  Moreover, delta power increased more for accurate than inaccurate responses in the

851  asynchronous conditions, independent of the presence or absence of a head-mask. Lastly,

852  participants perceived audio-visual synchrony less accurately when the speaker's facial

853  information was blurred. This was shown in weaker delta motor responses and that

854  synchronous and asynchronous conditions displayed not differences in delta power. Together,

855  these results confirm that left motor delta oscillations might reflect the successful detection of

856  audio-visual asynchrony, likely linked to the temporal analysis of multisensory speech.

857  Importantly, the responses found in the left inferior frontal gyrus align well with previous

858  research that established a role in cross-modal information integration between gestures and

859  speech (Park et al., 2018; Willems, Ozyürek & Hagoort, 2009; Zhao et al., 2018). Here,

860  participants perceived information carried by two modalities, and integrated gestures'

861  kinematics with auditory envelope modulations to perform an asynchrony detection task.

862  Further investigations will need to address whether the response modulations in the left IFG

863  were specific to the temporal integration of gesture and speech or could be reproduced using

864  moving dots following gestures' dynamics (Biau et al., 2016; Holle et al., 2012). In contrast, we

865  found no differential activation in further brain regions associated with multisensory speech

866  integration such as the left posterior superior temporal sulcus (Marstaller & Burianová, 2014).

867  Here, delta oscillations did not reflect multisensory integration *per se* but the temporal

868  alignment of multisensory information. It is worth noting that the present study focused on the

869  temporal analysis of visual and auditory prosodies and addressed how they (mis-)align. It could

870  be of further interest to look more closely into the temporal dynamics between sensory areas

871  in multisensory speech perception. For instance, comparing delta phase offsets between

31

872    synchronous and asynchronous conditions could help to understand whether the

873    synchronisation of delta oscillations between visual and auditory areas predicts delta responses

874    in the left motor cortex. Future studies may overcome the limitations of the current study to

875    perform source reconstruction analysis (e.g., including visual and auditory only conditions as

876    localizers), and address the role of delta synchronisation between sensory areas in multisensory

877    speech perception.

878        Finally, delta-beta coupling in the left motor cortex was larger in the no-mask conditions,

879    when listeners perceived audio-visual temporal alignment in both directions (i.e., synchronous

880    or asynchronous). Although somewhat speculative, delta-beta coupling might take place after

881    proper temporal analysis of visual and auditory prosodic features has occurred and might

882    support top-down predictions (e.g., auditory-motor coupling). For instance, Park et al. (2015)

883    showed that the left frontal areas modulated the phase of delta oscillations in the left auditory

884    cortex by means of top-down control in speech perception. Reciprocally, delta-beta PAC in the

885    auditory cortex respond to the modulations of rhythmic regularity in auditory speech

886    perception (Chang, Bosnyak and Trailor, 2019). Further, Keitel et al. (2018) reported that delta-

887    beta PAC in the left motor cortex predicted behavioural performance in speech comprehension.

888    Future research will need to unravel whether delta-beta coupling provides a ubiquitous means

889    of cross-regional communication to align temporally different dynamic input in sensory cortices

890    (Arnal, 2012; Fujioka, Ross & Trainor, 2015; Morillon et al., 2019). For example, Fontolan et al.

891    (2014) reported that delta-beta coupling in the associative auditory cortex modulated the

892    phase of gamma activity related to phonological processing in the primary auditory cortex in

893    auditory sentence perception (Giraud & Poeppel, 2012). Alternatively, delta-beta PAC may drive

894    the periodicity of attention to critical time-windows containing relevant accentuated speech

895    information, which fits with the dynamic attention theory (Large & Jones, 1999). It is important

896    to note that delta-beta PAC may increase in the head-mask conditions as well, but simply to a

897    lesser extent than in the no-mask conditions. If true, top-down predictions may be generated

898    during multisensory speech perception even when participants were less successful in detecting

899    audio-visual temporal (mis)alignment.

900    We propose that motor delta oscillations mirror the successful detection of asynchronous

901    multisensory prosodies, encoded separately in auditory and visual sensory cortices. The slow

902    timescale of delta (1-3Hz) may also offer the brain some flexibility to create a coherent

903    multisensory percept despite the natural delay between visual and auditory signal onsets in

904    speech (Chandrasekaran et al., 2009). In social interactions where conditions change quickly,

905    such a delta framework would help listeners to align speech information in a bottleneck fashion

906    to maintain stable synchronization in speech flow (Kotz, Ravignani & Fitch, 2018). When two

907    dynamic events cannot be integrated in a critical delta time-window due to their temporal

908    offsets, any effort to resolve such an audio-visual mismatch increases and shows in amplified

909    motor delta activity. At a certain point, i.e., when onsets of visual and auditory prosody onsets

910    mismatch ($\sim$ 400ms), delta power reaches a critical threshold, leading to the successful

911    detection of audio-visual asynchrony in speech. Further investigations will need to address

912    whether this with other timescales present in both the speech signal and brain oscillations. For

913    instance, we cannot fully discard that the prosodic contour in our stimuli still contained a

914    syllable structure embedded in it (e.g., at onsets and stress peaks). Further, lip movements and

915    auditory envelope convey syllabic information occurring at a theta rate (4-8 Hz) providing other

916    robust temporal information in the speech signal during face-to-face conversations

917    (Chandrasekaran et al., 2009; Giraud & Poeppel, 2012). Therefore, delta and theta activities

918    may actually couple to strengthen speaker-listener synchronization in social communicative

919    interactions.

920    **CONCLUSION**

921    Our findings show that left motor delta oscillations play a role in audio-visual asynchrony

922    detection of visual and auditory prosodies, and by extension contribute to the successful

923    temporal analysis of multisensory speech. We propose that a critical delta time window allows

924    for the (un-)successful temporal alignment of dynamic prosodic features, conveyed by distinct

925    sensory modalities in speech perception.

926    **ACKNOWLEDGMENT**

**RESOURCE SHARING**

Consent for sharing data at the level of the individual participant was received. Data for individual participants and associated scripts will be made available upon publication of the manuscript. Further information or requests should be directed to the corresponding authors.

**REFERENCES**

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

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

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

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

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

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

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

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

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

956 Crosse, M. J., Di Liberto, G. M., & Lalor, E. C. (2016). Eye can hear clearly now: inverse
957     effectiveness in natural audiovisual speech processing relies on long-term crossmodal
958     temporal integration. *The Journal of Neuroscience*, *36*(38), 9888–9895.
959     https://doi.org/10.1523/JNEUROSCI.1396-16.2016

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

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

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

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

970 Fujioka, T., Ross, B., & Trainor, L. J. (2015). Beta-Band Oscillations Represent Auditory Beat and
971     Its Metrical Hierarchy in Perception and Imagery. *Journal of Neuroscience*, *35*(45), 15187–
972     15198

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

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

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

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

984    Gunter, T. C., & Douglas Weinbrenner, J. E. (2017). When to Take a Gesture Seriously: On How
985         We Use and Prioritize Communicative Cues. *Journal of Cognitive Neuroscience, 29*(8),
986         1355-1367

987    Holle, H., Obermeier, C., Schmidt-Kassow, M., Friederici, A. D., Ward, J., & Gunter, T. C. (2012).
988         Gesture facilitates the syntactic analysis of speech. *Frontiers in psychology*, *3*, 74.
989         https://doi.org/10.3389/fpsyg.2012.00074

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

992    Keitel, A., Ince, R. A. A., Gross, J., & Kayser, C. (2017). Auditory cortical delta-entrainment
993         interacts with oscillatory power in multiple fronto-parietal networks. *NeuroImage*, *147*,
994         32–42. https://doi.org/10.1016/j.neuroimage.2016.11.062

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

997    Kösem, A., Gramfort, A., & van Wassenhove, V. (2014). Encoding of event timing in the phase of
998         neural                          oscillations. *NeuroImage*, *92*,                          274–284.
999         https://doi.org/10.1016/j.neuroimage.2014.02.010

1000   Kösem, A., & van Wassenhove, V. (2017). Distinct contributions of low- and high-frequency
1001        neural oscillations to speech comprehension. *Language, Cognition and Neuroscience*,
1002        *32*(5), 536–544

1003   Kotz S.A., Ravignani A., Fitch W.T. The Evolution of Rhythm Processing. *Trends Cogn Sci*.
1004        2018;22(10):896-910. doi:10.1016/j.tics.2018.08.002

1005   Large, E.W., & Jones, M.R. (1999). The dynamics of attending: how people track time-varying
1006        events. *Psychol. Rev.* 106, 119

1007   Mai, G., Minett, J. W., & Wang, W. S. Y. (2016). Delta, theta, beta, and gamma brain oscillations
1008        index levels of auditory sentence processing. *NeuroImage*, *133*, 516–528.
1009        https://doi.org/10.1016/j.neuroimage.2016.02.064

1010  Maris, E., & Oostenveld, R. (2007). Nonparametric statistical testing of EEG- and MEG-data.
1011       *Journal of Neuroscience Methods*, *164*(1), 177–190

1012  Marstaller L, Burianov_a H.(2014). The multisensory perception of co-speech gestures - a
1013       review and meta-analysis of neuroimaging studies. *J Neurolinguistics, 30*, 69-77.

1014  Macmillan, N. A., & Kaplan, H. L. (1985). Detection theory analysis of group data: Estimating
1015       sensitivity from average hit and false-alarm rates. *Psychological Bulletin*, 98, 185–199.

1016  Mercier MR, Molholm S, Fiebelkorn IC, Butler JS, Schwartz TH, Foxe JJ. (2015). Neuro-oscillatory
1017       phase alignment drives speeded multisensory response times: an electro-corticographic
1018       investigation. J Neurosci 35: 8546–8557.

1019  Meyer, L., Sun, Y., & Martin, A. E. (2019). Synchronous, but not entrained: exogenous and
1020       endogenous cortical rhythms of speech and language processing. *Language, Cognition and*
1021       *Neuroscience*, 1–11. https://doi.org/10.1080/23273798.2019.1693050

1022  Morillon, B., & Schroeder, C. E. (2015). Neuronal oscillations as a mechanistic substrate of
1023       auditory temporal prediction. *Annals of the New York Academy of Sciences*, *1337*(1), 26–
1024       31. https://doi.org/10.1111/nyas.12629

1025  Morillon, B., & Baillet, S. (2017). Motor origin of temporal predictions in auditory attention.
1026       *Proceedings of the National Academy of Sciences of the United States of America*, *114*(42),
1027       E8913–E8921

1028  Morillon, B., Arnal, L. H., Schroeder, C. E., & Keitel, A. (2019). Prominence of delta oscillatory
1029       rhythms in the motor cortex and their relevance for auditory and speech perception.
1030       *Neuroscience and Biobehavioral Reviews*, Vol. 107, pp. 136–142.
1031       https://doi.org/10.1016/j.neubiorev.2019.09.012

1032  Munhall, K. G., Jones, J. A., Callan, D. E., Kuratate, T., & Vatikiotis-Bateson, E. (2004). Visual
1033       prosody and speech intelligibility: head movement improves auditory speech perception.
1034       *Psychological Science*, *15*(2), 133–137

1035  Obermeier, C., Dolk, T., & Gunter, T. (2012). The benefit of gestures during communication:
1036       Evidence from hearing and hearing-impaired individuals. *Cortex, 48*, 857-870

37

1037   Obermeier, C., & Gunter, T. C. (2014). Multisensory Integration: The Case of a Time Window of
1038        Gesture-Speech Integration. *Journal of Cognitive Neuroscience*, 1–16

1039   Oostenveld R, Fries P, Maris E, Schoffelen JM (2011) FieldTrip: open-source software for
1040        advanced analysis of MEG, EEG, and invasive electrophysiological data. Comput Intell
1041        Neurosci 2011:156869

1042   Park, H., Ince, R. A. A., Schyns, P. G., Thut, G., & Gross, J. (2015). Frontal Top-Down Signals
1043        Increase Coupling of Auditory Low-Frequency Oscillations to Continuous Speech in Human
1044        Listeners. *Current Biology*, *25*(12), 1649–1653

1045   Park, H., Kayser, C., Thut, G., & Gross, J. (2016). Lip movements entrain the observers' low-
1046        frequency brain oscillations to facilitate speech intelligibility. *ELife*, *5*

1047   Park, H., Ince, R., Schyns, P. G., Thut, G., & Gross, J. (2018). Representational interactions during
1048        audiovisual speech entrainment: Redundancy in left posterior superior temporal gyrus and
1049        synergy      in      left      motor      cortex. *PLoS      biology*, *16*(8),      e2006558.
1050        https://doi.org/10.1371/journal.pbio.2006558

1051   Peelle, J. E., & Davis, M. H. (2012). Neural Oscillations Carry Speech Rhythm through to
1052        Comprehension. *Frontiers in Psychology*, *3*, 320

1053   Puzzo, I., Cooper, N. R., Vetter, P., & Russo, R. (2010). EEG activation differences in the pre-
1054        motor cortex and supplementary motor area between normal individuals with high and
1055        low traits of autism. *Brain Research*, *1342*, 104–110

1056   Saleh, M., Reimer, J., Penn, R., Ojakangas, C. L., & Hatsopoulos, N. G. (2010). Fast and Slow
1057        Oscillations in Human Primary Motor Cortex Predict Oncoming Behaviorally Relevant Cues.
1058        *Neuron*, *65*(4), 461–471

1059   Schultz, B. G., Biau, E., & Kotz, S. A. (2020). An open-source toolbox for measuring dynamic
1060        video framerates and synchronizing video stimuli with neural and behavioral responses.
1061        *Journal of Neuroscience Methods*, 108830

1062   Stegemöller, E. L., Allen, D. P., Simuni, T., & MacKinnon, C. D. (2017). Altered premotor cortical
1063        oscillations during repetitive movement in persons with Parkinson's disease. *Behavioural
1064        Brain Research*, *317*, 141–146

1065   Sumby, W. H., & Pollack, I. (1954). Visual Contribution to Speech Intelligibility in Noise. *Journal
1066        of the Acoustical Society of America*, *26*(2), 212–215

1067    Tadel F, Baillet S, Mosher JC, Pantazis D, Leahy RM (2011) Brainstorm: a user-friendly
1068        application for MEG/EEG analysis. Comput Intell Neurosci 2011 :879716

1069    Thut, G., Veniero, D., Romei, V., Miniussi, C., Schyns, P., & Gross, J. (2011). Rhythmic TMS
1070        causes local entrainment of natural oscillatory signatures. *Current Biology*, *21*(14), 1176–
1071        1185. https://doi.org/10.1016/j.cub.2011.05.049

1072    Tort, A. B. L., Komorowski, R., Eichenbaum, H., & Kopell, N. (2010). Measuring phase-amplitude
1073        coupling between neuronal oscillations of different frequencies. *Journal of*
1074        *Neurophysiology*, *104*(2), 1195–1210. https://doi.org/10.1152/jn.00106.2010

1075    van Veen, B. D., van Drongelen, W., Yuchtman, M., & Suzuki, A. (1997). Localization of brain
1076        electrical activity via linearly constrained minimum variance spatial filtering. *IEEE*
1077        *Transactions on Bio-Medical Engineering*, *44*(9), 867–880

1078    van Wassenhove, V., Grant, K. W., & Poeppel, D. (2005). Visual speech speeds up the neural
1079        processing of auditory speech. *Proceedings of the National Academy of Sciences of the*
1080        *United States of America*, *102*(4), 1181–1186

1081    Wagner, P., Malisz, Z., & Kopp, S. (2014). Gesture and speech in interaction: An overview.
1082        *Speech Communication*, *57*, 209–232

1083    Wang, D., Clouter, A., Chen, Q., Shapiro, K. L., & Hanslmayr, S. (2018). Single-trial phase
1084        entrainment of theta oscillations in sensory regions predicts human associative memory
1085        performance. *Journal of Neuroscience*, *38*(28), 6299–6309

1086    Willems, R. M., Ozyürek, A., & Hagoort, P. (2009). Differential roles for left inferior frontal and
1087        superior temporal cortex in multimodal integration of action and
1088        language. *NeuroImage*, *47*(4), 1992–2004.
1089        https://doi.org/10.1016/j.neuroimage.2009.05.066

1090    Zhao, W., Riggs, K., Schindler, I., & Holle, H. (2018). Transcranial Magnetic Stimulation over Left
1091        Inferior Frontal and Posterior Temporal Cortex Disrupts Gesture-Speech Integration. *The*
1092        *Journal of Neuroscience*, *38*(8), 1891–1900. https://doi.org/10.1523/JNEUROSCI.1748-
1093        17.2017

1094    Zoefel, B., Archer-Boyd, A., & Davis, M. H. (2018). Phase Entrainment of Brain Oscillations
1095        Causally Modulates Neural Responses to Intelligible Speech. *Current Biology*, *28*(3), 401-
1096        408.e5. https://doi.org/10.1016/j.cub.2017.11.071

1097

1098

1099

1100

1101

1102
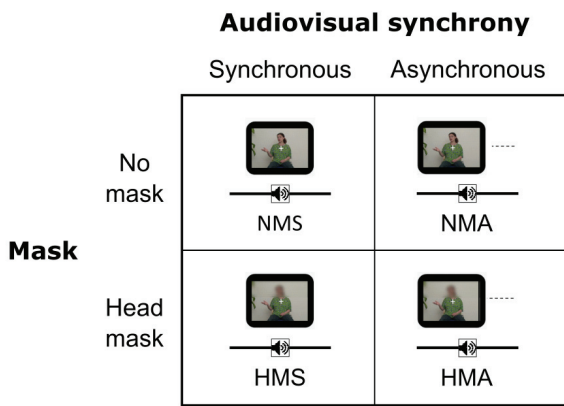
1103

1104

1105

1106

1107     Movie Legends

1108     Movie 1: Example of an audio-visual stimulus presented in the no-mask asynchronous condition
1109     (NMA). In this videoclip, video and audio information were presented in asynchrony (i.e., video
1110     onset led audio onset by 400 milliseconds), and the face of the speaker was fully visible.

1111     Movie 2: Example of an audio-visual stimulus presented in the no-mask synchronous condition
1112     (NMS). In this videoclip, video and audio information were presented in synchrony, and the face
1113     of the speaker was fully visible.

1114     Movie 3: Example of an audio-visual stimulus presented in the head-mask asynchronous
1115     condition (HMA). In this videoclip, video and audio information were presented in asynchrony
1116     (i.e., video onset led audio onset by 400 milliseconds), and the face of the speaker was blurred.

1117     Movie 4: Example of an audio-visual stimulus presented in the no-mask synchronous condition
1118     (HMS). In this videoclip, video and audio information were presented in synchrony, and the face
1119     of the speaker was blurred.

1120     Movie1_still: Example of a frame taken from a videoclip in the no-mask synchronous condition (NMS).
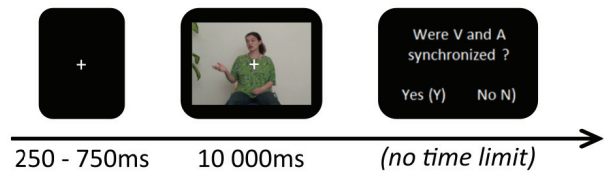
1121    Movie2_still: Example of a frame taken from a videoclip in the no-mask asynchronous condition (NMA).

1122    Movie3_still: Example of a frame taken from a videoclip in the head-mask synchronous condition (HMS).

1123    Movie4_still: Example of a frame taken from a videoclip in the head-mask asynchronous condition
1124    (HMA).

1125

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

**A**

**Audiovisual synchrony**

Synchronous          Asynchronous

Mask

No mask          NMS          NMA

Head mask          HMS          HMA

**B**

250 - 750ms          10 000ms          (no time limit)

Were V and A synchronized ?

Yes (Y)          No N)

**C**

**A**



**B**

**A**

**No-mask**
NMA-NMS

**Head-mask**
HMA-HMS

Frequencies [Hz]

Time [s]

Power difference

**B**

**No-mask**
NMA-NMS

**Head-mask**
HMA-HMS

Power difference

**C**

Normalized power

NMS  NMA        HMS  HMA

**No-mask**    **Head-mask**

**A**



NMA-NMS (ROI)　　NMA-NMS (RONI)

**B**

**A**

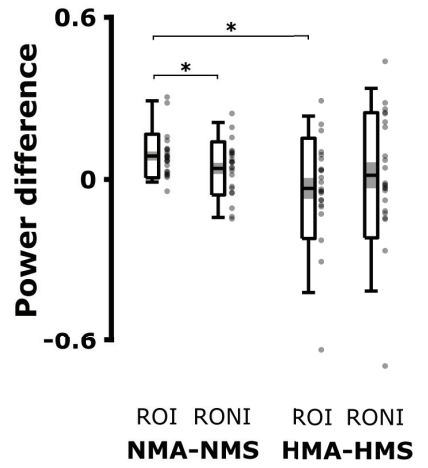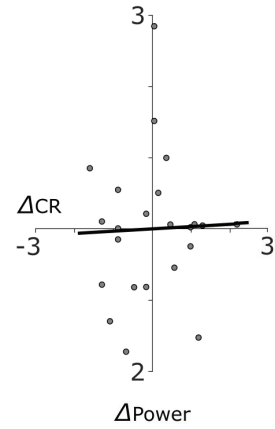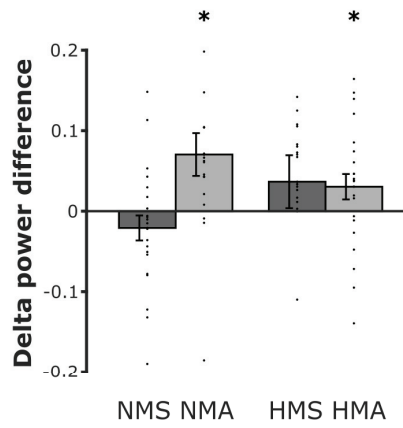**No-mask condition: NMA-NMS**



T values

3

2.1

**B**

**No-mask**
NMA-NMS



ΔCR

ΔPower

3

-3

3

2

**Head-mask**
HMA-HMS



ΔCR

ΔPower

3

-3

3

2

**C**



Delta power difference

0.2

0.1

0

-0.1

-0.2

*

*

NMS NMA    HMS HMA

**D**



Synchronous
Asynchronous

Absolute difference

2.2

1.8

1.4

1

0.6

full    head    body    audio

full    head    body    audio

No-mask    Head-mask



Neural peak frequency [Hz]

No-mask

3

2

1

Head-mask

3

2

1

**Ordered stimuli**

**A**



**B**