Research Articles: Behavioral/Cognitive

# Dissociable neural correlates of multisensory coherence and selective attention

1 **Dissociable neural correlates of multisensory coherence and selective**
2 **attention**

3 **Abbreviated title:** Multisensory coherence and selective attention

4 Fei Peng[1], Jennifer K. Bizley[2], Jan W. Schnupp[1*], Ryszard Auksztulewicz[1,3*]

5 [1] Department of Neuroscience, City University of Hong Kong, Hong Kong, China
6 [2] Ear Institute, University College London, London, United Kingdom
7 [3] Center for Cognitive Neuroscience Berlin, Department of Education and Psychology, Free
8 University Berlin, Germany
9
10 Corresponding author: wschnupp@cityu.edu.hk, ryszard.auksztulewicz@fu-berlin.de

11 43 pages and 6 figures

12 203 words in abstract, 641 words in introduction and 1751 words in discussion.

13 **Conflict of interest statement:**

14 All authors declare that they have no conflicts of interest.

15 **Acknowledgements**

23

24

25

26

27

28

29

30

31

32 **Abstract**

33 Previous work has demonstrated that performance in an auditory selective attention task can be

34 enhanced or impaired, depending on whether a task-irrelevant visual stimulus is temporally

35 coherent with a target auditory stream or with a competing distractor. However, it remains unclear

36 how audiovisual (AV) temporal coherence and auditory selective attention interact at the

37 neurophysiological level. Here, we measured neural activity using electroencephalography (EEG)

38 while human participants (men and women) performed an auditory selective attention task,

39 detecting deviants in a target audio stream. The amplitude envelope of the two competing auditory

40 streams changed independently, while the radius of a visual disc was manipulated to control the

41 audiovisual coherence. Analysis of the neural responses to the sound envelope demonstrated that

42 auditory responses were enhanced independently of the attentional condition: both target and

43 masker stream responses were enhanced when temporally coherent with the visual stimulus. In

44 contrast, attention enhanced the event-related response (ERP) evoked by the transient deviants,

45 independently of AV coherence. Finally, in an exploratory analysis, we identified a spatiotemporal

46 component of ERP, in which temporal coherence enhanced the deviant-evoked responses only in

47 the unattended stream. These results provide evidence for dissociable neural signatures of bottom-

48 up (coherence) and top-down (attention) effects in AV object formation.

49

50 **Keywords**: temporal coherence, selective auditory attention, audio-visual binding, object

51 formation

52

53 **Significance Statement**

54   Temporal coherence between auditory stimuli and task-irrelevant visual stimuli can enhance

55   behavioral performance in auditory selective attention tasks. However, how audiovisual temporal

56   coherence and attention interact at the neural level has not been  established. Here, we measured

57   EEG during a behavioral task designed to independently manipulate AV coherence and auditory

58   selective attention. While some auditory features (sound envelope) could be coherent with visual

59   stimuli, other features (timbre) were independent of visual stimuli. We find that audiovisual

60   integration can be observed independently of attention for sound envelopes temporally coherent

61   with visual stimuli, while the neural responses to unexpected timbre changes are most strongly

62   modulated by attention. Our results provide evidence for dissociable neural mechanisms of

63   bottom-up (coherence) and top-down (attention) effects on AV object formation.

64   **Introduction (650 words)**

65   In many real world sound environments, sounds originate from multiple sources – the auditory

66   system needs to appropriately segregate and group sound features to efficiently process the entire

67   scene (Maddox & Shinn-Cunningham, 2012; Middlebrooks et al., 2017; Shamma et al., 2011).

68   Several psychoacoustic studies have demonstrated that visual cues which are temporally coherent

69   with sounds can modulate auditory processing. For example, a synchronous, task-irrelevant light

70   flash improves the detection of weak auditory signals (Lovelace et al., 2003). Similarly, task-

71   irrelevant visual stimuli which are temporally coherent with a speech envelope enhance speech

72   intelligibility in background babble noise (Yuan et al., 2020). Furthermore, performance in an

73   auditory selective attention task can be enhanced or impaired, depending on whether the task-

74   irrelevant visual stimulus is temporally coherent with a target sound stream or a competing masker

75    stream (Maddox et al., 2015). However, the neural mechanisms mediating the interactions between

76    temporal coherence and selective attention in facilitating AV integration remain unknown.

77    Several previous studies have identified potential neural correlates of attentional modulation of

78    AV integration. For example, a study using simple tone pips and visual gratings demonstrated that

79    ERPs related to multisensory integration were amplified by selective attention (Talsma &

80    Woldorff, 2005). When both visual and auditory stimuli were attended, the ERP peak amplitude

81    showed superadditive AV effects; however, subadditive effects were observed for unattended

82    stimuli (Talsma et al., 2007). Some EEG and MEG studies have employed the analysis of "neural

83    envelope-tracking responses" to speech, by modeling the relationship between neural activity and

84    the auditory envelope (Crosse et al., 2015; Golumbic et al., 2013), and have found that congruent

85    audio-visual speech enhances the envelope tracking response relative to auditory speech alone or

86    the linear summation of auditory and visual speech. Other studies have used auditory selective

87    attention tasks to show that attention is necessary for AV speech integration. For example, Morís

88    Fernández et al. (2015) measured fMRI data and showed that multisensory integration occurred

89    almost exclusively only when the congruent AV speech was attended. However, Ahmed et al.,

90    (2021) measured EEG and found some evidence for early AV integration in the unattended stream,

91    consistent with the idea that distinct audiovisual computations emerge at different processing

92    stages (Kayser & Shams, 2015; Talsma et al., 2007; Talsma & Woldorff, 2005; Zumer et al., 2020).

93    One potential difficulty with interpreting findings from AV speech processing is that it can be hard

94    to know the extent to which they generalize to other continuous AV stimuli, given that speech

95    processing can be heavily influenced by linguistic knowledge and expectations. Thus these speech

96    specific studies might not represent more general mechanisms of visual influences on auditory

97    processing. Consistent with AV integration occurring independently of attention for non-speech

98    stimuli, neural correlates of AV integration were observed in single neurons in the auditory cortex

99    of passively exposed ferrets. This included enhancement of the neural representation of the

100   temporally coherent features (i.e., envelope), but also of the other (i.e., timbre) sound features

101   (Atilgan et al., 2018). Together, from these findings it remains unclear whether such bottom-up

102   effects modulate the cortical representation of auditory streams independently of attentional top-

103   down enhancement, or whether these effects are synergistic.

104   Here, we use EEG to investigate the electrophysiological correlates of AV temporal coherence and

105   auditory selective attention on sound processing in an auditory selective attention task. Listeners

106   were required to detect short timbre deviants in an attended audio stream, while a visual stimulus

107   was paired with either the target, masker or neither sound through coherent size/amplitude

108   fluctuations. First, we focused our analysis on how AV coherence and attention affected the neural

109   signatures of continuous stream processing, as manifest in the envelope-tracking response. Second,

110   we focused on the transient auditory deviants, whose timing was independent of the features of the

111   visual stream, and compared deviant-evoked ERPs between conditions. Our goal was to test the

112   hypothesis that attention and audiovisual integration operate independently.

## Materials and Methods

### Participants

115   Twenty volunteers were recruited for this experiment (median ± standard deviation (SD) age, 22

116   ± 2 years; 12 males; 19 right-handed). All participants were healthy, had self-reported normal

117   hearing and normal or corrected-to-normal vision. Prior to the experiment, each participant gave

118   written informed consent. All procedures were approved by the Human Subjects Ethics Sub-

119   Committee of the City University of Hong Kong.

4

120 **Stimuli**

121 We adapted the behavioral paradigm from previous psychoacoustics studies (Atilgan & Bizley,

122 2021; Maddox et al., 2015). Stimuli included two simultaneously presented auditory streams and

123 one visual stream. One auditory stream was meant to be attended, and will be referred to as the

124 target sound (At), while the other one was meant to be unattended, and will be referred to as the

125 masker stream (Am). Finally, stimulation included a concurrently presented visual stream (V)

126 which comprised a radius-modulated disc. Auditory streams were independently amplitude-

127 modulated and the modulation of the visual disc could be temporally coherent either with the

128 amplitude of the target stream (AtAmVt), the masker stream (AtAmVm), or independent of both

129 (AtAmVi) (Figure 1A).

130 The envelopes below 7 Hz were generated using the same methods as in Maddox et al. (2015).

131 Briefly, frequency domain synthesis was used to generate the envelopes. In the frequency domain,

132 amplitudes of frequency bins between 0-7 Hz were set to one and, for other frequency bins, to zero,

133 The non-zero bins were given a random phase from a uniform distribution between 0 and $2\pi$, at

134 an audio sampling rate of 24414 Hz. To maintain Hermitian symmetry, the corresponding

135 frequency bins across Nyquist frequency were set to the respective complex conjugates. The

136 inverse Fourier transform was calculated to create the time domain envelope. Three envelopes

137 of each trial were computed using the same method, and they were orthogonalized using the

138 Gram-Schmidt procedure. Visual envelopes were generated by downsampling the auditory

139 envelope at the monitor frame-rate of 60 Hz, where the disc radius of the first frame was

140 corresponding to the first auditory sample. Each auditory stream consisted of one continuous

141 amplitude modulated synthetic vowel, either /u/ or /a/, which were generated by filtering a click

142 train at four "formant" frequencies (F1-F4). The fundamental frequency (F0) of vowel /u/ was 175

5

143 Hz, and the formant peaks were 460, 1105, 2975, 4263 Hz, while the F0 of vowel /a/ was 195 Hz,

144 and the formant peaks were 936, 1551, 2975, 4263 Hz. Auditory deviants were embedded in the

145 auditory streams by temporarily changing the timbre of the vowel. The deviant in vowel /u/

146 transitioned (in F1/F2 space) towards the vowel /ɛ/, with the maximum timbre change resulting in

147 formant peaks at 525, 1334, 2975, 4263 Hz, while the deviant in vowel /a/ transitioned towards /i/

148 with formant peaks at 860, 1725, 2975, 4263 Hz. The duration of each deviant was 200 ms, which

149 included a linear change of the formants towards the deviant for 100 ms and then back for 100 ms.

150 The visual stimulus was a grey disc surrounded by a white ring presented at the center of the black

151 screen. The radius of the visual stimulus was modulated by the visual envelope, such that the disc

152 subtended between 1° and 2.5°, and the white ring extended 0.125° beyond the grey disc.

153 Each trial lasted 14 s and comprised three streams. A target audio stream and the visual stream

154 were each 14 s in duration while the masker stream, although also generated to be 14 s in duration,

155 was silenced for the first second. The initial 1 s, during which only the target stream was audible,

156 provided the cue for the listener which was the to-be-attended target stream. Auditory deviants

157 could occur at any time during a window beginning 2 s after the onset of the target audio stream

158 and ending 1 s before the end of the trial, subject to the constraint that the minimum interval

159 between auditory deviants was 1.2 s. On average each stream contained 2 deviants (range 1-3

160 across trials). Unlike Maddox et al. (2015), the visual stream did not contain any colour deviants.

161 **Procedure**

162 Participants were seated in a sound-attenuated room. Auditory stimuli were presented binaurally

163 via earphones (ER-3, Etymotic Research, Elk Grove Village, IL, USA), using an RZ6 signal

164 processor at a sampling rate of 24414 Hz (Tucker-Davis Technologies, Alachua, FL, USA). The

165 sound level was calibrated at 65 dB SPL. Visual stimuli were presented on a 24-inch computer

166　monitor using the Psychophysics Toolbox for MATLAB. Participants were asked to pay attention

167　to the target auditory stream and to detect the embedded auditory deviants by pressing a keyboard

168　button. They were instructed to refrain from pressing buttons in response to any events in the

169　masker stream.

170　Before the actual task, all participants completed a training session to verify that they were able to

171　detect the auditory deviants. The training session included four blocks, and each block included 9

172　trials. The feedback of performance was given after each block, and all participants showed they

173　could perform the experiment (d' > 0.8) in at least one block of four.

174　Participants were instructed to minimize eye blinks and body movements during the EEG

175　recording. Continuous EEG signals were collected using an ANT Neuro EEGo Sports amplifier

176　from 64 scalp channels at a sampling rate of 1024 Hz. The EEG signals were grounded at the

177　nasion and referenced to the Cpz electrode. Each participant completed 12 blocks of the task, with

178　18 trials (6 trials x 3 conditions) in each block. Trials belonging to different conditions were

179　presented in a randomly interleaved order. In total, each participant completed 216 trials (72 trials

180　x 3 conditions). Feedback on behavioral performance was provided after each block. Triggers

181　corresponding to trial and deviant onset were recorded along with the EEG signal.

182　**Behavioral data analysis**

183　A 'hit' was defined as the response to the deviant in the target auditory stream within 1 s following

184　the onset of the deviant, and a 'false alarm' was defined as the response to a deviant that occurred

185　in the masker stream. To study how visual coherence affects auditory deviant detection, we

186　conducted a one-way repeated measures ANOVA on the sensitivity measure d' with a within-

187　subjects factor of AV condition (visual coherent with the target, AtAmVt, visual coherent with the

188　masker, AtAmVm, and independent visual AtAmVi).

7

189 **EEG signal pre-processing**

190 EEG signals were pre-processed using the SPM12 Toolbox (Wellcome Trust Centre for

191 Neuroimaging, University College London) for MATLAB. Continuous data were downsampled

192 to 500 Hz, high-pass filtered at a cut-off frequency of 0.01 Hz, notch-filtered between 48 Hz and

193 52 Hz, and then low-pass filtered at 90 Hz. All filters were fifth-order zero-phase Butterworth.

194 Eyeblink artifacts were removed by the use of the principal component analysis (PCA) based on a

195 "preselection" spatial filtering technique described by Ille et al., (2002). Specifically, eyeblink

196 artifacts were detected by computing the principal components of the signal in the channel Fpz,

197 and removed by subtracting the first two spatiotemporal components associated with each eyeblink

198 from all channels (Ille et al., 2002). The EEG data were then re-referenced to the average of all

199 channels. The preprocessed data were further analyzed in two ways: For the response to the sound

200 amplitude envelope, the pre-processed data were bandpass filtered between 0.3 and 30 Hz (Crosse

201 et al., 2015), downsampled to 64 Hz, and subjected to a calculation of the TRF, or used for stimulus

202 reconstruction (see below). For the deviant evoked response analysis, the pre-processed data were

203 epoched from -100 ms to 500 ms relative to deviant onset. Epoched EEG signals were then

204 denoised using the "Dynamic Separation of Sources" (DSS) algorithm (de Cheveigné & Simon,

205 2008), which is commonly used to maximize reproducibility of stimulus-evoked response across

206 trials and maintain the differences across the different stimulus types (here: 2 vowel types × 3

207 experimental conditions). Epoched data were linearly detrended, and the first seven DSS

208 components were preserved and applied to project the data back into sensor space. The SD of the

209 voltage over time was computed for each trial, and we excluded the noisy trials whose SD

210 exceeded the median ± 2SD over trials. Across participants roughly 30 trials were excluded for

211 each participant (the included trials were 829 ± 31 (median ± SD) out of 864 trials). Denoised data

212 were averaged across the good trials.

213 **EEG response to sound amplitude envelopes**

214 **Stimulus reconstruction**

215 To investigate how visual temporal coherence and attention affect multisensory integration, we

216 quantified the accuracy of neural tracking of the sound amplitude envelope. We reconstructed

217 amplitude envelopes of different elements of the AV scene (Crosse et al., 2015) based on the EEG

218 data using a linear model as follows:

219
$$\check{s}(t) = \sum_{n=1}^{64} \sum_{\tau=0}^{500\ ms} r(t + \tau, n) g(\tau, n) \tag{1}$$

220 where $\check{s}(t)$ is the reconstructed envelope; $r(t + \tau, n)$ is the EEG data at channel $n$; and $g$ is the

221 linear decoder representing the linear mapping from the response to stimulus amplitude envelope

222 at time lag $\tau$. The time lag $\tau$ ranged from 0 to 500 ms post-stimulus. The decoder was obtained

223 separately for each condition using ridge regression as follows:

224
$$g = (R^T R + \lambda I)^{-1} R^T s \tag{2}$$

225 where $R$ is the lagged time series of the EEG response, $\lambda$ is the ridge parameter, $I$ is the

226 regularization term, and $s$ is the sound amplitude envelope. The decoder is a multivariate impulse

227 response function computed from all channels and all time-lags simultaneously. Decoders

228 corresponding to the AV, A-only, and V-only streams were generated separately as follows:

229
$$g_{AtVt}(\tau, n) = (R_{AtAmVt}^T R_{AtAmVt} + \lambda I)^{-1} R_{AtAmVt}^T s_{AtVt} \tag{3}$$

230
$$g_{AmVm}(\tau, n) = (R_{AtAmVm}^T R_{AtAmVm} + \lambda I)^{-1} R_{AtAmVm}^T s_{AmVm} \tag{4}$$

231
$$g_{At}(\tau, n) = (R_{AtAmVi}^T R_{AtAmVi} + \lambda I)^{-1} R_{AtAmVi}^T s_{At} \tag{5}$$

232
$$g_{Am}(\tau, n) = (R_{AtAmVi}^T R_{AtAmVi} + \lambda I)^{-1} R_{AtAmVi}^T s_{Am} \tag{6}$$

9

233
$$g_{Vi}(\tau, n) = (R_{AtAmVi}^T R_{AtAmVi} + \lambda I)^{-1} R_{AtAmVi}^T s_{Vi} \tag{7}$$

234 Since in the condition AtAmVi, the envelope of At, Am, and Vi are independent of each other, we

235 could obtain the decoder of the envelopes of the auditory target only, auditory masker only, and

236 visual only, respectively. To obtain the decoder for each condition, we used leave-one-trial-out

237 cross-validation to select the $\lambda$ value (from the set of $10^{-6}$, $10^{-5}$, …, $10^5$, $10^6$) for which the

238 correlation between $\check{s}(t)$ and $s(t)$ is maximized. To assess the effect of AV integration, we

239 reconstructed the sound envelopes (both target and masker sound) using the integration AV

240 decoder and the algebraic sum of the A and V decoder (A+V), separately, based on the following

241 formulas:

242
$$\widetilde{s_{At\_(AV)}}(t) = \sum_{n=1}^{64} \sum_{\tau=0}^{500\ ms} r_{AtAmVt}(t+\tau, n) g_{AtVt}(\tau, n) \tag{8}$$

243
$$\widetilde{s_{At\_(A+V)}}(t) = \sum_{n=1}^{64} \sum_{\tau=0}^{500\ ms} r_{AtAmVt}(t+\tau, n)(g_{At}(\tau, n) + g_{Vi}(\tau, n)) \tag{9}$$

244
$$\widetilde{s_{Am\_(AV)}}(t) = \sum_{n=1}^{64} \sum_{\tau=0}^{500\ ms} r_{AtAmVm}(t+\tau, n) g_{AmVm}(\tau, n) \tag{10}$$

245
$$\widetilde{s_{Am\_(A+V)}}(t) = \sum_{n=1}^{64} \sum_{\tau=0}^{500\ ms} r_{AtAmVm}(t+\tau, n)(g_{Am}(\tau, n) + g_{Vm}(\tau, n)) \tag{11}$$

246 The reconstruction accuracy (r) was defined as the Pearson correlation coefficient between the

247 actual stimulus envelope and the estimated envelope.

248 Based on our main research question - namely whether the effects of attention and coherence are

249 independent or synergistic - the possible scenarios of combining the effects of coherence and

250 attention were considered in the context of two main models of AV coherence: an integration

251 model and a summation model (Figure 1B). To test whether the reconstruction accuracy using

252 either the AV decoder ("integration model") and/or A+V decoder ("summation model") was

253 significantly larger than chance, we conducted a nonparametric permutation test. The null

254 distribution of 1000 Pearson's r values was created for each subject by calculating the correlation

255 between randomly shuffled response trials of estimated sound envelopes and actual sound

256 envelopes. We estimated sound envelopes using each decoder separately, and generated the null

257 distribution for each condition.

258 To test for the interaction of attention and AV integration, we computed a repeated-measures

259 ANOVA on reconstruction accuracy with two main within-subjects factors, attention (target vs.

260 masker) and integration decoder ("integration model": AV vs. "summation model": A+V).

261 **Temporal response function (TRF) estimation**

262 To investigate how the visual temporal coherence and attention affect AV integration across the

263 EEG channels, we estimated the linear temporal response function (TRF) (Crosse, Di Liberto,

264 Bednar, et al., 2016) which links the EEG response at each channel and the sound envelope. The

265 TRF is the linear filter that best describes the brain's transformation of the sound envelope to the

266 continuous neural response at each EEG channel location (Haufe et al., 2014). TRFs were

267 estimated separately for each experimental condition (AtAmVt, AtAmVm, AtAmVi) as follows:

268
$$r_{AtAmVt}(t,n) = \sum_\tau s_{AtVt}(t-\tau)\, TRF_{AtVt}(\tau,n) + \sum_\tau s_{A2}(t-\tau)\, TRF_{Am\prime}(\tau,n) + \varepsilon(t,n) \qquad (12)$$

269
$$r_{AtAmVm}(t,n) = \sum_\tau s_{At}(t-\tau)\, TRF_{At\prime}(\tau,n) + \sum_\tau s_{AmVm}(t-\tau)\, TRF_{AmVm}(\tau,n) + \varepsilon(t,n) \qquad (13)$$

270
$$r_{AtAmVi}(t,n) = \sum_\tau s_{At}(t-\tau)\, TRF_{At}(\tau,n) + \sum_\tau s_{Am}(t-\tau)\, TRF_{Am}(\tau,n) + \sum_\tau s_{Vi}(t-\tau)\, TRF_{Vi}(\tau,n) + \varepsilon(t,n) \quad (14)$$

271 where $r_{AtAmVt}$, $r_{AtAmVm}$, and $r_{AtAmVi}$ are the EEG response in each of the 3 conditions respectively;

272 $t$ is time, $n$ is the index of the EEG channel under consideration; $s_{At}$, $s_{Am}$, and $s_{Vi}$ are the stimulus

273 envelopes of At, Am, and Vi, respectively; $\tau$ represents the convolution time lag (-100 ms to 500

274 ms), and $\varepsilon(t,n)$ is the residual "error", that is, the part of the EEG recording not explained by the

275 TRF model. We use the term $TRF_{At1}$ to describe the TRF in the AtAmVm condition, and $TRF_{Am1}$

276 in the AtAmVt condition, to differentiate them from the $TRF_{At}$ and $TRF_{Am}$ estimated from the

277 AtAmVi condition, this being the only condition in which all three streams were fully independent

278 (Equations 13-15). The TRF for each condition was calculated at time lags from -100 ms to 500

279 ms relative to the stimulus as follows:

$$TRF_{AtVt} = (S_{AtVt}^T S_{AtVt} + \lambda I)^{-1} S_{AtVt}^T r_{AtAmVt} \tag{15}$$

$$TRF_{AmVm} = (S_{AmVm}^T S_{AmVm} + \lambda I)^{-1} S_{AmVm}^T r_{AtAmVm} \tag{16}$$

$$TRF_{At} = (S_{At}^T S_{At} + \lambda I)^{-1} S_{At}^T r_{AtAmVi} \tag{17}$$

$$TRF_{Am} = (S_{Am}^T S_{Am} + \lambda I)^{-1} S_{Am}^T r_{AtAmVi} \tag{18}$$

$$TRF_{Vi} = (S_{Vi}^T S_{Vi} + \lambda I)^{-1} S_{Vi}^T r_{AtAmVi} \tag{19}$$

285 where $S$ is the lagged time series of the stimulus envelope; $I$ is the regularization term used to

286 prevent overfitting; and $\lambda$ is the ridge parameter. $TRF_{AtVt}$, $TRF_{AmVm}$, $TRF_{At}$, $TRF_{Am}$, and $TRF_{Vi}$

287 were fitted separately for each condition using the MATLAB toolbox adapted from a previous

288 study by Crosse et al. (2016). The TRF of each channel was estimated using leave-one-out cross-

289 validation. The best $\lambda$ (in the range of $2^{10}$, $2^{11}$, …, $2^{21}$) was selected based on the maximum

290 correlation coefficient between the predicted response with the actual neural response for each

291 channel. The EEG signal of each trial (13 s long) was used to estimate the TRF, modeling the

292 neural response to the simultaneous presentation of both At and Am.

293 To test whether AV integration is affected by attention, we compared the TRF amplitude between

294 the temporally coherent and independent conditions across EEG channels and time points. Single-

295 participant TRF data were converted into three-dimensional images (2D: spatial topography, 1D:

296 time) and entered into a repeated-measures ANOVA with two within-subjects factors: attention

297 (attended: $TRF_{AtVt}$ and $TRF_{At} + TRF_{Vi}$, unattended: $TRF_{AmVm}$ and $TRF_{Am} + TRF_{Vi}$) and

298 integration (integration model: $TRF_{AtVt}$ and $TRF_{AmVm}$, linear summation model: $TRF_{At} + TRF_{Vi}$

299 and $TRF_{Am} + TRF_{Vi}$). The two-way repeated-measures ANOVA was implemented as a GLM in

300 SPM12. The resulting statistical parametric maps, representing the main and interaction effects,

301  were thresholded at p < 0.05 (two tailed) and corrected for multiple comparisons across

302  spatiotemporal voxels at a family-wise error (FWE)-corrected p = 0.05 (cluster-level) under

303  random field theory assumptions (Kilner et al., 2005).

**Auditory Deviant-evoked ERP**

305  To assess how attention and visual coherence affect deviant-evoked activity, the EEG data were

306  first subject to a traditional channel-by-channel mass-univariate analysis. Epoched data were

307  averaged over trials, separately for the deviants in At and Am and for each visual condition (Vt,

308  Vm, Vi). Single-subject ERP data were converted into three-dimensional images (two spatial

309  dimensions and one temporal dimension) and entered into a repeated-measures ANOVA with two

310  within-subjects factors: attention (attended: deviant in the At stream, unattended: deviant in the

311  Am stream) and visual coherence (coherent with the sound containing deviants: deviants in AtVt

312  and AmVm; visual condition independent of the sound: AtVm and AmVt). The two-way repeated-

313  measures ANOVA was implemented as a GLM in SPM12. The resulting statistical parametric

314  maps, representing the main and interaction effects, were thresholded at p < 0.05 (two-tailed) and

315  corrected for multiple comparisons across spatiotemporal voxels at a family-wise error (FWE)-

316  corrected p = 0.05 (cluster-level).

317  In a follow-up attempt to isolate dissociable neural signatures of attention and visual coherence,

318  we concatenated the ERP data across participants and used PCA to reduce the EEG data

319  dimensionality and obtain spatial principal components (PCs, representing the weight of channel

320  topographies) and temporal principal components (representing voltage time-series). The EEG

321  data were concatenated across participants before being subjected to PCA, in order to obtain the

322  same PCs across participants. The PCs quantified independent contributions to whole-scalp data,

323  such that the sensitivity to those isolated components increased (relative to original data,

324   containing a mixture of components). The first four PCs (explaining 80% of the original variance

325   across participants) were used to extract single-participant ERP components for further analysis.

326   Each PC was then converted into one-dimensional images (time) and subject to statistical inference

327   using repeated-measures ANOVAs, as above. Significance thresholds were kept identical to the

328   traditional univariate analysis, but correction for multiple comparisons was implemented across

329   time points (rather than spatiotemporal voxels).

330   **Correlating timbre deviant evoked ERP magnitude with behavioral performance**

331   Since the behavioral task was to detect deviants in the target auditory stream, we extracted the

332   EEG responses to deviants in At and measured the peak to peak amplitude of the PCs of ERP

333   identified above. We then calculated the Pearson correlation coefficients between the behavioral

334   mean d' and the mean PC amplitude over conditions (AtAmVt, AtAmVm, and AtAmVi). To

335   reduce the number of comparisons, we limited our correlation analyses to those ERP components

336   and factors which showed significant effects. Specifically, for the $1^{st}$ PC for which we have

337   identified the significant main effect of attention (see Results), the negative and positive peaks

338   were measured between 100 to 200 ms, and 220 to 300 ms. respectively. For the $3^{rd}$ PC for which

339   we have identified significant main and interaction effects of attention and coherence (see Results),

340   the positive and negative peaks were measured between 50 to 160 ms, and 220 to 400 ms,

341   respectively. Prior to calculating the correlations, we fitted the behavioral performance d' with the

342   PC peak-to-peak amplitude using a linear regression model, and detected the outliers in each

343   condition using Cook's distance (threshold: 3 means of Cook's distance).

344

# Results

**Behavioral results**

First, we investigated whether behavioral performance was stable over time, which would warrant pooling data from all blocks. To this end, we calculated the single-participant hit rate separately for each of the 12 blocks, and fitted the data using a linear regressor representing the block number. The resulting regression coefficient (slope) was not statistically different from zero across participants (one-sample t-test, $t = 1.11$, $p = 0.28$), suggesting that there were neither significant learning nor fatigue effects during the experiment.

To investigate the effect of the visual temporal coherence on behavioral performance, we performed one-way repeated measures ANOVAs on d' ($F = 0.15$, $p = 0.85$) (Figure 1C), hit rates ($F = 0.42$, $p = 0.66$) and false alarm rates ($F = 2.12$, $p = 0.13$). The hit rates were 69% ± 2.7%, 70% ± 2.6% and 70% ± 2.9% (mean ± SEM), and the mean false alarm rates were 4% ± 0.6%, 5% ± 0.9%, and 5% ± 1% for the three conditions (AtAmVt, AtAmVm, AtAmVi), respectively. No significant effect of visual coherence on deviant detection was observed, likely due to large variability and heterogeneity of response patterns across participants. For instance, while some participants showed behavioral benefits of visual coherence (e.g., larger d' in AtAmVt condition than AtAmVm), others showed the opposite effects (Figure 1C). Two previous studies using similar stimulus paradigms (Atilgan & Bizley, 2021; Maddox et al., 2015) reported enhanced task performance when the target stream and visual stimulus were temporally coherent. Our failure to replicate these data may be attributable to small but perhaps important differences in the details of the experimental paradigms, especially the manipulation of visual attention (see Discussion). Furthermore, our behavioral results are consistent with the general framework of the possible effects of attention and coherence (Figure 1B), in which the relative contribution of the integration

15

368  term might be small compared to the summation term. However, the aims of this study were to

369  identify effects of auditory selective attention and AV coherence on physiological measures of

370  neural stimulus representations, and the timbre deviants primarily served as a device for

371  controlling and monitoring our participants' attention. The relatively high hit rates and low false

372  alarm rates indicate that the deviants had fulfilled that purpose.

373  **Stimulus reconstruction reveals temporal coherence mediated audiovisual**

374  **integration**

375  To investigate the occurrence of AV integration at both attended and unattended conditions, we

376  reconstructed an estimation of the sound envelope from the recorded EEG waveforms. We used

377  the condition in which the visual stimulus was independent of both auditory streams to estimate

378  unimodal reconstructions for the target auditory stream (At), the masker stream (Am) and the

379  visual stream (Vi) (Figure 2B). From this condition we could independently estimate unisensory

380  response elements, without introducing some of the confounds inherent in comparing activity

381  across multisensory and unisensory trials, where prestimulus expectation and attention may differ

382  (Mishra et al., 2007, Rohe et al. 2019). We first confirmed that the unimodal reconstructions for

383  all conditions were significantly better than the chance estimated using a permutation test. From

384  the unisensory reconstructions,  we estimated the response to stimuli in which the visual stimulus

385  was coherent with one or the other stream by linear summation. This linear summation model was

386  compared to an integration model in which audiovisual envelopes were reconstructed based on the

387  responses to conditions in which the visual stimulus was temporally coherent with one or the other

388  stream (i.e., AtVt and AmVm) (Figure 2A). Testing for the interaction of attention and integration

389  in a two-way repeated measures ANOVA, we only found that the main effect of integration was

390  significant (F = 491.8, p <0.001). In post-hoc comparisons, we observed that the average

391   reconstruction accuracy of the AV decoder was significantly higher than that of the A+V decoder

392   for both the target stream (Figure 2C, Wilcoxon signed-rank test p < 0.001) and the masker stream

393   (Figure 2D, Wilcoxon signed-rank test p < 0.001), consistent with AV integration occurring

394   independently of attention.

395   **Forward models highlight attentional modulation of auditory responses**

396   We next asked how temporal coherence and attention affect AV integration across the different

397   EEG channels by estimating temporal response functions (TRFs) of each channel. While stimulus

398   reconstruction predicts the accuracy of cortical tracking of the amplitude envelope by using

399   multichannel EEG response (and may therefore be dominated by visual responses), TRFs reflect

400   the linear transformation of the sound envelope to the neural responses at each EEG channel. We

401   first explored whether we could observe similar evidence of audiovisual integration from the TRF

402   estimations as we did with the stimulus reconstruction. We estimated unisensory TRFs for the

403   auditory target stream ($TRF_{At}$), the auditory masker stream ($TRF_{Am}$), and the visual stimulus

404   ($TRF_{Vi}$), separately, from the response in the condition AtAmVi, in which temporal envelopes of

405   all three streams were independent. We then estimated the $TRF_{AtVt}$ and $TRF_{AmVm}$ using the

406   responses in the condition AtAmVt and AtAmVm, respectively.

407   To investigate how the cortical representation of amplitude envelopes was influenced by attention

408   and AV integration, we used a two-way repeated measures ANOVA to assess the influence of AV

409   integration and attention on the TRF amplitudes across all EEG channels. We observed a

410   significant main effect of attention (Figure 3A anterior cluster, 78 ms to 219 ms, $F_{max} = 26.68$, $Z_{max}$

411   $= 4.51$, $p_{FWE} < 0.001$; Figure 3B anterior and central cluster, 297 to 391 ms, $F_{max} = 27.98$, $Z_{max} =$

412   $4.61$, $p_{FWE} = 0.008$) and integration (Figure 3C anterior cluster, 219 to 250 ms, $F_{max} = 14.12$, $Z_{max}$

413   $= 3.35$, $p_{FWE} = 0.009$).

17

414  In summary, we observed evidence that AV integration occurred both in the target and masker

415  auditory stream when measures of stimulus reconstruction accuracy were used to analyse the

416  neural responses to the sound envelopes. Analysis of TRFs amplitude across all EEG channels

417  showed that attention modulated the magnitude of the TRF. AV integration was observed for the

418  masker stream in central and frontal channels. The attention effect was observed for a subset of

419  channels in the TRF analysis but not in the stimulus reconstruction, which utilized the responses

420  across all channels. The other possible reason that attentional effects were observed with the TRF

421  and not the stimulus reconstruction, is that the latter might be dominated by the responses to visual

422  stimulus (Figure 2B). Taken together, our results suggest audio-visual integration occurs

423  automatically, prior to attentional modulation.

424  **Effects of audiovisual temporal coherence and selective attention on deviant-evoked**

425  **responses**

426  The analysis so far has focused on the neural responses to the amplitude envelopes of the

427  audiovisual scene, and has revealed evidence for both attentional modulation of acoustic responses,

428  and AV integration of temporally coherent cross-modal sources. Since, in the temporally coherent

429  conditions, the visual and auditory streams convey redundant information, this integration falls

430  short of reaching the stricter definition of binding proposed by Bizley et al., (2016) which requires

431  an enhancement of independent features that are not those which link the stimuli across modalities.

432  Here, the presence or timing of the auditory timbre deviants that listeners detected in the selective

433  attention task are not predicted by the amplitude changes of the audio or visual envelopes, and they

434  thus provide a substrate with which to explore binding.

435  To investigate how AV temporal coherence and attention affect the deviant-evoked responses, we

436  compared the ERPs evoked by deviants embedded in At and Am streams, and, in order to look for

437    evidence of binding, asked how audiovisual temporal coherence modulated these responses

438    (Figure 4). As shown in scalp topographies which visualize the response change over time for each

439    condition (Figure 4A), the deviant-evoked response in the target stream was clearly stronger than

440    that in the masker stream.

441    Accordingly, in a traditional channel-by-channel mass-univariate analysis, correcting for multiple

442    comparisons across all channels and time points, we observed a significant main effect of attention

443    (anterior cluster, 196 to 302 ms, $F_{max}$ = 16.87, $Z_{max}$ = 3.65, $p_{FWE}$ < 0.001; posterior cluster, 210 to

444    320 ms, $F_{max}$ = 21.32, $Z_{max}$ = 4.08, $p_{FWE}$ < 0.001) and a significant interaction effect of attention

445    and temporal coherence (anterior cluster, 62 to 146 ms, $F_{max}$ = 11.26, $Z_{max}$ = 2.99, $p_{FWE}$ = 0.036,

446    posterior cluster, 58 to 146 ms, $F_{max}$ = 16.92, $Z_{max}$ = 3.66, $p_{FWE}$ = 0.03 ). No main effect of temporal

447    coherence was observed.

448    Significant post-hoc comparisons between conditions were consistent with the main effect of

449    attention: for both temporally coherent and temporally independent streams, the deviant response

450    in the target always exceeded that of the masker. The amplitude of the ERP evoked by timbre

451    deviants presented in the target stream (AtVm) was significantly larger than that in the masker

452    stream (AmVt) in two clusters: negative peak enhancement was observed over anterior channels

453    (Figure 4B the first row, 210-300 ms after deviant onset, $p_{FWE}$ < 0.001, $T_{max}$ = 4.23), and positive

454    peak enhancement over posterior channels (Figure 4B the second row, 212-302 ms after deviant

455    onset, $p_{FWE}$ < 0.001, $T_{max}$ = 4.68). In the AV coherent stream, we observed that ERP amplitude

456    evoked by the timbre deviants in the attended coherent stream (AtVt) was significantly stronger

457    than in the unattended coherent stream (AmVm) in two clusters: one over the central and frontal

458    channels between time lag 236 to 310 ms (Figure 4C the first row, cluster level $p_{FWE}$ < 0.001, $T_{max}$

459    = 3.8), and one over posterior channels between time lag 234 to 350 ms (Figure 4C the second row,

19

460    cluster level $p_{FWE} = 0.007$, $T_{max} = 4.18$).

461    Post-hoc comparisons also allowed us to examine the interaction between temporal coherence and

462    attentional condition. We observed that the amplitude of ERP evoked by deviants in the masker

463    stream was significantly smaller when this was accompanied by a temporally coherent visual

464    stimulus (Figure 4E). The deviant induced ERP was smaller in the AmVm condition than in the

465    AmVt condition in two clusters: one over the central and frontal channels between time lag 74 to

466    186 ms (Figure 4E the first row, cluster level $p_{FWE} = 0.011$, $T_{max} = 4.38$), and one over left temporal

467    and posterior channels between time lag 94 to 180 ms (Figure 4E the second row, cluster level

468    $p_{FWE} = 0.005$, $T_{max} = 3.72$). In contrast, audiovisual temporal coherence did not influence the size

469    of the deviant response in the target stream (Figure 4D).

470    From the mass-univariate ERP data analysis (i.e., when analysing all channels and correcting for

471    multiple comparisons across channels and time points), attention was the main modulator of the

472    size of the deviant response, with temporal coherence only influencing the deviant responses in

473    the masker stream. In a follow-up exploratory analysis, we investigated whether effects of visual

474    coherence, as well as attention, can be identified when EEG channels are grouped into principal

475    spatiotemporal components explaining different sources of variance. To this end, we performed a

476    principal component analysis to extract the spatiotemporal components of the ERP, and performed

477    separate two-way repeated measures of ANOVAs with two main factors: attention (attended and

478    unattended) and visual coherence (coherent and incoherent), on the first four principal components

479    (PCs) in the time domain. These four PCs together explained 80% of the original variance. The

480    analysis of the 1st PC (Figure 5A, explaining 67% of the original variance) only showed a main

481    effect of attention (time lag between 208 to 284 ms, $F_{max} = 32.53$, $Z_{max} = 4.92$, $p_{FWE} < 0.001$). No

482    main or interaction effects were found to be significant for the 2nd and 4th PC (Figure 5B and Figure

483    5D, explaining 6% and 3% of the original variance, respectively). However, the analysis of the 3rd

484    PC (explaining 4% of the original variance) showed a main effect of attention (time lag between

485    8 to 84 ms, $F_{max} = 43.33$, $Z_{max} = 5.53$, $p_{FWE} < 0.001$; 134 to 170 ms, $F_{max} = 77.98$, $Z_{max} = 6.88$, $p_{FWE}$

486    < 0.001; 260 ms, $F_{max} = 26.54$, $Z_{max} = 4.50$, $p_{FWE} < 0.001$), coherence (time lag at 346 ms, $F_{max} =$

487    9.98, $Z_{max} = 2.81$, $p_{FWE} < 0.001$), and the interaction effect between attention and visual coherence

488    (time lag between 214 to 238 ms, $F_{max} = 14.82$, $Z_{max} = 3.43$, $p_{FWE} < 0.001$). We therefore subjected

489    the 3rd PC to further analyses described below.

490    Post-hoc tests on this principal component supported the idea that attention dominates the neural

491    response, but that temporal coherence can modulate it. In keeping with the main ERP results, the

492    main effect of audiovisual temporal coherence was apparent in the unattended stream, suggesting

493    that the effect of attention may be strong enough to elicit a ceiling effect. Specifically, we observed

494    main effect of attention (AtVm > AmVt: 86 - 244 ms, cluster-level $p_{FWE} < 0.001$, $T_{max} = 8.16$;

495    AtVt > AtVm at 38 ms, cluster-level $p_{FWE} < 0.001$, $T_{max} = 4.60$; at 178 ms, cluster-level $p_{FWE} =$

496    0.032, $T_{max} = 3.45$; Figure 5C). The effect of attention on the incoherent stream extends over more

497    time points than the effect of attention on the coherent stream. Consistent with this being due to a

498    temporal coherence mediated enhancement of the masker stream, the deviant-evoked responses in

499    the masker were significantly greater when accompanied by a temporally coherent visual stimulus

500    (AmVm>AmVt: 100-132 ms, $T_{max} = 3.79$, cluster-level $p_{FWE} < 0.001$; 240 to 268 ms, cluster-level

501    $p_{FWE} < 0.001$, $T_{max} = 3.55$; Figure 5C). The PC was dominated by the responses from the left

502    temporal and right frontal channels (Figure 5C, last column).

503    **Correlations between behavioral performance and EEG**

504    To examine the relationship between the EEG responses and behavioral performance, we

505    calculated Pearson correlation coefficients between measures of behavioural performance and

21

506   neural activity. Outliers were deleted using Cook's distance if the distance was larger than 3 times

507   the means of Cook's distance. We first considered whether the magnitude of the deviant response

508   in the target stream correlated with overall behavioural performance (mean d' across all visual

509   conditions), reasoning that participants with a stronger deviant response might be better able to

510   accurately report timbre deviants. For both PC1 and PC3, the peak-to-peak PCs of ERP amplitudes

511   obtained for the deviants in the target stream (At) correlated with overall d' performance (PC1

512   peak-to-peak amplitude: Figure 6A, r = 0.55, p = 0.019; PC3: Figure 6B, r = 0.61, p = 0.005).

513   The auditory selective attention task required that participants not only detect timbre deviants, but

514   that they successfully differentiated target and masker events. We therefore hypothesised that

515   listeners who more successfully engaged selective attention mechanisms might show larger

516   differences in the magnitude of deviant response to target and masker deviants. To test this we

517   subtracted the peak to peak amplitude of EEG responses for masker deviants from the peak to peak

518   amplitude to target deviants, and then measured the correlation between the EEG responses

519   difference with the behavioral performance (d'). This relationship was observed for PC3 (Figure

520   6C, r=0.67, p=0.001), but not PC1 (r=-0.01, p=0.971)..

521   Finally, while the visual condition did not significantly influence behavioural performance at the

522   group level, there was significant heterogeneity within our listeners. To determine whether

523   modulation of behavioural performance by the visual stimulus correlated with the magnitude of

524   the attention × visual condition interaction in PC3, we considered the difference in the normalised

525   d' performance for target-coherent and masker-coherent trials (i.e. the difference in target-coherent

526   d' and masker-coherent performance d' / overall d') and correlated this with the difference in the

527   attentional modulation of the 3$^{rd}$ PC peak-to-peak amplitude across visual conditions, i.e. AtVt-

528   AmVt vs AtVm-AmVm (Figure 6D, r = 0.51, p = 0.031). While the correlation was significant

529    and in the predicted direction (i.e. participants who showed a benefit for target-coherent trials had

530    a greater attentional modulation in the target-coherent condition), we note that it's principally

531    driven by a single participant whose removal renders the correlation non-significant.

## Discussion

533    This study used an auditory selective attention task, performed in the presence of a temporally

534    modulated visual stimulus, to dissect the neural signatures of selective attention and audiovisual

535    temporal coherence. Our EEG data of envelope responses reveal evidence for audiovisual

536    integration of temporally coherent audiovisual envelopes which occurred independently of

537    selective attention. Meanwhile, selective attention had a strong effect on the amplitude of TRFs

538    derived from the envelope responses, with TRFs corresponding to target streams yielding higher

539    amplitudes than those corresponding to masker streams. To further investigate audiovisual binding

540    we examined the EEG responses to the timbre deviants which occurred independently of the

541    amplitude envelopes of the audio(visual) streams. The fact that the EEG responses elicited by the

542    timbre deviants were affected by the visual coherence of the stimulus can be interpreted as

543    evidence that temporal coherence in the audiovisual streams favored the emergence of a fused

544    audiovisual percept, which contrasts more strongly against the deviants  than a purely auditory

545    stream would. In direct support of this notion, we observed that, in some spatiotemporal

546    components of the neural response, audiovisual temporal coherence interacted with selective

547    attention.

**Temporal coherence based AV integration occurs independently of attention**

549    Based on the stimulus envelope reconstruction analysis, we found that the cortical responses to the

550    AV amplitude envelope were better explained by an AV integration model than by a linear

23

551  summation (A+V) model in both the attended and unattended streams, suggesting attention was

552  not required to link audio and visual streams. Our study thus provides evidence that AV integration

553  based on temporal coherence between the auditory and visual stream can occur independently of

554  attention. This result is in contrast to previous studies using speech as stimuli. Ahmed et al., (2021)

555  found AV integration was only observed for attended speech stream, demonstrating that responses

556  to attended speech were better explained by an AV model, while the responses to unattended

557  speech were better explained by the A+V model. However, their integration model outperformed

558  the linear summation model for unattended speech at very short (0-100 ms lag) latencies,

559  suggesting that distinct multisensory computations occur at different processing stages. In contrast

560  to studies utilizing natural speech and videos of faces, our visual disc was much simpler. One

561  possibility, which is already noted in Atilgan et al. (2018), is that bottom up audiovisual integration

562  does occur independently of attention for simple non-speech stimuli. Another possibility is that

563  watching a competing talker is more distracting than watching an uninformative disc, perhaps

564  leading to observers actively suppressing a competing face in the context of a selective attention

565  task. A final difference might be that subjects in Ahmed et al. (2021) were instructed to look at the

566  eyes of the face, whereas our listeners fixated on the disc itself; potentially the radius changes of

567  the disk, presented at the fovea, provide a more salient temporal cue. In support of this possibility,

568  we note that the stimulus reconstruction accuracy of the visual-only decoder in the independent

569  condition was quite high, and significantly larger than that of the audio-only decoder.

570  We used a forward model to examine the cortical representation of the sound amplitude envelope

571  across all EEG channels. Two-way repeated measures ANOVA indicated significant main effects

572  of attention and integration. In the unattended sound stream, the $TRF_{AV}$ amplitude was

573  significantly stronger than the summation of $TRF_A$ and $TRF_V$ amplitude, which suggests that AV

574 integration occurs independently of attention. This result is consistent with our results from the

575 envelope reconstruction (Figure 2), as well as a previous study from Crosse et al (2015), both in

576 terms of the direction of the effect (AV vs. A+V) and its latency in the ~200 ms range. Furthermore,

577 attention strongly modulated the TRF, with the $TRF_{AV}$ amplitude for the target stream being

578 significantly larger than that for the masker stream. This finding is consistent with previous studies,

579 demonstrating an enhancement of attended speech streams (Ding & Simon, 2012; Mesgarani &

580 Chang, 2012) and audiovisual streams (Zion Golumbic et al., 2013). An open question is why

581 audiovisual temporal coherence did not influence the attended stream $TRF_{AV.}$ Perhaps the

582 enhancement of the TRF by attention generated a ceiling effect, or possibly if we had required

583 subjects to attend to the visual stimulus we might have observed stronger audiovisual interactions.

584 Nevertheless, our TRF results reveal the effects of both audiovisual temporal coherence and

585 attention on the TRF amplitude.

586 **Attention and coherence effects on the deviant evoked responses**

587 In this study, we adapted the behavioral paradigm of previous studies (Atilgan & Bizley, 2021;

588 Maddox et al., 2015), however, we failed to replicate the behavioural findings. Two key

589 differences may explain this: first, the magnitude of the timbre deviants was increased, which

590 effectively rendered the task easier. The overall d' scores are higher in the current dataset than in

591 previous ones. A recent study (Cappelloni et al., 2022) also suggested that the temporal coherence

592 of the visual stream might not provide additional benefit if the two auditory streams were easily

593 segregated. Second, in these previous studies, listeners were also required to detect occasional

594 colour deviants in the visual stimulus, which required them to maintain some level of attention

595 towards the visual modality. In our experiment, the visual stimulus neither contained deviants of

596 its own, nor did it provide cues that might facilitate the detection of auditory deviants. Within the

597    framework of the model included in Figure 1B, attending to the visual stream would lead to further

598    enhancement. It is possible that this difference explains why, at the group level, we did not observe

599    a significant effect of audiovisual temporal coherence on auditory deviant detection.

600    A whole-scalp analysis of deviant-evoked ERPs brought evidence for a main effect of attention,

601    with the latency of the effect corresponding to a P300 time window. The P300 is a later component

602    in response to novelty occurring between 200-600 ms relative to deviant onset, and has been

603    previously described for the auditory and visual modalities (Friedman et al., 2001). Previous

604    studies showed that the P300 is attention-dependent (Polich et al., 2007), consistent with our

605    findings. The anterior-posterior topography of the effect shown on Figure 4 is due to our choice of

606    re-referencing to the average of all channels. In addition to this robust modulation of the deviant

607    response by attention, a further PCA based on the timbre deviant elicited ERPs revealed

608    interactions between attention and audiovisual temporal coherence. For specific principal

609    components, there was an attention-dependent enhancement of the deviant-evoked responses in

610    the target stream independent of the visual coherent. This suggests that the attentional modulation

611    of the target stream is sufficiently strong that temporal coherence exerts no additional effect. We

612    found the main effect of attention to modulate activity at very early latencies (8 - 84 ms), although

613    cluster-based statistics do not indicate that all time points within this time window show significant

614    effects, but rather that there are some time points within the cluster that show significant effects.

615    The post-hoc test showed that the early peak of the attention effect was at 38 ms (Figure 5C, AtVt

616    vs AtVm). Previous studies has shown similarly early attention effects on auditory responses, e.g.

617    in a previous MEG study (Auksztulewicz et al., 2015), a main effect of attention was observed

618    around 27-40 ms after tone onset. Such early latencies are consistent with earlier results obtained

619    in attentional paradigms based on auditory filtering (Rif et al., 1991) and could be interpreted as

26

620    evidence of attentional gating (Lange 2013). However, for the unattended stream, temporal

621    coherence does enhance the deviant evoked response in the masker stream. One possibility

622    therefore is that in this paradigm the attentional modulation was sufficiently strong that, for target

623    sounds, there was a ceiling effect preventing any further modulation by audiovisual temporal

624    coherence (equivalent in the model in Figure 1B to the magnitude of attentional enhancement

625    rendering small changes due to audiovisual integration as irrelevant to the eventual summed

626    activity). Some caution is required in interpreting these results given that the $3^{rd}$ PC accounted for

627    only 4% of the variance in the EEG data, but it is noteworthy that this PC also correlated with

628    differences in task performance. The magnitude of attentional modulation scaling with overall

629    behavioural performance d' (Figure 6C). There was some evidence for a correlation between the

630    extent to which the visual condition influenced behavioural performance and the magnitude of the

631    temporal coherence dependent attentional effects (Figure 6D), although this requires replication,

632    preferably in the context of task parameters that more reliably elicit a modulation of task

633    performance by audiovisual temporal coherence. That we see significant audiovisual integration

634    in the envelope tracking responses, but not in behaviour or in the main ERP analysis (Figure 4) of

635    the timbre deviants, potentially suggests that both behaviour and timbre deviant responses are

636    dominated by attentional effects. Future experiments could make attentional selection harder, for

637    example by making the pitch or timbre of the two streams more similar, in order to determine

638    whether it is possible to unmask audiovisual temporal coherence effects that are hinted at by our

639    PCA of the timbre deviant responses.

640    Our results are consistent with previous studies on 'cocktail party effect' speech stream segregation,

641    in which congruent visual stimuli enhanced the cortical representation of the speech envelope of

642    attended speech streams relative to unattended streams (Crosse, Di Liberto, & Lalor, 2016;

643 Golumbic et al., 2013). However, unlike in these previous studies, where visual speech provided

644 temporal and contextual information about the auditory envelope, we used a simple disc as a visual

645 stimulus, which provided no information about the auditory deviant. While previous studies have

646 demonstrated that attention dedicated to one feature of an object may enhance the responses to

647 other features of the object in both auditory (Alain & Arnott, 2000; Maddox & Shinn-Cunningham,

648 2012; Shamma et al., 2011; Shinn-Cunningham, 2008) and visual modalities (Blaser et al., 2000;

649 O'Craven et al., 1999), our results provide new evidence that temporal coherence modulates the

650 attentional enhancement of the neural response to the timbre deviants ("other" features) of the AV

651 object.

652 In summary, we examined the temporal coherence and attention effect on neural responses to the

653 continuous sound envelope and the deviant evoked response, respectively. Temporal coherence

654 facilitated the audiovisual integration independent of attention, and attention further enhanced the

655 audiovisual integration of the coherent audiovisual stream. Attention amplified a large portion of

656 the deviant-evoked response independent of temporal coherence, while coherence only modulated

657 deviant-evoked responses in the unattended auditory stream. These results provide evidence for

658 partly dissociable neural signatures of bottom-up (coherence) and top-down (attention) effects in

659 AV object formation.

## Acknowledgements

## References

Ahmed, F., Nidiffer, A. R., O'sullivan, A. E., Zuk, N. J., & Lalor, E. C. (2021). The integration of continuous audio and visual speech in a cocktail-party environment depends on attention. *BioRxiv*. https://doi.org/10.1101/2021.02.10.430634

Alain, C., & Arnott, S. R. (2000). Selectively attend to auditory objects. In *Frontiers in bioscience : a journal and virtual library* (Vol. 5). https://doi.org/10.2741/a505

Atilgan, H., & Bizley, J. K. (2021). Training enhances the ability of listeners to exploit visual information for auditory scene analysis. *Cognition*, *208*, 104529. https://doi.org/10.1016/j.cognition.2020.104529

Atilgan, H., Town, S. M., Wood, K. C., Jones, G. P., Maddox, R. K., Lee, A. K. C., & Bizley, J. K. (2018). Integration of Visual Information in Auditory Cortex Promotes Auditory Scene Analysis through Multisensory Binding. *Neuron*, *97*(3), 640-655.e4. https://doi.org/10.1016/j.neuron.2017.12.034

Auksztulewicz, R., & Friston, K. (2015). Attentional enhancement of auditory mismatch responses: a DCM/MEG study. Cerebral cortex, 25(11), 4273-4283.

Beauchamp, M. S., Lee, K. E., Argall, B. D., & Martin, A. (2004). Integration of auditory and visual information about objects in superior temporal sulcus. *Neuron*. https://doi.org/10.1016/S0896-6273(04)00070-4

686     Besle, J., Fort, A., Delpuech, C., & Giard, M. H. (2004). Bimodal speech: Early suppressive

687        visual effects in human auditory cortex. *European Journal of Neuroscience*, *20*(8),

688        2225–2234. https://doi.org/10.1111/j.1460-9568.2004.03670.x

689     Bizley, J. K., Maddox, R. K., & Lee, A. K. C. (2016). Defining Auditory-Visual Objects:

690        Behavioral Tests and Physiological Mechanisms. *Trends in Neurosciences*, *39*(2), 74–

691        85. https://doi.org/10.1016/j.tins.2015.12.007

692     Blaser, E., Pylyshyn, Z. W., & Holcombe, A. O. (2000). Tracking an object through feature

693        space. *Nature 2000 408:6809*, *408*(6809), 196–199. https://doi.org/10.1038/35041567

694     Busse, L., Roberts, K. C., Crist, R. E., Weissman, D. H., & Woldorff, M. G. (2005). The

695        spread of attention across modalities and space in a multisensory object. *Proceedings of*

696        *the National Academy of Sciences*, *102*(51), 18751–18756.

697        https://doi.org/10.1073/PNAS.0507704102

698     Calvert, G. A., & Campbell, R. (2003). Reading speech from still and moving faces: the neural

699        substrates of visible speech. *Journal of Cognitive Neuroscience*, *15*(1), 57–70.

700     Cappelloni, M. S., Mateo, V. S., & Maddox, R. K. (2022). Humans rely more on talker

701        identity than temporal coherence in an audiovisual selective attention task using speech-

702        like stimuli. bioRxiv.

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

704        The Natural Statistics of Audiovisual Speech. *PLoS Computational Biology*, *5*(7),

705        e1000436. https://doi.org/10.1371/journal.pcbi.1000436

706     Crosse, M. J., Butler, J. S., & Lalor, E. C. (2015). Congruent visual speech enhances

707        cortical entrainment to continuous auditory speech in noise-free conditions. *Journal of*

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

709       15.2015

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

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

712       continuous stimuli. *Frontiers in Human Neuroscience*, *10*(NOV2016), 1–14.

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

714    Crosse, M. J., Di Liberto, G. M., & Lalor, E. C. (2016). Eye can hear clearly now: Inverse

715       effectiveness in natural audiovisual speech processing relies on long-term crossmodal

716       temporal integration. *Journal of Neuroscience*, *36*(38), 9888–9895.

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

718    de Cheveigné, A., & Simon, J. Z. (2008). Denoising based on spatial filtering. *Journal of*

719       *Neuroscience Methods*, *171*(2), 331–339.

720    Ding, N., & Simon, J. Z. (2012). Emergence of neural encoding of auditory objects while

721       listening to competing speakers. *Proceedings of the National Academy of Sciences of the*

722       *United States of America*, *109*(29), 11854–11859.

723       https://doi.org/10.1073/pnas.1205381109

724    Friedman, D., Cycowicz, Y. M., & Gaeta, H. (2001). The novelty P3: an event-related brain

725       potential (ERP) sign of the brain's evaluation of novelty. Neurosci Biobehav Rev, 25(4),

726       355-373. https://doi.org/10.1016/s0149-7634(01)00019-7

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

728       Enhances Selective Speech Envelope Tracking in Auditory Cortex at a "Cocktail Party."

729   *The Journal of Neuroscience*, *January*. https://doi.org/10.1523/JNEUROSCI.3675-

730   12.2013

731   Haufe, S., Meinecke, F., Görgen, K., Dähne, S., Haynes, J. D., Blankertz, B., & Bießmann,

732   F. (2014). On the interpretation of weight vectors of linear models in multivariate

733   neuroimaging.   *NeuroImage*,   *87*,   96–110.

734   https://doi.org/10.1016/J.NEUROIMAGE.2013.10.067

735   Kayser, C., & Logothetis, N. K. (2009). Directed interactions between auditory and superior

736   temporal cortices and their role in sensory integration. *Frontiers in Integrative*

737   *Neuroscience*, *3*, 7.

738   Kayser, C., & Shams, L. (2015). Multisensory Causal Inference in the Brain. *PLOS Biology*,

739   *13*(2), e1002075. https://doi.org/10.1371/JOURNAL.PBIO.1002075

740   Kilner, J. M., Kiebel, S. J., & Friston, K. J. (2005). Applications of random field theory to

741   electrophysiology.   *Neuroscience*   *Letters*,   *374*(3),   174–178.

742   https://doi.org/10.1016/j.neulet.2004.10.052

743   Kondo, H., Saleem, K. S., & Price, J. L. (2003). Differential connections of the temporal pole

744   with the orbital and medial prefrontal networks in macaque monkeys. *Journal of*

745   *Comparative Neurology*, *465*(4), 499–523.

746   Kondo, H., Saleem, K. S., & Price, J. L. (2005). Differential connections of the perirhinal and

747   parahippocampal cortex with the orbital and medial prefrontal networks in macaque

748   monkeys. *Journal of Comparative Neurology*, *493*(4), 479–509.

749   Lange, K. (2013). The ups and downs of temporal orienting: a review of auditory temporal

750   orienting studies and a model associating the heterogeneous findings on the auditory N1

751        with opposite effects of attention and prediction. Frontiers in human neuroscience, 7,

752        263.

753    Litvak, V., & Friston, K. (2008). Electromagnetic source reconstruction for group studies.

754        *Human        Brain        Mapping        Journal*,        *42*,        1490–1498.

755        https://doi.org/10.1016/j.neuroimage.2008.06.022

756    Lovelace, C. T., Stein, B. E., & Wallace, M. T. (2003). An irrelevant light enhances auditory

757        detection in humans: a psychophysical analysis of multisensory integration in stimulus

758        detection. *Cognitive Brain Research*, *17*(2), 447–453. https://doi.org/10.1016/S0926-

759        6410(03)00160-5

760    Maddox, R. K., Atilgan, H., Bizley, J. K., & Lee, A. K. C. (2015). Auditory selective attention

761        is enhanced by a task-irrelevant temporally coherent visual stimulus in human listeners.

762        *ELife*, *4*, 1–11. https://doi.org/10.7554/eLife.04995

763    Maddox, R. K., & Shinn-Cunningham, B. G. (2012). Influence of task-relevant and task-

764        irrelevant feature continuity on selective auditory attention. *JARO - Journal of the*

765        *Association    for    Research    in    Otolaryngology*,    *13*(1),    119–129.

766        https://doi.org/10.1007/s10162-011-0299-7

767    Mcgurk, H., & Macdonald, J. (1976). Hearing lips and seeing voices. *Nature*.

768        https://doi.org/10.1038/264746a0

769    Mesgarani, N., & Chang, E. F. (2012). Selective cortical representation of attended speaker

770        in multi-talker speech perception. In *Nature* (Vol. 485, Issue 7397, pp. 233–236). NIH

771        Public Access. https://doi.org/10.1038/nature11020

772  Michel, M., & Morales, J. (2020). Minority reports: Consciousness and the prefrontal cortex.

773  *Mind and Language*, *35*(4), 493–513. https://doi.org/10.1111/mila.12264

774  Middlebrooks, J. C., Simon, J. Z., Popper, A. N., & Fay, R. R. (2017). *The auditory system*

775  *at the cocktail party* (Vol. 60). Springer.

776  Mishra, J., Martinez, A., Sejnowski, T. J., & Hillyard, S. A. (2007). Early cross-modal

777  interactions in auditory and visual cortex underlie a sound-induced visual illusion.

778  *Journal of Neuroscience*, *27*(15), 4120–4131.

779  Molholm, S., Sehatpour, P., Mehta, A. D., Shpaner, M., Gomez-Ramirez, M., Ortigue, S.,

780  Dyke, J. P., Schwartz, T. H., & Foxe, J. J. (2006). Audio-visual multisensory integration

781  in superior parietal lobule revealed by human intracranial recordings. *Journal of*

782  *Neurophysiology*, *96*(2), 721–729.

783  Morís Fernández, L., Visser, M., Ventura-Campos, N., Ávila, C., & Soto-Faraco, S. (2015).

784  Top-down attention regulates the neural expression of audiovisual integration.

785  *NeuroImage*, *119*, 272–285. https://doi.org/10.1016/J.NEUROIMAGE.2015.06.052

786  Ille, N., Berg, P., & Scherg, M. (2002). Artifact correction of the ongoing EEG using spatial

787  filters based on artifact and brain signal topographies. *Journal of clinical*

788  *neurophysiology*, *19*(2), 113-124. https://doi.org/10.1097/00004691-200203000-00002

789  O'Craven, K. M., Downing, P. E., & Kanwisher, N. (1999). fMRI evidence for objects as the

790  units of attentional selection. *Nature 1999 401:6753*, *401*(6753), 584–587.

791  https://doi.org/10.1038/44134

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

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

34

794      Cocktail Party Environment Can Be Decoded from Single-Trial EEG. *Cerebral Cortex*,

795      *25*(7), 1697–1706. https://doi.org/10.1093/cercor/bht355

796      Petrides, M, & Pandya, D. N. (1999). Dorsolateral prefrontal cortex: comparative

797      cytoarchitectonic analysis in the human and the macaque brain and corticocortical

798      connection patterns. *European Journal of Neuroscience*, *11*(3), 1011–1036.

799      Petrides, Michael, & Pandya, D. N. (2002). Comparative cytoarchitectonic analysis of the

800      human and the macaque ventrolateral prefrontal cortex and corticocortical connection

801      patterns in the monkey. *European Journal of Neuroscience*, *16*(2), 291–310.

802      Polich, J. (2007). Updating P300: an integrative theory of P3a and P3b. *Clinical*

803      *neurophysiology*, *118*(10), 2128-2148.

804      Rif, J., Hari, R., Hämäläinen, M. S., & Sams, M. (1991). Auditory attention affects two

805      different areas in the human supratemporal cortex. Electroencephalography and clinical

806      Neurophysiology, 79(6), 464-472.

807      Rohe, T., Ehlis, A.-C., & Noppeney, U. (2019). The neural dynamics of hierarchical Bayesian

808      causal inference in multisensory perception. *Nature Communications*, *10*(1), 1–17.

809      Shamma, S. A., Elhilali, M., & Micheyl, C. (2011). Temporal coherence and attention in

810      auditory scene analysis. *Trends in Neurosciences*, *34*(3), 114–123.

811      https://doi.org/10.1016/j.tins.2010.11.002

812      Shinn-cunningham, B. G. (2008). *Object-based auditory and visual attention*. *April*, 182–186.

813      https://doi.org/10.1016/j.tics.2008.02.003

814      Talsma, D., Doty, T. J., & Woldorff, M. G. (2007). Selective attention and audiovisual

815          integration: Is attending to both modalities a prerequisite for early integration? *Cerebral*

816          *Cortex*, *17*(3), 679–690. https://doi.org/10.1093/cercor/bhk016

817      Talsma, D., Senkowski, D., Soto-Faraco, S., & Woldorff, M. G. (2010). The multifaceted

818          interplay between attention and multisensory integration. *Trends in Cognitive Sciences*,

819          *14*(9), 400–410. https://doi.org/10.1016/J.TICS.2010.06.008

820      Talsma, D., & Woldorff, M. G. (2005). Selective attention and multisensory integration:

821          Multiple phases of effects on the evoked brain activity. *Journal of Cognitive*

822          *Neuroscience*, *17*(7), 1098–1114. https://doi.org/10.1162/0898929054475172

823      Van der Burg, E., Olivers, C. N. L., Bronkhorst, A. W., & Theeuwes, J. (2008). Pip and Pop:

824          Nonspatial Auditory Signals Improve Spatial Visual Search. *Journal of Experimental*

825          *Psychology: Human Perception and Performance*, *34*(5), 1053–1065.

826          https://doi.org/10.1037/0096-1523.34.5.1053

827      Van Wassenhove, V., Grant, K. W., & Poeppel, D. (2005). Visual speech speeds up the neural

828          processing of auditory speech. *Proceedings of the National Academy of Sciences of the*

829          *United States of America*, *102*(4), 1181–1186. https://doi.org/10.1073/pnas.0408949102

830      Yuan, Y., Wayland, R., & Oh, Y. (2020). Visual analog of the acoustic amplitude envelope

831          benefits speech perception in noise. *The Journal of the Acoustical Society of America*,

832          *147*(3), EL246–EL251. https://doi.org/10.1121/10.0000737

833      Zion Golumbic, E. M., Ding, N., Bickel, S., Lakatos, P., Schevon, C. A., McKhann, G. M.,

834          Goodman, R. R., Emerson, R., Mehta, A. D., Simon, J. Z., Poeppel, D., & Schroeder, C.

835    E. (2013). Mechanisms underlying selective neuronal tracking of attended speech at a

836    "cocktail party." *Neuron*, *77*(5), 980–991. https://doi.org/10.1016/j.neuron.2012.12.037

837  Zumer, J. M., White, T. P., & Noppeney, U. (2020). The neural mechanisms of audiotactile

838    binding depend on asynchrony. *European Journal of Neuroscience*, *52*(12), 4709–4731.

839    https://doi.org/10.1111/EJN.14928

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855 **Figure captions**

856 *Figure 1. Experimental paradigm, diagram of the possible effects of attention and coherence, and*

857 *behavioral performance. (A) Schematic plot of auditory and visual stimuli in the behavioral task.*

858 *Amplitude envelopes of target/attended sound (grey solid line), masker/unattended sound (red solid line),*

859 *and visual radius envelope (blue dashed line). (B) Diagram of the possible effects of attention and*

860 *coherence. Attention and coherence effects can be mapped onto four main scenarios: for non-coherent*

861 *stimuli in the absence of selective attention, the response is a linear summation of the three streams; for*

862 *coherent stimuli, there is an additional term representing audiovisual integration. In the context of a task*

863 *requiring selective attention to one sound, attention can enhance (illustrated by larger terms) either the*

864 *auditory stream only (for incoherent stimuli) or additionally the integration term (for coherent stimuli).*

865 *This model assumes that when temporal coherence is absent the relevant integration term becomes zero,*

866 *and these are shown in gray. (C) Behavioral sensitivity (d') for each visual condition. Each line shows data*

867 *of one participant. Solid lines indicate participants with higher d' in the AtAmVt vs. AtAmVm condition,*

868 *and dashed lines indicate participants with lower d' in the AtAmVt vs. AtAmVm condition. Black squares*

869 *represent group averages, and error bars indicate the standard error of the mean (SEM).*

870 *Figure 2. Stimulus reconstruction. (A) Examples of the original sound envelope (grey) with the grand-*

871 *average neural reconstruction (black) overlapped. The mean reconstruction accuracy over subjects is*

872 *indicated to the right. (B) The stimulus reconstruction accuracy for each stream in the independent*

873 *condition AtAmVi was significantly better than chance (permutation test). Each dot represents one subject.*

874 *(C, D) The stimulus reconstruction accuracy using the AV decoder and A+V decoder for the target and*

875 *masker sound was significantly better than chance (permutation test), respectively.*

876 *Figure 3. Temporal response function analysis. (A, B) Left panel, the TRF estimated for coherent target*

877 *stream (AtVt) had a stronger amplitude than that for the masker stream (AmVm). Right panel, the*

878 *summation of TRFs estimated for the target stream (At + Vi) was significantly stronger than that for the*

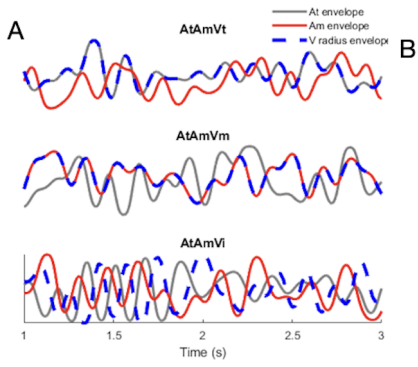879 *masker stream (Am + Vi). (C) Left panel, for the target sound (At) condition, the TRF estimated for coherent*

880    *AV streams (AtVt) was not significantly different from the summation of TRFs estimated for independent*

881    *AV streams (At + Vi). Right panel, for the masker sound (Am) condition, the TRF estimated for coherent*

882    *AV streams (AmVm) had a stronger amplitude than the summation of TRFs estimated for independent AV*

883    *streams (Am + Vi). Shaded areas indicate SEM (standard error of the mean) across subjects. The*

884    *topographical plot shows the EEG channel locations with a significant difference. Black horizontal bars:*

885    *$p_{FWE}$ < 0.05 (based on the main effects in the ANOVA; see Results).*

886    ***Figure 4. Grand-average deviant-evoked ERPs over participants and channels across conditions. (A)***

887    *Scalp topographies from deviant onset to 0.5 s after onset. Each row represents one condition (from top to*

888    *bottom, the condition corresponds to AtVm, AmVt, AtVt, and AmVm, respectively), each column represents*

889    *one 50-ms time window. **(B)** Deviants presented in the incoherent target stream (red dashed lines) and*

890    *masker stream (grey solid lines); **(C)** Deviants presented in the AV coherent target (red solid lines) and*

891    *masker stream (black dashed lines); **(D,E)** Deviants presented in the target and masker stream in each of*

892    *the two attentional conditions (left: masker stream; right: target stream); The topographical plots in panels*

893    *(A-C) show the EEG channel locations where a significant ERP amplitude difference between the two*

894    *conditions (as indicated at the top of each plot) was observed (FWE-corrected) except. The topographical*

895    *plot in (D) shows the EEG channel locations same as the locations in (C). The black bar represents the*

896    *time segment with a significant difference between the deviants in two different conditions. Shaded areas*

897    *represent SEM across subjects.*

898    ***Figure 5. Attentional enhancement of deviant-evoked ERPs: principal component analysis. (A,B,C,D)***

899    *Deviant-evoked response for the $1^{st}$ PC, $2^{nd}$ PC, $3^{rd}$ PC, and $4^{th}$ PC of ERP, respectively. The first two*

900    *columns represent the attention effect on the AV incoherent conditions (AtVm in red dashed lines and AmVt*

901    *in black solid lines) and AV incoherent conditions (AtVt in red solid lines and AmVm in black dashed lines).*

902    *Black bars indicate the time periods with a significant difference between the two conditions, and black*

903    *asterisks indicate the time points with a significant difference between the two conditions. Shaded areas*

904    *indicate SEM across subjects. The third and fourth columns represent the AV coherence effect on the target*

39

905 *conditions (AtVt in blue solid lines and AtVm in purple solid lines) and masker conditions (AmVt in purple*

906 *dashed lines and AmVm in blue dashed lines). The fifth column represents the spatial topography map of*

907 *the principal component weights across channels. Color indicates the weight (warm: high, cool: low).*
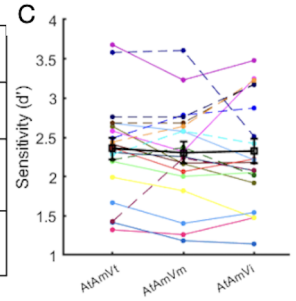
908 ***Figure 6. Correlations between the behavioral performance and EEG responses. (A, B)*** *The correlation*

909 *between mean d' and the mean 1$^{st}$ PC and 3$^{rd}$ PC peak-to-peak amplitude over conditions (AtAmVt,*

910 *AtAmVm, and AtAmVi), respectively.* ***(C)*** *The correlation between mean d' and the mean 3$^{rd}$ PC peak-to-*

911 *peak amplitude (At - Am).* ***(D)*** *Visual coherence modulation of behaviour performance with EEG responses.*

912 *The correlation between the hit rate difference (AtAmVt - AtAmVm) and the 3$^{rd}$ PC peak-to-peak amplitude*

913 *(AtVt – AmVt vs AtVm - AmVm). The unfilled circles represent outliers. (P value corrected for multiple*
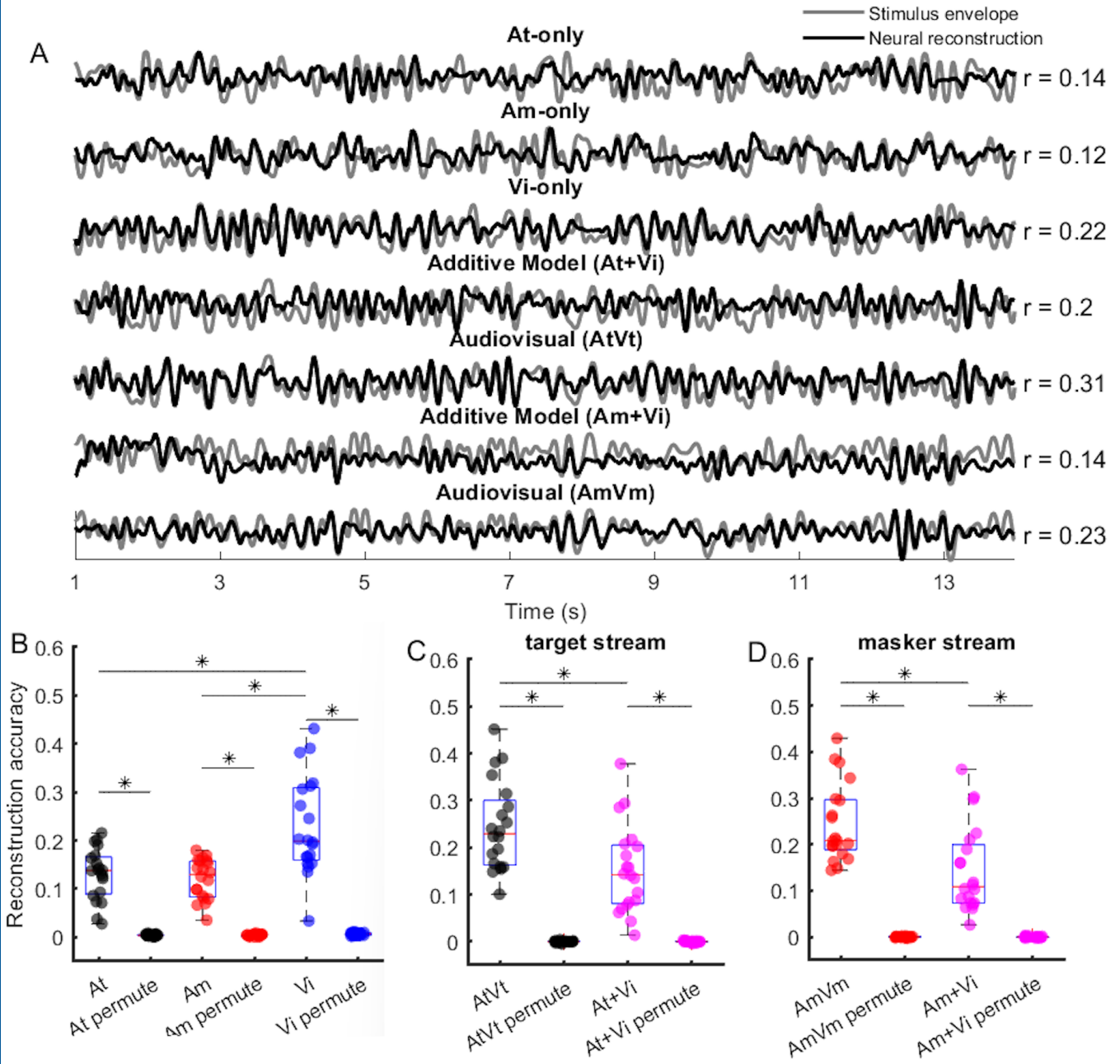
914 *comparison.)*

**A**

AtAmVt

AtAmVm

AtAmVi

Legend:
- At envelope
- Am envelope
- V radius envelope

Time (s)

**B**

$$\overline{\sum(.)} \;=\; \overline{(At+Am+V)^2} \;=\; \underbrace{\overline{(At)^2}+\overline{(Am)^2}+\overline{(V)^2}}_{\text{Summation}} \;+\; \underbrace{2\,(\overline{AtV}+\overline{AmV}+\overline{AtAm})}_{\text{Integration}}$$
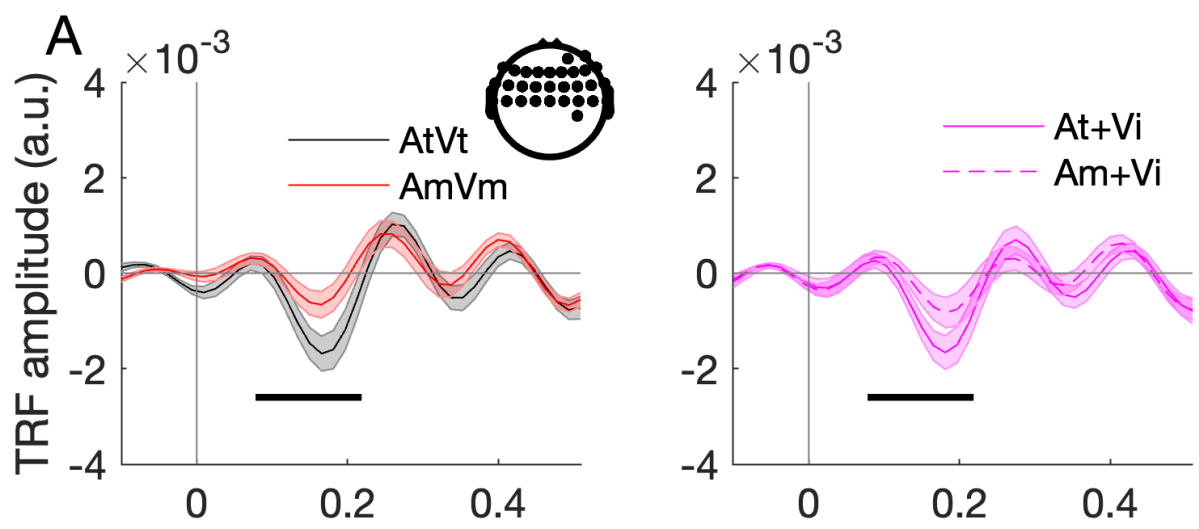
| | | |
|---|---|---|
| Independent AV, no attention | $\overline{(A_1)^2}+\overline{(A_2)^2}+\overline{(V)^2}\;+$ <br> $2\,(\overline{A_1V}+\overline{A_2V}+\overline{A_1A_2}\,)$ | Summation |
| Visual temporal coherent to A1, no attention | $\overline{(A_1)^2}+\overline{(A_2)^2}+\overline{(V)^2}\;+$ <br> $2\,(\overline{A_1V}+\overline{A_2V}+\overline{A_1A_2}\,)$ | Summation <br> + Integration |
| Visual temporal coherent to At (attended to At) | $\overline{(At)^2}+\overline{(Am)^2}+\overline{(Vt)^2}\;+$ <br> $2\,(\overline{AtV}+\overline{AmV}+\overline{AtAm}\,)$ | Summation <br> + Integration |
| Visual temporal coherent to Am (attended to At) | $\overline{(At)^2}+\overline{(Am)^2}+\overline{(Vt)^2}\;+$ <br> $2\,(\overline{AtV}+\overline{AmV}+\overline{AtAm}\,)$ | Summation <br> + Integration |

Gray indicates term = 0

**C**

Sensitivity (d')

AtAmVt    AtAmVm    AtAmVi

A

At-only
r = 0.14

Am-only
r = 0.12

Vi-only
r = 0.22

Additive Model (At+Vi)
r = 0.2

Audiovisual (AtVt)
r = 0.31

Additive Model (Am+Vi)
r = 0.14

Audiovisual (AmVm)
r = 0.23

Stimulus envelope
Neural reconstruction

Time (s)

B

Reconstruction accuracy

At  At permute  Am  Am permute  Vi  Vi permute

C  target stream

AtVt  AtVt permute  At+Vi  At+Vi permute

D  masker stream

AmVm  AmVm permute  Am+Vi  Am+Vi permute

**A** r = 0.55 p = 0.005

Mean PC1 peak-to-peak amplitude vs Behavioural performance Mean d'

**B** r = 0.61 p = 0.001

Mean PC3 peak-to-peak amplitude vs Behavioural performance Mean d'

**C** r = 0.67 p < 0.001

Difference of PC3 peak-to-peak amplitude At - Am vs Behavioural performance Mean d'

**D** r = 0.51 p = 0.008

Difference of PC3 peak-to-peak amplitude (AtVt-AmVt) - (AtVm-AmVm) vs Visual modulation of behavioural performance (AtAmVt-AtAmVm)