

1 **Integration of visual information in auditory cortex promotes auditory scene analysis through**  
2 **multisensory binding**

3

4 Huriye Atilgan<sup>1</sup>, Stephen M. Town<sup>1</sup>, Katherine C. Wood<sup>1</sup>, Gareth P. Jones<sup>1</sup>, Ross K. Maddox<sup>2,3</sup>, Adrian  
5 K.C. Lee<sup>3</sup> and Jennifer K. Bizley<sup>1</sup>

6 <sup>1</sup>The Ear Institute, University College London, UK

7 <sup>2</sup>Department of Biomedical Engineering, University of Rochester, Department of Neuroscience,  
8 University of Rochester, Rochester, USA

9 <sup>3</sup>Institute for Learning and Brain Sciences and Department of Speech and Hearing Sciences, University  
10 of Washington, Seattle, USA

11

12

13

14 Corresponding Author: Jennifer Bizley [j.bizley@ucl.ac.uk](mailto:j.bizley@ucl.ac.uk)

15 UCL Ear Institute, 332 Gray's Inn Road, London, WC1X 8EE.

## 16 Abstract

17 How and where in the brain audio-visual signals are bound to create multimodal objects remains  
18 unknown. One hypothesis is that temporal coherence between dynamic multisensory signals provides  
19 a mechanism for binding stimulus features across sensory modalities in early sensory cortex. Here we  
20 report that temporal coherence between auditory and visual streams enhances spiking  
21 representations in auditory cortex. We demonstrate that when a visual stimulus is temporally  
22 coherent with one sound in a mixture, the neural representation of that sound is enhanced.  
23 Supporting the hypothesis that these changes represent a neural correlate of multisensory binding,  
24 the enhanced neural representation extends to stimulus features other than those that bind auditory  
25 and visual streams. These data provide evidence that early cross-sensory binding provides a bottom-  
26 up mechanism for the formation of cross-sensory objects and that one role for multisensory binding  
27 in auditory cortex is to support auditory scene analysis.

28 When listening to a sound of interest, we frequently look at the source. However, how auditory and  
29 visual information are integrated into a coherent perceptual object is unknown. The temporal  
30 properties of a visual stimulus can be exploited to detect correspondence between auditory and visual  
31 streams <sup>1-3</sup>, can bias the perceptual organisation of a sound scene <sup>4</sup>, and can enhance or impair  
32 listening performance depending on whether the visual stimulus is temporally coherent with a target  
33 or distractor sound stream <sup>5</sup>. Together, these behavioural results suggest that temporal coherence  
34 between auditory and visual stimuli can promote binding of cross-modal features to enable the  
35 formation of an auditory-visual (AV) object <sup>6</sup>.

36 Visual stimuli can both drive and modulate neural activity in primary and non-primary auditory cortex  
37 <sup>7-11</sup> but the contribution that visual activity in auditory cortex makes to auditory perception remains  
38 unknown. One hypothesis is that the integration of cross-sensory information into early sensory cortex  
39 provides a bottom-up substrate for the binding of multisensory stimulus features into a single  
40 perceptual object <sup>6</sup>. We have recently argued that binding is a distinct form of multisensory integration  
41 that underpins perceptual object formation, and can be separated from other sorts of integration by  
42 demonstrating a benefit in the behavioural or neural discrimination of a stimulus feature orthogonal  
43 to the features that link crossmodal stimuli (Fig. 1a). Therefore, in order to demonstrate binding an  
44 appropriate crossmodal stimulus should elicit not only enhanced neural encoding of the stimulus  
45 features that bind auditory and visual streams, but that there should be enhancement in the  
46 representation of other stimulus features associated with the source (Fig. 1c).

47 Here we test the hypothesis that the incorporation of visual information into auditory cortex can  
48 determine the neuronal representation of an auditory scene through multisensory binding (Fig.1). We  
49 demonstrate that when visual luminance changes coherently with the amplitude modulations of one  
50 sound stream in a mixture, the neural representation of that sound stream is enhanced in the auditory  
51 cortex. Consistent with these effects reflecting cross-modal binding, the encoding of auditory timbre,

52 an orthogonal stimulus dimension, is subsequently enhanced in the temporally coherent auditory  
53 stream.

## 54 Results

55 We recorded neuronal responses in the auditory cortex of awake (n=9 ferrets, 221 single units, 311  
56 multi-units) and medetomidine-ketamine anesthetised ferrets (n=5 ferrets, 426 single units, 772 multi  
57 units) in response to naturalistic time-varying auditory and visual stimuli adapted from Maddox et al<sup>5</sup>.  
58 Recordings in anesthetized animals allowed us to isolate bottom-up attention-independent  
59 processing, permitted longer recording durations for additional control stimuli and enabled  
60 simultaneous characterization of neural activity across cortical laminae. Recordings in awake animals  
61 while they held their head at a drinking spout but were not engaged in a behavioral task allowed us  
62 to measure neural activity free from any confounds associated with pharmacological manipulation  
63 and in the absence of task-directed attention. The auditory streams were two vowels with a distinct  
64 pitch and timbre (denoted A1: [u], fundamental frequency (F0) = 175 Hz and A2: [a], F0 = 195 Hz) each  
65 of which was independently amplitude modulated with a low-pass (<7 Hz) envelope (Fig. 1d). A full-  
66 field visual stimulus accompanied the auditory stimuli, the luminance of which was temporally  
67 modulated with a modulation envelope from one of the two auditory streams (Fig. 1e). We tested  
68 stimulus conditions in which a single AV stimulus pair was presented ('single stream' stimuli), where  
69 the auditory and visual streams could be temporally coherent (A1V1, A2V2) or independently  
70 modulated (A1V2, A2V1). We also tested a dual auditory stream condition in which both auditory  
71 streams were presented and the visual stimulus was temporally coherent with one of the auditory  
72 streams (A12V1 or A12V2, Fig. 1e).

### 73 **Spike patterns in auditory cortex differentiate dynamic auditory-visual stimuli**

74 Before exploring the impact of temporal coherence between auditory and visual stimuli on auditory  
75 cortical neurons, we first used the responses to single stream stimuli to classify neurons according to  
76 whether they were modulated primarily by auditory or visual temporal dynamics. To determine

77 whether the auditory amplitude envelope reliably modulated spiking we used a spike-pattern classifier  
78 to decode the auditory stream identity, collapsed across visual stimuli (i.e. we decoded auditory  
79 stream identity from the combined responses to A1V1 and A1V2 stimuli and the combination of A2V1  
80 and A2V2 responses). An identical approach was taken to determine if neuronal responses reliably  
81 distinguished visual modulation (i.e. we decoded visual identity from the combined responses to A1V1  
82 and A2V1 stimuli and the combined responses elicited by A1V2 and A2V2). Neuronal responses which  
83 were informative about auditory or visual stimulus identity at a level better than chance (estimated  
84 with a bootstrap resampling) were classified as auditory-discriminating (Fig. 2a-b) and / or visual-  
85 discriminating (Fig. 2c-d) respectively.

86 In awake animals, 39.5% (210/532) of units were auditory-discriminating, 11.1% (59/532) were visual-  
87 discriminating, and only 0.38 (2/532) discriminated both auditory and visual stimuli. Overall a smaller  
88 proportion of units represented the identity of auditory or visual streams in the anaesthetised dataset:  
89 20.2% (242/1198) were auditory-discriminating, 6.8% (82/1198) were visual-discriminating, and 0.58%  
90 (7/1198) discriminated both.

91 During recordings made under anaesthesia, we also recorded responses to noise bursts and white-  
92 light flashes (both 100 ms) presented separately and together to map AV responsiveness in auditory  
93 cortex (Bizley et al., 2007). Specifically, responsiveness was defined using a two-way ANOVA (factors:  
94 auditory stimulus [on/off] and visual stimulus [on/off]) on spike counts measured during stimulus  
95 presentation. We defined units as being sound-driven (main effect of auditory stimulus, no effect of  
96 visual stimulus or interaction), light-driven (main effect of visual stimulus, no effect of auditory  
97 stimulus or interaction) or both (main effect of both auditory and visual stimuli or significant  
98 interaction;  $p < 0.05$ ). Using such stimuli revealed that the classification of units by visual / auditory  
99 discrimination of single stream stimuli selected a subset of light or sound driven units and that the  
100 proportions of auditory, visual and AV units recorded in our sample were in line with previous studies

101 (Bizley et al, 2007: 65.1% (328/504) of units were driven by noise bursts, 16.1% (81/504) by light  
102 flashes and 14.1% (71/504) by both.

103 We hypothesised that temporal coherence between auditory and visual stimuli would enhance the  
104 discriminability of neural responses, irrespective of a unit's classification as auditory or visual  
105 discriminating. We confirmed this prediction by comparing discrimination of temporally coherent  
106 (A1V1 vs. A2V2) and temporally independent (A1V2 vs. A2V1) stimuli (Fig. 2e, f): Temporally coherent  
107 AV stimuli produced more discriminable spike patterns than those elicited by temporally independent  
108 ones in both awake (Fig. 2e, pairwise t-test, auditory-discriminating  $t_{418} = 34.277$ ,  $p < 0.001$ ; visual-  
109 discriminating  $t_{116} = 13.327$ ,  $p < 0.001$ ; All  $t_{540} = 35.196$ ,  $p < 0.001$ ) and anaesthetised recordings (Fig. 2f,  
110 auditory-discriminating  $t_{482} = 27.631$ ,  $p < 0.001$ ; visual-discriminating  $t_{162} = 22.907$ ,  $p < 0.001$ ; All  $t_{664} =$   
111  $33.149$ ,  $p < 0.001$ ).

112 What might underlie the enhanced discriminability observed for temporally coherent cross-modal  
113 stimuli? The phase of low frequency oscillations determines the excitability of the surrounding cortical  
114 tissue<sup>12-14</sup>, is reliably modulated by naturalistic stimulation<sup>15-19</sup> and has been implicated in  
115 multisensory processing<sup>20,21</sup>. We hypothesised that sub-threshold visual inputs could modulate spiking  
116 activity by modifying the phase of the local field potential such that phase coupling to temporally  
117 coherent sounds is enhanced. This in turn would provide a mechanism by which neuronal spiking was  
118 enhanced when auditory and visual streams are temporally coherent.

### 119 **Dynamic visual stimuli elicit reliable changes in LFP phase**

120 Stimulus evoked changes in the local field potential (LFP) were evident from the recorded voltage  
121 traces and analysis of cross-trial phase coherence demonstrated that there were reliable changes in  
122 phase across repetitions of identical AV stimuli (Fig. 3 a, b). To isolate the influence of visual activity  
123 on the LFP for each unit, and address the hypothesis that visual stimuli elicited reliable changes in the  
124 LFP, we calculated phase and power dissimilarity functions for stimuli with identical auditory signals  
125 but differing visual stimuli<sup>17</sup>. Briefly, this analysis assumes that if the phase (or power) within a

126 particular frequency band differs systematically between responses to two different stimuli, then  
127 inter-trial phase coherence (ITPC) across repetitions of a single stimulus will be greater than across  
128 randomly selected stimuli. For each frequency band in the LFP, we therefore compared “within-  
129 stimulus” ITPC for responses to each stimulus (A1 stream Fig. 3c; A2 stream Fig. 3d) with “across-  
130 stimulus” ITPC calculated from stimuli with identical auditory stimuli but randomly selected visual  
131 stimuli (e.g. randomly drawn from A1V1 and A1V2; Fig. 3c). The difference between within-stimulus  
132 and across-stimulus ITPC was then calculated across frequency and described as the phase  
133 dissimilarity index (PDI) (Fig. 3e, f) with positive PDI values indicating reliable changes in phase  
134 coherence elicited by the visual component of the stimulus.

135 We calculated PDI values for each of the four single stream stimuli and grouped conditions by  
136 coherency (coherent: A1V1 / A2V2, or independent: A1V2 / A2V1). To determine at what frequencies  
137 the across-trial phase reliability was significantly positive, we compared within-stimulus and across-  
138 stimulus PDI for each frequency band (paired t-test with Bonferroni correction for 43 frequencies,  $\alpha =$   
139 0.0012). In awake subjects we identified a restricted range of frequencies between 10 and 20 Hz where  
140 visual stimuli enhanced the phase reliability (Fig. 4a, b). In anaesthetised animals, average PDI values  
141 were larger than in awake animals and all frequencies tested had single stream PDI values that were  
142 significantly positive (Fig. 4d, e). We therefore conclude that visual stimulation elicited reliable  
143 changes in the LFP phase in auditory cortex. In contrast to LFP phase, a parallel analysis of across trial  
144 power reliability showed no significant effect of visual stimuli on LFP power in any frequency band  
145 (Supplementary Fig. 1).

146 Next we asked whether there were any frequencies at which phase coherence was increased by AV  
147 temporal coherence by performing a pairwise comparison of single stream PDI values, yielded from  
148 temporally coherent and independent stimuli, for all frequency points. In anaesthetised animals, the  
149 single stream PDI did not differ between coherent and independent stimuli at any frequency (Fig. 4f).  
150 In awake animals, PDI values were similar for temporally coherent and temporally independent

151 stimuli, except in the 11-14 Hz band where coherent stimuli elicited significantly greater phase  
152 coherence (Fig. 4c). Together these data suggests that visual inputs modulate the phase of the low  
153 frequency field potential in auditory cortex independently of temporal coherence with auditory  
154 stimuli, and are consistent with auditory cortical neurons integrating visual and auditory information  
155 such that discriminability of spiking responses to temporally coherent auditory visual signals are  
156 enhanced (Fig.2e, f).

157 **Visual information enhances the representation of the temporally coherent auditory stream in a**  
158 **sound mixture**

159 Arguably the greatest challenge for the auditory brain is to reconstruct sound sources in the world  
160 from their overlapping cochlear representations. Having demonstrated that temporal coherence  
161 between auditory and visual stimuli enhances discriminability of auditory spiking responses, we next  
162 asked whether the temporal dynamics of a visual stimulus could enhance the representation of one  
163 sound in a mixture. We therefore recorded responses to auditory scenes composed of two sounds  
164 (A1 and A2) presented simultaneously with a visual stimulus that was temporally coherent with one  
165 or other auditory stream (A12V1 or A12V2). To test if a visual stimulus could enhance the  
166 representation of the temporally coherent auditory stream in such dual stream stimuli we then  
167 compared dual stream responses with responses to temporally coherent single stream stimuli.

168 Figure 5 illustrates this approach for a single unit: responses to the single stream AV stimuli (Fig. 5a)  
169 formed templates against which we judged the similarity of responses to the dual stream stimuli (Fig.  
170 5b). Responses to the dual stream stimuli more closely resembled A1V1 when the visual stimulus was  
171 V1, and A2V2 when the visual stimulus was V2. In our analysis of single-stream encoding, this unit was  
172 classified as visual-discriminating, but many auditory-discriminating units showed similar response  
173 properties (e.g. Supplementary Fig. 2). Enhancement of the coherent auditory stimulus representation  
174 was visible at the population level (Fig. 5c-f): Auditory cortical responses to dual-stream stimuli most  
175 closely resembled responses to the single stream stimulus with the same visual component. This

176 finding was robust in both awake (Fig. 5d, pairwise t-test:  $t_{540} = 6.073$ ,  $p < 0.001$ ) and anaesthetised  
177 animals (Fig. 5f,  $t_{660} = 9.5137$ ,  $p < 0.001$ ) suggesting that these effects were not mediated by attention.

178 Modulation of dual stream responses by visual stimulus identity was not simply a consequence of the  
179 shared visual component of single stream and dual stream stimuli (Fig. 6). To show this we decoded  
180 responses to dual stream stimuli (A12V1 and A12V2) using responses to auditory-only stimuli (Fig. 6a;  
181 A1 or A2). We also analysed responses to mixed auditory streams with no visual stimulus (A12) using  
182 responses either to coherent single stream stimuli (A1V1, A2V2). A two-way repeated measures  
183 ANOVA on the decoder responses with factors of visual stream (V1, V2, no visual), and template type  
184 (AV or A) demonstrated a significant effect of visual stream identity on dual stream decoding (Fig. 6d,  
185  $F(2, 528) = 19.320$ ,  $p < 0.001$ ), but there was no effect of template type ( $F(1, 528) = 0.073$ ,  $p = 0.787$ )  
186 or interaction between factors ( $F(2, 528) = 0.599$ ,  $p = 0.550$ ). Post-hoc comparison across units  
187 revealed that without visual stimulation there was no tendency to respond preferentially to either  
188 stream but that visual stream identity significantly influenced classification of dual stream responses.

189 The ability of a visual stimulus to modulate auditory representation in the dual stream condition was  
190 observed across all cortical layers (defined by current source density analysis, see online methods,  
191 Supplemental Fig. 3a,b) and across three tonotopic fields (Supplemental Fig. 3c,d) of anaesthetized  
192 subjects. While present in all effects, the influence of the visual stimulus was strongest in the supra-  
193 granular layers (supplemental Fig.3b). Thus cross-modal modulation by temporal coherence was a  
194 general phenomenon across auditory cortex. Separating auditory, visual, and auditory-visual units  
195 according to responsiveness to classic neurophysiological stimuli (noise bursts and light flashes)  
196 revealed that the impact of a visual stimulus on the representation of a sound mixture was present  
197 across functional sub-populations but strongest in visual and auditory-visual units (Supplemental Fig.  
198 4). Finally, we observed these effects in both single and multi-units (Supplemental Fig.5).

199 **Visual stimuli elicit changes in LFP phase in the context of an auditory scene**

200 Our findings indicate that visual stimuli can shape the representation of auditory mixtures and that  
201 temporal coherence between auditory and visual stimuli enhances across-trial phase coherence. To  
202 understand whether changes in phase coherence could provide a mechanism for visual modulation of  
203 auditory representations, we again generated within-stimulus ITPC for each dual-stream stimulus (Fig.  
204 7a, A12V1 and A12V2) and across-stimulus ITPC by randomly selecting responses across visual  
205 conditions (Fig. 7b). We then expressed the difference as the dual stream PDI (dual stream phase  
206 dissimilarity index, Fig. 7c). Since the auditory components were identical in each dual stream  
207 stimulus, the influence of the visual component on LFP phase could be isolated as non-zero dual  
208 stream PDI values (paired t-test, bonferoni corrected,  $\alpha = 0.0012$ ). In awake animals, dual stream PDI  
209 was significantly greater than zero at 11-14Hz and 16-19 Hz (Fig. 7d, e) whereas in anaesthetised  
210 animals, we found positive dual stream PDI values across all frequencies tested (Fig. 7f, g). In  
211 anaesthetised animals where we were able to use the responses of units to noise and light flashes to  
212 categorise units as auditory, visual or auditory-visual, we also confirmed significant PDI values in each  
213 of these subpopulations of units (Supplemental Fig. 4c). In awake animals, we confirmed the  
214 significance of PDI values in the 11-14 Hz range across different rates of amplitude modulation of the  
215 auditory stimulus (Supplemental Fig. 6).

### 216 **Neural responses to auditory timbre deviants are enhanced when changes in visual luminance and** 217 **auditory intensity are temporally coherent**

218 A hall-mark of an object-based rather than feature-based representation is that all stimulus features  
219 are bound into a unitary perceptual construct, including those features which do not directly mediate  
220 binding<sup>22</sup>. We predicted that binding across modalities would be promoted via synchronous changes  
221 in auditory intensity and visual luminance (Fig. 1b) and observed that the temporal dynamics of the  
222 visual stimulus enhanced the representation of temporally coherent auditory streams (Fig. 2e-f and  
223 5d-f). To determine whether temporal synchrony of visual and auditory stimulus components also  
224 enhanced the representation of orthogonal stimulus features (Fig. 1c) and thus fulfil a key prediction

225 of binding, we introduced brief timbre perturbations into our dual stream stimuli ( $n = 4$  deviants, two  
226 in A1 and two in A2). Such deviants could be detected by human listeners and were better detected  
227 when the auditory stream in which they were embedded was temporally coherent with an  
228 accompanying visual stimulus<sup>5</sup>. We hypothesised that, despite containing no information about the  
229 occurrence of deviants, a temporally coherent visual stimulus would enhance the representation of  
230 changes in timbre in the responses of auditory cortical neurons.

231 To isolate neural responses to the timbre change from those elicited by the on-going amplitude  
232 modulation, we extracted the 200ms epochs of the neuronal response during which the timbre  
233 deviant occurred and compared these to epochs from responses to otherwise identical stimuli without  
234 deviants. We observed that the spiking activity of many units differed between deviant and no-deviant  
235 trials (e.g. Fig 8a) and so we used a pattern-classifier approach to estimate the presence/absence of a  
236 timbre deviant in a given response window. We first considered the influence of temporal coherence  
237 between auditory and visual stimuli on the representation of timbre deviants in the single stream  
238 condition (A1V1, A1V2 etc.). We found that a greater proportion of units detected at least one deviant  
239 when the auditory stream in which deviants occurred was temporally coherent with the visual  
240 stimulus relative to the temporally independent condition. This was true both for awake (Fig. 8b;  
241 Pearson chi-square statistic,  $\chi^2 = 322.617$ ,  $p < 0.001$ ) and anaesthetised animals (Fig. 8e;  $\chi^2 = 288.731$ ,  $p$   
242  $< 0.001$ ). For units that discriminated at least one deviant, discrimination scores were significantly  
243 higher when accompanied by a temporally coherent visual stimulus (Fig.8c, awake dataset, pairwise  
244 t-test  $t_{300} = 3.599$   $p < 0.001$ ; Fig. 8f, anaesthetised data  $t_{262} = 4.444$   $p < 0.001$ ).

245 Across the population of units, we performed a two-way repeated measures ANOVA on discrimination  
246 performance with visual condition (V1/V2) and the auditory stream in which the deviants occurred  
247 (A1/A2) as factors. We predicted that enhancement of the representation of timbre deviants in the  
248 temporally coherent auditory stream would be revealed as a significant interaction term. Significant  
249 interactions were seen in both the awake (Fig. 8d,  $F(1, 600) = 29.138$ ,  $p < 0.001$ ) and anaesthetised

250 datasets (Fig. 8g,  $F(1, 524) = 16.652$ ,  $p < 0.001$ ). We also observed significant main effects of auditory  
251 and visual conditions in awake (main effect of auditory stream,  $F(1, 600) = 4.565$ ,  $p = 0.033$ ; main effect  
252 of visual condition,  $F(1, 600) = 2.650$ ,  $p = 0.010$ ) but not anaesthetised animals (main effect of auditory  
253 stream,  $F(1, 524) = 0.004$ ,  $p = 0.948$ ; main effect of visual condition,  $F(1, 524) = 1.355$ ,  $p = 0.245$ ). Thus  
254 we concluded a temporally coherent visual stimulus can enhance the representation of features (here  
255 auditory timbre) orthogonal to those that promote binding between auditory and visual streams. This  
256 finding is consistent with our model of cross-modal binding (Fig. 1a, c) and so these data fulfil our  
257 definition of binding.

## 258 Discussion

259 Here we provide mechanistic insight into how auditory and visual information could be bound  
260 together to form coherent perceptual objects. Visual stimuli elicit reliable changes in the phase of the  
261 local field potential in auditory cortex that result in an enhanced spiking representation of auditory  
262 information. These results are consistent with the binding of cross-modal information to form a  
263 multisensory object. When two sounds are presented together within an auditory scene, the  
264 representation of the stream that is temporally coherent with the visual stimulus is enhanced.  
265 Importantly, this enhancement is not restricted to the encoding of the amplitude changes that bind  
266 auditory and visual information but extends to the encoding of auditory timbre, a stimulus dimension  
267 orthogonal to the dimensions that link auditory and visual stimuli. Thus our results meet the  
268 requirements for a strict neural test of cross-modal binding that was laid out in Bizley et al.<sup>6</sup>. These  
269 data provide a physiological underpinning for the pattern of performance observed in human listeners  
270 performing an auditory selective attention task in which detection of a perturbation in a stimulus  
271 stream is enhanced or suppressed, when a visual stimulus is temporally coherent with the target or  
272 masker auditory stream respectively<sup>5</sup>. The electrophysiological data presented here suggest that the  
273 temporally coherent auditory stream would be represented more effectively, making the task easier  
274 when this stream was the target and making the task more challenging when it was the masker.

275 Surprisingly, the effects of the visual stimulus on the representation of an auditory scene can be  
276 observed in anesthetised animals ruling out any top-down effect of attentional modulation.

277 Previous investigations of the impact of visual stimuli on auditory scene analysis have frequently used  
278 speech stimuli. In order to probe more general principles that might relate to both speech and non-  
279 speech processing we chose to employ non-speech stimuli, but utilized modulation rates that fell  
280 within the range of syllable rates in human speech<sup>23</sup>. Previous work has demonstrated that a visual  
281 stimulus can enhance the neural representation of the speech envelope both in quiet and in noise<sup>3</sup>,  
282<sup>24,25</sup>. Being able to see a talker's mouth provides listeners with rhythm information and information  
283 about the amplitude of the speech waveform which may help listeners by cueing them to pay  
284 attention to the auditory envelope<sup>26</sup> as well as information about the place of articulation that can  
285 disambiguate different consonants<sup>27</sup>. Visual speech information is hypothesised to be relayed in  
286 parallel to influence the processing of auditory speech: Our data support the idea that early  
287 integration of visual information occurs<sup>26,28-30</sup> and is likely to reflect a general phenomenon whereby  
288 visual stimuli can cause phase-entrainment in the local field potential. Our data support the  
289 contention that such early integration is unlikely to be specific to speech. Indeed low-frequency  
290 entrainment to modulations in an on-going stimulus are observed in the human brain and have been  
291 shown to optimize listening performance for non-speech stimuli<sup>31</sup>. In contrast, later integration is  
292 likely to underlie information about speech gestures that might be used to constrain lexical identity<sup>26</sup>.

293 Consistent with previous studies, our analysis of local field potential activity revealed that visual  
294 information reliably modulated the phase of oscillatory activity in auditory cortex independently of  
295 the modulation frequency of the stimulus<sup>8-11</sup>. Neuronal excitability varies with LFP phase<sup>32-35</sup> and may  
296 be the physiological mechanism through which cross-sensory information is integrated. Our analysis  
297 allowed us to isolate changes in LFP phase that were directly attributable to the visual stimulus and  
298 identified reliable changes in the LFP phase irrespective of whether the visual stimulus was temporally  
299 coherent with the auditory stimulus. Such a finding is consistent with the idea that the LFP phase

300 synchronization arises from fluctuating inputs to cortical networks<sup>14,21,36</sup>. Our finding that visual  
301 stimulation elicited reliable phase modulation in both awake and anaesthetised animals suggests that  
302 bottom-up cross-modal integration interacts with selective attention, which also modulates phase  
303 information in auditory cortex<sup>20</sup>. While our data suggest that cross-modal binding can occur in the  
304 absence of attention, it is likely that the effects we observe in auditory cortex are the substrates on  
305 which selective attention acts to further boost the representation of cross-modal objects.

306 In the awake animal the impact of visual stimulation on LFP phase reliability was smaller than in the  
307 anaesthetised animal and was restricted to a narrower range of frequencies, consistent with a  
308 dependence of oscillatory activity on behavioural state<sup>37-39</sup>. Since the neural correlates of multisensory  
309 binding are evident in the anaesthetised animal, the specific increase in alpha phase reliability that  
310 occurred in awake animals in response to temporally coherent auditory-visual stimulus pairs (Fig. 4c  
311 & 7e) may indicate an attention-related signal triggered by temporal coherence between auditory and  
312 visual signals. Phase resetting or synchronisation of alpha phase has been associated both with  
313 enhanced functional connectivity<sup>38</sup> and as a top-down predictive signal for upcoming visual  
314 information<sup>40</sup>. Disambiguating these possibilities would require simultaneous recordings in auditory  
315 and visual cortex and/or recording during the performance of a task designed to explicitly manipulate  
316 attention.

317 Temporal coherence between sound elements has been proposed as a fundamental organising  
318 principle for auditory cortex<sup>41, 42</sup> and here we extend this principle to the formation of cross-modal  
319 constructs. Our data provide evidence that one role for the early integration of visual information into  
320 auditory cortex is to resolve competition between multiple sound sources within an auditory scene.  
321 While previous studies have demonstrated a role for visual information in conveying lip movement  
322 information to auditory cortex<sup>3, 8, 9, 20</sup>, here we suggest a more general phenomenon whereby visual  
323 temporal cues facilitate auditory scene analysis through the formation of cross-sensory objects. The  
324 origin of the visual inputs is an open question but both visual cortical and sub-cortical structures

325 innervate tonotopic auditory cortex<sup>7, 43</sup>. Identifying which of these inputs is responsible for the  
326 physiological effects we observe requires experiments that manipulate defined neural circuits.

327 In summary, activity in auditory cortex was reliably affected by visual stimulation in a manner that  
328 enhanced the representation of temporally coherent auditory information. Enhancement of auditory  
329 information was observed for sounds presented alone or in a mixture and for sound features that  
330 were related to (amplitude) and orthogonal to (timbre) variation in visual input. Such processes  
331 provide mechanistic support for a coherence based model of cross-modal binding in object formation.

332

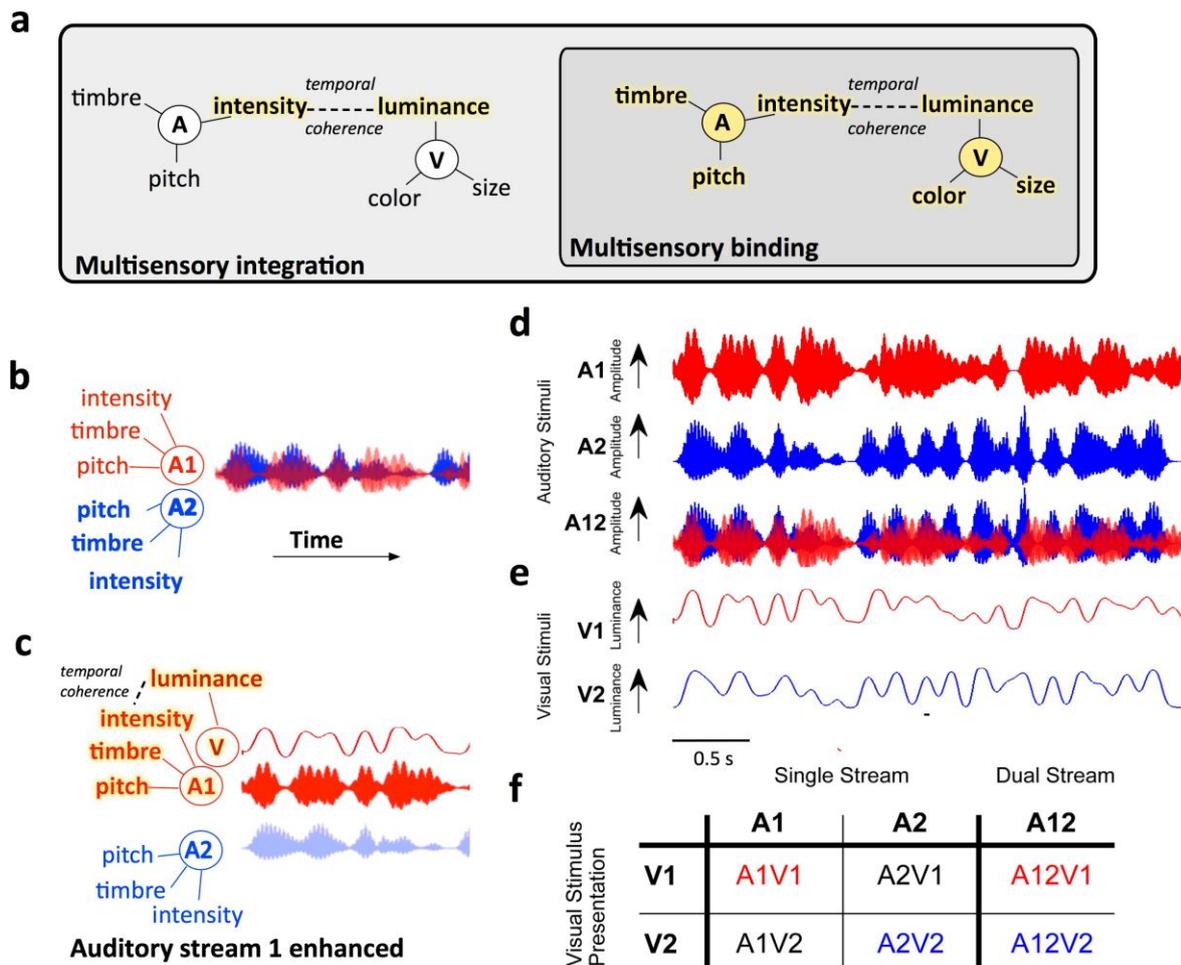
333 [Acknowledgments](#)

334 This work was funded by grants to each author: JKB: Wellcome Trust / Royal Society WT098418MA;  
335 Biotechnology and Biological Sciences Research Council (BB/H016813/1), and an Action on Hearing  
336 Loss Studentship (596: UEI: JB); RKM: NIH K99DC014288 and Hearing Health Foundation Emerging  
337 Research Grant; AKCL: NIH R01DC013260; and an International Exchanges Scheme award from the  
338 Royal Society to JKB and AKCL.

339 [Author contributions](#)

340 HA, RKM, AKCL, JKB Conception and design, HA, SMT, KCW, GPJ, JKB Acquisition of data, HA, JKB  
341 Analysis and interpretation of data, HA, SMT, RKM, AKCL, JKB Drafting or revising the article.

342 Figures

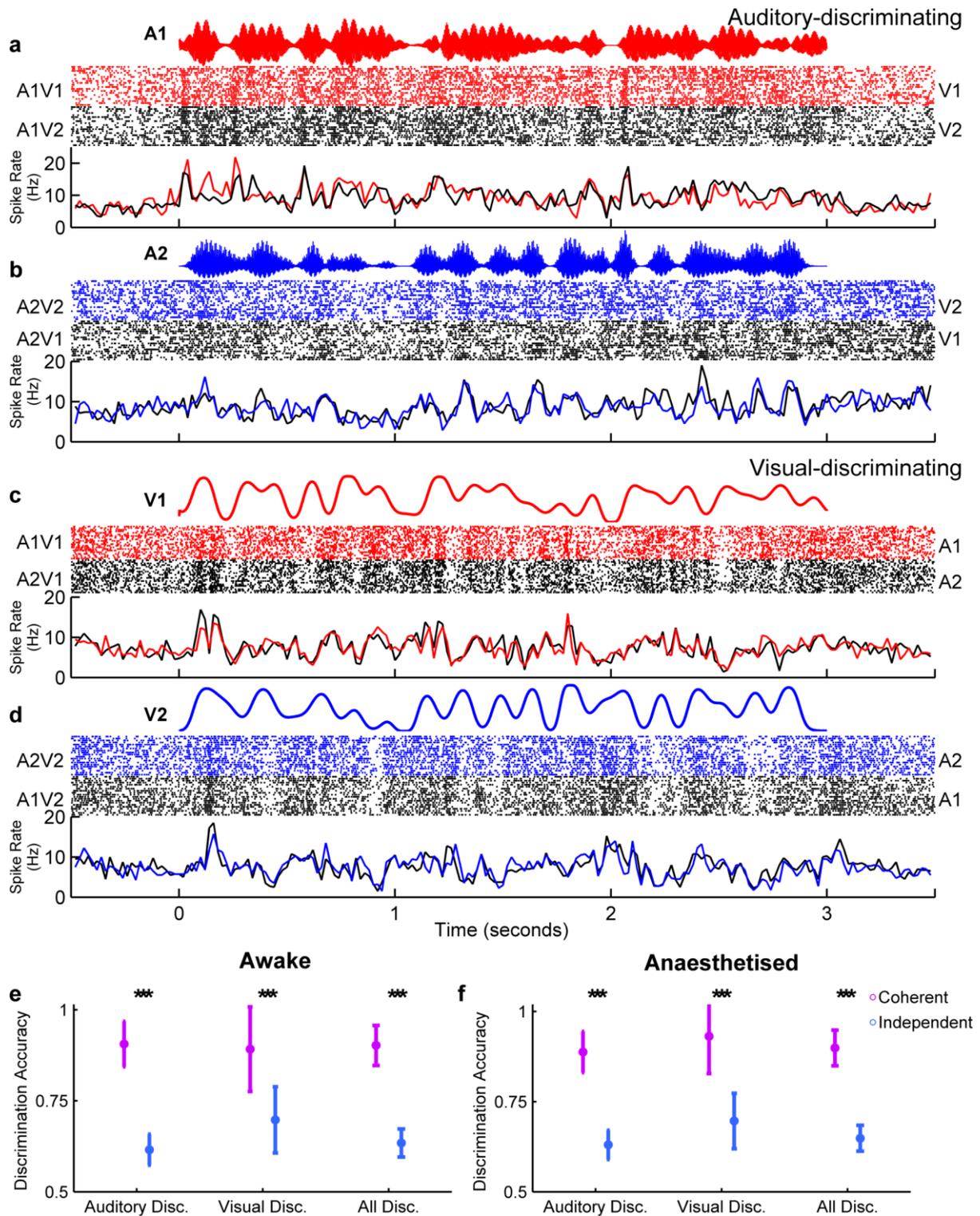


343

344 **Figure 1: Hypothesis and experimental design**

345 **a** Conceptual model illustrating how binding can be identified as a distinct form of multisensory  
 346 integration. Multisensory binding is defined as a subset of multisensory integration which results in  
 347 the formation of a crossmodal object. During binding, all features of the audio-visual object are linked  
 348 and enhanced including both those features that bind the stimuli across modalities (here temporal  
 349 coherence between auditory (A) intensity and visual (V) luminance) and orthogonal features such as  
 350 auditory pitch and timbre- and visual color and size. Other forms of multisensory integration would  
 351 result in only the features that promote binding being enhanced - here auditory intensity and visual  
 352 luminance. To identify binding therefore requires a demonstration that non-binding features (e.g.  
 353 here pitch, timbre, color or size) are enhanced **b** When two competing sounds (red and blue

354 waveforms) are presented they can be separated on the basis of their features, but may elicit  
355 overlapping neuronal representations in auditory cortex. **c** Hypothesised enhancement in auditory  
356 stream segregation when a temporally coherent visual stimulus enables multisensory binding. When  
357 the visual stimulus changes coherently with the red sound (A1, top) this sound is enhanced and the  
358 two sources are better segregated. Perceptually this would result in enhanced auditory scene analysis  
359 and an enhancement of the non-binding features. **d** Stimuli design: Auditory stimuli were two artificial  
360 vowels (denoted A1 and A2), each with distinct pitch and timbre and independently amplitude  
361 modulated with a noisy low pass envelope. **e** Visual stimulus: a luminance modulated white light was  
362 presented with one of two temporal envelopes derived from the amplitude modulations of A1 and  
363 A2. **f** illustrates the stimulus combinations that were tested experimentally in *single stream* (a single  
364 auditory visual pair) and *dual stream* (two sounds and one visual stimulus) conditions.



365

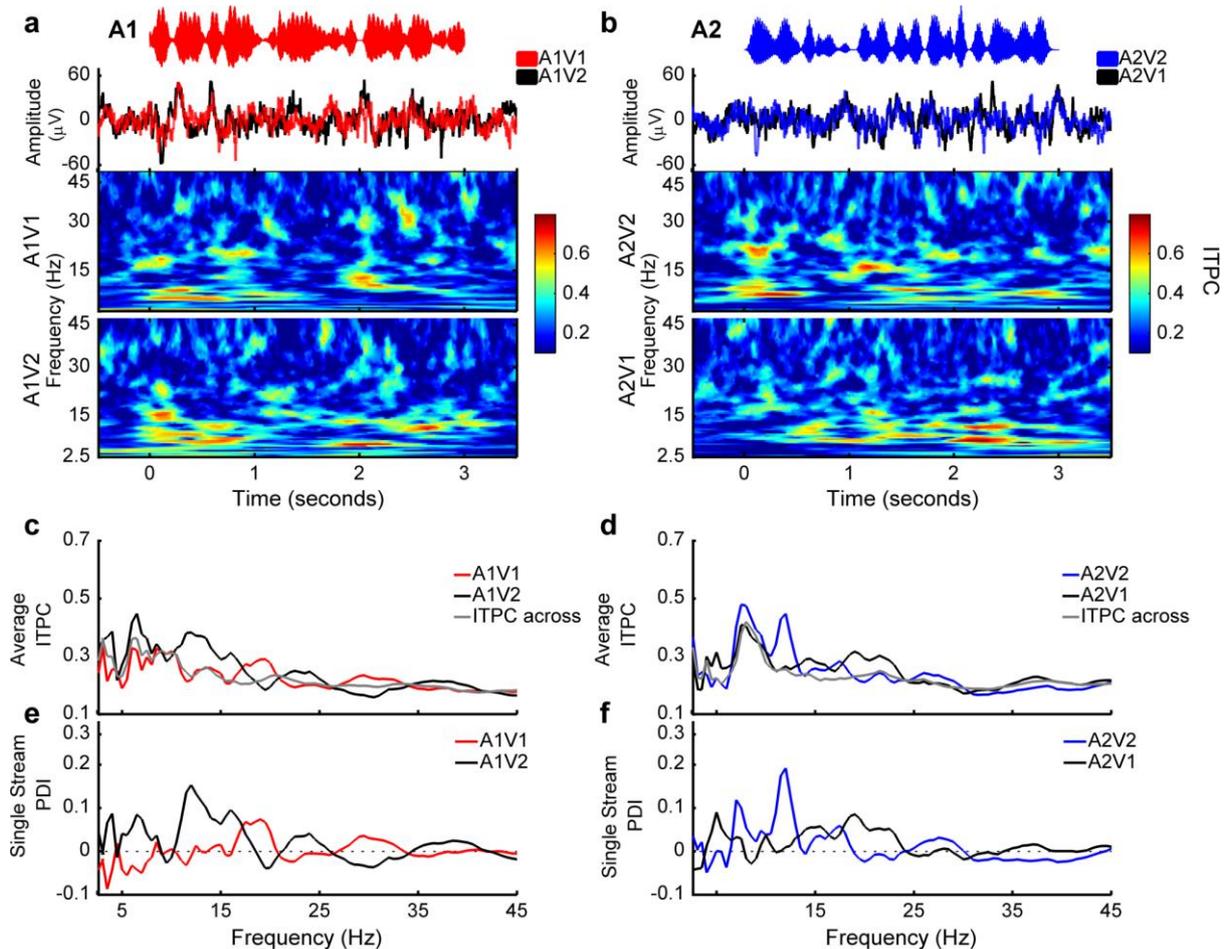
366 **Figure 2: Auditory-visual temporal coherence enhances neural coding in auditory cortex.**

367 A pattern classifier was used to determine whether neuronal responses were informative about

368 auditory or visual stimuli. The responses to single stream stimuli are shown for two example units,

369 with responses grouped according to the identity of the auditory (**a**, **b**, auditory discriminating unit)

370 or visual stream (**c, d**, visual discriminating unit). In each case the stimulus amplitude (a,b) / luminance  
371 (c,d) waveform is shown in the top panel with the resulting raster plots (20 trials per condition) and  
372 peri-stimulus time histogram (20 ms bin) below. **e, f**: Decoder performance (mean  $\pm$  SEM) for  
373 discriminating stimulus identity (coherent: A1V1 vs. A2V2; independent: A1V2 vs. A2V1) in auditory  
374 and visual classified units recorded in awake (e) and anaesthetised (f) ferrets. Pairwise comparisons  
375 for decoding of coherent versus independent stimuli:  $p < 0.001$  (\*\*\*, see results).

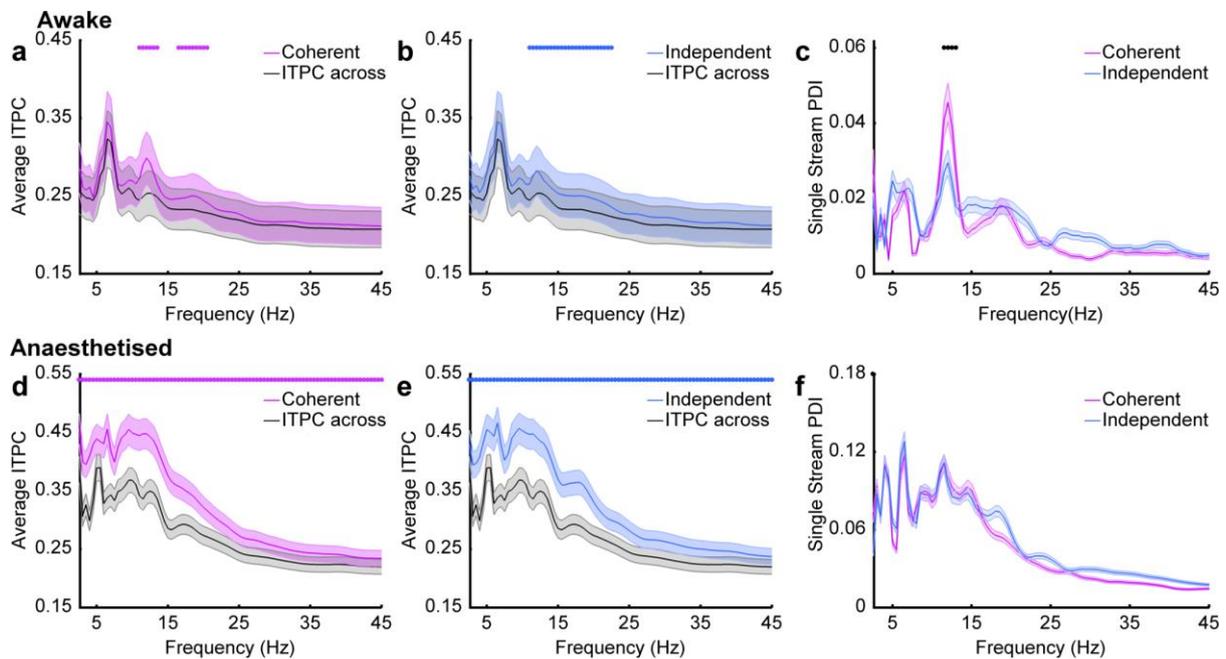


376

377 **Figure 3: Visual stimuli elicit reliable changes in the phase of the local field potential**

378 **a, b** Example LFP responses to single stream stimuli across visual conditions. Data obtained from the  
379 recording site at which multiunit spiking activity discriminated auditory stream identity in Fig. 2 a and  
380 b. The amplitude waveforms of the stimuli are shown in the top panel, with the evoked LFP  
381 underneath (mean across 21 trials). The resulting inter-trial phase coherence (ITPC) values are shown

382 in the bottom two panels. **c, d** ITPC was calculated for coherent and independent AV stimuli separately  
383 and compared to a null distribution (ITPC across). Single stream phase dissimilarity values (PDI) were  
384 calculated by comparing ITPC values to the ITPC across condition (**e, f**).



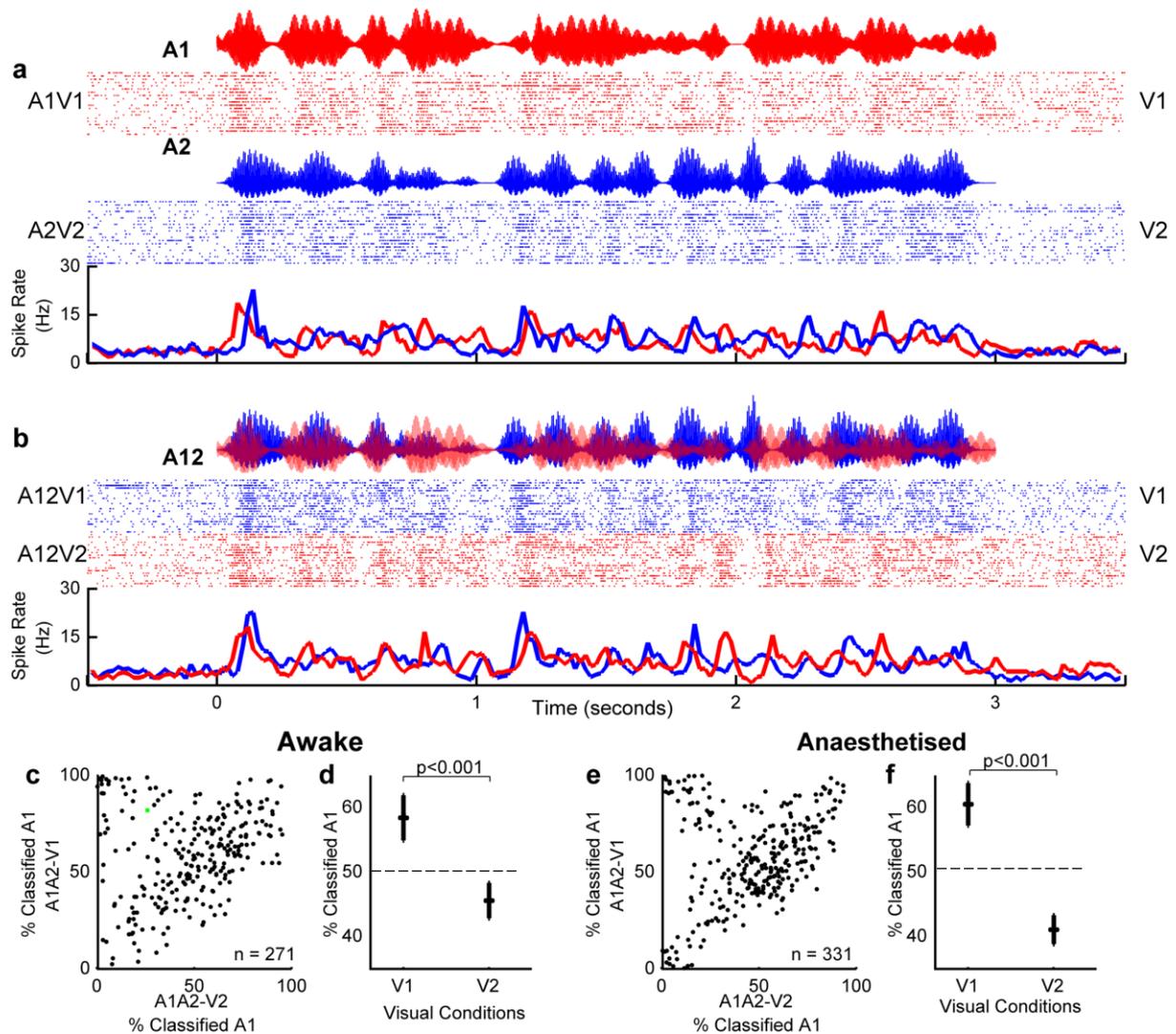
385

386 **Figure 4: Visual stimuli elicit reliable changes in LFP phase in awake and anaesthetised animals.**

387 Mean inter-trial phase coherence (ITPC) values across frequency for coherent (**a, d**) and independent  
388 (**b, e**) conditions. Dots indicate frequencies at which the ITPC values were significantly greater than  
389 chance (Pairwise ttest,  $\alpha = 0.0012$ , Bonferroni corrected for 43 frequencies). **c f**: Mean ( $\pm$ SEM) single  
390 stream phase dissimilarity index (PDI) values for coherent and independent stimuli in awake (**c**) and  
391 anaesthetised (**f**) animals. Black dots indicate frequencies at which the coherent stream PDI is  
392 significantly greater than in the independent conditions ( $p < 0.001$ ).

393

394

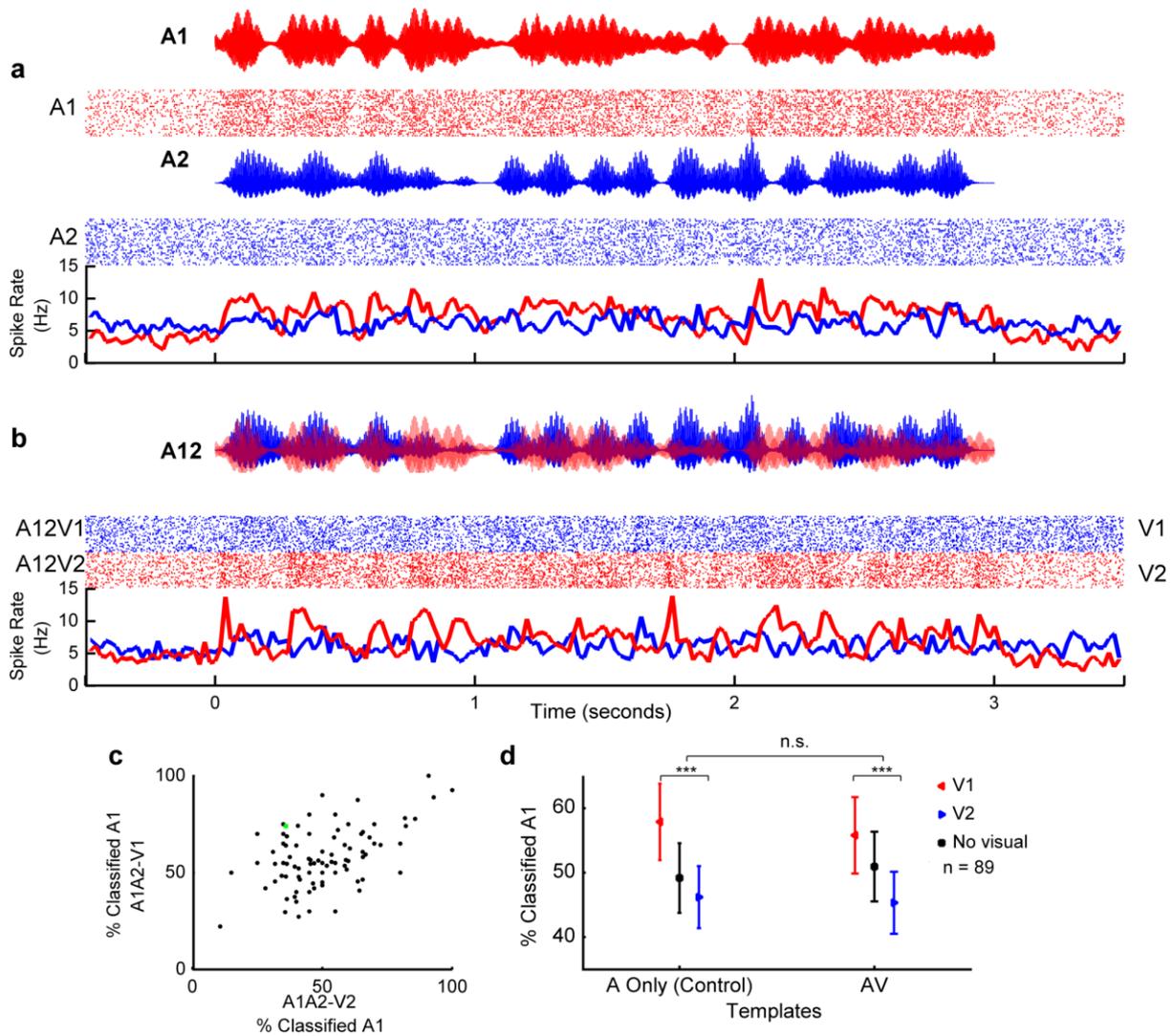


395

396 **Figure 5: Visual stimuli can determine which sound stream auditory cortical neurons follow in a**  
397 **mixture.**

398 Spiking responses from an example unit (visual classified unit from the awake dataset) in response to  
399 **a**, single stream AV stimuli used as decoding templates and **b**, dual stream stimuli, rasters and PSTH  
400 responses stimuli from which single trial responses were classified. When the visual component of the  
401 dual stream was V1, the majority of trials were classified as A1V1 (82% (19/23) of trials), and A2V2  
402 when the visual stimulus was V2 (26% (6/23) of responses classified as A1V1 (see also green datapoint  
403 in c).). **c-f** data for awake (c,d) and anaesthetised (e,f) datasets. In each case the left panel (c,e) shows  
404 the distribution of decoding values according to the visual condition and the right panel (d,f) shows

405 the population mean ( $\pm$  SEM) Pairwise comparisons revealed significant effect of visual condition on  
 406 decoding in both datasets ( $p < 0.001$ ).

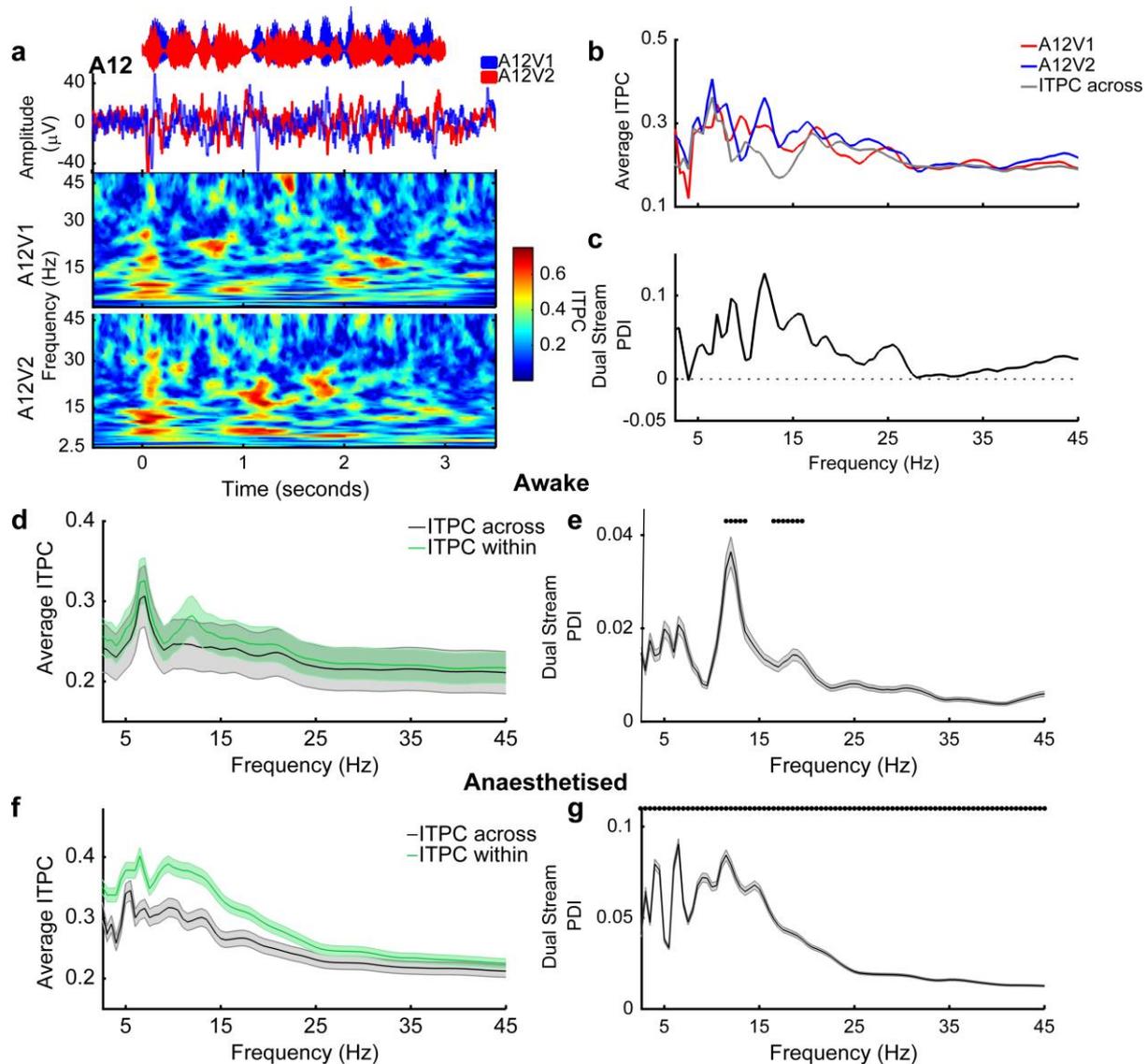


407

408 **Figure 6: Visual stimuli shape the neural representation of an auditory scene.**

409 In an additional control experiment ( $n=89$  units recorded in awake animals), the responses to coherent  
 410 AV and auditory-only (A Only) single stream stimuli were used as templates to decode dual stream  
 411 stimuli either accompanied by visual stimuli (V1/V2) or in the absence of visual stimulation (no visual).  
 412 Spiking responses from an example unit in response to **a**, single stream stimuli (A only, no visual) used  
 413 as decoding templates and **b**, dual stream stimuli, in each case the auditory waveform, rasters and  
 414 PSTHs are shown. A two-way ANOVA with template type and visual condition as factors revealed  
 415 significant effects of visual condition but not template type. Post-hoc comparisons demonstrated that

416 classification was significantly influenced by visual stimulus identity when both A only and AV  
 417 templates were used. The proportion of responses classified as A1 when the visual stimulus was V1 or  
 418 V2 are shown in **c** and **d**, Mean ( $\pm$  SEM) values for these units. Pairwise comparisons revealed  
 419 significant effect across visual conditions in both datasets ( $p < 0.001$ ).

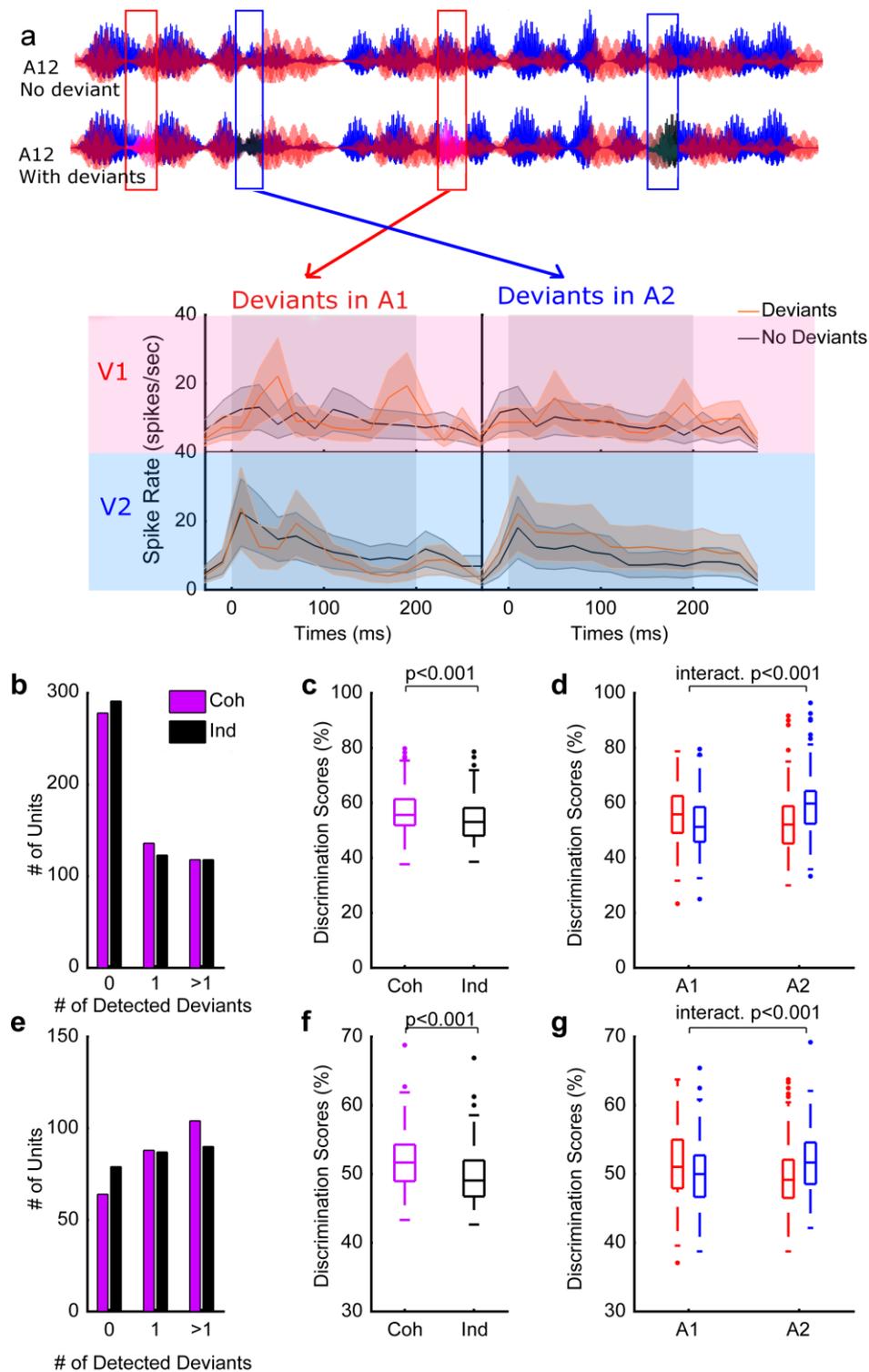


420

421 **Figure 7: Visual stimuli elicits reliable changes in LFP phase to shape auditory scene analysis**

422 **a**, Auditory stimulus waveform (top), evoked LFP (middle) and resulting inter-trial phase coherency  
 423 plots (bottom) for a typical recording site in response to dual stream stimuli. **b**, ITPC values were  
 424 calculated across frequency for responses to dual stream stimuli with identical visual stimuli ('ITPC  
 425 within') and across randomly drawn visual stimuli ('ITPC-across', grey). Dual stream phase selectivity

426 index (PDI) values were then calculated as the difference between shuffled and dual stream ITPC (**c**).  
427 **d, f** Average ITPC ( $\pm$  SEM) for dual stream ITPC-within and ITPC-across for awake and anaesthetised  
428 animals. Symbols indicate where the dual stream PDI was significant (pairwise t-test for ITPC within  
429 versus ITPC across,  $\alpha = 0.0012$  with correction). **e, g** Mean ( $\pm$  SEM) dual stream PDI values for awake  
430 and anaesthetised animals.



432 **Figure 8: Temporally coherent changes in visual luminance and auditory intensity enhance the**  
 433 **coding of a non-binding auditory feature.**

434 **a** Example unit (from the awake dataset) showing the influence of visual temporal coherence on  
 435 spiking responses to dual stream stimuli with or without deviants embedded. Shaded rectangles

436 indicate the 200 ms window over which the timbre deviant occurred and over which analysis was  
437 conducted. **b-d** timbre deviant discrimination in the awake dataset. Two deviants were included in  
438 each auditory stream giving a possible maximum of 4 per unit **b**, Histogram showing the number of  
439 deviants (out of 4) that could be discriminated from spiking responses **c**, Box plots showing the  
440 average timbre deviant discrimination scores in the single stream condition across different visual  
441 conditions (Coh: coherent, ind: independent). The boxes show the upper and lower quartile values,  
442 and the horizontal lines at their “waist” indicate the median. **d**, Discrimination scores for timbre  
443 deviant detection in dual stream stimuli. Discrimination scores are plotted according to the auditory  
444 stream in which the deviant occurred and the visual stream that accompanied the sound mixture. **e-g**  
445 show the same as **b-d** but for the anesthetised dataset.

446

447

448

449

450

451

452

453

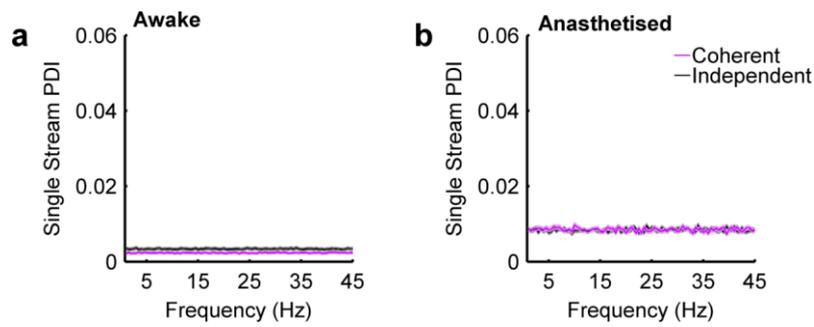
454

455

456

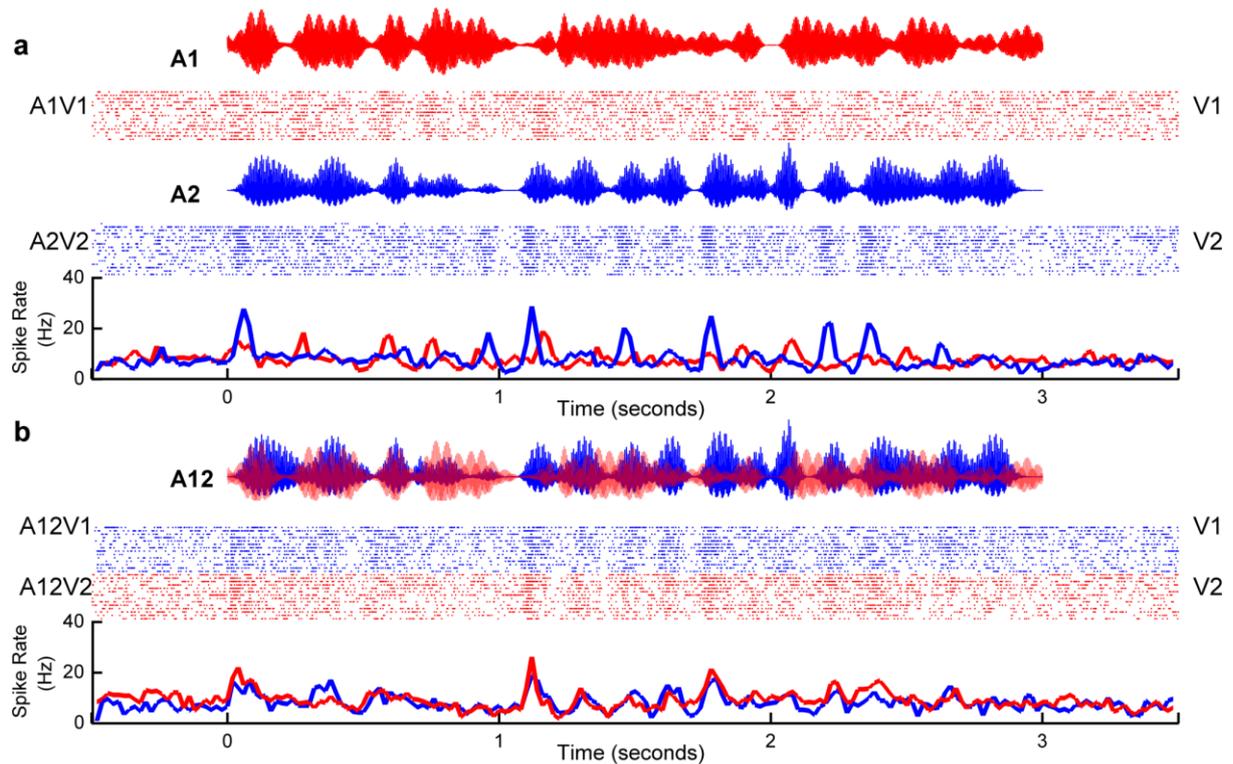
457

458 **Supplementary Figures**



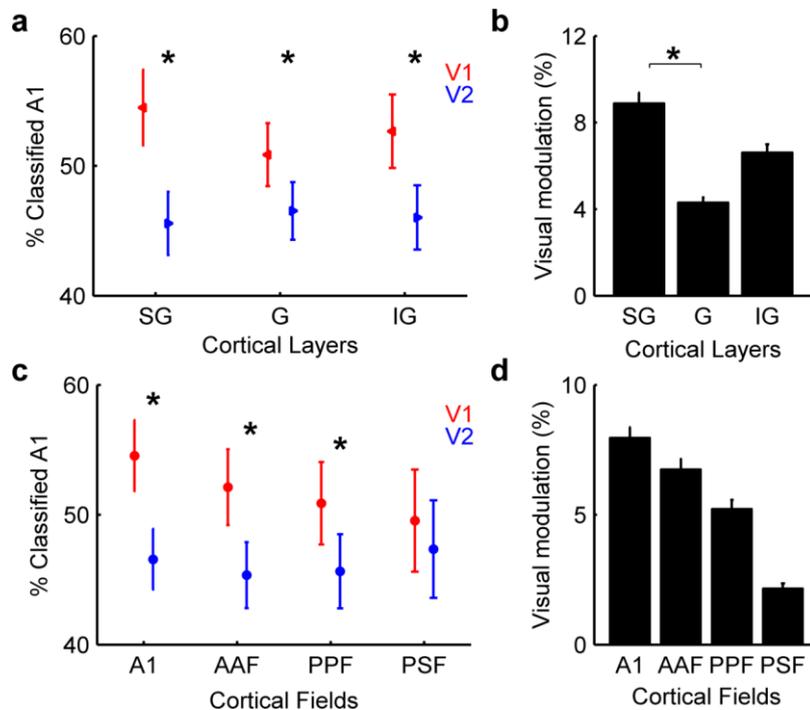
460 **Supplementary figure 1: Dynamic visual stimuli do not elicit reliable changes in LFP power.**

461 **a, b:** Mean ( $\pm$ SEM) single stream phase dissimilarity index values for coherent (A1V1 and A2V2) and  
462 independent (A1V2 and A2V1) stimuli at all recording sites at which there was a driven spiking  
463 response for awake (a) and anaesthetised (b) animals.



465 **Supplementary figure 2: Visual stimuli can determine which sound stream auditory cortical neurons**  
466 **follow in a mixture in auditory-discriminating units**

467 The spiking responses of an example unit are shown to coherent single stream (a) and dual stream  
 468 stimuli (b). This example unit was an auditory discriminating unit recorded from an awake animal.  
 469 When the visual stimulus was V1, 67% of trials were classified as A1V1, when the visual stimulus was  
 470 V2 only 29% of trials were classified as A1V1.



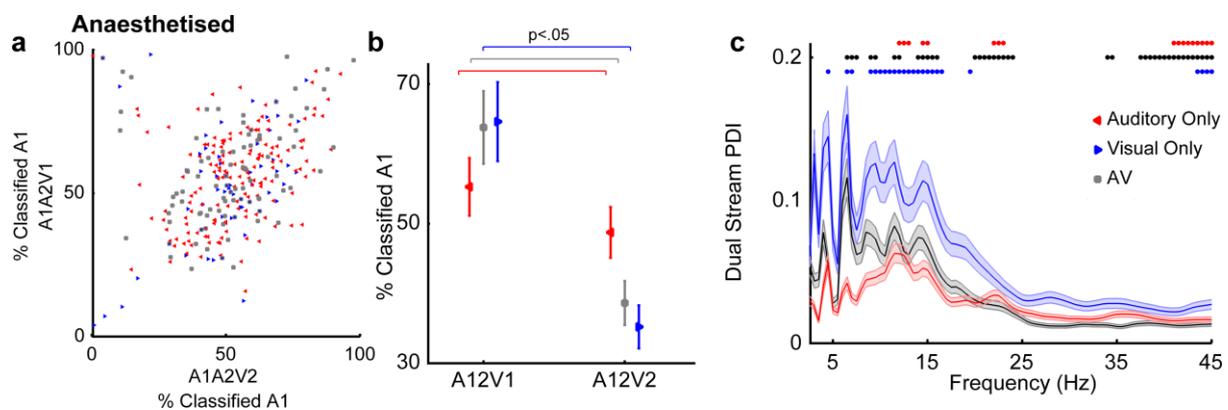
471

472 **Supplementary figure 3: Preference to the auditory stream with which the visual stimulus was**  
 473 **coherent observed across all cortical layers and in three tonotopic fields.**

474 In the anesthetised animal, we determined tonotopic gradients from responses to pure tone stimuli  
 475 to accurately confirm the cortical field in which any given recording was made. Recordings were made  
 476 with linear shank electrodes, facilitating current source density analysis to identify the cortical layers.

477 **a** Proportion of responses classified as A1 when the visual stimulus was V1 or V2 are shown for  
 478 supragranular layers (SG,  $t_{700}=5.686$ ,  $p<0.001$ ), granular layer (G,  $t_{878}=3.481$   $p<0.001$ ) and intragranular  
 479 layer (IG,  $t_{690}=4.418$ ,  $p<0.001$ ). Data shown as mean ( $\pm$  SEM) across units. A two-way ANOVA across  
 480 visual condition and cortical layers showed only a significant effect of visual condition ( $F(1, 2273)$   
 481  $=64.288$ ,  $p<0.001$ ) but not of layers ( $F(1, 2273) = 0.91$ ,  $p = 0.404$ ) and no interaction was observed ( $F$   
 482  $(1, 2273) = 2.679$ ,  $p = 0.068$ ). Significant post-hoc comparisons ( $p<0.05$ ) are shown with an asterisks. **b**,

483 In order to compare the magnitude of the effect of the visual stimulus across cortical layers we  
484 calculated visual modulation as the decoding score when the visual stimulus was V1 – the decoder  
485 score when the visual stimulus was V2. Plotted are the mean ( $\pm$  SEM) visual modulation across  
486 different layers. A one-way ANOVA across cortical layers showed a significant effect of cortical layer  
487 ( $F(2, 1134) = 3.15, p = 0.043$ ). Post hoc comparison revealed that visual modulation in SG is greater  
488 than G. c, as for a but for different cortical fields (A1,  $t_{798} = 5.435, p < 0.001$ ; AAF,  $t_{636} = 4.302, p < 0.001$ ;  
489 PPF,  $t_{510} = 3.609, p < 0.001$ ; PSF,  $t_{317} = 0.932, p = 0.352$ ). A two way ANOVA across visual condition and  
490 cortical field showed only a significant effect of visual condition (Visual condition:  $F(1, 2267) = 40.301,$   
491  $p < 0.001$ ; fields:  $F(1, 2267) = 1.937, p = 0.121$ ; visual condition  $\times$  fields:  $F(1, 2267) = 1.804, p = 0.144$ ).  
492 Significant post-hoc comparisons ( $p < 0.05$ ) are shown with an asterisks. In order to compare the  
493 magnitude of the effects across cortical fields we again calculated a visual modulation value. **b**, Mean  
494 ( $\pm$  SEM) visual modulation across different fields. A one-way ANOVA across cortical layers showed a  
495 no significant effect of cortical field ( $F(3, 1130) = 2.19, p = 0.0877$ ).



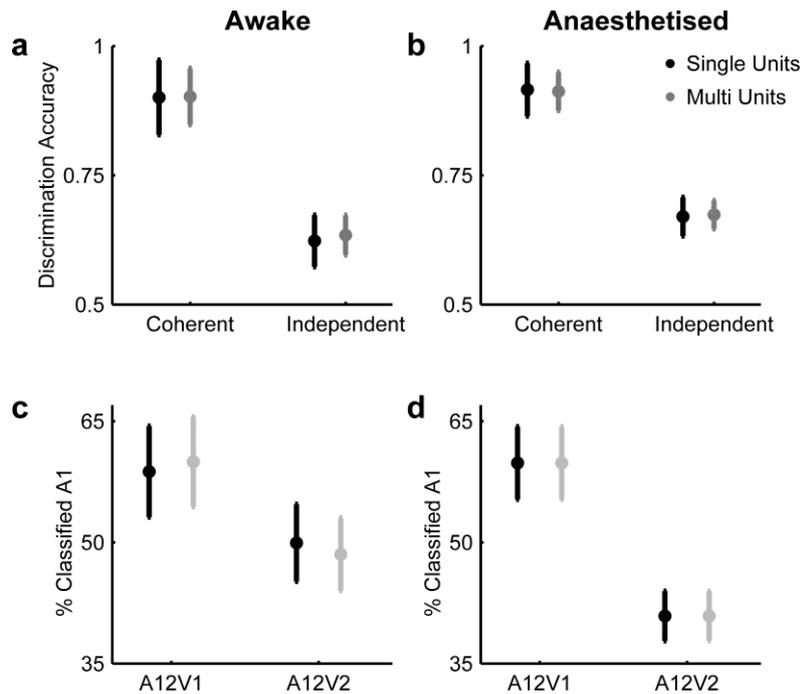
496

497 **Supplementary figure 4: Auditory, visual and auditory-visual units are influenced by visual stimuli.**

498 In anaesthetised animals we also recorded responses to 100 ms noise bursts and/or LED flashes in order  
499 to apply conventional classification of units as either auditory, visual or auditory-visual. Unit  
500 classification was based on a two way ANOVA of spike counts calculated over a 200ms window  
501 following stimulus onset with auditory and visual stimuli as factors. Units in which both auditory and  
502 visual stimuli significantly modulated spiking or in which there was a significant interaction between

503 stimuli were classified as auditory-visual (AV). This analysis yielded recording sites at which units were  
504 classified as auditory only (A; n=177 units), visual only (V, n=130) and AV (grey, n=150). This allowed  
505 us to explore the visual modulation effects according to independently determined functional  
506 definitions of units as auditory, visual or auditory visual.

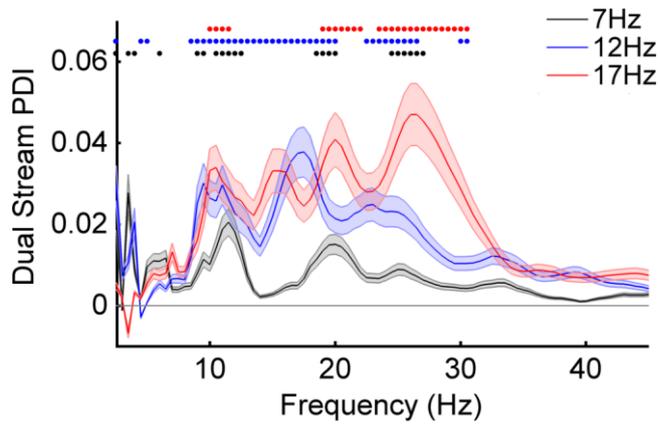
507 **a**, The proportion of responses classified as A1 when visual stimuli as V1 or V2 is plotted with each  
508 unit color coded according classification as A, V or AV. **b**, Mean ( $\pm$  SEM) values across the population.  
509 Two way ANOVA revealed a significant effect of visual condition ( $F(1, 1212) = 173.463, p < 0.001$ ) and  
510 unit type ( $F(2, 1212) = 1780.803, p < 0.001$ ), and interaction between factors ( $F(2, 1212) = 88.66,$   
511  $p < 0.001$ ) on the proportion of responses classified as A1. While the effect was strongest in AV and V  
512 units, pairwise post-hoc comparisons revealed a significant effect of visual stimulus in all subgroups  
513 ( $p < 0.05$ ). **c**, Mean ( $\pm$  SEM) dual stream phase dissimilarity index (PDI) values for recording sites  
514 categorised according to the spiking responses recorded there. Symbols indicate where the dual  
515 stream PDI index was significant (pairwise t-test,  $p < 0.001$  with correction). While the phase effects  
516 are greatest at the sites where visual activity was recorded, significant dual stream PDI values were  
517 observed in all three unit types. In all three cases significant phase coherence was seen at 12Hz,  
518 13.5Hz-14.5Hz and 42.5-44.5Hz. Modulation at 10-12 Hz was only observed at sites in which AV and V  
519 responses were recorded.



520

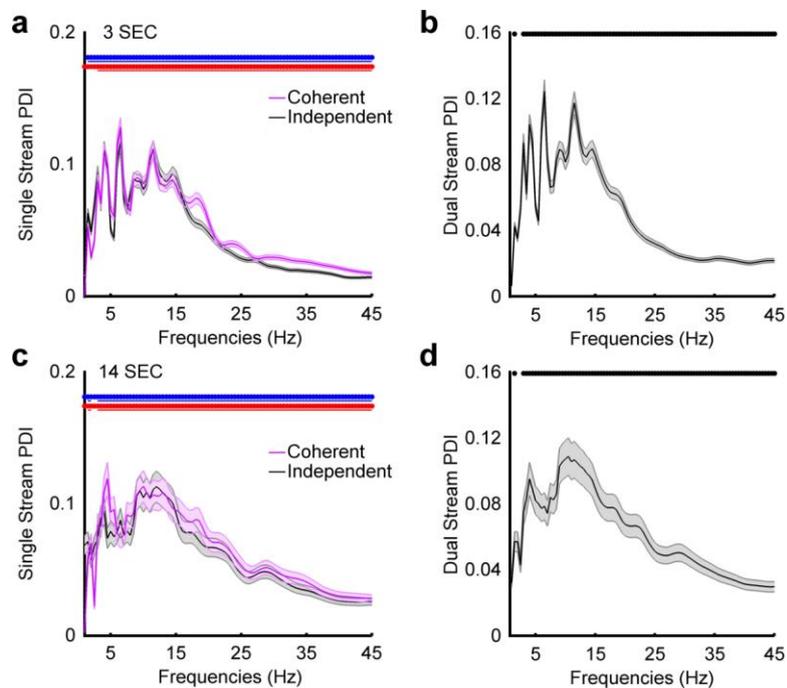
521 **Supplementary figure 5: The effects of temporal coherence on single stream decoding and of**  
522 **visual identity on dual stream decoding are evident in both single and multi-units.**

523 **a,b** Discrimination accuracy for single stream decoding is indistinguishable for single units and multi-  
524 unit activity. Two way ANOVA revealed a significant effect of visual condition and no effect of unit  
525 type in discrimination accuracy values in awake recording (a, Visual condition:  $F(1,851) = 14.6607$   
526  $p < 0.01$ , Unit Type:  $F(1,851) = 1.02$ ,  $p = 0.312$ ; interaction:  $F(1,851) = 0.005$ ,  $p = 0.449$ ) and anaesthetised  
527 recording. (b, Visual condition:  $F(1,1847) = 24.514$ ,  $p < 0.01$ , Unit Type:  $F(1,1847) = 0$ ,  $p = 0.986$ ;  
528 interaction:  $F(1,1847) = 0.66$ ,  $p = 0.418$ ) **c, d.** Impact of visual stimulus on dual stream decoding in single  
529 units and multi units. A two way ANOVA revealed only significant effect of visual condition and no  
530 effect of unit type in discrimination accuracy values in awake recording (c, Visual condition:  $F(1,449)$   
531  $= 19.41$   $p < 0.01$ , Unit Type:  $F(1,851) = 1.07$ ,  $p = 0.414$ ; interaction:  $F(1,851) = 0.33$ ,  $p = 0.568$ ) and  
532 anaesthetised recording. (d, Visual condition:  $F(1,1847) = 15.41$ ,  $p < 0.01$ , Unit Type:  $F(1,1847) = 1.1$ ,  
533  $p = 0.506$ ; interaction:  $F(1,1847) = 0.66$ ,  $p = 0.418$ ).



535 **Supplementary figure 6: Dual stream PDI values in the alpha range are independent of amplitude**  
536 **modulation rate.**

537 Stimuli were generated with three different amplitude modulation rates (<7Hz, as in the main  
538 experiment, <12Hz and <17Hz) and the responses to these were recorded in 92 units. Symbols  
539 indicate where the dual stream phase selectivity index was significant (pairwise t-test,  $p < 0.05$  with  
540 correction). In all three cases significant phase coherence is seen between 10Hz-11.5Hz, 19Hz-20Hz  
541 and 24-26 Hz.



543 **Supplementary figure 7: Stimulus duration does not influence the frequency range across which**  
544 **significant single stream PDI values and dual stream PDI values are observed.**

545 We observed that phase coherence values were higher for longer duration stimuli and therefore in  
546 order to directly compare anaesthetised and awake datasets, all analysis in the anaesthetised dataset  
547 was restricted to the first three seconds of stimuli. At longer stimulus durations the ITPC profile and  
548 resulting PDI varied more smoothly with frequency (a, c); However at both durations phase values  
549 were significantly different from zero at all frequencies. The pattern of significant phase selectivity  
550 values was also preserved across stimulus durations. (**b, d**). Frequency points at which the single  
551 stream PDI value and dual stream PDI values were similar in 3 second length (**a, b**) and 14 second  
552 length (**c, d**) Blue, red and black symbols indicate where the PDI was significant (pairwise t-test,  $\alpha =$   
553 0.0012 with correction).

## 554 Online Methods

555 **Animal preparation:** The experiments were approved by the Committee on Animal Care and Ethical  
556 Review of University College London and The Royal Veterinary College, and performed under license  
557 from the UK Home Office and in accordance with the Animals Scientific Procedures Act 1986. Neural  
558 responses were recorded in a total of 11 awake pigmented adult female ferrets (*Mustela putorius*  
559 *furo*; 1-5 years old). Data from 9 of these animals was used for the main experiment (532 units), data  
560 from 8 other animals (6/9 and two additional ferrets, 89 units) was collected for additional control  
561 analysis (Figures 6e, supplemental figure 6). Animals were chronically implanted with recording  
562 electrodes and passively listening/watching stimuli while holding their head at a water spout and  
563 receiving continuous water reward. These animals were trained in various listening tasks for other  
564 studies. An additional 5 adult females were used to record responses under anaesthesia.

565

566 Full methods for recording under anesthesia can be found in Bizley et al.,<sup>44</sup>. Briefly, ferrets were  
567 anesthetized with medetomidine (Domitor; 0.022mg/kg/h; Pfizer, Sandwich, UK) and ketamine  
568 (Ketaset; 5mg/kg/h; Fort Dodge Animal Health, Southampton, UK). The animal was intubated and the  
569 left radial vein was cannulated in order to provide a continuous infusion (5 ml/h) of a mixture of  
570 medetomidine and ketamine in lactated ringers solution augmented with 5% glucose, atropine sulfate  
571 (0.06 mg/kg/h; C-Vet Veterinary Products) and dexamethasone (0.5 mg/kg/h, Dexadreson; Intervet,  
572 UK). The ferret was placed in a stereotaxic frame in order to implant a bar on the skull, enabling the  
573 subsequent removal of the stereotaxic frame. The left temporal muscle was largely removed, and the  
574 suprasylvian and pseudosylvian sulci were exposed by a craniotomy, revealing auditory cortex (Kelly  
575 et al., 1986). The dura was removed over auditory cortex and the brain protected with 3% agar  
576 solution. The eyes were protected with zero-refractive power contact lenses. The animal was then  
577 transferred to a small table in a sound-attenuating chamber. Body temperature, end-tidal CO<sub>2</sub>, and  
578 the electrocardiogram were monitored throughout the experiment. Experiments typically lasted

579 between 36 and 56 h. Neural activity was recorded with multisite silicon electrodes (Neuronexus  
580 Technologies) in a 1x 16, 2x 16 or 4x 8 (shank x number of sites) configuration.

581

582 Full surgical methods for recording implanting electrode arrays to facilitate recording from awake  
583 animals are available in Bizley et al. <sup>45</sup>. Briefly, animals were bilaterally implanted with WARP-16 drives  
584 (Neuralynx, Montana, USA) loaded with high impedance tungsten electrodes (FHC, Bowdoin, USA)  
585 under general anaesthesia (medetomidine and ketamine induction, as above, isoflurane maintenance  
586 1-3%). Craniotomies were made over left and right auditory cortex, a small number of screws were  
587 inserted into the skull for anchoring and grounding the arrays, and the WARP-16 drive was anchored  
588 with dental acrylic and protected with a capped well. Animals were allowed to recover for a week  
589 before the electrodes were advanced into auditory cortex. Pre-operative, peri-operative and post-  
590 operative analgesia were provided to animals under veterinary advice.

591

592 **Stimulus presentation:** All stimuli were created using TDT System 3 hardware (Tucker-Davis  
593 Technologies, Alachua, FL) and controlled via MATLAB (Mathworks, USA). For recordings in awake  
594 animals, sounds were presented over two loud speakers (Visaton FRS 8). Water deprived ferrets were  
595 placed in a dimly lit testing box (69 x 42 x 52 cm length x width x height) and received water from a  
596 central reward spout located between the two speakers. Sound levels were calibrated using a Brüel  
597 and Kjær (Norcross, GA) sound level meter and free-field ½-inch microphone (4191). Auditory streams  
598 were presented at 65 dB SPL (Fig. 1a). Visual stimuli were delivered by illuminating the spout with a  
599 white LED which provided full field illumination (Precision Gold N76CC Luxmeter, 0 to 36.9 lux). The  
600 animals were not required to do anything other than maintain their heads in position at the spout  
601 where they were freely rewarded. Recording was terminated when animals were sated.

602 For anaesthetised recordings, acoustic stimuli were presented using Panasonic headphones (Panasonic  
603 RP-HV297, Bracknell, UK) at 65 dB SPL. Visual stimuli were presented with a white Light Emitting Diode  
604 (LED) which was placed in a diffuser at a distance of roughly 10 cm from the contralateral eye so that  
605 it illuminated virtually the whole contralateral visual field.

606

607 *Stimuli and data acquisition:* *Auditory stimuli* were artificial vowel sounds that were created in  
608 Matlab (MathWorks, USA). In the behavioural experiment that motivated this study<sup>5</sup>, stimuli were 14  
609 seconds in duration. However, we adapted the stimulus duration in awake recordings to 3 seconds in  
610 order to collect sufficient repetitions of all stimuli, and to ensure animals maintained their head  
611 position facing forwards for the whole trial duration. In the anaesthetised recording stimulus streams  
612 were 14 seconds long, as in the human psychophysics<sup>5</sup> but we only analysed the first 3 seconds to  
613 ensure datasets were directly comparable (Supplemental Figure 7).

614 Stimulus A1 was the vowel [u] (formant frequencies F1-4: 460, 1105, 2857, 4205 Hz, F0= 195Hz), A2  
615 was [a] (F1-4: 936, 1551, 2975, 4263 Hz, F0= 175Hz). Streams were amplitude modulated with a noisy  
616 lowpass (7 Hz cutoff) envelope. Unless specifically noted, the timbre of the auditory stream remained  
617 fixed throughout the trial. However, we also recorded responses to auditory streams that included  
618 brief timbre deviants. As in our previous behavioural study, deviants were 200ms epochs in which the  
619 identity of the vowel was varied by smoothly changing the first and second formant frequencies to  
620 and from those identifying another vowel. Stream A1 was morphed to/from [ε] (730, 2058, 2857,  
621 4205 Hz) and A2 to/from [i] (437, 2761, 2975, 4263 Hz).

622 *Visual stimuli* were generated using an LED whose luminance was modulated with dynamics that  
623 matched the amplitude modulation applied to A1 or A2. In single stream conditions a single auditory  
624 and single visual stream were presented (e.g. A1V1, A1V2, A2V1, or A2V2) whereas in dual stream  
625 conditions both auditory streams were presented simultaneously, accompanied by a single visual  
626 stimulus (A12V1, A12V2, A12V1 A12V2) (Fig. 1e). Auditory streams were always presented from both

627 speakers so that spatial cues could not facilitate segregation, and stimulus order was varied pseudo-  
628 randomly. In the anaesthetised recordings each stimulus was presented 20 times. In the awake dataset,  
629 where recording duration was determined by how long the ferret remained at the central location  
630 (mean repetitions: 20, minimum: 14, maximum: 34).

631 During anaesthetised recordings, pure tone stimuli (150 Hz to 19 kHz in 1/3-octave steps, from 10 to  
632 80 dB SPL in 10 dB, 100 ms in duration, 5 ms cosine ramped) were also presented. These allowed us  
633 to characterize individual units and determine tonotopic gradients, so as to confirm the cortical field  
634 in which any given recording was made. Additionally broadband noise bursts and diffuse light flashes  
635 (100 ms duration, 70 dB SPL) were presented and used to classify a stimulus as auditory, visual or  
636 auditory visual. LFPs were subjected to current source density analysis to identify sources and sinks as  
637 described by Kaur et al.<sup>46</sup>

638

639 **Data Analysis:** Electrophysiological data were analysed offline. Spiking activity and local field potential  
640 signals were extracted from the broadband voltage waveform by filtering at 0.3-5kHz and 1-150 Hz  
641 respectively. Spikes were detected, extracted and then sorted with a spike-sorting algorithm  
642 (WaveClus) (Quiroga et al., 2004).

643 We used a Euclidean distance based pattern classifier (Schnupp et al., 2006) with leave-one-out cross  
644 validation to determine whether the neuronal responses to different stimuli could be discriminated.  
645 Spiking responses to a given stimulus were binned into a series of spike counts from stimulus onset (0  
646 s) to offset (3s) in 20 ms bins. The average across-repetition response to each stimulus (excluding the  
647 to-be-classified response) were used as templates and the response to a single stimulus presentation  
648 was classified by calculating the Euclidean distance between itself and the template sweeps and  
649 assigning it to the closest template. To determine whether the classifier performed significantly better  
650 than expected by chance, a 1000 iteration permutation test was performed where trials were drawn  
651 (with replacement) from the observed data and randomly assigned to a stimulus that was then used

652 for template formation / decoding. A neural response was considered to be significantly informative  
653 about stimulus identity if the observed value exceeded the 95th percentile of the randomly-drawn  
654 distribution.

655 This approach allowed us to classify units according to their functional properties: auditory units  
656 discriminated two auditory stimuli based on the amplitude modulation of sound (A1 versus A2)  
657 regardless of visual dynamics, (Fig. 2a, b), visual units discriminated visual presentations based on  
658 temporal envelope of visual stimuli (V1 versus V2) regardless of auditory presentation (Fig. 2c, d) and  
659 AV units could do both. This approach was extended to classify dual stream responses by using the  
660 average response to each of the temporally coherent single stream stimuli (A1V1 or A2V2) as  
661 templates. Performance was (arbitrarily) expressed as the proportion of responses classified as being  
662 from the A1, and compared for the two dual stream stimuli with different visual conditions (Figure 5).  
663 To be considered in this analysis the response of a unit had to be informative about the single stream  
664 stimuli (i.e. classified as either auditory and/or visual discriminating).

665

666 **Phase/power dissimilarity analysis:** Local field potential recordings were considered for all sites at  
667 which there was a significant driven spiking response, irrespective of whether that response could  
668 discriminate auditory or visual stream identity. For the single stream trials, we computed a single  
669 Stream Phase Dissimilarity Index (PDI), which characterizes the consistency and uniqueness of the  
670 temporal phase/power pattern of neural responses to continuous auditory stimuli (Luo and Poeppel,  
671 2007). This analysis compares the phase (or power) consistency across repetitions of the same  
672 stimulus with a baseline of phase-consistency across trials in which different stimuli were presented.

673 In the first stage of PDI analysis, we obtained a time-frequency representation of each response using  
674 wavelet decomposition with complex 7-cycle Morlet wavelets in 0.5 steps between 2.5–45 Hz,  
675 resulting in 86 frequency points. Next, we calculated the inter-trial phase-coherence value (ITPC;  
676 Equ.1) at each time-frequency point, across all trials in which the same stimulus was presented. For

677 each frequency band, the ITPC time-course was averaged over the duration of the analysis window  
678 and across all repetitions to obtain the average *within-stimulus ITPC*.

679

$$680 \quad ITPC_{t,f} = \left| \frac{\sum_{k=1}^N e^{i\theta_{k,t,f}}}{N} \right| \quad \text{Equ.1}$$

681 In which N is equal to the number of trials, and  $\theta$  is the phase of trial  $k$  at a given frequency ( $f$ ) and  
682 time ( $t$ ). The *across-stimuli ITPC* was estimated using the same approach but using shuffled data, such  
683 that the ITPC was computed across trials with the same auditory stimulus but randomly drawn visual  
684 stimuli. The single stream phase dissimilarity index (Single stream PDI) was computed as the difference  
685 between the ITPC value calculated for *within* trials and the ITPC values calculated *across* visual trials  
686 (Equ.2).  $\gamma$ . The dissimilarity function for each frequency bin  $i$  was defined as;

687

$$688 \quad \text{Single Stream PDI}_i = \frac{\sum_{j=1}^N ITPC_{ij} \text{ within}_{vis}}{N} - \frac{\sum_{j=1}^N ITPC_{ij} \text{ across}_{vis}}{N} \quad \text{Equ.2}$$

689

690 Large positive PDI indicate that responses to individual stimuli have a highly consistent response on  
691 single trials. Single stream PDI values were calculated for each stimulus type and then averaged across  
692 stimuli to calculate values for temporally coherent and temporally independent auditory visual stimuli.  
693 Single stream PDI was positive if within stimulus ITPC was larger than across-stimulus ITPC (pairwise  
694 t-test,  $p < 0.05$  Bonferroni correction for 86 frequencies points) and was considered significant if a  
695 minimum of 2 adjacent bins exceeded the corrected threshold.

696 Dual stream phase dissimilarity index (dual stream PDI) values were calculated by extending this  
697 approach for dual stream stimuli with the goal of determining how the temporal envelope of the visual  
698 stimulus influences the neural response to a sound mixture. To this end, we calculated the *within-dual*  
699 *ITPC* from the A12V1 trials and A12V2 trials separately and *across-dual ITPC* by randomly selecting

700 trials from both stimuli (Equ.3). The within-dual and across-dual ITPCs were then averaged over time  
701 and subtracted to yield the dual stream PDI (Equ.3).

702

$$703 \quad \text{Dual Stream PDI}_i = \frac{\sum_{j=1}^N \text{ITPC}_{ij} \text{ within}_{dual}}{N} - \frac{\sum_{j=1}^N \text{ITPC}_{ij} \text{ across}_{dual}}{N} \quad \text{Equ.3}$$

704

705 Positive dual stream PDI values indicate that the time course of the neural responses was influenced  
706 by visual input, despite the identical acoustic input. We determined whether the dual stream PDI was  
707 greater if the *within\_dual ITPC* was significantly larger than *across\_dual ITPC* (pairwise t-test,  $p < 0.05$   
708 Bonferroni correction, as above).

709

710 **Timbre deviant analysis:** In order to determine how a visual stimulus influenced the ability to decode  
711 timbre deviants embedded within the auditory streams we used the cross-validated pattern classifier  
712 described above for analysing single stream stimuli to discriminate deviant from no-deviant trials.  
713 Responses were considered over the 200 ms time window that the deviant occurred (or the equivalent  
714 point in the no-deviant stimulus) binned with a 10 ms resolution. Significance was assessed by a 1000  
715 iteration permutation test in which trials were randomly drawn with replacement from deviant and  
716 no-deviant responses. The discrimination score was calculated as the proportion of correctly classified  
717 trials.

718

719

720

## 721 **References**

722 1. Denison, R.N., Driver, J. & Ruff, C.C. Temporal structure and complexity affect audio-visual  
723 correspondence detection. *Front Psychol* **3** (2013).

- 724 2. Rahne, T. *et al.* A multilevel and cross-modal approach towards neuronal mechanisms of  
725 auditory streaming. *Brain Research* **1220**, 118-131 (2008).
- 726 3. Crosse, M.J., Butler, J.S. & Lalor, E.C. Congruent Visual Speech Enhances Cortical Entrainment  
727 to Continuous Auditory Speech in Noise-Free Conditions. *The Journal of Neuroscience* **35**,  
728 14195-14204 (2015).
- 729 4. Brosch, M., Selezneva, E. & Scheich, H. Neuronal activity in primate auditory cortex during the  
730 performance of audiovisual tasks. *European Journal of Neuroscience* **41**, 603-614 (2015).
- 731 5. Maddox, R.K., Atilgan, H., Bizley, J.K. & Lee, A.K. Auditory selective attention is enhanced by a  
732 task-irrelevant temporally coherent visual stimulus in human listeners. *Elife* **4**, e04995 (2015).
- 733 6. Bizley, J.K., Maddox, R.K. & Lee, A.K. Defining Auditory-Visual Objects: Behavioral Tests and  
734 Physiological Mechanisms. *Trends in Neurosciences* (2016).
- 735 7. Bizley, J.K., Nodal, F.R., Bajo, V.M., Nelken, I. & King, A.J. Physiological and anatomical  
736 evidence for multisensory interactions in auditory cortex. *Cereb Cortex* **17**, 2172-2189 (2007).
- 737 8. Chandrasekaran, C., Lemus, L. & Ghazanfar, A.A. Dynamic faces speed up the onset of auditory  
738 cortical spiking responses during vocal detection. *Proceedings of the National Academy of  
739 Sciences* **110**, E4668-E4677 (2013).
- 740 9. Ghazanfar, A.A., Maier, J.X., Hoffman, K.L. & Logothetis, N.K. Multisensory integration of  
741 dynamic faces and voices in rhesus monkey auditory cortex. *The Journal of Neuroscience* **25**,  
742 5004-5012 (2005).
- 743 10. Kayser, C., Petkov, C.I. & Logothetis, N.K. Visual modulation of neurons in auditory cortex.  
744 *Cerebral Cortex* **18**, 1560-1574 (2008).
- 745 11. Perrodin, C., Kayser, C., Logothetis, N.K. & Petkov, C.I. Natural asynchronies in audiovisual  
746 communication signals regulate neuronal multisensory interactions in voice-sensitive cortex.  
747 *Proceedings of the National Academy of Sciences* **112**, 273-278 (2015).
- 748 12. Azouz, R. & Gray, C.M. Cellular mechanisms contributing to response variability of cortical  
749 neurons in vivo. *The Journal of neuroscience* **19**, 2209-2223 (1999).
- 750 13. Okun, M., Naim, A. & Lampl, I. The subthreshold relation between cortical local field potential  
751 and neuronal firing unveiled by intracellular recordings in awake rats. *The Journal of  
752 neuroscience* **30**, 4440-4448 (2010).
- 753 14. Szymanski, F.D., Rabinowitz, N.C., Magri, C., Panzeri, S. & Schnupp, J.W. The laminar and  
754 temporal structure of stimulus information in the phase of field potentials of auditory cortex.  
755 *The Journal of Neuroscience* **31**, 15787-15801 (2011).
- 756 15. Chandrasekaran, C., Turesson, H.K., Brown, C.H. & Ghazanfar, A.A. The influence of natural  
757 scene dynamics on auditory cortical activity. *The Journal of Neuroscience* **30**, 13919-13931  
758 (2010).
- 759 16. Kayser, C., Petkov, C.I. & Logothetis, N.K. Multisensory interactions in primate auditory cortex:  
760 fMRI and electrophysiology. *Hearing Res* **258**, 80-88 (2009).
- 761 17. Luo, H. & Poeppel, D. Phase patterns of neuronal responses reliably discriminate speech in  
762 human auditory cortex. *Neuron* **54**, 1001-1010 (2007).
- 763 18. Ng, B.S.W., Schroeder, T. & Kayser, C. A precluding but not ensuring role of entrained low-  
764 frequency oscillations for auditory perception. *The Journal of Neuroscience* **32**, 12268-12276  
765 (2012).
- 766 19. Schyns, P.G., Thut, G. & Gross, J. Cracking the code of oscillatory activity. *PLoS Biol* **9**, e1001064  
767 (2011).
- 768 20. Golumbic, E.Z., Cogan, G.B., Schroeder, C.E. & Poeppel, D. Visual input enhances selective  
769 speech envelope tracking in auditory cortex at a "cocktail party". *The Journal of Neuroscience*  
770 **33**, 1417-1426 (2013).
- 771 21. Lakatos, P., Chen, C.-M., O'Connell, M.N., Mills, A. & Schroeder, C.E. Neuronal oscillations and  
772 multisensory interaction in primary auditory cortex. *Neuron* **53**, 279-292 (2007).
- 773 22. Desimone, R. & Duncan, J. Neural mechanisms of selective visual attention. *Annual review of  
774 neuroscience* **18**, 193-222 (1995).

- 775 23. Chandrasekaran, C., Trubanova, A., Stillittano, S., Caplier, A. & Ghazanfar, A.A. The natural  
776 statistics of audiovisual speech. *Plos Comput Biol* **5**, e1000436 (2009).
- 777 24. Luo, H., Liu, Z. & Poeppel, D. Auditory cortex tracks both auditory and visual stimulus dynamics  
778 using low-frequency neuronal phase modulation. *PLoS Biol* **8**, e1000445 (2010).
- 779 25. Crosse, M.J., Di Liberto, G.M. & Lalor, E.C. Eye Can Hear Clearly Now: Inverse Effectiveness in  
780 Natural Audiovisual Speech Processing Relies on Long-Term Crossmodal Temporal  
781 Integration. *The Journal of Neuroscience* **36**, 9888-9895 (2016).
- 782 26. Peelle, J.E. & Sommers, M.S. Prediction and constraint in audiovisual speech perception.  
783 *Cortex* **68**, 169-181 (2015).
- 784 27. Sumbly, W.H. & Pollack, I. Visual contribution to speech intelligibility in noise. *The journal of*  
785 *the acoustical society of america* **26**, 212-215 (1954).
- 786 28. Schroeder, C.E., Lakatos, P., Kajikawa, Y., Partan, S. & Puce, A. Neuronal oscillations and visual  
787 amplification of speech. *Trends in cognitive sciences* **12**, 106-113 (2008).
- 788 29. Okada, K., Venezia, J.H., Matchin, W., Saberi, K. & Hickok, G. An fMRI study of audiovisual  
789 speech perception reveals multisensory interactions in auditory cortex. *Plos One* **8**, e68959  
790 (2013).
- 791 30. Möttönen, R., Schürmann, M. & Sams, M. Time course of multisensory interactions during  
792 audiovisual speech perception in humans: a magnetoencephalographic study. *Neuroscience*  
793 *letters* **363**, 112-115 (2004).
- 794 31. Henry, M.J. & Obleser, J. Frequency modulation entrains slow neural oscillations and  
795 optimizes human listening behavior. *Proceedings of the National Academy of Sciences* **109**,  
796 20095-20100 (2012).
- 797 32. Jacobs, J., Kahana, M.J., Ekstrom, A.D. & Fried, I. Brain oscillations control timing of single-  
798 neuron activity in humans. *The Journal of neuroscience* **27**, 3839-3844 (2007).
- 799 33. Klimesch, W., Sauseng, P. & Hanslmayr, S. EEG alpha oscillations: The inhibition-timing  
800 hypothesis. *Brain Research Reviews* **53**, 63-88 (2007).
- 801 34. Lakatos, P. *et al.* The spectrotemporal filter mechanism of auditory selective attention. *Neuron*  
802 **77**, 750-761 (2013).
- 803 35. Lőrincz, M.L., Kékesi, K.A., Juhász, G., Crunelli, V. & Hughes, S.W. Temporal framing of thalamic  
804 relay-mode firing by phasic inhibition during the alpha rhythm. *Neuron* **63**, 683-696 (2009).
- 805 36. Mazzone, A., Panzeri, S., Logothetis, N.K. & Brunel, N. Encoding of naturalistic stimuli by local  
806 field potential spectra in networks of excitatory and inhibitory neurons. *PLoS Comput Biol* **4**,  
807 e1000239 (2008).
- 808 37. Tukker, J.J., Fuentealba, P., Hartwich, K., Somogyi, P. & Klausberger, T. Cell type-specific tuning  
809 of hippocampal interneuron firing during gamma oscillations in vivo. *The journal of*  
810 *neuroscience* **27**, 8184-8189 (2007).
- 811 38. Voloh, B. & Womelsdorf, T. A Role of Phase-Resetting in Coordinating Large Scale Neural  
812 Networks During Attention and Goal-Directed Behavior. *Frontiers in systems neuroscience* **10**  
813 (2016).
- 814 39. Wang, X.-J. Neurophysiological and computational principles of cortical rhythms in cognition.  
815 *Physiological reviews* **90**, 1195-1268 (2010).
- 816 40. Samaha, J., Bauer, P., Cimaroli, S. & Postle, B.R. Top-down control of the phase of alpha-band  
817 oscillations as a mechanism for temporal prediction. *Proceedings of the National Academy of*  
818 *Sciences* **112**, 8439-8444 (2015).
- 819 41. O'Sullivan, J.A., Shamma, S.A. & Lalor, E.C. Evidence for neural computations of temporal  
820 coherence in an auditory scene and their enhancement during active listening. *The Journal of*  
821 *Neuroscience* **35**, 7256-7263 (2015).
- 822 42. Elhilali, M., Ma, L., Micheyl, C., Oxenham, A.J. & Shamma, S.A. Temporal Coherence in the  
823 Perceptual Organization and Cortical Representation of Auditory Scenes. *Neuron* **61**, 317-329  
824 (2009).

- 825 43. Budinger, E., Heil, P., Hess, A. & Scheich, H. Multisensory processing via early cortical stages:  
826 connections of the primary auditory cortical field with other sensory systems. *Neuroscience*  
827 **143**, 1065-1083 (2006).
- 828 44. Bizley, J.K., Walker, K.M.M., Silverman, B.W., King, A.J. & Schnupp, J.W.H. Interdependent  
829 Encoding of Pitch, Timbre, and Spatial Location in Auditory Cortex. *Journal of Neuroscience*  
830 **29**, 2064-2075 (2009).
- 831 45. Bizley, J.K., Walker, K.M.M., King, A.J. & Schnupp, J.W.H. Spectral timbre perception in ferrets:  
832 Discrimination of artificial vowels under different listening conditions. *J Acoust Soc Am* **133**,  
833 365-376 (2013).
- 834 46. Kaur, S., Rose, H., Lazar, R., Liang, K. & Metherate, R. Spectral integration in primary auditory  
835 cortex: laminar processing of afferent input, in vivo and in vitro. *Neuroscience* **134**, 1033-1045  
836 (2005).

837

838