

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, Department of Neuroscience, Del Monte Institute for  
8 Neuroscience, University of Rochester, Rochester, NY, USA

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

11

12 7 Figures

13 68072 characters

14

15

16 Corresponding Author and Lead Contact: Jennifer Bizley [j.bizley@ucl.ac.uk](mailto:j.bizley@ucl.ac.uk)

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

## 18 Summary

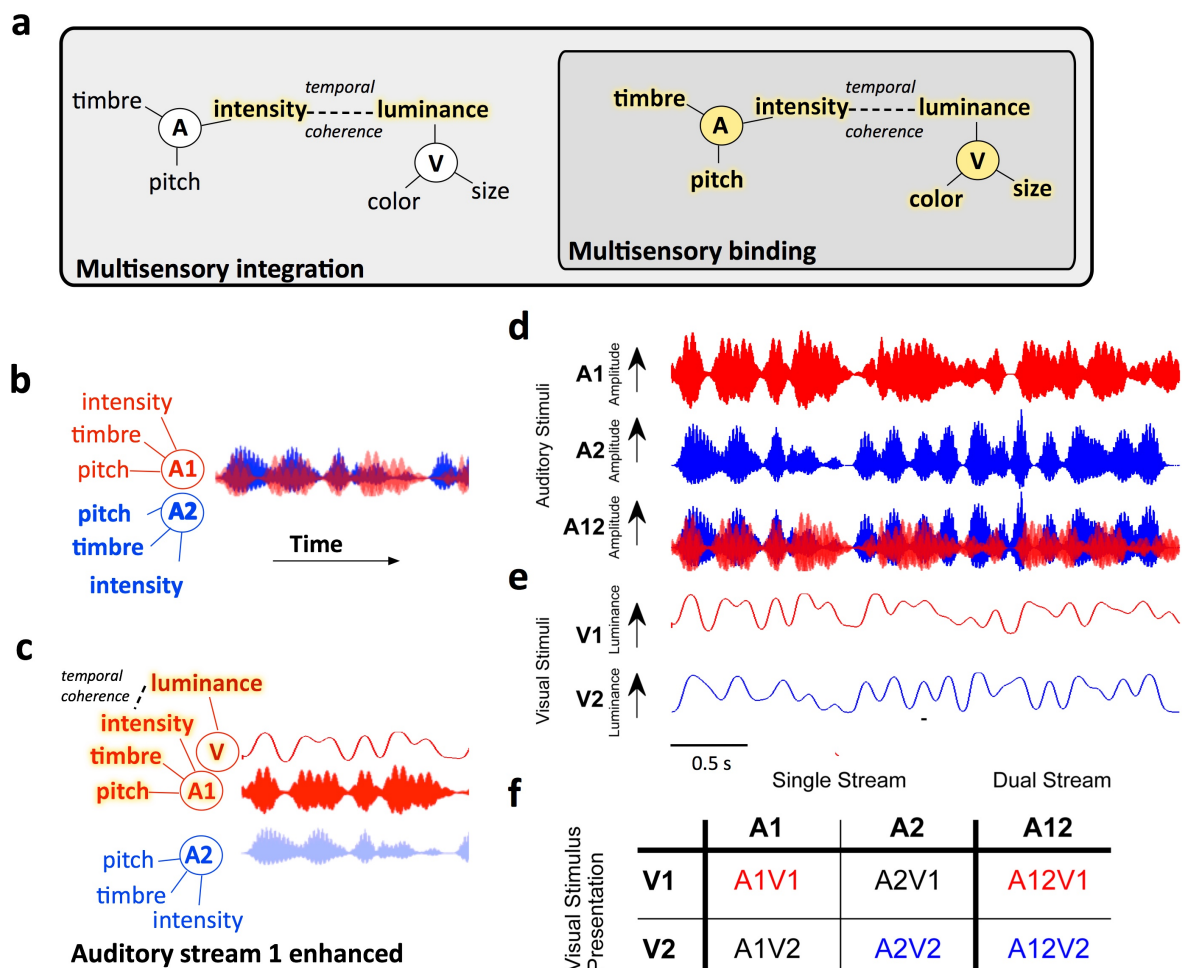
19 How and where in the brain audio-visual signals are bound to create multimodal objects remains  
20 unknown. One hypothesis is that temporal coherence between dynamic multisensory signals  
21 provides a mechanism for binding stimulus features across sensory modalities. Here we report that  
22 when the luminance of a visual stimulus is temporally coherent with the amplitude fluctuations of  
23 one sound in a mixture, the representation of that sound is enhanced in auditory cortex. Critically,  
24 this enhancement extends to include both binding and non-binding features of the sound. We  
25 demonstrate that visual information conveyed from visual cortex, via the phase of the local field  
26 potential is combined with auditory information within auditory cortex. These data provide evidence  
27 that early cross-sensory binding provides a bottom-up mechanism for the formation of cross-sensory  
28 objects and that one role for multisensory binding in auditory cortex is to support auditory scene  
29 analysis.

## 30 Introduction

31 When listening to a sound of interest, we frequently look at the source. However, how auditory and  
32 visual information are integrated to form a coherent perceptual object is unknown. The temporal  
33 properties of a visual stimulus can be exploited to detect correspondence between auditory and  
34 visual streams (Crosse et al., 2015; Denison et al., 2013; Rahne et al., 2008), can bias the perceptual  
35 organisation of a sound scene (Brosch et al., 2015), and can enhance or impair listening performance  
36 depending on whether the visual stimulus is temporally coherent with a target or distractor sound  
37 stream (Maddox et al., 2015). Together, these behavioural results suggest that temporal coherence  
38 between auditory and visual stimuli can promote binding of cross-modal features to enable the  
39 formation of an auditory-visual (AV) object (Bizley et al., 2016b).

40 Visual stimuli can both drive and modulate neural activity in primary and non-primary auditory  
41 cortex (Bizley et al., 2007a; Chandrasekaran et al., 2013; Ghazanfar et al., 2005; Kayser et al., 2008;  
42 Kayser et al., 2010, Perrodin et al., 2015), but the contribution that visual activity in auditory cortex  
43 makes to auditory function remains unknown. One possibility is that the integration of cross-sensory  
44 information into early sensory cortex provides a bottom-up substrate for the binding of multisensory  
45 stimulus features into a single perceptual object (Bizley et al., 2016b). We have recently argued that  
46 binding is a distinct form of multisensory integration that underpins perceptual object formation.  
47 We hypothesise that binding is associated with a modification of the sensory representation and can  
48 be identified by demonstrating a benefit in the behavioural or neural discrimination of a stimulus  
49 feature orthogonal to the features that link crossmodal stimuli (Fig. 1a). Therefore, in order to  
50 demonstrate binding, an appropriate crossmodal stimulus should elicit not only enhanced neural  
51 encoding of the stimulus features that bind auditory and visual streams (the “binding features”), but  
52 that there should be enhancement in the representation of *other* stimulus features (“non-binding  
53 features” associated with the source (Fig. 1c).

54 Here we test the hypothesis that the incorporation of visual information into auditory cortex can  
 55 determine the neuronal representation of an auditory scene through multisensory binding (Fig.1).  
 56 We demonstrate that when visual luminance changes coherently with the amplitude of one sound in  
 57 a mixture, auditory cortex is biased towards representing the temporally coherent sound. Consistent  
 58 with these effects reflecting cross-modal binding, the encoding of sound timbre, a non-binding  
 59 stimulus feature, is subsequently enhanced in the temporally coherent auditory stream. Finally, we  
 60 demonstrate that the site of multisensory convergence is in auditory cortex and that visual  
 61 information is conveyed via the local field potential directly from visual cortex.



62

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

64 **a** Conceptual model illustrating how binding can be identified as a distinct form of multisensory  
 65 integration. Multisensory binding is defined as a subset of multisensory integration that results in  
 66 the formation of a crossmodal object. During binding, all features of the audio-visual object are  
 67 linked and enhanced - including both those features that bind the stimuli across modalities (here

68 temporal coherence between auditory (A) intensity and visual (V) luminance) and orthogonal  
69 features such as auditory pitch and timbre, and visual colour and size. Other forms of multisensory  
70 integration would result in enhancement of only the features that promote binding - here auditory  
71 intensity and visual luminance. To identify binding therefore requires a demonstration that non-  
72 binding features (e.g. here pitch, timbre, colour or size) are enhanced. Enhanced features are  
73 highlighted in yellow. **b** When two competing sounds (red and blue waveforms) are presented they  
74 can be separated on the basis of their features, but may elicit overlapping neuronal representations  
75 in auditory cortex. **c** Hypothesised enhancement in auditory stream segregation when a temporally  
76 coherent visual stimulus enables multisensory binding. When the visual stimulus changes coherently  
77 with the red sound (A1, top) this sound is enhanced and the two sources are better segregated.  
78 Perceptually this would result in more effective auditory scene analysis and an enhancement of the  
79 non-binding features. **d** Stimulus design: Auditory stimuli were two artificial vowels (denoted A1 and  
80 A2), each with distinct pitch and timbre and independently amplitude modulated with a noisy low  
81 pass envelope. **e** Visual stimulus: a luminance modulated white light was presented with one of two  
82 temporal envelopes derived from the amplitude modulations of A1 and A2. **f** illustrates the stimulus  
83 combinations that were tested experimentally in *single stream* (a single auditory visual pair) and  
84 *dual stream* (two sounds and one visual stimulus) conditions. See also supplemental figure 1.

85

## 86 Results

87 We recorded neuronal responses in the auditory cortex of awake passively listening ferrets (n=9  
88 ferrets, 221 single units, 311 multi-units) in response to naturalistic time-varying auditory and visual  
89 stimuli adapted from Maddox et al (2015). The stimuli are designed to share properties with natural  
90 speech; they are modulated at approximately syllable rate and, like competing voices, can be  
91 separated on the basis of their fundamental frequency (F0, the physical determinant of pitch). These  
92 sounds are devoid of any linguistic content permitting the separation of general sensory processing  
93 mechanisms from language-specific ones for human listeners. Maddox et al. (2015) used both pure  
94 tones and synthetic vowels as stimuli; here we use synthetic vowels as these robustly drive auditory  
95 cortical responses in the ferret in neurons with a wide range of characteristic frequencies (Bizley et  
96 al., 2009). Ferrets are also well able to distinguish the timbre of artificial vowels (Bizley et al., 2013,  
97 Town et al., 2015), and, like human listeners, both ferret behavioural and neural responses show  
98 invariant responses to vowel timbre across changes in sound level, location and pitch (Town et al.,  
99 2017). We additionally recorded neural responses in medetomidine-ketamine anaesthetised ferrets  
100 (n=5 ferrets, 426 single units, 772 multi units) which allowed us to entirely eliminate attentional

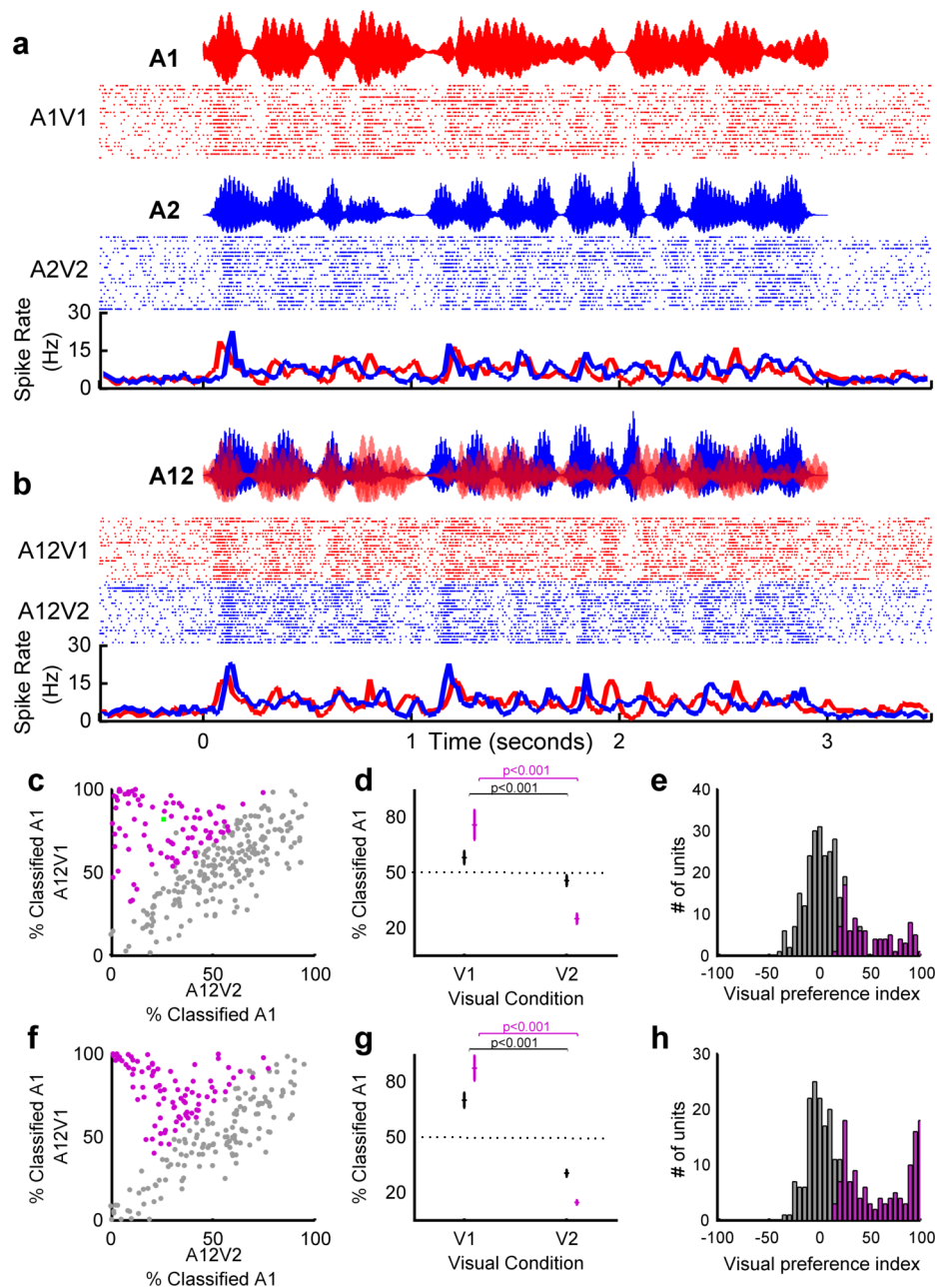
101 effects and limit the impact of top-down processing. These experiments also permitted longer  
102 recording durations for additional control stimuli and enabled simultaneous characterization of  
103 neural activity across cortical laminae. In a subset of these animals we were able to reversibly silence  
104 visual cortex during recording, in order to determine the origin of visual-stimulus elicited neural  
105 changes. Recordings were made in awake freely moving animals while they held their head at a  
106 drinking spout but were not engaged in a behavioural task, and allowed us to measure neural  
107 activity free from any confounds associated with pharmacological manipulation and in the absence  
108 of task-directed attention which would likely engage additional neural circuits.

109 The stimuli were two auditory streams comprised of two vowels, each with a distinct pitch and  
110 timbre (denoted A1: /u/, F0 = 175 Hz and A2: /a/, F0 = 195 Hz, Fig.1) and independently amplitude  
111 modulated with a low-pass (<7 Hz) envelope (Fig.1d). A full-field visual stimulus accompanied the  
112 auditory stimuli, the luminance of which was temporally modulated with the modulation envelope  
113 from one of the two auditory streams (Fig.1e). We tested stimulus conditions in which both auditory  
114 streams were presented (“dual stream”) and the visual stimulus was temporally coherent with one  
115 or other of the auditory streams (A12V1 or A12V2, Fig.1e). We also tested conditions in which a  
116 single AV stimulus pair was presented (‘single stream’ stimuli), where the auditory and visual  
117 streams could be temporally coherent (A1V1, A2V2) or independent (A1V2, A2V1), as well as no-  
118 visual control conditions.

## 119 **Auditory-visual temporal coherence shapes the representation of a sound scene in** 120 **auditory cortex**

121 We first asked whether the temporal dynamics of a visual stimulus could selectively enhance the  
122 representation of one sound in a mixture. We therefore recorded responses to auditory scenes  
123 composed of two sounds (A1 and A2), presented simultaneously, with a visual stimulus that was  
124 temporally coherent with one or other auditory stream (A12V1 or A12V2). A visual stimulus is known  
125 to enhance the representation of the envelope of an attended speech stream in auditory cortex

126 (Zion Golumbic et al., 2013; Park et al., 2016). To test whether we could observe a similar  
127 phenomenon in single neurons in the absence of selective attention, we used neural responses to  
128 temporally coherent single stream stimuli (i.e. A1V1 and A2V2) to determine to what extent the  
129 neural response to the sound mixture was specific to one or other sound stream.



130

131 **Figure 2: Visual stimuli can determine which sound stream auditory cortical neurons follow in a**  
132 **mixture.**

133 Spiking responses from an example unit in response to **a**, single stream AV stimuli used as decoding templates  
134 and **b**, dual stream stimuli. In each case rasters and PSTHs are illustrated. When the visual component of the  
135 dual stream was V1, the majority of trials were classified as A1V1 (82% (19/23 trials), and A2V2 when the

136 visual stimulus was V2 (26% 6/23, trials) of responses classified as A1V1 (see also green data point in c),  
137 yielding a visual preference score of 56%. **c-h** population data for awake (**c,d,e** 271 units) and anaesthetised  
138 (**f,g,h** 331 units) datasets. In each case the left panel (**c,f**) shows the distribution of decoding values according  
139 to the visual condition, the middle panel (**d,g**) shows the population mean ( $\pm$  SEM) projecting onto the vertical  
140 axis of panel c / f for V1 condition, and horizontal axis of panel c / f for the V2 condition. **e,h** shows the visual  
141 preference index (VPI). Units in which the VPI was significantly  $>0$  are coloured purple. Pairwise comparisons  
142 revealed significant effect of visual condition on decoding in all datasets: Awake: All:  $t_{540}=6.1, p=2.3e-09$   
143 ( $n=271$ ), Sig VPI:  $t_{180}=18.8, p=2.0e-44$  ( $n=91$ ). Anaesthetised: All:  $t_{660}=9.5, p=3.3e-20$  ( $n=331$ ), Sig. VPI:  $t_{348} =$   
144  $38.9, p=1.2e-128$  ( $n=175$ ) See also supplemental figures 2-4.

145

146 Figure 2 illustrates this approach for a single unit: responses to the temporally coherent single  
147 stream AV stimuli (Fig.2a) formed templates which were used to decode the responses to the dual  
148 stream stimuli (Fig.2b) using a Euclidean distance based spike pattern classifier. Such an approach is  
149 ideally suited for classifying neural responses to time-varying stimuli. Auditory cortical responses to  
150 the dual stream stimuli (A12V1 or A12V2) were more commonly decoded as A1V1 when the visual  
151 stimulus was V1, and A2V2 when the visual stimulus was V2. Performing this analysis for each  
152 neuron in our recorded population yielded similar observations: the coherent auditory stimulus  
153 representation was enhanced (Fig.2c,d,f,g) such that auditory cortical responses to dual-stream  
154 stimuli most closely resembled responses to the single stream stimulus with the shared visual  
155 component.

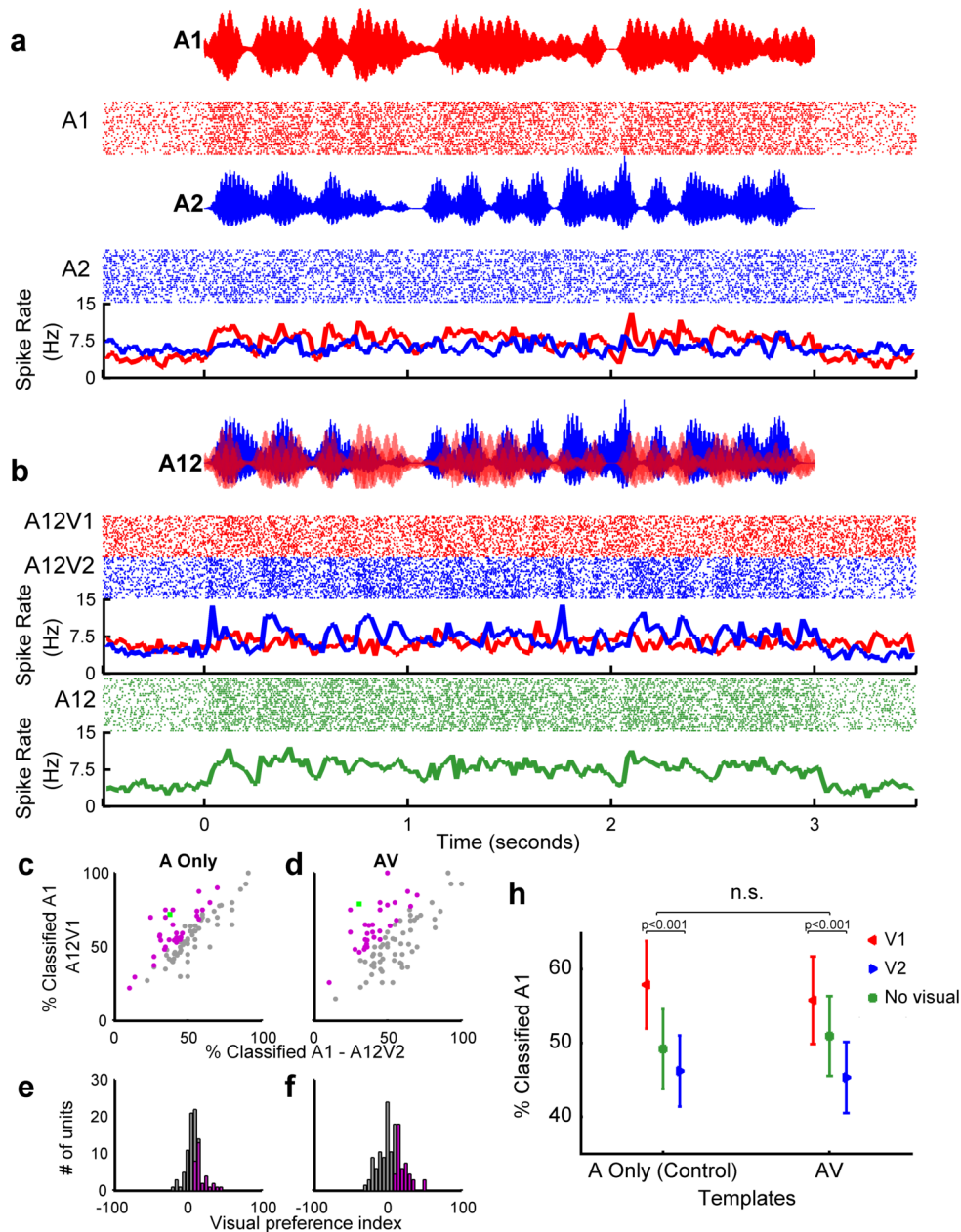
156 To quantify whether the responses of individual units were significantly influenced by the visual  
157 stimulus identity, we first calculated a visual preference index (VPI) as the difference between the  
158 percentage of A12V1 trials labelled A1 and the percentage of A12V2 trials labelled A1. Units which  
159 were fully influenced by the identity of the visual stimulus would have a visual preference score of  
160 100, while those in which the visual stimulus did not influence the response at all would have a score  
161 of 0 (Fig. 2e,h). We assessed the significance of observed VPI scores using a permutation test ( $p <$   
162  $0.05$ ).to revealed that 33.6% of driven units recorded in awake animals (91/271 units) and 52.9% of  
163 units in the anaesthetised dataset (175/331 units) had responses to dual stream stimuli significantly  
164 influenced by the visual stimulus.



165 Modulation of dual stream responses by visual stimulus identity was not simply a consequence of  
166 the shared visual component of single stream and dual stream stimuli and was observed in neurons  
167 in which visual or auditory identity could be decoded (example response from a unit in which only  
168 auditory stimulus identity could be decoded: Fig.S2). If this effect was only apparent in visual  
169 neurons in auditory cortex, then eliminating the visual element of the single stream stimuli should  
170 impair decoding performance for the dual stream stimuli. Additional control experiments (n=89  
171 driven units, awake animals) demonstrated that this was not the case: the enhancement of the  
172 temporally coherent sound in the sound mixture was evident whether dual stream stimuli (A12V1  
173 and A12V2) were decoded using responses to auditory-only stimuli (A1 or A2) or auditory-visual  
174 stimuli (A1V1, A2V2 etc.). Within this control data (Example unit: Fig. 3a,b, population data Fig 3c-h),  
175 32 units had a significant VPI scores when dual stream responses were decoded from an auditory-  
176 only single stream templates and 31 units when decoded with an auditory-visual template.  
177 Furthermore, the distribution of VPI values was statistically indistinguishable for decoding dual  
178 stream responses with A-only or AV templates (Kolmogorov–Smirnov test: all units,  $p = 0.9016$ ; units  
179 with visual preference scores significantly  $>0$ ,  $p > 0.9998$ ), and the distribution of values in Fig.3d was  
180 statistically indistinguishable from that in Fig.2e ( $p = 0.0864$ ). We also determined that removing the  
181 visual stimulus from the dual-stream condition eliminated any decoding difference in responses  
182 observed (Fig.3h). A two-way repeated measures ANOVA on decoded responses with factors of  
183 visual stream (V1, V2, no visual), and template type (AV or A) demonstrated a significant effect of  
184 visual stream identity on dual stream decoding ( $F(2, 528) = 19.320$ ,  $p < 0.001$ ), but there was no  
185 effect of template type ( $F(1,528) = 0.073$ ,  $p = 0.787$ ) or interaction between factors ( $F(2,528) =$   
186  $0.599$ ,  $p = 0.550$ ). Post-hoc comparisons revealed that without visual stimulation, there was no  
187 tendency to respond preferentially to either stream, but that visual stream identity significantly  
188 influenced the classification of dual stream responses.

189

190



191

192 **Figure 3: Visual stimuli shape the neural representation of an auditory scene.**

193 In an additional control experiment (n=89 units recorded in awake animals), the responses to  
 194 coherent AV and auditory-only (A Only) single stream stimuli were used as templates to decode dual  
 195 stream stimuli either accompanied by visual stimuli (V1/V2) or in the absence of visual stimulation  
 196 (no visual). Spiking responses from an example unit in response to **a**, single stream auditory stimuli  
 197 which were used as decoding templates to decode the responses to dual stream stimuli in **b**, in each  
 198 case the auditory waveform, rasters and PSTHs are shown. In this example, when decoded with AV  
 199 templates: 79% (22/28) of responses were classified as A1 when the visual stimulus was V1, and 32  
 200 % of responses (9/28) were classified as A1 when the visual stimulus was V2, yielding a VPI score of  
 201 47%. When decoded with A-only templates the values were 75% when V1 (22/28) and 35% when V2  
 202 (10/28), yielding a VPI of 40%. For comparison the auditory-only condition (A12) is shown in **c**, **d**,  
 203 population data showing the proportion of responses classified as A1 when the visual stimulus was  
 204 V1 or V2 when decoded with auditory-only templates or auditory visual templates. **e**, **f**, resulting VPI  
 205 scores. **h**, Mean ( $\pm$  SEM) values for these units when decoded with A-only templates, AV templates

206 (as in Fig.2) or in the absence of a visual stimulus. The green data point in **d** depicts the example in **a**,  
207 **b**.  
208

209 Analysis of recording site locations demonstrated that in the awake animals recordings in the  
210 Posterior Ectosylvian Gyrus (PEG, which contains two tonotopic secondary fields) were most strongly  
211 influenced by the visual stimulus (Fig.S3b). In anaesthetised animals the magnitude of the visual  
212 preference scores was similar to that of awake animals in the primary fields, but was not significantly  
213 different across cortical areas (Fig.S3e). In both awake and anaesthetised animals units that were  
214 classified as ‘visual-discriminating’ (see Fig.5/methods) and ‘auditory-discriminating’ were influenced  
215 by the visual stimulus, with the magnitude of the effects being greatest in the visual-discriminating  
216 units. In anaesthetised animals we confirmed using noise bursts and light flashes that a substantial  
217 proportion of visual-discriminating and auditory-discriminating units were auditory-visual (of 136  
218 visual discriminating units with a significant VPI, 19 were categorised as auditory, 39 as visual and 78  
219 as auditory-visual, of 39 auditory-discriminating units with significant VPI values 21 were auditory, 2  
220 were visual and 16 were auditory visual Fig. S3i). The ability of auditory-visual temporal coherence to  
221 enhance one sound in a mixture was observed across all cortical layers (anaesthetised dataset; layers  
222 defined by current source density analysis, see methods, Fig.S3f), but was strongest in the supra-  
223 granular layers (Fig.S3g). Finally, we observed these effects in both single and multi-units (Fig.S5a,b).

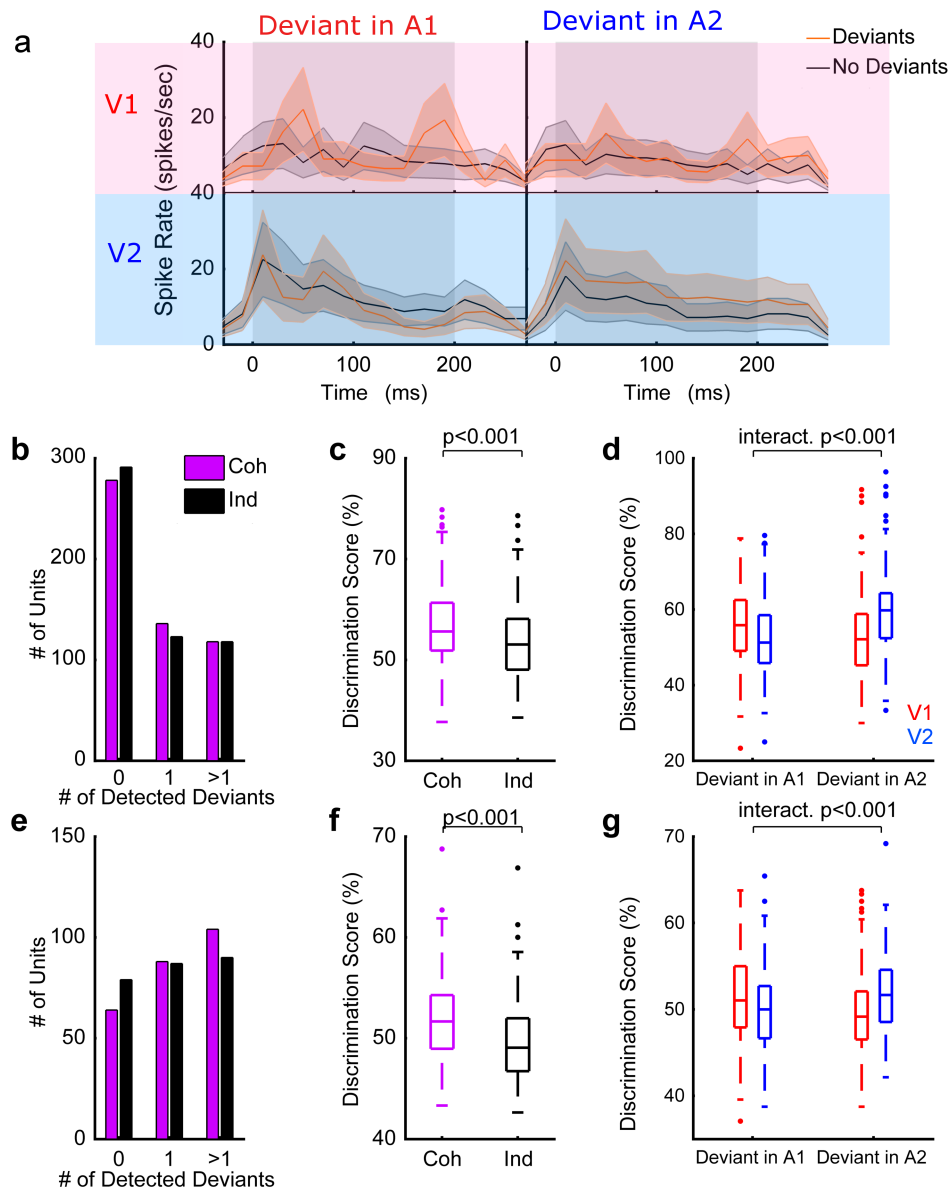
#### 224 **Auditory-visual temporal coherence enhances non-binding sound features**

225 A hallmark of an object-based rather than feature-based representation is that all stimulus features  
226 are bound into a unitary perceptual construct, including those features which do not directly  
227 mediate binding (Desimone and Duncan, 1995). We predicted that binding across modalities would  
228 be promoted via synchronous changes in auditory intensity and visual luminance (Fig.1b, S1) and  
229 observed that the temporal dynamics of the visual stimulus enhanced the representation of  
230 temporally coherent auditory streams (Fig.2c-h and 3d-f). To determine whether temporal  
231 synchrony of visual and auditory stimulus components also enhanced the representation of

232 orthogonal stimulus features and thus fulfil a key prediction of binding (Bizley et al., 2016b), we  
233 introduced brief timbre perturbations into our acoustic stimuli (two in each of the A1 and A2  
234 streams). Each deviant lasted for 200 ms during which the spectral timbre smoothly transitioned to  
235 the identity of another vowel and back to the original. It is important to note that neither the  
236 amplitude of the auditory envelope nor the visual luminance were informative about whether, or  
237 when, a change in sound timbre occurred (Fig.S1). Such timbre deviants could be detected by human  
238 listeners and were better detected when embedded in an auditory stream that was temporally  
239 coherent with an accompanying visual stimulus (Maddox et al., 2015). We hypothesised that a  
240 temporally coherent visual stimulus would enhance the representation of timbre deviants in the  
241 responses of auditory cortical neurons.

242 To isolate neural responses to the timbre change from those elicited by the on-going amplitude  
243 modulation, we extracted 200 ms epochs of the neuronal response during the timbre deviant and  
244 compared these to epochs from stimuli without deviants that were otherwise identical (Fig.S1). We  
245 observed that the spiking activity of many units differed between deviant and no-deviant trials (e.g.  
246 Fig.4a, Fig.S6) and we were able to discriminate deviant from no-deviant trials with a spike pattern  
247 classifier. For each neuron, our classifier reported both the number of deviants that could be  
248 detected (i.e. discriminated better than chance as assessed with a permutation test, the maximum is  
249 4, two per auditory stream), and a classification score (where 100% implies perfect discrimination,  
250 and 50% chance discrimination, averaged across all deviants for any unit in which at least one  
251 deviant was successfully detected). We first considered the influence of temporal coherence  
252 between auditory and visual stimuli on the representation of timbre deviants in the single stream  
253 condition (A1V1, A1V2 etc.). We found that a greater proportion of units detected at least one  
254 deviant when the auditory stream in which deviants occurred was temporally coherent with the  
255 visual stimulus relative to the temporally independent condition. This was true both for awake (Fig.  
256 4b; Pearson chi-square statistic,  $\chi^2 = 322.617$ ,  $p < 0.001$ ) and anaesthetised animals (Fig. 4e;  $\chi^2 =$   
257  $288.731$ ,  $p < 0.001$ ). For units that detected at least one deviant, discrimination scores were

258 significantly higher when accompanied by a temporally coherent visual stimulus (Fig.4c, awake  
 259 dataset, pairwise t-test  $t_{300} = 3.599$   $p < 0.001$ ; Fig. 4f, anesthetised data  $t_{262} = 4.444$   $p < 0.001$ ).



260

261 **Figure 4: Temporally coherent changes in visual luminance and auditory intensity enhance the**  
 262 **representation of auditory timbre**

263 **a** Example unit response (from the awake dataset) showing the influence of visual temporal  
 264 coherence on spiking responses to dual stream stimuli with (red PSTH) or without (black PSTH)  
 265 timbre deviants. **b-d** timbre deviant discrimination in the awake dataset. Two deviants were  
 266 included in each auditory stream giving a possible maximum of 4 per unit **b**, Histogram showing the  
 267 number of deviants (out of 4) that could be discriminated from spiking responses **c**, Box plots  
 268 showing the timbre deviant discrimination scores in the single stream condition across different  
 269 visual conditions (Coh: coherent, ind: independent). The boxes show the upper and lower quartile  
 270 values, and the horizontal lines indicates the median, the whiskers depict the most extreme data

271 points not considered to be outliers (which are marked as individual symbols). **d**, Discrimination  
272 scores for timbre deviant detection in dual stream stimuli. Discrimination scores are plotted  
273 according to the auditory stream in which the deviant occurred and the visual stream that  
274 accompanied the sound mixture. V1 stimuli are plotted in red, and V2 stimuli in blue; therefore for **d**  
275 and **g** the boxplots at the far left and right of the plot represent the cases in which the deviants  
276 occurred in an auditory stream which was temporally coherent with the visual stimulus while the  
277 central two boxplots represent the discrimination of deviants occurring in the auditory stream which  
278 was temporally independent of the visual stimulus. **e-g** show the same as **b-d** but for the  
279 anaesthetised dataset. See also supplemental figure 6.

280

281 We also observed an enhancement in the representation of timbre changes in the context of a  
282 sound scene (Fig 4d,g): timbre changes were more reliably encoded when the sound stream in which  
283 they were embedded was accompanied by a temporally coherent visual stimulus. We performed a  
284 two-way repeated measures ANOVA on deviant discrimination performance with visual condition  
285 (V1/V2) and the auditory stream in which the deviants occurred (A1/A2) as factors. We anticipated  
286 that enhancement of the representation of timbre deviants in the temporally coherent auditory  
287 stream would be revealed as a significant interaction term in the dual stream data. Significant  
288 interaction terms were seen in both the awake (Fig.4d,  $F(1,600) = 29.138$ ,  $p < 0.001$ ) and anaesthetised  
289 datasets (Fig.4g,  $F(1,524) = 16.652$ ,  $p < 0.001$ ). We also observed significant main effects of auditory  
290 and visual conditions in awake (main effect of auditory stream,  $F(1,600) = 4.565$ ,  $p = 0.033$ ; main  
291 effect of visual condition,  $F(1,600) = 2.650$ ,  $p = 0.010$ ) but not anaesthetised animals (main effect of  
292 auditory stream,  $F(1,524) = 0.004$ ,  $p = 0.948$ ; main effect of visual condition,  $F(1,524) = 1.355$ ,  $p =$   
293  $0.245$ ).

294 Finally, to determine whether a temporally coherent visual stimulus enhanced the representation of  
295 non-binding features relative to auditory-alone stimuli, we collected additional control data (3  
296 animals, 39 driven units) in which single stream stimuli were presented with, or without a temporally  
297 coherent visual stimulus. These data (Fig.S6a-c) confirmed that the presence of a visual stimulus  
298 enhanced the encoding of timbre deviants relative to the auditory-only condition. The magnitude of  
299 the influence of auditory-visual temporal coherence on timbre deviant encoding was equivalent in  
300 single and multi units (Fig.S5c,d).

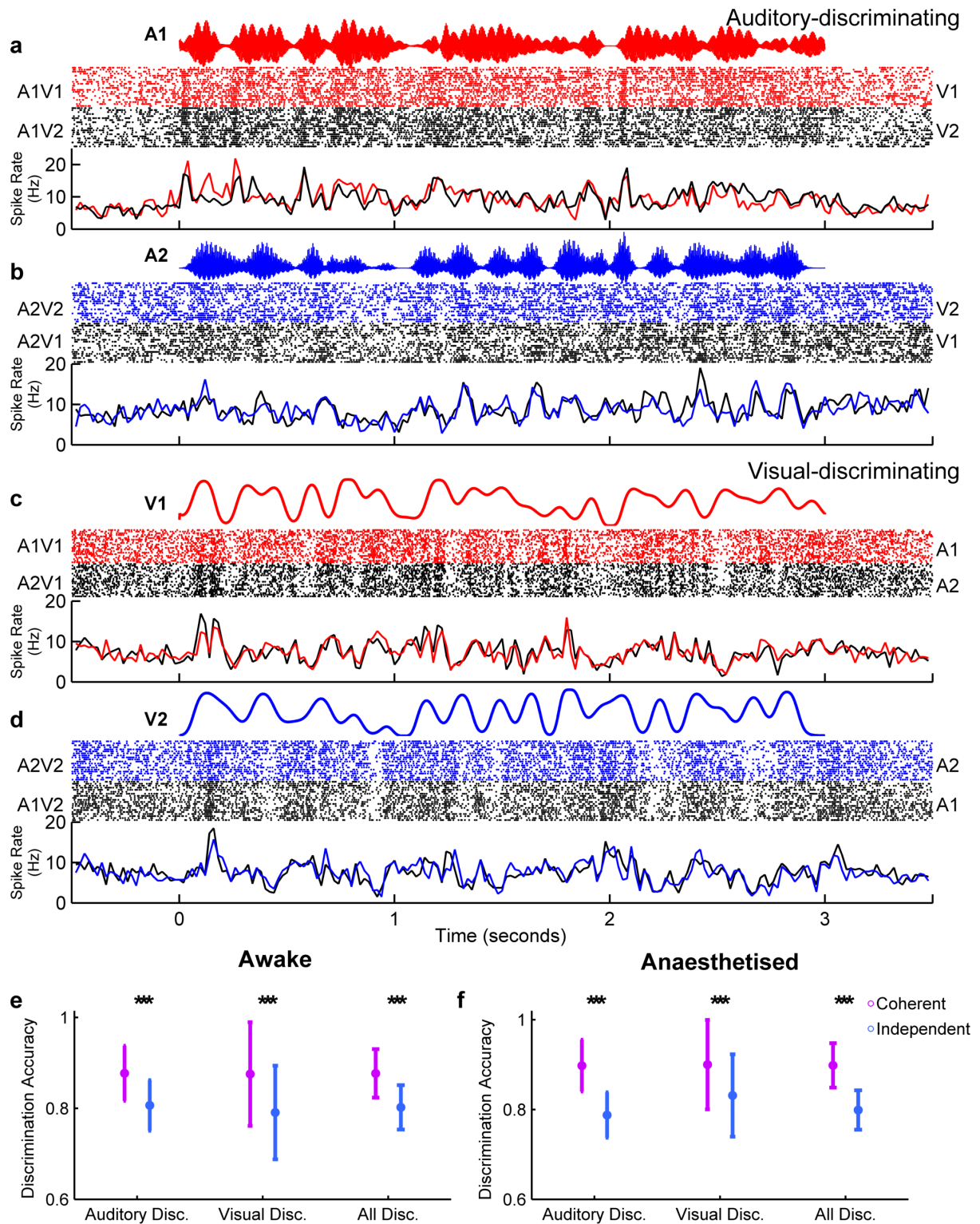
301

302 Together these data demonstrate the predicted enhancement in the neural representation of both  
303 binding (i.e. auditory amplitude) and non-binding features (here auditory timbre) that are  
304 orthogonal to those that promote binding between auditory and visual streams, meaning the effects  
305 we observe in auditory cortex fulfil our definition of multisensory binding. Next we turn to the  
306 question of how these effects are mediated, and whether they emerge within or outside of auditory  
307 cortex.

308

309 **Auditory cortical spike patterns differentiate dynamic auditory-visual stimuli more**  
310 **effectively when stimuli are temporally coherent**

311 We used the responses to single stream stimuli to classify neurons according to whether they were  
312 dominantly modulated by auditory or visual temporal dynamics. To determine whether the auditory  
313 amplitude envelope reliably modulated spiking, we used a spike-pattern classifier to decode the  
314 auditory stream identity, collapsed across visual stimulus (i.e. we decoded auditory stream identity  
315 from the combined responses to A1V1 and A1V2 stimuli and the combination of A2V1 and A2V2  
316 responses). An identical approach was taken to determine if neuronal responses reliably  
317 distinguished visual modulation (i.e. we decoded visual identity from the combined responses to  
318 A1V1 and A2V1 stimuli and the combined responses elicited by A1V2 and A2V2). Neuronal responses  
319 which were informative about auditory or visual stimulus identity at a level better than chance  
320 (estimated with a bootstrap resampling) were classified as auditory-discriminating (Fig. 5a-b) and /  
321 or visual-discriminating (Fig. 5c-d) respectively.



322

323 **Figure 5: Auditory-visual temporal coherence enhances neural coding in auditory cortex.**

324 A pattern classifier was used to determine whether neuronal responses were informative about  
 325 auditory or visual stimuli. The responses to single stream stimuli are shown for two example units,  
 326 with responses grouped according to the identity of the auditory (**a**, **b**, for an auditory discriminating  
 327 unit) or visual stream (**c**, **d**, for a visual discriminating unit). In each case the stimulus amplitude (a,b)  
 328 / luminance (c,d) waveform is shown in the top panel with the resulting raster plots and PSTHs  
 329 below. **e**, **f**: Decoder performance (mean  $\pm$  SEM) for discriminating stimulus identity (coherent: A1V1



330 vs. A2V2, purple; independent: A1V2 vs. A2V1, blue) in auditory and visual classified units recorded  
331 in awake (e) and anaesthetised (f) ferrets. Pairwise comparisons for decoding of coherent versus  
332 independent stimuli (\*\*\*) indicates  $p < 0.001$ ).

333

334 In awake animals, 39.5% (210/532) of driven units were auditory-discriminating, 11.1% (59/532)  
335 were visual-discriminating, and only 0.4% (2/532) discriminated both auditory and visual stimuli.

336 Overall a smaller proportion of units represented the identity of auditory or visual streams in the

337 anaesthetised dataset: 20.2% (242/1198) were auditory-discriminating, 6.8% (82/1198) were visual  
338 discriminating, and 0.6% (7/1198) discriminated both. Using simple noise bursts and light flashes in

339 anaesthetised animals revealed that the classification of units as visual / auditory discriminating

340 based on the single stream stimuli selected a subset of light and/or sound driven units and that the

341 proportions of auditory, visual and AV units recorded in our sample were in line with previous

342 studies from ferret auditory cortex (65.1% (328/504) of units were driven by noise bursts, 16.1%

343 (81/504) by light flashes and 14.1% (71/504) by both). When considering the units which were

344 classified as auditory or visual discriminating based on single stream stimuli, and for which we

345 recorded responses to noise bursts and light flashes, 53% (160/307) were classified as auditory, 17%

346 (53/307) as visual and 31% (94/307) as auditory-visual when classified with simple stimuli (see also

347 Fig.S3i).

348 We hypothesised that the effects we observed in the dual-stream condition might be a consequence

349 of temporal coherence between auditory and visual stimuli enhancing the discriminability of neural

350 responses. We confirmed this prediction by using the same spike pattern decoder to compare our

351 ability to discriminate temporally coherent (A1V1 vs. A2V2) and temporally independent (A1V2 vs.

352 A2V1) stimuli (Fig.5e,f): Temporally coherent AV stimuli produced more discriminable spike patterns

353 than those elicited by temporally independent ones in both awake (Fig. 5e, pairwise t-test, auditory-

354 discriminating  $t_{418} = 11.872$ ,  $p < 0.001$ ; visual-discriminating  $t_{116} = 6.338$ ,  $p < 0.001$ ; All  $t_{540} = 13.610$ ,

355  $p < 0.001$ ) and anaesthetised recordings (Fig.5f, auditory-discriminating  $t_{482} = 17.754$ ,  $p < 0.001$ ; visual-

356 discriminating  $t_{162} = 8.186$ ,  $p < 0.001$ ; All  $t_{664} = 19.461$ ,  $p < 0.001$ ). We further determined that neither

357 the mean nor maximum evoked spike rates were different between trials in response to temporally  
358 coherent and temporally independent auditory visual stimuli (Fig. S4). We also observed that the  
359 impact of auditory-visual temporal coherence was stronger in single units than multiunits in the  
360 awake dataset (Fig.S5e). Therefore the improved discrimination ability observed in response to  
361 temporally coherent auditory-visual stimuli is most likely to arise due to an increase in the reliability  
362 with which a spiking response occurred.

363

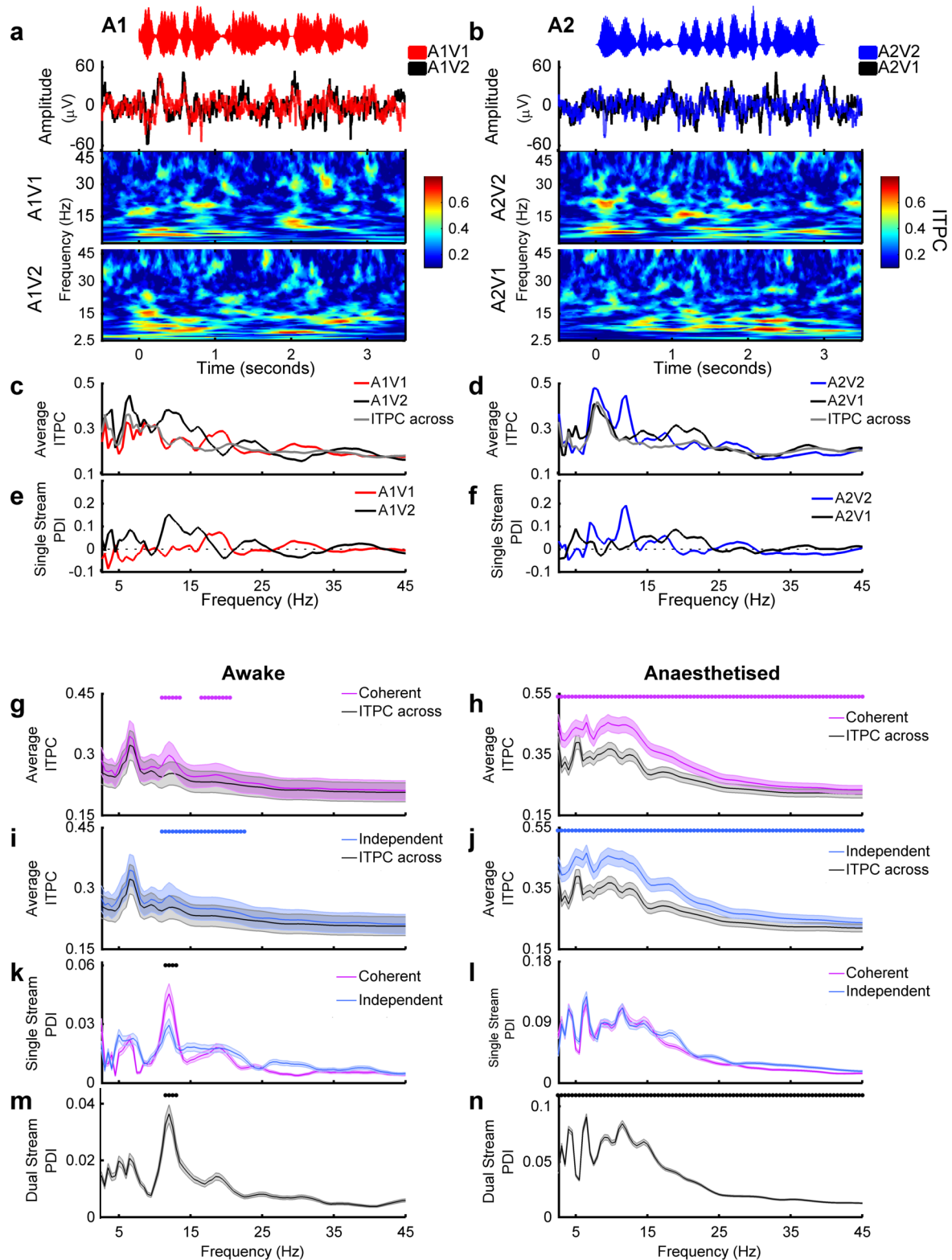
### 364 **Dynamic visual stimuli elicit reliable changes in LFP phase**

365 Temporal coherence between auditory and visual stimulus streams results in more discriminable  
366 spike trains in the single stream condition, and an enhancement of the representation of the  
367 temporally coherent sound when that sound forms part of an auditory scene. What might underlie  
368 the increased discriminability observed for temporally coherent cross-modal stimuli? The phase of  
369 on-going oscillations determines the excitability of the surrounding cortical tissue (Azouz and Gray,  
370 1999; Okun et al., 2010; Szymanski et al., 2011). LFP phase is reliably modulated by naturalistic  
371 stimulation (Chandrasekaran et al., 2010; Kayser et al., 2009; Luo and Poeppel, 2007b; Ng et al.,  
372 2012; Schyns et al., 2011) and has been implicated in multisensory processing (Golumbic et al., 2013;  
373 Lakatos et al., 2007). We hypothesised that sub-threshold visual inputs could modulate spiking  
374 activity by modifying the phase of the local field potential such that, when visual-stimulus induced  
375 changes in LFP phase coincided with auditory-stimulus evoked activity, the spiking precision in  
376 auditory cortex was enhanced.

377

378 Stimulus-evoked changes in the local field potential (LFP) were evident from the recorded voltage  
379 traces, and analysis of inter-trial phase coherence demonstrated that there were reliable changes in  
380 phase across repetitions of identical AV stimuli (Fig.6a,b). To isolate the influence of visual activity on  
381 the LFP at each recording site, and address the hypothesis that visual stimuli elicited reliable changes

382 in the LFP, we calculated phase and power dissimilarity functions for stimuli with identical auditory  
383 signals but differing visual stimuli (Luo and Poeppel, 2007b). Briefly, this analysis assumes that if the  
384 phase within a particular frequency band differs systematically between responses to two different  
385 stimuli, then inter-trial phase coherence (ITPC) across repetitions of the same stimulus will be  
386 greater than across randomly selected stimuli. For each frequency band in the LFP, we therefore  
387 compared “within-stimulus” ITPC for responses to each stimulus (A1 stream Fig. 6c; A2 stream Fig.  
388 6d) with “across-stimulus” ITPC calculated from stimuli with *identical auditory components* but  
389 randomly selected visual stimuli (e.g. randomly drawn from A1V1 and A1V2). The difference  
390 between within-stimulus and across-stimulus ITPC was then calculated across frequency and  
391 described as the phase dissimilarity index (PDI) (Fig.6e,f, single site example, k,l population data)  
392 with positive PDI values indicating reliable changes in phase coherence elicited by the visual  
393 component of the stimulus. Importantly, because both test distributions and the null distribution  
394 contain identical sounds any significant PDI value can be attributed directly to the visual component  
395 of the stimulus.



396

397 **Figure 6: Visual stimuli elicit reliable changes in the phase of the local field potential**

398 **a, b** Example LFP responses to each single stream auditory stimulus across visual conditions. Data  
 399 obtained from the recording site at which multiunit spiking activity discriminated auditory stream  
 400 identity in Fig. 5a and b. The amplitude waveforms of the stimuli are shown in the top panel, with  
 401 the evoked LFP underneath (mean across 21 trials). The resulting inter-trial phase coherence (ITPC)  
 402 values are shown in the bottom two panels (**c, d**). ITPC was calculated for coherent and independent

403 AV stimuli separately and compared to a null distribution (ITPC across). **e,f** Single stream phase  
404 dissimilarity values (PDI) were calculated by comparing ITPC within values to the ITPC across null  
405 distribution. **g,l** Population mean inter-trial phase coherence (ITPC) values across frequency for  
406 coherent (**g**, significant frequencies 10.5-13, 16-20 Hz) and independent (**i** significant frequencies  
407 10.5-22 Hz) conditions. Dots indicate frequencies at which the ITPC-within values were significantly  
408 greater than the ITPC-across values (Pairwise t-test,  $\alpha = 0.0012$ , Bonferroni corrected for 43  
409 frequencies). **k**: Mean ( $\pm$ SEM) single stream phase dissimilarity index (PDI) values for coherent and  
410 independent stimuli in animals. Black dots indicate frequencies at which the temporally coherent  
411 single stream PDI is significantly greater than in the independent conditions ( $p < 0.001$ , significant  
412 frequencies 10.5-12.5). **h,j,l** as **g,i,k** for anaesthetised dataset. **m, n**, mean  $\pm$ SEM dual stream PDI  
413 values for awake (**m**, significant frequencies 10.5-12.5) and anaesthetised (**n**) datasets.

414

415 We calculated PDI values for each of the four single stream stimuli and grouped conditions by  
416 coherency (coherent: A1V1 / A2V2, or independent: A1V2 / A2V1). To determine at which  
417 frequencies the across-trial phase reliability was significantly positive, we compared the within-  
418 stimulus values with the across-stimulus values for each frequency band (paired t-test with  
419 Bonferroni correction for 43 frequencies,  $\alpha = 0.0012$ ). In awake subjects we identified a restricted  
420 range of frequencies between 10.5 and 20 Hz where visual stimuli enhanced the phase reliability  
421 (Fig. 6g,i). In anaesthetised animals, average PDI values were larger than in awake animals and all  
422 frequencies tested had single stream PDI values that were significantly non-zero (Fig. 6h,j). We  
423 therefore conclude that visual stimulation elicited reliable changes in the LFP phase in auditory  
424 cortex. In contrast to LFP phase, a parallel analysis of across trial power reliability showed no  
425 significant effect of visual stimuli on LFP power in any frequency band (Fig.S7a,c).

426 If visual information was only conveyed in the case of temporally coherent stimuli, this might  
427 indicate that the locus of binding was outside of auditory cortex and that the information being  
428 provided to auditory cortex already reflected an integrated auditory-visual signal. The LFP is thought  
429 to reflect the combined synaptic inputs to a region (Viswanathan and Freeman 2007) and so  
430 significant PDI values for both temporally independent and coherent stimuli suggests that the  
431 correlates of binding observed in auditory cortex were not simply inherited from its inputs. Since  
432 there were significant PDI values for both temporally independent and coherent stimuli, we next  
433 asked whether there were any frequencies at which phase coherence was significantly greater in AV

434 stimuli which were temporally coherent compared to temporally independent. We performed a  
435 pairwise comparison of single stream PDI values obtained from temporally coherent and  
436 independent stimuli, for all frequency points. In awake animals, PDI values were similar for  
437 temporally coherent and temporally independent stimuli, except in the 10.5-12.5 Hz band where  
438 coherent stimuli elicited significantly greater phase coherence (Fig. 6k). In anaesthetised animals,  
439 the single stream PDI did not differ between coherent and independent stimuli at any frequency  
440 (Fig. 6l). Together these data suggests that visual inputs modulate the phase of the field potential in  
441 auditory cortex largely independently of any temporal coherence between auditory and visual  
442 stimuli. This finding supports the conjecture that multisensory binding occurs within auditory cortex.

443 To understand whether the same mechanisms could underlie the visual-stimulus induced  
444 enhancement of a temporally coherent sound in a mixture, we performed similar analyses on the  
445 data collected in response to the dual stream stimuli. We generated within-stimulus ITPC values for  
446 each dual-stream stimulus (i.e. A12V1 and A12V2) and across-stimulus ITPC by randomly selecting  
447 responses across visual conditions. We then expressed the difference as the dual stream phase  
448 dissimilarity index (dual stream PDI, Fig. 6m,n). Since the auditory components were identical in  
449 each dual stream stimulus, the influence of the visual component on LFP phase could be isolated as  
450 non-zero dual stream PDI values (paired t-test, Bonferroni corrected,  $\alpha = 0.0012$ ). In awake animals,  
451 the dual stream PDI was significantly non-zero at 10.5-12.5 (Fig.6m) whereas in anaesthetised  
452 animals, we found positive dual stream PDI values across all frequencies tested (Fig.6n). In  
453 anaesthetised animals, where we could use the responses of units to noise and light flashes to  
454 categorise units as auditory, visual or auditory-visual, we confirmed significant PDI values in the LFP  
455 recorded on the same electrode as units in each of these subpopulations (Fig.S3l). In awake animals,  
456 we tested auditory visual stimuli presented at three different modulation rates (7, 12 and 17Hz) and  
457 confirmed that significant PDI values were obtained at very similar LFP frequencies across these  
458 modulation rates - consistent with these being evoked phase alignments rather than stimulus-  
459 entrained oscillations (Fig. S7i). Additional evidence for that hypothesis comes from the fact that, in

460 the awake data, the frequencies at which the single and dual stream PDI values are significant are  
461 entirely non-overlapping with the modulation rate of the stimulus, which was band-limited to 7 Hz.

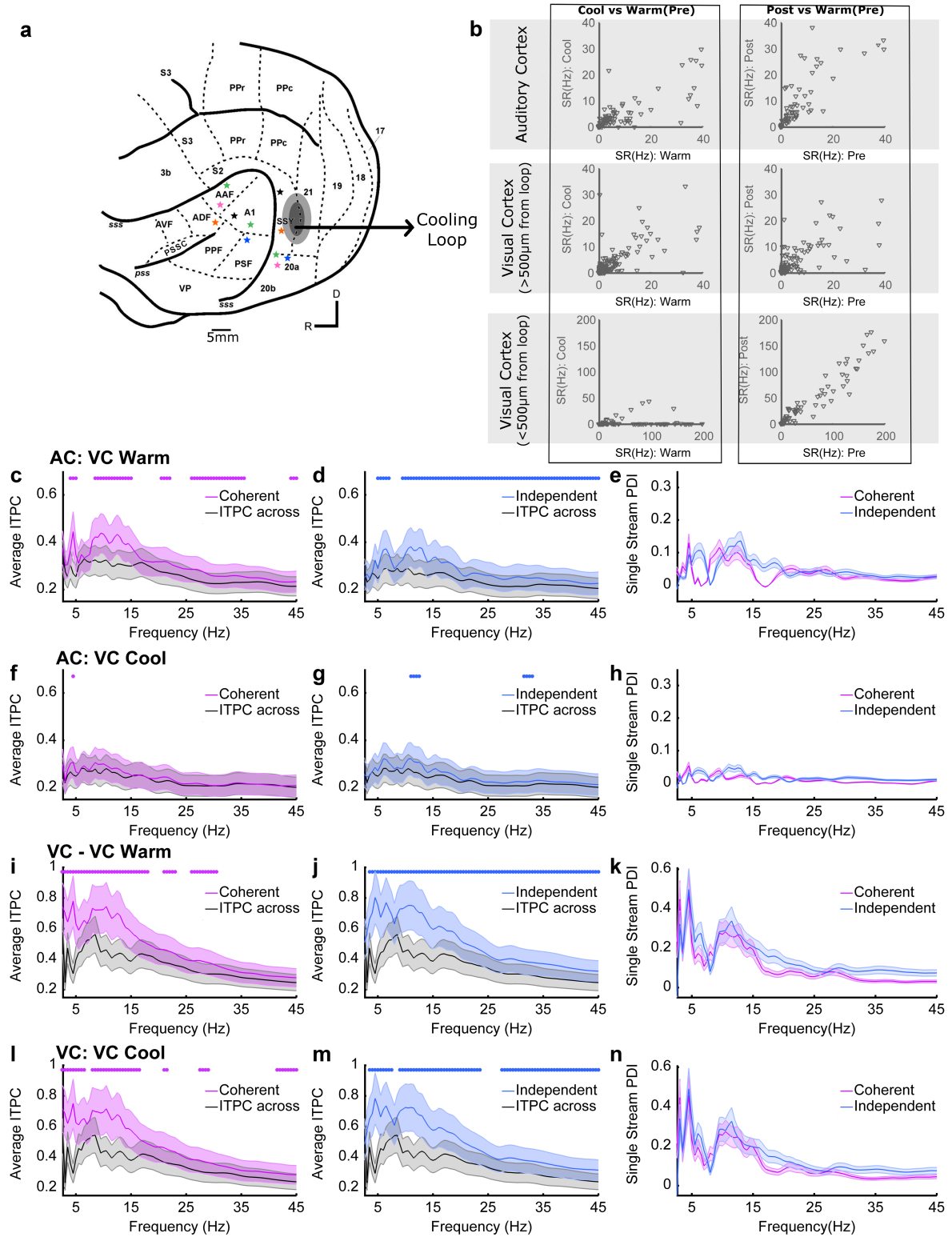
462

### 463 **Visual cortex mediates visual-stimulus induced LFP changes in auditory cortex**

464 Visual inputs to auditory cortex potentially originate from many sources: in the ferret, multiple visual  
465 cortical fields are known to innervate auditory cortex (Bizley et al., 2007), but frontal and thalamic  
466 areas are additional candidates for sources of top-down and bottom-up multisensory innervation.  
467 To determine the origin of the visual effects that we observe in auditory cortex, we performed an  
468 additional experiment in which we cooled the gyral surface of the posterior suprasylvian sulcus  
469 where visual cortical fields SSY (Cantone et al., 2006) and Area 21 (Innocenti et al., 2002) are located  
470 (Fig. 7a). Neural tracer studies have demonstrated that these areas directly project to auditory  
471 cortex in the ferret (Bizley et al., 2007). We used a cooling loop cooled to 9-10°C to reversibly silence  
472 neural activity within <500 μm of the loop (see Fig. 7b, see also Wood et al., 2017). Using simple  
473 noise bursts and light flashes at each site that we cooled, we verified that cooling visual cortex did  
474 not alter the response to noise bursts in auditory cortex (repeated measures ANOVA on spike rates  
475 in response to a noise burst pre-cooling, during cooling, after cooling,  $F(2,164) = 0.42$   $p=0.88$ ), but  
476 did reversibly attenuate the spiking response to light flashes in visual cortical sites >500 μm from the  
477 cooling loop (repeated measures ANOVA  $F_{(2,92)} = 6.83$   $p=0.001$ , post-hoc comparisons shows pre-cool  
478 and cool-post were significantly different, pre-post were not significantly different indicating the  
479 effects were reversible) and under the loop ( $F(2,210) = 30.2586$ ;  $p = 2.8350e-12$ , pre-cool vs. cooled,  
480 cooled vs. post-cooled significantly different, pre-post not significantly different). We measured  
481 responses to the single stream stimuli in auditory and visual cortex before and during cooling. From  
482 the LFP, we calculated the across-trial-phase coherence and phase dissimilarity indexes (as in Fig. 6).  
483 Cooling visual cortex significantly decreased the magnitude of the single stream PDI values in  
484 auditory cortex (Fig 7e,h). A 3-way repeated measures ANOVA with factors of visual condition

485 (coherent/independent), frequency, and cortical temperature (warm/cooled) on the SS PDI values  
486 obtained in auditory cortex showed a main effect of frequency ( $F_{(88,22605)} = 47.91, p < 0.1 \cdot 10^9$ )  
487 and temperature ( $F_{(1,22605)} = 1072, p < 0.1 \cdot 10^9$ ), but not visual condition ( $p = 0.49$ ). In contrast, LFP at  
488 recording sites in visual cortex away ( $>500 \mu\text{m}$ ) from the loop were unaffected by cooling (3-way  
489 ANOVA demonstrated that the magnitude of the SS PDI value was influenced by frequency ( $F_{(88,17265)}$   
490  $= 24.73, p < 1 \cdot 10^9$ ), but not temperature ( $p = 0.75$ ) or visual condition ( $p = 0.29$ ), Fig 7i-n). From these  
491 data, we conclude that the influence of visual stimuli on the auditory cortical field potential phase is  
492 mediated, at least in part, by inputs from visual cortical areas SSY and 21. While cooling does not  
493 allow us to confirm that visual inputs are direct mono-synaptic connections (Bizley et al., 2016a), the  
494 observation that the phase effects in other areas of visual cortex are unaffected suggests that  
495 cooling selectively influenced communication between auditory and visual cortices rather than  
496 suppressing visual processing generally.





497

498 **Figure 7: Visual stimulus induced LFP phase changes in auditory cortex are mediated in**

499 **cortex**

500 **a** Schematic showing the location of auditory cortical recording sites and the location of a cooling  
501 loop (black, grey line marks the 500  $\mu\text{m}$  radius over which cooling is effective (Wood et al., 2017)  
502 which was used to inactivate visual cortex. Individual recording sites contributing to c-n are shown  
503 with stars (simultaneous recordings are marked in the same colour). **b** Spike rate responses in  
504 Auditory Cortex (top row) and visual cortex (bottom row, sites  $>500 \mu\text{m}$  from the loop) in response  
505 to noise bursts or light flashes before, during and after cooling. **c,d**, inter-trial phase coherence  
506 values for the coherent (**c**) and independent stimuli (**d**) AV stimuli recorded in auditory cortex (AC)  
507 prior to cooling visual cortex (VC) compared to the shuffled null distribution (ITPC across). Asterisks  
508 indicate the frequencies at which the ITPC values are significantly different from the shuffled ITPC-  
509 across distribution **e** single stream Phase Dissimilarity Index values calculated from the ITPC values in  
510 **c** and **d**. **f-h** - as c-e but while visual cortex was cooled to 9 degrees. **i-n** as c-h but for sites in visual  
511 cortex  $>500\mu\text{m}$  from the cooling loop. c-h includes data from 83 sites from 6 electrode penetrations,  
512 i-n includes data from 47 sites from 5 penetrations.

513

## 514 Discussion

515 Here we provide insight into how and where auditory-visual binding occurs, and provide evidence  
516 that this effect is mediated by cortico-cortical interactions between visual and auditory cortex. While  
517 numerous studies have reported the incidence of auditory-visual interactions in auditory cortex over  
518 the past decade (Bizley et al., 2007a; Chandrasekaran et al., 2013; Ghazanfar et al., 2005; Kayser et  
519 al., 2008; Kayser et al., 2010, Perrodin et al., 2015), evidence for the functional role has remained  
520 less apparent. Here we show that one role for the early integration of auditory and visual  
521 information is to support auditory scene analysis. Visual stimuli elicit reliable changes in the phase of  
522 the local field potential in auditory cortex, irrespective of auditory-visual temporal coherence,  
523 indicating that the inputs to auditory cortex reflect the unisensory stimulus properties. When the  
524 visual and auditory stimuli are temporally aligned, phase resets elicited by the visual stimulus  
525 interacts with feed-forward sound-evoked activity and results in spiking output that more precisely  
526 represents the temporally coherent sound within auditory cortex. These results are consistent with  
527 the binding of cross-modal information to form a multisensory object because they result in a  
528 modification of the representation of the sound which is not restricted to the features that link  
529 auditory and visual signals but extends to other non-binding features. These data provide a  
530 physiological underpinning for the pattern of performance observed in human listeners performing

531 an auditory selective attention task, in which the detection of a perturbation in a stimulus stream is  
532 enhanced or impaired when a visual stimulus is temporally coherent with the target or masker  
533 auditory stream respectively (Maddox et al., 2015). The effects of the visual stimulus on the  
534 representation of an auditory scene can be observed in anaesthetised animals suggesting that these  
535 effects can occur independently of attentional modulation.

536

537 Previous investigations of the impact of visual stimuli on auditory scene analysis have frequently  
538 used speech stimuli. Being able to see a talker's mouth provides listeners with information about the  
539 rhythm and amplitude of the speech waveform which may help listeners by cueing them to pay  
540 attention to the auditory envelope (Pelle and Sommers, 2015) as well as by providing cues to the  
541 place of articulation that can disambiguate different consonants (Sumbly and Pollack, 1954).  
542 However, the use of speech stimuli makes it difficult to dissociate general multisensory mechanisms  
543 from speech-specific ones when testing in human subjects. Therefore, in order to probe more  
544 general principles across both human (Maddox et al., 2015) and non-human animals (here), we  
545 chose to employ continuous naturalistic non-speech stimuli that utilized modulation rates that fell  
546 within the range of syllable rates in human speech but lacked any linguistic content. Previous work  
547 has demonstrated that a visual stimulus can enhance the neural representation of the speech  
548 amplitude envelope both in quiet and in noise (Crosse et al., 2015; Crosse et al., 2016; Luo et al.,  
549 2010, Park et al., 2016), but functional imaging methods make it difficult to demonstrate enhanced  
550 neural encoding of features beyond the amplitude envelope. The implication of our findings is that  
551 representation of the spectro-temporal features that allow speech recognition such as voice pitch  
552 would be enhanced in auditory cortex when a listener views a talker's face, even though such  
553 spectro-temporal features may not be represented by the visual stimulus.

554 Visual speech information is hypothesised to be relayed to auditory cortex through multiple routes  
555 in parallel to influence the processing of auditory speech: Our data support the idea that early

556 integration of visual information occurs (Möttönen et al., 2004; Okada et al., 2013; Peelle and  
557 Sommers, 2015; Schroeder et al., 2008) and is likely to reflect a general phenomenon whereby visual  
558 stimuli can cause phase-entrainment in the local field potential. Within this framework, cross-modal  
559 binding potentially results from the temporal coincidence of evoked auditory responses and visual-  
560 stimulus elicited inputs that we observe as phasic changes of the LFP.

561

562 Consistent with previous studies, our analysis of local field potential activity revealed that visual  
563 information reliably modulated LFP phase in auditory cortex (Chandrasekaran et al., 2013; Ghazanfar  
564 et al., 2005; Kayser et al., 2008; Perrodin et al., 2015). This occurred independently of the  
565 modulation frequency of the stimulus suggesting that, rather than entraining oscillations at the  
566 stimulus modulation rate, relatively broadband phase resets are triggered by particular features  
567 within the stream (presumably points at which the luminance changed rapidly from low-high  
568 amplitude). LFP reflects the synaptic inputs to a region and LFP phase synchronization is thought to  
569 arise from fluctuating inputs to cortical networks (Lakatos et al., 2007; Mazzone et al., 2008;  
570 Szymanski et al., 2011). Since neuronal excitability varies with LFP phase (Jacobs et al., 2007;  
571 Klimesch et al., 2007; Lakatos et al., 2013; Lőrincz et al., 2009), synaptic inputs from visual cortex  
572 may provide a physiological mechanism through which temporally coincident cross-sensory  
573 information is integrated. Our analysis allowed us to isolate changes in LFP phase that were directly  
574 attributable to the visual stimulus and identify reliable changes in LFP phase irrespective of whether  
575 the visual stimulus was temporally coherent with the auditory stimulus. Such results suggest that the  
576 observed effects of cross-modal temporal coherence were not simply inherited within the inputs to  
577 auditory cortex. Moreover the effects that we observed in the LFP were lost when we silenced visual  
578 cortex, indicating that inputs from visual cortex are a key contributor to the effects of auditory-visual  
579 temporal coherence that we observed in auditory cortex. Our finding that visual stimulation elicited  
580 reliable phase modulation in both awake and anesthetised animals suggests that bottom-up cross-

581 modal integration interacts with selective attention, which has also been associated with modulation  
582 of phase information in auditory cortex (Golombic et al., 2013, Park et al., 2016). While our data  
583 suggest that cross-modal binding can occur in the absence of attention, it is likely that the additional  
584 neural pathways engaged during selective attention act to further enhance the representation of  
585 attended cross-modal objects.

586 In both awake and anaesthetised animals we observed that visual stimuli elicit robust effects on the  
587 LFP phase, that auditory-visual temporal coherence shapes the response to a sound mixture such  
588 that temporally coherent auditory-visual stimuli are more reliably represented in the spiking  
589 response, and that the spiking response to auditory timbre deviants (a non-binding feature) was  
590 enhanced. While these key findings were recapitulated in both states, there were some important  
591 differences: Firstly, in the awake animal the phase alignment in the LFP was generally smaller in  
592 magnitude and was only significantly modulated across a smaller range of frequencies (10.5-20 Hz as  
593 opposed to 4-45 Hz). Such differences are consistent with a dependence of oscillatory activity on  
594 behavioural state (Tukker et al., 2007; Voloh and Womelsdorf, 2016; Wang, 2010). Secondly, in the  
595 awake animal we observed a significant increase in the phase reliability (at 10.5-12.5 Hz) for  
596 temporally coherent auditory-visual stimuli when compared to temporally independent stimuli.  
597 Since the neural correlates of multisensory binding are evident in the anaesthetised animal, the  
598 specific increase in alpha phase reliability that occurred only in awake animals in response to  
599 temporally coherent auditory-visual stimulus pairs (Fig. 6k,m) may indicate an attention-related  
600 signal triggered by temporal coherence between auditory and visual signals, or an additional top-  
601 down signal conveying cross-modal information. Phase resetting or synchronisation of alpha phase  
602 has been associated both with enhanced functional connectivity (Voloh and Womelsdorf, 2016) and  
603 as a top-down predictive signal for upcoming visual information (Samaha et al., 2015).  
604 Understanding how attention engages additional brain networks and disambiguating these  
605 possibilities would require simultaneous recordings in auditory and visual cortex recording while  
606 trained animals performed the auditory selective attention task which motivated this study. Finally,

607 in awake and anaesthetised animals we observed that the impact of auditory-visual temporal  
608 coherence on the representation of sound mixtures (as assessed by visual preference scores) was of  
609 a similar magnitude in the primary areas (A1 and AAF, located in the MEG). In contrast, in the awake  
610 animal, neurons in the PEG, where secondary tonotopic fields PPF and PSF are located, had  
611 significantly higher VPI scores than those in the MEG, while in anaesthetised animals VPI scores were  
612 statistically indistinguishable across cortical fields. This suggests that in the awake animal additional  
613 mechanisms exist to enhance the effects that are present in the primary areas. These results were  
614 mirrored in the impact of auditory-visual temporal coherence on non-binding features (as assessed  
615 by the impact of auditory-visual temporal coherence on deviant detection ability) where the visual  
616 stimulus had a stronger influence in PEG than MEG in the awake animal, and did not differ across  
617 regions (and was overall of a smaller magnitude) in anaesthetised animals. Our cooling studies (in  
618 anaesthetised animals) do not allow us to determine whether this enhancement reflects the greater  
619 variety of inputs from visual cortex that terminate in secondary as opposed to primary auditory  
620 cortex (Bizley et al., 2007), top down inputs from higher areas (e.g. parietal or frontal cortex), or are  
621 a consequence of intracortical processing within auditory cortex.

622 Temporal coherence between sound elements has been proposed as a fundamental organising  
623 principle for auditory cortex (Elhilali et al., 2009; O'Sullivan et al., 2015b) and here we extend this  
624 principle to the formation of cross-modal constructs. Our data provide evidence that one role for the  
625 early integration of visual information into auditory cortex is to resolve competition between  
626 multiple sound sources within an auditory scene and that these neural computations occur pre-  
627 attentively. While some proponents of a temporal coherence based model for auditory streaming  
628 have stressed the importance of attention in auditory stream formation (Lu et al., 2017), neural  
629 signatures of temporal-coherence based streaming are present in passively listening subjects  
630 (O'Sullivan et al., 2015a; Teki et al., 2016). Previous studies have demonstrated a role for visual  
631 information in conveying lip movement information to auditory cortex (Chandrasekaran et al., 2013;  
632 Crosse et al., 2015; Ghazanfar et al., 2005; Golumbic et al., 2013), but such stimuli make it difficult to

633 separate sensory information from linguistic cues. Our data obtained using non-speech stimuli  
634 provide evidence that at least part of the boost provided by visualising a speaker's mouth arises  
635 from a more general (language-independent) phenomenon whereby visual temporal cues facilitate  
636 auditory scene analysis through the formation of cross-sensory objects. Our data are supportive of  
637 visual cortical areas providing at least one source of information. Other visual cortex fields and sub-  
638 cortical structures innervate tonotopic auditory cortex (Bizley et al., 2007; Budinger et al., 2006) and  
639 may potentially provide additional visual inputs to auditory cortex. Further dissecting the origin of  
640 visual innervation requires experiments that allow pathway specific manipulation of neuronal  
641 activity (for example by silencing the terminal fields of neurons that project from a candidate area  
642 into auditory cortex, Bizley et al., 2016a).

643 Finally, the neural correlates of multisensory binding were apparent in units which best  
644 discriminated either the auditory or visual characteristics of single auditory-visual streams, although  
645 the magnitude of the effects was stronger in visual-discriminating units. Nevertheless, both classes  
646 of neurons showed enhanced encoding of temporally coherent versus temporally independent  
647 auditory visual streams suggesting that both subgroups could be described as auditory-visual. This  
648 was confirmed in anaesthetised animals where neurons were additionally characterised with simple  
649 stimuli (noise bursts and light flashes) and revealed that 54% of visual-discriminating stimuli and 41%  
650 of auditory-discriminating neurons were classified as auditory-visual. Together, these results  
651 suggest that multisensory processing is prevalent throughout auditory cortex and that cross-sensory  
652 processing has the potential to have a significant impact on the representation of acoustic features  
653 in auditory cortex.

654 In summary, activity in auditory cortex was reliably affected by visual stimulation in a manner that  
655 enhanced the representation of temporally coherent auditory information. Enhancement of auditory  
656 information was observed for sounds presented alone or in a mixture and for sound features that  
657 were related to (amplitude) and orthogonal to (timbre) variation in visual input. Such processes

658 provide mechanistic support for a coherence-based model of cross-modal binding in object  
659 formation and indicate that one role for the early integration of visual information in auditory cortex  
660 is to support auditory scene analysis.



661

## 662 Acknowledgments

663 This work was funded by grants to each author: JKB: Wellcome Trust / Royal Society WT098418MA;  
664 Biotechnology and Biological Sciences Research Council (BB/H016813/1), and an Action on Hearing  
665 Loss Studentship (596: UEI: JB); RKM: NIH R00DC014288 and Hearing Health Foundation Emerging  
666 Research Grant; AKCL: NIH R01DC013260; and an International Exchanges Scheme award from the  
667 Royal Society to JKB and AKCL.

## 668 Author contributions

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

## 671 References

672 Azouz, R., and Gray, C.M. (1999). Cellular mechanisms contributing to response variability of cortical  
673 neurons in vivo. *The Journal of neuroscience* 19, 2209-2223.  
674 Bizley, J.K, Jones, G.P., and Town, S.M. (2016a) Where are multisensory signals combined for  
675 perceptual decision-making? *Current opinion in neurobiology*, 40, 31-37.  
676 Bizley, J.K., Maddox, R.K., and Lee, A.K. (2016b). Defining Auditory-Visual Objects: Behavioral Tests  
677 and Physiological Mechanisms. *Trends in Neurosciences*. 39(2), 74-85.  
678 Bizley, J.K., Nodal, F.R., Bajo, V.M., Nelken, I., and King, A.J. (2007). Physiological and anatomical  
679 evidence for multisensory interactions in auditory cortex. *Cereb Cortex* 17, 2172-2189.  
680 Bizley, J.K., Walker, K.M.M., King, A.J., and Schnupp, J.W.H. (2013). Spectral timbre perception in  
681 ferrets: Discrimination of artificial vowels under different listening conditions. *J Acoust Soc Am* 133,  
682 365-376.  
683 Bizley, J.K., Walker, K.M.M., Silverman, B.W., King, A.J., and Schnupp, J.W.H. (2009). Interdependent  
684 Encoding of Pitch, Timbre, and Spatial Location in Auditory Cortex. *Journal of Neuroscience* 29, 2064-  
685 2075.  
686 Brosch, M., Selezneva, E., and Scheich, H. (2015). Neuronal activity in primate auditory cortex during  
687 the performance of audiovisual tasks. *European Journal of Neuroscience* 41, 603-614.  
688 Budinger, E., Heil, P., Hess, A., and Scheich, H. (2006). Multisensory processing via early cortical  
689 stages: connections of the primary auditory cortical field with other sensory systems. *Neuroscience*  
690 143, 1065-1083.  
691 Cantone, G., Xiao, J., and Levitt, J.B. (2006). Retinotopic organization of ferret suprasylvian cortex.  
692 *Visual neuroscience* 23, 61-77.  
693 Chandrasekaran, C., Lemus, L., and Ghazanfar, A.A. (2013). Dynamic faces speed up the onset of  
694 auditory cortical spiking responses during vocal detection. *Proceedings of the National Academy of*  
695 *Sciences* 110, E4668-E4677.

- 696 Chandrasekaran, C., Trubanova, A., Stillittano, S., Caplier, A., and Ghazanfar, A.A. (2009). The natural  
697 statistics of audiovisual speech. *Plos Comput Biol* 5, e1000436.
- 698 Chandrasekaran, C., Turesson, H.K., Brown, C.H., and Ghazanfar, A.A. (2010). The influence of natural  
699 scene dynamics on auditory cortical activity. *The Journal of Neuroscience* 30, 13919-13931.
- 700 Crosse, M.J., Butler, J.S., and Lalor, E.C. (2015). Congruent visual speech enhances cortical  
701 entrainment to continuous auditory speech in noise-free conditions. *The Journal of Neuroscience* 35,  
702 14195-14204.
- 703 Crosse, M.J., Di Liberto, G.M., and Lalor, E.C. (2016). Eye Can Hear Clearly Now: Inverse Effectiveness  
704 in Natural Audiovisual Speech Processing Relies on Long-Term Crossmodal Temporal Integration. *The*  
705 *Journal of Neuroscience* 36, 9888-9895.
- 706 Denison, R.N., Driver, J., and Ruff, C.C. (2013). Temporal structure and complexity affect audio-visual  
707 correspondence detection. *Front Psychol* 3.
- 708 Desimone, R., and Duncan, J. (1995). Neural mechanisms of selective visual attention. *Annual review*  
709 *of neuroscience* 18, 193-222.
- 710 Elhilali, M., Ma, L., Micheyl, C., Oxenham, A.J., and Shamma, S.A. (2009). Temporal Coherence in the  
711 Perceptual Organization and Cortical Representation of Auditory Scenes. *Neuron* 61, 317-329.
- 712 Foffani, G., and Moxon, K.A. (2004). PSTH-based classification of sensory stimuli using ensembles of  
713 single neurons. *J Neurosci Methods* 135, 107-120.
- 714 Ghazanfar, A.A., Maier, J.X., Hoffman, K.L., and Logothetis, N.K. (2005). Multisensory integration of  
715 dynamic faces and voices in rhesus monkey auditory cortex. *The Journal of Neuroscience* 25, 5004-  
716 5012.
- 717 Golumbic, E.Z., Cogan, G.B., Schroeder, C.E., and Poeppel, D. (2013). Visual input enhances selective  
718 speech envelope tracking in auditory cortex at a “cocktail party”. *The Journal of Neuroscience* 33,  
719 1417-1426.
- 720 Henry, M.J., and Obleser, J. (2012). Frequency modulation entrains slow neural oscillations and  
721 optimizes human listening behavior. *Proceedings of the National Academy of Sciences* 109, 20095-  
722 20100.
- 723 Innocenti, G.M., Manger, P.R., Masiello, I., Colin, I., and Tettoni, L. (2002). Architecture and callosal  
724 connections of visual areas 17, 18, 19 and 21 in the ferret (*Mustela putorius*). *Cereb Cortex* 12, 411-  
725 422.
- 726 Jacobs, J., Kahana, M.J., Ekstrom, A.D., and Fried, I. (2007). Brain oscillations control timing of single-  
727 neuron activity in humans. *The Journal of neuroscience* 27, 3839-3844.
- 728 Kaur, S., Rose, H., Lazar, R., Liang, K., and Metherate, R. (2005). Spectral integration in primary  
729 auditory cortex: laminar processing of afferent input, in vivo and in vitro. *Neuroscience* 134, 1033-  
730 1045.
- 731 Kayser, C., Petkov, C.I., and Logothetis, N.K. (2008). Visual modulation of neurons in auditory cortex.  
732 *Cerebral Cortex* 18, 1560-1574.
- 733 Kayser, C., Petkov, C.I., and Logothetis, N.K. (2009). Multisensory interactions in primate auditory  
734 cortex: fMRI and electrophysiology. *Hearing Res* 258, 80-88.
- 735 Kayser, C., Logothetis, N.K., Panzeri, S. Visual enhancement of the visual information representation in  
736 auditory cortex **Curr Biol.** 2010 Jan 12;20(1):19-24.
- 737 Klimesch, W., Sauseng, P., and Hanslmayr, S. (2007). EEG alpha oscillations: The inhibition–timing  
738 hypothesis. *Brain Research Reviews* 53, 63-88.
- 739 Lakatos, P., Chen, C.-M., O'Connell, M.N., Mills, A., and Schroeder, C.E. (2007). Neuronal oscillations  
740 and multisensory interaction in primary auditory cortex. *Neuron* 53, 279-292.
- 741 Lakatos, P., Musacchia, G., O'Connell, M.N., Falchier, A.Y., Javitt, D.C., and Schroeder, C.E. (2013). The  
742 spectrotemporal filter mechanism of auditory selective attention. *Neuron* 77, 750-761.
- 743 Lőrincz, M.L., Kékesi, K.A., Juhász, G., Crunelli, V., and Hughes, S.W. (2009). Temporal framing of  
744 thalamic relay-mode firing by phasic inhibition during the alpha rhythm. *Neuron* 63, 683-696.
- 745 Lu, K., Xu, Y., Yin, P., Oxenham, A.J., Fritz, J.B., and Shamma, S.A. (2017). Temporal coherence  
746 structure rapidly shapes neuronal interactions. *Nature communications* 8, 13900.

- 747 Luo, H., Liu, Z., and Poeppel, D. (2010). Auditory cortex tracks both auditory and visual stimulus  
748 dynamics using low-frequency neuronal phase modulation. *PLoS Biol* 8, e1000445.
- 749 Luo, H., and Poeppel, D. (2007a). Phase patterns of neuronal responses reliably discriminate speech  
750 in human auditory cortex. *Neuron* 54, 1001-1010.
- 751 Luo, H., and Poeppel, D. (2007b). Phase patterns of neuronal responses reliably discriminate speech  
752 in human auditory cortex. *Neuron* 54, 1001-1010.
- 753 Maddox, R.K., Atilgan, H., Bizley, J.K., and Lee, A.K. (2015). Auditory selective attention is enhanced  
754 by a task-irrelevant temporally coherent visual stimulus in human listeners. *Elife* 4, e04995.
- 755 Mazzone, A., Panzeri, S., Logothetis, N.K., and Brunel, N. (2008). Encoding of naturalistic stimuli by  
756 local field potential spectra in networks of excitatory and inhibitory neurons. *PLoS Comput Biol* 4,  
757 e1000239.
- 758 Möttönen, R., Schürmann, M., and Sams, M. (2004). Time course of multisensory interactions during  
759 audiovisual speech perception in humans: a magnetoencephalographic study. *Neuroscience letters*  
760 363, 112-115.
- 761 Ng, B.S.W., Schroeder, T., and Kayser, C. (2012). A precluding but not ensuring role of entrained low-  
762 frequency oscillations for auditory perception. *The Journal of Neuroscience* 32, 12268-12276.
- 763 O'Sullivan, J.A., Shamma, S.A., and Lalor, E.C. (2015a). Evidence for Neural Computations of  
764 Temporal Coherence in an Auditory Scene and Their Enhancement during Active Listening. *J*  
765 *Neurosci* 35, 7256-7263.
- 766 O'Sullivan, J.A., Shamma, S.A., and Lalor, E.C. (2015b). Evidence for neural computations of temporal  
767 coherence in an auditory scene and their enhancement during active listening. *The Journal of*  
768 *Neuroscience* 35, 7256-7263.
- 769 Okada, K., Venezia, J.H., Matchin, W., Saberi, K., and Hickok, G. (2013). An fMRI study of audiovisual  
770 speech perception reveals multisensory interactions in auditory cortex. *Plos One* 8, e68959.
- 771 Okun, M., Naim, A., and Lampl, I. (2010). The subthreshold relation between cortical local field  
772 potential and neuronal firing unveiled by intracellular recordings in awake rats. *The Journal of*  
773 *neuroscience* 30, 4440-4448.
- 774 Park, H., Kayser, C., Thut, G., and Gross, J. (2016). Lip movements entrain the observers' low-  
775 frequency brain oscillations to facilitate speech intelligibility. *eLife* 5.
- 776 Peelle, J.E., and Sommers, M.S. (2015). Prediction and constraint in audiovisual speech perception.  
777 *Cortex* 68, 169-181.
- 778 Perrodin, C., Kayser, C., Logothetis, N.K., and Petkov, C.I. (2015). Natural asynchronies in audiovisual  
779 communication signals regulate neuronal multisensory interactions in voice-sensitive cortex.  
780 *Proceedings of the National Academy of Sciences* 112, 273-278.
- 781 Rahne, T., Deike, S., Selezneva, E., Brosch, M., König, R., Scheich, H., Bockmann, M., and Brechmann,  
782 A. (2008). A multilevel and cross-modal approach towards neuronal mechanisms of auditory  
783 streaming. *Brain Research* 1220, 118-131.
- 784 Samaha, J., Bauer, P., Cimaroli, S., and Postle, B.R. (2015). Top-down control of the phase of alpha-  
785 band oscillations as a mechanism for temporal prediction. *Proceedings of the National Academy of*  
786 *Sciences* 112, 8439-8444.
- 787 Schroeder, C.E., Lakatos, P., Kajikawa, Y., Partan, S., and Puce, A. (2008). Neuronal oscillations and  
788 visual amplification of speech. *Trends in cognitive sciences* 12, 106-113.
- 789 Schyns, P.G., Thut, G., and Gross, J. (2011). Cracking the code of oscillatory activity. *PLoS Biol* 9,  
790 e1001064.
- 791 Sumbly, W.H., and Pollack, I. (1954). Visual contribution to speech intelligibility in noise. *The journal*  
792 *of the acoustical society of america* 26, 212-215.
- 793 Szymanski, F.D., Rabinowitz, N.C., Magri, C., Panzeri, S., and Schnupp, J.W. (2011). The laminar and  
794 temporal structure of stimulus information in the phase of field potentials of auditory cortex. *The*  
795 *Journal of Neuroscience* 31, 15787-15801.
- 796 Teki, S., Barascud, N., Picard, S., Payne, C., Griffiths, T.D., and Chait, M. (2016). Neural Correlates of  
797 Auditory Figure-Ground Segregation Based on Temporal Coherence. *Cereb Cortex* 26, 3669-3680.

798 Town, S. M., Atilgan, H., Wood, K. C., & Bizley, J. K. (2015). The role of spectral cues in timbre  
799 discrimination by ferrets and humans. *J Acoust Soc Am* 137(5), 2870-2883.  
800 Town, S. M., Wood, K. C., & Bizley, J. K. (2017). Neural correlates of perceptual constancy in Auditory  
801 Cortex. *BioXriv* <https://www.biorxiv.org/content/early/2017/07/15/102889>  
802 Tukker, J.J., Fuentealba, P., Hartwich, K., Somogyi, P., and Klausberger, T. (2007). Cell type-specific  
803 tuning of hippocampal interneuron firing during gamma oscillations in vivo. *The journal of*  
804 *neuroscience* 27, 8184-8189.  
805 Viswanathan, A., and Freeman, R.D. (2007). Neurometabolic coupling in cerebral cortex reflects  
806 synaptic more than spiking activity. *Nat. Neurosci.* 10, 1308–1312.  
807 Voloh, B., and Womelsdorf, T. (2016). A Role of Phase-Resetting in Coordinating Large Scale Neural  
808 Networks During Attention and Goal-Directed Behavior. *Frontiers in systems neuroscience* 10.  
809 Wang, X.-J. (2010). Neurophysiological and computational principles of cortical rhythms in cognition.  
810 *Physiological reviews* 90, 1195-1268.  
811 Wood, K.C., Town, S.M., Atilgan, H., Jones, A.K., and Bizley, J.K. (2017). Acute Inactivation of Primary  
812 Auditory Cortex Causes a Sound Localisation Deficit in Ferrets. *PLoS One* DOI:  
813 10.1371/journal.pone.0170264.  
814 Zion Golumbic, E., Cogan, G.B., Schroeder, C.E., and Poeppel, D. (2013). Visual input enhances  
815 selective speech envelope tracking in auditory cortex at a "cocktail party". *J Neurosci* 33, 1417-1426.

816

817

## 818 **Figure Legends**

### 819 **Figure 1: Hypothesis and experimental design**

820 **a** Conceptual model illustrating how binding can be identified as a distinct form of multisensory  
821 integration. Multisensory binding is defined as a subset of multisensory integration that results in  
822 the formation of a crossmodal object. During binding, all features of the audio-visual object are  
823 linked and enhanced - including both those features that bind the stimuli across modalities (here  
824 temporal coherence between auditory (A) intensity and visual (V) luminance) and orthogonal  
825 features such as auditory pitch and timbre, and visual colour and size. Other forms of multisensory  
826 integration would result in enhancement of only the features that promote binding - here auditory  
827 intensity and visual luminance. To identify binding therefore requires a demonstration that non-  
828 binding features (e.g. here pitch, timbre, colour or size) are enhanced. Enhanced features are  
829 highlighted in yellow. **b** When two competing sounds (red and blue waveforms) are presented they  
830 can be separated on the basis of their features, but may elicit overlapping neuronal representations  
831 in auditory cortex. **c** Hypothesised enhancement in auditory stream segregation when a temporally  
832 coherent visual stimulus enables multisensory binding. When the visual stimulus changes coherently  
833 with the red sound (A1, top) this sound is enhanced and the two sources are better segregated.  
834 Perceptually this would result in more effective auditory scene analysis and an enhancement of the  
835 non-binding features. **d** Stimulus design: Auditory stimuli were two artificial vowels (denoted A1 and  
836 A2), each with distinct pitch and timbre and independently amplitude modulated with a noisy low  
837 pass envelope. **e** Visual stimulus: a luminance modulated white light was presented with one of two  
838 temporal envelopes derived from the amplitude modulations of A1 and A2. **f** illustrates the stimulus  
839 combinations that were tested experimentally in *single stream* (a single auditory visual pair) and  
840 *dual stream* (two sounds and one visual stimulus) conditions. See also supplemental figure 1.

841

842

843 **Figure 2: Visual stimuli can determine which sound stream auditory cortical neurons follow in a**  
844 **mixture.**

845 Spiking responses from an example unit in response to **a**, single stream AV stimuli used as decoding templates  
846 and **b**, dual stream stimuli. In each case rasters and PSTHs are illustrated. When the visual component of the  
847 dual stream was V1, the majority of trials were classified as A1V1 (82% (19/23 trials), and A2V2 when the  
848 visual stimulus was V2 (26% 6/23, trials) of responses classified as A1V1 (see also green data point in c),  
849 yielding a visual preference score of 56%. **c-h** population data for awake (**c,d,e** 271 units) and anesthetised  
850 (**f,g,h** 331 units) datasets. In each case the left panel (**c,f**) shows the distribution of decoding values according  
851 to the visual condition, the middle panel (**d,g**) shows the population mean ( $\pm$  SEM) projecting onto the vertical  
852 axis of panel c / f for V1 condition, and horizontal axis of panel c / f for the V2 condition. **e,h** shows the visual  
853 preference index (VPI). Units in which the VPI was significantly  $>0$  are coloured purple. Pairwise comparisons  
854 revealed significant effect of visual condition on decoding in all datasets: Awake: All:  $t_{540}=6.1, p=2.3e-09$   
855 (n =271), Sig VPI:  $t_{180}=18.8, p = 2.0e-44$  (n=91). Anesthetised: All:  $t_{660}=9.5, p=3.3e-20$  (n=331), Sig. VPI:  $t_{348} =$   
856  $38.9, p = 1.2e-128$  (n =175) See also supplemental figures 2-4.

857

858 **Figure 3: Visual stimuli shape the neural representation of an auditory scene.**

859 In an additional control experiment (n=89 units recorded in awake animals), the responses to  
860 coherent AV and auditory-only (A Only) single stream stimuli were used as templates to decode dual  
861 stream stimuli either accompanied by visual stimuli (V1/V2) or in the absence of visual stimulation  
862 (no visual). Spiking responses from an example unit in response to **a**, single stream auditory stimuli  
863 which were used as decoding templates to decode the responses to dual stream stimuli in **b**, in each  
864 case the auditory waveform, rasters and PSTHs are shown. In this example, when decoded with AV  
865 templates: 79% (22/28) of responses were classified as A1 when the visual stimulus was V1, and 32  
866 % of responses (9/28) were classified as A1 when the visual stimulus was V2, yielding a VPI score of  
867 47%. When decoded with A-only templates the values were 75% when V1 (22/28) and 35% when V2  
868 (10/28), yielding a VPI of 40%. For comparison the auditory-only condition (A12) is shown in **c, d**,  
869 population data showing the proportion of responses classified as A1 when the visual stimulus was  
870 V1 or V2 when decoded with auditory-only templates or auditory visual templates. **e,f**, resulting VPI  
871 scores. **h**, Mean ( $\pm$  SEM) values for these units when decoded with A-only templates, AV templates  
872 (as in Fig.2) or in the absence of a visual stimulus. The green data point in **d** depicts the example in **a**,  
873 **b**.

874 **Figure 4: Temporally coherent changes in visual luminance and auditory intensity enhance the**  
875 **representation of auditory timbre**

876 **a** Example unit response (from the awake dataset) showing the influence of visual temporal  
877 coherence on spiking responses to dual stream stimuli with (red PSTH) or without (black PSTH)  
878 timbre deviants. **b-d** timbre deviant discrimination in the awake dataset. Two deviants were  
879 included in each auditory stream giving a possible maximum of 4 per unit **b**, Histogram showing the  
880 number of deviants (out of 4) that could be discriminated from spiking responses **c**, Box plots  
881 showing the timbre deviant discrimination scores in the single stream condition across different  
882 visual conditions (Coh: coherent, ind: independent). The boxes show the upper and lower quartile  
883 values, and the horizontal lines indicates the median, the whiskers depict the most extreme data  
884 points not considered to be outliers (which are marked as individual symbols). **d**, Discrimination  
885 scores for timbre deviant detection in dual stream stimuli. Discrimination scores are plotted  
886 according to the auditory stream in which the deviant occurred and the visual stream that  
887 accompanied the sound mixture. V1 stimuli are plotted in red, and V2 stimuli in blue; therefore for d

888 and g the boxplots at the far left and right of the plot represent the cases in which the deviants  
889 occurred in an auditory stream which was temporally coherent with the visual stimulus while the  
890 central two boxplots represent the discrimination of deviants occurring in the auditory stream which  
891 was temporally independent of the visual stimulus. **e-g** show the same as **b-d** but for the  
892 anaesthetised dataset. See also supplemental figure 6.

893 **Figure 5: Auditory-visual temporal coherence enhances neural coding in auditory cortex.**

894 A pattern classifier was used to determine whether neuronal responses were informative about  
895 auditory or visual stimuli. The responses to single stream stimuli are shown for two example units,  
896 with responses grouped according to the identity of the auditory (**a, b**, for an auditory discriminating  
897 unit) or visual stream (**c, d**, for a visual discriminating unit). In each case the stimulus amplitude (a,b)  
898 / luminance (c,d) waveform is shown in the top panel with the resulting raster plots and PSTHs  
899 below. **e, f**: Decoder performance (mean  $\pm$  SEM) for discriminating stimulus identity (coherent: A1V1  
900 vs. A2V2, purple; independent: A1V2 vs. A2V1, blue) in auditory and visual classified units recorded  
901 in awake (e) and anaesthetised (f) ferrets. Pairwise comparisons for decoding of coherent versus  
902 independent stimuli (\*\*\*) indicates  $p < 0.001$ .

903

904 **Figure 6: Visual stimuli elicit reliable changes in the phase of the local field potential**

905 **a, b** Example LFP responses to each single stream auditory stimulus across visual conditions. Data  
906 obtained from the recording site at which multiunit spiking activity discriminated auditory stream  
907 identity in Fig. 5a and b. The amplitude waveforms of the stimuli are shown in the top panel, with  
908 the evoked LFP underneath (mean across 21 trials). The resulting inter-trial phase coherence (ITPC)  
909 values are shown in the bottom two panels (**c, d**). ITPC was calculated for coherent and independent  
910 AV stimuli separately and compared to a null distribution (ITPC across). **e, f** Single stream phase  
911 dissimilarity values (PDI) were calculated by comparing ITPC within values to the ITPC across null  
912 distribution. **g, i** Population mean inter-trial phase coherence (ITPC) values across frequency for  
913 coherent (**g**, significant frequencies 10.5-13, 16-20 Hz) and independent (**i** significant frequencies  
914 10.5-22 Hz) conditions. Dots indicate frequencies at which the ITPC-within values were significantly  
915 greater than the ITPC-across values (Pairwise t-test,  $\alpha = 0.0012$ , Bonferroni corrected for 43  
916 frequencies). **k**: Mean ( $\pm$ SEM) single stream phase dissimilarity index (PDI) values for coherent and  
917 independent stimuli in animals. Black dots indicate frequencies at which the temporally coherent  
918 single stream PDI is significantly greater than in the independent conditions ( $p < 0.001$ , significant  
919 frequencies 10.5-12.5). **h, j, l** as **g, i, k** for anaesthetised dataset. **m, n**, mean  $\pm$ SEM dual stream PDI  
920 values for awake (**m**, significant frequencies 10.5-12.5) and anaesthetised (n) datasets.

921

922 **Figure 7: Visual stimulus induced LFP phase changes in auditory cortex are mediated in visual**  
923 **cortex**

924 **a** Schematic showing the location of auditory cortical recording sites and the location of a cooling  
925 loop (black, grey line marks the 500  $\mu$ m radius over which cooling is effective (Wood et al., 2017)  
926 which was used to inactivate visual cortex. Individual recording sites contributing to c-n are shown  
927 with stars (simultaneous recordings are marked in the same colour). **b** Spike rate responses in  
928 Auditory Cortex (top row) and visual cortex (bottom row, sites  $>500 \mu$ m from the loop) in response  
929 to noise bursts or light flashes before, during and after cooling. **c, d**, inter-trial phase coherence  
930 values for the coherent (**c**) and independent stimuli (**d**) AV stimuli recorded in auditory cortex (AC)

931 prior to cooling visual cortex (VC) compared to the shuffled null distribution (ITPC across). Asterisks  
932 indicate the frequencies at which the ITPC values are significantly different from the shuffled ITPC-  
933 across distribution **e** single stream Phase Dissimilarity Index values calculated from the ITPC values in  
934 **c** and **d**. **f-h** - as c-e but while visual cortex was cooled to 9 degrees. **i-n** as c-h but for sites in visual  
935 cortex >500um from the cooling loop. c-h includes data from 83 sites from 6 electrode penetrations,  
936 i-n includes data from 47 sites from 5 penetrations.

937

938

## 939 **Methods**

### 940 **CONTACT FOR REAGENT AND RESOURCE SHARING**

941 Further information and requests for resources and reagents (data and MATLAB code) should be directed to and  
942 will be fulfilled by the Lead Contact, Jennifer Bizley ([j.bizley@ucl.ac.uk](mailto:j.bizley@ucl.ac.uk)).

943

### 944 **EXPERIMENTAL MODEL DETAILS**

945

946

947 The experiments were approved by the Animal Welfare and Ethical Review Board of University  
948 College London and The Royal Veterinary College, and performed under license from the UK Home  
949 Office (PPL 70/7267) and in accordance with the Animals Scientific Procedures Act 1986.

950 Neural responses were recorded in a total of 19 awake pigmented adult female ferrets (*Mustela*  
951 *putorius furo*; 1-5 years old). Fourteen of these animals contributed data to the awake dataset: Data  
952 from 9 of these animals was used for the main experiment (532 units), data from 11 other animals  
953 (6/9 in the main experiment, plus five other ferrets, totalling 128 units) was collected for additional  
954 control analysis (Fig. 3e, Fig. S5). Females (700-1500g, wild type) were co-housed in groups of 2-9.  
955 These animals were trained in a variety of psychoacoustic tasks unrelated to the current study prior  
956 to and after the implantation of recording electrodes. Animals were tested for this study on days  
957 when they were not participating in psychoacoustic testing. Five adult females were used to record  
958 responses under anaesthesia.

959 **METHOD DETAILS**

960

961 **Animal preparation**

962 Full methods for recording under anaesthesia can be found in Bizley et al., (2009). Briefly, ferrets  
963 were anesthetized with medetomidine (Domitor; 0.022mg/kg/h; Pfizer, Sandwich, UK) and ketamine  
964 (Ketaset; 5mg/kg/h; Fort Dodge Animal Health, Southampton, UK). The animal was intubated and  
965 the left radial vein was cannulated in order to provide a continuous infusion (5 ml/h) of a mixture of  
966 medetomidine and ketamine in lactated ringers solution augmented with 5% glucose, atropine  
967 sulfate (0.06 mg/kg/h; C-Vet Veterinary Products) and dexamethasone (0.5 mg/kg/h, Dexadreson;  
968 Intervet, UK). The ferret was placed in a stereotaxic frame in order to implant a bar on the skull,  
969 enabling the subsequent removal of the stereotaxic frame. The left temporal muscle was largely  
970 removed, and the suprasylvian and pseudosylvian sulci were exposed by a craniotomy, revealing  
971 auditory cortex (Kelly et al., 1986). The dura was removed over auditory cortex and the brain  
972 protected with 3% agar solution. The eyes were protected with zero-refractive power contact lenses.  
973 The animal was then transferred to a small table in a sound-attenuating chamber. Body  
974 temperature, end-tidal CO<sub>2</sub>, and the electrocardiogram were monitored throughout the experiment.  
975 Experiments typically lasted between 36 and 56 h. Neural activity was recorded with multisite silicon  
976 electrodes (Neuronexus Technologies) in a 1x 16, 2x 16 or 4x 8 (shank x number of sites)  
977 configuration. For experiments in which visual cortex was cooled we extended the craniotomy  
978 caudally to expose visual cortex and placed a cooling loop over the posterior suprasylvian gyrus.  
979 Details of the manufacture of the cooling loop and validation of its efficacy in the ferret animal  
980 model are provided in full in (Wood et al., 2017).

981

982 Full surgical methods for recording implanting electrode arrays to facilitate recording from awake  
983 animals are available in Bizley et al. (2013). Briefly, animals were bilaterally implanted with WARP-16  
984 drives (Neuralynx, Montana, USA) loaded with high impedance tungsten electrodes (FHC, Bowdoin,



985 USA) under general anaesthesia (medetomidine and ketamine induction, as above, isoflurane  
986 maintenance 1-3%). Craniotomies were made over left and right auditory cortex, a small number of  
987 screws were inserted into the skull for anchoring and grounding the arrays, and the WARP-16 drive  
988 was anchored with dental acrylic and protected with a capped well. Recording electrodes in awake  
989 animals targeted tonotopic auditory cortex (area MEG, containing fields A1 and AAF, and PEG,  
990 tonotopic belt areas PPF and PSF are located). Auditory fields were estimated prior to implantation  
991 based on known sulcal landmarks and confirmed with regular assessments of frequency tuning and  
992 post-mortem histology. Animals were allowed to recover for a week before the electrodes were  
993 advanced into auditory cortex. Pre-operative, peri-operative and post-operative analgesia were  
994 provided to animals under veterinary advice. Recordings were made over the next 1-2.5 years, with  
995 electrodes individually advanced every few weeks until the thickness of auditory cortex was  
996 traversed. Recordings were made while animals were passively listening/watching stimuli and  
997 holding their head at a water spout. During the recording a continuous stream of water was  
998 delivered from the spout.

999

## 1000 **Stimulus Presentation**

1001 All stimuli were created using TDT System 3 hardware (Tucker-Davis Technologies, Alachua, FL) and  
1002 controlled via MATLAB (Mathworks, USA). For recordings in awake animals, sounds were presented  
1003 over two loud speakers (Visaton FRS 8). Water deprived ferrets were placed in a dimly lit testing box  
1004 (69 x 42 x 52 cm length x width x height) and received water from a central reward spout located  
1005 between the two speakers. Sound levels were calibrated using a Brüel and Kjær (Norcross, GA)  
1006 sound level meter and free-field ½-inch microphone (4191). Auditory streams were presented at  
1007 65 dB SPL (Fig. 1a). Visual stimuli were delivered by illuminating the spout with a white LED which  
1008 provided full field illumination (Precision Gold N76CC Luxmeter, 0 to 36.9 lux). The animals were not

1009 required to do anything other than maintain their heads in position at the spout where they were  
1010 freely rewarded. Recording was terminated when animals were sated.

1011 For anaesthetised recordings, acoustic stimuli were presented using Panasonic headphones  
1012 (Panasonic RP-HV297, Bracknell, UK) at 65 dB SPL. Visual stimuli were presented with a white Light  
1013 Emitting Diode (LED) which was placed in a diffuser at a distance of roughly 10 cm from the  
1014 contralateral eye so that it illuminated virtually the whole contralateral visual field.

1015

### 1016 **Stimuli and data acquisition:**

1017 *Auditory stimuli* were artificial vowel sounds that were created in Matlab (MathWorks, USA). In the  
1018 behavioural experiment that motivated this study (Maddox et al., 2015), stimuli were 14 seconds in  
1019 duration. However, we adapted the stimulus duration in awake recordings to 3 seconds in order to  
1020 collect sufficient repetitions of all stimuli, and to ensure animals maintained their head position  
1021 facing forwards for the whole trial duration. The animals were observed constantly via a webcam  
1022 and recording was terminated / paused if the animal's head moved from the centre spout. In the  
1023 anaesthetised recording stimulus streams were 14 seconds long, as in the human psychophysics but  
1024 we only analysed the first 3 seconds to ensure datasets were directly comparable (see also Fig. S7 e-  
1025 h which replicates analysis for 3 second and 14 second stimuli).

1026 Stimulus A1 was the vowel [u] (formant frequencies F1-4: 460, 1105, 2857, 4205 Hz, F0= 195Hz), A2  
1027 was [a] (F1-4: 936, 1551, 2975, 4263 Hz, F0= 175Hz). Streams were amplitude modulated with a  
1028 noisy lowpass (7 Hz cut-off) envelope. Unless specifically noted, the timbre of the auditory stream  
1029 remained fixed throughout the trial. However, we also recorded responses to auditory streams that  
1030 included brief timbre deviants. As in our previous behavioural study, deviants were 200ms epochs in  
1031 which the identity of the vowel was varied by smoothly changing the first and second formant

1032 frequencies to and from those identifying another vowel. Stream A1 was morphed to/from [ε] (730,  
1033 2058, 2857, 4205 Hz) and A2 to/from [i] (437, 2761, 2975, 4263 Hz).

1034 *Visual stimuli* were generated using an LED whose luminance was modulated with dynamics that  
1035 matched the amplitude modulation applied to A1 or A2. In single stream conditions a single auditory  
1036 and single visual stream were presented (e.g. A1V1, A1V2, A2V1, or A2V2) whereas in dual stream  
1037 conditions both auditory streams were presented simultaneously, accompanied by a single visual  
1038 stimulus (A12V1, A12V2, A12V1 A12V2) (Fig. 1e). Auditory streams were always presented from both  
1039 speakers so that spatial cues could not facilitate segregation, and stimulus order was varied pseudo-  
1040 randomly. In the anaesthetised recordings each stimulus was presented 20 times. In the awake  
1041 dataset, where recording duration was determined by how long the ferret remained at the central  
1042 location (mean repetitions: 20, minimum: 14, maximum: 34).

1043 During anaesthetised recordings, pure tone stimuli (150 Hz to 19 kHz in 1/3-octave steps, from 10 to  
1044 80 dB SPL in 10 dB, 100 ms in duration, 5 ms cosine ramped) were also presented. These allowed us  
1045 to characterize individual units and determine tonotopic gradients, so as to confirm the cortical field  
1046 in which any given recording was made. Additionally broadband noise bursts and diffuse light flashes  
1047 (100 ms duration, 70 dB SPL) were presented and used to classify a stimulus as auditory, visual or  
1048 auditory visual. LFPs were subjected to current source density analysis to identify sources and sinks  
1049 as described by Kaur et al. (2005).

#### 1050 **Cortical cooling**

1051 During these experiments we made joint recordings in visual cortex (usually >500 μm from the  
1052 cooling loop, in order to determine whether visual cortical processing was impaired generally) and  
1053 auditory cortex simultaneously. We recorded responses to the single stream stimuli before and  
1054 during cooling, and, at each site additionally recorded responses to noise bursts and light flashes  
1055 before, during and after cooling. We used the responses to simple stimuli such as these to show that  
1056 we could recover the original spiking responses (and data was excluded from any recording sites in  
1057 which did not return to within 20% (a common criterion used in cooling studies: e.g. Antunes and  
1058 Malmeirca, 2011) of their pre-cooling spike rates). We did not record responses to the longer stimuli

1059 used in this study in the post-cooling condition as the additional recording time for these stimuli  
1060 would have compromised our ability to record across several different sites in each animal.

1061

## 1062 QUANTIFICATION AND STATISTICAL ANALYSIS

1063

1064 Electrophysiological data were analysed offline. Spiking activity and local field potential signals were  
1065 extracted from the broadband voltage waveform by filtering at 0.3-5kHz and 1-150 Hz respectively.  
1066 Spikes were detected, extracted and then sorted with a spike-sorting algorithm (WaveClus, Quiroga  
1067 et al., 2004).

### 1068 **Spiking responses**

1069 We used a Euclidean distance based pattern classifier (Foffani and Moxon, 2004) with leave-one-out  
1070 cross validation to determine whether the neuronal responses to different stimuli could be  
1071 discriminated. Spiking responses to a given stimulus were binned into a series of spike counts from  
1072 stimulus onset (0 s) to offset (3s) in 20 ms bins. The average across-repetition response to each  
1073 stimulus (excluding the to-be-classified response) were used as templates and the response to a  
1074 single stimulus presentation was classified by calculating the Euclidean distance between itself and  
1075 the template sweeps and assigning it to the closest template. To determine whether the classifier  
1076 performed significantly better than expected by chance, a 1000 iteration permutation test was  
1077 performed where trials were drawn (with replacement) from the observed data and randomly  
1078 assigned to a stimulus that was then used for template formation / decoding. A neural response was  
1079 considered to be significantly informative about stimulus identity if the observed value exceeded the  
1080 95th percentile of the randomly-drawn distribution.

1081 This approach allowed us to classify units according to their functional properties: auditory units  
1082 discriminated two auditory stimuli based on the amplitude modulation of sound (A1 versus A2)  
1083 regardless of visual dynamics, (Fig. 5a, b), visual units discriminated visual presentations based on  
1084 temporal envelope of visual stimuli (V1 versus V2) regardless of auditory presentation (Fig. 5c,d) and  
1085 AV units could do both. This approach was extended to classify dual stream responses by using the

1086 average response to each of the temporally coherent single stream stimuli (A1V1 or A2V2) as  
1087 templates (Fig. 2,3, S2-4). Performance was (arbitrarily) expressed as the proportion of responses  
1088 classified as being from the A1, and compared for the two dual stream stimuli with different visual  
1089 conditions (Figure 5). All units in which either auditory or visual stimulus identity could be decoded  
1090 were included in the dual-stream analysis. A VPI was derived from this measure as the difference  
1091 between the percentage of A1V2 trials labelled A1 and the percentage of A1V1 trials labelled A1  
1092 (multiplied by -1 to make the index positive for units that were strongly influenced by the visual  
1093 stimulus). Therefore units which were fully influenced by the identity of the visual stimulus would  
1094 have a visual preference score of 100, while those in which the visual stimulus did not influence the  
1095 response at all would have a score around 0 (Fig. 2e,h). We then assessed the significance of  
1096 observed VPI scores using a permutation test ( $p < 0.05$ ) in which the identity of single stream trials  
1097 used to generate classifier templates was shuffled and the VPI recalculated for 1000 iterations.

1098

#### 1099 **Timbre deviant analysis:**

1100 In order to determine how a visual stimulus influenced the ability to decode timbre deviants  
1101 embedded within the auditory streams we used the cross-validated pattern classifier described  
1102 above for analysing single stream stimuli to discriminate deviant from no-deviant trials. Responses  
1103 were considered over the 200 ms time window that the deviant occurred (or the equivalent time  
1104 point in the no-deviant stimulus) binned with a 10 ms resolution. Significance was assessed by a  
1105 1000 iteration permutation test in which trials were randomly drawn with replacement from deviant  
1106 and no-deviant responses. The discrimination score was calculated as the proportion of correctly  
1107 classified trials.

#### 1108 **Classification as auditory or visual with simple stimuli**

1109 During recordings made under anaesthesia, we also recorded responses to noise bursts and light  
1110 flashes (both 100 ms duration) presented separately and together to compare how the proportion of  
1111 auditory / visual discriminating units measured to naturalistic dynamic stimuli compared to more  
1112 traditional artificial stimuli. Specifically, responsiveness was defined using a two-way ANOVA  
1113 (factors: auditory stimulus [on/off] and visual stimulus [on/off]) on spike counts measured during  
1114 stimulus presentation. We defined units as being sound-driven (main effect of auditory stimulus, no  
1115 effect of visual stimulus or interaction), light-driven (main effect of visual stimulus, no effect of  
1116 auditory stimulus or interaction) or auditory-visual (main effect of both auditory and visual stimuli or  
1117 significant interaction;  $p < 0.05$ ) as in (Bizley et al., 2007).

#### 1118 **Phase/power dissimilarity analysis:**

1119 Local field potential recordings were considered for all sites at which there was a significant driven  
1120 spiking response, irrespective of whether that response could discriminate auditory or visual stream  
1121 identity. For the single stream trials, we computed a single stream Phase Dissimilarity Index (PDI),  
1122 which characterizes the consistency and uniqueness of the temporal phase/power pattern of neural  
1123 responses to continuous auditory stimuli (Luo and Poeppel, 2007a). This analysis compares the  
1124 phase (or power) consistency across repetitions of the same stimulus with a baseline of phase-  
1125 consistency across trials in which different stimuli were presented.

1126 In the first stage of PDI analysis, we obtained a time-frequency representation of each response  
1127 using wavelet decomposition with complex 7-cycle Morlet wavelets in 0.5 steps between 2.5–45 Hz,  
1128 resulting in 86 frequency points. Next, we calculated the inter-trial phase-coherence value (ITPC;  
1129 Equ.1) at each time-frequency point, across all trials in which the same stimulus was presented. For  
1130 each frequency band, the ITPC time-course was averaged over the duration of the analysis window  
1131 and across all repetitions to obtain the average *within-stimulus ITPC*.

1132

1133 
$$ITPC_{t,f} = \left| \frac{\sum_{k=1}^N e^{i\theta_{k,t,f}}}{N} \right|$$
 Equ.1

1134 In which N is equal to the number of trials, and  $\theta$  is the phase of trial  $k$  at a given frequency ( $f$ ) and  
1135 time ( $t$ ). The *across-stimuli ITPC* was estimated using the same approach but using shuffled data,  
1136 such that the ITPC was computed across trials with the same auditory stimulus but randomly drawn  
1137 visual stimuli. The single stream phase dissimilarity index (Single stream PDI) was computed as the  
1138 difference between the ITPC value calculated for *within* trials and the ITPC values calculated *across*  
1139 visual trials (Equ.2).  $\gamma$ . The dissimilarity function for each frequency bin  $i$  was defined as;

1140

1141 
$$Single\ Stream\ PDI_i = \frac{\sum_{j=1}^N ITPC_{ij, within\ vis}}{N} - \frac{\sum_{j=1}^N ITPC_{ij, across\ vis}}{N}$$
 Equ.2

1142

1143 Large positive PDI indicate that responses to individual stimuli have a highly consistent response on  
1144 single trials. Single stream PDI values were calculated for each stimulus type and then averaged  
1145 across stimuli to calculate values for temporally coherent and temporally independent auditory  
1146 visual stimuli. Single stream PDI was positive if within stimulus ITPC was larger than across-stimulus  
1147 ITPC (pairwise t-test,  $p < 0.05$  Bonferroni correction for 86 frequencies points) and was considered  
1148 significant if a minimum of 2 adjacent bins exceeded the corrected threshold. PDI magnitude values  
1149 were calculated by summing the PDI values across all significant frequencies.

1150 Dual stream phase dissimilarity index (dual stream PDI) values were calculated by extending this  
1151 approach for dual stream stimuli with the goal of determining how the temporal envelope of the  
1152 visual stimulus influences the neural response to a sound mixture. To this end, we calculated the  
1153 *within-dual ITPC* from the A12V1 trials and A12V2 trials separately and *across-dual ITPC* by randomly  
1154 selecting trials from both stimuli (Equ.3). The within-dual and across-dual ITPCs were then averaged  
1155 over time and subtracted to yield the dual stream PDI (Equ.3).

1156

1157

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

1158

1159 Positive dual stream PDI values indicate that the time course of the neural responses was influenced  
1160 by visual input, despite the identical acoustic input. We determined whether the dual stream PDI  
1161 was greater if the *within\_dual* ITPC was significantly larger than *across\_dual* ITPC (pairwise t-test,  
1162  $p < 0.05$  Bonferroni correction, as above). PDI magnitude values were calculated by summing the PDI  
1163 values across all significant frequencies.

#### 1164 Analysis of responses during cooling

1165 Our physiological recordings confirm that within the vicinity of the loop the inactivation spans all  
1166 cortical layers. As the temperature change drops off with distance, at distances further from the  
1167 loop the cooling is more restricted to superficial layers. These data are presented in full in Wood et  
1168 al., 2017.

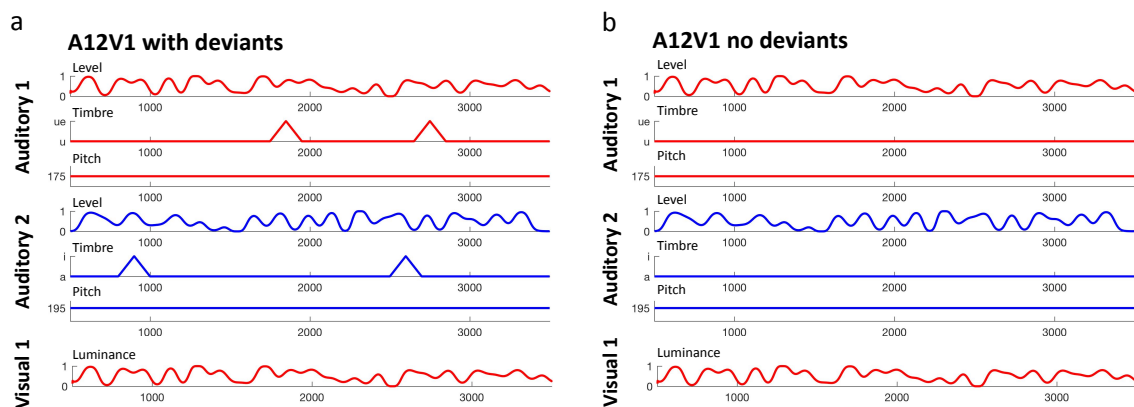
1169

1170

## 1171 Supplemental Results

1172

### 1173 Supplemental Figure 1 (related to Figure 1).



1174

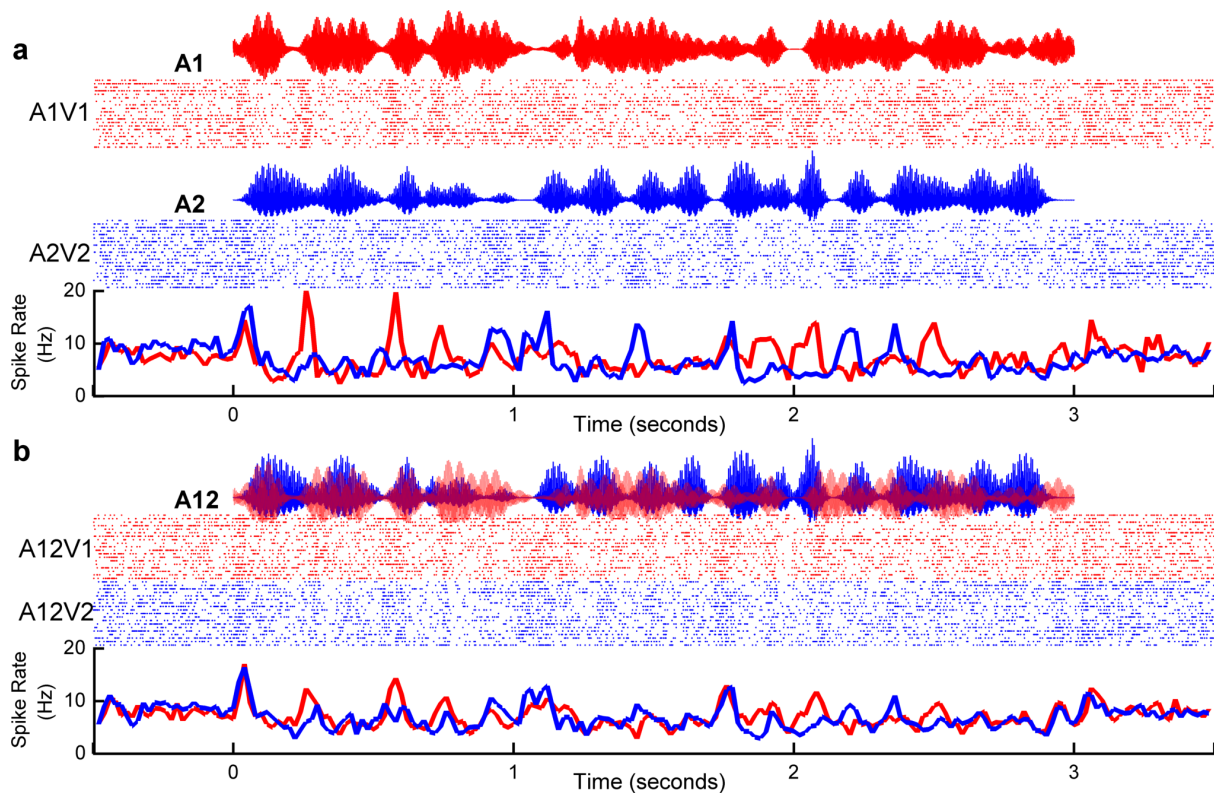
1175



1176 Schematic of the stimuli used in this study to illustrate the difference between stimuli with deviants  
1177 (a) and without (b). Each row depicts the time course of each feature within the stimulus over a  
1178 single trial. Importantly, the timing of the timbre deviants is not predicted by the temporally  
1179 coherent changes in the binding features: here sound level and visual luminance.

1180

1181 **Supplemental Figure 2 (related to Figure 2)**



1182

1183 **Visual stimuli can determine which sound stream auditory cortical neurons follow in a mixture:**  
1184 **example auditory-discriminating unit**

1185 The spiking responses of an example unit are shown to coherent single stream (a) and dual stream  
1186 stimuli (b). This example unit was an auditory discriminating unit recorded in an awake animal. In  
1187 this example 68% (15/22) of responses were classified as A1 when the visual stimulus was V1, and 40  
1188 % of responses (9/22) were classified as A1 when the visual stimulus was V2, yielding a VPI score of  
1189 28%.

1190

1191

1192

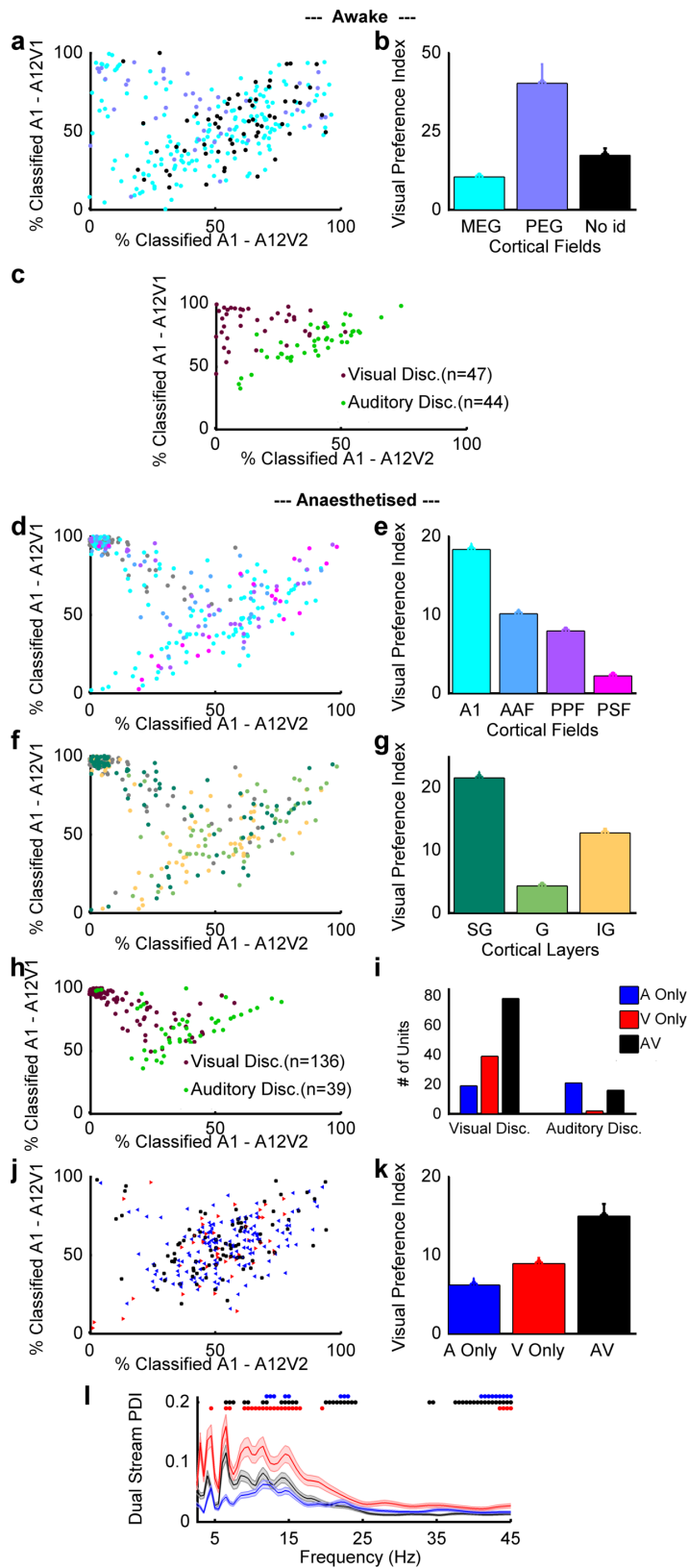
1193

1194

1195

1196 **Supplemental Figure 3 (related to Figures 2 and 3)**

1197 **Effects of cortical field, cortical lamina and response type on visual modulation of dual stream**  
 1198 **responses**



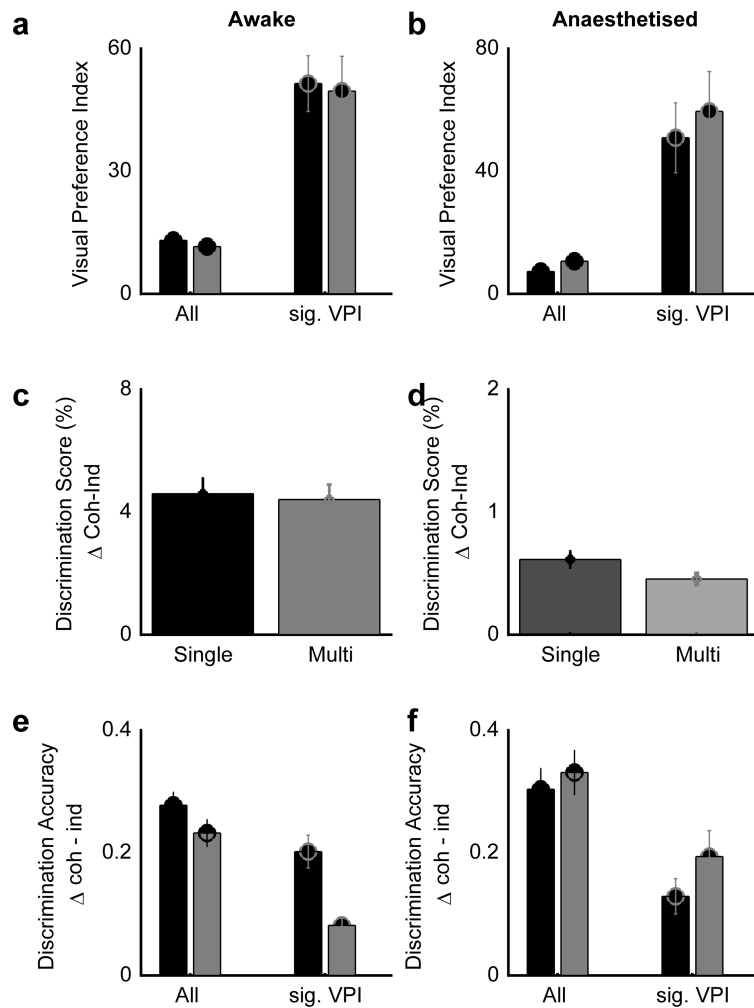
1199

1200

1201 **a** The distribution of decoding values for dual stream responses, according to the visual condition (as  
1202 in Figure 2C) colour coded according to recording location. Awake recordings were made from  
1203 animals in which recording electrodes targeted the MEG (where A1 and AAF are located) and PEG  
1204 (where fields PSF and PPF are located). The sampling density of our recording arrays (16 electrodes  
1205 in a 4 x 4 grid with electrodes separated by 800  $\mu\text{m}$ ) does not provide the high spatial resolution  
1206 necessary to determine recording locations precisely – particularly at the low frequency reversal that  
1207 separates PEG from MEG in the ferret. Therefore some recording electrodes are classified as  
1208 ‘unidentified’ as in these cases neither the frequency tuning, nor the post-mortem histology, allowed  
1209 us to unambiguously ascribe the recording site to MEG or PEG. **b** VPI scores from **a** by cortical area  
1210 (mean  $\pm$ SEM). A one way ANOVA revealed there to be a significant effect of field ( $F(2,278) = 7.1354$ ,  
1211  $p = 9.5072e-04$ ), with post-hoc comparisons indicating PEG was significantly higher than MEG or  
1212 unidentified sites. **c** re-plots the data in **a** showing only units with a significant VPI, colour-coded  
1213 according to whether they were visual classified or auditory classified. **d,e**, as **a,b** but for the  
1214 anaesthetised dataset where we were able to make sufficient penetrations (30-50) to generate a high  
1215 resolution tonotopic map and hence ascribe recording sites to cortical subfields. A one way ANOVA  
1216 revealed there to be no significant effect of field on VPI ( $F(3,1130) = 2.1886$ ,  $p = 0.0877$ ).

1217 Recordings in anaesthetised animals were made with linear shank electrodes, facilitating current source  
1218 density analysis to identify the cortical layers. **f** shows the distribution of decoding values in the dual  
1219 stream condition according to recording location and depth in the anaesthetised dataset. **g**  
1220 summarises the data in **c** by cortical field. A one-way ANOVA across cortical layers showed a  
1221 significant effect of layer ( $F(2,1134) = 3.1543$ ,  $p = 0.0430$ ) with post-hoc comparisons indicating that  
1222 the VPI scores were greater in the supra-granular than granular layers. **h** plots the distribution of dual  
1223 stream decoding values for only units with a significant VPI, colour-coded according to whether they  
1224 were classified as auditory-discriminating or visual-discriminating. In the anaesthetised animal we  
1225 additionally used simple noise bursts and light flashes to describe units as auditory (A;  $n=160$  units),  
1226 visual (V,  $n=53$ ) or auditory visual (AV; grey,  $n=94$ ). **i** shows the distribution of A, V and AV units  
1227 that were also classified as auditory-discriminating or visual-discriminating. Of 136 visual  
1228 discriminating units with a significant VPI, 19 were categorised as auditory, 39 as visual and 78 as  
1229 auditory-visual, of 39 auditory-discriminating units with significant VPI values 21 were auditory, 2  
1230 were visual and 16 were auditory visual Fig. S3i. **j,l**, as **d,e**, but with units colour coded according to  
1231 whether they were classified as A, V or AV with simple stimuli. **k** mean ( $\pm$  SEM) dual stream phase  
1232 dissimilarity index (PDI) values for recording sites categorised according to the spiking responses  
1233 recorded there. Symbols indicate the frequencies at which the dual stream PDI index was significant  
1234 (pairwise t-test,  $p < 0.001$  with correction). While the phase effects are greatest at the sites where visual  
1235 activity was recorded, significant dual stream PDI values were observed in all three unit types. In all  
1236 three cases significant phase coherence was seen at 12Hz, 13.5Hz-14.5Hz and 42.5-44.5Hz.  
1237 Modulation at 10-12 Hz was only observed at sites in which AV and V responses were recorded.

1238



1239

1240 **Supplementary figure 4: Related to Figures 2,3 and 5**

1241 **The effects of temporal coherence on single stream decoding and of visual identity on dual**  
 1242 **stream decoding are evident in both single and multi-units.**

1243

1244 **a,b** No effect of unit type (single versus multi-unit) was found for VPI values in awake recordings (all  
 1245 units:  $F(1,270) = 0.1595, p = 0.6899$ ; significant VPI:  $F(1,90) = 0.1048, p = 0.7469$ ) and anaesthetised  
 1246 recording (all units:  $F(1,332) = 0.4740, p = 0.4921$ ; significant VPI:  $F(1,174) = 0.6867, p = 0.4123$ )  
 1247 **c,d** Discrimination scores for timbre deviant detection in dual stream stimuli is indistinguishable for  
 1248 single units and multi units in awake recording  $F(1,167) = 0.0326, p = 0.8570$ ) and in anaesthetised  
 1249 recording ( $F(1,221) = 0.8339, p = 0.3625$ )

1250 **e,f** Single units had significantly higher influence of temporal coherence on discrimination accuracy  
 1251 for single stream stimuli (as in Figure 5e,f) in awake recordings (All units:  $F(1,270) = 4.9916, p =$   
 1252  $0.0263$ ; for units with significant VPI:  $F(1,90) = 10.1780, p = 0.0020$ ). Single and multiunits had  
 1253 equivalent performance in the anaesthetised dataset ( $F(1,332) = 1.2558, p = 0.2641$ ; significant VPI:  
 1254  $F(1,174) = 1.5121, p = 0.2262$ ).

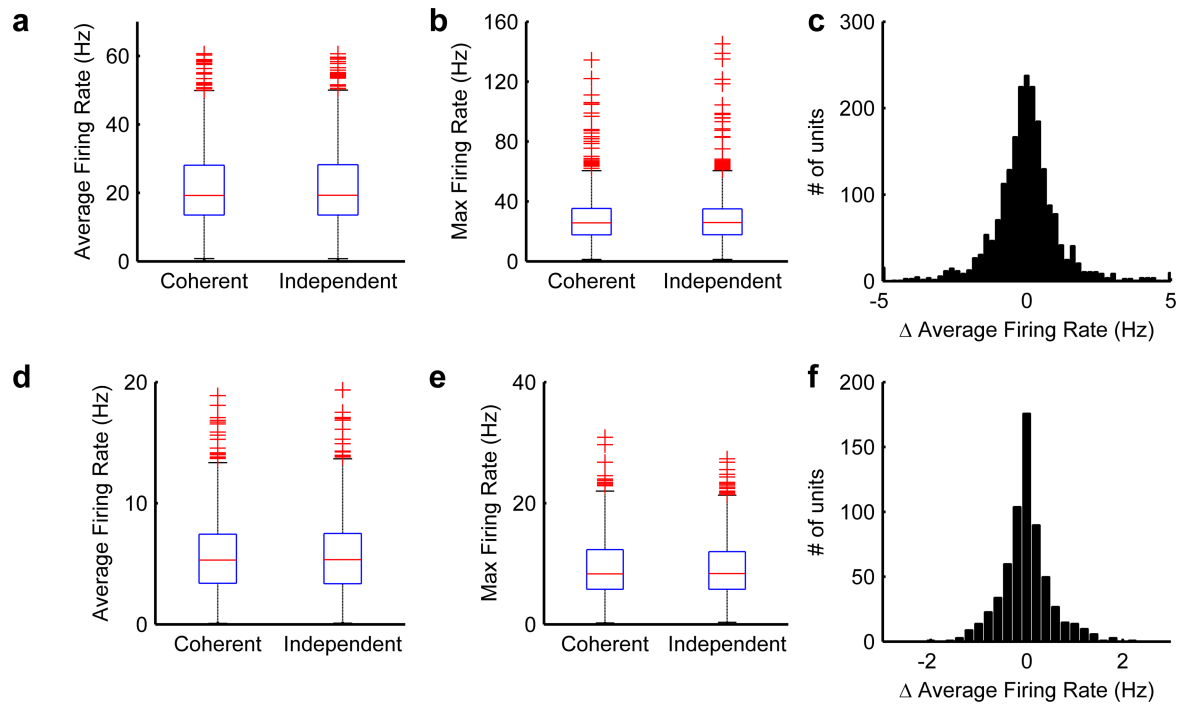
1255

1256

1257

1258

1259 **Supplemental Figure 5**



1260

1261 There were no statistically significant changes in mean (**a,d**) or max (**b,e**) firing rate between  
1262 temporally coherent and temporally independent datasets in either the awake (**a,b,c**) or anesthetised  
1263 (**d,e,f**) datasets. **Awake dataset:** For all units: Mean firing rate  $t_{540} = -0.0308$ ,  $p = 0.9754$ ; Max firing  
1264 rate:  $t_{540} = 0.4354$ ,  $p = 0.6636$ . For units with a significant VPI value: Mean:  $t_{180} = -0.0631$ ,  $p = 0.7694$ ;  
1265 Max:  $t_{180} = 0.8563$ ,  $p = 0.0939$ . **Anesthetised dataset:** all units: Mean firing rate  $t_{664} = -0.0638$ ,  $p =$   
1266  $0.9492$ . Max firing rate:  $t_{664} = 0.0047$ ,  $p = 0.9947$ . Significant VPI units Mean firing rate  $t_{348} = 0.0308$ ,  $p =$   
1267  $0.9498$ . Max firing rate:  $t_{348} = 0.0235$ ,  $p = 0.9912$ .

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

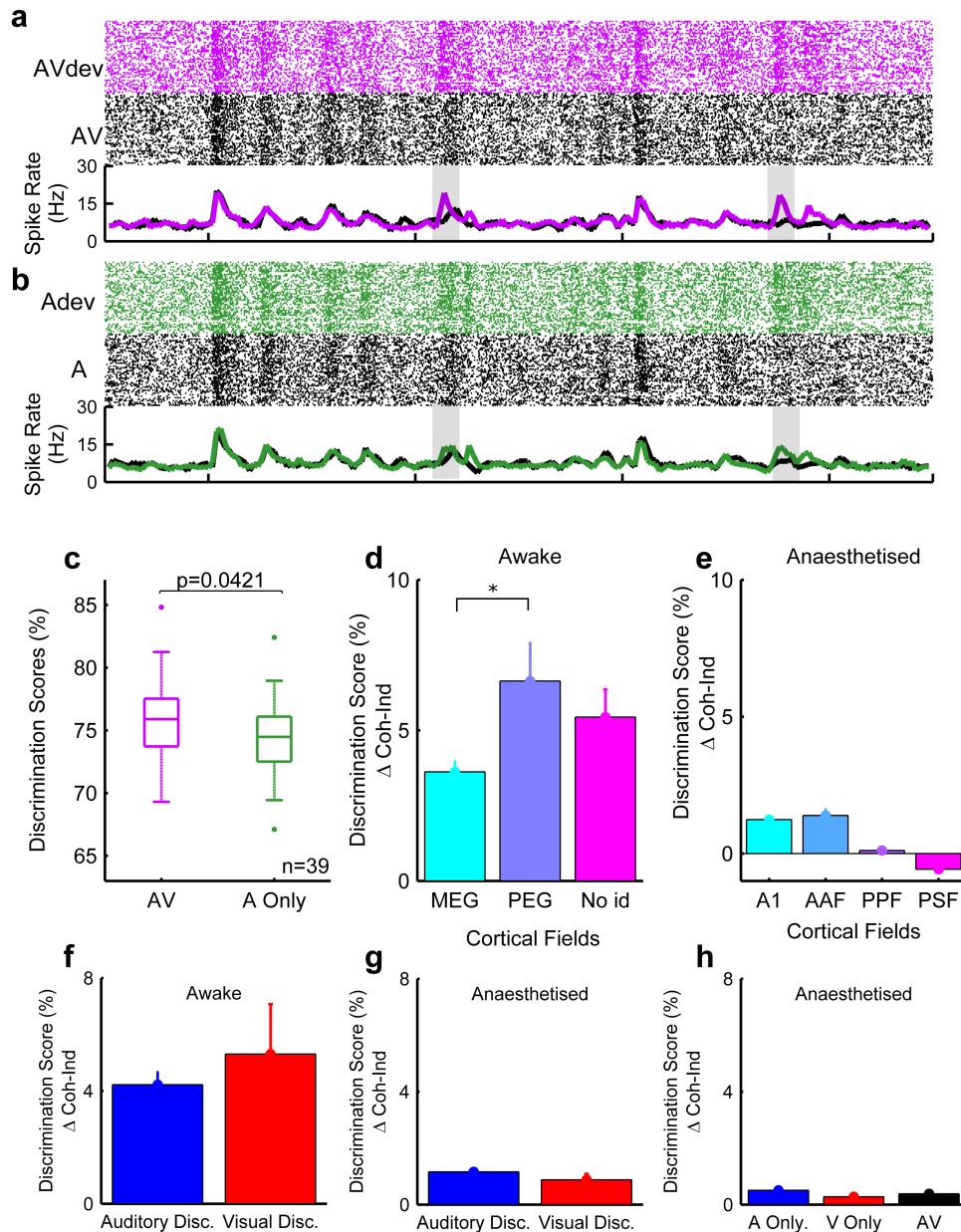
1281

1282

1283

1284 **Supplemental Figure 6 (related to Figure 4).**

1285



1286

1287 **a,b**, Rasters and PSTH for single stream stimuli containing deviants (top row, black) or without  
 1288 deviants (bottom row, purple/green). **a**, auditory stream presented with a temporally coherent visual  
 1289 stimulus, **b**, auditory stream presented in isolation. Grey panels indicate the timing of the timbre  
 1290 deviants. **c**, discrimination scores for detecting trials with deviants in them were significantly higher  
 1291 in AV trials than A only trials. Recordings were made in awake animals (# of animals =3, n=39 driven  
 1292 units.). Pairwise comparison  $t_{76} = 2.0676$ ,  $p = 0.0421$ )

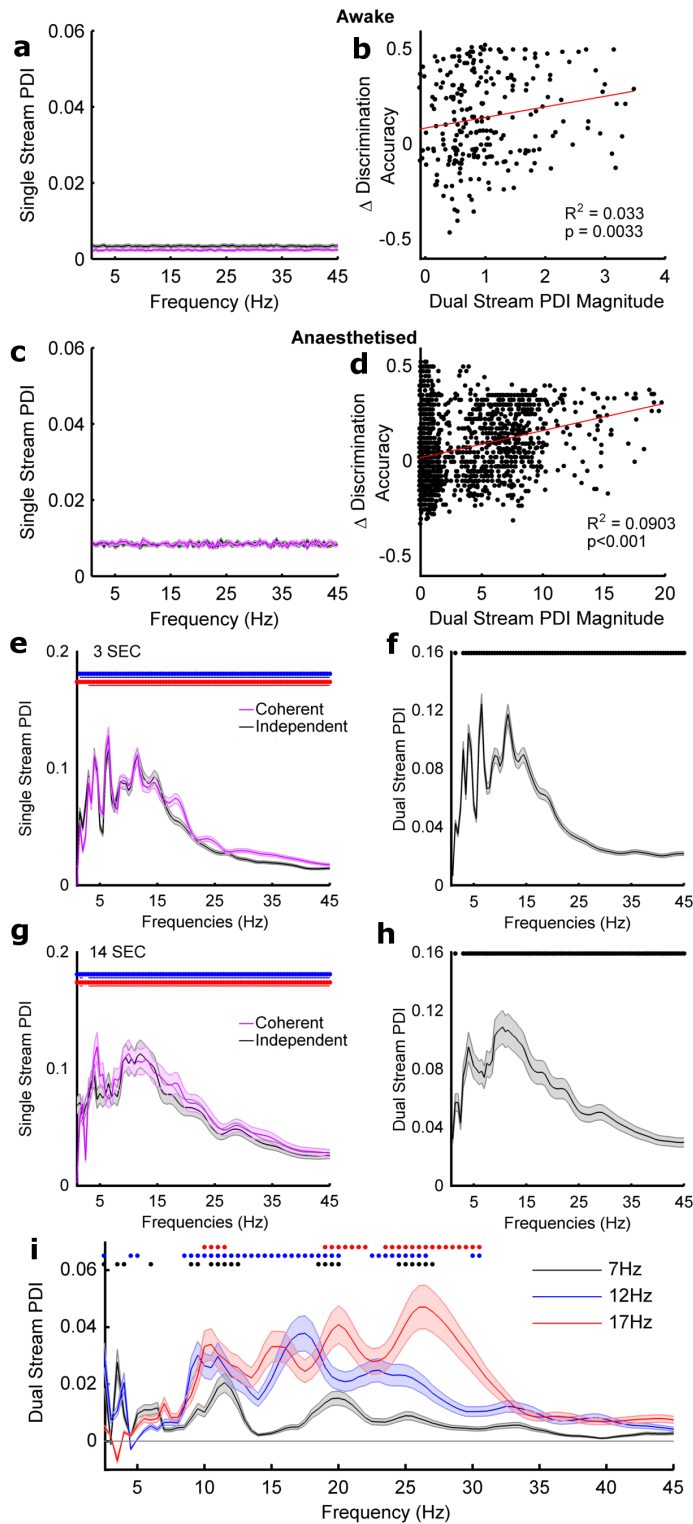
1293 **d,e**, comparison of how visual temporal coherence influenced deviant encoding across cortical fields  
 1294 in awake data and anaesthetised data. In awake data, a one-way ANOVA across cortical fields showed  
 1295 a significant effect of field ( $F(2,165) = 2.6710$ ,  $p = 0.0322$ ) with post-hoc comparisons indicating that  
 1296 the discrimination scores were greater in the PEG than MEG. However, there were no significant  
 1297 effect of field in anaesthetised data ( $F(3,245) = 2.0627$ ,  $p = 0.1057$ )

1298 f.g No differences were found in discrimination score between auditory discriminating and visual  
1299 discriminating units in awake recordings ( $F(1,100) = 0.2547$ ,  $p = 0.6149$ ) or in anaesthetised  
1300 recordings ( $F(1,119) = 0.0689$ ,  $p = 0.7933$ ).

1301 h There were also no differences across different unit type in anaesthetised recording ( $F(2,423) =$   
1302  $0.846$ ,  $p = 0.9169$ )

1303

1304 **Supplemental Figure 7 (related to Figure 6)**



1305

1306 **a,c** dual stream power discriminability index values for awake (a) and anaesthetised datasets (c). In  
 1307 neither case was there any frequency whose power was significantly influenced by the visual stimulus  
 1308 identity. **b,d**, relationship between phase discriminability index (PDI) values measured in the dual  
 1309 stream condition and the VPI score. There is a weak correlation between the magnitude of the dual  
 1310 stream PDI values and VPI values.

1311 During our initial analysis we observed that PDI values were higher in anaesthetised animals than  
 1312 awake animals. In order to determine whether this was a difference due to behavioural state or simply



1313 an artefact of stimulus length for all of the analysis reported in this paper we restricted analysis of the  
1314 anaesthetised responses to the first three seconds of the stimulus. In **e-h** we explicitly compare the PDI  
1315 values obtained for 3 second (**e,f**) and 14 second (**g,h**) single stream (**e,g**) and dual stream (**f,h**)  
1316 stimuli. to match that recorded in the awake dataset. While phase coherence values were slightly  
1317 higher for longer duration stimuli and hence at longer stimulus durations the ITPC profile and  
1318 resulting PDI varied more smoothly with frequency. However at both durations phase values were  
1319 significantly different from zero at all frequencies. The pattern of significant phase selectivity values  
1320 was also preserved across stimulus durations. (**b, d**). Frequency points at which the single stream PDI  
1321 value and dual stream PDI values were similar in 3 second length (**a, b**) and 14 second length (**c, d**)  
1322 Blue, red and black symbols indicate where the PDI was significant (pairwise t-test,  $\alpha = 0.0012$  with  
1323 bonferoni correction).

1324 **i** Dual stream stimuli were generated with three different amplitude modulation rates (<7Hz, as in the  
1325 main experiment, <12Hz and <17Hz, values picked to avoid harmonics of 7 Hz) and responses to  
1326 these were recorded in 92 units. Symbols indicate where the dual stream phase selectivity index was  
1327 significant (pairwise t-test,  $p < 0.05$  with correction). In all three cases significant phase coherence is  
1328 seen between 10Hz-11.5Hz, 19Hz-20Hz and 24-26 Hz.

1329

1330

1331

1332

1333