1    **Title:**

2    Multiplexed and multivariate representations of sound identity during perceptual constancy

3    **Authors:**

4    Stephen M. Town, Katherine C. Wood & Jennifer K. Bizley

5    **Affiliations:**

6    Ear Institute, University College London, 332 Gray's Inn Road, London, WC1X 8EE, UK

7    **Corresponding authors:**

8    Stephen Town (s.town@ucl.ac.uk)

9    Jennifer Bizley (j.bizley@ucl.ac.uk)

## Summary

Perceptual constancy requires neural representations that are selective for object identity, but also tolerant for identity-preserving transformations. How such representations arise in the brain and contribute to perception remains unclear. Here we studied tolerant representations of sound identity in the auditory system by recording multi-unit activity in tonotopic auditory cortex of ferrets discriminating the identity of vowels which co-varied across orthogonal stimulus dimensions (fundamental frequency, sound level, location and voicing). We found that neural decoding of vowel identity was most successful across the same orthogonal dimensions over which animals generalized their behavior. We also decoded orthogonal sound features and behavioral variables including choice and accuracy to show a behaviorally-relevant, multivariate and multiplexed representation of sound, with each variable represented over a distinct time-course. Finally, information content and timing of sound feature encoding was modulated by task-engagement and training, suggesting that tolerant representations during perceptual constancy are attentionally and experience-dependent.

## Keywords

Perceptual constancy, hearing, auditory cortex, ferret, behavior, electrophysiology, attention, learning, vowel

## Introduction

28

29    Perceptual constancy, also known as perceptual invariance, is the ability to recognize objects

30    across variations in sensory input such as a face from multiple angles or a word spoken by different

31    talkers (Bizley and Cohen, 2013; Logothetis and Sheinberg, 1996). Perceptual constancy requires that

32    sensory systems such as vision and hearing develop a level of tolerance to identity preserving

33    transformations (DiCarlo and Cox, 2007; DiCarlo et al., 2012). In hearing, the development of

34    tolerance is critical to the representation of sounds such as individual words or phonemes across

35    talkers, voice pitch, background noise and other acoustic transformations (Sharpee et al., 2011) and

36    is a key step in auditory object formation and scene analysis (Bizley and Cohen, 2013; Bregman,

37    1990; Griffiths et al., 2004).

38    Both humans and other animals perceive sound features constantly despite variation in

39    sensory input: we can recognize loudness across variation in location (Zahorik and Wightman, 2001),

40    frequency across sound level (Polley et al., 2006) and sound identity across talkers (Kojima and

41    Kiritani, 1989; Ohms et al., 2010), vocal tract length (Ghazanfar et al., 2007; Schebesch et al., 2010;

42    Smith et al., 2005) and fundamental frequency (F0)(Bizley et al., 2013a; Honorof and Whalen, 2010;

43    Town et al., 2015). At the neural level, tolerance emerges within auditory cortex for sounds including

44    vocalizations (Billimoria et al., 2008; Carruthers et al., 2015; Meliza and Margoliash, 2012), pure

45    tones (Sadagopan and Wang, 2008) and pulse trains (Bendor and Wang, 2007). Auditory cortical

46    neurons are modulated by multiple features of speech sounds, such as synthesized vowels (Bizley et

47    al., 2009), and when variables are considered in discrete time windows, tolerant responses of vowel

48    identity, as well as information about sound location and F0 can be recovered (Walker et al., 2011).

49    However, tolerance has yet to be shown in subjects actively demonstrating perceptual constancy,

50    and the behavioral relevance of tolerant representations in auditory cortex remains unclear.

51    Furthermore, although auditory cortical processing is modulated by attention and experience

52    (Osmanski and Wang, 2015), it is unknown how these processes affect tolerant representations.

53       Here we asked if tolerant representations exist in auditory cortex during perceptual

54      constancy, how tolerance was related to behavior, and modulated by attention and experience. To

55      address these questions, we recorded the activity of auditory cortical neurons in ferrets

56      discriminating synthesized vowel sounds across identity-preserving, orthogonal acoustic

57      transformations - including variations in F0, sound location, level and voicing.

58       We hypothesised that auditory cortical neurons would show tolerance across the same

59      range of orthogonal variables over which animals demonstrate perceptual constancy, and that such

60      tolerance would be degraded in cases where animals failed to generalize vowel identity. As auditory

61      cortex represents multiple stimulus variables, we expected tolerance to be accompanied by

62      information about both task-relevant and irrelevant sound features. Finally, we predicted that the

63      neural correlates of perceptual constancy should be dependent on an animal's behavioral

64      performance, attentional state and training. Our findings confirmed that neurons could represent

65      vowel identity across orthogonal variations and thus provide tolerant representations in perceptual

66      constancy. Furthermore, we also demonstrated these representations were sensitive to behavioral

67      performance, failures to perceive vowel constancy, attentional state and long-term experience.

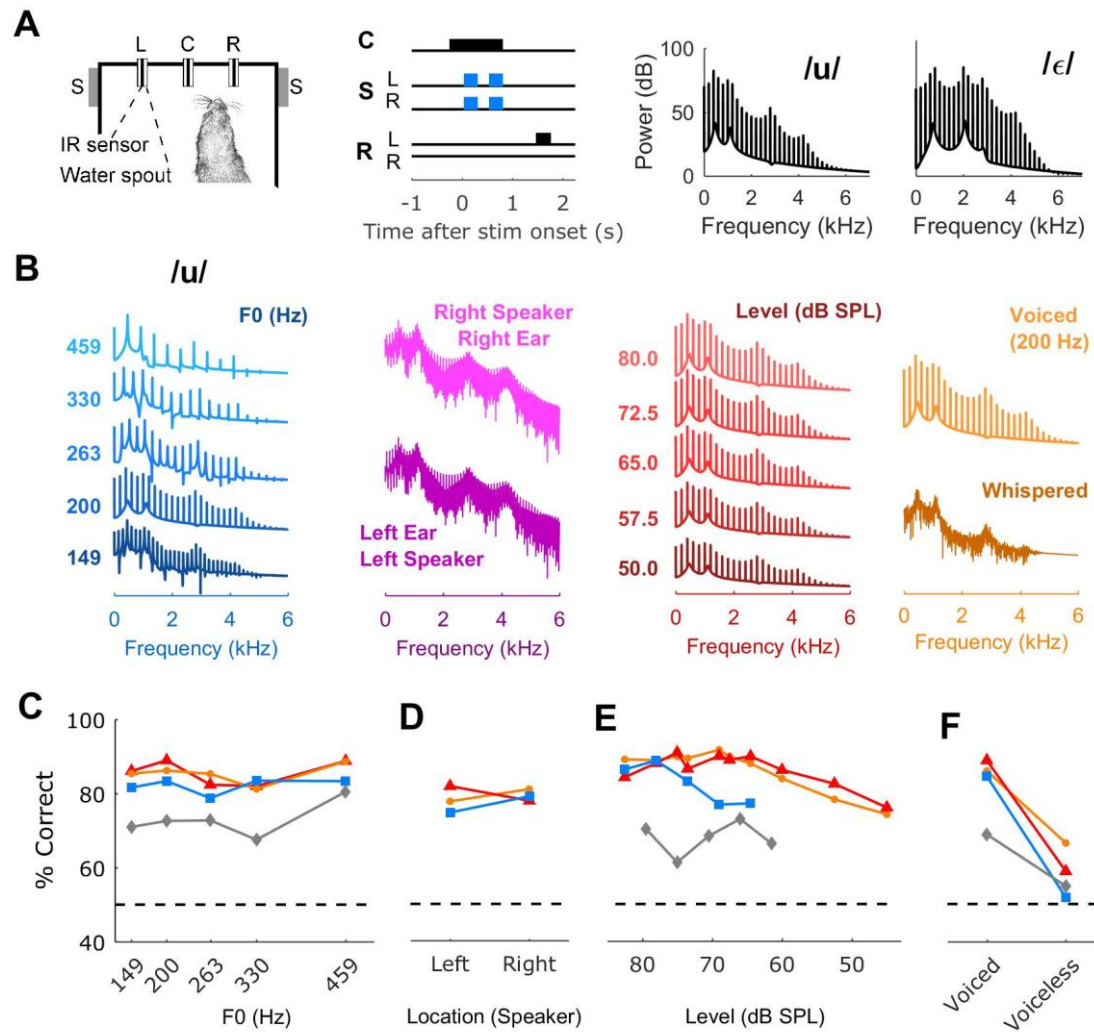## Results

### Perceptual constancy during vowel discrimination

To establish a behavioral model of perceptual constancy, ferrets were trained in a two-choice task (Fig 1A) to identify synthesized vowels varying in F0 (149 – 459 Hz), location (±90°), sound level (45 – 82.5 dB SPL), or voicing (in which vowels were generated to sound whispered and presented on 20% of trials as probe trials). Changes in these task-irrelevant orthogonal dimensions produced different spectra while preserving the formants peaks in the spectral envelope (Fig 1B) critical for vowel identification (Peterson and Barney, 1952; Town and Bizley, 2013). On each trial, the animal visited a central port to trigger presentation of the stimulus: two tokens of *the same* vowel, each lasting 250 ms with an inter-stimulus interval of 250 ms. Subjects then responded at a left / right spout depending on vowel identity, with correct responses rewarded with water and errors leading to a brief timeout (1-5 s). In each test session, vowels varied across only one orthogonal dimension (i.e. F0, level, location or voicing). Variation in each orthogonal dimension was sufficient that had the animals been discriminating these acoustic features performance would have been at, or close to, ceiling (Hine et al., 1994; Sinnott et al., 1992; Walker et al., 2011; Walker et al., 2009; Wood et al., 2017).

Ferrets discriminated vowels accurately across orthogonal dimensions: Performance was consistent across F0s (Fig 1C) and across locations (Fig 1D) (no effect of orthogonal dimension, logistic regression, $p > 0.05$, Table S1) and significantly better than chance at each F0 and location tested (binomial test vs. 50%, $p < 0.001$, Table S2). For all sound levels, performance was also better than chance (Fig 1E and Table S2), however performance increased significantly (but moderately) with sound level in 3 / 4 ferrets ($p < 0.01$; Table S1). Nevertheless performance was constant over a range of intensities and performance at lowest sound levels still exceeded chance. In contrast to the other orthogonal dimensions, ferrets failed to generalize across voicing: performance was significantly worse for whispered than voiced vowels (Fig 1F) and only two ferrets discriminated

93    whispered stimuli better than chance (Table S2). These results confirmed that ferrets perceived a

94    constant vowel identity across variations in acoustic input related to F0, sound location and sound

95    level but not voicing, and so we predicted that we would find tolerant representations of vowel

96    identity across changes in F0, location and sound level.

---

97    **Figure 1 Perceptual constancy during vowel discrimination**

98    **(A)** Schematic of task design: Animals initiated trials by visiting a central port (C) and waiting for a

99    variable period before stimulus presentation. Speakers (S) presented sounds (two tokens *of the*

100   *same* vowel; blue) to the left and right of the head in all conditions except when sound varied across

101   location - in which case they were presented from either left ($S_L$) or right ($S_R$) speaker only. Animals

102   responded at the left or right spout depending on vowel identity. **(B)** Spectra for 13 examples of one

103   vowel /u/ with varying F0, location, sound level and voicing. Spectra for sounds across location were

104   generated in virtual acoustic space (Schnupp et al., 2003) although sounds varied in free-field

105   location. **(C-F)** Behavioral performance when discriminating vowels across F0 (C), location (D), level

106   (E) and location (F). Individual subjects are shown as separate lines.

## Decoding sound features

109    We implanted arrays of independently moveable tungsten microelectrodes in left and right

110    auditory cortex, where electrodes targeted the low frequency reversal between tonotopic primary

111    and posterior fields (Bizley et al., 2005; See Fig S2 in Town et al., 2017). We recorded 471 sound-

112    responsive multi-units and, for each unit, measured responses to vowels across F0, sound location,

113    level and voicing during task performance (Fig 2A). We quantified the information available about

114    vowel identity, F0, etc., by decoding stimulus features in one dimension across changes in the

115    orthogonal dimension from single trial responses. Our decoder compared the Euclidean distances of

116    time-varying patterns of neural activity, with leave-one-out cross validation (Foffani and Moxon,

117    2004)(Fig S1A). The time window over which responses were decoded was variable and we searched

118    for those parameters (start time and window duration) that gave best decoding performance (Fig

119    S1B). Optimization significantly improved decoding performance (Fig S2, rank-sum, $p < 0.001$) and

120    enabled comparison of the time windows over which units were maximally informative. We decoded

121    responses from correct trials only as we reasoned these would provide the clearest demonstration

122    of auditory cortical encoding. For each unit, we reported decoding performance (Fig 2B) and

123    whether the unit could be classified as significantly informative as determined by a permutation test

124    ($p < 0.05$, Fig 2C and Fig S1C), indicating that the unit provided a tolerant representation of vowel

125    identity across variation in sensory input resulting from changes in orthogonal dimensions.

126          We found that the proportion of vowel informative units was highest across dimensions

127    over which behavioral performance was most constant. Across variation in F0, 42.1% of units

128    (154/366) were informative about vowel identity, 43.5% (50/115) were informative across varying

129    sound location and 40.6% (80/197) across varying sound level, whereas only 30.4% (63/207) were

130    informative about vowel identity across voicing (Fig 2C). Furthermore, when we decoded vowel

131    identity at each orthogonal value (Fig 2D) we found that the proportion of vowel informative units

132    was independent of variation in F0 (logistic regression, $\chi^2 = 0.776$, $p = 0.378$), location ($\chi^2 = 2.17$, $p =$

133    0.140), level ($\chi^2 = 0.447$, $p = 0.504$) and voicing ($\chi^2 = 0.983$, $p = 0.321$). Together this suggests that

134    auditory cortical neurons provide representations of vowel identity that are tolerant to variations in

135    acoustic input caused by changes in F0, sound location and level, and to a lesser extent, voicing.

### Conserved information content

137          If units that represent vowel identity across one orthogonal dimension provide a truly

138    tolerant representation, they should also represent vowel identity across multiple orthogonal

139    dimensions. To test this, we counted the number of sound-responsive units from the entire recorded

140    population that were vowel informative across F0, sound location, level and / or voicing.

141          While not every unit was tested across every orthogonal dimension, we found that units

142    remained informative about vowel identity across the dimensions (F0, sound location and level) over

143    which behavioral performance was constant (Fig 2E). Across sound level and location, 38.6% of units

144    (22/57) were vowel informative across both dimensions. This value was close to the proportion of

145    units sensitive to vowel identity across level or location (which were 40% and 43% respectively),

146    indicating that the majority of vowel informative units represented sound identity across both

147    orthogonal dimensions. Similarly, 34.8% of units (54/155) were informative about vowel identity

148    across F0 and level, 34.4% of units (32/93) across F0 and location and 37.0% of units (20/54) were

149    informative across all of F0, sound level and location. In contrast, notably fewer units (<22.5%) were

150    informative about vowel identity over the other combinations of orthogonal factors – of which all

151    included voicing, and across which, animals also generalized poorly. These findings suggest that a

152    sizeable subpopulation of units provide tolerant information about vowel identity across orthogonal

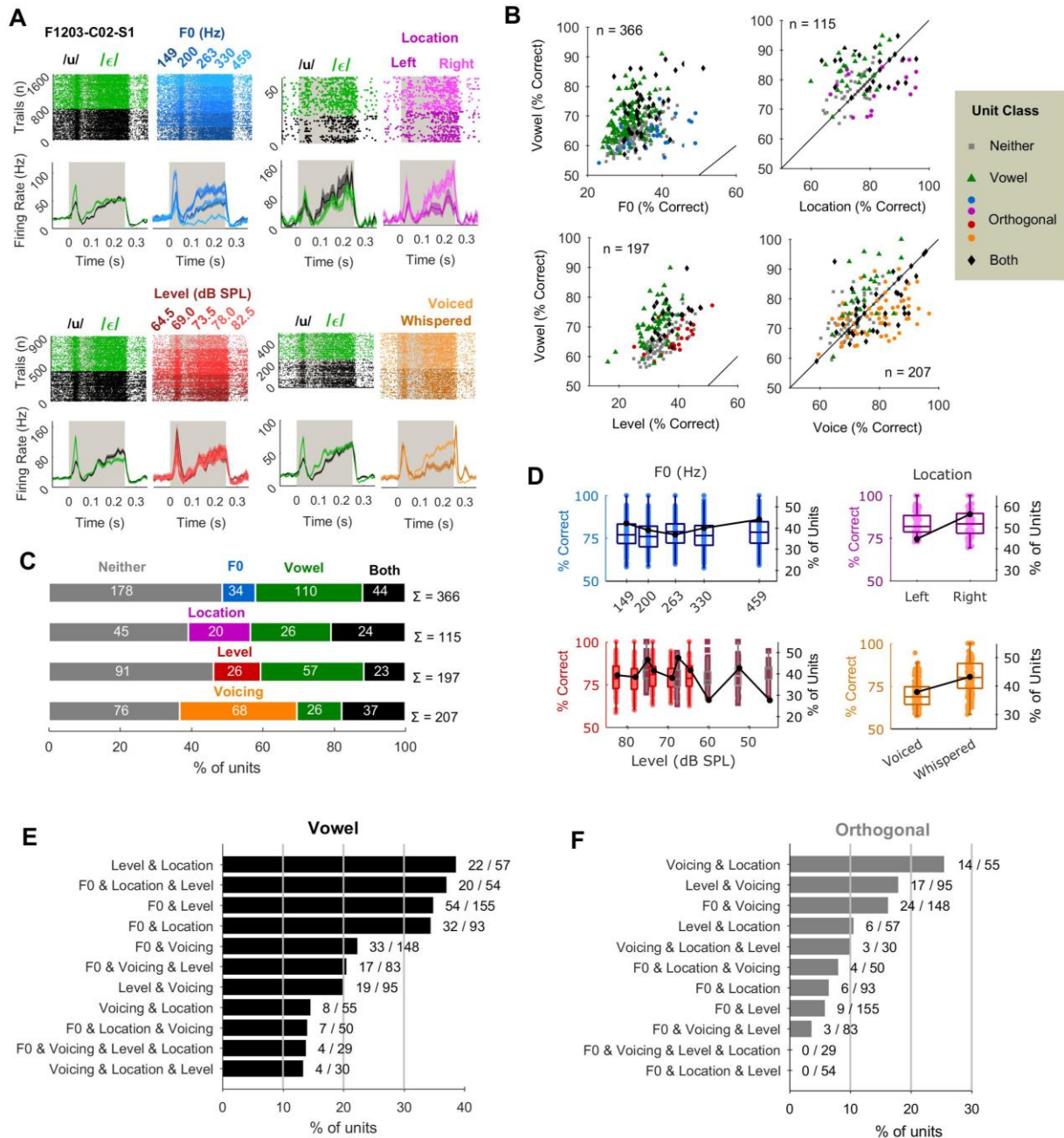153    dimensions during perceptual constancy.

154    **Encoding of orthogonal dimensions**

155        In addition to encoding vowel information during perceptual constancy, we also asked if

156    neural responses conveyed information about orthogonal features of sounds that were irrelevant for

157    task performance (Fig 2D). When considering all F0s or sound levels, we found 21.3% of units

158    (78/366) were informative about F0 across vowels, and 24.9% (49/197) about sound level. These

159    proportions increased to 38.5% (142/369) across F0; and 40.3% (29/72) across sound level when we

160    decoded across the most extreme orthogonal values tested (149 vs 459 Hz or 45 vs 75 dB SPL). While

161    a similar percentage of units (38.3%, 44/115) were informative about sound location, a greater

162    proportion of units (50.7%, 105/207) were informative about voicing. Thus, the balance of units

163    encoding task relevant and orthogonal dimensions was important for perceptual constancy: the

164    proportion of units conveying information about vowel identity was greater than (F0, level) or

165    similar to (location) the dimensions over which animals generalized, whereas across voicing, the

166    balance of informative units was shifted towards the orthogonal dimension (50% to 30%; Fig. 2C).

167     We also tested whether units that were informative about one orthogonal variable (e.g.

168     sound location) were also informative about other orthogonal variables (e.g. voicing). While we

169     observed that some units were informative about multiple orthogonal dimensions (Fig 2F), such

170     groups were significantly smaller than the corresponding analysis of vowel identity (sign-rank test on

171     proportion of conserved units, *p* = 0.0098). Thus, while information about vowel identity was

172     conserved *across* orthogonal dimensions, few units were informative about *multiple* orthogonal

173     dimensions.

---

174     **Figure 2 Neural responses and decoding acoustic features**

175     **(A)** Raster and peri-stimulus time histograms (PSTHs) of neural responses of one unit to vowels

176     across orthogonal variation in F0, location, level and voicing. Data plotted during presentation of the

177     first sound token (grey bar) by vowel identity and by the orthogonal variable. PSTHs show mean ±

178     s.e.m. firing rate. **(B)** Decoding performance when reconstructing vowel identity and orthogonal

179     values from single trial responses of individual units. Data points indicate the best decoding

180     performance of each unit. Chance performance for vowel identity, location and voicing was 50% and

181     20% for F0 and sound. **(C)** Number of units informative about vowel identity and / or orthogonal

182     values when considering responses across all orthogonal values tested. **(D)** Decoder performance

183     and proportion of vowel informative units at each orthogonal value. For sound level, different

184     shades reflect the distinct sound ranges over which units were tested. **(E)** Number of units classified

185     as vowel informative *across* multiple orthogonal dimensions. **(F)** As E but for units classified as being

186     informative *about* multiple orthogonal dimensions.

187

---

188

## Temporal multiplexing of sound features

190   Our data show that multiple sound features are represented in auditory cortex, in some

191   cases by the same neurons. Multivariate encoding in auditory cortex has been linked to temporal

192   multiplexing, where units encode information about different stimulus features at distinct time

193   points (Walker et al., 2011). We therefore asked whether multiplexing occurred during perceptual

194    constancy and if the representation of vowel identity across orthogonal features was matched by

195    conserved timing of information.

196        To study multiplexing, we compared the time windows that gave best performance decoding

197    each stimulus feature following optimization. A time window was defined by its start time and

198    duration, which we summarized as its midpoint (start time + duration/2). For each unit, we

199    measured the midpoint for best decoding vowel identity across each orthogonal dimension, and for

200    decoding each orthogonal dimension across vowel identity. We then compared cumulative

201    distribution functions (CDFs) of midpoints across units that were informative about multiple stimulus

202    features (dual-feature units) or only one feature (single-feature units).

203        We first confirmed the occurrence of multiplexing during perceptual constancy, finding that

204    for dual-feature units, information about each feature emerged over a distinct time-course.

205    Information about vowel identity arose significantly earlier than F0 (Fig 3A, Sign-rank test: z = -2.43,

206    $p$ = 0.015) and sound level (Fig 3C, z = -2.13, $p$ = 0.033) whereas information about sound location

207    arose significantly earlier than vowel identity (Fig 3B, z = 2.26, $p$ = 0.024). In contrast, there was no

208    significant difference in the timing of vowel identity and voicing (z = 1.33, $p$ = 0.184). Thus

209    multiplexing in these units only occurred for sounds that animals showed perceptual constancy.

210        For single-feature units, vowel identity was also best decoded earlier than F0 (Wilcoxon

211    rank-sum test: z = -2.31, $p$ = 0.021) but the differences between decoding of vowel identity and

212    sound level (z = -0.933, $p$ = 0.351), and vowel identity and location (z = 1.29, $p$ = 0.198) were not

213    significant. For single feature units, vowel identity was decoded significantly later than voicing (z
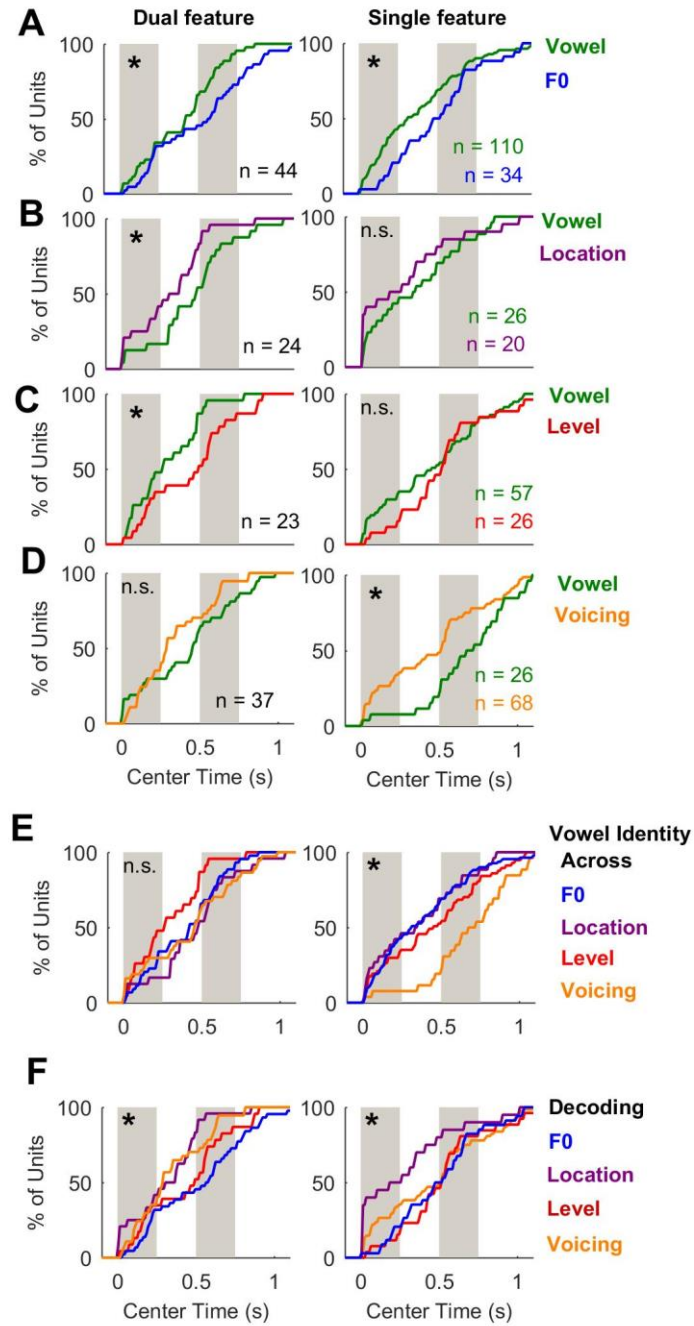
214    =2.79, $p$ = 0.005).

215        Across dual-feature units the timing of vowel information was conserved: CDFs did not differ

216    significantly across orthogonal dimensions (Fig 3E, Kruskal-Wallis test: $\chi^2$ = 5.76, $p$ = 0.124). In

217    contrast, single-feature units showed significant differences in timing of vowel identity information

218    across orthogonal dimensions ($\chi^2$ = 19.95, $p$ = 1.74 x 10$^{-4}$) with post-hoc comparisons showing that

219    information about vowel identity across voicing emerged significantly later than across every other

220    orthogonal factor (Tukey-Kramer corrected, F0: $p$ = 0.001, location: $p$ = 0.002, level: $p$ = 0.013). Dual-

221    feature units also encoded information about different orthogonal dimensions at significantly

222    different times (Fig 3F, $\chi^2$ = 9.77, $p$ = 0.012) with post-hoc comparisons showing F0 was decoded

223    significantly later than location ($p$ = 0.020) and voicing ($p$ = 0.047). Similar results were also found for

224    single feature units where encoding of orthogonal dimensions differed significantly in time ($\chi^2$ =

225    11.04, $p$ = 0.0206) with post-hoc comparisons showing location was decoded significantly earlier

226    than F0 ($p$ = 0.026) and sound level ($p$ = 0.028).

227        These results emphasise an important role for temporal multiplexing: perceptual constancy

228    only occurred when neurons that were sensitive to multiple stimulus features encoded information

229    about each dimension in distinct time windows. Moreover, while the relative timing of vowel

230    information and orthogonal dimensions was not important for generalisation, when vowel

231    information was shifted in time, as in the case of voicing, perceptual constancy failed. Very similar

232    results were observed when considering the start time or decoding window duration (Fig S4 and S5).

233 **Figure 3 Temporal multiplexing in auditory cortex**

234 **(A-D)** Cumulative distributions showing center times for best performance when decoding vowel

235 identity or orthogonal variables (**A:** F0, **B**: location, **C**: level and **D**: voicing). Units are shown

236 separately by classification as informative about vowel identity and orthogonal values (Dual feature

237 units), or only vowel identity or orthogonal values (Single feature units). Grey bars represent the

238 duration of each vowel token. **(E)** CDFs for decoding vowel identity across each orthogonal variable.

239 **(F)** CDFs for decoding orthogonal values across vowels. Asterisks show significant differences

240 between vowel and orthogonal (A-D, rank-sum or sign-rank depending on pairing, $p < 0.05$) or across

241 orthogonal variables (Kruskal-Wallis, $p < 0.05$).

242

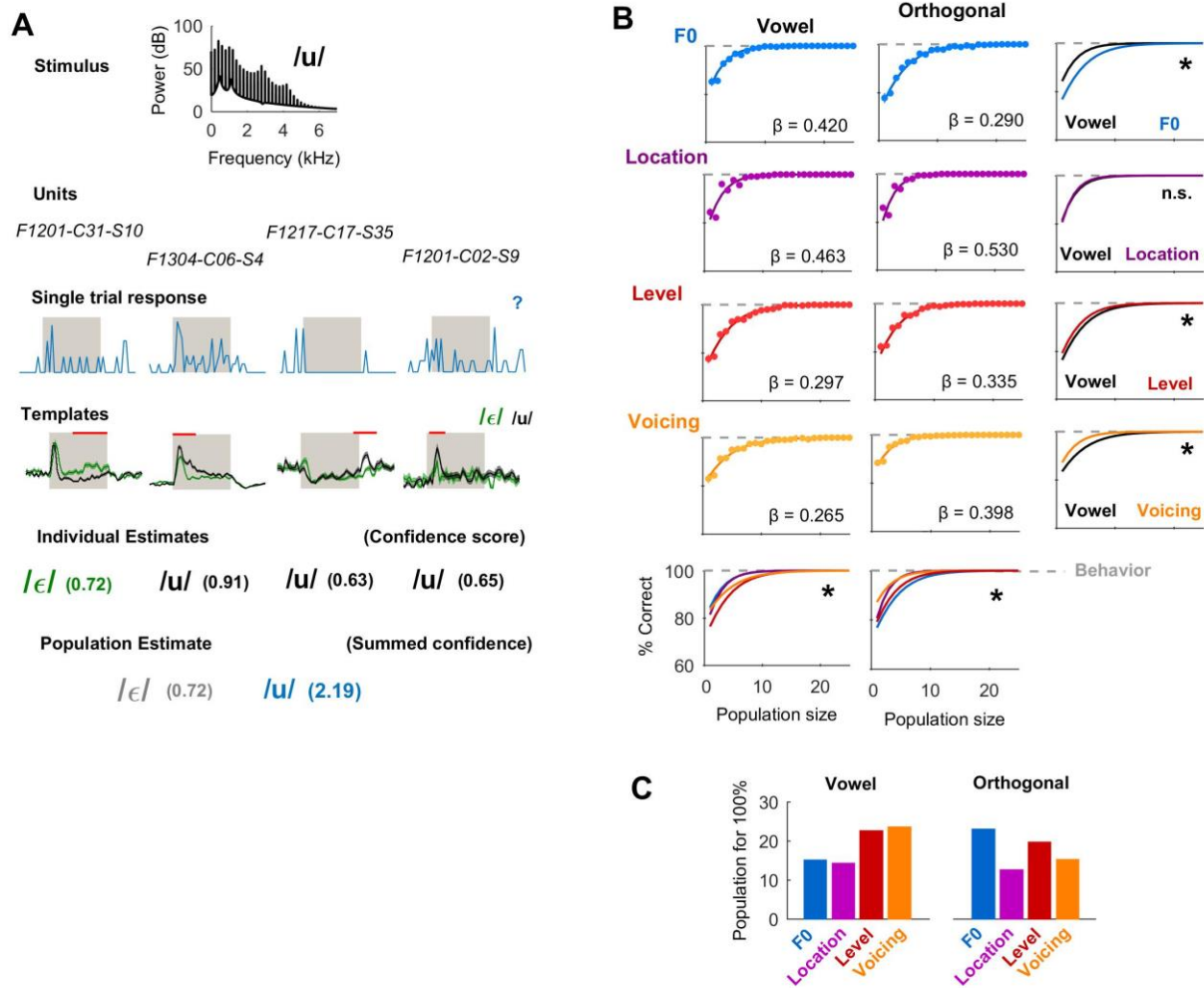**Population decoding matches behavioral performance**

244    Our results show a tolerant representation of vowel identity in auditory cortex during

245    perceptual constancy; however we also wanted to understand how neural encoding was related to

246    behavior. To match the subject's behavioral performance, decoding performance must reach 100%

247    as we only decoded neural responses on correct trials. While few individual units reached this level,

248    it was possible to decode sound features with 100% performance from small populations of units

249    (Fig 4A). Population decoding summed the number of individual units estimating each value of a

250    stimulus feature (e.g. vowel /u/ or /ɛ/) with a weighting based on the relative spike-distance

251    between decoding templates and test trials (see Methods). Decoding improved with population size,

252    following a logistic function (Fig 4B, *p*<0.001) that allowed us to find the minimum number of units

253    required reach 100% performance and compare decoding across conditions (logistic regression,

254    analysis of deviance on main effect of stimulus feature). To equate the number of stimulus features

255    decoded (n=2), we compared decoding across vowel, location and voicing (voiced and whispered)

256    with the F0s (149 and 459 Hz) and sound levels (45 and 75 dB SPL) of greatest separation.

257    Population decoding required fewer units to perfectly reconstruct vowel identity across the

258    orthogonal dimensions across which animals also showed perceptual constancy: To reach 100%

259    performance across F0 required 15 units, and across sound location required 14 units whereas

260    across sound level required 23 units and across voicing required 24 units (Fig 4B; see also Fig S3).

261    Comparison of population growth curves confirmed that the growth of performance with population

262    size was significantly different across the orthogonal dimensions across which vowel identity was

263    decoded ($\chi^2$ = 315.3, *p* < 0.001). Vowel decoding performance increased significantly faster with

264    population size than did decoding of F0 ($\chi^2$ = 319.7, *p* < 0.001), or equivalent to the orthogonal

265    dimension (location, $\chi^2$ = 1.87, *p* = 0.172). Conversely, growth curves for population decoding of

266    vowel identity rose significantly more slowly than for decoding of orthogonal features that animals'

267    failed to generalize across (voicing: $\chi^2$ = 174.9, *p* < 0.001) or incompletely generalised (sound level: $\chi^2$

268    = 40.9, $p$ = 1.62 x $10^{-10}$). Population decoding of orthogonal values also differed significantly across

269    dimensions ($\chi^2$ = 191.5, $p$ < 0.001), with 23 units required to decode F0 with 100% performance, 13

270    units for sound location, 20 units for sound level with 20 units and 15 units for voicing. These

271    findings, suggest that the dynamics of population decoding reflect the ability of animals to

272    generalise: a hallmark of perceptual constancy across a given dimension is that a performance can

273    be supported by a smaller number of units.

---

274    **Figure 4 Population decoding can match behavioral performance**

275    **(A)** Schematic illustration of population decoder in which individual unit estimates of acoustic

276    features (e.g. vowel identity) were weighted using a confidence score. Red lines above templates

277    indicate time window of response considered for each unit. **(B)** Decoding performance obtained with

278    increasing population sizes for decoding of vowel identity and orthogonal values as sounds varied in

279    F0, sound location, level or voicing. Data points show mean ± s.e.m. performance of individual

280    populations for each population size. Curves show logistic regression fits with coefficients (β).

281    Asterisks show significant differences (analysis of deviance, $p$ < 0.001) between vowel and

282    orthogonal curves (right column) or between orthogonal dimensions (bottom row). Behavioral

283    performance (grey lines) was 100% when considering only correct trials. **(C)** Number of units

284    required to decode variables with performance matching animal behavior.

285

---

## Error trials reveal behavioral role for auditory cortex

Population decoding showed that animal's behavioral performance could be matched by information sampled in the responses of small groups of neurons. We next asked how auditory cortical activity was related to behavior by analysing neural responses on error trials. We reasoned that, if activity was relevant for perception, decoding of sound features should be worse when animals made mistakes; whereas if activity was purely stimulus-driven and independent of behavior, decoding should be similar on correct and error trials, as the same stimuli were presented. Using the timing parameters optimized to decode vowel identity on correct trials, we found that, for vowel-informative units, decoding performance was significantly worse on error than correct trials (Fig 5A-B): this was true for sounds that varied across F0 (Wilcoxon sign-rank: $z = -8.64$, $p = 5.83 \times 10^{-18}$),
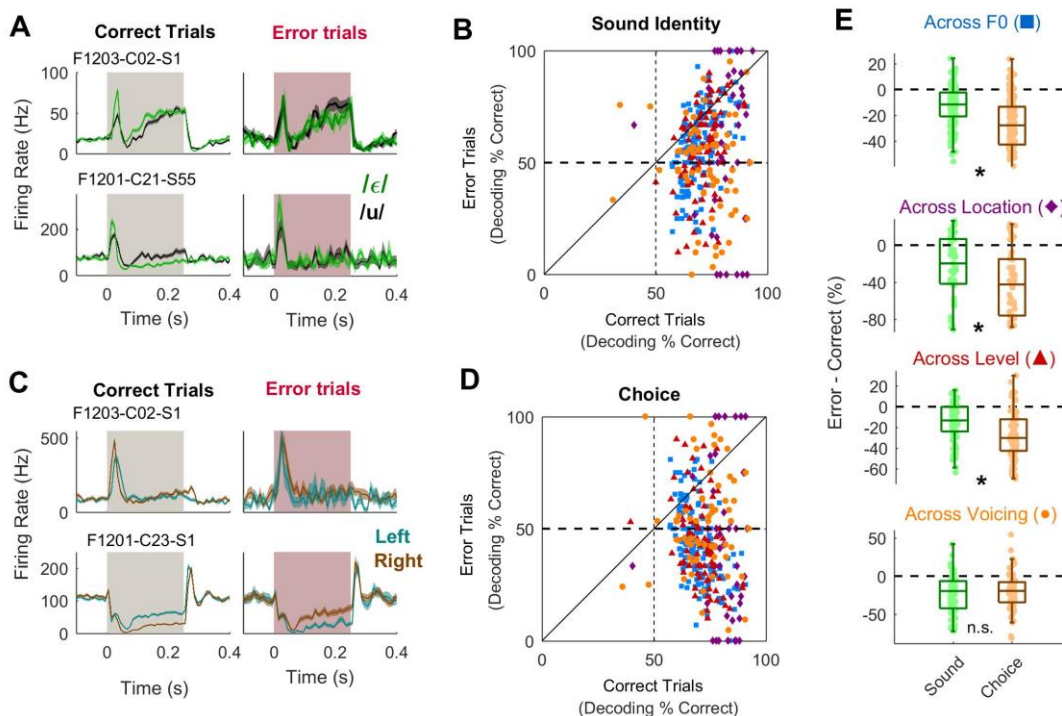
296  location ($z$ = 3.57, $p$ = 3.61 x 10$^{-4}$), sound level ($z$ = 6.07, $p$ = 1.30 x 10$^{-9}$), voicing ($z$ = 5.81, $p$ = 6.16 x

297  10$^{-9}$), and when decoding orthogonal values (Fig S6).

298  The decline in decoding performance we observed on error trials could reflect impairment in

299  the representation of the animals' choice (respond left/right) rather than vowel identity, as choice

300  and vowel were equivalent on correct trials – as correct trials are defined as those on which a

301  specific vowel produces a specific response (e.g. always respond left to /ɛ/). If units were purely

302  choice-driven, decoding of the animals' behavioral response (left or right) should be similar on error

303  and correct trials, while decoding of the stimulus should be significantly worse than chance (50%) as

304  the decoder systematically mis-categorizes trials. While there were many units in which stimulus

305  decoding was substantially below the 50% point, we also saw significantly worse decoding of choice

306  on error trials (Fig 5C-D and Fig S7) when sounds varied across F0 ($z$ = -10.1, $p$ = 7.50 x 10$^{-24}$), location

307  ($z$ = 5.22, $p$ = 1.82 x 10$^{-7}$), level ($z$ = 6.93, $p$ = 4.19 x 10$^{-12}$) and voicing ($z$ = 5.25, $p$ = 1.48 x 10$^{-7}$).

308  To contrast the influence of sensory and choice information on neural activity, we compared

309  the error-related decline in decoding of vowel identity and behavioral choice. Decline in decoding

310  performance was larger for choice than vowel identity when sounds varied across F0 (Fig 5E, rank-

311  sum test, z = 6.54, $p$ = 6.09 x 10$^{-11}$), location (z = 2.48, $p$ = 0.013) or level (z = 4.11, $p$ = 3.89 x 10$^{-5}$) but

312  not voicing (z = 0.473, $p$ = 0.636). These findings suggest that, across auditory cortex, neurons

313  provide a predominantly stimulus based representation whose quality determined the animals'

314  discrimination ability. However, the presence of units in which choice decoding on error trials was

315  maintained, and the observation of units in which decoding of the stimulus on error trials was

316  substantially worse than chance, indicates that the representation is not purely sensory and includes

317  choice information.

318     **Figure 5 Effects of task accuracy on auditory cortex**

319     (**A**) Discrimination of sound (vowel) identity by individual units on correct and error trials. Bars

320     represent the duration of the first vowel token after stimulus onset; neural responses shown as

321     mean ± s.e.m. firing rates. (**B**) Performance decoding sound identity on correct and error trials for all

322     units. Data presented separately for vowels varied across F0, location, sound level and voicing. (**C**)

323     Discrimination of behavioral choice when animals responded at left or right port by individual units

324     on correct and error trials. Data is shown as in (A). (**D**) Performance decoding behavioral choice on

325     correct and error trials for all units. Data shown as in B. (**E**) Comparison of the effects of task

326     accuracy on decoding sound identity and behavioral choice for vowels varied across each orthogonal

327     dimension. Asterisks show significant differences between sound and choice.



328

## Choice and accuracy related activity

330         Our analysis indicated the encoding of both sensory and behavioral variables during

331     perceptual constancy. To study this further, we subsampled neural responses to generate matched

332     datasets containing equal numbers of correct and error trials, vowel identities and choice directions

333     from data across pooled all orthogonal dimensions for which animals showed perceptual constancy

334     (Fig 6A; See Methods). This allowed us to determine modulation of the neural response by the

335     stimulus, the behavioral choice and accuracy (which might indicate confidence or inattention as

336     analysis time windows were restricted to the time before behavioral response). For matched data,

337     behavioral performance would correspond to 50% correct and modulation by one variable (e.g.

338     choice) could not be trivially explained by other variables (e.g. accuracy or sound).
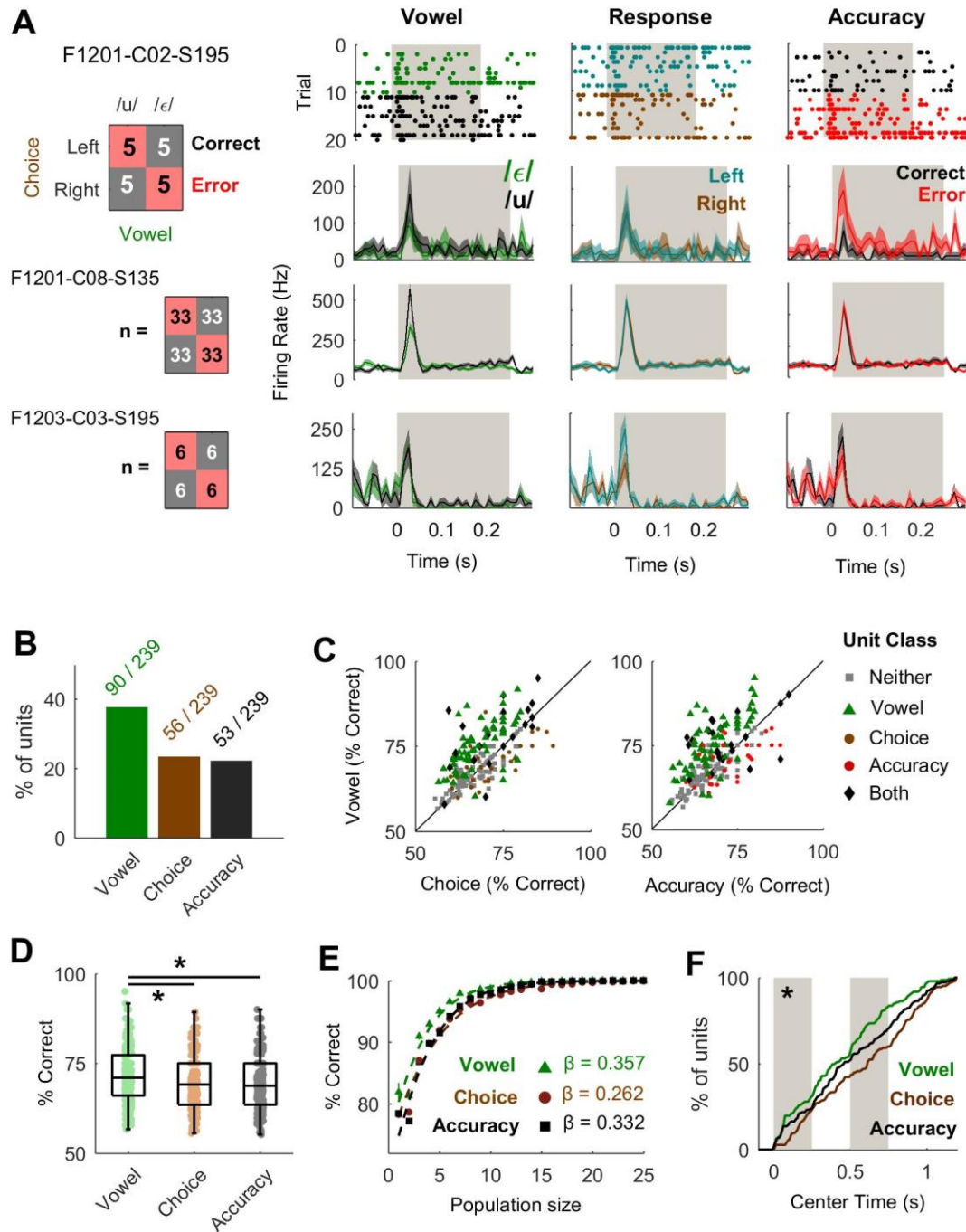
339         When decoding neural responses from matched data, we confirmed that while information

340     about stimulus identity was more widespread than about behavioral variables, units also conveyed

341     information about choice and accuracy: 37.7% of units (90/239) were significantly informative

342     (permutation test, $p < 0.05$) about sound identity, 23.4% (56/239) about choice and 22.2% (53/239)

343     informative about trial accuracy (Fig 6B). Decoding performance was significantly better for vowel

344     identity than for choice or accuracy (Fig. 6C-D, Kruskal-Wallis test: $\chi^2$ = 17.5, $p$ = 1.58 x $10^{-4}$; Tukey-

345     Kramer corrected pairwise comparisons: vowel vs. choice, $p$ = 0.0024; vowel vs. accuracy, $p$ = 3.43 x

346     $10^{-4}$; choice vs. accuracy, $p$ = 0.864). Population decoding plateaued with fewer units when decoding

347     vowel identity (Fig 6E, 18 units required for 100% correct) than decoding accuracy (21 units) or

348     choice (25 units) and performance growth curves differed significantly vowel identity and choice

349     (Bonferroni corrected analysis of deviance, $\chi^2$ = 154.3, $p < 0.001$), and vowel identity and accuracy ($\chi^2$

350     = 109.0, $p < 0.001$). We also compared decoding of choice and accuracy, but found no significant

351     differences ($p > 0.05$) in decoding performance of individual units (Fig 6D) or population decoding

352     functions (Fig 6E). Nonetheless, there were clear representations of behavioral, as well as stimulus,

353     variables, as we could decode the animal's behavioral choice and accuracy better than chance in a

354     substantial proportion (>20%) of individual units and with perfect decoding performance across

355     small populations.

356    Our results show that vowel identity was represented across behavioural variables as well as

357    orthogonal stimulus variations, as units represented vowel identity across the animals' behavioural

358    responses and performance accuracy. Given that auditory cortex multiplexed sound features, we

359    asked if sensory and non-sensory variables were also encoded at different times. We found that

360    information about sound identity emerged earliest, followed by task accuracy and then behavioral

361    choice (Fig 6F): For 147 units that were informative about sound identity, choice and/or accuracy,

362    the time of best decoding differed significantly between dimensions (Kruskal-Wallis test, $\chi^2$ = 15.07,

363    $p$ = 5.35 x $10^{-4}$) with choice represented later than sound identity (Tukey-Kramer corrected, $p$ = 3.07

364    x $10^{-4}$) but timing of information about accuracy not significantly different from either variable ($p$ >

365    0.1). Thus units multiplexed behavioral, as well as sensory variables, with a sequence consistent with

366    sensory-motor transformation, and provided a tolerant representation of sound identity across

367    orthogonal behavioral, as well as acoustic dimensions.

368    **Figure 6 Auditory cortical neurons encode sound identity, behavioral choice and task accuracy**

369    (**A**) Analysis design for matching equal numbers of neural responses to each sound identity, left and

370    right choices, and correct and error trials. Data shown as raster plots of spike times on each trial for

371    one unit and PSTHs representing mean ± s.e.m. firing rate across trials for three units. Grey bars

372    show the first stimulus token. Trial contingency (i.e. respond left for /ɛ/) shown as an example on

373    which one ferret was trained (F1217). (**B**) Percentage of units informative about sound, choice

374    and/or accuracy in matched data. (**C**) Performance decoding sound, choice and accuracy across all

375    units. (**D**) Comparison of performance decoding sound identity, behavioral choice and task accuracy;

376    boxplots show mean and interquartile range. Lines show significant pairwise comparisons ($p$ < 0.01).

377    (**E**) Performance decoding sound, choice or accuracy with populations of units. Data points show

378    mean ± s.e.m. population performance for each population size. (**G**) Cumulative distributions

379    showing center times for best performance when decoding vowel identity behavioral choice and

380    task accuracy. Grey bars represent the duration of the each token within the stimulus. Data shown

381    for all units informative about one or more variables; asterisk reflects significant difference between

382    variables ($p < 0.001$).



383

384    **Effects of Task Engagement**

385          The encoding of animal's choice and accuracy illustrates that auditory cortex processing

386    extends beyond the representation of acoustic input. Attentional state influences auditory cortical

387    activity (Dong et al., 2013; Kuchibhotla et al., 2017; Otazu et al., 2009) and receptive field properties

388    (Atiani et al., 2014; David et al., 2012; Fritz et al., 2003; Jaramillo and Zador, 2011; Lee and

389    Middlebrooks, 2011; Lu et al., 2017; Niwa et al., 2012) . We therefore asked if neural tolerance only

390    emerged during task engagement, by comparing unit responses (e.g. Fig S8) recorded during task

391    performance and during passive listening.

392          We observed that task engagement suppressed spiking responses in the first 100 ms after

393    stimulus onset (Fig 7A, sign-rank test, z = 3.62, $p$ = 2.93 x $10^{-4}$). In the same time window, we

394    decoded vowel identity significantly better from units recorded task-engaged than passively listening

395    animals (Fig 7B, z = -2.83, $p$ = 0.0047). We then expanded our analysis in time to consider effects of

396    engagement with a sliding window, finding that changes in spiking activity and decoding

397    performance were strongly time-dependent: Engagement-related suppression of firing rates

398    occurred throughout stimulus presentation, and contrasted with sustained enhancement of activity

399    in the anticipatory period before stimulus onset (Fig 7C). Furthermore, the *difference in firing rate*

400    between passive and task-engaged units differed significantly with time (one-way anova, $F_{30, 4743}$ =

401    7.08, $p$ = 1.38 x $10^{-28}$). Engagement-related enhancement of vowel decoding was observed at the

402    onset and offset of sounds but not in the sustained period of sound presentation (Fig 7D) and the

403    effect of task-engagement varied significantly with time ($F_{30, 4650}$ = 1.57, $p$ = 0.0247).
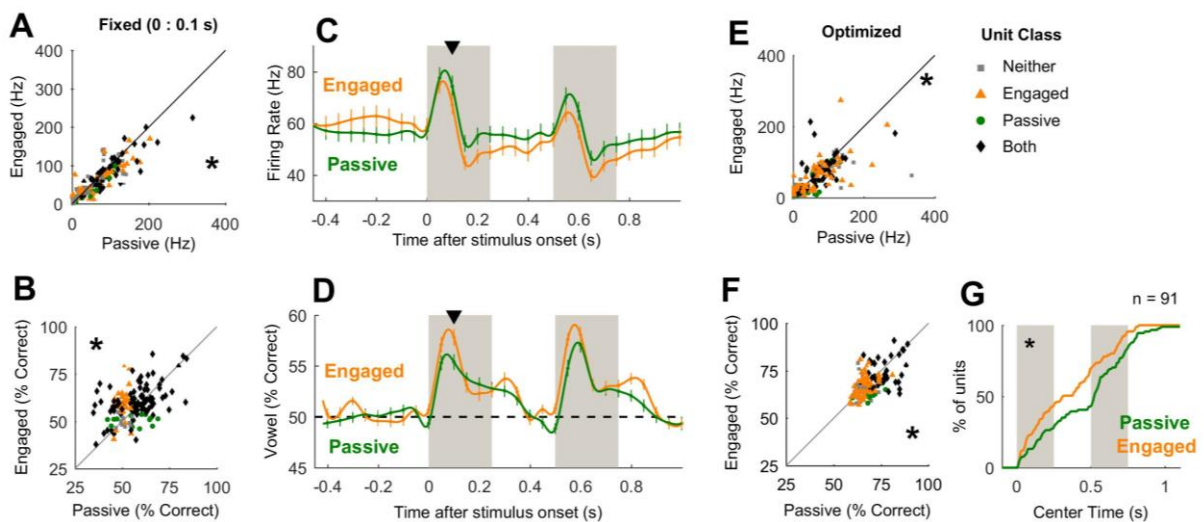
404          To understand how the time-dependent effects of task-engagement modulated overall

405    information content, we compared spiking and vowel decoding in the time window that gave best

406    decoding performance, optimised for each condition independently. Consistent with fixed window

407    analyses, firing rates in optimized windows were lower in the engaged than passive condition (Fig

408    7E; Wilcoxon sign-rank test: z = 3.20, $p$ = 0.0014). However, in contrast to findings with fixed time

409     windows, engagement did not improve optimal decoding performance: for units that were

410     significantly vowel informative during active listening, decoding performance was statistically

411     indistinguishable (z = -0.55, p = 0.582, note firing rate difference was still significant for these units, z

412     = 2.41, $p$ = 0.016) while if all units were considered, there was a small but significant drop in

413     performance (Fig 7F; z = 2.15, $p$ = 0.032). Task engagement similarly affected representation of F0,

414     by supressing spiking activity (z = 2.98, $p$ = 0.003) and decoding performance (z = 4.45, $p$ = 8.47 x 10$^{-}$

415     $^{6}$) in optimized time windows (Fig S9).

416             To understand the origin of differences in fixed-window and optimised analyses, we

417     compared the timing parameters that gave best performance decoding vowel identity, with a focus

418     on units that were significantly vowel informative during task performance. This revealed that the

419     optimized time window for vowel-informative units was significantly earlier during task performance

420     than passive listening (sign-rank test on center time: z = 2.79, $p$ = 0.015). The effects of task-

421     engagement were therefore not to enhance the degree of tolerance of vowel informative units, as

422     judged by optimized decoding of vowel identity across F0, which remained similar across states, but

423     rather to enhance the speed and efficiency of tolerant representations, encoding vowel identity

424     faster and using fewer spikes.

425    **Figure 7 Modulation of auditory processing by task engagement**

426    (**A**) Paired comparison of mean firing rate in the 100 ms after stimulus presentation for units (n =

427    154) recorded during task performance (engaged) and passive listening conditions. Data points show

428    individual units labelled by classification as informative about vowel identity in engaged and passive

429    conditions. (**B**) Paired comparison of performance decoding vowel identity using neural responses

430    measured in the 100 ms after stimulus onset. Individual data points shown each unit (n = 151). Data

431    is shown as in (A). (**C-D**) Paired comparison of firing rate (C) and vowel decoding performance (D) in

432    time windows fixed relative to stimulus onset. Data points show mean ± s.e.m. Black triangles

433    indicate the comparison at 0 – 100 ms in A-B. (**E**) Firing rate in the time window that gave best

434    performance decoding vowel identity (optimized independently for each unit in each experimental

435    condition [passive/ engaged]). Data is shown as in (A). (**F**) Paired comparison of best performance

436    decoding vowel identity in optimized time windows. Data is shown as in (A). (**G**) Cumulative density

437    distributions showing center times giving best decoding performance. Data is shown for units

438    informative about vowel identity during task performance.



439

## Effects of Training

440

441    In addition to task engagement, long-term experience can also affect auditory cortical

442    processing (Bao et al., 2004; Ohl et al., 2001; Polley et al., 2004; Polley et al., 2006; Schnupp et al.,

443    2006; Whitton et al., 2014)(Atilgan et al. Unpublished) and so we also asked if training to

444    discriminate vowels altered auditory representations. We recorded sound-evoked responses (Fig.

445    S8) to vowels in four naïve ferrets (86 units), and in two trained animals presented with untrained

446    vowels (56 units) and compared these with units recorded in trained animals responding to trained

447    vowels (230 units). As we could not pair units across trained and naïve animals, we conducted

448    unpaired comparisons of neural activity (normalized relative to a pre-stimulus baseline period) and

449    decoding performance.

450    Training suppressed neural activity in both comparisons of unit responses to trained and

451    untrained stimuli (Fig 8A) and comparisons of units from trained and naïve animals (Fig 8B): Using a

452    roving analysis window and ANOVA to compare normalized firing rates, with time bin and stimulus

453    training as factors, we found significant effects of time ($F_{30, 8804}$ = 29.0, $p$ < 0.001), training ($F_{1, 8804}$ =

454    24.1, $p$ < 0.001), and a time x training interaction ($F_{30, 8804}$ = 1.89, $p$ = 0.0024). When comparing firing

455    rates in units recorded from trained and naïve subjects, we also found significant effects of time ($F_{30,}$

456    $_{9734}$ = 51.3, $p$ < 0.001), training ($F_{1, 9734}$ = 25.3, $p$ < 0.001), and a time x training interaction ($F_{30, 9734}$ =

457    3.83, $p$ < 0.001).

458    Training also reduced performance decoding vowel identity across F0 in both comparisons of

459    unit responses to trained and untrained stimuli (Fig 8C), and of units recorded in trained and naïve

460    animals (Fig 8D): Comparisons across time (two-way ANOVA) showed significant effects of stimulus

461    training ($F_{1, 8804}$ = 7.69, $p$ = 0.006), time ($F_{30, 8804}$ = 13.6, $p$ < 0.001), and a significant time x training

462    interaction ($F_{30, 8804}$ = 1.65, $p$ = 0.014). Similarly, subject training ($F_{1, 9734}$ = 12.4, $p$ < 0.001) and time

463    ($F_{30, 9734}$ = 16.3, $p$ < 0.001) significantly affected vowel decoding – although we found no significant

464    interaction ($F_{30, 9734}$ = 0.73, $p$ = 0.857). Thus when assessing the effects of training by comparison of

465    stimuli or subjects, neural activity and vowel decoding performance was suppressed.
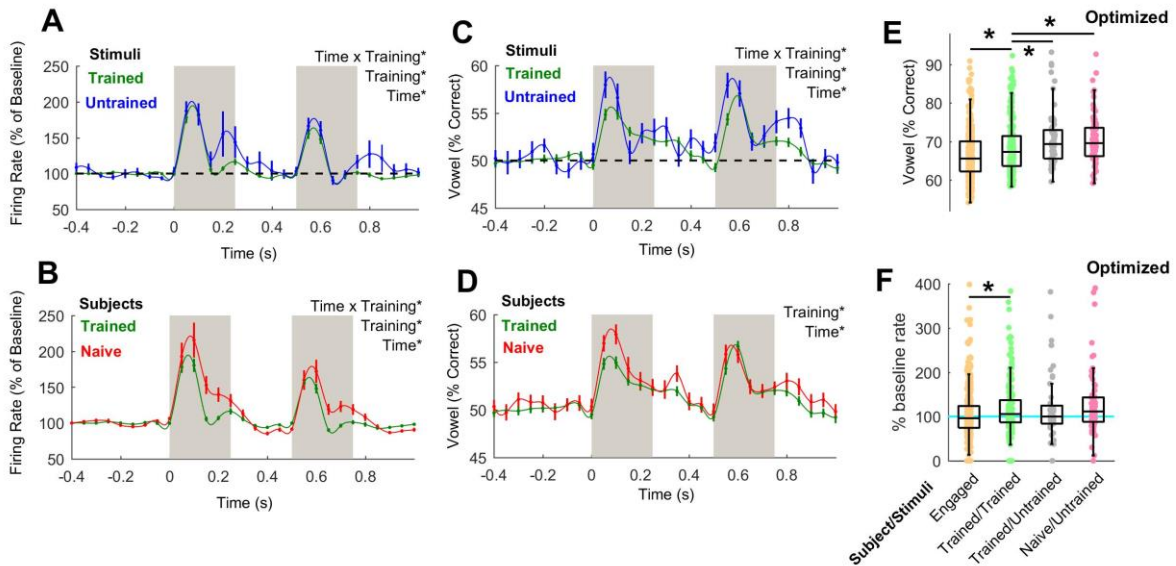
466        Training-related suppression of vowel decoding was also observed when the time

467    parameters of decoding were optimized for each unit (Fig. 8E): Comparing decoding performance

468    across all units recorded in passive conditions revealed a significant effect of experimental group

469    (Kruskal-Wallis test, $\chi^2$ = 12.08, $p$ = 0.002), with pairwise comparisons revealing significant

470    differences between decoding of responses to trained and untrained sounds (Tukey-Kramer

471    corrected, $p$ = 0.046), and between units in trained and naïve animals responding to the same

472    physical stimuli ($p$ = 0.007) but not between units in trained and naïve animals responding to

473    unfamiliar sounds ($p$ = 0.987). In contrast to fixed time window analysis however, we saw no

474    significant effects of training on firing rates in optimized time windows (Fig 8F, $p$ > 0.1). Thus both

475    fixed-window and optimized decoding show training-related reduction in information about vowel

476    identity across F0. Similar results were also found for decoding of F0 across vowels (Fig S10),

477    suggesting that training has broad effects on auditory processing and that information about sound

478    features was, paradoxically, more robust in naïve than trained animals.

479

480    **Figure 8 Modulation of auditory processing by training**

481    (**A-B**) Firing rates of units evoked by trained and untrained sounds (A) and in units recorded from

482    trained and naïve animals (B). Data is shown as mean ± s.e.m. in 100 ms windows at 50 ms intervals

483    with spline interpolation across means. (**C-D**) Unpaired comparison of performance decoding vowel

484    identity from unit responses to trained and untrained sounds (C) and units recorded in trained and

485    naïve animals (D). Data is shown as in A-B. (**E**) Comparison of best performance decoding vowel

486    identity in optimized time window. Individual data points show individual units; box plots show

487    median and inter-quartile ranges. Asterisks show significant comparisons between experimental

488    groups (Tukey correction for multiple comparisons, $p < 0.05$). Effect of task engagement shown for

489    reference. (**F**) Normalized firing rate in the time window giving best performance decoding vowel

490    identity. Data is shown as in F.



491

## Discussion

493    Here we demonstrate that auditory cortical neurons reliably represent vowel sounds across

494    a range of orthogonal acoustic transformations that mirror those preserved in perceptual constancy.

495    The neural representation provided by auditory cortex was multivariate, as units represented

496    multiple stimulus features, and multiplexed, as variables were best represented at different times.

497    Multivariate encoding extended to behavioral dimensions as units represented subjects' choice and

498    accuracy and decoding performance differed between correct and error trials. Consistent with a shift

499    from stimulus-related to task-related neural representation, we found that both task-engagement

500    and long-term training significantly affected the representation of vowel identity in auditory cortex.

501    Together our findings demonstrate that auditory cortical neurons provide a degree of tolerance

502    across variation in sensory input and behavior that was sufficient to represent the identity of target

503    sounds during perceptual constancy.

504    Ferrets identified vowels by their spectral timbre while sounds varied across the major

505    acoustic dimensions key to real-world hearing, including F0 that determines voice pitch, sound

506    location and sound level. Both animals and neurons generalized across the same acoustic

507    dimensions (F0, space etc.). Encoding of multiple features of speech-like sounds, sometimes by the

508    same units, supports previous reports of distributed coding in auditory cortex (Bizley et al., 2009;

509    Griffiths et al., 2010; Ortiz-Rios et al., 2017) and shows that even when potentially disruptive to

510    behavior, orthogonal variables (e.g. F0) are encoded. Furthermore, the encoding of vowel identity by

511    even small populations of units was sufficient to account for, or exceed, animal's behavioral

512    performance. This suggests that auditory cortex provides a multivariate representation of sounds

513    from which downstream neurons may select behaviorally relevant dimensions during perceptual

514    constancy.

515    The point at which multivariate encoding of stimulus features might transition to a

516    univariate representation of a task-relevant dimension is unclear. Here we recorded from a

517    combination of primary and secondary tonotopic areas of auditory cortex; however the limited

518    density of electrodes in our recording array prevented us from mapping the precise boundaries

519    between regions necessary to determine if neural tolerance differed between fields as suggested

520    elsewhere (Carruthers et al., 2015). In future it will be important to record using denser arrays with

521    which tolerant representations can be mapped across auditory cortex and beyond: Neurons in

522    prefrontal cortex (PFC) and higher-order auditory cortex (dPEG) are selective for behaviorally

523    relevant sounds (Atiani et al., 2014; Ding and Simon, 2012; Fritz et al., 2010; Russ et al., 2008;

524    Tsunada et al., 2016) and so we might expect that in such areas, tolerant representations of sound

525    identity are preserved while encoding of orthogonal task-irrelevant dimensions is lost.


526        We decoded vowel identity and orthogonal variables independently and without *a priori*

527    selection of neural response time windows. This approach showed that responses of units

528    informative about both vowel identity and orthogonal features were best decoded in distinct time

529    windows. Temporal multiplexing by units mirrored the time-course of sound perception: Decoding

530    of vowel identity and sound location earlier than voicing or F0 is consistent with perception of sound

531    location and vowel identity at sound onset (Litovsky et al., 1999; Stecker and Hafter, 2002), while

532    listeners require longer to estimate F0 (Gray, 1942; Mckeown and Patterson, 1995; Walker et al.,

533    2011). Best decoding of sound level after vowel identity, sound location or voicing may reflect the

534    time course of temporal integration by the auditory system when assessing moderate level sounds

535    (Buus et al., 1997; Glasberg and Moore, 2002). However, information about voicing was decoded

536    earlier than other stimulus attributes, suggesting that information about harmonicity is available

537    earlier in the neural response, and that temporal multiplexing occurs even when perceptual

538    constancy does not.


539        The order in which acoustic feature representations emerged during perceptual constancy

540    also matched the encoding of vowel identity, F0 and location under anaesthesia (Walker et al.,

541    2011), indicating that multiplexing is a general principle of encoding in auditory cortex. Our work

542    extends these findings to additional acoustic features (voicing and sound level) as well as non-

543    sensory variables (choice and accuracy). Furthermore, our comparison of engaged and passively

544    listening conditions showed that the time-course of multiplexing was plastic and depended on

545    behavioral state. By accelerating the encoding of acoustic variables during task performance,

546    neurons may create time for integration of motor and motivational signals, as well as taught

547    associations (Fritz et al., 2003; Fritz et al., 2010; McGinley et al., 2015; Schneider et al., 2014) in

548    order to coordinate behavioral responses. We would therefore predict that delaying the encoding of

549    acoustic features, but preserving the overall information content of auditory cortical responses,

550    would either disrupt or retard sound discrimination.

551    An open question is why training animals to discriminate sounds reduced information about

552    stimulus features. Such effects are consistent with independent findings that training animals to

553    discriminate vowel identity leads to a reduction in the variation in auditory cortical responses

554    attributable to vowel identity and F0 (Atilgan et al., Unpublished). Furthermore, changes in decoding

555    performance could not be explained trivially by changes in firing rate, as we observed both

556    suppression of neural activity and enhancement of decoding performance (Fig 7), and suppression of

557    decoding performance in the absence of changes in neural activity (Fig 8). One possibility could be

558    that responses to untrained sounds reflect purely feedforward information about sound features

559    extracted earlier in the auditory pathway, but that the association of sounds with non-sensory

560    dimensions in auditory cortex comes at the cost of representing acoustic information.

561    Our findings confirm the importance of behavioral variables in auditory cortical processing

562    (Bizley et al., 2013b; Dong et al., 2013; Niwa et al., 2012): decoding of sound features was impaired

563    on error trials, and we found many units that encoded information about the animals' choice and /or

564    accuracy. The significant drop in decoding performance on error trials, and the sensitivity of units to

565    accuracy, shows that auditory cortical activity is predictive of upcoming mistakes. Given this

566    information, and the finding that stimulus identity could be decoded perfectly from small

567  populations of units, why do animals make errors? One possibility is that errors arise from

568  inattention, which has a distinct neural signature (Lakatos et al., 2016) that our decoder uses to

569  distinguish correct and error trials. At present it is unclear whether the accuracy signal we decode

570  reflects such an attentional lapse or arises as an interaction between representations of sound

571  identity and behavioral choice, or a representation of confidence in auditory processing, or

572  anticipation of reward (Metzger et al., 2006). Future experiments in which confidence or reward

573  value are systematically explored may explain the precise nature of accuracy information reported

574  here.

575  In summary, our results show that during perceptual constancy, neurons in auditory cortex

576  provide tolerant representations of vowel identity and that small populations of units can represent

577  sounds as well as, or better than animal's behavior. Auditory cortical units also encoded information

578  about F0, sound location, level and voicing, as well as the animal's choice and accuracy in the task,

579  each with a specific temporal profile that shows a multivariate and multiplexed system. Task-

580  engagement and training modulated auditory processing, demonstrating a role for attention and

581  long-term experience in perceptual constancy. Across all these variables and experimental

582  conditions, auditory cortical responses showed sufficient tolerance to unambiguously represent

583  vowel identity in the same conditions that animals successfully generalized behavioral performance,

584  and thus provided a neural correlate of perceptual constancy.

## Author contributions:

586  SMT and JKB designed the experiments and wrote the paper; all authors were involved in data

587  collection; SMT analysed the data.

## Acknowledgements

# References

591    Atiani, S., David, S.V., Elgueda, D., Locastro, M., Radtke-Schuller, S., Shamma, S.A., and Fritz, J.B.

593    (2014). Emergent selectivity for task-relevant stimuli in higher-order auditory cortex. Neuron 82,

594    486-499.

595    Bao, S., Chang, E.F., Woods, J., and Merzenich, M.M. (2004). Temporal plasticity in the primary

596    auditory cortex induced by operant perceptual learning. Nat Neurosci 7, 974-981.

597    Bendor, D., and Wang, X.Q. (2007). Differential neural coding of acoustic flutter within primate

598    auditory cortex. Nat Neurosci 10, 763-771.

599    Billimoria, C.P., Kraus, B.J., Narayan, R., Maddox, R.K., and Sen, K. (2008). Invariance and sensitivity

600    to intensity in neural discrimination of natural sounds. Journal of Neuroscience 28, 6304-6308.

601    Bizley, J.K., and Cohen, Y.E. (2013). The what, where and how of auditory-object perception. Nat Rev

602    Neurosci 14, 693-707.

603    Bizley, J.K., Nodal, F.R., Nelken, I., and King, A.J. (2005). Functional organization of ferret auditory

604    cortex. Cerebral cortex 15, 1637-1653.

605    Bizley, J.K., Walker, K.M., King, A.J., and Schnupp, J.W. (2013a). Spectral timbre perception in ferrets:

606    discrimination of artificial vowels under different listening conditions. J Acoust Soc Am 133, 365-376.

607    Bizley, J.K., Walker, K.M., Nodal, F.R., King, A.J., and Schnupp, J.W. (2013b). Auditory cortex

608    represents both pitch judgements and the corresponding acoustic cues. Current Biology 23, 620-625.

609    Bizley, J.K., Walker, K.M., Nodal, F.R., King, A.J., and Schnupp, J.W. (2013c). Auditory cortex

610    represents both pitch judgments and the corresponding acoustic cues. Curr Biol 23, 620-625.

611    Bizley, J.K., Walker, K.M., Silverman, B.W., King, A.J., and Schnupp, J.W. (2009). Interdependent

612    encoding of pitch, timbre, and spatial location in auditory cortex. J Neurosci 29, 2064-2075.

613    Bregman, A.S. (1990). Auditory scene analysis (Cambridge, MA: MIT Press).

614    Buus, S., Florentine, M., and Poulsen, T. (1997). Temporal integration of loudness, loudness

615    discrimination, and the form of the loudness function. J Acoust Soc Am 101, 669-680.

616    Carruthers, I.M., Laplagne, D.A., Jaegle, A., Briguglio, J.J., Mwilambwe-Tshilobo, L., Natan, R.G., and

617    Geffen, M.N. (2015). Emergence of invariant representation of vocalizations in the auditory cortex. J

618    Neurophysiol 114, 2726-2740.

619    David, S.V., Fritz, J.B., and Shamma, S.A. (2012). Task reward structure shapes rapid receptive field

620    plasticity in auditory cortex. Proc Natl Acad Sci U S A 109, 2144-2149.

621    DiCarlo, J.J., and Cox, D.D. (2007). Untangling invariant object recognition. Trends Cogn Sci 11, 333-

622    341.

623    DiCarlo, J.J., Zoccolan, D., and Rust, N.C. (2012). How does the brain solve visual object recognition?

624    Neuron 73, 415-434.

625    Ding, N., and Simon, J.Z. (2012). Emergence of neural encoding of auditory objects while listening to

626    competing speakers. Proceedings of the National Academy of Sciences of the United States of

627    America 109, 11854-11859.

628    Dong, C., Qin, L., Zhao, Z., Zhong, R., and Sato, Y. (2013). Behavioral modulation of neural encoding

629    of click-trains in the primary and nonprimary auditory cortex of cats. J Neurosci 33, 13126-13137.

630    Foffani, G., and Moxon, K.A. (2004). PSTH-based classification of sensory stimuli using ensembles of

631    single neurons. J Neurosci Methods 135, 107-120.

632    Fritz, J., Shamma, S., Elhilali, M., and Klein, D. (2003). Rapid task-related plasticity of spectrotemporal

633    receptive fields in primary auditory cortex. Nat Neurosci 6, 1216-1223.

634    Fritz, J.B., David, S.V., Radtke-Schuller, S., Yin, P., and Shamma, S.A. (2010). Adaptive, behaviorally

635    gated, persistent encoding of task-relevant auditory information in ferret frontal cortex. Nat

636    Neurosci 13, 1011-1019.

637    Ghazanfar, A.A., Turesson, H.K., Maier, J.X., van Dinther, R., Patterson, R.D., and Logothetis, N.K.

638    (2007). Vocal tract resonances as indexical cues in rhesus monkeys. Current Biology 17, 425-430.

639    Glasberg, B.R., and Moore, B.C.J. (2002). A model of loudness applicable to time-varying sounds. J

640    Audio Eng Soc 50, 331-342.

641   Gray, G.W. (1942). Phonemic Microtomy: The Minimum Duration of Perceptible Speech Sounds.

642   Speech Monogr 9, 75-90.

643   Griffiths, T.D., Kumar, S., Sedley, W., Nourski, K.V., Kawasaki, H., Oya, H., Patterson, R.D., Brugge,

644   J.F., and Howard, M.A. (2010). Direct recordings of pitch responses from human auditory cortex.

645   Current biology : CB 20, 1128-1132.

646   Griffiths, T.D., Warren, J.D., Scott, S.K., Nelken, I., and King, A.J. (2004). Cortical processing of

647   complex sound: a way forward? Trends Neurosci 27, 181-185.

648   Hine, J.E., Martin, R.L., and Moore, D.R. (1994). Free-field binaural unmasking in ferrets. Behavioral

649   neuroscience 108, 196-205.

650   Honorof, D.N., and Whalen, D.H. (2010). Identification of speaker sex from one vowel across a range

651   of fundamental frequencies. J Acoust Soc Am 128, 3095-3104.

652   Jaramillo, S., and Zador, A.M. (2011). The auditory cortex mediates the perceptual effects of acoustic

653   temporal expectation. Nat Neurosci 14, 246-251.

654   Kojima, S., and Kiritani, S. (1989). Vocal-Auditory Functions in the Chimpanzee - Vowel Perception.

655   Int J Primatol 10, 199-213.

656   Kuchibhotla, K.V., Gill, J.V., Lindsay, G.W., Papadoyannis, E.S., Field, R.E., Sten, T.A., Miller, K.D., and

657   Froemke, R.C. (2017). Parallel processing by cortical inhibition enables context-dependent behavior.

658   Nat Neurosci 20, 62-71.

659   Lakatos, P., Barczak, A., Neymotin, S.A., McGinnis, T., Ross, D., Javitt, D.C., and O'Connell, M.N.

660   (2016). Global dynamics of selective attention and its lapses in primary auditory cortex. Nat Neurosci

661   19, 1707-1717.

662   Lee, C.C., and Middlebrooks, J.C. (2011). Auditory cortex spatial sensitivity sharpens during task

663   performance. Nat Neurosci 14, 108-114.

664   Litovsky, R.Y., Colburn, H.S., Yost, W.A., and Guzman, S.J. (1999). The precedence effect. Journal of

665   the Acoustical Society of America 106, 1633-1654.

666     Logothetis, N.K., and Sheinberg, D.L. (1996). Visual object recognition. Annu Rev Neurosci 19, 577-

667     621.

668     Lu, K., Xu, Y., Yin, P., Oxenham, A.J., Fritz, J.B., and Shamma, S.A. (2017). Temporal coherence

669     structure rapidly shapes neuronal interactions. Nat Commun 8, 13900.

670     McGinley, M.J., David, S.V., and McCormick, D.A. (2015). Cortical Membrane Potential Signature of

671     Optimal States for Sensory Signal Detection. Neuron 87, 179-192.

672     Mckeown, J.D., and Patterson, R.D. (1995). The Time-Course of Auditory Segregation - Concurrent

673     Vowels That Vary in Duration. Journal of the Acoustical Society of America 98, 1866-1877.

674     Meliza, C.D., and Margoliash, D. (2012). Emergence of Selectivity and Tolerance in the Avian

675     Auditory Cortex. Journal of Neuroscience 32, 15158-15168.

676     Metzger, R.R., Greene, N.T., Porter, K.K., and Groh, J.M. (2006). Effects of reward and behavioral

677     context on neural activity in the primate inferior colliculus. J Neurosci 26, 7468-7476.

678     Musial, P.G., Baker, S.N., Gerstein, G.L., King, E.A., and Keating, J.G. (2002). Signal-to-noise ratio

679     improvement in multiple electrode recording. J Neurosci Methods 115, 29-43.

680     Niwa, M., Johnson, J.S., O'Connor, K.N., and Sutter, M.L. (2012). Activity related to perceptual

681     judgment and action in primary auditory cortex. J Neurosci 32, 3193-3210.

682     Ohl, F.W., Scheich, H., and Freeman, W.J. (2001). Change in pattern of ongoing cortical activity with

683     auditory category learning. Nature 412, 733-736.

684     Ohms, V.R., Gill, A., Van Heijningen, C.A.A., Beckers, G.J.L., and ten Cate, C. (2010). Zebra finches

685     exhibit speaker-independent phonetic perception of human speech. P R Soc B 277, 1003-1009.

686     Ortiz-Rios, M., Azevedo, F.A., Kusmierek, P., Balla, D.Z., Munk, M.H., Keliris, G.A., Logothetis, N.K.,

687     and Rauschecker, J.P. (2017). Widespread and Opponent fMRI Signals Represent Sound Location in

688     Macaque Auditory Cortex. Neuron 93, 971-983 e974.

689     Osmanski, M.S., and Wang, X. (2015). Behavioral dependence of auditory cortical responses. Brain

690     Topogr 28, 365-378.

691    Otazu, G.H., Tai, L.H., Yang, Y., and Zador, A.M. (2009). Engaging in an auditory task suppresses

692    responses in auditory cortex. Nat Neurosci 12, 646-654.

693    Peterson, G.E., and Barney, H.L. (1952). Control methods used in a study of vowels. J Acoust Soc Am

694    24, 175-184.

695    Polley, D.B., Heiser, M.A., Blake, D.T., Schreiner, C.E., and Merzenich, M.M. (2004). Associative

696    learning shapes the neural code for stimulus magnitude in primary auditory cortex. Proc Natl Acad

697    Sci U S A 101, 16351-16356.

698    Polley, D.B., Steinberg, E.E., and Merzenich, M.M. (2006). Perceptual learning directs auditory

699    cortical map reorganization through top-down influences. J Neurosci 26, 4970-4982.

700    Russ, B.E., Orr, L.E., and Cohen, Y.E. (2008). Prefrontal neurons predict choices during an auditory

701    same-different task. Current biology : CB 18, 1483-1488.

702    Sadagopan, S., and Wang, X. (2008). Level invariant representation of sounds by populations of

703    neurons in primary auditory cortex. Journal of Neuroscience 28, 3415-3426.

704    Schebesch, G., Lingner, A., Firzlaff, U., Wiegrebe, L., and Grothe, B. (2010). Perception and neural

705    representation of size-variant human vowels in the Mongolian gerbil (Meriones unguiculatus).

706    Hearing Res 261, 1-8.

707    Schneider, D.M., Nelson, A., and Mooney, R. (2014). A synaptic and circuit basis for corollary

708    discharge in the auditory cortex. Nature 513, 189-194.

709    Schnupp, J.W., Booth, J., and King, A.J. (2003). Modeling individual differences in ferret external ear

710    transfer functions. J Acoust Soc Am 113, 2021-2030.

711    Schnupp, J.W.H., Hall, T.M., Kokelaar, R.F., and Ahmed, B. (2006). Plasticity of temporal pattern

712    codes for vocalization stimuli in primary auditory cortex. Journal of Neuroscience 26, 4785-4795.

713    Sharpee, T.O., Atencio, C.A., and Schreiner, C.E. (2011). Hierarchical representations in the auditory

714    cortex. Curr Opin Neurobiol 21, 761-767.

715    Sinnott, J.M., Brown, C.H., and Brown, F.E. (1992). Frequency and intensity discrimination in

716    Mongolian gerbils, African monkeys and humans. Hear Res 59, 205-212.

717    Smith, D.R.R., Patterson, R.D., Turner, R., Kawahara, H., and Irino, T. (2005). The processing and

718    perception of size information in speech sounds. Journal of the Acoustical Society of America 117,

719    305-318.

720    Stecker, G.C., and Hafter, E.R. (2002). Temporal weighting in sound localization. Journal of the

721    Acoustical Society of America 112, 1046-1057.

722    Town, S.M., Atilgan, H., Wood, K.C., and Bizley, J.K. (2015). The role of spectral cues in timbre

723    discrimination by ferrets and humans. J Acoust Soc Am 137, 2870-2883.

724    Town, S.M., and Bizley, J.K. (2013). Neural and behavioral investigations into timbre perception.

725    Front Syst Neurosci 7, 88.

726    Town, S.M., Brimijoin, W.O., and Bizley, J.K. (2017). Egocentric and allocentric representations in

727    auditory cortex. Plos Biol 15, e2001878.

728    Tsunada, J., Liu, A.S., Gold, J.I., and Cohen, Y.E. (2016). Causal contribution of primate auditory

729    cortex to auditory perceptual decision-making. Nat Neurosci 19, 135-142.

730    Walker, K.M., Bizley, J.K., King, A.J., and Schnupp, J.W. (2011). Multiplexed and robust

731    representations of sound features in auditory cortex. J Neurosci 31, 14565-14576.

732    Walker, K.M., Schnupp, J.W., Hart-Schnupp, S.M., King, A.J., and Bizley, J.K. (2009). Pitch

733    discrimination by ferrets for simple and complex sounds. The Journal of the Acoustical Society of

734    America 126, 1321-1335.

735    Whitton, J.P., Hancock, K.E., and Polley, D.B. (2014). Immersive audiomotor game play enhances

736    neural and perceptual salience of weak signals in noise. Proc Natl Acad Sci U S A 111, E2606-2615.

737    Wood, K.C., Town, S.M., Atilgan, H., Jones, G.P., and Bizley, J.K. (2017). Acute Inactivation of Primary

738    Auditory Cortex Causes a Sound Localisation Deficit in Ferrets. PLoS One 12, e0170264.

739    Zahorik, P., and Wightman, F.L. (2001). Loudness constancy with varying sound source distance. Nat

740    Neurosci 4, 78-83.

741

## Methods

### Animals

Subjects were four pigmented female ferrets (1-5 years old) trained to discriminate vowels across fundamental frequency, sound level, voicing and location (Bizley et al., 2013a; Town et al., 2015). Each ferret was chronically implanted with Warp-16 microdrives (Neuralynx, MT) housing sixteen independently moveable tungsten microelectrodes (WPI Inc., FL) positioned over primary and posterior fields of left and right auditory cortex (Fig S2). Details of the surgical implantation procedures and histological confirmation of electrode position are described elsewhere (Bizley et al., 2013c). A further six ferrets (also pigmented females) implanted with the same microdrives were used as naïve animals for passive recording. These animals were trained in a variety of psychophysical tasks that did not involve the vowel sounds presented here.

Subjects were water restricted prior to testing; on each day of testing, subjects received a minimum of 60ml/kg of water either during testing or supplemented as a wet mash made from water and ground high-protein pellets. Subjects were tested in morning and afternoon sessions on each day for up to five days in a week. Test sessions lasted between 10 and 50 minutes and ended when the animal lost interest in performing the task.

The weight and water consumption of all animals was measured throughout the experiment. Regular otoscopic examinations were made to ensure the cleanliness and health of ferrets' ears. Animals were housed in groups of two or more animals in enriched housing conditions. All experimental procedures were approved by a local ethical review committee and performed under license from the UK Home Office and in accordance with the Animals (Scientific Procedures) Act 1986.

## Apparatus

764

765       Ferrets were trained to discriminate sounds in a customized pet cage (80 cm x 48 cm x 60

766   cm, length x width x height) within a sound-attenuating chamber (IAC) lined with sound-attenuating

767   foam. The floor of the cage was made from plastic, with an additional plastic skirting into which

768   three spouts (center, left and right) were inserted. Each spout contained an infra-red sensor (OB710,

769   TT electronics, UK) that detected nose-pokes and an open-ended tube through which water could be

770   delivered.

771       Sound stimuli were presented through two loud speakers (Visaton FRS 8) positioned on the

772   left and right sides of the head at equal distance and approximate head height. These speakers

773   produce a smooth response (±2 dB) from 200Hz to 20 kHz, with an uncorrected 20 dB drop-off from

774   200 to 20 Hz when measured in an anechoic environment using a microphone positioned at a height

775   and distance equivalent to that of the ferrets in the testing chamber. An LED was also mounted

776   above the center spout and flashed (flash rate: 3 Hz) to indicate the availability of a trial. The LED

777   was continually illuminated whenever the animal successfully made contact with the IR sensor

778   within the center spout until a trial was initiated. The LED remained inactive during the trial to

779   indicate the expectation of a peripheral response and was also inactive during a time-out following

780   an incorrect response.

781       The behavioral task, data acquisition, and stimulus generation were all automated using

782   custom software running on personal computers, which communicated with TDT real-time signal

783   processors (RZ2 and RZ6, Tucker-Davis Technologies, Alachua, FL).

## Task Design, Stimuli and Behavioral Testing

784

785       Ferrets discriminated vowel identity in a two-alternative forced choice task described

786   elsewhere (Town et al., 2015). Briefly, on each trial the animal was required to approach the center

787   spout and hold head position for a variable period (0 – 500 ms) before stimulus presentation. Each

788    stimulus consisted of a 250 ms artificial vowel sound repeated once with an interval of 250 ms.

789    Animals were required to maintain contact with the center spout until the end of the interval

790    between repeats (i.e. 500 – 1000 ms after initial nose-poke) and could then respond at either left or

791    right spout. Correct responses were rewarded with water delivery whereas incorrect responses led

792    to a variable length time-out (3 - 8 s). To prevent animals from developing biases, incorrect

793    responses were also followed by a correction trial on which animals were presented with the same

794    stimuli. Correction trials and trials on which the animal failed to respond within the trial window (60

795    s) were not analysed. The only exception to this protocol was for whispered sounds, which we

796    presented as probe sounds in 10 – 20% of trials on which any response was rewarded and correction

797    trials did not follow.

798        We also tested subjects under passive listening conditions, in which animals were provided

799    with water at the center port to recreate the head position and motivational context occurring

800    during task performance. Sounds were presented with the same two-token stimulus structure as

801    during task performance, with a minimum of 1 second between stimuli. During test sessions, sound

802    presentation began once the animal approached the center spout and began licking and ended

803    when the animal became sated and lost interest in remaining at the spout.

804        Stimuli were artificial vowel sounds synthesized in MATLAB (MathWorks, USA) based on an

805    algorithm adapted from Malcolm Slaney's Auditory Toolbox

806    (https://engineering.purdue.edu/~malcolm/interval/1998-010/). The adapted algorithm simulates

807    vowels by passing a sound source (either a click train, to mimic a glottal pulse train for voiced

808    stimuli, or broadband noise for whispered stimuli) through a biquad filter with appropriate

809    numerators such that formants are introduced in parallel. Four formants (F1-4) were modelled:

810    three subjects were trained to discriminate /u/ (F1-4: 460, 1105, 2857, 4205 Hz) from /ɛ/ (730, 2058,

811    2857, 4205 Hz) while one subject was trained to discriminate /a/ (936, 1551, 2975, 4263 Hz) from /i

812    (437, 2761, 2975, 4263 Hz). Selection of formant frequencies was based on previously published

813    data (Peterson and Barney, 1952; Town et al., 2015) and synthesis produced sounds consistent with

814    the intended phonetic identity. Formant bandwidths were kept constant at 80, 70, 160 and 300 Hz

815    (F1-4 respectively) and all sounds were ramped on and off with 5 ms cosine ramps.

816         To test perceptual constancy, we varied the rate of the pulse train to generate different

817    fundamental frequencies and used broadband noise rather than pulse train to generate whispered

818    vowel. For sound level we simply attenuated signals in software prior to stimulus generation. For

819    sound location, we presented vowels only from the left or right speaker whereas for all other tests

820    sounds were presented from both speakers. Across variations in F0, voicing and space, we fixed

821    sound level at 70 dB SPL. For tests across sound level and location, voiced vowels were generated

822    with 200 Hz fundamental frequency. Sound levels were calibrated using a Brüel & Kjær (Norcross,

823    USA) sound level meter and free-field [1/2] inch microphone (4191) placed at the position of the

824    animal's head during trial initiation.

## Neural Recording

826         Neural activity in auditory cortex was recorded continuously throughout task performance.

827    On each electrode, voltage traces were recorded using TDT System III hardware (RX8 and RZ2) and

828    OpenEx software (Tucker-Davis Technologies, Alachua, FL) with a sample rate of 25 kHz. For

829    extraction of action potentials, data were bandpass filtered between 300 and 5000 Hz and motion

830    artefacts were removed using a decorrelation procedure applied to all voltage traces recorded from

831    the same microdrive in a given session (Musial et al., 2002). For each channel within the array, we

832    identified spikes (putative action potentials) as those with amplitudes between -2.5 and -6 times the

833    RMS value of the voltage trace and defined waveforms of events using a 32-sample window

834    centered on threshold crossings.

835         In the current study, waveform shapes were not sorted and data from multiple test sessions

836    combined across days. The activity for each unit thus represents the unsorted multi-unit activity of a

837    small population of cells at the recording site. We identified sound responsive units in task-engaged

838    animals as those whose stimulus evoked response within the 300 ms after onset of first token

839    differed significantly from spontaneous activity in the 300 ms before making contact with the spout

840    (Sign-rank test, $p < 0.05$). In passive conditions, we identified responsive units using a similar

841    comparison, but spontaneous activity was measured in the 300 ms before stimulus presentation.

## Decoding procedure

843        We decoded stimulus features (e.g. vowel identity, F0 etc.) on single trials using a simple

844    spike-distance decoder with leave-one-out cross-validation (LOCV). For every trial over which an

845    individual unit was tested in a given dataset (e.g. vowels varied across F0 during task performance),

846    we calculated template responses for each stimulus class (e.g. each vowel or each F0) as the mean

847    PSTH of responses on all other trials. We then estimated the stimulus feature on the test trial as the

848    template with the smallest Euclidean distance to the test trial (Fig S1A). Where equal distances were

849    observed between test trial and multiple templates, we randomly estimated (i.e. guessed) which of

850    the equidistant templates was the true stimulus feature. This procedure was repeated for all trials

851    and decoding performance was measured as the percentage of trials on which the stimulus feature

852    was correctly recovered. Although this approach was simple and did not account for the variance of

853    neural activity, it provided a simple and intuitive relationship between neural activity and

854    information content that we could use with small data sets (sample sizes down to five trials per

855    condition). Robustness to sample size was particularly important because the animal's behavior

856    determined the number of trials in each condition and we aimed to analyse as many units as

857    possible rather than develop a more sophisticated decoder.

858        Auditory cortical units showed a wide variety of response profiles that made it difficult to

859    select a single fixed time window over which to decode neural activity. To accommodate the

860    heterogeneity of auditory cortical neurons and identify the time at which stimulus information

861    arose, we repeated our decoding procedure using different time windows (n = 1550) varying in start

862    time (-0.5 to 1 s after stimulus onset, varied at 0.1 s intervals) and duration (10 to 500 ms, 10 ms

863     intervals) (Fig S1B). Within this parameter space we then reported the parameters that gave best

864     decoding performance. Where several parameters gave best performance we reported the time

865     window with earliest start time and shortest duration.

866         To assess the significance of decoding performance, we conducted a permutation test in

867     which the decoding procedure (including temporal optimization) was repeated 100 times but with

868     vowel identity randomly shuffled between trials to give a null distribution of decoder performance

869     (Fig S1C). The null distribution of shuffled decoding performance was then parameterized and fitting

870     a Gaussian probability density function, for which we then calculated the probability of observing

871     the real decoding performance. Units were identified as informative when the probability of

872     observing the real performance was less than 0.05. Parameterization of the null distribution was

873     used to reduce the number of shuffled iterations over which decoding was repeated. This was

874     necessary because the optimization search for best timing parameters dramatically increased the

875     computational demands of decoding.

876     **Population Decoding**

877         To decode vowel identity from the single trial responses of populations of units, we simply

878     the summed the number of units that estimated each stimulus, weighted by the confidence of each

879     unit's estimate, and took the stimulus with the maximum value as the population estimate. Weights

880     for individual unit (w) estimates were calculated as

$$w = \frac{d_{min}}{\sum_{j=1}^{n} d_j}$$

881     Where *n* was the number of stimulus classes (e.g. vowel identities) and *d* was the spike distance

882     between a test trial response and response templates generated for each stimulus class. Here, $d_{min}$

883     represents the minimum spike distance that indicated the estimated stimulus for that unit.

884     We tested populations of up to 35 units, by which point decoder performance had typically

885     saturated at 100% (with the exception of decoding F0 and sound level across larger [n = 5] numbers

886     of feature classes [e.g. 149, 200, 263, 330 and 459 Hz]). Populations were constructed first by

887     selecting the top 35 units that performed best at decoding the relevant parameter at the individual

888     unit level. Within this subpopulation, we randomly sampled 100 combinations of units without

889     replacement from the large number of possible combinations of units available.

890     **Data Analysis**

891     **Behavior:** Perceptual constancy was reported when the orthogonal factor did not significantly affect

892     task performance, i.e. the likelihood of responding correctly. To test this, we analysed the proportion

893     of correct trials as a function of each orthogonal dimension (e.g. F0) using a logistic regression (Table

894     S1). Regressions were performed separately for each animal, and each orthogonal dimension, and

895     any significant effect ($p < 0.05$) was reported as a failure of constancy. We also asked if an animal's

896     performance at specific orthogonal values was better than chance (50%) using a binomial test ($p <$

897     $0.001$, Table S2).

898     **Neural activity:** The times of spikes was referenced to the onset of the stimulus on each trial and

899     used to create raster and peri-stimulus time histograms. In our analysis of task engagement and

900     training, we measured on each trial the firing rate in 100 ms bins after stimulus onset at 50 ms

901     intervals. For paired comparisons, firing rates in engaged and passively listening animals was

902     compared using a Wilcoxon sign-rank test. For unpaired analyses, we normalized firing rates in these

903     bins relative to the firing rate in a pre-stimulus baseline period in the 450 ms before stimulus onset

904     (passively listening animals) or before the animal began waiting at the center spout (task-engaged

905     animals). Across passively listening groups presented with familiar / unfamiliar sounds (Fig 8E), we

906     compared normalized firing rates and baseline firing rates (i.e. the normalization factors in each

907     condition) across groups using a Kruskal-Wallis test with pairwise post-hoc comparisons performed

908     with Tukey-Kramer correction for multiple comparisons.

909    **Individual unit decoding:** In addition to classifying whether units were informative about a particular

910    stimulus feature (permutation test, $p < 0.05$), we also compared decoding performances (Fig 5B, 5D,

911    6D, 7B, 7F, 8E, S6E, S7, S9A and S10A). When comparing decoding performance across more than

912    two conditions (i.e. in passively listening animals; Fig 8E), data were analysed using a Kruskal-Wallis

913    test with Tukey-Kramer corrected post-hoc comparisons where relevant. When comparing two

914    conditions directly, we used a Wilcoxon sign-rank test for paired data (e.g. comparing performance

915    on correct and error trials; Fig 5B). For comparison of changes in decoding performance between

916    conditions (e.g. decoding sound identity vs. choice on correct and error trials; Fig. 5E), we used a

917    Wilcoxon rank-sum comparison for unpaired data.

918    **Timing:** For each unit, we determined the timing window for which we achieved best decoding

919    performance (Fig S1B) and took the window center (Fig 3), start time (Fig S4) or window duration

920    (Fig S5). We then compared the change in parameter value (e.g. change in center time) for best

921    decoding of between vowel identity and orthogonal dimensions using a Wilcoxon rank-sum test (Fig

922    3A-D). The same approach was used when comparing the timing of decoding vowel identity and F0

923    in task-engaged and passively listening animals (Fig 7G). We also compared the times of best

924    decoding of vowel identity across orthogonal dimensions using a Kruskal-Wallis anova with Tukey-

925    Kramer correction for post-hoc comparisons (Fig 3E). We used the same approach to compare the

926    decoding of orthogonal dimensions (Fig 3F) and decoding of vowel identity, behavioral choice and

927    accuracy (Fig 6F).

928    **Population decoding:** For each unit in a given population, we generated estimates of the target

929    value on each trial based on the minimum spike distance from templates generated on all other

930    trials (i.e. the same LOCV method as for individual unit decoding – see above). Templates were

931    generated using the timing parameters that gave best decoding in the individual unit case and thus

932    each unit's response was sampled independently. In addition to an estimated target value, we also

933    retained a confidence score for that estimate: the spike distance from test trial to the closest

934    template, expressed as a proportion of the sum of spike distances between test trial and all

935    templates. Across the population, we then summed confidence scores for each possible target value

936    and selected the value with the largest sum as the population estimate for a given trial. We then

937    repeated the procedure across trials to get the decoding performance of a given population.

938        We summarized the relationship between population size and decoding performance by

939    fitting a logistic regression model to the proportion each population scored correct, with population

940    size as a predictor. To compare population decoding across conditions (e.g. decoding vowel identity

941    or sound location; Fig 4B) we fitted a logistic model with and without the condition as an additional

942    predictor and assessed significance of improvement in model fit using an analysis of deviance.

943    **Error trial analysis:** We trained the decoder on correct trials using the LOCV procedure to estimate

944    vowel identity on each individual correct trial from templates built on all-other correct trials. For

945    error trials, we used the training templates calculated across all correct trials and estimated vowel

946    identity on each error trial. Only units that were informative about vowel identity were analysed,

947    with the exception of three units recorded when the animal performed perfectly (i.e. made no

948    errors) when vowels varied across sound location and thus error trials could not be studied. We

949    repeated the same procedure for decoding orthogonal variables using only units informative about

950    the relevant dimension. Decoding performance was compared for vowel identity (Fig 5B),

951    orthogonal values (Fig S6) and for behavioral choice (Fig. 5D) using a Wilcoxon sign-rank test. We

952    compared the change in decoding performance between correct and error trials when decoding

953    vowel identity and behavioral choice using a Wilcoxon rank-sum test (Fig 5E).

954    **Datasets matched for vowel, choice and accuracy:** To study the tolerance of a given unit to

955    behavioral as well as acoustic variables, we subsampled neural responses from all conditions in

956    which animals showed perceptual constancy: Specifically we included sounds varied across F0,

957    sound location and sound level above 60 (three ferrets) or 70 dB SPL (one ferret). We excluded all

958    data when sounds were whispered. To prevent trial outcome (water reward or timeout) from

959    confounding accuracy signals, we also excluded trials on which animals responded within one

960    second of stimulus onset. Following pooling and exclusion, we balanced data sets for the number of

961    each vowel, choice and trial outcome by randomly selecting *N* trials, where N was the minimum

962    number of trials in which any one condition (e.g. left responses to /u/) was tested. As with our

963    earlier decoding analysis, we only considered units for which N ≥ 5. We then decoded vowel identity,

964    behavioral choice and accuracy using the same LOCV decoding procedure described above. We

965    compared decoding performance for vowel identity, choice and accuracy across all units with a

966    Kruskal-Wallis anova and post-hoc comparisons using the Tukey-Kramer correction (Fig 6D).

967 **Supplemental Tables**

968 **Table S1**

969 Results of logistic regressions comparing performance across orthogonal variables.

| Ferret | Orthogonal Dimension | | | |
| | F0 (149, 200, 263, 330 & 459 Hz) | Voicing (Voiced / Whispered) | Sound Level (45 – 82.5 dB SPL) | Location (±90°) (i.e. Left / Right) |
|---|---|---|---|---|
| F1201 | $df = 12034$, $p = 0.7731$ | $df = 5802$, $p < 0.001$ | $df = 3490$, $p < 0.001$ | $df = 879$, $p = 0.219$ |
| F1203 | $df = 9945$, $p = 0.764$ | $df = 4212$, $p < 0.001$ | $df = 3063$, $p = 0.005$ | $df = 701$, $p = 0.185$ |
| F1217 | $df = 4790$, $p = 0.368$ | $df = 3784$, $p < 0.001$ | $df = 2214$, $p < 0.001$ | $df = 145$, $p = 0.523$ |
| F1304 | $df = 1485$, $p = 0.388$ | $df = 617$, $p = 0.002$ | $df = 455$, $p = 0.882$ | Not tested |

970

971 **Table S2**

972 Comparison of observed vowel discrimination against chance performance (50%); data shown as

973 fraction of trials correct and probability of observed performance (binomial test). Orthogonal values

974 tested separately for voicing and sound level when a significant main effect of orthogonal value was

975 observed on behavioral performance (Table S1).

976

| Dimension | Data Set | Ferret | | | |
|---|---|---|---|---|---|
| | | F1201 | F1203 | F1217 | F1304 |
| F0 | All | 10214 / 12036 84.9% p < 0.001 | 8502 / 9947 85.5% p < 0.001 | 3946 / 4792 82.4% p < 0.001 | 1051 / 1487 70.7% p < 0.001 |
| Voicing | Voiced | 3830 / 4472 85.6% $p$ < 0.001 | 2833 / 3207 88.3% $p$ < 0.001 | 2585 / 3038 85.1% $p$ < 0.001 | 333 / 473 70.4% $p$ < 0.001 |
| | Whispered | 870 / 1332 65.3% $p$ < 0.001 | 596 / 1007 59.2% $p$ < 0.001 | 386 / 748 51.6% $p$ = 0.400 | 83 / 146 56.9% $p$ = 0.116 |
| Sound level (dB SPL) | 45 | 138 / 186 74.2% p < 0.001 | 149 / 196 76.0% p < 0.001 | Not tested | Not tested |
| | 52.5 | 151 / 193 78.2% p < 0.001 | 169 / 205 82.4% p < 0.001 | Not tested | Not tested |
| | 60 | 161 / 192 83.9% p < 0.001 | 162 / 188 86.2% p < 0.001 | Not tested | Not tested |
| | 64.5 | 438 / 498 88.0% p < 0.001 | 372 / 414 89.9% p < 0.001 | 341 / 442 77.2% p < 0.001 | Not tested |
| | 67.5 | 175 / 195 89.7% p < 0.001 | 176 / 198 88.9% p < 0.001 | Not tested | Not tested |
| | 69 | 479 / 523 91.6% p < 0.001 | 357 / 397 89.9% p < 0.001 | 345 / 449 76.8% p < 0.001 | Not tested |
| | 73.5 | 460 / 515 89.3% p < 0.001 | 352 / 407 86.5% p < 0.001 | 380 / 457 83.2% p < 0.001 | Not tested |
| | 75 | 168 / 187 89.8% p < 0.001 | 179 / 197 90.9% p < 0.001 | Not tested | Not tested |
| | 78 | 445 / 501 88.8% p < 0.001 | 382 / 434 88.0% p < 0.001 | 392 / 442 88.7% p < 0.001 | Not tested |
| | 82.5 | 445 / 500 89.0% p < 0.001 | 357 / 424 84.2% p < 0.001 | 366 / 424 86.3% p < 0.001 | Not tested |
| | All | Not tested | Not tested | Not tested | 309 / 455 67.9% p < 0.001 |
| Location | All | 701 / 879 79.8% $p$ < 0.001 | 564 / 703 80.2% $p$ < 0.001 | 112 / 145 77.2% $p$ < 0.001 | Not tested |

## Supplemental Figures

**Fig S1 Decoder structure**

**(A)** Schematic showing the decoding of trial parameters (e.g. vowel identity) from single trial neural responses for one unit. We used a leave-one-out cross validation method in which templates were calculated as the mean response to each stimulus class (e.g. vowel) on all but one test trial of the data set. Mean responses were averaged across trials from spike times within a decoding window binned at 10 ms intervals. For the test trial, the decoded estimate of stimulus class was assigned as the template class with the smallest Euclidean distance to the test response. Every trial in the dataset was decoded as a test trial with templates recalculated from all other trials.

987    **(B)** Optimizing timing parameters. To accommodate potential variation in timing of information

988    content, we varied the temporal parameters (start time and duration) that defined the decoding

989    window. Start time was varied from -0.5 to 1 s after stimulus onset in 50 ms intervals; duration was

990    varied between 10 and 500 ms in 10 ms intervals. For each combination of start time and duration,

991    we calculated decoding performance across trials and mapped temporal parameter space using a

992    simple grid search. While this search protocol may not find the true optimal parameters for best

993    decoding performance, it nonetheless enabled us to improve decoding performance and estimate

994    those times in the trial at which information about a given feature was most strongly represented.



995

996    **(C)** Permutation testing of decoder performance. Each test variable (e.g. vowel identity) was shuffled

997    and the decoding was repeated with the full optimization search. For each unit, we repeated this

998    shuffling procedure on 100 iterations (we used a relatively small number of iterations and

999    parameterized the permutation distribution to compromise for the computational cost of

1000   optimization). When shuffled, PSTH responses to each vowel were virtually identical. To determine

1001   whether a unit was informative, we fitted the distribution of best performance values obtained for

1002   each shuffle and calculated the probability of measuring the observed decoding performance.



1003

1004 **Fig S2 Improvement in decoding performance with optimization of time window parameters**

1005 Comparison of performance decoding vowel identity using neural responses in a fixed duration (100

1006 ms centered at different times after stimulus onset) or using optimized time window. Data show

1007 mean ± s.e.m. with individual data points showing individual units for optimized data. For each time

1008 point, optimized decoding performance was significantly better than fixed window performance

1009 (Bonferroni corrected for multiple comparisons, $p < 1 \times 10^{-10}$).



1010

1011

1012     **Fig S3 Decoder performance vs. feature set size**

1013     Population decoding performance for two and five way classification. Data shown as logistic

1014     regression model fits for all populations tested (see Fig 5 of the main text for examples of original

1015     data).



1016

1017 **Fig S4 Temporal profiles for the onset (start time) of best decoding windows**

1018 **(A-D)** Cumulative distributions showing start time for best performance when decoding vowel

1019 identity or orthogonal variables (**A:** F0, **B**: location, **C**: level and **D**: voicing). Units are shown

1020 separately by classification as informative about vowel identity and orthogonal values (Dual feature

1021 units), or only vowel identity or orthogonal values (Single feature units). **(E)** CDFs for decoding vowel

1022 identity across each orthogonal variable. **(F)** CDFs for decoding orthogonal values across vowels.

1023 Asterisks show significant differences between vowel and orthogonal (A-D, rank-sum or sign-rank

1024 depending on pairing, $p < 0.05$) or across orthogonal variables (Kruskal-Wallis, $p < 0.05$).



1025

1026 **Fig S5 Temporal profiles for the duration of best decoding windows**
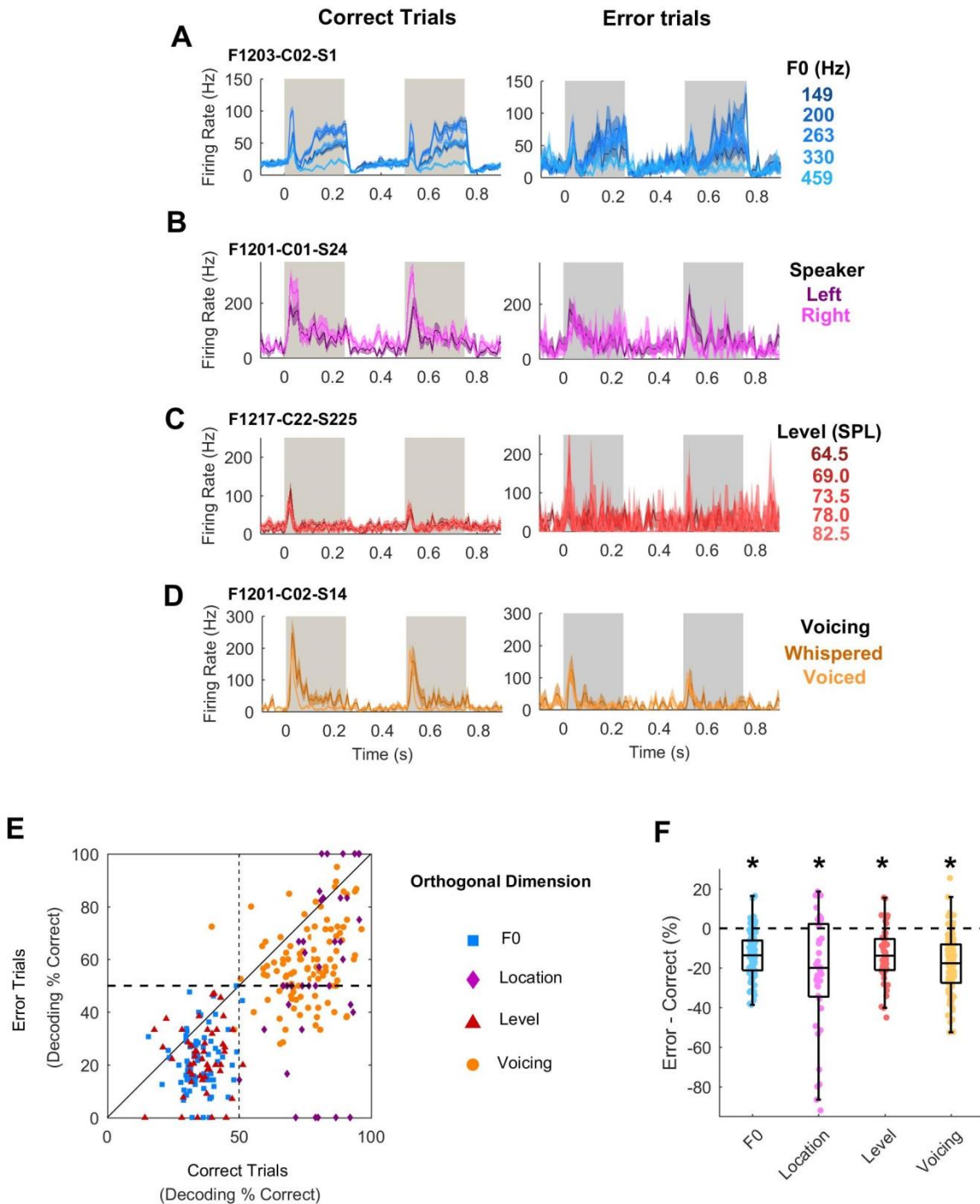
1027 **(A-D)** Cumulative distributions showing duration for best performance when decoding vowel identity

1028 or orthogonal variables (**A:** F0, **B**: location, **C**: level and **D**: voicing). Units are shown separately by

1029 classification as informative about vowel identity and orthogonal values (Dual feature units), or only

1030 vowel identity or orthogonal values (Single feature units). **(E)** CDFs for decoding vowel identity

1031 across each orthogonal variable. **(F)** CDFs for decoding orthogonal values across vowels. Asterisks

1032 show significant differences between vowel and orthogonal (A-D, rank-sum or sign-rank depending

1033 on pairing, $p < 0.05$) or across orthogonal variables (Kruskal-Wallis, $p < 0.05$).



1034

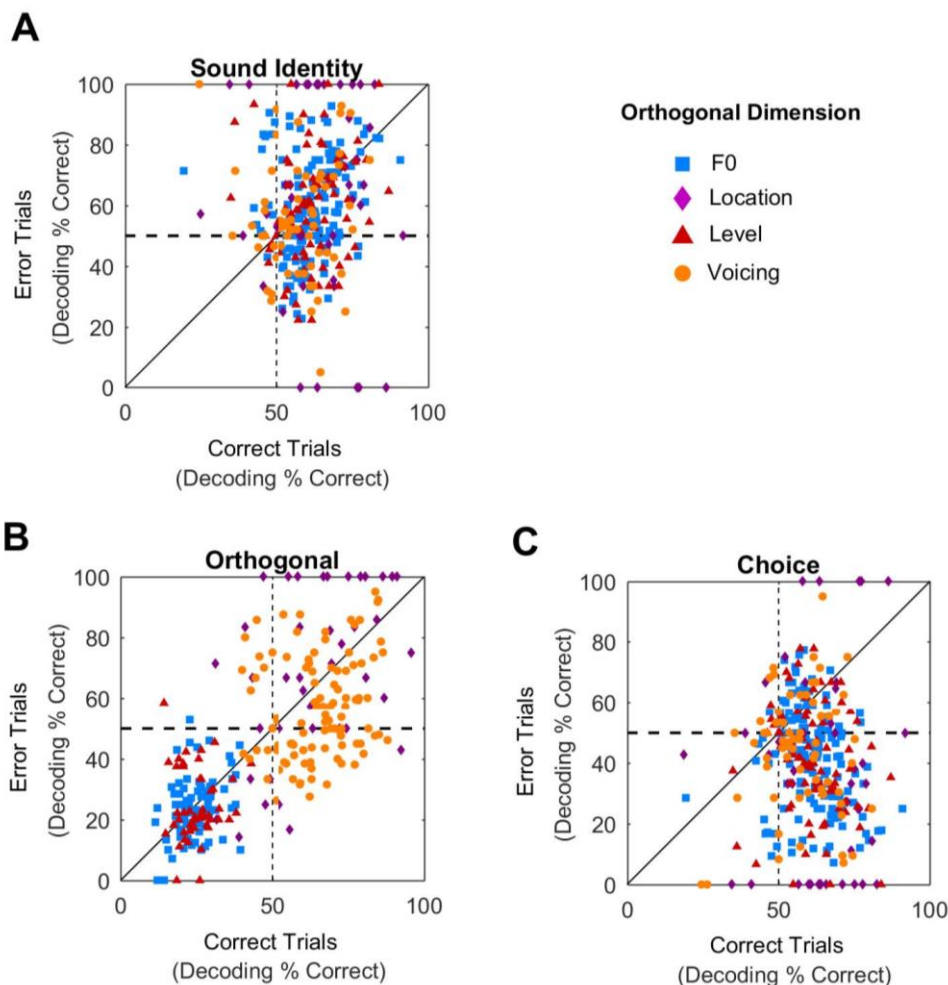1035 **Fig S6 Encoding of orthogonal variables on correct and incorrect trials**

1036 (**A-D**) Example units show encoding of orthogonal variables on correct and error trials. Data is shown

1037 as mean ± s.e.m. (**E-F**) Change in decoding performance from correct to error trials. Asterisks show

1038 significant comparisons (paired t-test: F0: $t_{77}$ = -10.3, $p$ = 3.42 x $10^{-16}$, across location: $t_{40}$ = -4.82, $p$ =

1039 2.12 x $10^{-5}$, across sound level: $t_{48}$ = -6.98, $p$ = 7.85 x $10^{-9}$, across voicing: $t_{104}$ = -11.9, $p$ = 4.52 x $10^{-21}$)



1040

**Fig S7 Error trial decoding performance using a fixed decoding window**

Decoding of vowel identity (**A**), orthogonal values such as F0 or sound level (**B**) and behavioral

response (**C**) on correct and error trials using a fixed time window in the first stimulus presentation

(0 to 250 ms). There was no consistent difference between correct and error trials when decoding

vowel identity or orthogonal values ($p > 0.1$ for all comparisons). However decoding the animal's

response direction was worse on error than correct trials (Across F0: $t_{153} = -12.1$, $p = 3.40 \times 10^{-24}$,

across location: $t_{45} = -5.5$, $p = 1.93 \times 10^{-6}$, across sound level: $t_{79} = -8.1$, $p = 5.77 \times 10^{-12}$, across voicing:

$t_{62} = -4.7$, $p = 1.45 \times 10^{-5}$). The effects of trial accuracy were greater on behavioral response than

stimulus identity when sounds varied across F0 ($t_{306} = 8.4$, $p = 2.15 \times 10^{-15}$), across location ($t_{90} = 3.8$,

$p = 2.94 \times 10^{-4}$), across sound level ($t_{158} = 5.3$, $p = 3.90 \times 10^{-7}$) and across voicing ($t_{124} = 2.9$, $p = 0.004$).
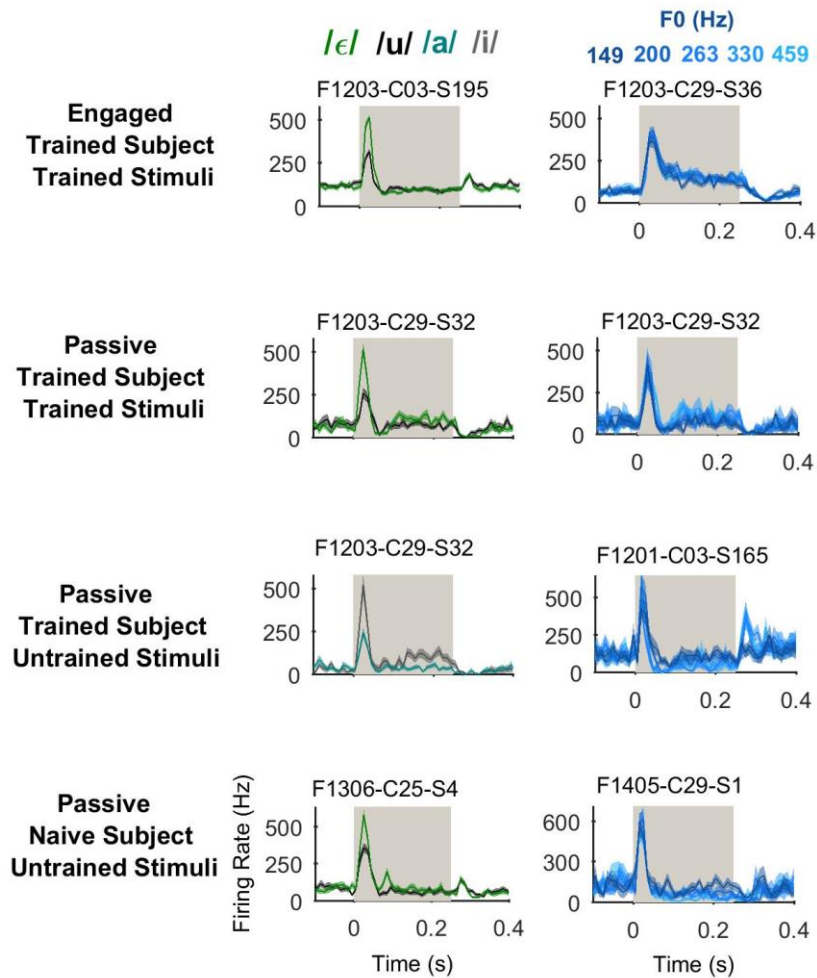
1052    **Fig S8 Example responses during task engagement and passive listening**

1053    Sound-evoked responses of individual unit examples, in task-engaged and passively listening, trained

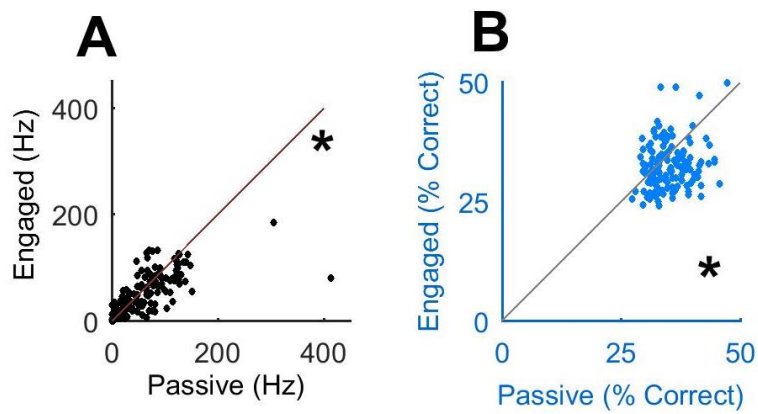1054    and untrained animals to trained and untrained vowels. Plots show mean ± s.e.m. firing rate.



1055

1056

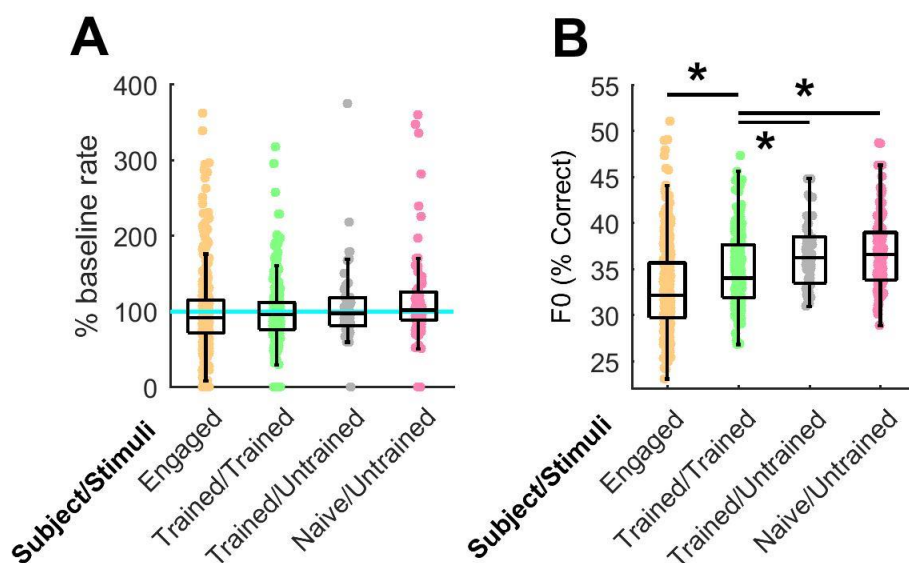1057    **Fig S9 Example responses during task engagement and passive listening**

1058    (**A**) Firing rate in the time window that gave best performance decoding F0 (optimized

1059    independently for each unit in each experimental condition [passive/ engaged]). Data points indicate

1060    individual units (**B**) Paired comparison of best performance decoding F0 in optimized time windows.

1061    Data is shown as in (A). Asterisks show significant engagement-related suppression ($p < 0.05$).



1062

1063 **Fig S10 Effects of training on spiking activity and performance in time windows optimized for**

1064 **decoding F0**

1065 (**A-B**) Comparison of spiking activity normalized to baseline firing (A) and best performance decoding

1066 F0 (B) in optimized time window. Individual data points show individual units; box plots show

1067 median and inter-quartile ranges. Data also shown for task engaged responses for reference.

1068 Asterisks show significant comparisons between experimental groups: Normalized firing rates did

1069 not differ significantly between neurons recorded in any passive conditions, or between units

1070 responding to trained sounds during task engagement and passive listening. Decoding performance

1071 across all units differed significantly between groups (Kruskal-Wallis anova, $\chi^2$ = 21.0, $p$ = 2.76 x $10^{-5}$)

1072 with decoding being significantly worse in units recorded from trained than naïve animals (Tukey-

1073 Kramer corrected for multiple comparisons, $p$ = 1.0 x $10^{-4}$), and worse for units responding to trained

1074 than untrained sounds ($p$ = 0.007). Performance decoding F0 of untrained sounds in naïve and

1075 trained animals was not significantly different ($p$ = 0.935). Decoding performance of units responding

1076 to trained sounds during task engagement was significantly worse than when passively listening ($z$ =

1077 9.81, $p$ = 1.06 x $10^{-8}$).



1078