1 **Title page**
2 Title: Visual modulation of spectrotemporal receptive fields in mouse auditory cortex
3
4 Author names and affiliations:
5 James Bigelow[1,3], Ryan J. Morrill[1,2,3], Timothy Olsen[1,3], Stephanie N. Bazarini[1,3], Andrea R.
6 Hasenstaub[1,2,3]
7 [1]Coleman Memorial Laboratory
8 [2]Neuroscience Graduate Program
9 [3]Department of Otolaryngology–Head and Neck Surgery, University of California, San
10 Francisco, 94143
11
12 Corresponding author:
13 Dr. Andrea R. Hasenstaub
14 Sandler Neurosciences Center
15 Department of Otolaryngology–Head and Neck Surgery
16 675 Nelson Rising Lane
17 Box 0444
18 University of California, San Francisco
19 San Francisco, CA 94143
20

26 **Abstract**

27

28 Recent studies have established significant anatomical and functional connections between
29 visual areas and primary auditory cortex (A1), which may be important for perceptual processes
30 such as communication and spatial perception. However, much remains unknown about the
31 microcircuit structure of these interactions, including how visual context may affect different cell
32 types across cortical layers, each with diverse responses to sound. The present study examined
33 activity in putative excitatory and inhibitory neurons across cortical layers of A1 in awake male
34 and female mice during auditory, visual, and audiovisual stimulation. We observed a
35 subpopulation of A1 neurons responsive to visual stimuli alone, which were overwhelmingly
36 found in the deep cortical layers and included both excitatory and inhibitory cells. Other neurons
37 for which responses to sound were modulated by visual context were similarly excitatory or
38 inhibitory but were less concentrated within the deepest cortical layers. Important distinctions in
39 visual context sensitivity were observed among different spike rate and timing responses to
40 sound. Spike rate responses were themselves heterogeneous, with stronger responses evoked
41 by sound alone at stimulus onset, but greater sensitivity to visual context by sustained firing
42 activity following transient onset responses. Minimal overlap was observed between units with
43 visual-modulated firing rate responses and spectrotemporal receptive fields (STRFs) which are
44 sensitive to both spike rate and timing changes. Together, our results suggest visual information
45 in A1 is predominantly carried by deep layer inputs and influences sound encoding across
46 cortical layers, and that these influences independently impact qualitatively distinct responses to
47 sound.

48

49 **Significance statement**

50

51 Multisensory integration is ubiquitous throughout the brain, including primary sensory cortices.
52 The present study examined visual responses in primary auditory cortex, which were found in
53 both putative excitatory and inhibitory neurons and concentrated in the deep cortical layers.
54 Visual-modulated responses to sound were similarly observed in excitatory and inhibitory
55 neurons but were more evenly distributed throughout cortical layers. Visual modulation
56 moreover differed substantially across distinct sound response types. Transient stimulus onset
57 spike rate changes were far less sensitive to visual context than sustained spike rate changes
58 during the remainder of the stimulus. Spike timing changes were often modulated independently
59 of spike rate changes. Audiovisual integration in auditory cortex is thus diversely expressed
60 among cell types, cortical layers, and response types.

61 **Introduction**

62

63 Evidence accumulated within recent decades demonstrates that primary auditory cortex (A1) is
64 not exclusively involved in processing sound. Instead, A1 integrates information carried by
65 projections from multiple sensory and motor areas with input from the ascending auditory
66 pathway (Schneider and Mooney, 2018; King et al., 2019). For many species including humans,
67 monkeys, and mice, visual projections comprise a particularly dense source of input to A1
68 (Banks et al., 2011). This likely reflects the tendency of environmental features and events to be
69 simultaneously transduced by the auditory and visual modalities, thus giving rise to audiovisual
70 perceptual processes such as spatial localization and communication. Consistent with these
71 observations, physiological studies have found that sound-evoked responses may be modulated
72 by simultaneously presented visual stimuli, either increasing or decreasing firing rates relative to
73 sound alone (Bizley and King, 2008; Kayser et al., 2009).
74 Most physiological studies of audiovisual integration in A1 have not investigated
75 potential differences among cortical layers and cell types (e.g., excitatory vs. inhibitory). This is
76 surprising, as an essential aspect of cortical organization is its division into layers, each with
77 distinct cell type compositions and connectivity patterns with cortical and subcortical structures.
78 Consistent with these anatomical differences, numerous studies in A1 have reported differences
79 in sound encoding properties of neurons across cortical layers (Atencio et al., 2009) and cell
80 types (Atencio and Schreiner, 2008; Phillips et al., 2017b). These findings raise the possibility
81 that multisensory integrative properties of A1 might similarly vary by cortical layer and cell type.
82 Indeed, a recent study from our lab found that a subset of neurons in mouse A1 were
83 responsive to visual flash stimuli, and that these neurons were concentrated in the infragranular
84 layers (Morrill and Hasenstaub 2018). However, this study left open the question of whether
85 neurons with audiovisual integrative responses, such as visual-modulated responses to sound
86 or responses to both modalities, are distributed in parallel with the infragranular visual-
87 responsive neurons. Similarly, whether unimodal visual responses or audiovisual integrative
88 responses differ between neuron types (excitatory, inhibitory) remains to be investigated.
89 With few exceptions (Kayser et al., 2010; Atilgan et al., 2018), studies examining
90 multisensory integration in A1 and elsewhere have relied on changes in time-averaged spike
91 rates to quantify integrative effects. However, neurons throughout the auditory pathway may
92 encode sound features and other events through changes in spike timing, rate, or both
93 (deCharms and Merzenich 1996; Malone et al., 2010; Insanally et al., 2019). Capturing both
94 spike rate and timing changes in A1 may be fundamental to understanding audiovisual
95 integration for two reasons. First, neurons with spike timing changes alone are common in A1,
96 in some preparations reflecting the majority (Insanally et al., 2019). Thus, focusing on spike-rate
97 changes alone may underestimate the prevalence or strength of multisensory integrative
98 activity. Second, downstream targets of A1 may be differently influenced by spike rate and
99 timing changes. Capturing spike timing effects may therefore provide insight into multisensory
100 processes in structures receiving projections from A1.
101 Spike-rate changes are themselves multifaceted and may include transient firing
102 changes at stimulus onset or offset, as well as sustained changes throughout the stimulus
103 period (Lu et al., 2001; Wang et al., 2005; Malone et al., 2015). These diverse response types
104 may reflect distinct network states (Churchland et al., 2010) and sources of information from the

105 ascending auditory pathway (Liu et al., 2019). Resolving potential differences in multisensory
106 integrative properties among response types may be similarly fundamental to understanding the
107 nature and extent of multisensory integration in A1.
108     The current study examined single unit responses in awake mouse A1 to auditory,
109 visual, and audiovisual stimulation. High-density multichannel electrode arrays enabled cortical
110 depth estimation for each neuron and physiological features permitted classification of putative
111 excitatory and inhibitory units. Broadband receptive field estimation stimuli delivered in
112 segments enabled measurement of both transient onset and sustained firing rate responses, as
113 well spectrotemporal receptive fields (STRFs), which are sensitive to both spike rate and timing
114 changes.
115
116 **Materials and Methods**
117
118 *Subjects and surgical preparation*
119
120 All procedures were approved by the Institutional Animal Care and Use Committee at the
121 University of California, San Francisco. A total of 15 adult mice (6 female) served as subjects
122 (median age 99 days, range 58–169 days). All mice had a C57BL/6 background and expressed
123 optogenetic effectors targeting interneuron subpopulations, which were not manipulated in the
124 current experiment. Mice were housed in groups of two to five under a 12H-12H light-dark cycle.
125 Surgical procedures were performed under isoflurane anesthesia with perioperative analgesics
126 (lidocaine, meloxicam, and buprenorphine) and monitoring. A custom stainless steel headbar
127 was affixed to the cranium above the right temporal lobe with dental cement, after which
128 subjects were allowed to recover for at least two days. Prior to electrophysiological recording, a
129 small craniotomy (~1–2 mm diameter) centered above auditory cortex (~2.5–3.5 mm posterior
130 to bregma and under the squamosal ridge) was made within a window opening in the headbar.
131 The craniotomy was then sealed with silicone elastomer (Kwik-Cast, World Precision
132 Instruments). The animal was observed until ambulatory (~5–10 min) and allowed to recover for
133 a minimum of 2 h prior to electrophysiological recording. The craniotomy was again sealed with
134 silicone elastomer at the conclusion of recording, and the animal was housed alone thereafter.
135 Electrophysiological recordings were conducted for each animal on up to five consecutive days
136 following the initial craniotomy procedure.
137
138 *Auditory and visual stimuli*
139
140 All stimuli were generated in MATLAB (Mathworks) and delivered using Psychophysics Toolbox
141 Version 3 (Kleiner et al., 2007). Sounds were delivered through a free-field electrostatic speaker
142 (ES1, Tucker-Davis Technologies) approximately 15–20 cm from the left (contralateral) ear
143 using an external soundcard (Quad Capture, Roland) at a sample rate of 192 kHz. Sound levels
144 were calibrated to 60 ± 5 dB at ear position (Model 2209 meter, Model 4939 microphone, Brüel
145 & Kjær). Visual stimuli were presented on a 19-inch LCD monitor with a 60 Hz refresh rate
146 (ASUS VW199 or Dell P2016t) centered 25 cm in front of the mouse. Monitor luminance was
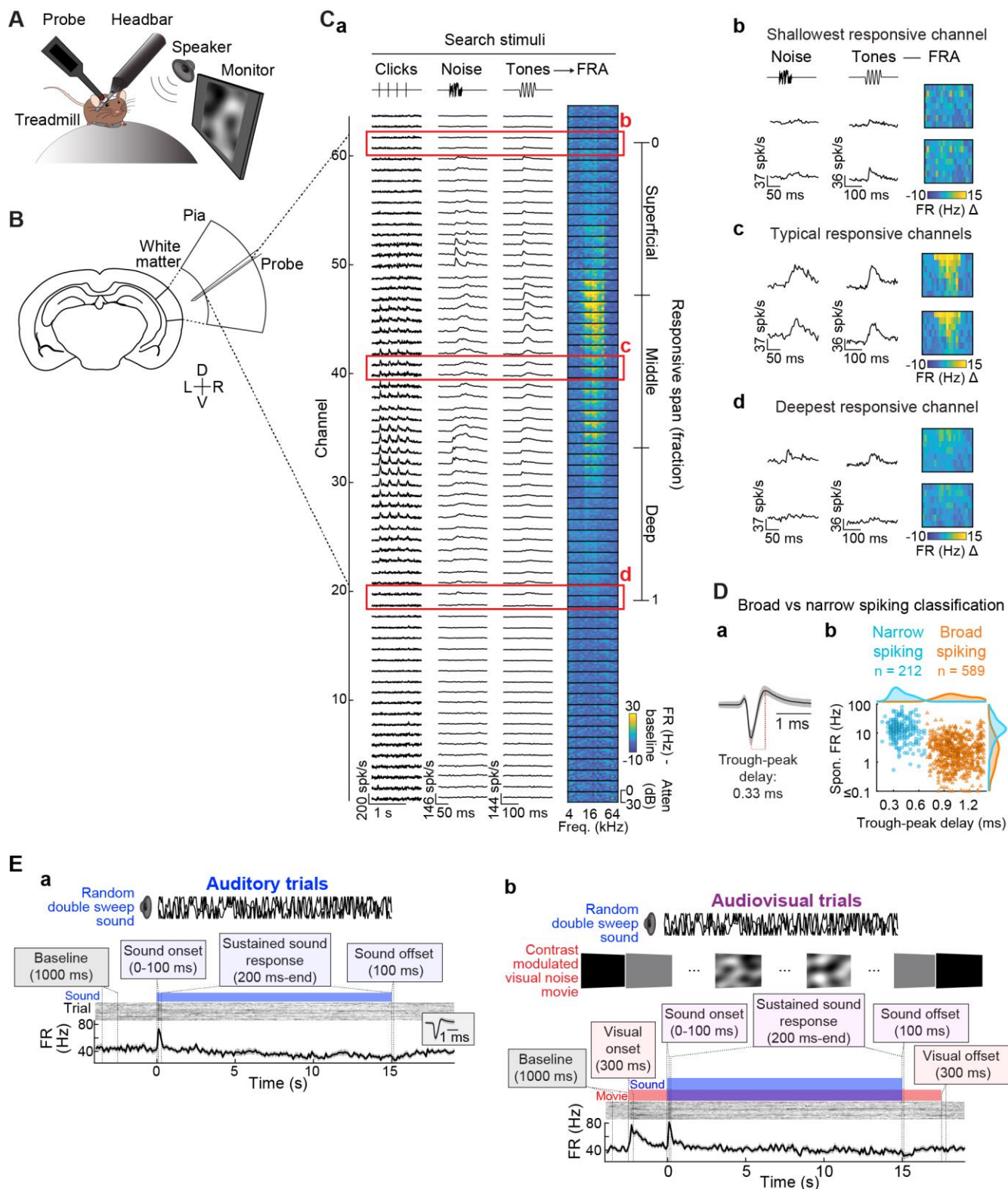147 calibrated to 25 cd/m$^2$ for 50% gray at eye position.

148    For the majority of recordings, search stimuli used for cortical depth estimation included
149  click trains, noise bursts, and pure tone pips, plus the experimental stimuli described below. For
150  a small minority of recordings, only tone pips and experimental stimuli were presented due to
151  time constraints. In some recordings, additional search stimuli were presented, such as
152  frequency-modulated sweeps. Click trains comprised broadband 5 ms non-ramped white noise
153  pulses presented at 4 Hz for 1 s at 60 dB with a ~1 s interstimulus interval (ISI), with 20–50
154  repetitions. Noise bursts consisted of 50 ms non-ramped band-passed noise with a uniform
155  spectral distribution between 4 and 64 kHz presented at 60 dB in 500 unique trials with a ~350
156  ms ISI. Pure tones consisted of 100 ms sinusoids with 5-ms cosine-squared onset/offset ramps
157  presented at a range of frequencies (4–64 kHz, 0.2 octave spacing) and attenuation levels (30–
158  60 dB, 5 dB steps). Three repetitions of each frequency-attenuation combination were
159  presented in pseudorandom order with an ISI of ~550 ms. Peristimulus-time histograms
160  (PSTHs) quantifying time-binned multi-unit firing rates were constructed for each stimulus. For
161  tone pips, frequency-response area (FRA) functions were constructed from baseline-subtracted
162  spike counts during the stimulus period averaged across trials at each frequency-attenuation
163  combination. PSTHs and FRAs from an example recording are shown in **Figure 1C**.
164    As depicted by **Figure 1E**, experiments comprised two trial types: (a) Auditory trials,
165  which presented sound only, and (b) Audiovisual trials, which included both sound and visual
166  stimuli. For both trial types, the auditory stimulus was a random double sweep (RDS), a
167  continuous, spectrally sparse receptive field estimation stimulus capable of effectively driving
168  activity across diversely tuned neurons in A1 (Gourévitch et al., 2015). The RDS comprised two
169  uncorrelated random sweeps, each varying continuously and smoothly over time between 4 and
170  64 kHz, with a maximum sweep modulation frequency of 20 Hz. Sample RDS frequency vectors
171  are depicted in **Figure 1E** and **Figure 4A, a**. The RDS was delivered in 15-s non-repeating
172  segments (40 trials, 10 minutes total stimulation; cf. Rutkowski et al., 2002). The inter-sound
173  interval was ~9 s, with visual stimuli trailing and leading sounds within this interval. Thus,
174  intertrial intervals were ~9 s for consecutive auditory trials, ~4 s for consecutive audiovisual
175  trials, and ~6.5 seconds for mixed trial type sequences. The same 40 RDS segments were used
176  for Auditory and Audiovisual trials to maintain identical stimulus statistics between conditions.
177  Audiovisual trials were thus identical to Auditory trials with the exception of an additional
178  contrast modulated visual noise stimulus.
179    As described in detail elsewhere (Niell and Stryker, 2008), CMN is a broadband stimulus
180  designed to drive as many primary visual cortical neurons as possible. The stimulus is
181  generated by first creating a random frequency spectrum in the Fourier domain. The temporal
182  frequency spectrum was flat with a low-pass cutoff at 10 Hz. The spatial frequency spectrum
183  dropped off as $A(f) \sim 1/(f+f_c)$, with $f_c$ = 0.05 cycles/°. A spatiotemporal movie was then created
184  by inverting the three-dimensional spectrum. Finally, contrast modulation was imposed by
185  multiplying the movie by a sinusoidally variable contrast function. The CMN stimulus was
186  generated at 60×60 pixels, then interpolated to 900×900 pixels. The first and last frames of the
187  CMN movie were uniform 50% gray, providing abrupt luminance changes at stimulus onset and
188  offset from black during the intertrial interval. The CMN stimulus led and trailed the RDS sounds
189  by 2.5 s to allow ample time for potential visual-evoked spiking responses to reach an adapted
190  state prior to sound onset responses and persist throughout sound offset responses.

191
192

**Figure 1.** Single unit recording and audiovisual stimulation in awake mouse auditory cortex. (**A**) Mice were head fixed atop a spherical treadmill. A headbar window provided access to primary auditory cortex (A1) of the right hemisphere for extracellular recording with translaminar probes. Sounds were presented to the contralateral ear through an electrostatic speaker and visual stimuli were presented via a monitor centered in front of subjects at 25 cm distance. (**B**) Coronal mouse brain section with magnification of A1 and linear multichannel electrode arrays (64-channels, 20μm spacing) used to simultaneously record neuronal activity across all cortical layers. (**C**) Auditory cortical depth estimation. (**a**) Multiunit responses evoked by search stimuli (e.g., click trains, noise bursts, pure tones) were used to guide visual demarcation of the span of responsive channels which served as an estimate of the putative cortical span and was used to assign a fractional depth value to each recorded neuron. Fractional responsive span was further divided into Superficial, Middle, and Deep bins. (**b–d**) Example multiunit responses from (b) the shallowest channel of the responsive span, (c) responsive channels from the middle of the probe, and (d) deepest channel of the responsive span. (**D**) Identification of putative excitatory and inhibitory neurons by waveform morphology clustering. (**a**) Example single-unit waveform (black line: median, gray shading: median absolute deviation) showing trough-peak delay calculation. (**b**) The distribution of trough-peak delay times was sharply bimodal, permitting straightforward identification of broad spiking (BS; putative excitatory) and narrow spiking (NS; putative inhibitory) neurons. BS and NS unit populations were further distinguished by differences in spontaneous firing rate. (**E**) Auditory and visual stimulation paradigm. (**a**) Auditory trials comprised non-repeating 15-s segments of a random double sweep (RDS) stimulus, comprising two continuously frequency-modulated pure tones which varied independently of one another between 4 and 64 kHz. Trials were separated by silent intertrial intervals (~4–9 s), permitting calculation of spontaneous firing rates. Sound onset firing rate responses were defined by a 100-ms window post-stimulus onset. Sustained firing rates were quantified within a window 200-ms post-stimulus onset to the end of the stimulus (15 s). Inset shows the example unit spike waveform (median ± MAD). (**b**) Audiovisual trials were identical to Auditory trials (including the same RDS segments) with the addition of a visual contrast modulated noise (CMN) stimulus. The CMN stimulus led and trailed the RDS stimulus by 2.5 s, permitting unambiguous assessment of visual onset and offset firing rate responses and allowing adaptation to the visual stimulus prior to sound onset. The auditory (RDS) and visual (CMN) stimuli were uncorrelated with each other. Auditory and Audiovisual trials were interleaved in pseudorandom order.

193 *Electrophysiology*

194

195    Recordings were conducted inside a sound attenuation chamber (Industrial Acoustics

196    Company). Anesthesia has well known and profound influences on auditory cortical encoding,

197    including the relative prevalence of onset and sustained firing rate responses (Wang et al.,

198    2005). Recordings were thus conducted in awake, headfixed animals moving freely atop a

199    spherical treadmill in **Figure 1A** (Dombeck et al., 2007; Niell and Stryker, 2010; Phillips and

200    Hasenstaub, 2016; Phillips et al., 2017a, 2017b; Morrill and Hasenstaub, 2018; Bigelow et al.,

201    2019). The silicone elastomer filling the craniotomy was removed and a single shank, linear

202    multichannel electrode array (Cambridge Neurotech) was slowly lowered into cortex using a

203    motorized microdrive (FHC). Arrays with 64 channels (20 μm site spacing, 1260 μm total span)

204    were used for all recordings except one, which used a 32-channel array (25 μm site spacing,

205    775 μm total span). Prior to lowering the probe, the craniotomy was filled with 2% agarose to

206    stabilize the brain surface. After reaching depths of approximately 800–1000 μm below the first

207    observation of action potentials, probes were allowed to settle for at least 20 mins before

208    initiating recording. Continuous extracellular voltage traces were collected using an RHD2000

209    (Intan Technologies) at a sample rate of 30 kHz. Other experimental events such as stimulus

210    event times were stored concurrently by the same system.

211        The span of recording channels (1260 μm) exceeded mouse cortical depth (~800 μm;
212    Paxinos and Franklin, 2019), resulting in a majority of sound-responsive channels plus an
213    additional subset of channels recorded outside of A1. As depicted in **Figure 1C**, multi-unit
214    responses evoked by search stimuli (e.g., click trains, noise bursts, tone pips) were used to
215    guide visual demarcation of the range of sound-responsive channels, which served as an
216    estimate of cortical span. Although penetrations were approximately perpendicular to the
217    cortical surface, it was impractical to achieve perfect orthogonality. Thus, the responsive span of
218    each recording was normalized such that each channel was expressed as a fraction of total
219    depth (Morrill and Hasenstaub, 2018). We note that this cortical depth estimation procedure is
220    less precise than our prior study, in which Di-I was applied to the probe for histologically
221    referenced depth estimation (Morrill and Hasenstaub, 2018). Nevertheless, we observed parallel
222    depth distributions of visual responsive neurons in the current and prior studies, suggesting the
223    current method achieved a rough approximation to the histological approach. However, due to
224    the lack of histological verification, a more conservative depth categorization approach was
225    adopted, dividing the responsive span into three equal bins reflecting superficial, middle, and
226    deep-layer neuron populations.

227        Recordings targeted A1 using stereotaxic coordinates and anatomical landmarks such
228    as characteristic vasculature patterns (Joachimsthaler et al., 2014). Previous studies have
229    reported significant differences in tone onset latencies between primary and non-primary
230    auditory cortical fields (Joachimsthaler et al., 2014), with latencies between 5 and 18 ms for
231    primary fields (median ~9 ms), and 8–32 ms for non-primary fields (median ~12–16 ms). Thus,
232    tone onset latencies were used to support designations of putative primary recording sites.
233    PSTHs were constructed from multi-unit activity (negative threshold crossings exceeding 4.5
234    median absolute deviations of the continuous voltage trace distribution) using 2-ms bins and
235    smoothed with a Savitzky-Golay filter (3rd order, 10-ms window). Onset latency was defined as
236    the first post-stimulus bin in which the firing rate exceeded 2.5 standard deviations of the pre-
237    stimulus firing rate bins (Morrill and Hasenstaub, 2018). Recordings for which the median
238    latency across the responsive channel span was 14 ms or less were considered putative
239    primary sites and retained for further analysis (49 of 60 total recordings). The retained putative
240    primary recordings universally exhibited robust multi-unit responses to click trains, noise bursts,
241    and tone pips, and clear evidence of frequently-level tuning in the FRA plots (**Figure 1C**) as well
242    as spectrotemporal tuning in single-unit responses to RDS stimuli (**Figure 4**).

243        Single-unit activity was isolated from continuous multichannel traces using Kilosort 2.0
244    (Pachitariu et al., 2016; available: https://github.com/MouseLand/Kilosort) and further validated
245    by auto- and cross-correlation analysis, refractory period analysis, and cluster isolation
246    statistics. Although the majority of isolated units were held throughout the entire recording (~28
247    minutes), isolation of individual neurons was occasionally disrupted and lost partway through
248    the experiment. Thus, the active timespan for each unit was estimated by visual demarcation of
249    unit activity plots over time. Inactive trials were discarded from further analysis. For the
250    remaining subset of active trials, RDS stimuli were matched between conditions by only using
251    available RDS segments common to both conditions. This ensured strict equivalence of auditory
252    stimuli between conditions, isolating any observed differences to the presence of the visual
253    CMN stimulus. Only units with 10 or more active trials (5 per condition) were retained for final
254    analysis. A total of 801 units were included in the analyses below. As in previous publications

255 (Phillips et al., 2017; Bigelow et al., 2019), units were classified as narrow-spiking (NS;
256 expected to be overwhelmingly inhibitory) or broad-spiking (BS; expected to be mainly
257 excitatory) on the basis of a clear bimodal distribution of waveform trough-peak delays (**Figure
258 1D**; NS, <600 μs, n = 212; BS, ≥600 μs, n = 589). Consistent with this classification, NS units
259 had characteristically higher spontaneous firing rates than BS units (estimated from baseline
260 period shown in **Figure 1E, a**; F = 391.31, p < $10^{-70}$, $\eta^2$ = 0.329).

262 *Spectrotemporal receptive field estimation*

264 Spectrotemporal fields (STRFs) were estimated using standard reverse-correlation techniques
265 (Wu et al., 2006) as depicted in **Figure 4A, a–c**. RDS stimuli were discretized in 1/8 oct
266 frequency bins and 5-ms time bins, which is sufficient resolution for modeling response
267 properties in the majority of A1 neurons (Thorson et al., 2015). The spike-triggered average
268 (STA) was obtained by adding the discretized stimulus segment preceding each spike to a
269 cumulative total, and then dividing by the total spike count. For all data analyses, the peri-spike
270 time analysis window spanned 0–100 ms prior to spike event times, sufficient for capturing the
271 full latency-adjusted temporal response periods of the majority of A1 neurons (Atencio and
272 Schreiner 2013, See et al., 2018). A broader window was used for display purposes, spanning
273 200 ms before and 50 ms after spike event times. The 50 ms post-spike window was included
274 as a visualization of estimated acausal values, i.e., those that would be expected by chance
275 given the finite recording time, stimulus and spike timing statistics, and smoothing parameters
276 (Gourévitch et al. 2015). The first 200 ms of the RDS response from each trial were dropped
277 from all STA calculations analyses to minimize bias reflecting strong onset transients.
278 The STA thus reflects the average binned stimulus segment preceding spike events and
279 can be viewed as a linear approximation to the optimal stimulus for driving neuronal firing
280 (deCharms and Merzenich, 1998). As discussed by Rutkowski et al. (2002), the STA can be
281 formalized as the probability (*P*) of a stimulus frequency *f* occurring at time $t_i$-*τ* given that a spike
282 occurred, as expressed by the equation

284 $$P(S[f,\tau]|i) = \left(\sum_i^n S(f, t_i - \tau) \bullet \delta(t_i)\right)/n \qquad \text{Eq. (1)}$$

286 where *i* indicates a spike, $t_i$ is a spiketime, *τ* is the time analysis window, *n* is the spike count,
287 and Σ indicates summing across spikes. *S(f,t)* is the stimulus value at a given time-frequency
288 bin, equaling one if an RDS frequency intersects the bin, two if both RDS frequencies coincide
289 with the bin, and zero otherwise. $S(f, t_i$-*τ*) represents the windowed stimulus aligned to a spike
290 time. δ($t_i$) is equal to one if a spike occurs at time $t_i$ and zero otherwise. With a slight
291 modification, STRF time-frequency bins can be expressed in terms of deviation from mean
292 driven firing rate (spikes/s - mean) using the terms defining the STA and Bayes' theorem,

294 $$P(i|S[f,\tau]) = P(S[f,\tau]|i) \bullet P(i)/P(S[f]) \qquad \text{Eq. (2)}$$

296 where *P(i)* is the probability of a spike occurring in a bin, equal to $n_i/T$, where *T* is the total
297 stimulus time and *P(S[f])* is the probability of a tone frequency occurring in a bin. The mean
298 driven firing rate is then subtracted from the STRF such that individual time-frequency bins

299  reflect increases (positive, red) or decreases (negative, blue) from the mean driven rate (**Figure**
300  **4A, c**). Finally, the STRF is smoothed by a uniform 3✕3 bin window to reduce overfitting to
301  finite-sampled stimulus statistics. The STA and STRF expressed in units of spikes/s are
302  multiples of each other since terms in the expression $P(i)/P(S[f])$ are constant and thus nearly
303  identical for the purposes of all data analyses, including comparisons between conditions.
304  However, we opt to report STRFs represented in firing-rate change units to facilitate
305  interpretation of stimulus driven changes in neuronal activity.
306       In addition, 'Null' STRFs were calculated using identical procedures to those described
307  above except that the stimulus was reversed in time, while preserving the original spike event
308  times. This modification breaks the temporal relationship between the stimulus and spike times,
309  but preserves spike count and timing statistics (e.g., interspike interval distribution), as well as
310  the statistical distributions of the stimulus (Bigelow and Malone, 2017). The resulting STRF was
311  used to estimate time-frequency bin values expected by chance within the constraints of finite
312  spike counts and stimulus time.
313       STRFs were calculated independently for each condition. As in previous studies (Fritz et
314  al., 2003), a difference STRF, which we refer to throughout as ΔSTRF, was calculated by
315  subtracting the STRF obtained from Auditory trials from the STRF obtained from Audiovisual
316  trials. A Null ΔSTRF was similarly obtained by subtracting Null Auditory STRF from the Null
317  Audiovisual STRF.
318
319  *Mutual information analysis*
320
321  Although the STRF effectively captures the spectrotemporal tuning of a neuron, it does not
322  reveal the degree to which driven spiking activity is dependent upon similarity between the
323  stimulus and receptive field. For instance, it is not apparent how consistently spiking activity is
324  observed when the stimulus closely approximates the STRF (e.g., RDS frequency vectors
325  intersecting STRF excitatory subfields) and whether spiking is inhibited when the stimulus is
326  anticorrelated or uncorrelated with the STRF (e.g., RDS frequency vectors intersecting STRF
327  inhibitory subfields or regions with values near zero). Thus, mutual information was calculated to
328  quantify the relationship between probability of firing and similarity between the stimulus and
329  STRF according to previously published methodology (Atencio et al., 2008; Atencio and
330  Schreiner, 2016).
331       As depicted by **Figure 12A**, mutual information reflects a scaled ratio of two
332  distributions: $P(x)$, which reflects the 'similarity' between the STRF and all possible RDS
333  segments in the experiment, and the other, $P(x|spike)$, reflecting STRF-stimulus 'similarity'
334  values at time bins in which a spike occurred. Stimulus-STRF 'similarity' (x) is operationally
335  defined as the inner product (i.e., projection value) between the STRF and the stimulus segment
336  of equivalent dimensions. Put another way, the STRF is convolved with the stimulus, yielding a
337  projection value for each 5-ms time bin in the stimulus (**Figure 12A, a–c**). For ease of
338  interpretation, raw projection values were standardized by subtracting the mean of the raw $P(x)$
339  distribution and dividing by its standard deviation. Thus, highly positive and negative
340  standardized values reflect stimulus segments highly similar and dissimilar to the STRF,
341  respectively. Values near zero imply a random relationship between the stimulus and STRF.
342  The continuously valued projection values were separated into nine linearly spaced bins (**Figure**

343    **12A, d**). As in prior studies, the most extreme positive bin was dropped from the information
344    calculation due to undersampling (the most extreme STRF-stimulus matches are rare), resulting
345    in eight analyzed bins. Finally, mutual information (bits/spike) is calculated by the calculating
346    log2-transformed ratio of $P(x|spike)$ divided by $P(x)$, multiplying the result by $P(x|spike)$, and
347    summing across bins. In equation form,

349    $$I = \int dx P(x|spike) log2\left[\frac{P(x|spike)}{P(x)}\right]$$    Eq. (3)

351    The information estimation procedure occasionally produced values of zero for a small
352    minority of units, which were excluded from further information analyses. Information values
353    were not calculated for units with <200 total spikes due to undersampling concerns. Intuitively,
354    mutual information increases as the $P(x)$ and $P(x|spike)$ distributions diverge. For instance, if a
355    given binned $P(x)$ value is low but $P(x|spike)$ is high, we conclude that the stimulus (and its
356    associated 'similarity' value to the STRF) occurs infrequently but usually evokes a spiking
357    response. Similarly, if $P(x)$ is high but $P(x|spike)$ is low, we infer that a regularly occurring
358    stimulus typically inhibits spiking. In each scenario, the stimulus and response are thus mutually
359    informative of one another. By contrast, if the binned $P(x)$ and $P(x|spike)$ values are similar
360    (both high or low), the stimulus is not informative about spiking behavior.

362    *Statistical analysis*

364    Our stimulus design allowed measurement of unimodal and bimodal evoked firing-rate
365    responses as well as stimulus-driven changes in spike timing (e.g., spike alignment with RDS
366    stimulus features). As indicated by **Figure 1E, a**, baseline (spontaneous) firing rates were
367    measured from a window spanning 3.5 to 2.5 s prior to sound onset (immediately preceding
368    visual stimulus onset for audiovisual trials). Sound onset responses were measured from spikes
369    occurring within the first 100 ms of the stimulus. Sustained firing rates were estimated within the
370    same window used for STRF estimation, from 200 ms post stimulus onset to the end of the
371    stimulus (15 s). All sound-evoked firing rate response windows and STRF calculations were
372    identical between conditions. Visual onset responses were analyzed in a window spanning 300
373    ms after the visual stimulus began. The wider window used to capture onset transients for visual
374    compared to auditory stimuli accommodated the substantially longer response latencies typical
375    of visual responses in A1 (Morrill and Hasenstaub, 2018).
376    Auditory and visual offset responses were similarly analyzed using 100- and 300-ms
377    windows following stimulus offset, respectively. Unlike onset and sustained responses, offset
378    responses required comparison against two baselines: first, the spontaneous window described
379    above, and second, a window of equivalent duration (1 s) preceding the offset response (i.e.,
380    the final second of the stimulus). Differences from the second baseline ensured offset
381    responses were not merely carryover from a sustained response. In addition, 'significant' offset
382    responses were required to be above or below both baselines (but not between them) to ensure
383    offset responses did not simply reflect the return of a sustained firing rate change to
384    spontaneous activity. The larger of the two p-values was used to assess significance of the
385    offset response, and its associated baseline was used for calculating effect size described
386    below.

387     For each unit, the significance of all firing rate responses was assessed by comparison
388     to baseline Wilcoxon signed-rank tests (paired) using $\alpha$ = 0.05. For all tests of individual unit
389     significance, false discovery rate (FDR) was limited by implementing the Benjamini–Hochberg
390     procedure with $q$ = 0.05 across the unit population (Benjamini and Hochberg, 1995). The
391     adjusted p-values produced by the FDR procedure are used to indicate significance for each
392     response type for all example unit plots throughout the manuscript. For a standardized measure
393     of response strength that accommodated both increases and decreases in firing rate from
394     baseline, we defined effect sizes as the absolute difference between the evoked and
395     spontaneous firing rate means, divided by the standard deviation of the spontaneous firing rate.
396     This facilitated interpretation of response deviations from chance reflecting different analysis
397     windows, unit types.
398     Significant differences in firing rate between conditions (auditory, audiovisual) were
399     similarly assessed with Wilcoxon signed-rank tests ($\alpha$ = 0.05, Benjamini–Hochberg FDR
400     correction with $q$ = 0.05), and effect sizes were estimated as the absolute difference between
401     conditions divided by the standard deviation of the auditory condition.
402     To assess the significance of STRFs, we used a reliability index in which the correlation
403     coefficient was calculated between two STRFs computed from random trial halves (**Figure 4A,**
404     **c**; Escabí et al., 2014). Reliability was defined as the mean across 1000 iterations for both 'data'
405     (time-preserved stimulus) and 'null' (time-reversed stimulus) STRFs. A p-value was calculated
406     reflecting the proportion of the null distribution exceeding the mean of the data distribution
407     (**Figure 4A, e**). P-values equal to zero (cases where none of the null correlations exceeded the
408     data mean) were adjusted to 0.000999 in reflection of the resolution permitted by the number of
409     subsample iterations. Finally, the p-values were multiplied by two for an estimate of two-tailed
410     significance and adjusted using Benjamini–Hochberg FDR correction ($q$ = 0.05). Because the
411     null subsampled STRF distribution skewed negative for some units (i.e., null reliability index >
412     0), we further required STRF reliability >0.2 in order to be considered 'significant'. Similar to
413     firing rate responses, STRF reliability effect sizes were estimated as the absolute difference
414     between the data and null distributions, divided by the standard deviation of the null distribution.
415     Significance of ΔSTRFs was assessed with the same approach, except using subsampled
416     ΔSTRF correlation distributions (**Figure 9A, a–b**).
417     Except where otherwise noted, tests of population-level differences (e.g., among units in
418     superficial, middle, deep cortical depth bins) were assessed by independent one-way analysis
419     of variance (ANOVA), using effect sizes described above as the dependent variable. For
420     uniformity in presenting the results, we use the same approach for testing between two
421     variables (e.g., differences between NS and BS units), wherein ANOVA and the Student's $t$-test
422     produce equivalent p-values with $F = t^2$.

**Results**

426     The current study examined visual-modulation of sound-evoked responses in awake mouse A1.
427     By delivering a continuous auditory receptive field estimation stimulus in 15-s segments
428     separated by silent intervals, we were able to capture both transient onset firing rate responses
429     as well as sustained firing rate responses throughout the duration of the stimulus. By using
430     reverse correlation methods, we were further able to estimate STRFs, which are simultaneously

431    sensitive to both spike rate and timing. Half of the trials included a continuous visual stimulus,
432    enabling direct comparison of firing rate responses and STRFs between auditory alone and
433    visual-modulated conditions. The visual stimulus both led and trailed the sound stimulus by 2.5
434    seconds, further allowing measurement of averaged firing rate responses evoked by the visual
435    stimulus alone.
436
437    *Some A1 neurons respond to visual stimulation alone*
438
439    Example visual-responsive units are shown in **Figure 2A–B**. Consistent with our earlier study
440    examining visual-flash evoked responses in A1 (Morrill and Hasenstaub, 2018), we found that
441    significant responses to the visual CMN stimulus were most prevalent in the deepest cortical
442    depth bin (**Figure 2C**). Extending our previous work, we found that for a minority of units, visual
443    responses reflected decreases in firing rate relative to baseline (e.g., **Figure 2B**). For a
444    standardized measure of visual response strength that avoided overweighting units with high
445    firing rates and accommodated both increases and decreases in firing rate from baseline, we
446    defined effect sizes as the absolute difference between the evoked and spontaneous firing rate
447    means, divided by the standard deviation of the spontaneous firing rate. One-way ANOVA
448    confirmed that effect size of visual-evoked onset firing rate changes was significantly dependent
449    upon cortical depth for both unit types, with the strongest responses in the deepest bin (NS: F =
450    13.95, p < $10^{-5}$, $\eta^2$ = 0.118; BS: F = 7.62, p < $10^{-3}$, $\eta^2$ = 0.025).
451        Visual responses were detected in a larger percentage of NS units (15.6%) than BS
452    units (7.8%), and visual responses were significantly stronger for NS units in the deepest
453    cortical depth bin (F = 4.32, p = 0.039, $\eta^2$ = 0.016). Differences were non-significant for the
454    remaining depth bins (shallow: F = 0.25, p = 0.619, $\eta^2$ = 0.002; middle: F = 0.29, p = 0.593, $\eta^2$ =
455    0.001). We conducted follow-up analyses to test whether the larger percentage of visual
456    responsive NS units could be explained by their characteristically higher firing rates – and thus
457    statistical power for detecting firing rate changes. Mean firing rates were calculated for each unit
458    across audiovisual trials (including the full stimulus period plus spontaneous activity 1 s before
459    and after the stimulus), producing population mean rates of 18.11 Hz for NS units and 4.45 Hz
460    for BS units (**Extended Data Figure 2-1A, a**). By randomly subsampling spikes from NS units
461    with firing rates above the BS mean, we created a pseudo-population of NS units with mean
462    firing rate equivalent to BS units (4.45 Hz; **Extended Data Figure 2-1A, b**). Visual onset
463    response data were recalculated, including the p-value distribution which was readjusted by the
464    Benjamini–Hochberg FDR procedure. As seen in **Extended Data Figure 2-1B**, differences in
465    visual response effect sizes were non-significant for all depth bins (all F-ratios < 1.4, all p-values
466    > 0.24). Thus, visual responses were observed in both NS and BS units, but any differences
467    between these unit subpopulations may have reflected inherent differences in firing rate.
468        Because stimulus offset responses have been observed in primary visual cortex (Liang
469    et al., 2008), we also examined the possibility of significant firing rate changes following
470    termination of the visual CMN stimulus in the current study. However, after FDR correction, we
471    found that there were no units with significant visual offset responses. We therefore relied on
472    visual onset responses throughout the remainder of the analyses to define visual responsive
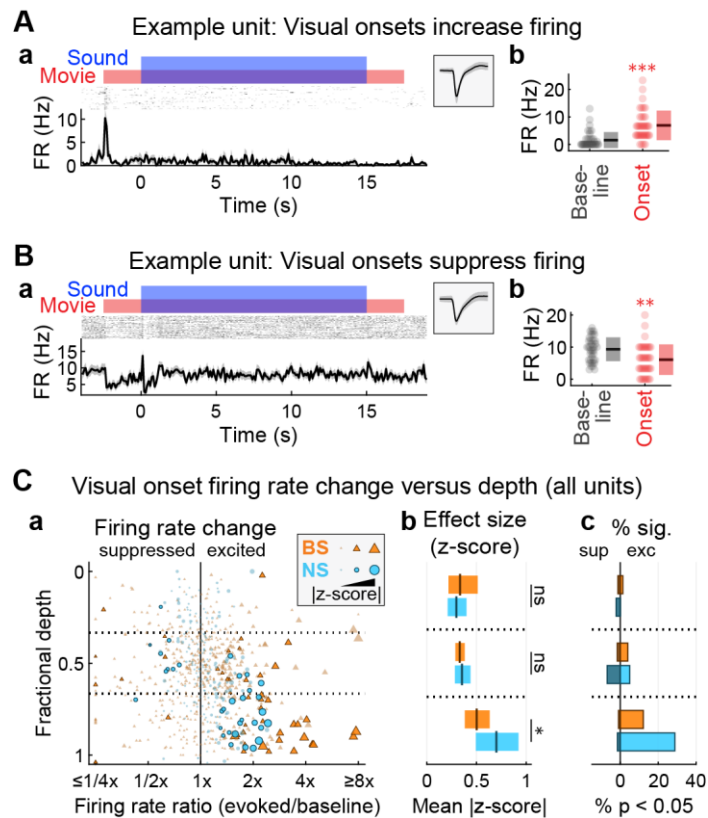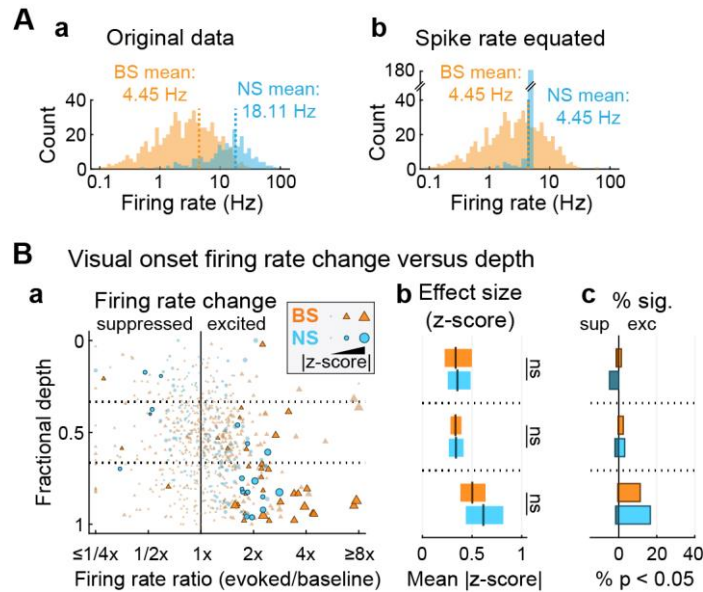473    units.

**Figure 2.** Some primary auditory cortical neurons respond to visual stimulation alone. (**A**) Example unit with an excitatory visual onset response. (**a**) Spiking responses quantified by peristimulus time histograms (lower) and binned-spike count matrices (upper) with blue and red bars indicating auditory and visual stimulus intervals, respectively (temporal binning: 100 ms). Inset shows the unit spike waveform (median ± MAD). (**b**) Summary of visual onset firing rate responses compared to baseline. Each dot represents mean firing rate for a single trial, with mean ± SD across trials indicated to the right. Wilcoxon signed-rank tests (paired): $*p<0.05$, $**p<0.01$, $***p<0.001$, ns $p>0.05$. (**B**) Example unit with suppressed visual onset response. (**C**) Summary of visual onset firing rate responses by unit type and cortical depth. (**a**) Scatter plot depicting visual onset responses for each unit at its estimated fractional depth. Firing rate ratio values above and below 1 (x-axis) indicate excited and suppressed responses, respectively. Outlined markers indicate statistically significant responses ($p < 0.05$, Benjamini–Hochberg FDR correction). Marker sizes are scaled by effect size (absolute difference between onset and baseline means, divided by baseline SD). (**b**) Mean effect size (plus 99% confidence interval) across all recorded units (significant and non-significant visual onset responses included) by unit type and depth. Visual onset effect sizes for NS units were significantly stronger than BS units in the deepest cortical depth bin (but see Extended Data Figure 2-1 for spike-equated analysis). One-way ANOVA: $*p<0.05$, $**p<0.01$, $***p<0.001$, ns $p>0.05$. (**c**) Histograms indicating percentages of all recorded units with significant excited and suppressed responses by unit type and depth. Significant visual responses were most concentrated in the deepest cortical depth bin.
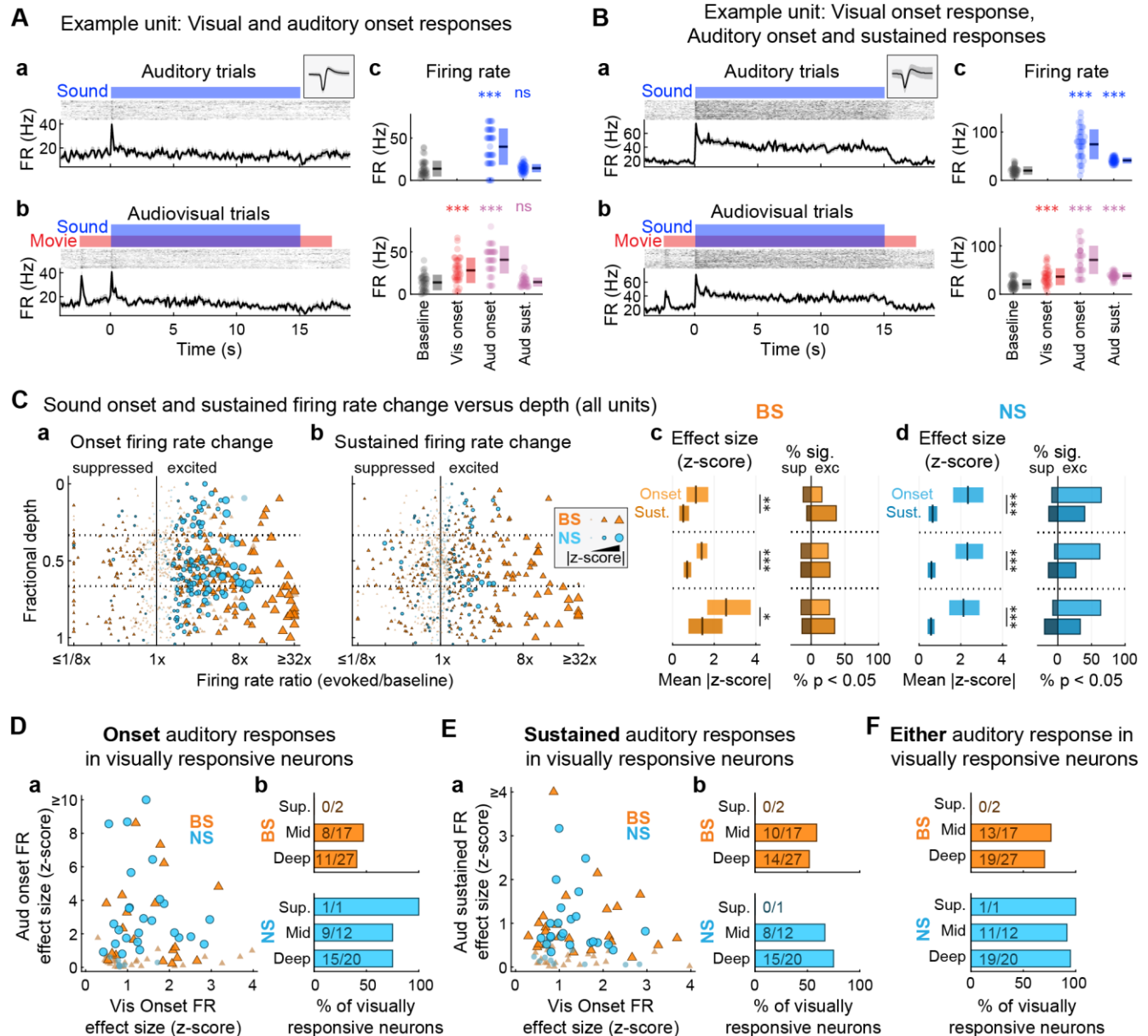
474

**Extended Data Figure 2-1.** Differences in visual responses between NS and BS units may reflect inherent differences in firing rate. (**A**) Firing rate histograms for each unit type inclusive of both spontaneous and stimulated periods. (**a**) NS units tend to have higher mean firing rates than BS units. (**b**) By subsampling spikes from NS units with firing rates above the mean, a pseudopopulation was created with equivalent mean firing rate to BS units. (**B**) Summary of visual onset firing rate responses by cortical depth for the spike-equated NS and BS unit populations. (**a**) Scatter plot depicting visual onset responses for each unit at its estimated fractional depth. (**b**) Mean effect size (plus 99% confidence interval) across all recorded units (significant and non-significant visual onset responses included) by unit type and depth. No significant differences were observed between BS and spike-equated NS units. One-way ANOVA: ns p>0.05. (**c**) Histograms indicating percentages of units with significant excited and suppressed responses by unit type and depth.

*Visual-responsive neurons in A1 also typically respond to sound*

The visual-responsive example unit shown in **Figure 2B** also had an apparent increase in firing rate aligned to sound onset, whereas the example in **Figure 2A** lacked any discernible change in firing rate during sound presentation. To quantify the proportions of unimodal and bimodal visual-responsive units, we examined intersections of significant firing rate changes evoked by unimodal visual and auditory stimuli. Because numerous prior studies reported differences between auditory transient spike rate changes at stimulus onset and sustained responses throughout the remainder of the stimulus, we separately analyzed firing rate responses averaged within the first 100 ms of the stimulus (onset) and from 200-ms to stimulus end (sustained). Example visual-responsive neurons with significant auditory onset and sustained responses are shown in **Figure 3A and B**, respectively. Consistent with previous studies, we found that onset responses were substantially stronger than sustained responses (**Figure 3C**) for both unit types (BS units: shallow: F = 8.25, p = 0.005, $\eta^2$ = 0.053; middle: F = 32.68, p < 10$^{-7}$, $\eta^2$ = 0.051; deep: F = 5.13, p = 0.025, $\eta^2$ = 0.012; NS units: shallow: F = 33.15, p < 10$^{-6}$, $\eta^2$ = 0.274; middle: F = 40.98, p < 10$^{-8}$, $\eta^2$ = 0.170; deep: F = 26.17, p < 10$^{-5}$, $\eta^2$ = 0.168). As shown in **Figure 3D–E**, roughly half of all visual responsive units had either significant onset or sustained firing rate responses to sound. Considering onset and sustained responses together, well over two-thirds of visual-responsive neurons exhibited significant sound-evoked firing rate changes (**Figure 3F**).

495

**Figure 3**. Visually responsive neurons in A1 also typically exhibit sound-evoked firing rate responses. (**A**) Example unit with visual and auditory onset evoked firing rate responses. (**a**) Auditory and (**b**) Audiovisual trials. Spiking responses quantified by peristimulus time histograms (lower) and binned-spike count matrices (upper) with red and blue bars indicating auditory and visual stimulus intervals, respectively (temporal binning: 100 ms). Inset in (a) shows the unit spike waveform (median ± MAD). (**c**) Window-averaged firing rate responses, with dots representing single trials and mean ± SD across trials indicated to the right. Wilcoxon signed-rank tests (paired): *p<0.05, **p<0.01, ***p<0.001, ns p>0.05. (**B**) Example unit with visual onset response as well as both auditory onset and sustained responses, with subplot organization as in (A). (**C**) Sound onset firing rate responses are stronger than sustained responses. (**a, b**) Summary of sound onset and sustained firing rate changes from baseline (including both visual responsive and non-responsive units) separated by unit type and cortical depth. (**c**) Comparison of onset and sustained responses for BS units. Left: mean effect size (plus 99% confidence interval) across all recorded units (significant and non-significant responses included). Effect sizes were significantly greater for onset responses (upper bars, lighter coloring) than sustained responses at each depth bin. One-way ANOVA: *p<0.05, **p<0.01, ***p<0.001, ns p>0.05. Right: Percentages of all recorded units with significant onset and sustained responses by unit type and depth. (**d**) Comparison of onset and sustained responses for NS units, with subplot organization as in (c). (**D**) Approximately half of visually responsive units have significant sound onset firing rate responses. (**a**) Scatter plot of visual and sound onset response effect sizes. Large markers with outlines reflect units with significant sound onset firing rate responses. (**b**) Bar plot showing the percentages of visually responsive units with significant sound onset firing rate responses. (**E**) Approximately half of visually responsive units have significant sound sustained firing rate responses. Subplot organization as in (D). (**F**) The majority of visually responsive units have significant sound-evoked firing rate responses including either onset or sustained responses. Bars represent the union of sound-responsive units E, b and D, b. Effect sizes were similarly greater for onset responses across depth bins.

496       In addition to onset and sustained responses, we observed transient firing rate changes
497 at the offset of RDS segments in only a small minority of units (NS units: 15/212 [7.1%]; BS
498 units: 27/589 [4.6%]). We thus focused on onset and sustained responses to define sound-
499 evoked firing rate changes in the current report. Indeed, as in previous studies (Scholl et al.,
500 2010), we found that the majority of these units also had significant onset and/or sustained
501 responses (NS units: 12/15 [80.0%]; BS units: 14/27 [51.9%]).
502       The distinction between onset and sustained responses observed in the present study
503 reinforces the conclusions of numerous previous studies suggesting spike rate changes driven
504 by different stimulus phases (e.g., onset, offset, sustained) are at least partially dissociable and
505 may contain different information about sounds (Lu et al., 2001; Wang et al., 2005; Malone et
506 al., 2015). We pursued a related question as to whether sound-evoked responses reflecting
507 spike rate and timing changes were similarly partially independent by examining intersections of
508 sound-evoked firing rate responses (onset and sustained) and STRFs, which are sensitive to
509 both spike rate and timing changes. As depicted by **Figure 4A, a–c**, STRFs are calculated by
510 averaging the windowed stimulus segments preceding each spike. Thus, spectrotemporal
511 tuning depends strictly upon temporal alignment between spike events and stimulus features.
512 Importantly, such alignment may occur with or without changes from the spontaneous spike
513 rate. We used a response reliability metric to determine whether observed STRF structure was
514 statistically different from chance, as defined by STRFs calculated using time-reversed RDS
515 segments (**Figure 4A, d–e**; Escabí et al., 2014).
516       An example unit with significant STRF reliability, which also had clear firing rate changes
517 from baseline, is shown in **Figure 4B**. By contrast, **Figure 4C** shows an example unit with
518 significant STRF reliability but neither significant onset nor sustained firing rate changes.
519 Significantly reliable STRFs were observed in the majority of both NS and BS units (**Figure 4D**).
520 Notably, significant changes in firing rate were not observed in many of the units with significant
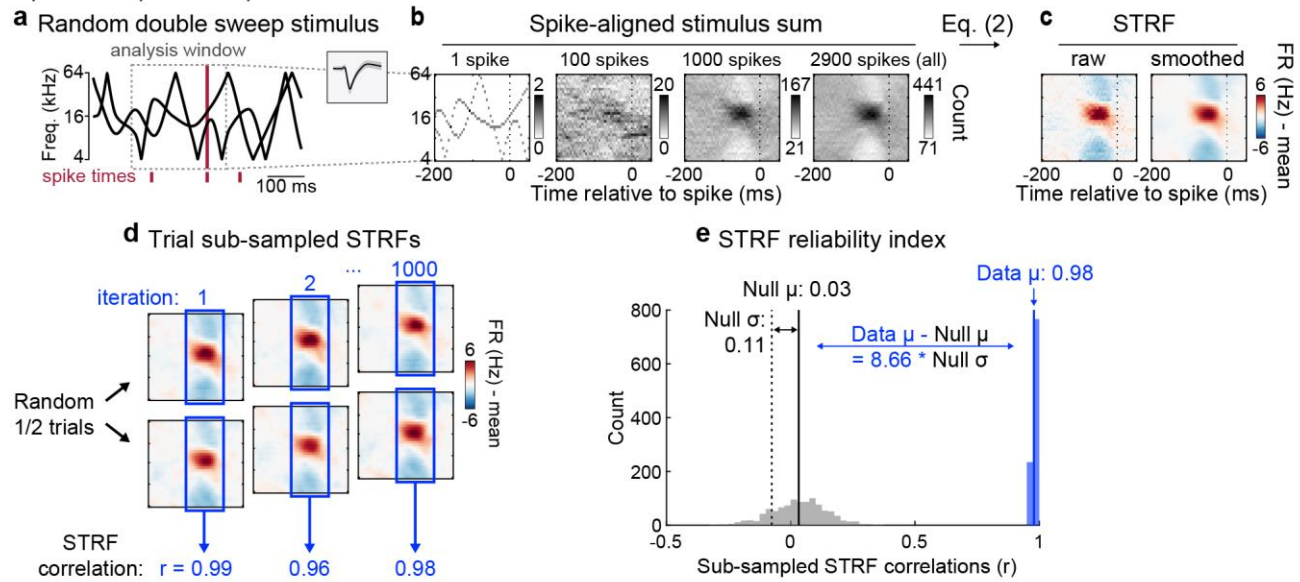
521    STRFs (**Figure 4E**). For BS units, significant STRFs without onset responses were slightly more
522    common than both STRF and onset responses together, whereas the reverse was true for NS
523    units (**Figure 4E, a**). For both unit types, sustained responses were observed in roughly half of
524    units with significant STRFs (**Figure 4E, b**). Significant firing rate changes (onset or sustained)
525    without STRFs were rare. Because sustained responses and STRFs were calculated from the
526    same analysis windows, spikes, and stimulus distributions, our results underscore the important
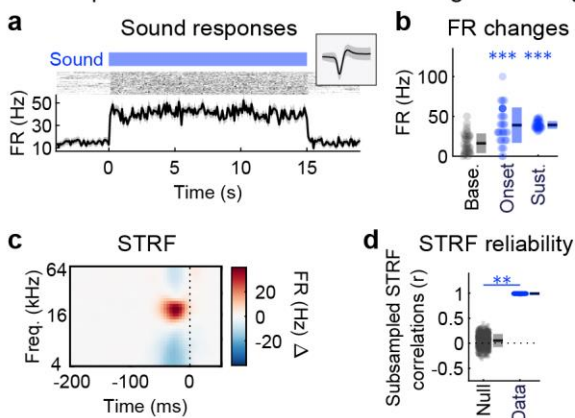527    distinction between changes in spike rate and timing in response to sound.
528            We quantified how many visual-responsive neurons were also responsive to sound
529    features as revealed by STRF calculation, e.g., as shown by the visual-responsive example unit
530    with significant spectrotemporal tuning in **Figure 5A**. Group data indicated that visual-
531    responsive neurons with significant STRFs were the rule rather than the exception, with
532    significant STRFs observed in all but a handful of BS units (**Figure 5B**). Extending the definition
533    of 'sound responsive' to include STRFs as well as onset or sustained firing rate responses
534    indicated that, with just a few exceptions, visual-responsive neurons in A1 also respond to
535    sound (**Figure 5C**).

**A** Spectrotemporal receptive field estimation

**a** Random double sweep stimulus

**b** Spike-aligned stimulus sum — Eq. (2) →

**c** STRF

**d** Trial sub-sampled STRFs

**e** STRF reliability index

**B** Example unit: STRF with sustained firing rate change

**a** Sound responses **b** FR changes

**c** STRF **d** STRF reliability

**C** Example unit: STRF without sustained firing rate change

**a** Sound responses **b** FR changes

**c** STRF **d** STRF reliability

**D** STRF reliability versus depth (all units)

**a** Reliability index **b** Effect size (z-score) **c** % sig.

**E** Overlap between significant STRFs and FR responses

**a** **Onset** responses **b** **Sustained** responses



536

19

**Figure 4.** Auditory cortical neurons may encode spectrotemporal features with or without significant spike rate changes. (**A**) Spectrotemoral receptive field (STRF) estimation procedure. (**a**) For each spike, a segment of the RDS stimulus was stored using a window 200 ms before and 50 ms after the spike time. (**b**) RDS segments aligned to each spike were added to a cumulative sum. Structure in the time-frequency bins typically only emerges after several hundred spikes or more. (**c**) Transforming the spike-aligned stimulus sum according to Eq. (2) yielded an STRF estimate expressed in firing rate (Hz) deviations from the mean driven rate. Red and blue regions indicate stimulus energy at time-frequency bins associated with increases and decreases in firing rate, respectively. (**d**) A subsampling procedure was used to determine the statistical significance of time-frequency bin structure in the STRFs. The correlation coefficient between STRFs calculated from random trial halves (without replacement) was calculated across 1000 iterations. (**e**) The STRF reliability index was defined as the mean of the subsampled correlation coefficient distribution. A null distribution was obtained from STRFs calculated using time-reversed stimulus RDS segments, which breaks the temporal relationship between spikes and stimulus features but preserves spike count and timing statistics. A p-value was obtained by dividing the number of null STRF correlations exceeding the reliability index (data) by the number of iterations and multiplying by two for two-tailed significance. Effect size reflected the absolute difference between null and data means, divided by the null standard deviation. (**B**) Example unit with significant STRF reliability as well as onset and sustained firing rate responses. (**a**) Spiking responses quantified by peristimulus time histograms (lower) and binned-spike count matrices (upper) with blue bar indicating auditory stimulus interval (temporal binning: 100ms). Inset shows the unit spike waveform (median ± MAD). (**b**) Summary of sound onset and sustained firing rate responses compared to baseline. Each dot represents a single trial, with mean ± SD across trials indicated to the right. Wilcoxon signed-rank tests (paired): *p<0.05, **p<0.01, ***p<0.001, ns p>0.05. (**c**) STRF as calculated in (A, a–c). (**d**) STRF reliability as calculated in (A, d–e). Each dot represents the correlation between STRFs for a single subsample iteration, with mean ± SD across trials indicated to the right. Subsampling test: *p<0.05, **p<0.01, ***p<0.001, ns p>0.05. (**C**) Example unit significant STRF reliability but non-significant onset and sustained firing rate responses. Subplot organization as in (B). (**D**) Summary of STRF reliability by unit type and cortical depth. (**a**) STRF reliability for each unit at its estimated cortical depth. Marker sizes are scaled by effect size, with outlined markers indicating units with significant STRF reliability (p < 0.05, Benjamini–Hochberg FDR correction). (**b**) Mean effect size (plus 99% confidence interval) across all recorded units (units with significant and non-significant reliability included) by unit type and depth. A small difference between unit types was observed in the deepest cortical bin. One-way ANOVA: *p<0.05, **p<0.01, ***p<0.001, ns p>0.05. (**c**) Histograms indicating percentages of all recorded units with significant STRF reliability. (**E**) Significant STRF reliability may occur with or without significant firing rate changes. (**a**) Intersections of significant onset firing rate responses alone, significant STRF reliability alone, or both. (**b**) Intersections of significant sustained firing rate responses alone, significant STRF reliability alone, or both.
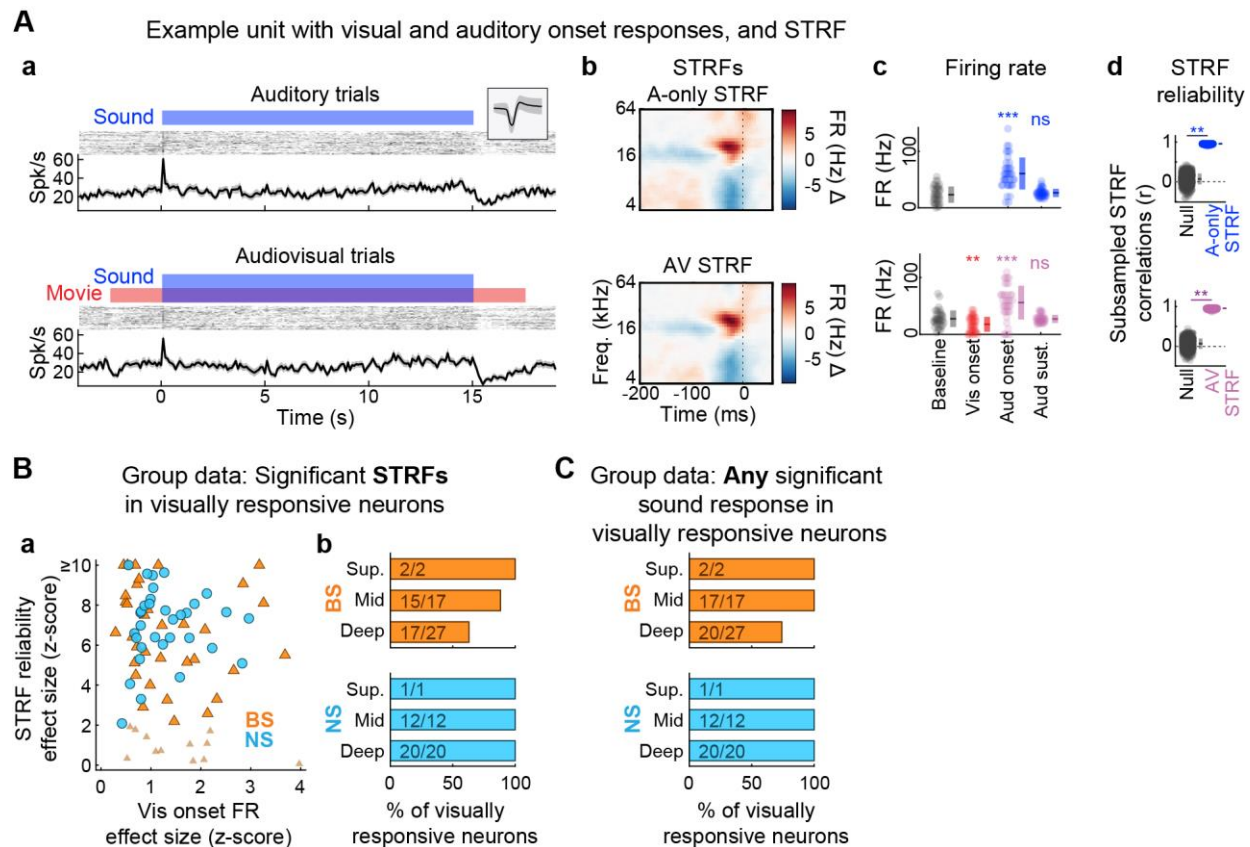
537

**Figure 5**. Visually responsive A1 neurons are typically tuned for spectrotemporal features. (**A**) Example unit with both visual and auditory onset firing rate responses as well as significant STRF reliability. (**a**) Spiking responses quantified by peristimulus time histograms (lower) and binned-spike count matrices (upper) with blue and red bars indicating auditory and visual stimulus intervals, respectively (temporal binning: 100ms). Inset shows the unit spike waveform (median ± MAD). (**b**) STRFs for each condition. (**c**) Summary of firing rate responses compared to baseline. Each dot represents a single trial, with mean ± SD across trials indicated to the right. Wilcoxon signed-rank tests (paired): *p<0.05, **p<0.01, ***p<0.001, ns p>0.05. (**d**) STRF reliability for each condition. Each dot represents the correlation between STRFs for a single subsample iteration, with mean ± SD across trials indicated to the right. Subsampling test: *p<0.05, **p<0.01, ***p<0.001, ns p>0.05. (**B**) The majority of visually responsive units have significant STRF reliability. (**a**) Scatter plot showing effect size for visual onset and STRF reliability. The larger, outlined markers indicate units with significant STRF reliability. (**b**) Bar plots showing percentages of visually responsive units with significant STRF reliability. (**C**) The majority of visually responsive units also respond to sound as defined by either significant firing rate responses (onset or sustained) or STRF reliability. Bars represent the union of units with significant STRF reliability and units with significant firing rate responses (onset or sustained).

538

539    The preceding results indicated that if a neuron was responsive to visual stimuli, it was
540    likely also responsive to sound. We addressed a corollary question of whether the presence or
541    absence of an auditory response predicted whether a unit was responsive to visual stimuli. Chi-
542    squared tests were used to compare percentages of visual-responsive units within unit
543    subpopulations separated by significant and non-significant auditory responses. As depicted in
544    **Figure 6**, visual responsiveness was not strongly predicted by the presence or absence of any
545    auditory response type for either unit subpopulation. Small but statistically significant effects
546    were observed for NS units in the deepest bin for sustained spike rate changes (**Figure 6B**; $X^2$
547    = 4.63, p = 0.031) and STRFs (**Figure 6C**; $X^2$ = 4.28, p = 0.039). A similar trend was observed
548    for sustained responses in the middle depth bin (**Figure 6B**; $X^2$ = 3.84, p = 0.050). These results
549    suggested that the absence of either sound response was associated with decreased probability
550    of a visual response, recapitulating analyses above suggesting responses to the two modalities
551    tended to occur together. Differences for all other depth bins, and across depth bins for BS units
552    were non-significant (all $X^2$ < 1.90, all p-values > 0.16). Thus, with minor exceptions, visual
553    responses were approximately evenly distributed among sound response types, leaving cortical
554    depth as the most meaningful predictor of visual responsiveness.
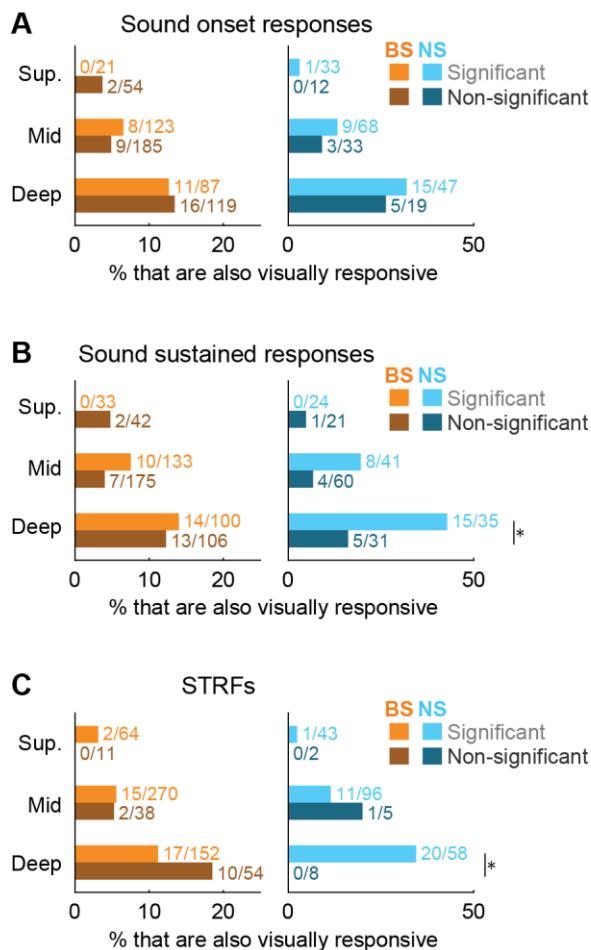555



**Figure 6.** Sound response archetypes are not strongly predictive of visual responsiveness. Each subplot shows the proportion of BS and NS units (left and right) separated by cortical depth with and without significant sound responses (light and dark bars) as defined by (**A**) onset firing rate responses, (**B**) sustained firing rate responses, and (**C**) STRFs. NS units in the deepest cortical bin with significant sustained sound responses and STRFs were more likely to have significant visual responses (Chi-square tests: *p<0.05). The presence or absence of any other sound response at any other depth bin did not predict visual responsiveness.

556    *Visual context differentially modulates sound onset and sustained firing rate responses*
557

558   The foregoing results documented unimodal visual responses in a subpopulation of A1 neurons,
559   and further that most of these units were bimodal, responding to both visual and auditory stimuli.
560   We investigated further the extent of audiovisual integration in A1 by examining the frequency
561   with which sound onset and sustained firing rate responses were themselves modulated by the
562   presence or absence of the visual stimulus. **Figure 7A** shows an example neuron for which both
563   sound onset and sustained responses significantly increased on audiovisual trials. For the
564   example unit in **Figure 7B**, the sustained response decreased significantly on audiovisual trials,
565   but the onset was unaffected. The latter outcome was typical of units with visual-modulated
566   responses to sound. Indeed, the example neuron in **Figure 7A** was the only unit in the entire
567   population for which sound onset response changed significantly with visual stimulation (**Figure
568   7C, a**). By comparison, sustained sound firing rate responses were significantly modulated by
569   visual context in 7.1% of NS units and 5.1% of BS units (**Figure 7C, b**). ANOVA confirmed
570   visual modulation effect sizes were larger for sustained than onset responses for the middle and
571   deep cortical depth bins for both unit types (BS units: shallow: $F = 0.16$, $p = 0.687$, $\eta^2 = 0.001$;
572   middle: $F = 8.62$, $p = 0.003$, $\eta^2 = 0.014$; deep: $F = 9.79$, $p = 0.002$, $\eta^2 = 0.024$; NS units: shallow:
573   $F = 0.10$, $p = 0.756$, $\eta^2 = 0.001$; middle: $F = 6.60$, $p = 0.011$, $\eta^2 = 0.032$; deep: $F = 9.46$, $p =$
574   $0.003$, $\eta^2 = 0.068$). Notably, these outcomes were exactly opposite of responses driven by
575   sound alone, in which case onsets were substantially stronger than sustained firing rate
576   changes (**Figure 3**).
577   Consistent with the concentration of unimodal visual responses in the deepest cortical
578   bin (**Figure 2**), the majority of units with visual-modulated sustained responses to sound were
579   observed in either the middle or deep cortical bin, with the strongest mean effect size in the
580   deepest bin (**Figure 7C, c–d**). ANOVA confirmed small but significant effects of cortical depth
581   on visual modulated sustained sound responses for both unit types (BS units: $F = 3.22$, $p =$
582   $0.041$, $\eta^2 = 0.011$; NS units: $F = 4.69$, $p = 0.010$, $\eta^2 = 0.043$). Visual modulated onset responses
583   were not significantly dependent upon depth (BS units: $F = 2.35$, $p = 0.096$, $\eta^2 = 0.008$; NS
584   units: $F = 0.25$, $p = 0.780$, $\eta^2 = 0.002$). All differences between unit type were similarly non-
585   significant (all F-ratios < 1.5, p-values > 0.23), with the exception of borderline effect suggesting
586   stronger sustained response modulation for BS units in the shallowest bin ($F = 3.91$, $p = 0.051$,
587   $\eta^2 = 0.032$).
588   Sound-evoked sustained firing rates were suppressed below baseline for the example
589   units in **Figure 7A–B**. However, the effect of visual stimulation was opposite for each neuron,
590   elevating firing near the spontaneous rate for the example in **Figure 7A** and further decreasing
591   the suppressed response for the example in **Figure 7B**. These examples raised the possibility
592   that effects of visual stimulation on sound-evoked firing rate responses might differ depending
593   on whether responses to sound alone reflected an increase or decrease in spike rate relative to
594   baseline (spontaneous). We thus expanded the analyses above by dividing onset and sustained
595   responses into subgroups for which the sound-evoked response was greater or less than
596   spontaneous. Distributions of sound-evoked firing rate responses divided by baseline rates are
597   shown in **Extended Data Figure 7-1A–B**, indicating both increases and decreases in firing rate
598   relative to baseline were common for both onset and sustained responses in both unit
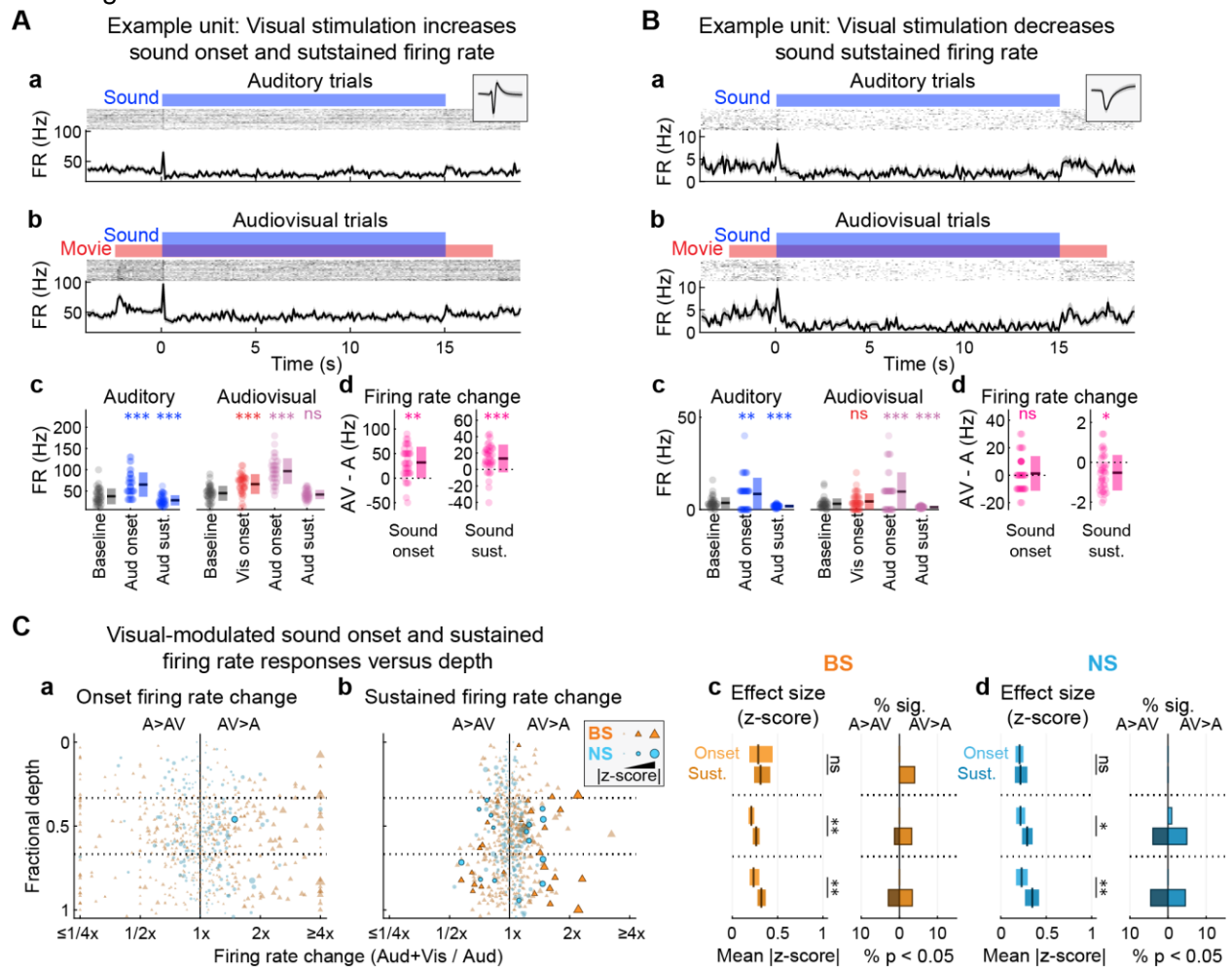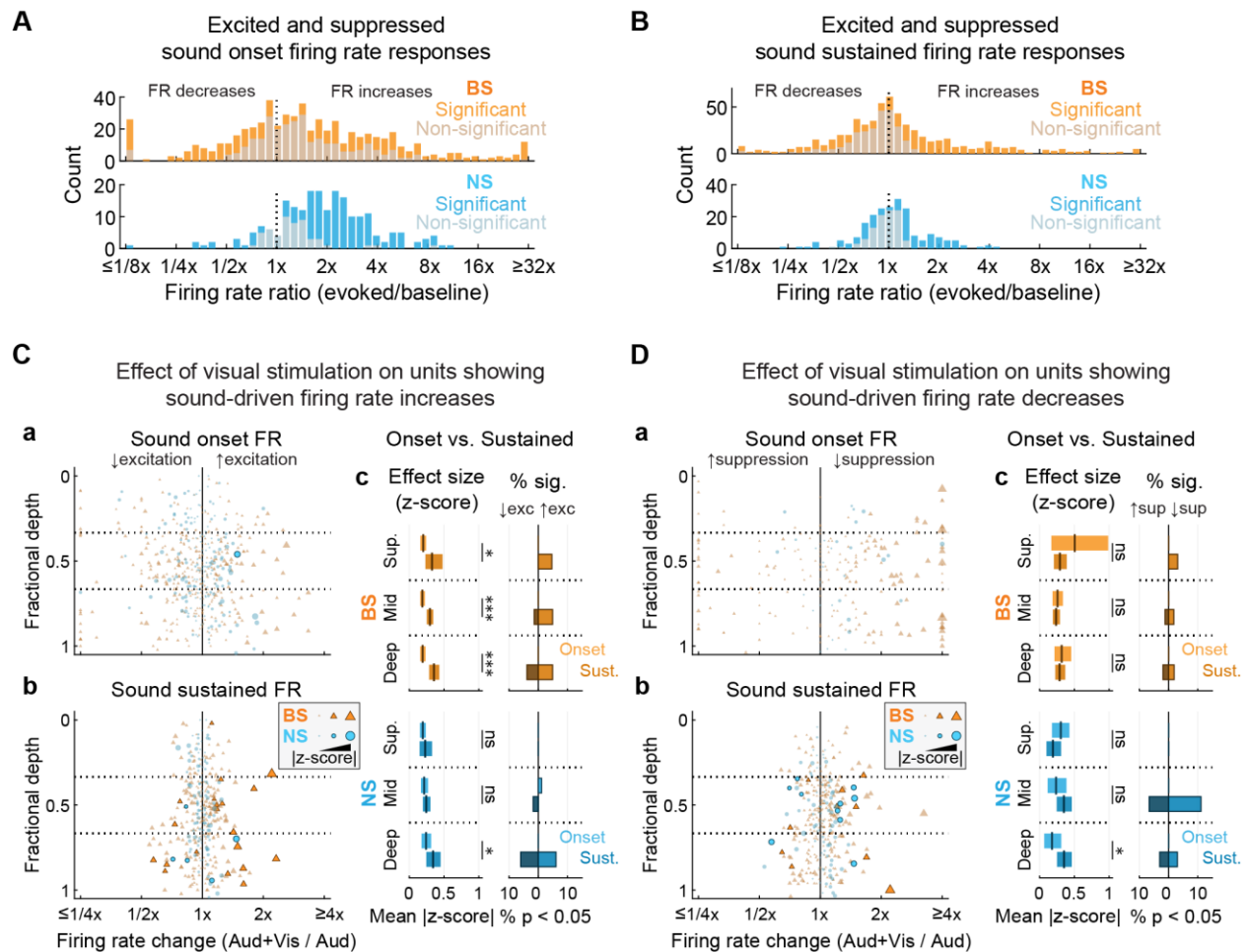599   subpopulations. Designations of 'excited' and 'suppressed'

**Figure 7.** Visual context differentially modulates sound onset and sustained firing rate responses. (**A**) Example unit for which visual stimulation significantly increased onset and sustained firing rate responses to sound. (**a**) Auditory and (**b**) Audiovisual trials. Spiking responses quantified by peristimulus time histograms (lower) and binned-spike count matrices (upper) with red and blue bars indicating auditory and visual stimulus intervals, respectively (temporal binning: 100 ms). Inset in (a) shows the unit spike waveform (median ± MAD). (**c**) Window-averaged firing rate responses, with dots representing single trials and mean ± SD across trials indicated to the right. Wilcoxon signed-rank tests (paired): ∗p<0.05, ∗∗p<0.01, ∗∗∗p<0.001, ns p>0.05. (**d**) Difference in averaged firing rate responses between conditions (Audiovisual - Auditory trials), with dots representing single trials and mean ± SD across trials indicated to the right. Wilcoxon signed-rank tests (paired): ∗p<0.05, ∗∗p<0.01, ∗∗∗p<0.001, ns p>0.05. (**B**) Example unit for which visual stimulation significantly increased sustained but not onset firing rate responses to sound. Subplot organization as in (A). (**C**) Visual modulation effects are stronger for sound sustained firing rate responses than onset responses. (**a, b**) Summary of visual-modulated sound onset and sustained firing rate changes (Audiovisual divided by Auditory) separated by unit type and cortical depth for all units. Scatter plots depict firing rate changes between conditions for each unit by cortical depth. Firing rate ratio values above and below 1 (x-axis) indicate increases and decreases in firing on Audiovisual trials relative to Auditory trials, respectively. Outlined markers indicate statistically significant responses. Marker sizes are scaled by effect size (absolute difference between Audiovisual and Auditory means, divided by Auditory SD). (**c**) Comparison of visual-modulated onset and sustained responses for BS units. Left: mean effect size (plus 99% confidence interval) across all recorded units (significant and non-significant responses included). Effect sizes were significantly greater for sustained responses (lower bars, darker coloring) than onset responses at the middle and deep bins. One-way ANOVA: ∗p<0.05, ∗∗p<0.01, ∗∗∗p<0.001, ns p>0.05. Right: Percentages of all recorded units with significant visual-modulated onset and sustained responses. (**d**) Comparison of onset and sustained responses for NS units, with subplot organization as in (c). Effect sizes were similarly greater for sustained responses (lower bars, darker coloring) at the middle and deep bins. See Extended Data Figure 7-1 for results separated by units for which sound responses reflected increases and decreases from baseline.

600

**Extended Data Figure 7-1.** Visual modulation of excited and suppressed sound onset and sustained firing rate responses. (**A**) Sound may evoke increases or decreases in onset firing rate. Histograms show the numbers of units (BS, top; NS, bottom) with increases and decreases in firing rate relative to baseline, with significant (p<0.05) and non-significant (p>=0.05) responses indicated by dark and light bars, respectively. (**B**) Sound may evoke increases or decreases in sustained firing rate. Histograms as in (A). (**C**) Summary of visual stimulation effects on onset and sustained firing rate responses that were excited relative to baseline. (**a, b**) Summary of visual-modulated sound onset and sustained firing rate changes (Audiovisual divided by Auditory) separated by unit type and cortical depth for all units with excited responses. Scatter plots depict firing rate changes between conditions for each unit by cortical depth. Firing rate ratio values above and below 1 (x-axis) indicate increases and decreases in excitation on Audiovisual trials relative to Auditory trials, respectively. Outlined markers indicate statistically significant responses. Marker sizes are scaled by effect size (absolute difference between Audiovisual and Auditory means, divided by Auditory SD). (**c**) Comparison of visual-modulated onset and sustained responses for BS (top) and NS units (bottom). Left: mean effect size (plus 99% confidence interval) across units with excited sound responses. Visual modulation effect sizes were significantly greater for sustained responses (lower bars, darker coloring) than onset responses across depth bins for BS units, and at the deepest bin for NS units. One-way ANOVA: *p<0.05, **p<0.01, ***p<0.001, ns p>0.05. Right: Percentages of all recorded units with significant visual-modulated onset and sustained responses. (**D**) Summary of visual stimulation effects on onset and sustained firing rate responses that were suppressed relative to baseline. Subplot organization as in (C).

601

602    units were inclusive of both significant and non-significant increases and decreases in
603    firing rate, respectively, to accommodate potential cases in which sound-evoked responses only
604    reached significance in one or the other condition (e.g., sustained response in **Figure 7A**). By
605    separating visual-modulatory influences into excited and suppressed subgroups, the analysis
606    included four possible forms of modulation: increases and decreases in excited and suppressed
607    responses.
608    Outcomes of this analysis were generally consistent with the pooled results reported
609    above. For BS units with excited responses (**Extended Data Figure 7-1C, c**), visual modulation
610    effects were stronger for sustained than onset responses at all depth bins (shallow: F = 5.99, p
611    = 0.016, $\eta^2$ = 0.062; middle: F = 27.20, p < $10^{-6}$, $\eta^2$ = 0.074; deep: F = 29.82, p < $10^{-6}$, $\eta^2$ =
612    0.115). For BS units with suppressed responses (**Extended Data Figure 7-1D, c**), differences
613    in visual modulation between sustained and onset responses were non-significant at all depth
614    bins (all F-ratios < 2.15, p-values > 0.14). These non-significant outcomes may reflect floor
615    effects limiting suppression of the already relatively low spontaneous rates in BS units. For NS
616    units with both excited and suppressed responses, visual modulation effects were stronger for
617    sustained than onset responses at the deepest bin only (excited: F = 4.59, p = 0.035, $\eta^2$ =
618    0.051; suppressed: F = 4.79, p = 0.034, $\eta^2$ = 0.100; all other F-ratios < 2.0, p-values > 0.16).
619    For the example neuron in **Figure 7A**, but not the example in **Figure 7B**, visual
620    modulated responses to sound coincided with a significant response to visual stimulation alone.
621    We examined the consistency of these outcomes by calculating intersections between unimodal
622    visual responses and visual-modulated sustained responses to sound. **Figure 8A** shows an
623    example unit for which a sustained, excitatory response to sound was further elevated by visual
624    stimulation even though the response to visual stimulation alone was not significant. This was
625    true for the majority of units with visual-modulated sustained sound responses, as seen in
626    **Figure 8B**, which did not respond outright to unimodal visual stimuli. This finding suggests only
627    partial overlap between visual responsive and visual modulated unit subpopulations.
628
629    *Visual stimulation may modify STRFs independently of firing rate changes*
630
631    As reported above, many units for which sound-evoked firing rate responses were significantly
632    modulated by visual stimulation were not significantly responsive to visual stimulation alone.
633    This suggested that STRFs might likely also be modulated by the presence of visual stimuli,
634    possibly for units without responses to unimodal visual stimuli. We therefore examined visual
635    influences on spectrotemporal encoding by calculating difference STRFs, referred to hereafter
636    as ΔSTRFs, reflecting the STRF in the auditory condition subtracted from the STRF in the
637    audiovisual condition. The significance of ΔSTRFs was determined by iteratively subsampling
638    STRFs from each condition and calculating the difference, then comparing the ΔSTRF
639    correlation distribution to an equivalent null distribution (**Figure 9A, a–b**). An example neuron
640    with significant ΔSTRF reliability is shown in **Figure 9B**. In total, 7.1% of NS units and 6.3% of
641    BS units had STRFs that were significantly modulated by visual stimulation (**Figure 9C**). There
642    was no evidence that ΔSTRF reliability depended significantly on cortical depth or unit type (all
643    F-ratios < 2.6, all p-values > 0.11). As with sustained firing rate responses to sound, visual
644    modulation of STRFs most often occurred without significant responses to visual stimulation
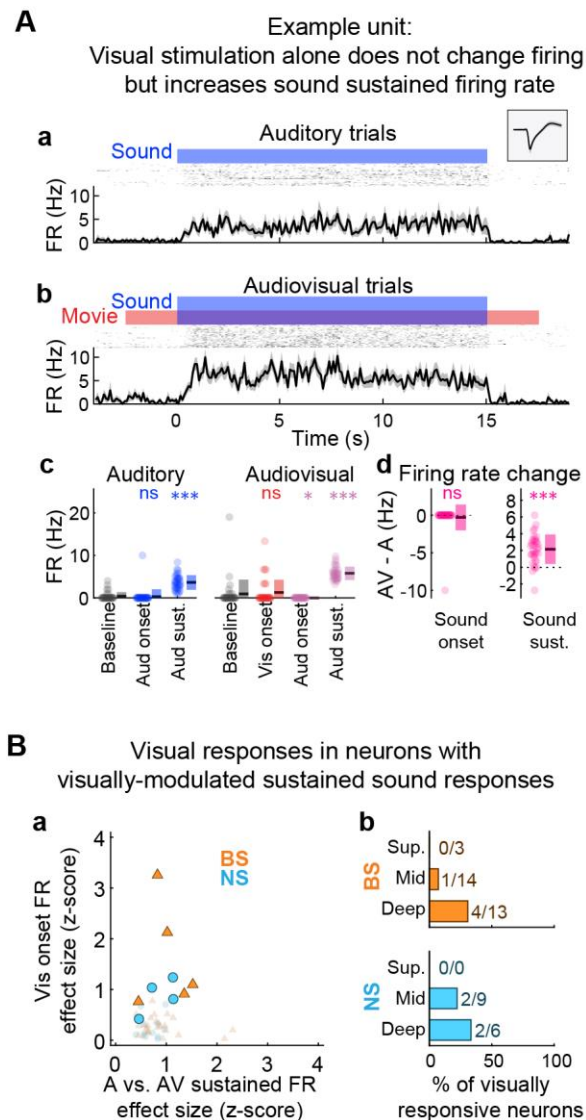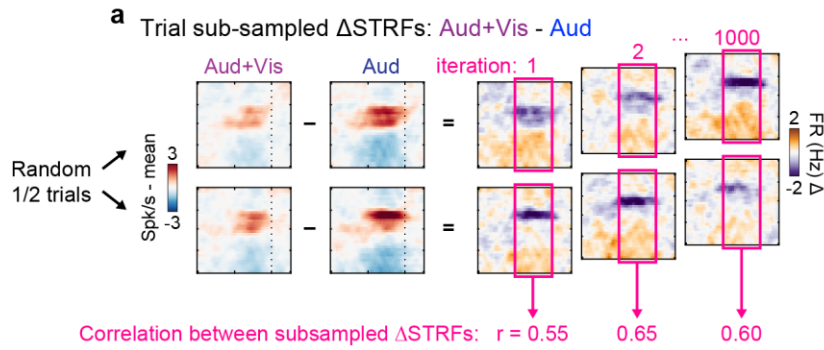645    alone (**Figure 9D**).

**Figure 8.** Sound-evoked firing rate responses may be modulated by visual inputs even without responses to visual stimulation alone. (**A**) Example unit with non-significant visual onset response but significant visual-modulated increases in sustained firing rate responses to sound. (**a**) Auditory and (**b**) Audiovisual trials. Spiking responses quantified by peristimulus time histograms (lower) and binned-spike count matrices (upper) with red and blue bars indicating auditory and visual stimulus intervals, respectively (temporal binning: 100 ms). Inset in (a) shows the unit spike waveform (median ± MAD). (**c**) Window-averaged firing rate responses, with dots representing single trials and mean ± SD across trials indicated to the right. Wilcoxon signed-rank tests (paired): ∗p<0.05, ∗∗p<0.01, ∗∗∗p<0.001, ns p>0.05. (**d**) Difference in averaged firing rate responses between conditions (Audiovisual - Auditory trials), with dots representing single trials and mean ± SD across trials indicated to the right. Wilcoxon signed-rank tests (paired): ∗p<0.05, ∗∗p<0.01, ∗∗∗p<0.001, ns p>0.05. (**B**) Many units with significant visual-modulated sustained sound firing rate responses are not responsive to visual stimulation alone. (**a**) Scatter plot of effect sizes for visual onset responses and visual-modulated sustained sound responses. Large markers with outlines reflect units with significant visual onset firing rate responses. (**b**) Bar plot showing the percentages of units with significant visual-modulated sustained sound responses that are also visual responsive.
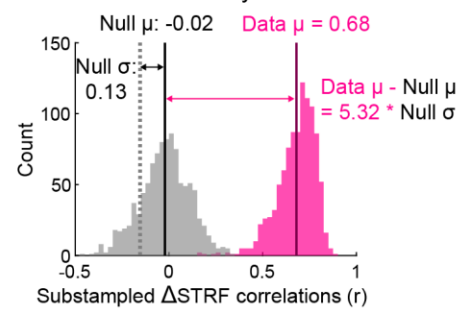
646        Two outcomes reported above further raised the possibility that STRFs may be
647    modulated by visual stimulation independently of firing rate changes. First, significant STRFs
648    were frequently obtained on auditory trials even in the absence of time-averaged firing rate
649    changes (**Figure 4**). Second, sound-evoked firing rate changes (onset vs. sustained) were
650    themselves independently modulated by visual context (**Figure 7**). To address this possibility,
651    we examined overlap between units with significant ΔSTRF reliability and visual-modulated
652    sustained firing rate responses. For the example unit in **Figure 10A**, the STRF was significantly
653    modulated by visual context, even though the sustained firing rate response did not significantly
654    change from baseline, let alone between conditions. Group data presented in **Figure 10B**
655    revealed that similar independent modulation of sustained firing rate responses and STRFs
656    occurred in the majority of units sensitive to visual context.
657
658

27

659

**Figure 9.** Visual context may modify spectrotemporal receptive fields even in units without significant responses to visual stimulation alone. (**A**) Procedure for estimating differences in STRFs between conditions (ΔSTRFs). (**a**) ΔSTRFs were estimated by subtracting the Audiovisual STRF from the Auditory alone STRF. A subsampling test was used to determine the statistical significance of time-frequency bin structure in the ΔSTRFs. The correlation coefficient between ΔSTRFs calculated from random trial halves (without replacement) was calculated across 1000 iterations. (**b**) ΔSTRF reliability was defined as the mean of the subsampled correlation coefficient distribution. A null ΔSTRF distribution was obtained from ΔSTRFs calculated using time-reversed stimulus RDS segments. A p-value was obtained by dividing the number of null STRF correlations exceeding the reliability index (data) by the number of iterations and multiplying by two for two-tailed significance. Effect size reflected the absolute difference between null and data means, divided by the null standard deviation. (**B**) Example unit with significant ΔSTRF reliability. (**a**) Spiking responses quantified by peristimulus time histograms (lower) and binned-spike count matrices (upper) with blue and red bars indicating auditory and visual stimulus intervals, respectively (temporal binning: 100ms). Inset shows the unit spike waveform (median ± MAD). (**b**) STRFs for each condition (**c**) ΔSTRF (**d**) STRF reliability for each condition. (**e**) ΔSTRF reliability. For (d) and (e), each dot represents the correlation between STRFs or ΔSTRFs for a single subsample iteration, with mean ± SD across trials indicated to the right. Subsampling test: ∗p<0.05, ∗∗p<0.01, ∗∗∗p<0.001, ns p>0.05. (**C**) Summary of ΔSTRF reliability by unit type and cortical depth. (**a**) ΔSTRF reliability for each unit at its estimated cortical depth. Marker sizes are scaled by effect size, with outlined markers indicating units with significant ΔSTRF reliability (p < 0.05, Benjamini–Hochberg FDR correction). (**b**) Mean effect size (plus 99% confidence interval) across all recorded units (units with significant and non-significant reliability included) by unit type and depth. No differences between unit types were observed. One-way ANOVA: ns p>0.05. (**c**) Histograms indicating percentages of all recorded units with significant ΔSTRF reliability. (**D**) Significant ΔSTRF reliability often occurs without significant visual onset firing rate changes. (**a**) Scatter plot of visual onset and ΔSTRF reliability response effect sizes. Large markers with outlines reflect units with significant visual onset firing rate responses. (**b**) Bar plot showing the percentages of units with significant ΔSTRF reliability that are also visually responsive.

660

**A**    Example unit: visual inputs alter STRF but not sound-evoked firing rates
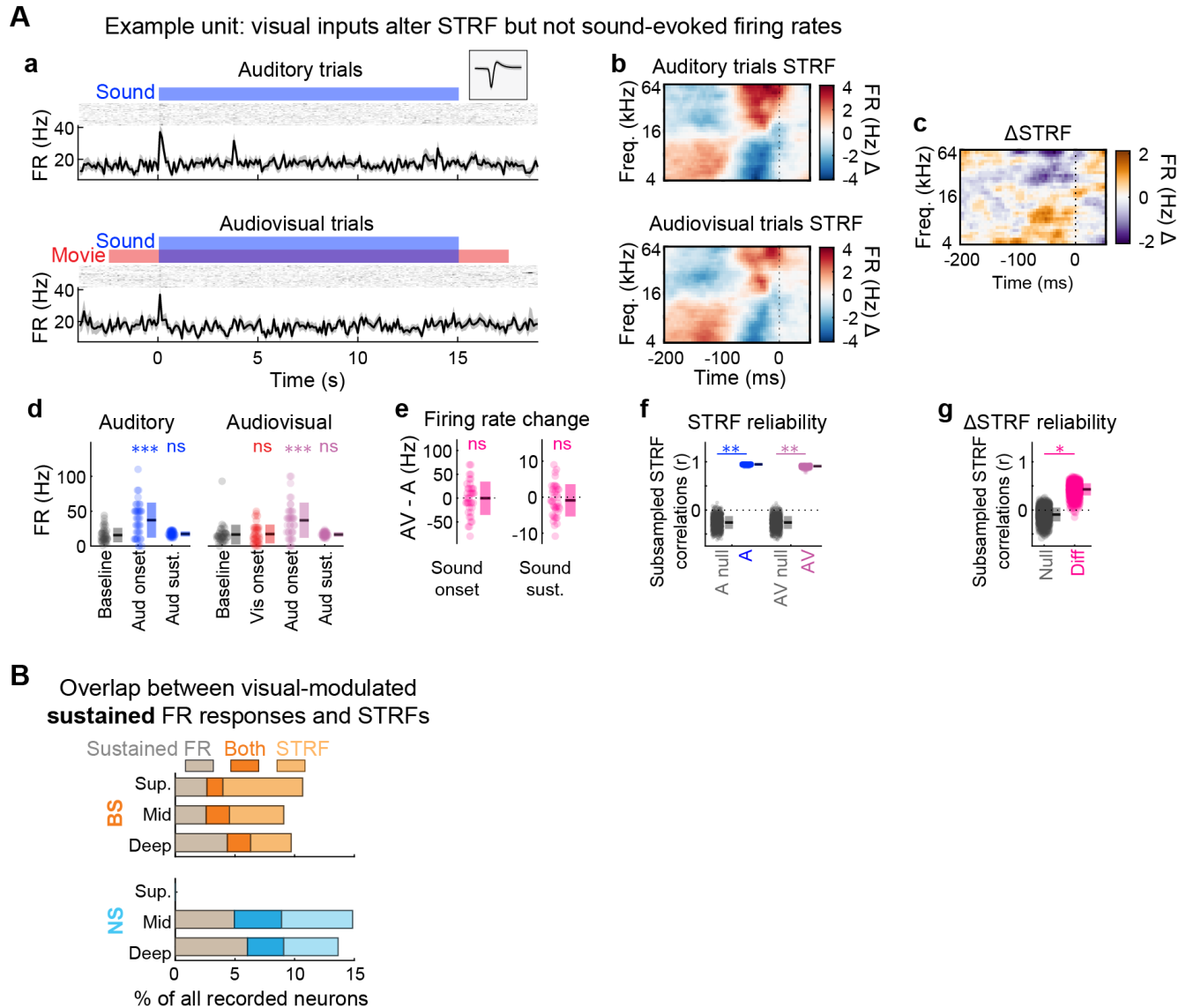


**Figure 10.** Visual stimulation may modify spectrotemporal receptive fields independently of averaged firing rates. (**A**) Example unit with significant ΔSTRFs reliability but non-significant visual-modulated sustained firing rate. (**a**) Spiking responses quantified by peristimulus time histograms (lower) and binned-spike count matrices (upper) for each condition with blue and red bars indicating auditory and visual stimulus intervals, respectively (temporal binning: 100ms). Inset shows the unit spike waveform (median ± MAD). (**b**) STRFs for each condition (**c**) ΔSTRF (**d**) Window-averaged firing rate responses, with dots representing single trials and mean ± SD across trials indicated to the right. Wilcoxon signed-rank tests (paired): ∗p<0.05, ∗∗p<0.01, ∗∗∗p<0.001, ns p>0.05. (**e**) Difference in averaged firing rate responses between conditions (Audiovisual - Auditory trials), with dots representing single trials and mean ± SD across trials indicated to the right. Wilcoxon signed-rank tests (paired): ∗p<0.05, ∗∗p<0.01, ∗∗∗p<0.001, ns p>0.05. (**f**) STRF reliability for each condition. (**g**) ΔSTRF reliability. For (f) and (g), each dot represents the correlation between STRFs or ΔSTRFs for a single subsample iteration, with mean ± SD across trials indicated to the right. Subsampling test: ∗p<0.05, ∗∗p<0.01, ∗∗∗p<0.001, ns p>0.05. (**B**) Significant ΔSTRF reliability may occur with or without significant visual modulation of sustained sound-evoked firing rate changes. Bar plots indicate intersections of significant visual modulated sustained firing rate responses alone, significant ΔSTRF reliability alone, or both.

661

662 *Visual context may alter STRF gain but preserves spectrotemporal tuning*

663

664 In the previous sections, we found that STRFs are sensitive to visual context for some units as
665 indicated by the reliability of ΔSTRFs. This metric confirmed that the STRFs differed significantly
666 between conditions but did not reveal what was different about the STRFs, qualitatively or
667 quantitatively. We therefore examined the degree to which differences between conditions
668 reflected gain (quantitative) or tuning changes (qualitative) by calculating slope and correlation
669 coefficient parameters for time-frequency bin values between conditions. For this analysis, we
670 included only STRFs for which the reliability index in the auditory condition was >0.5 to ensure
671 reliable baseline spectrotemporal tuning prior to estimating differences between conditions.
672 Example units with relatively strong and weak correlations between STRFs in the auditory and
673 audiovisual conditions are shown in **Figure 11A, a and b**, respectively. As indicated by **Figure
674 11A, c**, STRFs were usually highly correlated between conditions (NS median: r = 0.95; BS
675 median: r = 0.90), implying largely preserved spectrotemporal tuning between conditions.
676 Absolute correlations between auditory STRFs and ΔSTRFs were similarly high (NS median: r =
677 0.73; BS median: r = 0.62). This implied audiovisual STRFs were similarly structured, but with
678 either larger or smaller time-frequency bin values than auditory STRFs for units with positive
679 and negative correlations between auditory STRFs and ΔSTRFs, respectively.
680          Changes in STRF gain were modeled using standardized major axis regression, which is
681 designed for capturing bivariate relationships with estimation error affecting both variables
682 (Warton et al., 2006). With the auditory and audiovisual conditions plotted on the x and y axes,
683 respectively, positive and negative slopes for the regression model implied that visual
684 stimulation increased and decreased gain, respectively. Slopes were only analyzed for units for
685 which the relationship between conditions could be modeled with high accuracy ($r^2 > 0.5$).
686 STRFs were first downsampled by a factor of three in both dimensions to avoid inflating
687 correlations between conditions resulting from the smoothing operation (see Materials and
688 Methods). Example units with gain decreases and increases are shown in **Figure 11B, a–b** for
689 BS units and **Figure 11B, d–e** for NS units. Group data in **Figure 11B, c–f** indicated that gain
690 changes were diverse, showing both increases and decreases. In extreme cases, gain was
691 nearly doubled or halved. However, STRFs were highly correlated between conditions even for
692 units with large gain changes, implying STRF differences were more quantitative than
693 qualitative. As indicated by the y-axis marginal histogram in **Figure 11B, c**, the median gain
694 change was not significantly different from one for BS units (Wilcoxon signed-rank test: p =
695 0.888), implying units with gain increases and decreases were approximately balanced. A small
696 effect of borderline significance suggesting an overall tendency toward decreased gain (median
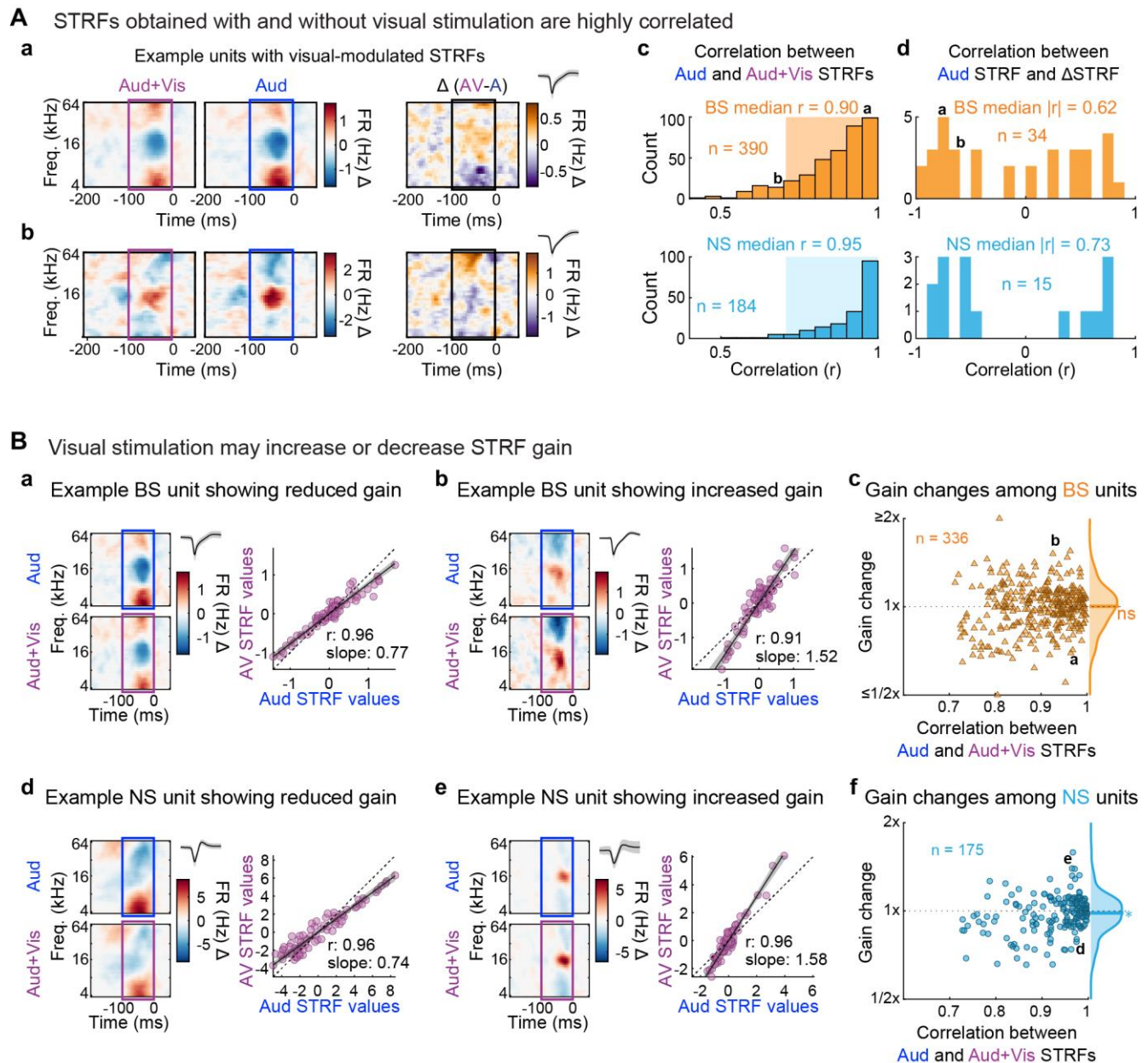697 0.984) was observed for the NS unit population (Wilcoxon signed-rank test: p = 0.044).

698

699

**A**  STRFs obtained with and without visual stimulation are highly correlated



**Figure 11.** Visual stimulation may change STRF response magnitude but preserves spectrotemporal tuning. (**A**) STRFs obtained with and without visual stimulation are highly correlated. (**a**) Example unit with highly correlated STRFs between conditions. (**b**) Example unit with moderately correlated STRFs between conditions. (**c**) Distributions of STRF correlations between conditions for BS (top) and NS units (bottom), indicating similar spectrotemporal tuning for the majority of units. (**d**) Distributions of correlations between Auditory and ΔSTRFs for each BS (top) and NS units (bottom). ΔSTRFs are typically highly correlated or anti-correlated with Auditory STRFs, suggesting visual stimulation produces global increases or decreases in time-frequency bin values without substantial modifications in STRF structure. (**B**) Visual inputs may alter STRF gain without altering spectrotemporal tuning. (**a**) Example BS unit with gain decrease in the Audiovisual condition. Left: STRFs for each condition with unit waveform (median ± MAD). Right: best fit lines to STRF time-frequency bins from each condition (shading indicates 95% confidence intervals). (**b**) Example BS unit with gain increase in the Audiovisual condition. Subplot organization as in (b). (**c**) Summary of STRF gain changes between conditions as a function of STRF correlations between conditions for BS units. Values above and below 1 (y-axis) indicate higher or lower gain in the Audiovisual condition, respectively. Deviations from 1 were observed even for units with highly correlated STRFs between conditions. Units with increases and decreases in gain were in approximate balance, such that no significant group-level difference from 1 was observed. (**d–e**) Example NS units with decreased and increased gain and in the Audiovisual condition, with subplot organization as in (a–b). (**f**) Summary of STRF gain changes between conditions as a function of STRF correlations between conditions for NS units. Decreases in gain were significantly more common than increases, such that the unit population median was significantly below 1. Wilcoxon signed-rank tests (paired): *p<0.05.

700

701     *Visual stimulation may increase or decrease auditory information*
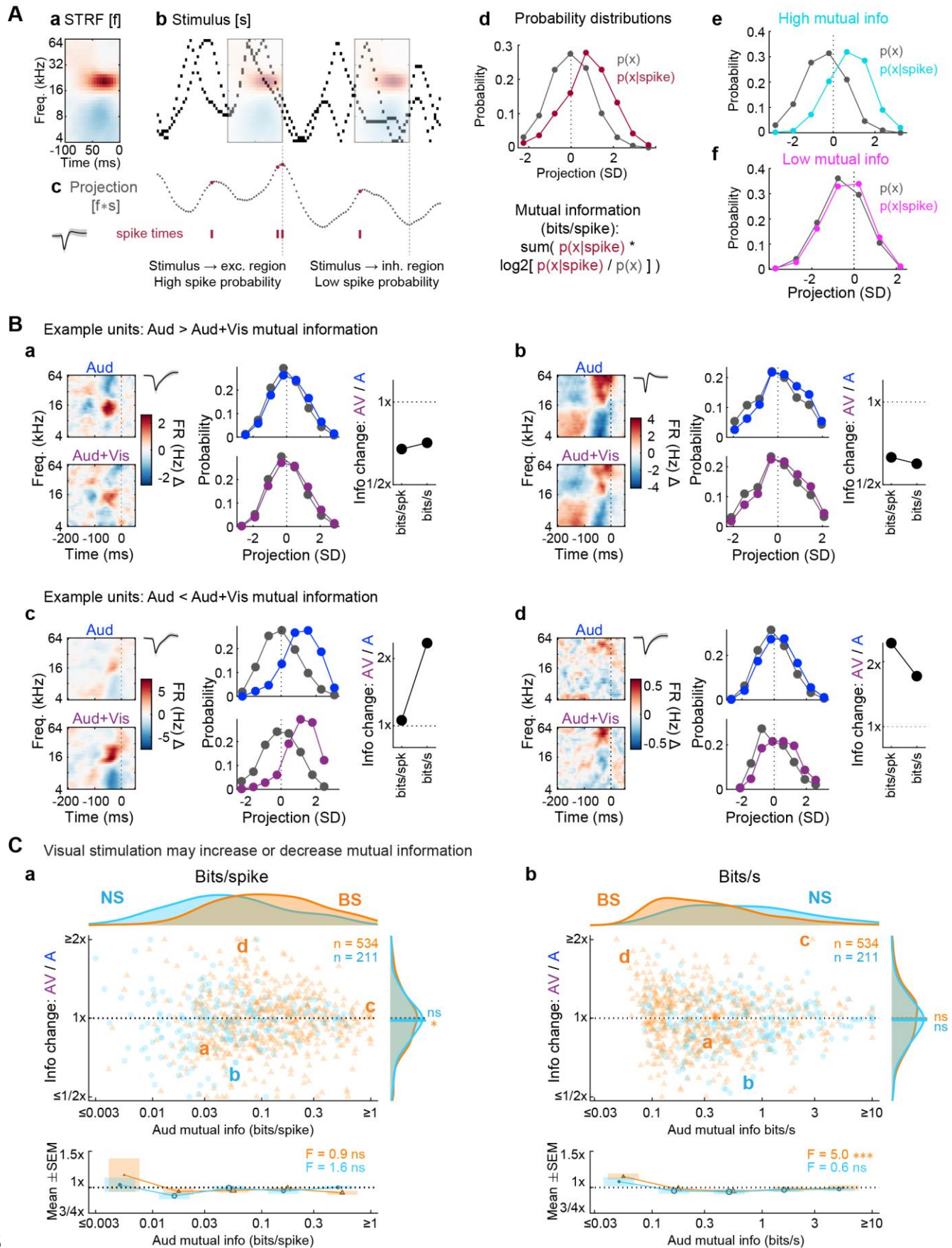
702

703     The analysis above showed that STRFs for some neurons either increased or decreased in gain

704     on audiovisual trials, i.e., showed greater firing rate deviations from the mean driven rate. This

705     suggested that concurrent visual stimulation may have 'helped' or 'hurt' encoding of the RDS

706     sound features by individual neurons. However, gain change analysis was restricted to units for

707     which STRFs from each condition were highly correlated ($r^2 > 0.5$), which also implied these

708     units had highly STRF reliability (since an unreliable STRF in one condition would not be well

709     correlated with the STRF from the other condition). We further explored the idea that visual

710     stimulation might facilitate or impede encoding of sounds in A1 using an information theoretic

711     approach to quantify visual-mediated changes in information between the auditory stimulus and

712     spike events. This analysis was inclusive of all units with at least 200 driven spikes (to avoid

713     undersampling concerns) regardless of STRF reliability or significance, including potential units

714     with poorly or randomly structured STRFs in one condition or the other. Because previous

715     studies have found multisensory interactions are often dependent upon baseline (unimodal)

716     responses (Allman and Meredith 2007), we further broke down information changes according

717     to baseline information rates in the auditory condition.

718     As depicted in **Figure 12A, a–c**, similarity between the STRF and stimulus was

719     quantified by convolution, producing a projection value (x) for each time bin used in STRF

720     calculation. The full sample of projection values defined the distribution p(x), which was

721     standardized by its mean and standard deviation. Relatively high and low projections values

722     represent stimulus segments that are similar and dissimilar to the STRF, respectively (examples

723     of each are superimposed over the binned stimulus in **Figure 12A, b**). For auditory responsive

724     units with significant STRFs, spikes are more likely to occur in time bins with high projection

725     values, and least likely to occur in bins with low projection values. This can be demonstrated by

726     identifying the subset of projection values that coincide with spikes, i.e., the distribution

727     p(x|spike). As shown in **Figure 12A, d**, values of p(x|spike) tend to be higher than p(x) for

728     positive standardized projection values, implying increased spike probability for closer matches

729     between the STRF and stimulus. Values of p(x|spike) moreover tend to be lower than p(x) for

730     negative projection values, implying poor matches between the STRF and stimulus result in

731     spiking probability below the mean firing rate. Mutual information is proportional to the ratio of

732     these distributions, as defined by pseudocode in **Figure 12A, d**. Intuitively, p(x) and p(x|spike)

733     become asymptotically equivalent over time for randomly a spiking neuron but diverge to the

734     degree to which a neuron is selectively activated by close stimulus approximations to its

735     receptive field (and suppressed otherwise). Thus, spikes from the example unit in **Figure 12A**

736     are said to be informative about the stimulus because we can infer that the structure of the

737     stimulus is relatively similar to the STRF when the ratio between p(x|spike) and p(x) is high, and

738     relatively dissimilar when the ratio is low.

739     A total of 211 NS units and 534 BS units were included in the analysis, reflecting the

740     subset of the full dataset with at least 200 driven spikes and information values greater than

741     zero. Example units with visual-mediated decreases and increases in mutual information are

742     presented in **Figure 12B, a–d**. At the population level, we found that BS units had significantly

743     higher bits per spike than NS units (**Figure 12C, a**; F = 72.17, $p < 10^{-15}$, $\eta^2 = 0.088$). However,

744     considering BS units tend to have lower firing rates than NS units, we also analyzed information

**A**

**a** STRF [f]

**b** Stimulus [s]

**c** Projection [f∗s]

spike times

Stimulus → exc. region
High spike probability

Stimulus → inh. region
Low spike probability

**d** Probability distributions

p(x)
p(x|spike)

Mutual information (bits/spike):
sum( p(x|spike) *
log2[ p(x|spike) / p(x) ] )

**e** High mutual info

p(x)
p(x|spike)

**f** Low mutual info

p(x)
p(x|spike)

**B** Example units: Aud > Aud+Vis mutual information

**a** Aud / Aud+Vis

**b** Aud / Aud+Vis

Example units: Aud < Aud+Vis mutual information

**c** Aud / Aud+Vis

**d** Aud / Aud+Vis

**C** Visual stimulation may increase or decrease mutual information

**a** Bits/spike

NS        BS

n = 534
n = 211

F = 0.9 ns
F = 1.6 ns

**b** Bits/s

BS        NS

n = 534
n = 211

F = 5.0 ***
F = 0.6 ns

**Figure 12.** Visual stimulation may increase or decrease information carried about an auditory stimulus. (**A**) Mutual information estimation procedure. (**a**) Example STRF, which serves as filter [f]. (**b**) Example RDS stimulus segment [s]. The overlaid STRFs highlight timepoints at which the RDS frequency vectors intersect excitatory and inhibitory regions of the STRF. (**c**) Projection values (A.U.) for each time point reflecting the convolution of the STRF and stimulus [f*s]. Relatively high and low projection values result from intersections of stimulus energy with excitatory and inhibitory regions of the STRF. Unit waveform shown at left (median ± MAD). (**d**) Distributions of all normalized projection values (gray) and those at time points for which spikes occurred (red). High and low projection values are associated with increases and decreases in spiking probability, respectively. Mutual information between the stimulus (projection value) and response (spike) is estimated from the relationship between these distributions as indicated by the pseudocode below the distributions. (**e–f**) Examples of projection value distributions (hypothetical) which would produce high (**e**) and low (**f**) mutual information values. (**B**) Example units for which visual context decreases and increases mutual information between sound stimuli and spiking responses. (**a–b**) Example units with visual-driven decreases in auditory information. Each example plot includes STRFs (left) and projection value distributions (middle) for each condition, plus an information change summary plot (right) depicting rate changes in both bits/spike and bits/s. Unit waveform shown at to the right of the Auditory STRF (median ± MAD). (**c–d**) Example units with visual-driven increases in information, with subplots organized as in (a–b). (**C**) Summary of changes in mutual information with the addition of visual stimulation, as a function of information in Auditory trials. (**a**) Summary of information changes expressed in bits/spike. Across unit populations, bits/spike tended to be higher for BS units (x-axis marginal). As depicted by the scatter plot (top), information changes were highly heterogeneous across both unit types, with both increases and decreases. Letters correspond to example units in (a–d). As indicated by the marginal histograms for the y-axis, a significant group level bias was observed for BS units only (medians shown by colored lines), with a tendency toward decreased information (Wilcoxon signed-rank tests [paired]: ∗p<0.05, ns p>0.05). As indicated by binned means below the x-axis of the scatter plot (mean ± SEM), increases and decreases in bits/spike were not significantly dependent upon baseline Auditory information values for either unit type (one-way ANOVA: ∗p<0.05, ∗∗p<0.01, ∗∗∗p<0.001, ns p>0.05). (**b**) Summary of information changes expressed in bits/s. Across unit populations, bits/s tended to be higher for NS units (x-axis marginal). No significant group-level biases were observed for either unit type (y-axis marginal). As indicated by binned means below the x-axis of the scatter plot (mean ± SEM), increases in bits/s were more likely for BS units with extremely low Auditory information values, with decreases more likely for moderate to high values. No significant relationship was observed for NS units.

746

747    in terms of bits per second by multiplying the bits/spike value for each unit by its mean firing rate
748    obtained from the sustained response window (spikes/s). The result indicated significantly
749    higher bits/s values for NS units (**Figure 12C, b**; F = 58.33, p < $10^{-13}$, $\eta^2$ = 0.072).
750          Changes in information rates for individual units were generally consistent with the gain
751    change results reported above. Information rates in the audiovisual condition could be either
752    greater or less than in the auditory condition, for both BS and NS units and for both bits/s and
753    bits/spike. As indicated by y-axis marginals in **Figure 12C, a**, a small but significant decrease in
754    median bits/spike was observed for BS units (Wilcoxon signed-rank test: p = 0.013) but a similar
755    trend for NS units was not significant (Wilcoxon signed-rank test: p = 0.072). Changes in
756    bits/spike did not significantly depend upon baseline auditory information rates for either unit
757    type (BS units: F = 0.874, p = 0.479, $\eta^2$ = 0.007; NS units: F = 1.61, p = 0.174, $\eta^2$ = 0.030).
758    Population level decreases in bits/s were non-significant for both unit types (Wilcoxon signed-
759    rank tests: BS units: p = 0.385; NS units: p = 0.067). For BS units, changes in bits/spike
760    depended significantly upon baseline auditory information rates (F = 5.01, p < $10^{-3}$, $\eta^2$ = 0.037),

35

761 such that units with very low auditory information (<0.1 bits/s) on average showed increases in
762 information transfer with visual stimulation, whereas information for most other units decreased.
763 A similar relationship was not observed for NS units was not significant (F = 0.613, p = 0.654, $\eta^2$
764 = 0.012).
765     Considered together, gain and information change analyses produced three insights
766 about the influence of visual context on sound encoding in A1. First, effects are highly
767 heterogeneous at the individual unit level, including units with both increases and decreases in
768 gain and information. Second, population level changes in gain and information in the
769 audiovisual condition are either subtle or non-significant, and generally reflect decreases
770 relative to the auditory condition. Third, whether visual context facilitates or interferes with sound
771 encoding may partially depend on baseline auditory responsiveness.
772
773 **Discussion**
774
775 Perception is inherently multisensory under natural conditions. Environmental events are often
776 encoded by more than one sensory modality, such as correlated visual and auditory cues
777 supporting spatial perception and conspecific communication (Sugihara et al., 2006; Allman and
778 Meredith 2007; Bigelow and Poremba, 2016). Even unisensory events must be encoded within
779 the context of uncorrelated sensory processing by other modalities. Growing recognition of the
780 ubiquity of such phenomena has resulted in increased attention to multisensory context within
781 the sensory physiology community, which has traditionally been dominated by unimodal
782 paradigms. One of the most important insights from these efforts is that multisensory interaction
783 is pervasive throughout cortex, including direct anatomical connections and functional
784 interaction between primary sensory cortices (Ghazanfar and Schroeder, 2006; Bizley et al.,
785 2007; Banks et al., 2011; Iurilli et al., 2012; Bizley et al., 2016; McClure and Polack 2019). For
786 instance, a recent study from our lab found that a subset of neurons in awake mouse A1 were
787 responsive to visual flash stimuli alone, and that these neurons were concentrated in the
788 infragranular layers (Morrill and Hasenstaub, 2018). These outcomes were replicated and
789 extended in the present study using the CMN visual stimulus (**Figure 2**), which also confirmed
790 that both putative inhibitory (NS) and excitatory neurons (BS) can be visually responsive. Of
791 these, all but a few were also significantly responsive to the RDS sound stimuli (**Figures 3, 5–
792 6**), suggesting purely visual neurons in A1 are very rare.
793     The current study further identified neurons in A1 for which responses to sound were
794 modulated by visual stimulation (**Figures 7–12**). A major question motivating the present study
795 was whether there was close overlap between units with such visual-modulated responses and
796 units responsive to visual stimulation alone. If so, the expected depth distribution of visual-
797 modulated responses to sound would be similarly concentrated in the deep cortical layers, as
798 was recently observed in a human imaging study (Gau et al., 2020). We obtained only mixed
799 support for this hypothesis. Fewer than half of neurons with visual-modulated sound responses
800 showed significant responses to visual stimulation alone. These outcomes replicate the findings
801 of previous studies that unimodal visual responses are neither necessary nor sufficient for
802 visual-modulated responses to sound. Consistent with these observations, the depth
803 distributions of units with visual-modulated sound responses and units with unimodal visual
804 responses were, at best, only weakly similar. For instance, in parallel with unimodal visual

805 responses, we found that visual influences on sound sustained firing responses depended on
806 cortical depth, with the majority of significant effects in the deepest cortical bin (**Figure 7**).
807 However, similar depth dependencies were not observed for modulated onset responses or
808 STRFs. Together, these outcomes support a model of audiovisual integration wherein visual
809 projections to infragranular A1 drive responses to visual stimulation alone, after which such
810 activity may modulate sound-evoked responses by propagating throughout cortical layers.
811 By delivering segmented RDS stimuli separated by intertrial intervals, we were able to
812 estimate time-averaged firing rate changes from baseline reflecting both transient onset
813 responses and sustained responses throughout the sound period. Moreover, averaging the
814 binned stimulus values within a window preceding each spike event time enabled estimation of
815 STRFs, which are sensitive to both spike count and timing with respect to the stimulus.
816 Capturing each response type turned out to be important both for characterizing responses to
817 the sounds themselves and for the dependence of these responses on visual context. Notably,
818 a double dissociation was observed for onset and sustained firing responses, wherein much
819 stronger baseline deviation effects were observed for onset responses on auditory trials (**Figure**
820 **3**), but much stronger visual modulation effects were observed for sustained responses on
821 audiovisual trials (**Figure 7**). We further found that neither firing rate response was necessarily
822 present in units with significantly reliable STRFs (**Figure 4**). Similarly, only partial overlap was
823 observed among units with visual-modulated sustained firing rate responses and units with
824 visually modulated STRFs, even though the spikes entering each analysis were identical
825 (**Figure 10**). Collectively, these outcomes reinforce previous studies concluding that spike rate
826 and timing changes in A1 carry non-redundant information about the stimulus (Brugge and
827 Merzenich, 1973; deCharms and Merzenich 1996; Recanzone, 2000; Lu et al., 2001; Wang et
828 al., 2005; Malone et al., 2010; Malone et al., 2015; Insanally et al., 2019; Liu et al., 2019).
829 An important caveat regarding onset and sustained firing rate responses in the current
830 study is that they were elicited by complex, non-repeating stimuli. This aspect of our study
831 design has at least two implications for the results. First, averaged spike rates within any given
832 time bin, and across full analysis windows, reflect mixed responses to diverse spectrotemporal
833 features produced by random draws from the RDS stimulus distribution. Thus, they likely reflect
834 a combination of weak and strong excitatory responses, as well as suppressed responses
835 resulting from intersections of spectral energy with receptive field inhibitory sidebands,
836 refractory periods, and from simultaneous and forward masking effects the stimulus was
837 designed to produce (Gourévitch et al., 2015). By contrast, binned spike averages elicited by a
838 repeating RDS stimulus segment would reflect specific spectrotemporal features that drive or
839 inhibit spiking activity. A neuron without a significant change in averaged sustained firing in a
840 non-repeating stimulus paradigm may thus exhibit regular temporal patterns of excitation and
841 suppression in a repeating stimulus paradigm. Second, onset and sustained responses
842 observed in the current study are not directly comparable to many previous studies examining
843 diverse firing rate responses in A1 (onset, sustained, offset), the majority of which presented
844 repeated tone pips, often at the best frequency of the neuron (Brugge and Merzenich, 1973;
845 Recanzone, 2000; Wang et al., 2005; Malone et al., 2015; Liu et al., 2019). Thus, it cannot be
846 assumed that onset and sustained responses observed in previous studies would necessarily
847 be subject to similar asymmetries in terms of visual context modulation. The diverse averaged
848 firing rate response types observed in the present study are thus qualitatively different from

849  previous results obtained with repeated tone pips, but nonetheless underscore the same
850  conclusion that different firing rate responses contain unique information about the stimulus and
851  are not equivalently sensitive to sensory and other contextual variables.
852      Previous studies of multisensory integration across multiple modalities in a wide range of
853  cortical and subcortical stations have reported that stimulus encoding may be either facilitated
854  or impeded by the addition of a second modality (Perrault et al., 2003; Romanski, 2007; Stein
855  and Stanford 2008; Sugihara et al., 2010; Kobayasi and Riquimaroux, 2013). Similar outcomes
856  have been obtained with both correlated and uncorrelated bimodal stimulation (Dahl et al.,
857  2010), with several studies suggesting benefits are more common with correlated stimulation
858  (Diehl and Romanski, 2014; Meijer et al., 2017). Costs are thought to possibly reflect
859  competition by each modality over limited spike resources within a given neuron. For example,
860  in a neuron responsive to two modalities, a subset of the spikes otherwise under the control of a
861  single modality may instead be driven by a second modality during bimodal stimulation, thus
862  potentially decreasing the mutual information between spiking and the stimulus of the first
863  modality. In correlated bimodal stimulus paradigms, such as natural audiovisual speech,
864  encoding benefits of bimodal stimulation likely reflect synergistic interaction of temporally and/or
865  spatially coincident feature selectivity within each respective modality. In uncorrelated bimodal
866  stimulus paradigms, including the present study, improved encoding with the addition of an
867  uncorrelated signal from another modality may reflect stochastic resonance, in which sensitivity
868  to a weakly detectable stimulus may be increased even by random noise within or between
869  modalities (Wiesenfeld and Moss, 1995; Shu et al. 2003; Hasenstaub et al. 2005; Crosse et al.,
870  2016; Malone et al., 2017). Consistent with this possibility, some evidence obtained in the
871  present study suggested that auditory information 'costs' and 'benefits' produced by visual
872  stimulation were non-randomly distributed with respect to baseline information rates obtained
873  with auditory stimulation alone. Specifically, information increases for BS units (bits/s) were
874  most likely for units with the lowest auditory information rates, suggesting weak and inconsistent
875  responses to sound features may be strengthened and regularized by the additional visual
876  stimulus. These and similar findings by previous studies establish neural phenomena parallel to
877  psychophysical experiments reporting significantly improved detection of weakly perceptible
878  events following the addition of both correlated and uncorrelated stimulation within a second
879  modality (Ward et al., 2010; Huang et al., 2017; Gleiss and Kayser, 2014; Krauss et al., 2018).

**References**

880 Allman BL, Meredith MA (2007) Multisensory processing in "unimodal" neurons: cross-modal subthreshold auditory effects in cat extrastriate visual cortex. J Neurophys 98:545-549.

Atencio CA, Schreiner CE (2008) Spectrotemporal processing differences between auditory cortical fast-spiking and regular-spiking neurons. J Neurosci 28(15):3897-3910.

Atencio CA, Schreiner CE (2013) Stimulus choices for spike-triggered receptive field analysis. In: Depireux DA, Elhilali M, editors. Handbook of modern techniques in auditory cortex. New York: Nova Biomedical, pp. 61-100.

Atencio CA, Schreiner CE (2016) Functional congruity in local auditory cortical microcircuits. Neuroscience 316:402-419.

Atencio CA, Sharpee TO, Schreiner CE (2008) Cooperative nonlinearities in auditory cortical neurons. Neuron 58:956-966.

Atencio CA, Sharpee TO, Schreiner CE (2009) Hierarchical computation in the canonical auditory cortical circuit. Proc Natl Acad Sci USA 106:21894-21899.

Atilgan H, Town SM, Wood KC, Jones GP, Maddox RK, Lee AK, Bizley JK (2018) Integration of visual information in auditory cortex promotes auditory scene analysis through multisensory binding. Neuron 97:640-655.

Banks MI, Uhlrich DJ, Smith PH, Krause BM, Manning KA (2011) Descending projections from extrastriate visual cortex modulate responses of cells in primary auditory cortex. Cereb Cortex 21:2620-2638.

Benjamini Y, Hochberg Y (1995) Controlling the false discovery rate: a practical and powerful approach to multiple testing. J R Stat Soc B 57:289–300.

Bigelow J, Malone BJ (2017) Cluster-based analysis improves predictive validity of spike-triggered receptive field estimates. PloS ONE 12:e0183914.

Bigelow J, Morrill RJ, Dekloe J, Hasenstaub AR (2019) Movement and VIP interneuron activation differentially modulate encoding in mouse auditory cortex. eNeuro 6(5).

Bigelow J, Poremba A (2016) Audiovisual integration facilitates monkeys' short-term memory. Anim Cogn 19:799-811.

Bizley JK, Jones GP, Town SM (2016) Where are multisensory signals combined for perceptual decision-making?. Curr Opin Neurobiol 40:31-37.

Bizley JK, King AJ (2008) Visual–auditory spatial processing in auditory cortical neurons. Brain Res 1242:24-36.

Bizley JK, Nodal FR, Bajo VM, Nelken I, King AJ (2007) Physiological and anatomical evidence for multisensory interactions in auditory cortex. Cereb Cortex 17:2172-2189.

Brugge JF, Merzenich MM (1973) Responses of neurons in auditory cortex of the macaque monkey to monaural and binaural stimulation. J Neurophys 36:1138-1158.

Churchland MM, Byron MY, Cunningham JP, Sugrue LP, Cohen MR, Corrado GS, Newsome WT, Clark AM, Hosseini P, Scott BB, Bradley DC (2010) Stimulus onset quenches neural variability: a widespread cortical phenomenon. Nat Neurosci 13:369-378.

Crosse MJ, Di Liberto GM, Lalor EC (2016) Eye can hear clearly now: inverse effectiveness in natural audiovisual speech processing relies on long-term crossmodal temporal integration. J Neurosci 36:9888-9895.

923 Dahl C, Logothetis NK, Kayser C (2010) Modulation of visual responses in the superior temporal
924     sulcus by audio-visual congruency. Front Integr Neurosci 4:10.
925 deCharms RC, Blake DT, Merzenich MM (1998) Optimizing sound features for cortical neurons.
926     Science 280:1439-1444.
927 deCharms RC, Merzenich MM (1996) Primary cortical representation of sounds by the
928     coordination of action-potential timing. Nature 381:610-613.
929 Diehl MM, Romanski LM (2014) Responses of prefrontal multisensory neurons to mismatching
930     faces and vocalizations. J Neurosci 34:11233-11243.
931 Dombeck DA, Khabbaz AN, Collman F, Adelman TL, Tank DW (2007) Imaging large-scale
932     neural activity with cellular resolution in awake, mobile mice. Neuron 56:43–57.
933 Escabí MA, Read HL, Viventi J, Kim DH, Higgins NC, Storace DA, Liu AS, Gifford AM, Burke
934     JF, Campisi M, Kim YS (2014) A high-density, high-channel count, multiplexed µECoG
935     array for auditory-cortex recordings. J Neurophys 112:1566-1583.
936 Fritz J, Shamma S, Elhilali M, Klein D (2003) Rapid task-related plasticity of spectrotemporal
937     receptive fields in primary auditory cortex. Nat Neurosci 6:1216-1223.
938 Gau R, Bazin PL, Trampel R, Turner R, Noppeney U (2020) Resolving multisensory and
939     attentional influences across cortical depth in sensory cortices. Elife 9:e46856.
940 Ghazanfar AA, Schroeder CE (2006) Is neocortex essentially multisensory? Trends Cogn Sci
941     10:278-285.
942 Gleiss S, Kayser C (2014) Acoustic noise improves visual perception and modulates occipital
943     oscillatory states. J Cogn Neurosci 26:699-711.
944 Gourévitch B, Occelli F, Gaucher Q, Aushana Y, Edeline JM (2015) A new and fast
945     characterization of multiple encoding properties of auditory neurons. Brain Topogr
946     28:379-400.
947 Hasenstaub A, Shu Y, Haider B, Kraushaar U, Duque A, McCormick DA (2005) Inhibitory
948     postsynaptic potentials carry synchronized frequency information in active cortical
949     networks. Neuron 47:423-435.
950 Huang J, Sheffield B, Lin P, Zeng FG (2017) Electro-tactile stimulation enhances cochlear
951     implant speech recognition in noise. Sci Rep 7:1-5.
952 Insanally MN, Carcea I, Field RE, Rodgers CC, DePasquale B, Rajan K, DeWeese MR,
953     Albanna BF, Froemke RC (2019) Spike-timing-dependent ensemble encoding by non-
954     classically responsive cortical neurons. Elife 8:e42409.
955 Iurilli G, Ghezzi D, Olcese U, Lassi G, Nazzaro C, Tonini R, Tucci V, Benfenati F, Medini P
956     (2012). Sound-driven synaptic inhibition in primary visual cortex. Neuron 73(4):814-828.
957 Joachimsthaler B, Uhlmann M, Miller F, Ehret G, Kurt S (2014) Quantitative analysis of neuronal
958     response properties in primary and higher-order auditory cortical fields of awake house
959     mice (*Mus musculus*). Eur J Neurosci 39:904-918.
960 Kayser C, Logothetis NK, Panzeri S (2010) Visual enhancement of the information
961     representation in auditory cortex. Curr Biol 20:19-24.
962 Kayser C, Petkov CI, Logothetis NK. Multisensory interactions in primate auditory cortex: fMRI
963     and electrophysiology. Hear Res 258:80-88.
964 King AJ, Hammond-Kenny A, Nodal FR (2019) Multisensory processing in the auditory cortex.
965     In: Lee A, Wallace M, Coffin A, Popper A, Fay R, editors. Multisensory Processes.
966     Springer Handbook of Auditory Research. Springer, Cham. pp. 105-133

967  Kleiner M, Brainard D, Pelli D (2007) What's new in Psychtoolbox-3? ECVP Abstr Suppl 14.

968  Kobayasi KI, Riquimaroux H (2013) Audiovisual integration in the primary auditory cortex of an
969      awake rodent. Neurosci Lett 534:24-29.

970  Kopp-Scheinpflug C, Sinclair JL, Linden JF. When sound stops: offset responses in the auditory
971      system. Trends Neurosci 41:712-728.

972  Liang Z, Shen W, Sun C, Shou T (2008) Comparative study on the offset responses of simple
973      cells and complex cells in the primary visual cortex of the cat. Neuroscience 156(2):365-
974      373.

975  Liu J, Whiteway MR, Sheikhattar A, Butts DA, Babadi B, Kanold PO (2019) Parallel processing
976      of sound dynamics across mouse auditory cortex via spatially patterned thalamic inputs
977      and distinct areal intracortical circuits. Cell Rep 27:872-885.

978  Lu T, Liang L, Wang X (2001) Temporal and rate representations of time-varying signals in the
979      auditory cortex of awake primates. Nat Neurosci 4:1131-1138.

980  Malone BJ, Heiser MA, Beitel RE, Schreiner CE (2017) Background noise exerts diverse effects
981      on the cortical encoding of foreground sounds. J Neurophys 118:1034-1054.

982  Malone BJ, Scott BH, Semple MN. Temporal codes for amplitude contrast in auditory cortex. J
983      Neurosci 30:767-784.

984  Malone BJ, Scott BH, Semple MN (2015) Diverse cortical codes for scene segmentation in
985      primate auditory cortex. J Neurophys 113:2934-2952.

986  McClure Jr JP, Polack PO (2019) Pure tones modulate the representation of orientation and
987      direction in the primary visual cortex. J Neurophys 121:2202-2214.

988  Meijer GT, Montijn JS, Pennartz CM, Lansink CS (2017) Audiovisual modulation in mouse
989      primary visual cortex depends on cross-modal stimulus configuration and congruency. J
990      Neurosci 37:8783-8796.

991  Morrill RJ, Hasenstaub AR (2018) Visual information present in infragranular layers of mouse
992      auditory cortex. J Neurosci 38:2854-2862.

993  Niell CM, Stryker MP (2008) Highly selective receptive fields in mouse visual cortex. J Neurosci
994      28:7520-7536.

995  Niell CM, Stryker MP (2010) Modulation of visual responses by behavioral state in mouse visual
996      cortex. Neuron 65:472-479.

997  Pachitariu M., Steinmetz N, Kadir S, Carandini M, Kenneth DH (2016) Kilosort: realtime spike-
998      sorting for extracellular electrophysiology with hundreds of channels. bioRxiv: 061481.

999  Paxinos G, Franklin KBJ (2019) Paxinos and Franklin's the mouse brain in stereotaxic
1000     coordinates. Academic press.

1001 Perrault Jr TJ, Vaughan JW, Stein BE, Wallace MT (2003) Neuron-specific response
1002     characteristics predict the magnitude of multisensory integration. J Neurophys 90:4022-
1003     4026.

1004 Phillips EA, Hasenstaub AR (2016) Asymmetric effects of activating and inactivating cortical
1005     interneurons. Elife 5:e18383.

1006 Phillips EA, Schreiner CE, Hasenstaub AR (2017a) Cortical interneurons differentially regulate
1007     the effects of acoustic context. Cell Rep 20(4):771-8.

1008 Phillips EA, Schreiner CE, Hasenstaub AR (2017b) Diverse effects of stimulus history in waking
1009     mouse auditory cortex. J Neurophys 118:1376-1393.

Recanzone GH (2000) Response profiles of auditory cortical neurons to tones and noise in behaving macaque monkeys. Hear Res 150:104-118.

Romanski LM (2007) Representation and integration of auditory and visual stimuli in the primate ventral lateral prefrontal cortex. Cereb Cortex 17:i61-69.

Rutkowski RG, Shackleton TM, Schnupp JW, Wallace MN, Palmer AR (2002) Spectrotemporal receptive field properties of single units in the primary, dorsocaudal and ventrorostral auditory cortex of the guinea pig. Audiol Neurootol 7:214-227.

Schneider DM, Mooney R (2018) How movement modulates hearing. Annu Rev Neurosci 41:553-572.

Scholl B, Gao X, Wehr M (2010) Nonoverlapping sets of synapses drive on responses and off responses in auditory cortex. Neuron 65(3):412-421.

See JZ, Atencio CA, Sohal VS, Schreiner CE (2018) Coordinated neuronal ensembles in primary auditory cortical columns. Elife 7:e35587.

Shu Y, Hasenstaub A, Badoual M, Bal T, McCormick DA (2003) Barrages of synaptic activity control the gain and sensitivity of cortical neurons. J Neurosci 23:10388-10401.

Stein BE, Stanford TR (2008) Multisensory integration: current issues from the perspective of the single neuron. Nat Rev Neurosci 9:255-266.

Stringer C, Pachitariu M, Steinmetz N, Reddy CB, Carandini M, Harris KD (2019) Spontaneous behaviors drive multidimensional, brainwide activity. Science 364(6437).

Sugihara T, Diltz MD, Averbeck BB, Romanski LM (2006) Integration of auditory and visual communication information in the primate ventrolateral prefrontal cortex. J Neurosci 26:11138-11147.

Thorson IL, Liénard J, David SV (2015) The essential complexity of auditory receptive fields. PLoS Comput Biol 11:e1004628.

Wang X, Lu T, Snider RK, Liang L (2005) Sustained firing in auditory cortex evoked by preferred stimuli. Nature 435:341-346.

Ward LM, MacLean SE, Kirschner A (2010) Stochastic resonance modulates neural synchronization within and between cortical sources. PloS ONE 5:e14371.

Warton DI, Wright IJ, Falster DS, Westoby M (2006) Bivariate line-fitting methods for allometry. Biol Rev 81:259-291.

Wiesenfeld K, Moss F (1995) Stochastic resonance and the benefits of noise: from ice ages to crayfish and SQUIDs. Nature 373:33-36.

Wu MC, David SV, Gallant JL (2006) Complete functional characterization of sensory neurons by system identification. Annu Rev Neurosci. 29:477-505.