

Humans use local spectrotemporal correlations to detect rising and falling pitch

Parisa A. Vaziri¹, Samuel D. McDougle^{*2,3}, Damon A. Clark^{*3,4,5,6,7}

1 – Yale College, Yale University, New Haven, CT 06511

2 – Dept of Psychology, Yale University, New Haven, CT 06511

3 – Wu Tsai Institute, Yale University, New Haven, CT 06511

4 – Dept of Molecular Cellular and Developmental Biology, Yale University, New Haven, CT 06511

5 – Dept of Physics, Yale University, New Haven, CT 06511

6 – Dept of Neuroscience, Yale University, New Haven, CT 06511

7 – Quantitative Biology Institute, Yale University, New Haven, CT 06511

* Equal contributors and lead contacts: samuel.mcdougle@yale.edu, damon.clark@yale.edu

Abstract

To discern speech or appreciate music, the human auditory system detects how pitch increases or decreases over time. However, the algorithms used to detect changes in pitch, or pitch motion, are incompletely understood. Here, using psychophysics, computational modeling, functional neuroimaging, and analysis of recorded speech, we ask if humans detect pitch motion using computations analogous to those used by the visual system. We adapted stimuli from studies of vision to create novel auditory correlated noise stimuli that elicited robust pitch motion percepts. Crucially, these stimuli possess no persistent features across frequency or time, but do possess positive or negative local spectrotemporal correlations in intensity. In psychophysical experiments, we found clear evidence that humans judge pitch direction based on both positive and negative spectrotemporal correlations. The observed sensitivity to negative correlations is a direct analogue of illusory “reverse-phi” motion in vision, and thus constitutes a new auditory illusion. Our behavioral results and computational modeling led us to hypothesize that human auditory processing employs pitch direction opponency. fMRI measurements in auditory cortex supported this hypothesis. To link our psychophysical findings to real-world pitch perception, we analyzed recordings of English and Mandarin speech and discovered that pitch direction was robustly signaled by the same positive and negative spectrotemporal correlations used in our psychophysical tests, suggesting that sensitivity to both positive and negative correlations confers ecological benefits. Overall, this work reveals that motion detection algorithms sensitive to local correlations are deployed by the central nervous system across disparate modalities (vision and audition) and dimensions (space and frequency).

1 Introduction

2

3 From discriminating phonemes to being moved by Bach's *Partitas*, detecting changes in pitch
4 over time, or pitch motion, is fundamental to human audition. Indeed, in everyday speech we use
5 both intonation and lexical tones — including complex rising and falling pitches — to signify
6 meaning (1-3). In English, for instance, rising pitch at the end of a sentence signifies a question.
7 In Mandarin Chinese, changes of pitch within words conveys fundamental differences in
8 meaning. But how does the human auditory system detect changes in pitch?

9

10 Changes in pitch can, in principle, be detected in at least two ways. First, listeners could identify
11 an auditory “object” corresponding to the pattern of frequencies made by a voice or any other
12 sound source (e.g., a friend's speech, a violin, etc.). If at the next instant in time the object
13 moved to higher frequencies, listeners would infer a rising pitch, or the opposite if the object
14 moved to lower frequencies (**Fig. 1A**). By identifying and tracking auditory objects, listeners can
15 perceive changes in the object's pitch over time. In vision, humans use this “feature tracking”
16 approach as one mechanism for detecting motion (4).

17

18 An alternative method for detecting changes in pitch would be to compute local correlations in
19 sound volume over time at nearby frequencies. These local correlations would enable listeners to
20 infer whether pitches are rising or falling without the added burden of first identifying auditory
21 objects. Methods like these are the basis of canonical models for spatial motion detection in
22 vision (5, 6). They can be dramatically revealed by visual illusions involving negative
23 correlations, including “reverse phi” phenomena (5-7). Thus, at least in vision, humans use both
24 object-tracking and intensity correlations to detect motion in the environment (8, 9).

25

26 Object tracking is a plausible method for detecting changes in pitch. Humans are clearly adept at
27 identifying and tracking auditory objects: In the well-known “cocktail party” effect, guests at a
28 noisy party can pick out and track a single voice in a sea of other voices (10-12). More generally,
29 listeners can group nearby frequencies into auditory objects, which strongly influences the
30 perception of rising and falling pitch (13). Likewise, the perception of continuity with rising and
31 falling tones is also consistent with tracking auditory objects (14), and psychophysical studies of
32 frequency change detection have tended to use isolated frequencies or persistent sound spectra in
33 which auditory object tracking is possible (15-18). Studies also show that pitch change
34 discrimination can occur over seconds, suggestive of object tracking (19).

35

36 What are the neural correlates of detecting rising and falling tones? Neurophysiological studies
37 have shown that both subcortical neurons (20, 21) and cortical neurons (22, 23), including in
38 primates (24), respond selectively to rising or falling tones in a narrow range of frequencies.
39 They achieve this selectivity by nonlinearly combining different frequency inputs at different
40 delays. Moreover, studies of many cortical auditory neurons have characterized complex
41 spectrotemporal receptive fields, which show how responses depend on different frequencies
42 over time (25, 26). Thus, although neural responses to auditory stimuli with local
43 spectrotemporal correlations have not been measured to date, neurons with appropriate
44 spectrotemporal tuning could detect such correlations. Neurons that detect rising or falling tones
45 could in principle support algorithms that detect pitch motion by object tracking but could also,
46 crucially, support those that work by sensing spectrotemporal correlations. It thus remains

47 unclear whether the human auditory system can use spectrotemporal correlations to perceive
48 directed changes in pitch. In this study, we hypothesize that detecting local spectrotemporal
49 correlations is a fundamental computation of the human auditory system.

50

51 **Results**

52

53 *Spectral motion without features*

54

55 We set out to test whether humans can detect auditory motion based on local spectrotemporal
56 correlations. To do this, we adapted a stimulus used to study visual motion detection (27, 28) to
57 develop new correlated noise auditory stimuli that use increments and decrements in volume to
58 generate local correlations in volume at specific offsets in frequency and time (see **Methods**).

59 We designed four stimuli with positive or negative correlations in volume at an offset of 1/6
60 second, with the frequency directed either upward or downward by 1/15 octave (**Fig. 1B, S1**).

61 These sounds were inharmonic, so that fundamental frequencies could not be used to judge pitch
62 changes (29, 30). We presented these stimuli to participants for 2 seconds and asked them to
63 report whether they perceived the sound as having a rising or falling pitch profile over time.

64

65 Participants reported that upward-directed positive correlations rose in pitch over time, while
66 those with downward directed positive correlations fell in pitch over time (**Fig. 1C, Supp. Movie**
67 **1**). This psychophysical result demonstrates that humans can identify rising or falling pitch based
68 on local correlations alone, without persistent auditory objects.

69

70 Remarkably, when we presented stimuli with negative correlations in frequency and time,
71 participants reported the opposite percepts (**Fig. 1C, Supp. Movies 1 and 2**). That is, the
72 upward-directed negative correlations sounded like they were falling in pitch, while the
73 downward directed negative correlations sounded like they were rising in pitch. Participants who
74 consistently perceived rising or falling pitch in the stimuli with positive correlations also
75 consistently perceived rising or falling pitch in the stimuli with negative correlations (**Figure**
76 **S1**). This striking illusion demonstrates that humans are sensitive not just to positive
77 spectrotemporal correlations, but to negative ones as well. This result is a direct analog to
78 illusory reverse-phi visual motion percepts, which have been reported across many species and
79 phyla (5, 7, 31-33).

80

81 How does the strength of these spectrotemporal correlations relate to perception? To answer this
82 question, we varied the coherence of the stimulus and again asked participants to judge whether
83 tones were rising or falling in pitch (**Fig. 1D**). We titrated the coherence of the stimuli from 1 to
84 0 by randomly replacing correlated time-frequency elements with random ones, such that the
85 coherence represented the fraction of original correlations remaining (see **Methods**). With high
86 coherence, participants perceived rising and falling pitches in a pattern similar to the first
87 experiment (**Fig. 1C**). As coherence decreased, however, the probability of judging a sound as
88 rising tended towards chance (0.5). There were no significant differences between the curves for
89 ($\uparrow +$) and ($\downarrow -$) or ($\downarrow +$) and ($\uparrow -$) ($p > 0.05$ for each, as measured by a two-way, repeated
90 measures ANOVA), meaning that inverting the stimulus correlation and direction led to
91 indistinguishable percepts. These results reveal a clear monotonic relationship between the

92 strength of spectrotemporal correlations and the strength of pitch change percepts, both for
93 positive and for negative correlations.

94
95 In vision, object tracking can integrate information between the two eyes, while correlation
96 based algorithms rely on correlations within each eye (9). We next asked if spectrotemporal
97 correlations for pitch motion detection are computed monaurally or binaurally. The structure of
98 our correlated noise stimulus is created by summing a random binary mask with itself at a
99 frequency-time offset (**Fig. 1E, Methods**). This allowed us to play one binary mask to the left
100 ear and a shifted one to the right ear, so that neither ear alone would be presented with any
101 correlations. In this context, detecting spectrotemporal correlations can only proceed by
102 integrating information across the two ears. We played all four types of binaural correlations to
103 participants and asked them to judge whether they heard rising or falling sounds. They reported
104 the same pattern of percepts as in the monaural stimuli, though with average reported directions
105 somewhat closer to chance. This demonstrates that the perception of rising or falling pitch can
106 use information from both ears to integrate volume information to compute spectrotemporal
107 correlations. This is consistent with data showing that many cortical auditory neurons integrate
108 signals from both ears (34).

109
110 *Tuning of human spectrotemporal correlation detectors*

111
112 Our next step was to characterize the spectral and temporal tuning of the correlation sensitivity
113 we had observed. To do this, we designed a different kind of stimulus, one inspired by random
114 dot kinetograms in visual neuroscience (35). In these stimuli, a medium intensity sound that
115 played at all frequencies was interrupted by brief pips at different frequencies, 50 ms in duration
116 (24). These pips either increased the volume of a specific frequency or decreased it to zero (**Fig.**
117 **2A**, see **Methods**). After an initial set of pips were placed randomly in frequency and time, we
118 added a second set of pips with a specific delay in time and change in frequency, yielding
119 correlated pip pairs. These pairs had positive correlations when both pips were loud or both were
120 silent, and negative correlations when one was loud and one was silent. This allowed us to create
121 auditory stimuli with upward and downward-directed pairs of pips with positive or negative
122 correlations (**Fig. 2B**). Like the stimuli used in **Fig. 1**, these stimuli had no auditory objects that
123 persisted in time or frequency, but crucially, they allowed us to vary the delay continuously
124 between correlated pips.

125
126 We first used these stimuli to map out the sensitivity to different delays between individuated
127 tones. We kept the frequency change at 1/15 octave and swept values of the delay between
128 correlated pips while asking participants to judge whether the pitch was rising or falling over
129 time (**Fig. 2C**). For both negative and positive correlations and upward and downward-directed
130 displacements, we found that peak directional sensitivity occurred at a delay of around 40 ms.
131 This peak did not change appreciably when the pip duration was shortened to 20 ms (**Fig. S2**).
132 According to models for visual motion estimation, this peak sensitivity value reflects the typical
133 relative delays in the circuits detecting local motion signals (27). The delay seen here is on a
134 similar timescale, though is slightly longer than delays measured by similar experiments in
135 human and fly visual systems (27, 36).

136

137 We then measured sensitivity to the magnitude of displacements in frequency space. Using a
138 similar method, we set the delay to 40 ms and varied the frequency displacement within a pip
139 pair (**Fig. 2D**). We found that peak sensitivity occurred for tone displacements of 1/15th octave,
140 though there was still significant direction-selectivity at 2/15th octave displacements ($p < 0.05$
141 for both positive and negative correlations by a paired t-test). This result shows that correlation-
142 based motion detectors in the human auditory system are most sensitive to small shifts in
143 frequency in the vicinity of 1/15th of an octave (4.7% changes in frequency) or less. This result
144 is consistent with peak sensitivity for changes in complex sounds (37) and with the smaller
145 values of frequency discrimination thresholds in humans (15).

146 147 *Sensitivity to spectrotemporal volume patterns*

148
149 Our positive and negative correlation stimuli each consist of multiple patterns in volume over
150 frequency and time. Upward-directed positive correlation ($\uparrow +$) stimuli consist of both loud-loud
151 and soft-soft combinations, whereas the negative versions ($\uparrow -$) consist of both loud-soft and
152 soft-loud combinations. Prior work using long-lasting spectrotemporal correlations in auditory
153 stimuli has suggested that humans are selectively sensitive to loud-loud combinations (38). Are
154 humans sensitive to all four pairwise combinations, or to just a subset of them? To address this
155 question, we generated new correlated pip auditory stimuli (**Fig. 3A**) where each stimulus had
156 paired pips of only one of the four types: loud-loud, soft-soft, loud-soft, or soft-loud (see
157 **Methods**). We asked participants to judge whether these different stimuli were rising or falling
158 and recorded their responses (**Fig. 3B**). Participants were sensitive to all four different pairings
159 with both upward and downward displacements.

160
161 In visual motion detection, one generalization beyond pairwise correlations involves so-called
162 triplet correlations (39, 40). In vision, triplet correlations are patterns that contain spatiotemporal
163 correlations over three points in space and time, but no pairwise correlations, and can elicit
164 visual motion percepts in humans (39, 41), flies (41, 42), and fish (43). Visual motion detection
165 algorithms are sensitive to this higher-order relative structure, but is the same true in
166 audition? When participants were presented with auditory analogs of visual triplet correlation
167 stimuli (see **Methods**), they did indeed perceived auditory motion (**Figure S3**) and did so in a
168 pattern much like that found in fly and fish visual perception. This correspondence across both
169 species and modalities points to significant similarities in the neural algorithms used by animals
170 in processing auditory and visual motion.

171 172 *Psychophysical and cortical signatures of opponent subtraction of spectral motion signals*

173
174 When we presented positively and negatively correlated stimuli, we observed a striking
175 symmetry: Tuning of negative correlation percepts matched the tuning of positive correlation
176 percepts, but in the opposite direction (**Fig. 2**). This clear symmetry is highly suggestive of an
177 opponent architecture. To investigate this, we first built a simple motion energy model unit to
178 describe a hypothetical directionally tuned auditory neuron (**Fig. 4A**). The model unit filtered
179 sound intensity linearly over frequency and time in a pattern that enhanced upward-directed
180 spectral motion, similar to prior suggestions (44), before sending the signal through a quadratic
181 nonlinearity (6). When we presented this model with correlated pip stimuli (**Fig. 2**), it responded
182 at an elevated baseline level but with deviations that depended on the direction and sign of the

183 stimulus correlation (**Fig. 4B**). As designed, it responded more to upward-directed positive
184 correlations than to downward-directed ones. Since this model relies solely on pairwise
185 correlations, it was also expected that negative correlation stimuli elicited equal and opposite
186 deviations to positive correlation stimuli. Crucially, however, in this model, negatively correlated
187 stimuli exhibit a different tuning from oppositely directed positive stimuli; that is, inverting the
188 correlation is not equivalent to inverting the direction (i.e., the temporal delay).

189
190 We next created an opponent signal by subtracting signals from two model units with opposite
191 directional tuning (**Fig. 4C**). This opponent signal responded to positively correlated stimuli with
192 positive and negative values when they were directed upward and downward (**Fig. 4C, green**).
193 Critically, this opponent signal has an important symmetry: Responses to negatively correlated
194 stimuli have the same tuning as positively correlated stimuli in the opposite direction. Thus,
195 upward-directed negative correlation stimuli yield the same responses as downward-directed
196 positive correlation stimuli. We also derived this result analytically (see **Methods**): When
197 motion energy signals are opponently subtracted, negative correlation stimuli elicit mean
198 responses that match oppositely directed positive correlation stimuli.

199
200 To demonstrate that our data contained this symmetry, we compared percepts of negative
201 correlation stimuli to percepts of positive correlation stimuli in the opposite direction, for both
202 frequency change and delay time tuning (**Fig. 4D, E**, replotting data from **Fig. 2**). The curves
203 appeared to fully superimpose. ANOVA tests confirmed that there was no measurable difference
204 between the positive correlation curves and the flipped negative correlation curves (see figure
205 legends for statistics). This robust symmetry between positive and negative correlation stimuli
206 has also been found in visual motion detection in fruit flies (27) and in humans (36).

207
208 In primate vision, opponent subtraction occurs in visual area V5, also called MT (45, 46), which
209 has been shown to be causally involved in visual motion percepts (47). Similarly, flies also
210 subtract visual motion signals with opposing preferred directions (48). Motivated by our
211 psychophysical results, analogies with vision, proposals for opponent subtraction to determine
212 spectral direction (16), and by spectral direction opponent auditory cells found in bats (49), we
213 reasoned that human auditory cortex might possess signatures of opponent processing.

214
215 We followed the logic of previous functional magnetic resonance imaging (fMRI) studies that
216 identified opponent signals in human cortical area MT and used visual stimuli that summed
217 motion in opposite directions (50). To start, we assume that cortical voxels involved in detecting
218 spectral motion contain units that respond preferentially to rising tones and units that respond
219 preferentially to falling tones, but none that respond to both (**Fig. 4F**) (51). Such a voxel should
220 thus respond reliably to stimuli containing either rising or falling tones. The key distinction
221 between a system with or without opponency lies in its response to a summed stimulus that
222 contains superimposed rising and falling tones: If units are opponent, then the summed stimulus
223 should cause a decrease in voxel activity due to a net suppression of signals in units with
224 opponent responses (50). We therefore designed simple stimuli consisting of rising tones, falling
225 tones, or their sum (**Fig. 4G, S4**) and presented them to subjects while measuring blood-oxygen-
226 level-dependent (BOLD) signals via fMRI.

227

228 We searched within a broad auditory cortex mask for voxels that responded more to the non-
229 summed (rising or falling) stimuli than to the summed (opponent) stimulus (see **Methods**).
230 Strikingly, at both group and individual levels, a bilateral region within superior temporal cortex
231 was significantly more activated by the non-summed stimuli than by the summed stimulus (**Fig.**
232 **4H, I**), consistent with opponency. The group map extended over multiple bilateral functional
233 subregions of the human auditory cortex (52), including core regions A1 and RI, Area 52, and
234 lateral and medial belt regions (**Fig. S4**). According to the opponency hypothesis, activity in
235 opponent voxels should be similar in magnitude for rising and falling stimuli and suppressed for
236 the summed stimulus. Thus, we wanted to ensure that our result followed this symmetry and was
237 not biased by either the rising or falling stimulus alone (see **Methods**). Activity in putative
238 opponent regions was indeed comparable for rising and falling tones (**Fig. 4J**). Overall, our
239 fMRI findings demonstrate that a key result from our behavioral studies — the clear symmetry
240 between positive and negative correlation percepts — lead to a specific neural hypothesis that
241 was borne out in neuroimaging data. To our knowledge, this general region of human auditory
242 cortex has not previously been identified as a potential locus for opponent spectral motion
243 signals.

244
245 *Positive and negative correlation spectrotemporal cues signal tone modulation in speech*

246
247 Is there an ecological advantage in detecting both positive and negative spectrotemporal
248 correlations? To address this question, we chose to look at human speech, where tone modulation
249 contains critical semantic information in both tonal and non-tonal languages (1-3). Since humans
250 are sensitive to both positive and negative pairwise correlations in frequency and time, we
251 hypothesized that these correlations could convey information about the direction and speed of
252 tone modulation in human speech. Following in the tradition of relating auditory processing to
253 natural sounds (53), we analyzed corpora of spoken English and Mandarin and examined how
254 tone modulation is related to underlying positive and negative pairwise spectrotemporal
255 correlations in volume (**Fig. 5, Methods**).

256
257 Our analysis took several steps. First, we computed spectrograms for each of the speech
258 recordings (**Fig. 5A, top**). We then used an optical flow algorithm to estimate the change in tone
259 at each point in time – that is, the degree to which the sound was rising or falling in frequency at
260 each time (**Fig. 5A, bottom, see Methods**). Next, we binarized the spectrogram and looked for
261 specific patterns of volume in frequency and time, examining all four combinations of loud and
262 soft: loud-loud, soft-soft, loud-soft, and soft-loud (**Fig. 5B**). We next computed the local net
263 signal for each pattern at each frequency and time by subtracting the downward directed patterns
264 from the upward directed ones (**Figs. 5C, D**). Finally, we averaged these local net signals over all
265 frequencies to obtain a net pattern signal (**Figs. 5C, D**). Computing net pattern signals is
266 consistent with the opponency we observed psychophysically and in fMRI (**Fig. 4**). For the loud-
267 loud patterns, there was a positive correlation between the time trace of the net pattern signal and
268 the tone change. For the loud-soft patterns, the correlation was negative. The clear suggestion is
269 that negative correlations contain information about tone changes that could be useful to listeners
270 in detecting rising and falling tones in speech.

271
272 To see whether this result generalized, we analyzed hundreds of speech snippets that totaled over
273 90 minutes in English and 40 minutes in Mandarin Chinese (**Fig. 5E, F, see Methods**). In

274 English, the tone changes should be dominated by intonation, while in Mandarin Chinese, the
275 tone changes should reflect both intonational and within-syllable changes in tone (1-3). We
276 reproduced the analysis of the different volume patterns, and then correlated the net signal for
277 each pattern with the computed change in tone. In both English and Mandarin Chinese, the two
278 positive correlation patterns (loud-loud and soft-soft) produced a strong positive correlation ($r >$
279 0.5) with the intonation velocity, whereas the two negative patterns (loud-soft and soft-loud)
280 produced a strong negative correlation ($r < -0.5$) (**Fig. 5E, F**). These results show that all four
281 patterns could be useful in estimating tone changes in speech. The negative stimulus correlation
282 produced an *anti*-correlation with tone changes, which explains why they elicit percepts in the
283 opposite direction: upward directed negative correlations indicate downward directed tone
284 changes. We obtained similar results when we processed with the spectrograms with continuous
285 rather than digital operations to obtain positive and negative spectrotemporal correlations in the
286 speech data (see **Methods, Fig. S5**). Thus, this analysis provides an ecological explanation of the
287 observed inverted percepts to negative auditory correlations.

288

289 **Discussion**

290

291 In the studies reported here, we have demonstrated that humans are sensitive to local
292 spectrotemporal correlations in volume over frequency and time as they discern whether a sound
293 is rising or falling in pitch (**Figs. 1-3**). Participants were equally sensitive to both negative and
294 positive spectrotemporal correlations, a pattern that mirrors a powerful visual phenomenon, the
295 reverse-phi illusion, in a different modality (audition) and over a different dimension of motion
296 (frequency). Inspired by our behavioral results showing symmetry between inverting correlation
297 and inverting direction, we hypothesized that the human auditory system might implement
298 opponent subtraction, echoing a similar operation in visual motion detection. Using fMRI, we
299 discovered that, like visual cortex, regions within human auditory cortex show signatures of
300 opponency (**Fig. 4**). Finally, we demonstrated that negative spectrotemporal correlations likely
301 act as reliable cues to assess tone changes in speech (**Fig. 5**).

302

303 The stimuli we developed here (**Figs. 1-3**) in some ways resemble Shepard tones (54), which
304 were designed to sound like they are unceasingly rising or falling. However, Shepard tones
305 consist of periodic auditory features that persist over frequency and time (similar to **Fig. 4F**).
306 Thus, the rising or falling of a Shepard tone could be assessed by simply tracking auditory
307 features over time. The auditory stimuli we developed and investigated here, however, have no
308 such persistent features — a rising or falling percept must instead depend on the detection of
309 positive and negative pairwise spectrotemporal correlations within the stimulus. Thus, the strong
310 percepts of rising and falling tones, which depended on the sign of the correlation, reflect an
311 authentic auditory illusion in which there is no true rising or falling tone but only the imposition
312 of specific spectrotemporal correlations in volume.

313

314 Sensitivity to spectrotemporal correlations in judging pitch direction likely acts in coordination
315 with other algorithms for judging changes in pitch. In particular, changes in frequency can be
316 judged over gaps of seconds (37), which points to a different system for such judgements.
317 Similarly, judgements about relative pitch can be made using fundamental frequencies in
318 harmonic sounds (29, 55). These examples suggest that auditory spectral motion processing is

319 similar to visual motion processing, where positional changes can be detected by both local
320 correlational algorithms and by slower, longer range object-tracking algorithms (8).

321
322 We found regions in both primary and non-primary auditory cortex across both Heschl's gyrus
323 and the superior temporal gyrus (STG) that may perform opponent computations to resolve net
324 pitch direction (**Fig. 4H, I**). How might this relate to the neural underpinnings of speech
325 perception? Human auditory cortex displays regional specialization, with areas that selectively
326 encode different aspects of speech, primarily in the STG (56-59). Our results are broadly
327 consistent with findings that regions within STG encode variability in speaker intonation and
328 lexical tone (59, 60). Moreover, our observation that significant portions of Heschel's gyrus also
329 showed spectral motion sensitivity is broadly consistent with other work (61), though we saw the
330 effects bilaterally (**Fig. S4**) and, critically, with an opponent signature. Our results thus suggest
331 that opponency may be a signature of pitch direction processing in circuits involved in simple
332 pitch computations (in primary areas) and in more complex perceptual tasks like speech
333 processing (in non-primary areas).

334
335 Canonical algorithmic models for motion detection are sensitive to negative correlations (5, 6),
336 and more neurophysiologically-inspired models for motion detection are similarly sensitive to
337 negative correlations (62, 63). At the single neuron level, units in rodent (22), bat (23), and
338 primate (24) auditory cortex display spectrotemporally oriented receptive fields, which should
339 confer sensitivity to both positive and negative spectrotemporal correlations (**Fig. 4**) (6). Our
340 results suggest that neurons with this type of sensitivity could underlie spectrotemporal
341 correlation detection in humans. Meanwhile, our psychophysical and fMRI results also suggest
342 that units in multiple regions of auditory cortex exhibit directional opponency, a property
343 observed in bat auditory neurons (49). This direction opponency could arise in primary motion
344 detectors (64), or be the result of subtracting opposing cortical or subcortical motion signals (21).

345
346 There are well-established similarities in the processing of visual motion between invertebrates
347 and vertebrates (65-67), phyla that diverged hundreds of millions of years ago. Our study shows
348 that local correlational algorithms for motion detection also span modalities, since human
349 audition and vision appear to employ similar computational motifs. Audition thus joins olfaction
350 (68) as a non-visual sense where pairwise, local correlations can generate rich motion percepts.
351 In these experiments, sensitivity to pairwise stimulus correlations also includes sensitivity to
352 negative correlations. This sensitivity to negative correlations is due in part to the mathematics of
353 computing correlations (see **Methods**) (6, 40), providing a conceptual framework for
354 understanding the neural detection of motion that spans modality and species.

355
356 Lastly, negative correlations sensed in audition likely act as useful cues to infer real-world
357 changes in the frequency domain (**Fig. 5**), just as they may help in visual motion detection (69,
358 70). Thus, the illusory pitch motion described here is not just an interesting laboratory
359 epiphenomenon. Rather, it reflects neural sensitivity to the statistics of the auditory world, with
360 direct implications for everyday speech and music perception.

361 **Contributions**

362 PAV and DAC designed auditory stimuli. PAV and SDM acquired data. PAV, SDM, and DAC
363 analyzed and interpreted data. PAV, SDM, and DAC wrote the paper.

365

366 **Acknowledgements**

367 This work was funded by a grant from the Wu Tsai Institute at Yale University. DAC and this

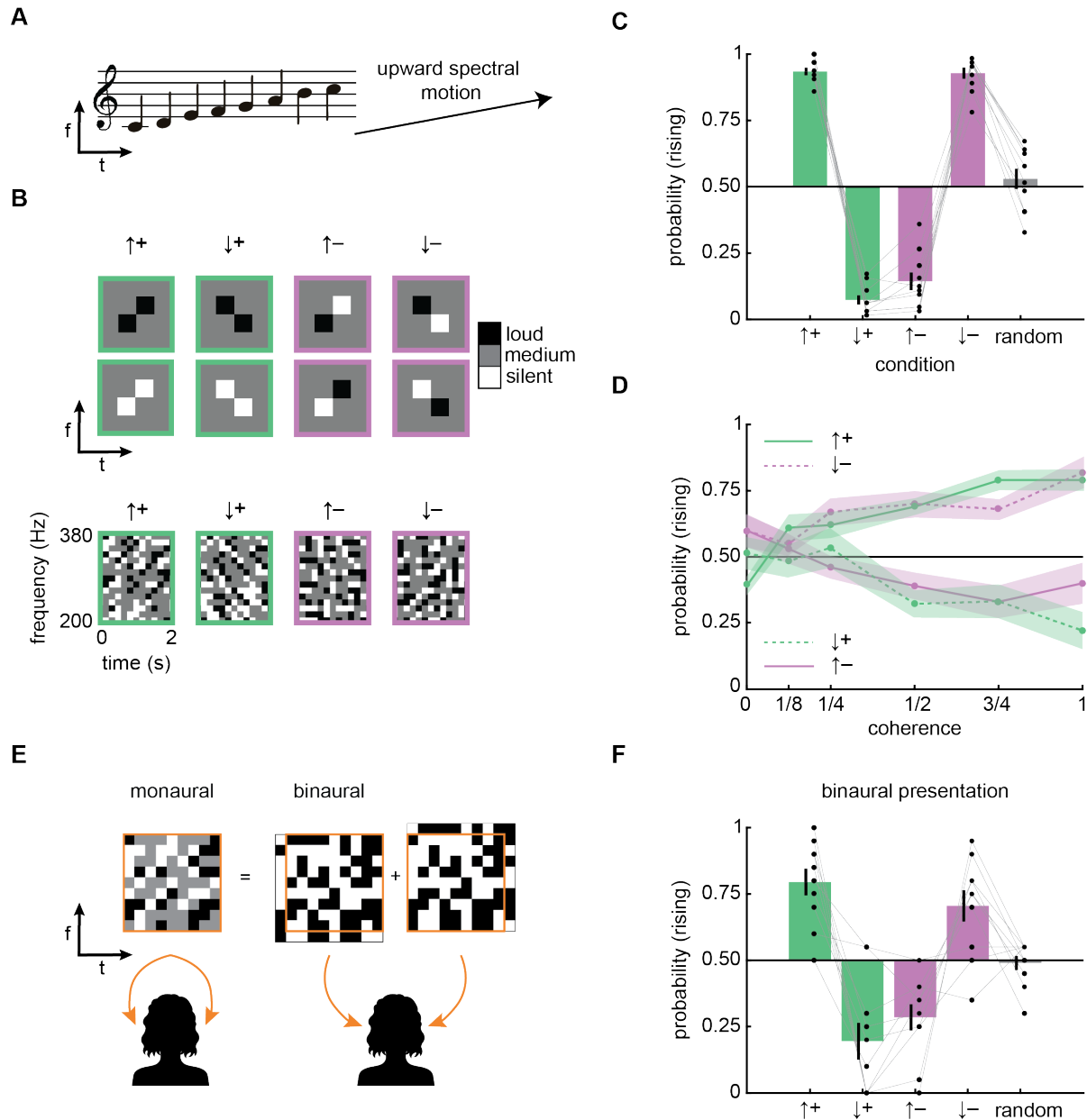
368 work were funded by NIH R01 EY026555. We thank R. Aslin, E. V. Clark, H. H. Clark, I.

369 Yildirim, and J. Zavatone-Veth for helpful discussions and comments on this project.

370

371

372 **Figures**
373

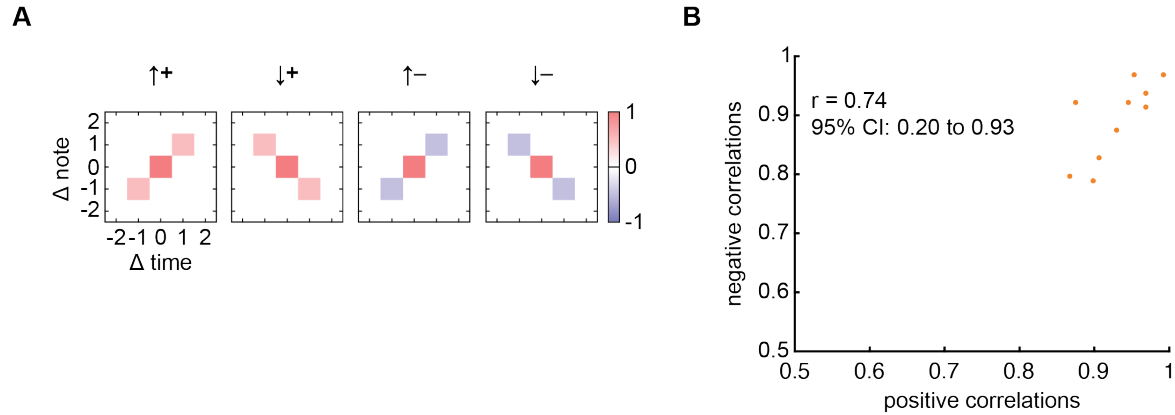


374
375
376
377
378
379
380
381
382
383
384

Figure 1. Humans detect auditory motion in pairwise frequency-time correlations.

- A) Simple schematic of a rising sound written on a music staff and in frequency-time.
- B) Diagrams showing sample (*top*) and actual (*bottom*) stimuli. Frequency-time correlations can be directed either upward or downward and be either positively or negatively correlated.
- C) Perceived direction of stimuli with varying direction and correlation. Mean \pm SEM over $N=10$ subjects. One-sample t-tests revealed significant deviations from chance (0.50) in pitch direction judgements in all four stimulus conditions (all $ps < 10^{-5}$). Pitch direction judgements in the random stimulus condition were not significantly different from chance ($p = 0.45$). Error bars represent mean \pm SEM ($N = 10$).

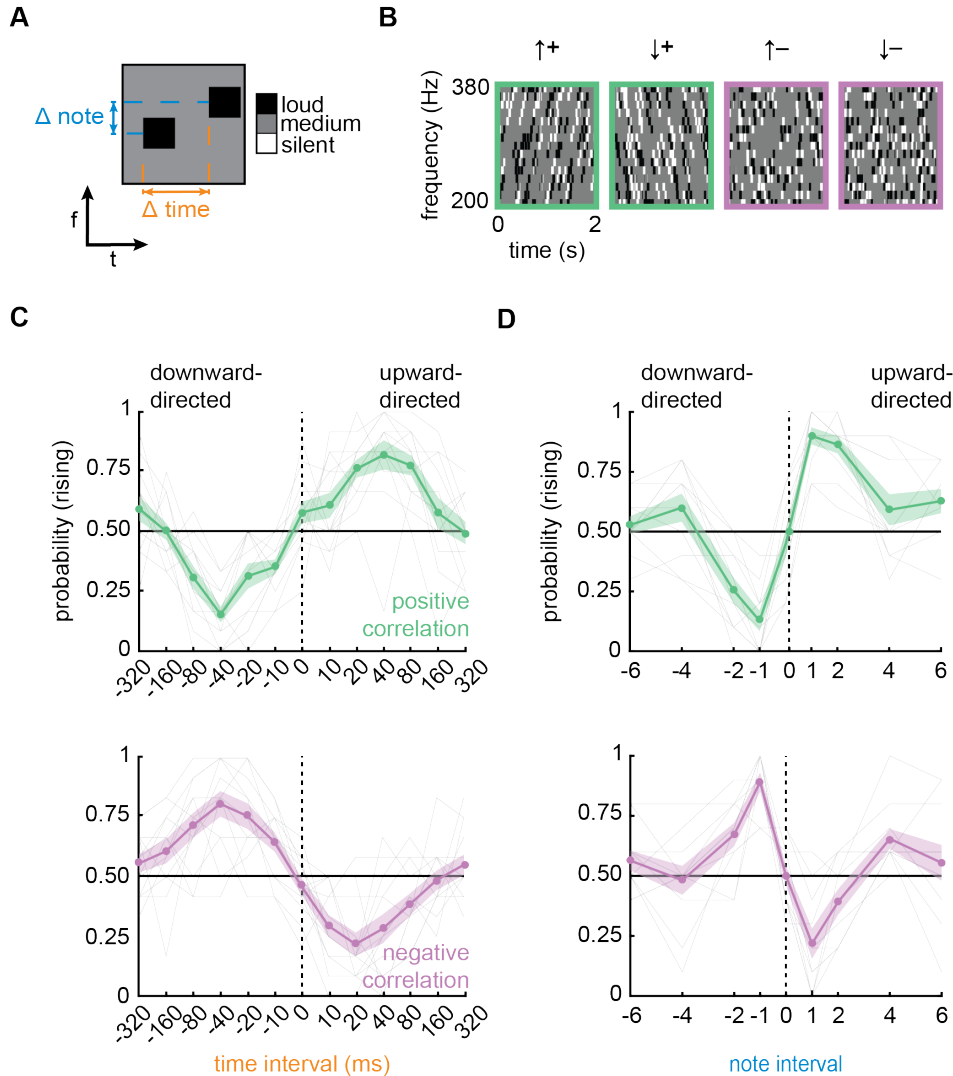
- 385 D) Perceived direction of stimuli with varying degrees of correlation (coherence) in the
386 stimulus. The upward directed positive and downward directed negative curves were not
387 significantly different ($p > 0.05$ by a two-way, repeated measures ANOVA); similarly,
388 the downward directed positive and upward directed negative curves were also not
389 significantly different ($p > 0.05$, same test). Both ANOVAs revealed significant main
390 effects of coherence on pitch direction judgements (all $ps < 10^{-5}$). Error shading
391 represents \pm SEM (N = 10).
- 392 E) Diagram showing how binaural stimuli were presented to each ear.
- 393 F) Perceived direction of stimuli with varying directions and correlations using binaural
394 presentation. One-sample t-tests revealed significant deviations from chance (0.50) in
395 pitch direction judgements in all four stimulus conditions (all $ps < 10^{-3}$). Pitch direction
396 judgements in the random stimulus condition were not significantly different from chance
397 ($p = 0.72$). Error bars represent mean \pm SEM (N = 10).
- 398



399
400
401
402
403
404
405
406
407
408
409
410
411
412

Supplemental Figure S1.

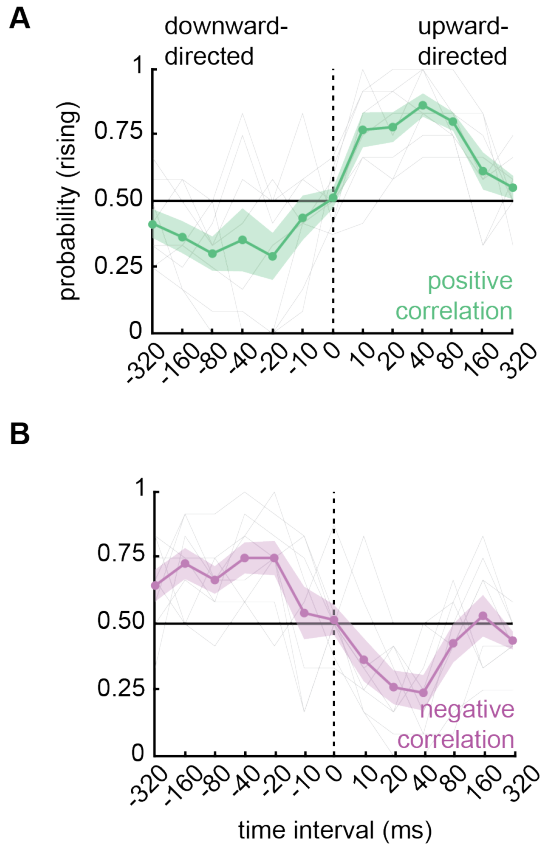
- A) Stimulus autocorrelation plots at different note and time offsets for the stimuli in **Figure 1B**. The stimuli have positive or negative correlations at a single spectrotemporal offset, directed either upward or downward in frequency over time. These plots are normalized so that the origin has correlation of 1.
- B) Correlation between perception of positively correlated and negatively correlated stimuli. To obtain the positive correlation values, we averaged $P(\text{rising})$ for the upward directed, positive correlation stimuli with $1-P(\text{rising})$ for the downward directed, positive correlation stimuli. To obtain the negative correlation values, we averaged $P(\text{rising})$ for the downward directed, negative correlation stimuli with $1-P(\text{rising})$ for the upward directed, negative correlation stimuli. Correlation coefficient is the Pearson correlation, and a 95% confidence interval is noted.



413

414 **Figure 2.** Correlation detection is tuned to small frequency changes and short delays in time.

- 415 A) Diagram showing a correlated pip pair with a frequency displacement (Δ note) and a
 416 delay between pips.
 417 B) Spectrotemporal diagrams of 4 different correlated pip stimuli directed upward and
 418 downward with positive and negative correlations. Pip duration in these experiments was
 419 50 ms.
 420 C) Perceived direction of stimuli with Δ note = +1 and varying pip delays; positive pip
 421 correlations (*top*) and negative pip correlations (*bottom*). One-way, repeated measures
 422 ANOVAs for the positive and negative correlation curves revealed significantly different
 423 responses across pip delays (all $ps < 10^{-21}$). Gray lines are individual participant curves.
 424 Error shading represents \pm SEM (N = 13).
 425 D) Perceived direction of stimuli using varying note intervals and 40 ms pip delays; positive
 426 pip correlations (*top*) and negative pip correlations (*bottom*). One-way, repeated measures
 427 ANOVAs for the positive and negative correlation curves revealed significantly different
 428 responses across note intervals (all $ps < 10^{-12}$). Gray lines are individual participant
 429 curves. Error shading represents \pm SEM (N = 13).



430

431 **Supp. Fig. 2.** Interval sweep with a different pip duration.

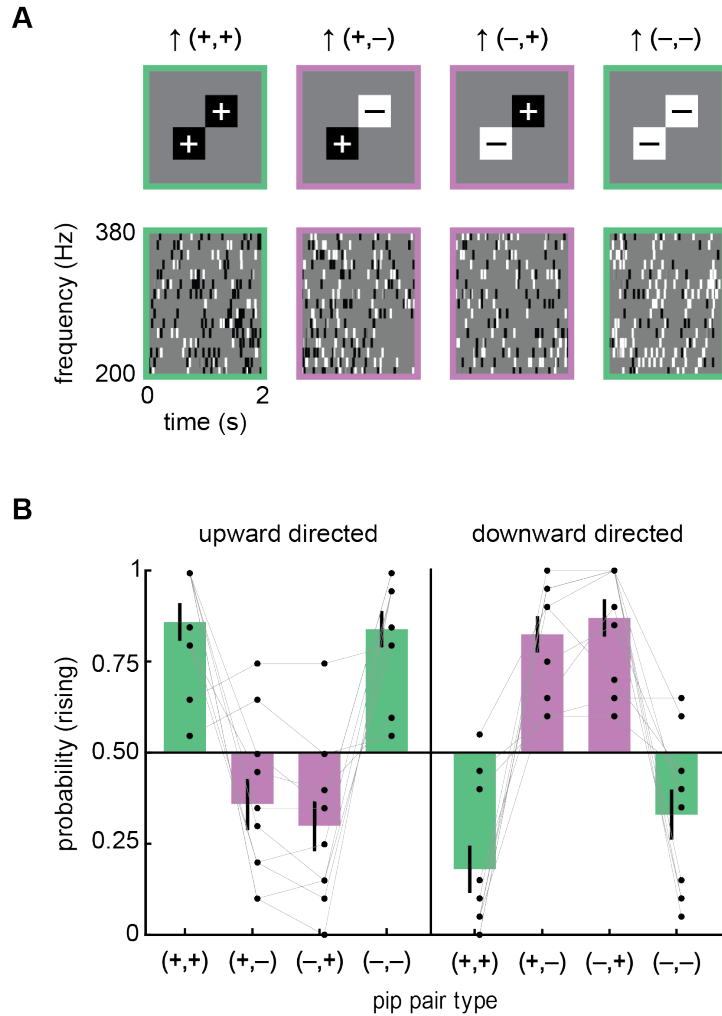
432

433 A) Perceived direction of positively correlated stimuli with varying pip delays and 20 ms
434 pips. Sensitivity tends to peak around 40 ms delays, similar to the data in **Figure 2C**. A
435 one-way, repeated measures ANOVA revealed significantly different responses across
436 pip delays ($p < 10^{-10}$). Error shading represents \pm SEM (N = 9).

436

437 B) Perceived direction of negatively correlated stimuli with varying pip delays and 20 ms
438 pips. Sensitivity tends to peak around 40 ms delays, similar to the data in **Figure 2C**. A
439 one-way, repeated measures ANOVA revealed significantly different responses across
440 pip delays ($p < 10^{-8}$). Error shading represents \pm SEM (N = 9).

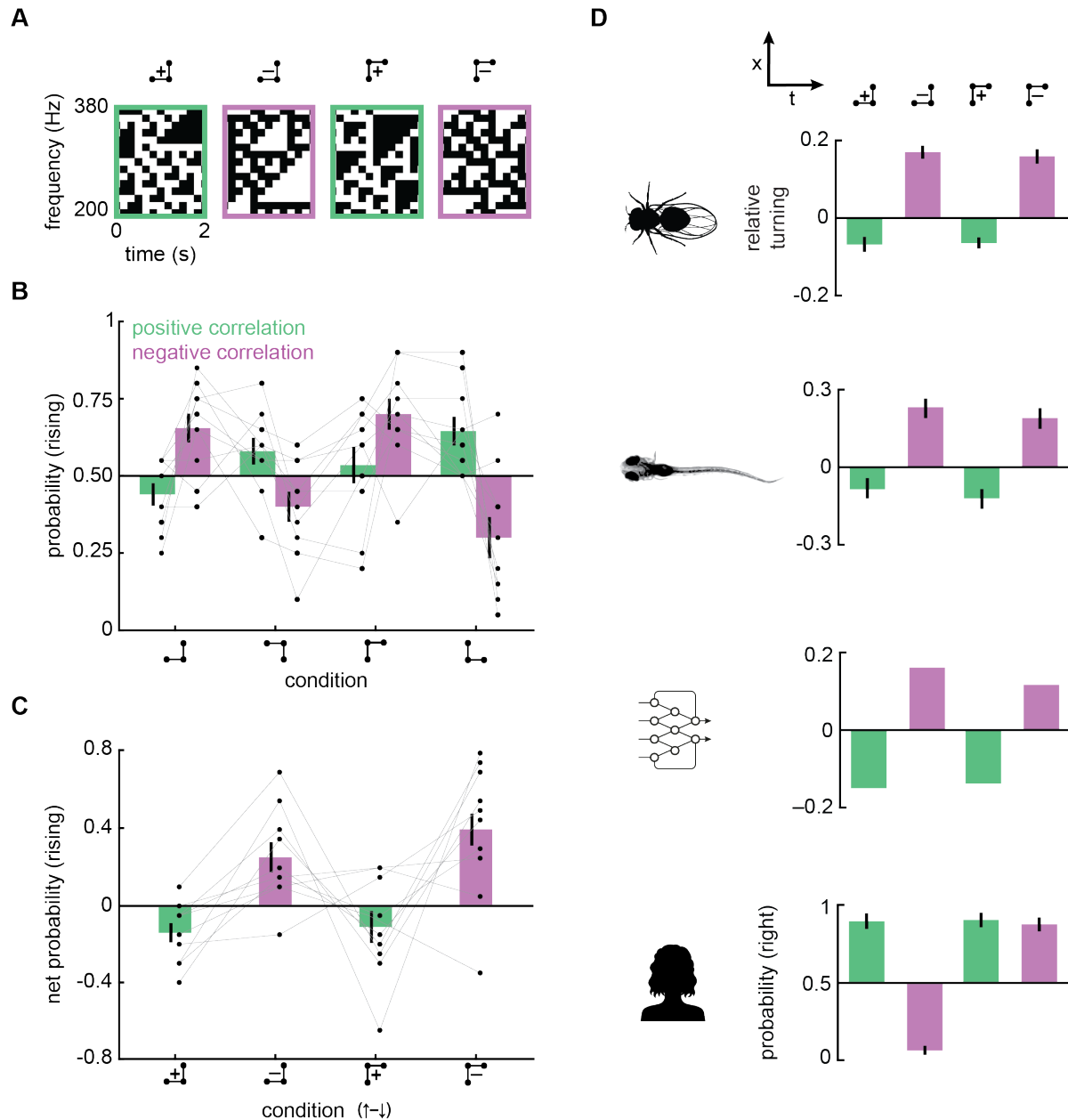
439



440
441
442
443
444
445
446
447

Figure 3. Sensitivity to all four pairwise loudness combinations contribute to rising and falling pitch perception.

- A) Frequency-time diagram of 4 different pip combinations, presented with 40 ms delays.
B) Probability of perceiving rising pitch for each of the four loudness combinations directed upward (*left*) and downward (*right*). Paired t-tests comparing upward- versus downward-directed stimuli for each matched pair revealed significant direction selectivity across all pitch direction judgements (all $ps < 10^{-4}$). Error bars represent mean \pm SEM (N = 10).



448

449 **Supplemental Figure S3.** Human auditory sensitivity to 3-point glider stimuli resembles visual
 450 sensitivity in different species.

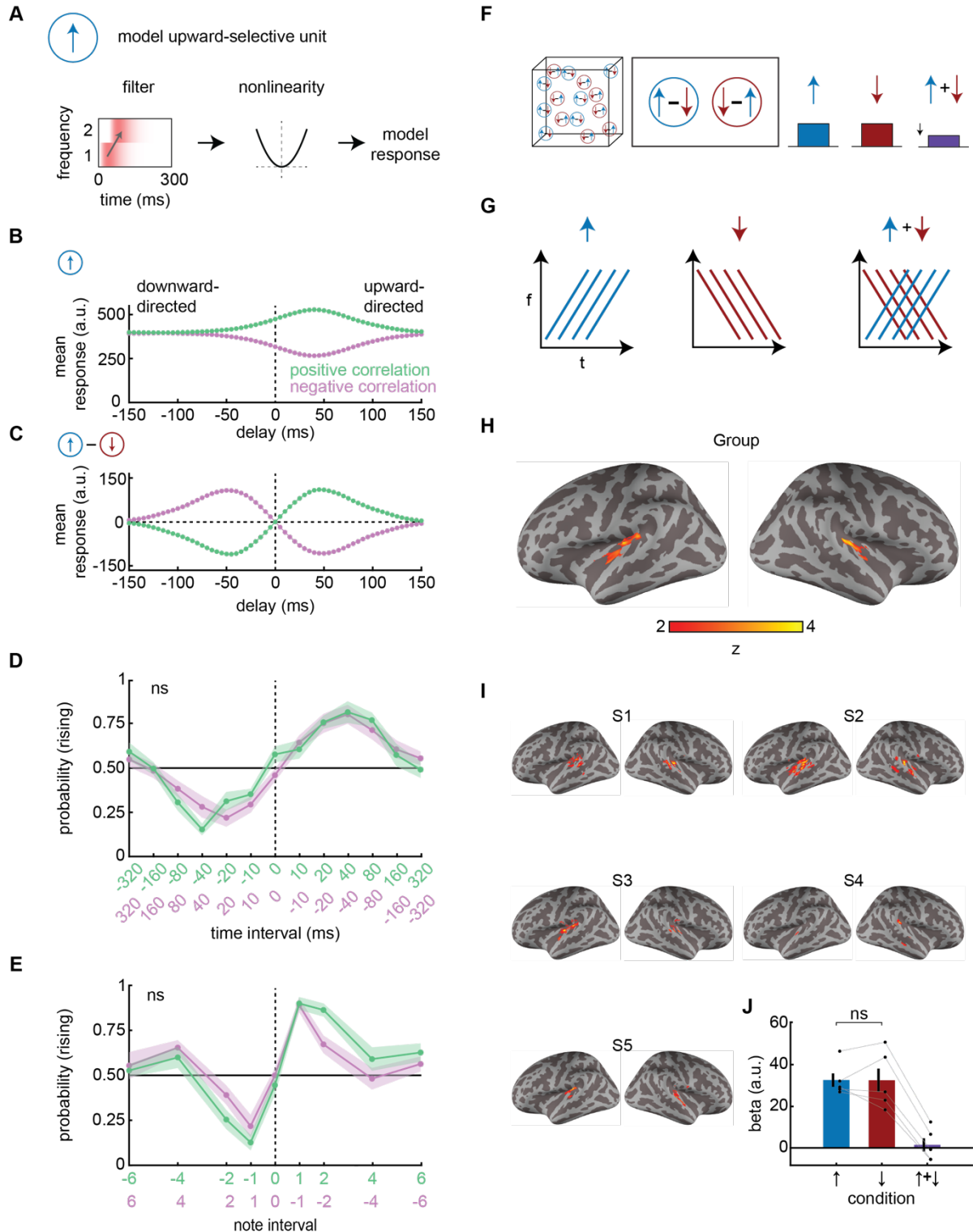
451 A) Diagram of 3-point glider stimuli in frequency and time (39). 3-point glider stimuli
 452 contain correlations between triplets of points as denoted by the barbell diagrams, and
 453 contain no pairwise correlations. Thus, motion percepts with these stimuli would have to
 454 rely on correlations beyond pairwise ones.

455 B) Perceived direction of 3-point glider stimuli. Participants heard rising and falling tones in
 456 these triplet correlation stimuli. Error bars represent mean \pm SEM (N = 10).

457 C) Net perceived direction of 3-point glider stimuli with positive and negative correlations.
 458 The net probability rising is computed by subtracting the downward directed P(rising)
 459 from the upward directed P(rising) in panel (B). Positively correlated stimuli were
 460 perceived as falling, while negatively correlated stimuli were perceived as rising. Paired

461 t-tests revealed significantly different responses to positively and negatively correlated
462 diverging gliders, and to positively and negatively correlated converging gliders (all $ps <$
463 10^{-3}). Error bars represent mean \pm SEM ($N = 10$).

464 D) Net perceived direction of 3-point glider stimuli across various visual systems. Data is
465 replotted from prior publications for fruit flies (41), larval zebrafish (43), a machine
466 learning algorithm (69), and human visual psychophysics (41). Human auditory percepts
467 resemble fruit fly and zebrafish visual percepts and machine learning responses, but not
468 the human visual percepts.
469

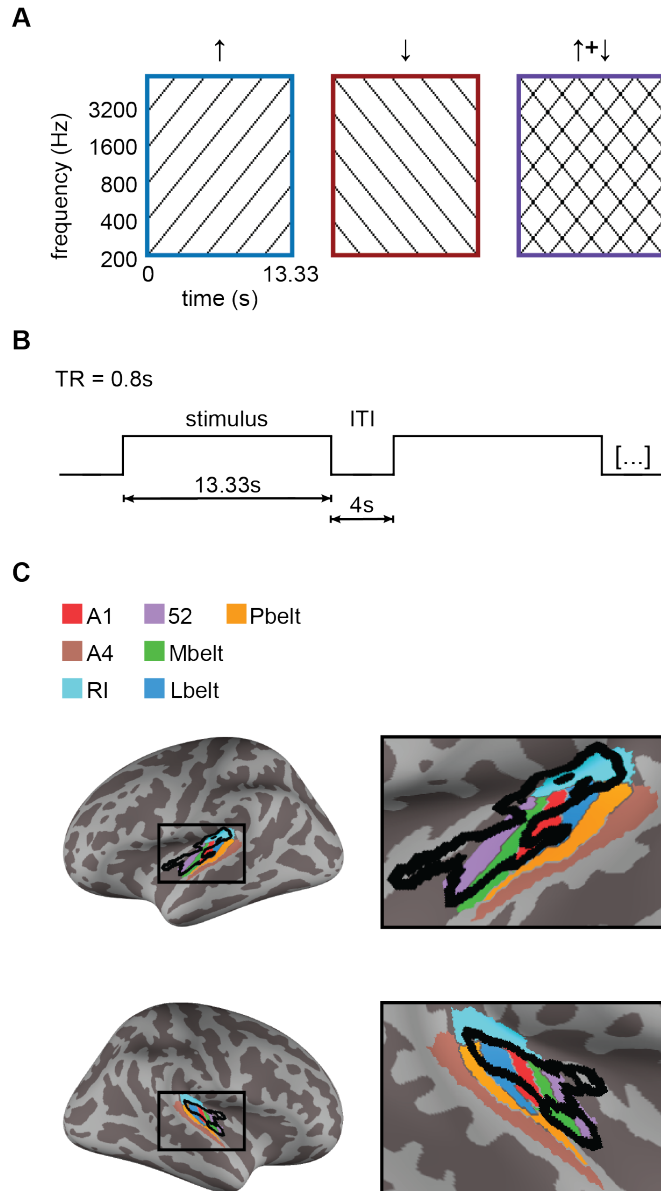


470
471
472
473
474
475

Figure 4. Bilateral regions of human auditory cortex show signatures of opponency.

A) A simple model auditory unit that responds more to upward direction spectral motion than downward directed spectral motion. The stimulus spectrogram is convolved with an upward-oriented spectrotemporal filter before the result is squared, as in a motion energy model (6).

- 476 B) Mean response of the unit to correlated pip stimuli with different delays and correlation
477 signs, corresponding to upward and downward directed positive and negative
478 correlations.
- 479 C) As in (B), but for an opponent signal, consisting of an upwardly tuned unit response
480 minus an identical unit tuned to downward motion.
- 481 D) Comparison of P(rising) for positive and negative correlation stimuli sweeping time
482 interval, aligning upward directed positive correlation stimuli with downward directed
483 negative correlation stimuli. Data replotted from Figure 2. The curves were not
484 significantly different ($p > 0.05$ by a two-way, repeated measures ANOVA).
- 485 E) As in (D) but for sweeping the tone difference. The curves were not significantly
486 different ($p > 0.05$ by a two-way, repeated measures ANOVA).
- 487 F) Conceptual schematic of opponency in brain regions. An opponent voxel/region would
488 respond strongly to rising and falling tones but be suppressed by the sum of the two
489 stimuli.
- 490 G) Stimulus design. Stimuli were rising, falling, or summed rising and falling.
- 491 H) Group level analysis. A bilateral region within auditory cortex responded less to summed
492 stimuli than non-summed stimuli. Cluster-corrected with false positive rate at $p < 0.05$
493 with a cluster-forming threshold of 20 voxels.
- 494 I) Individual level analysis. Regions in auditory cortex across subjects responded less to
495 summed stimuli than non-summed stimuli. Cluster-corrected with false positive rate at p
496 < 0.05 with a cluster-forming threshold of 20 voxels.
- 497 J) Control analysis showing symmetric beta values in response to rising and falling stimuli
498 in individually defined opponent ROIs ($p > 0.05$ via one-sample t-test). (Note that all beta
499 values are relative to an implicit baseline that includes responses to ambient scanner
500 noise.) Error bars represent mean \pm SEM (N = 5).
501

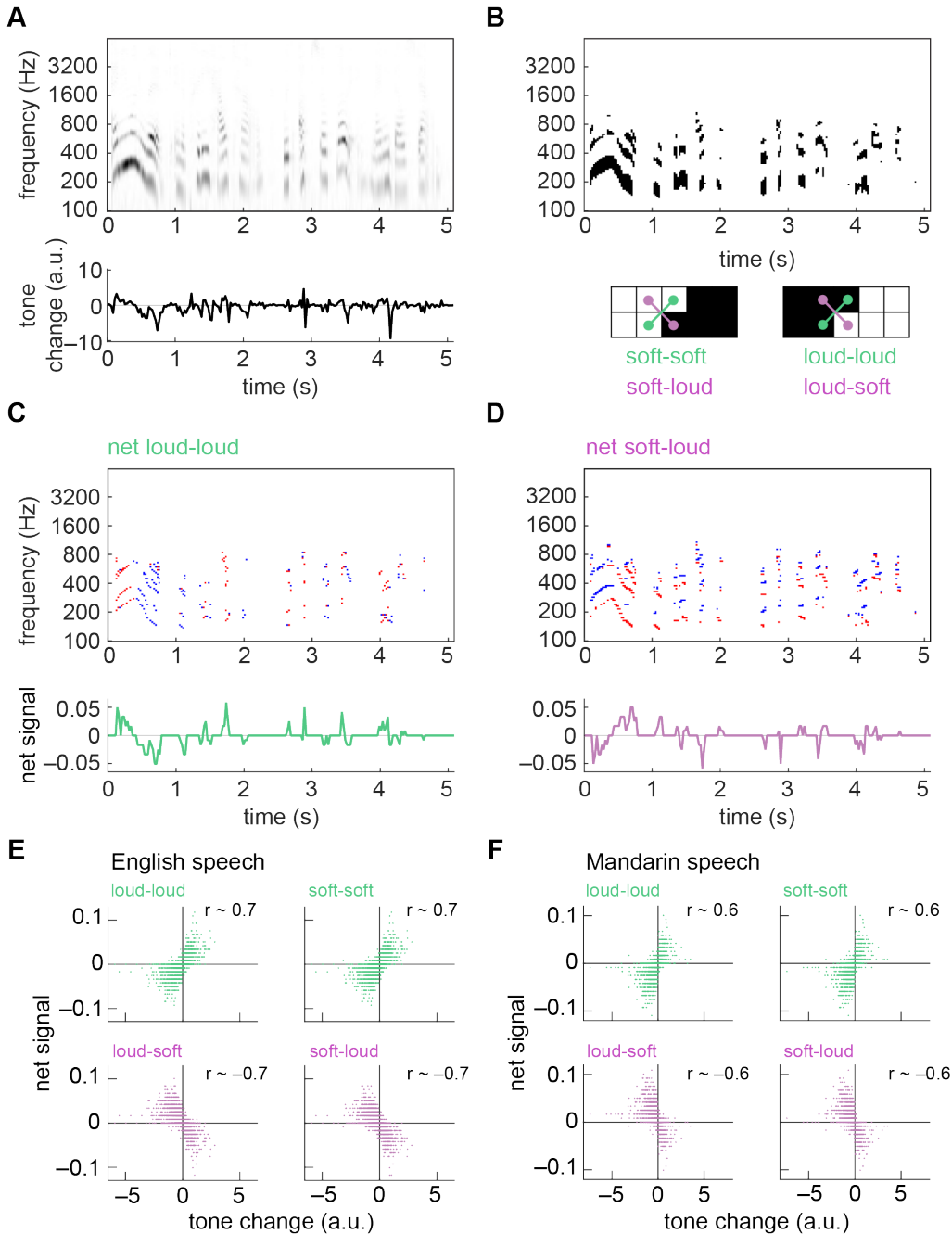


Supplemental Figure S4. Further method and result details from fMRI opponency experiment.

A) Depiction of actual stimuli used for the opponency experiment.

B) Time course of fMRI trial structure.

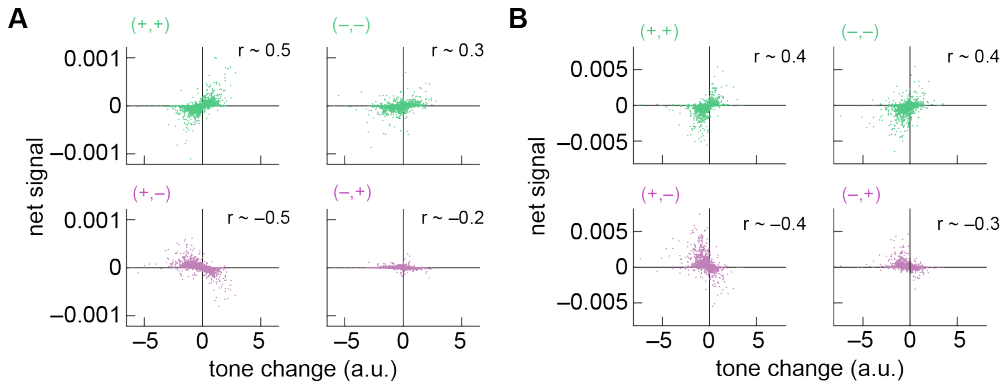
C) Group level analysis showing bilateral regions within auditory cortex that demonstrate significant opponent properties. Black outline reflects significant clusters from **Figure 4H**. Colored patches show cortical regions in accordance with (52). *RI* = retroinsular cortex; *Mbelt* = medial belt of auditory cortex; *Lbelt* = lateral belt of auditory cortex; *Pbelt* = parabelt region.



511
512 **Figure 5.** Rising and falling tone in spoken language can be detected through both positive and
513 negative pairwise correlations.

- 514 A) Spectrogram of voice saying, “Anyone lived in a pretty how town (with up so falling
515 many bells down)” (*top*). Intonation velocity estimate from spectrogram (*bottom*, see
516 **Methods**). Positive tone changes correspond to rising frequencies in the sound.
517 B) Binarized spectrogram from (A) (*top*). Four distinct loud and soft frequency-time
518 combinations in the binarized spectrogram (*bottom*).
519 C) Net loud-loud instances at each frequency and time in the binarized spectrogram in (B)
520 (*top*). Red is +1, blue is -1, white is 0. Frequency-averaged net loud-loud signal (*bottom*).

- 521 D) Net soft-loud instances at each frequency and time in the binarized spectrogram in (B)
522 (*top*). Red is +1, blue is -1, white is 0. Frequency-averaged net soft-loud signal (*bottom*).
523 E) Correlations between the tone change estimate at each time and the frequency-averaged
524 net signals for loud-loud, soft-soft, loud-soft, and soft-loud patterns. Data from English
525 speech corpus (71).
526 F) As in (E) but for Mandarin speech corpus (72).
527



528

529

530

Supplemental Figure S5. Multiplicative interactions of amplitude derivatives are informative about intonation direction.

531

532

533

534

535

536

537

538

539

540

541

542

A) Correlations between the tone change estimate at each time and a continuous correlator model using only positive signals (+,+), only negative signals (-,-), and mixtures of the two (+,- and -,+ (see **Methods**). The correlations comprising the net signal were obtained by taking the derivative of the spectrogram amplitude in time, then multiplying derivatives of neighboring frequencies with a time-step delay and subtracting a mirror image product. Signals were rectified before multiplication to obtain the four pairs of multiplied signals, which together add up to a full correlator model. The net signals computed from (+,+) and (-,-) pairs correlated positively with tone change, while the net signals from (+,-) and (-,+) pairs correlated negatively with tone change. Data from English speech corpus (71).

B) As in (A) but with data from Mandarin speech corpus (72).

543 **Supp. Movie 1.** Demonstration of positive and negative pairwise correlations using ternary
544 correlated noise stimuli, as in Figure 1.

545

546 **Supp. Movie 2.** Demonstration of positive and negative pairwise correlations using ternary
547 correlated noise stimuli, analogous to the stimuli in **Supp. Movie 1** but in visual motion
548 detection (27).

549

550 **Methods**

551

552 *Psychophysical measurements*

553

554 All participants (N = 33; 12 female; mean age: 23.3 years, range of 18 years to 32 years)
555 provided informed, written consent in accordance with procedures approved by the Yale
556 University Institutional Review Board. To measure human psychophysical curves (**Figures 1-3**),
557 we recruited participants with self-reported normal hearing from within the university
558 population. Participants were seated in a quiet room, wearing headphones (Model DT 770 PRO,
559 Beyerdynamic, Heilbronn, Germany) to listen to various sound stimuli and make perceptual
560 judgments. The sounds were created in Matlab and presented using Psychtoolbox (73-75) on a
561 Macbook Pro, using its native soundcard. Participants adjusted the volume to a comfortable
562 level, which we estimated to typically be around 60 dB. Each sound was played for 2 seconds,
563 after which participants were cued to judge, to the best of their ability, whether it sounded like a
564 rising or falling tone. To ensure they understood the task, participants went through several
565 example sounds with the researcher before beginning the experiment. Participants usually
566 completed two experiments lasting approximately 15 minutes each. The data was analyzed using
567 custom code written in Matlab. The code to produce the sounds, all anonymized data, and the
568 code used to analyze the data and produce Figure 1-3 are all publicly available at: [GitHub
569 repository here, to be made available on publication].

570

571 *Creating correlated sounds*

572

573 We created complex sounds containing multiple frequencies, following the design of visual
574 stimuli that have been informative in that field. To do this, we created a comb of constant carrier
575 frequencies, with frequencies ranging over 6 octaves from 200 Hz to 6400 Hz, with 15
576 frequencies per octave, equally spaced in log-space. The sampling frequency was chosen to be
577 20 kHz for all experiments. Each carrier frequency was then multiplied by a slower, time varying
578 envelope, before the frequencies were summed to make the overall waveform for that sound.
579 Mathematically, the sound waveform, $w(t)$, looks like:

580

$$w(t) = \sum_{i=1}^N \theta_i m_i(t) \sin(2\pi f_i t)$$

581 Where the f_i is the indexed carrier frequencies, t is sampled at 20 kHz, and the value θ_i was
582 chosen to roughly equalize the perceptual salience of the different frequencies, using the ISO
583 standard 226 at 60 dB. (We note that in various tests in lab, this perceptual salience scaling was
584 not critical for the percepts we measured; since we included it in initial experiments, we included
585 it for all stimuli in this study.) It remains to compute the suite of $m_i(t)$ envelope functions to
586 create each sound. The envelope functions were computed as outlined below. All envelope
587 functions are computed to have non-negative binary or ternary values, and were filtered with a
588 25 ms low-pass filter in the ternary stimuli (Fig. 1) and at 0.5 ms low-pass filter in the pip stimuli
589 (Figs. 2 and 3) to eliminate sharp transitions. After all waveforms $w(t)$ were created, they were
590 scaled to have a minimum value of -1 and maximum value of $+1$.

591

592 Ternary pairwise correlations (Figure 1B-D).

593 To create sounds with only local, pairwise correlations between specific frequency and time
594 offsets, we followed a protocol used in prior visual experiments (27, 28, 76). Based on informal
595 experiments attempting to optimize our own percepts, we discretized frequencies into 15 notes
596 per octave and time into 1/6 second frames. This change in frequency is similar to the most
597 salient change in frequency in a prior study (37). We then created an initial binary mask in this
598 coarse-time representation, B_{ij} , where i indexed the frequency and j the time step in 1/6 second
599 intervals. In each trial, each element of B was chosen from a Bernoulli distribution with
600 probability 0.5, then centered to have values of $\pm 1/2$ instead of 0 and 1. A ternary mask, M , was
601 created by the following formula:

602

$$603 \quad M_{i,j} = B_{i,j} + PB_{i+d,j+1}$$

604

605 The mask is thus the binary matrix added back to itself with a displacement in frequency of $d =$
606 ± 1 for upward and downward directed correlations. The mask is ternary, with values of 0, 1, and
607 -1 . The correlation parity is chosen by $P = \pm 1$, so that the offset matrices are added to create
608 positive correlations and subtracted to create negative correlations. The discrete autocorrelation
609 function of this mask M is equal to:

610

$$611 \quad C_{m,n} = \frac{1}{2} \delta_{m,0} \delta_{n,0} + \frac{1}{4} P (\delta_{m,d} \delta_{n,1} + \delta_{m,-d} \delta_{n,-1})$$

612

613 Where the $\delta_{i,j}$ terms are Kronicker delta functions (see **Fig. S1**). Importantly, the elements in the
614 mask are not deterministically the same or different at the spectrotemporal offset of the
615 correlated displacement, so that spectral patterns vary substantially at each temporal update of
616 the stimulus.

617

618 A continuous time expression for the autocorrelation function is available in a prior work
619 describing similar stimuli in vision (28).

620

621 The coarse-time matrix M was recentered to have values of 0, 0.5, and 1, then up-sampled to the
622 sampling frequency F_s to create $m_i(t)$ at each frequency. The masks were filtered with a 25 ms
623 low-pass filter to eliminate sharp transitions.

624

625 To create the stimuli with varying coherence, we replaced a fraction of mask elements with
626 random ternary stimuli, drawn from the values (0, 0.5, 1) with probabilities (0.25, 0.5, 0.25). The
627 fraction replaced was equal to $(1 - C)$ where C is the coherence value.

628

629 Binaural pairwise correlations (Figure 1D, E).

630 To play sounds such that correlations only existed by integrating across the ears, we simply
631 played $B_{i,j}$ in one ear and $PB_{i+d,j+1}$ in the other ear, for the correlations as described above to
632 create the ternary pairwise correlations. To play these binary masks, we created two masks
633 $M_{i,j} = B_{i,j}$ and $M_{i,j} = PB_{i\pm d,j+1}$ to play to the two ears. The matrices were recentered to have
634 values of 0 and 1, then up-sampled to the sampling frequency. The masks were filtered with a 0.5
635 ms low-pass filter to eliminate sharp transitions.

636

637 Correlated pips with time and frequency offsets (Figure 2).

638 To create the correlated pip stimulus, we discretized frequency space into 15 tones per octave.
639 We first initialized our masks $m_i(t)$ to be 0 for all times, sampled at the sampling frequency F_s .
640 We then placed initial delta-function pips in a Poisson distribution across all frequencies and
641 times in our sound, at a rate of 4 pips per frequency per second. Positive and negative pips were
642 equally probable, represented by mask values of ± 1 . We then created a second set of pips offset
643 by the selected change in frequency and delay time, according to the two different correlation
644 types. After imposing the correlations, the overall pip rate became 8 pips per frequency per
645 second. We then convolved this event-trace with a boxcar function with the length of the pip
646 duration to create the mask at F_s . Pips had a duration of 40 ms in Figure 2 and 20 ms in Figure
647 S2. Last, the masks were linearly transformed to be between 0 and 1 and filtered with a 0.5 ms
648 low-pass filter to eliminate sharp transitions. The loud values corresponded to values of 1 in the
649 mask, the soft to values of 0, and the background to values of 0.5.

650

651 Correlations between loud and soft pips (Figure 3).

652 These stimuli were generated similarly to the correlated pips stimulus above. However, only two
653 thirds of all pips were in correlated pairs of loud-loud, soft-soft, loud-soft, or soft-loud. In the
654 case of the loud-loud correlated pips, the remaining third of pips consisted of randomly placed
655 soft pips. In the case of soft-soft correlated pips, the remaining third of pips consisted of
656 randomly placed loud pips. And in the cases of soft-loud and loud-soft, the remaining third were
657 equally distributed between soft and loud pips. Thus, the four types had equal numbers correlated
658 pairs in each stimulus. The overall rate of pips for all stimuli was 6 pips per frequency per
659 second.

660

661 Triplet correlations (Figure S3).

662 We made triplet correlation binary masks, discretized in frequency at time, following prior
663 procedures (39, 41). The frequency was discretized in 15 tones per octave and time was
664 discretized into 1/6 second frames. The frequencies began at 200 Hz and ranged over 5 octaves.
665 The masks $m_i(t)$ were linearly transformed to have values of 0 and 1 and were filtered in time
666 with a 0.5 ms low-pass filter to eliminate sharp transitions.

667

668 Rising, falling, and opponent tones (Figure 5).

669 To create the rising, falling, and opponent tones used in our fMRI experiment, we used
670 frequencies discretized into 1/16 octave steps and time discretized into 1/6 second steps.
671 Ascending tones were created from a binary mask equal to an ascending line of time-frequency
672 elements in this discretized space (Fig. S5) and descending tones consisted of a descending line
673 of time-frequency elements. The summed ascending plus descending was the sum of the two
674 masks. All masks were filtered in time with a 0.5 ms low-pass filter to eliminate sharp
675 transitions. We switched to 16 steps per octave for this experiment so that the ascending and
676 descending stimuli never played the same frequency simultaneously, making the addition of the
677 stimuli more straightforward.

678

679 Code to generate the sounds used in these experiments is available at [GitHub repository on
680 publication].

681

682 *Model motion energy unit (Figure 4)*

683

684 We created a model motion energy unit by convolving a linear filter with a sound spectrogram,
685 then squaring the result. That is:

686

$$687 \quad r(t) = ((f_1 * S_1)(t) + (f_2 * S_2)(t))^2$$

688

689 The filters were chosen to be:

690

$$691 \quad f_1(t) = \frac{t}{\tau^2} e^{-t/\tau}$$
$$692 \quad f_2(t) = f_1(t - T)\Theta(t - T)$$

693

694 Where f_2 is just a time-shifted version of f_1 with a time shift of $T = 40$ ms. The function Θ is a
695 Heaviside step function. The two filters are applied to adjacent frequencies in the spectrogram,
696 $S_1(t)$ and $S_2(t)$, so that the filter enhances signals directed upward over time.

697

698 We computed the mean of $r(t)$ over time to get the mean response for a given stimuli. Stimuli
699 were created to match the correlated pip-style stimuli in Figure 2. The opponent response was
700 computed as

701

$$702 \quad r_{\text{opp}}(t) = ((f_1 * S_1)(t) + (f_2 * S_2)(t))^2 - ((f_2 * S_1)(t) + (f_1 * S_2)(t))^2$$

703

704 The second, negative term is the same as the first term but with the filter flipped in frequency
705 space, so that it corresponds to a downward selective unit. This response was likewise averaged
706 over time to produce the plots in Figure 4.

707

708 Matlab code to create **Figures 4B, C** is available at [Github repository on publication].

709

710 *Speech analysis*

711

712 Spoken language databases were analyzed to ask how spectrotemporal correlations could act as
713 indicators for rising and falling tones in speech. Using Matlab, we first loaded short snippets of
714 speech from two databases: 438 snippets constituting a total of 91 minutes of data from
715 Librispeech, a corpus of read English (71); and 749 snippets constituting a total of 52 minutes of
716 data from Magicdata Mandarin Chinese Read Speech Corpus (72), a corpus of read Mandarin.
717 We computed a spectrogram for each snippet of speech using the Matlab command
718 `spectrogram`; we extracted the spectral amplitude at a resolution of 40 samples per second
719 with no overlap between samples, at 20 evenly spaced frequencies per octave from 100 Hz to
720 6400 Hz (**Figure 5A**). We estimated the rising/falling intonation change of the sound at each
721 point using the Matlab command `opticalFlowHS`, which uses the Horn-Schunck method (77)
722 to estimate directional local flow (typically optic flow) between frames. We averaged the
723 calculated flow over frequencies to compute an estimate of the frequency “flow” with arbitrary
724 units, which we termed tone change (**Figure 5A**). This method does not make strong
725 assumptions about how changes in speech tone or frequency should be computed. It should work

726 to extract tone changes from most complex sounds. We then examined estimators of this tone
727 change as follows:

728

- 729 1) To compute binary correlations in frequency and time, we first binarized the spectrogram
730 using Otsu's method (Matlab command `imbinarize`) (78), which maximizes the
731 variance between the binarized time-frequency element amplitudes while minimizing
732 variance within each of the two categories (**Figure 5B**). We made 8 new binary
733 frequency-time data arrays, containing Boolean values at each point in time and
734 frequency, $V_{t,f,\uparrow,\pm,\pm} := (\{A_{t,f}, A_{t+1,f+1}\} = \{\pm 1, \pm 1\})$ and an equivalent one for
735 downward directed volume patterns. These matrices are records of the existence of each
736 pattern of sound intensity at each time and frequency. From these, we computed the net
737 signal of each pattern at each frequency by subtracting the downward directed matrix
738 from the upward directed one. We last found the mean net signal over all frequencies for
739 each pattern (**Figure 5C, D**). We computed the correlation between these mean net
740 signals at each time point with the calculated upward or downward flow velocity (**Figure**
741 **5E, F**). Note that the sum of these net pattern signals sum to 0 over the four different
742 patterns (\pm, \pm), so that the 4 signals are not independent.
- 743 2) To generate non-binarized correlation plots, we first linearly filtered the spectrogram
744 amplitudes, $A_{t,f}$, to take temporal derivatives: $F_{t,f} = A_{t,f} - A_{t-1,f}$. We then used these
745 derivatives, $F_{t,f}$, which have positive and negative values, as inputs to a Hassenstein-
746 Reichardt correlator model (Hassenstein and Reichardt 1956, Fitzgerald and Clark 2015).
747 We then computed the net (+,+) correlations, for instance, as $N_{t,f,+,+} =$
748 $[F_{t,f}]_+ [F_{t+1,f+1}]_+ - [F_{t+1,f}]_+ [F_{t,f+1}]_+$, where $[x]_+ = x$ when $x > 0$ and $[x]_+ = 0$
749 otherwise. A similar process computed the net (-,-), (+,-) and (-,+) correlations. We
750 averaged these signals over frequency to obtain a single indicator of velocity at each
751 point in time. These indicators were then correlated with the estimated tone change of the
752 sound snippet at that point in time (**Figure S5**).

753

754 Code to analyze the spoken language databases and produce the panels in Figure 5 is available at
755 [GitHub repository].

756

757 *fMRI recordings and analysis*

758

759 Whole-brain imaging was performed at the Brain Imaging Center at Yale University, on a
760 Siemens 3 T Prisma MRI scanner using a 32-channel head coil. Functional data were acquired
761 with a gradient-echo echoplanar pulse sequence (TR = 0.80 s, TE = 30 ms, flip angle = 52°,
762 voxel size = 2.4 mm × 2.4 mm × 2.4 mm, MB acc. factor = 6). T1-weighted MP-RAGE
763 anatomical images were collected as well (TR = 2.5 s, TE = 2.0 ms, flip angle = 8°, 208 slices,
764 voxel size = 1.0 mm isotropic). Functional imaging in our sample (N=5; 1 female; mean age:
765 26.2 years; authors PAV and SDM were participants in the fMRI study) was performed in ~5-
766 minute runs, with the total number of functional runs per participant ranging from 3-5. Fifteen
767 auditory stimuli were presented per run in an event-related design (5 each of three stimulus
768 types: rising, falling, and summed). Each stimulus lasted for 13.33 s, separated by an inter-trial
769 interval (ITI) of 4 s. The order of the three stimulus types was randomized in each run.
770 Participants passively listened to the tones and were not required to render any responses. MRI-

771 optimized noise-canceling headphones (Optoacoustics OptoACTIVE III) were used to limit
772 effects of background scanner noise and the noise-cancelling software was trained on the EPI
773 sequence sound features before each session using a brief calibration run.

774
775 The fMRI-Prep toolbox was used for preprocessing (79). The anatomical image was corrected
776 for intensity non-uniformity (INU) with N4BiasFieldCorrection (80) and used as T1w-reference.
777 The T1w-reference was then skull-stripped with a Nipype implementation of the
778 antsBrainExtraction.sh workflow in ANTs, and tissue segmentation of cerebrospinal fluid (CSF),
779 white-matter (WM), and gray-matter (GM) was performed on the brain-extracted T1w using
780 FFAST (FSL 6.0.5) (81). Volume-based spatial normalization to standard (MNI) space was
781 performed through nonlinear registration with antsRegistration (ANTs 2.3.3). For each of the
782 BOLD runs, a reference volume and its skull-stripped version were generated using a custom
783 methodology of fMRIPrep. Head-motion parameters were estimated using MCFLIRT (FSL
784 6.0.5) (82) and BOLD time-series were resampled into native space by applying the transforms
785 to correct for head-motion, and the BOLD reference was co-registered to the anatomical
786 reference using mri_coreg (FreeSurfer) followed by FLIRT. Co-registration was configured with
787 6 DOF. Several confounding time-series were calculated based on the preprocessed BOLD:
788 framewise displacement (FD), DVARS and three region-wise global signals. The BOLD time-
789 series were resampled into standard space, and volumetric resamplings were performed using
790 ANTs.

791
792 Our main analyses involved constructing general linear models (GLMs) to quantify the effects of
793 the three stimulus types within auditory cortex. GLM analyses were performed using Nilearn
794 (83). Confound regressors of no interest (generated using fMRIPrep, see above) were entered
795 into each GLM. These included six standard motion regressors, the framewise displacement time
796 course, and white matter and global signal time courses. Each stimulus type (rising, falling, and
797 summed) was modeled using boxcar regressors over the entire stimulus presentation phase
798 (13.33 s) of the relevant trials, and was convolved with the canonical double-gamma
799 hemodynamic response function. The main contrast of interest at the group and individual levels
800 compared BOLD responses to the non-summed directional stimuli (i.e., rising and falling) to the
801 summed stimuli (i.e., superimposed rising + falling). The contrast was designed to highlight
802 deviations from a null hypothesis of equivalent responses between directional and opponent
803 stimuli. Individual subject runs were combined in a fixed effects analysis and then brought to the
804 group level for mixed-effect analyses, where we controlled the false positive rate at $p < 0.05$ with
805 a cluster-forming threshold of 20 voxels. Critically, individual-level results for all subjects were
806 also analyzed and displayed, using the same thresholding parameters. All contrasts were
807 performed within an *a priori* anatomical mask that consisted of any voxels crossing the 50%
808 probability threshold within a combined bilateral probabilistic atlas (Harvard-Oxford) that
809 included both the STG and Heschel's gyrus. Individual and group results were projected onto
810 the standard (MNI) cortical surface (FreeSurfer) for visualization.

811
812 A simple control analysis was also performed to ensure that the non-summed > summed results
813 were not driven by a single non-summed stimulus (e.g., rising or falling) having a proportionally
814 larger response, but rather by symmetric responses to the rising and falling stimuli. To perform
815 this control analysis, we first extracted individualized regions of interest (ROIs) from the non-
816 summed > summed contrast (using the threshold described above), and then extracted average

817 beta values within that ROI for each stimulus type. We note that while this was of course not an
 818 unbiased ROI relative to the hypothesis that non-summed stimuli would on average show
 819 stronger activity than summed, it was unbiased relative to the hypothesis of symmetric responses
 820 to rising versus falling tones.

821
 822 *Opponency implies a symmetry in responses with opposite correlations in opposite directions*
 823

824 The motion energy model uses pairwise correlations to extract motion information from input
 825 stimuli and seems to accurately represent important aspects of cellular physiology (6). In the
 826 motion energy model, stimuli over space and time, $S(x, t)$, are convolved with a space-time
 827 oriented linear filter, $H(x, t)$. (In this section, we will derive results in space, but a frequency
 828 variable f could substitute for x and this approach would apply sound intensity over frequency
 829 rather than light intensity over space.) The result of the convolution is squared to obtain a
 830 response:

$$831$$

$$832 \quad r(x, t) = \left(\iint dx' dt' H(x', t') S(x - x', t - t') \right)^2$$

$$833$$

834 This response is stronger, on average, to stimuli with motion in the preferred direction than in the
 835 null direction. The preferred direction corresponds to the orientation of the filter H in space time,
 836 which amplifies signals when the motion direction aligns with the filter orientation. When the
 837 response is averaged over time and space, it yields a pleasing form in Fourier space, such that the
 838 mean response is the dot product of the stimulus power with a weighting function (6):

$$839$$

$$840 \quad \langle r \rangle = \iint dk d\omega |\tilde{H}(k, \omega)|^2 |\tilde{S}(k, \omega)|^2$$

$$841$$

842 Where \tilde{H} and \tilde{S} are the Fourier transforms of H and S . Therefore, to understand responses of this
 843 model, it is useful to compute the power spectrum of the stimulus.

844
 845 For a random dot kinetogram in which the dots are displaced by Δx in space and Δt in time, the
 846 covariance density, C , of the stimulus is a function of the offsets in time and space, x and t :

$$847$$

$$848 \quad C(x, t) = \beta \delta(x, t) + \alpha \delta(x - \Delta x, t - \Delta t) + \alpha \delta(x + \Delta x, t + \Delta t)$$

$$849$$

850 Where the first term is the stimulus autocovariance and the remaining two terms correspond to
 851 correlations in the stimulus at offsets of $(\Delta x, \Delta t)$ and $(-\Delta x, -\Delta t)$. For random dot kinetograms,
 852 $\beta < 1$ and α can take on positive or negative values for positively and negative correlated
 853 random dot kinetograms. This derivation is in continuous space, using Dirac delta function
 854 correlations; a similar result with discrete time and frequencies was found earlier in the methods
 855 for the ternary stimuli. The power spectrum of the stimulus is the Fourier transform of this
 856 covariance function:

$$857$$

$$858 \quad |\tilde{S}(k, \omega)|^2 = \iint dx dt e^{ikx} e^{i\omega t} C(x, t) = \beta + \alpha \cos(\omega \Delta t + k \Delta x)$$

$$859$$

860 The power is highest/lowest along lines of constant phase in cosine, or when $\omega\Delta t + k\Delta x = n\pi$.
 861 When the α is negative, for negative correlation stimuli, this effectively changes the phase of the
 862 cosine by 180 degrees. The motion energy model says the mean response to such a stimulus, for
 863 a unit with filter H , is:

864

$$\langle r \rangle = \iint dk d\omega |\tilde{H}(k, \omega)|^2 (\beta + \alpha \cos(\omega\Delta t + k\Delta x))$$

865

866 This is the type of curve shown in **Figure 4B**, in which there is a baseline response determined
 867 by β and the integral of $|\tilde{H}(k, \omega)|^2$. There is a modulatory term that depends on α and the dot
 868 product of $|\tilde{H}(k, \omega)|^2$ with $\cos(\omega\Delta t + k\Delta x)$, which gives the modulation a directional tuning.
 869 This form means that the modulation inverts when the sign of the correlation (sign of α) inverts.
 870 If there is a peak response to a stimulus with correlation α at a specific Δt and Δx , then the peak
 871 will be equal and opposite when α is inverted. Importantly, however, the peak is not the same
 872 when the direction of the stimulus is inverted, that is when $\Delta x \rightarrow -\Delta x$.

873

874 However, if we compute an opponent response, in which we subtract the response with one filter
 875 orientation from the response with the opposite filter orientation (inverting the k in the Fourier
 876 domain), then we find:

877

$$\langle r_{opp} \rangle = \iint dk d\omega (|\tilde{H}(k, \omega)|^2 - |\tilde{H}(-k, \omega)|^2) (\beta + \alpha \cos(\omega\Delta t + k\Delta x))$$

878

$$\langle r_{opp} \rangle = \alpha \iint dk d\omega (|\tilde{H}(k, \omega)|^2 - |\tilde{H}(-k, \omega)|^2) (\cos(\omega\Delta t + k\Delta x))$$

879

880 Here, we see that the opponent subtraction causes the β term to drop out entirely so that the
 881 remaining term is just proportional to α , the correlation in the stimulus. The mean opponent
 882 response can be computed for correlation stimuli with parameters α , Δt , and Δx :
 883 $\langle r_{opp}(\alpha, \Delta t, \Delta x) \rangle$. Because of the directional opponency, the response inverts when the stimulus
 884 is reversed in space:

885

$$\langle r_{opp}(\alpha, \Delta t, \Delta x) \rangle = -\langle r_{opp}(\alpha, \Delta t, -\Delta x) \rangle$$

886

887 And because of the proportionality with the correlation, the response inverts when the stimulus
 888 correlation is inverted:

889

$$\langle r_{opp}(\alpha, \Delta t, \Delta x) \rangle = -\langle r_{opp}(-\alpha, \Delta t, \Delta x) \rangle$$

890

891 Therefore, for an opponent signal, inverting the correlation is equivalent to inverting the
 892 direction of the signal:

893

$$\langle r_{opp}(-\alpha, \Delta t, \Delta x) \rangle = \langle r_{opp}(\alpha, \Delta t, -\Delta x) \rangle$$

894

895 For any set of filters, as long as they are opponently subtracted, inverting the sign of the
 896 correlation is identical to inverting the direction of the stimulus, when computing the

897

902 spatiotemporal average response. So when stimuli can be generated that have autocovariance
903 structures like those in the ternary scintillator (**Fig. 1**) or in a random dot kinetogram (**Fig. 2**), if
904 the computation is based on pairwise correlations and is opponent, the equations above show that
905 the response will always be inverted when the stimulus correlation is inverted, and always be
906 equivalent to inverting the direction of the stimulus. Therefore, opponency implies the sort of
907 inversion symmetries we observed in our data, where inverting the correlation sign generates
908 percepts with the same tuning as inverting the direction of the stimulus (**Fig. 4D, E**, but also
909 visible in **Figs. 1-3**). Opponency also implies the sort of consistent symmetries between positive
910 and negative correlation stimuli observed in human motion perception (36). We note that it is
911 also possible to achieve this kind of symmetry using precisely defined filters that lead to
912 opponent properties in single units, without a subtractive step (64).

913
914
915
916

917 **References**

918

- 919 1. D. Hirst, A. Di Cristo, Intonation systems. *A survey of Twenty Languages* (1998).
- 920 2. J. Gandour, Tone perception in Far Eastern languages. *Journal of phonetics* **11**, 149-175
921 (1983).
- 922 3. M. J. W. Yip, *Tone* (Cambridge University Press, 2002).
- 923 4. Z.-L. Lu, L. A. Lesmes, G. Sperling, Perceptual motion standstill in rapidly moving
924 chromatic displays. *Proc. Natl. Acad. Sci. USA* **96**, 15374-15379 (1999).
- 925 5. B. Hassenstein, W. Reichardt, Systemtheoretische Analyse der Zeit-, Reihenfolgen- und
926 Vorzeichenauswertung bei der Bewegungspertzeption des Rüsselkäfers *Chlorophanus*.
927 *Zeits. Naturforsch.* **11**, 513–524 (1956).
- 928 6. E. Adelson, J. Bergen, Spatiotemporal energy models for the perception of motion. *JOSA*
929 *A* **2**, 284-299 (1985).
- 930 7. S. M. Anstis, B. J. Rogers, Illusory reversal of visual depth and movement during
931 changes of contrast. *Vision Res.* **15**, 957-IN956 (1975).
- 932 8. Z. L. Lu, G. Sperling, Three-systems theory of human visual motion perception: review
933 and update. *JOSA A* **18**, 2331-2370 (2001).
- 934 9. Z.-L. Lu, G. Sperling, The functional architecture of human visual motion perception.
935 *Vision Res.* **35**, 2697-2722 (1995).
- 936 10. J. H. McDermott, The cocktail party problem. *Curr. Biol.* **19**, R1024-R1027 (2009).
- 937 11. E. C. Cherry, Some experiments on the recognition of speech, with one and with two
938 ears. *The Journal of the acoustical society of America* **25**, 975-979 (1953).
- 939 12. J. K. Bizley, Y. E. Cohen, The what, where and how of auditory-object perception. *Nat.*
940 *Rev. Neurosci.* **14**, 693-707 (2013).
- 941 13. A. S. Bregman, J. Campbell, Primary auditory stream segregation and perception of order
942 in rapid sequences of tones. *Journal of experimental psychology* **89**, 244 (1971).
- 943 14. V. Ciocca, A. S. Bregman, Perceived continuity of gliding and steady-state tones through
944 interrupting noise. *Perception & Psychophysics* **42**, 476-484 (1987).
- 945 15. J. M. Sinnott, R. N. Aslin, Frequency and intensity discrimination in human infants and
946 adults. *The Journal of the Acoustical Society of America* **78**, 1986-1992 (1985).
- 947 16. L. Demany, C. Ramos, On the binding of successive sounds: Perceiving shifts in
948 nonperceived pitches. *The Journal of the Acoustical Society of America* **117**, 833-841
949 (2005).

- 950 17. R. N. Aslin, Discrimination of frequency transitions by human infants. *The Journal of the*
951 *Acoustical Society of America* **86**, 582-590 (1989).
- 952 18. K. Siedenburg, J. Graves, D. Pressnitzer, A unitary model of auditory frequency change
953 perception. *PLoS Comp. Biol.* **19**, e1010307 (2023).
- 954 19. L. Demany, C. Semal, Automatic frequency-shift detection in the auditory system: A
955 review of psychophysical findings. *Neuroscience* **389**, 30-40 (2018).
- 956 20. H.-W. Lu, P. H. Smith, P. X. Joris, Mammalian octopus cells are direction selective to
957 frequency sweeps by excitatory synaptic sequence detection. *Proc. Natl. Acad. Sci. USA*
958 **119**, e2203748119 (2022).
- 959 21. R. I. Kuo, G. K. Wu, The generation of direction selectivity in the auditory system.
960 *Neuron* **73**, 1016-1027 (2012).
- 961 22. C.-q. Ye, M.-m. Poo, Y. Dan, X.-h. Zhang, Synaptic mechanisms of direction selectivity
962 in primary auditory cortex. *J. Neurosci.* **30**, 1861-1868 (2010).
- 963 23. S. Andoni, N. Li, G. D. Pollak, Spectrotemporal receptive fields in the inferior colliculus
964 revealing selectivity for spectral motion in conspecific vocalizations. *J. Neurosci.* **27**,
965 4882-4893 (2007).
- 966 24. R. C. DeCharms, D. T. Blake, M. M. Merzenich, Optimizing sound features for cortical
967 neurons. *science* **280**, 1439-1444 (1998).
- 968 25. F. E. Theunissen, K. Sen, A. J. Doupe, Spectral-temporal receptive fields of nonlinear
969 auditory neurons obtained using natural sounds. *J. Neurosci.* **20**, 2315-2331 (2000).
- 970 26. L. M. Miller, M. A. Escabí, H. L. Read, C. E. Schreiner, Spectrotemporal receptive fields
971 in the lemniscal auditory thalamus and cortex. *J. Neurophysiol.* **87**, 516-527 (2002).
- 972 27. E. Salazar-Gatzimas *et al.*, Direct measurement of correlation responses in *Drosophila*
973 elementary motion detectors reveals fast timescale tuning. *Neuron* **92**, 227-239 (2016).
- 974 28. S. Roy, R. d. R. van Steveninck, Bilocal visual noise as a probe of wide field motion
975 computation. *J. Vis.* **16**, 8-8 (2016).
- 976 29. M. J. McPherson, J. H. McDermott, Relative pitch representations and invariance to
977 timbre. *Cognition* **232**, 105327 (2023).
- 978 30. M. J. McPherson, J. H. McDermott, Diversity in pitch perception revealed by task
979 dependence. *Nature human behaviour* **2**, 52-66 (2018).
- 980 31. D. A. Clark, L. Bursztyn, M. A. Horowitz, M. J. Schnitzer, T. R. Clandinin, Defining the
981 computational structure of the motion detector in *Drosophila*. *Neuron* **70**, 1165-1177
982 (2011).

- 983 32. M. B. Orger, M. C. Smear, S. M. Anstis, H. Baier, Perception of Fourier and non-Fourier
984 motion by larval zebrafish. *Nat. Neurosci.* **3**, 1128-1133 (2000).
- 985 33. B. Krekelberg, T. Albright, Motion mechanisms in macaque MT. *J. Neurophysiol.* **93**,
986 2908 (2005).
- 987 34. J. B. Kelly, P. W. Judge, Binaural organization of primary auditory cortex in the ferret
988 (*Mustela putorius*). *J. Neurophysiol.* **71**, 904-913 (1994).
- 989 35. K. H. Britten, M. N. Shadlen, W. T. Newsome, J. A. Movshon, The analysis of visual
990 motion: a comparison of neuronal and psychophysical performance. *J. Neurosci.* **12**,
991 4745-4765 (1992).
- 992 36. R. Bours, M. Kroes, M. Lankheet, Sensitivity for reverse-phi motion. *Vision Res.* **49**, 1-9
993 (2009).
- 994 37. L. Demany, D. Pressnitzer, C. Semal, Tuning properties of the auditory frequency-shift
995 detectors. *The Journal of the Acoustical Society of America* **126**, 1342-1348 (2009).
- 996 38. J. Allik, E. Dzhafarov, A. Houtsma, J. Ross, N. Versfeld, Pitch motion with random
997 chord sequences. *Perception & Psychophysics* **46**, 513-527 (1989).
- 998 39. Q. Hu, J. D. Victor, A set of high-order spatiotemporal stimuli that elicit motion and
999 reverse-phi percepts. *J. Vis.* **10** (2010).
- 1000 40. J. E. Fitzgerald, A. Y. Katsov, T. R. Clandinin, M. J. Schnitzer, Symmetries in stimulus
1001 statistics shape the form of visual motion estimators. *Proc. Natl. Acad. Sci. USA* **108**,
1002 12909-12914 (2011).
- 1003 41. D. A. Clark *et al.*, Flies and humans share a motion estimation strategy that exploits
1004 natural scene statistics. *Nat. Neurosci.* **17**, 296-303 (2014).
- 1005 42. J. Chen, H. B. Mandel, J. E. Fitzgerald, D. A. Clark, Motion estimates in flies include
1006 higher-order correlations that cancel noise induced by the structure of natural scenes. *In*
1007 *preparation* (2018).
- 1008 43. T. Yildizoglu, C. Riegler, J. E. Fitzgerald, R. Portugues, A Neural Representation of
1009 Naturalistic Motion-Guided Behavior in the Zebrafish Brain. *Curr. Biol.* (2020).
- 1010 44. G. J. Brown, M. Cooke, Computational auditory scene analysis. *Computer Speech &*
1011 *Language* **8**, 297-336 (1994).
- 1012 45. R. J. Snowden, S. Treue, R. G. Erickson, R. A. Andersen, The response of area MT and
1013 V1 neurons to transparent motion. *J. Neurosci.* **11**, 2768-2785 (1991).
- 1014 46. N. Qian, R. A. Andersen, Transparent motion perception as detection of unbalanced
1015 motion signals. II. Physiology. *J. Neurosci.* **14**, 7367-7380 (1994).

- 1016 47. C. D. Salzman, C. M. Murasugi, K. H. Britten, W. T. Newsome, Microstimulation in
1017 visual area MT: effects on direction discrimination performance. *J. Neurosci.* **12**, 2331-
1018 2355 (1992).
- 1019 48. A. S. Mauss *et al.*, Neural circuit to integrate opposing motions in the visual field. *Cell*
1020 **162**, 351-362 (2015).
- 1021 49. S. Andoni, G. D. Pollak, Selectivity for spectral motion as a neural computation for
1022 encoding natural communication signals in bat inferior colliculus. *J. Neurosci.* **31**, 16529-
1023 16540 (2011).
- 1024 50. D. J. Heeger, G. M. Boynton, J. B. Demb, E. Seidemann, W. T. Newsome, Motion
1025 opponency in visual cortex. *J. Neurosci.* **19**, 7162-7174 (1999).
- 1026 51. B. Tian, J. P. Rauschecker, Processing of frequency-modulated sounds in the lateral
1027 auditory belt cortex of the rhesus monkey. *J. Neurophysiol.* **92**, 2993-3013 (2004).
- 1028 52. M. F. Glasser *et al.*, A multi-modal parcellation of human cerebral cortex. *Nature* **536**,
1029 171-178 (2016).
- 1030 53. N. C. Singh, F. E. Theunissen, Modulation spectra of natural sounds and ethological
1031 theories of auditory processing. *The Journal of the Acoustical Society of America* **114**,
1032 3394-3411 (2003).
- 1033 54. R. N. Shepard, Circularity in judgments of relative pitch. *The journal of the acoustical*
1034 *society of America* **36**, 2346-2353 (1964).
- 1035 55. M. J. McPherson, J. H. McDermott, Time-dependent discrimination advantages for
1036 harmonic sounds suggest efficient coding for memory. *Proc. Natl. Acad. Sci. USA* **117**,
1037 32169-32180 (2020).
- 1038 56. H. G. Yi, M. K. Leonard, E. F. Chang, The encoding of speech sounds in the superior
1039 temporal gyrus. *Neuron* **102**, 1096-1110 (2019).
- 1040 57. G. Hickok, D. Poeppel, The cortical organization of speech processing. *Nat. Rev.*
1041 *Neurosci.* **8**, 393-402 (2007).
- 1042 58. W. A. de Heer, A. G. Huth, T. L. Griffiths, J. L. Gallant, F. E. Theunissen, The
1043 hierarchical cortical organization of human speech processing. *J. Neurosci.* **37**, 6539-
1044 6557 (2017).
- 1045 59. Y. Li, C. Tang, J. Lu, J. Wu, E. F. Chang, Human cortical encoding of pitch in tonal and
1046 non-tonal languages. *Nature communications* **12**, 1161 (2021).
- 1047 60. C. Tang, L. Hamilton, E. Chang, Intonational speech prosody encoding in the human
1048 auditory cortex. *Science* **357**, 797-801 (2017).

- 1049 61. I. S. Johnsrude, V. B. Penhune, R. J. Zatorre, Functional specificity in the right human
1050 auditory cortex for perceiving pitch direction. *Brain* **123**, 155-163 (2000).
- 1051 62. J. A. Zavatone-Veth, B. A. Badwan, D. A. Clark, A minimal synaptic model for direction
1052 selective neurons in *Drosophila*. *J. Vis.* **20**, 2-2 (2020).
- 1053 63. C.-H. Mo, C. Koch, Modeling reverse-phi motion-selective neurons in cortex: double
1054 synaptic-veto mechanism. *Neural Comput.* **15**, 735-759 (2003).
- 1055 64. B. A. Badwan, M. S. Creamer, J. A. Zavatone-Veth, D. A. Clark, Dynamic nonlinearities
1056 enable direction opponency in *Drosophila* elementary motion detectors. *Nat. Neurosci.*
1057 **22**, 1318-1326 (2019).
- 1058 65. D. A. Clark, J. B. Demb, Parallel computations in insect and mammalian visual motion
1059 processing. *Curr. Biol.* **26**, R1062-R1072 (2016).
- 1060 66. J. R. Sanes, S. L. Zipursky, Design principles of insect and vertebrate visual systems.
1061 *Neuron* **66**, 15 (2010).
- 1062 67. A. Borst, M. Helmstaedter, Common circuit design in fly and mammalian motion vision.
1063 *Nat. Neurosci.* **18**, 1067-1076 (2015).
- 1064 68. N. Kadakia *et al.*, Odour motion sensing enhances navigation of complex plumes. *Nature*
1065 **611**, 754-761 (2022).
- 1066 69. J. E. Fitzgerald, D. A. Clark, Nonlinear circuits for naturalistic visual motion estimation.
1067 *eLife*, e09123 (2015).
- 1068 70. E. Salazar-Gatzimas, M. Agrochao, J. E. Fitzgerald, D. A. Clark, The Neuronal Basis of
1069 an Illusory Motion Percept Is Explained by Decorrelation of Parallel Motion Pathways.
1070 *Curr. Biol.* **28**, 3748-3762 (2018).
- 1071 71. V. Panayotov, G. Chen, D. Povey, S. Khudanpur (2015) Librispeech: an asr corpus based
1072 on public domain audio books. in *2015 IEEE international conference on acoustics,*
1073 *speech and signal processing (ICASSP)* (IEEE), pp 5206-5210.
- 1074 72. MagicData, Retrieved from <https://www.openslr.org/123/> on March 28, 2023. (2019).
- 1075 73. M. Kleiner *et al.*, What's new in Psychtoolbox-3. *Perception* **36**, 1 (2007).
- 1076 74. D. H. Brainard, The psychophysics toolbox. *Spatial vision* **10**, 433-436 (1997).
- 1077 75. D. G. Pelli, The VideoToolbox software for visual psychophysics: Transforming numbers
1078 into movies. *Spatial vision* **10**, 437-442 (1997).
- 1079 76. R. van Steveninck, W. Bialek, M. Potters, R. Carlson, G. Lewen (1996) Adaptive
1080 movement computation by the blowfly visual system. in *Natural & Artificial Parallel*
1081 *Computation: Proceedings of the Fifth NEC Research Symposium* (SIAM), p 21.

- 1082 77. B. K. Horn, B. G. Schunck, Determining optical flow. *Artif. Intell.* **17**, 185-203 (1981).
- 1083 78. N. Otsu, A threshold selection method from gray-level histograms. *Automatica* **11**, 23-27
1084 (1975).
- 1085 79. O. Esteban *et al.*, fMRIPrep: a robust preprocessing pipeline for functional MRI. *Nat.*
1086 *Methods* **16**, 111-116 (2019).
- 1087 80. N. J. Tustison *et al.*, N4ITK: improved N3 bias correction. *IEEE Trans. Med. Imaging*
1088 **29**, 1310-1320 (2010).
- 1089 81. Y. Zhang, M. Brady, S. Smith, Segmentation of brain MR images through a hidden
1090 Markov random field model and the expectation-maximization algorithm. *IEEE Trans.*
1091 *Med. Imaging* **20**, 45-57 (2001).
- 1092 82. M. Jenkinson, P. Bannister, M. Brady, S. Smith, Improved optimization for the robust
1093 and accurate linear registration and motion correction of brain images. *NeuroImage* **17**,
1094 825-841 (2002).
- 1095 83. A. Abraham *et al.*, Machine learning for neuroimaging with scikit-learn. *Frontiers in*
1096 *neuroinformatics* **8**, 14 (2014).
1097