# A structural theory of pitch

A structural theory of pitch

**Jonathan Laudanski**[1,5,†], **Yi Zheng**[1,2,3,4] **and Romain Brette**[1,2,3,4]

[1]*Institut d'Etudes de la Cognition, Ecole Normale Supérieure, Paris, France*

[2]*Sorbonne Universités, UPMC Univ. Paris 06, UMR_S 968, Institut de la Vision, Paris, F-75012, France*

[3]*INSERM, U968, Paris, F-75012, France*

[4]*CNRS, UMR_7210, Paris, F-75012, France*

[5]*Scientific and Clinical Research Department, Neurelec, Vallauris, France*

[†]deceased, 11 May 2014

Corresponding Author: Romain Brette, INSERM, 17 rue Moreau, Paris, France 75012, voice: +33153462536, Email: romain.brette@inserm.fr

eNeuro

A structural theory of pitch

1  **Abstract**

2  Musical notes can be ordered from low to high along a perceptual dimension called « pitch ». A
3  characteristic property of these sounds is their periodic waveform, and periodicity generally
4  correlates with pitch. Thus pitch is often described as the perceptual correlate of the periodicity
5  of the sound's waveform. However, the existence and salience of pitch also depends in a complex
6  way on other factors, in particular harmonic content: for example, periodic sounds made of high-
7  order harmonics tend to have a weaker pitch than those made of low-order harmonics. Here we
8  examine the theoretical proposition that pitch is the perceptual correlate of the regularity
9  structure of the vibration pattern of the basilar membrane, across place and time - a
10 generalization of the traditional view on pitch. While this proposition also attributes pitch to
11 periodic sounds, we show that it predicts differences between resolved and unresolved
12 harmonic complexes and a complex domain of existence of pitch, in agreement with
13 psychophysical experiments. We also present a possible neural mechanism for pitch estimation
14 based on coincidence detection, which does not require long delays, in contrast with standard
15 temporal models of pitch.

16

17  **Significance statement**

18 Melodies are composed of sounds that can be ordered on a musical scale. "Pitch" is the
19 perceptual dimension on that scale. To a large extent, the periodicity of the sound's waveform
20 can be mapped to pitch. However, the existence and strength of pitch also depends on the
21 harmonic content sounds, i.e., their timbre, which does not fit with this simple view. We propose
22 to explain these observations by the fact that the input to the auditory system is the spatio-
23 temporal vibration of the basilar membrane in the cochlea, rather than the acoustic waveform.
24 We show that defining pitch as the regularity structure of that vibration can explain some
25 aspects of the complexity of pitch perception.

26

27  **Introduction**

28  A musical note played by a piano or a trumpet has a perceptual attribute called "pitch", which
29  can be low or high. The same key played on different instruments produces sounds with
30  different spectral content but identical pitch. To a large extent, pitch can be mapped to the
31  periodicity, or repetition rate (f0), of the acoustic waveform (Oxenham, 2012). For this reason,
32  theories of pitch perception have focused on how the auditory system extracts periodicity. In the
33  cochlea, the mechanical response of the basilar membrane (BM) to sounds has both a spatial and
34  a temporal dimension. The BM vibrates in response to tones, following the frequency of the tone.
35  The place of maximal vibration along the BM also changes gradually with tone frequency, from
36  the base (high frequency) to the apex (low frequency). Accordingly, there are two broad types of
37  theories of pitch, emphasizing either time or place (de Cheveigné, 2010).

38  Place theories (or pattern recognition theories) propose that the spatial pattern of BM vibration
39  is compared to internal templates, consisting of harmonic series of fundamental frequencies
40  (Terhardt, 1974). Pitch is then estimated from the fundamental frequency of the best matching
41  template. This mechanism requires that harmonics of the sound produce clear peaks in the
42  spatial pattern of BM vibration, i.e., that harmonics are "resolved" by the cochlea, but this is
43  typically not the case for high-order harmonics because the bandwidth of cochlear filters
44  increases with center frequency. In contrast, tone complexes with only unresolved harmonics
45  can elicit a pitch (Ritsma, 1962; Oxenham et al., 2011). In addition, the firing rate of auditory
46  nerve fibers as well as most neurons in the cochlear nucleus saturates at high levels, but pitch
47  perception does not degrade at high levels (Cedolin and Delgutte, 2005).

48  Temporal theories propose that periodicity is estimated from the temporal waveform in each
49  auditory channel (cochlear place), and estimates are then combined across channels (Licklider,
50  1951; Meddis and O'Mard, 1997; de Cheveigné, 2010). Sound periodicity is indeed accurately
51  reflected in the patterns of spikes produced by auditory nerve fibers (Cariani and Delgutte,
52  1996a, 1996b; Cedolin and Delgutte, 2005). Resolvability plays little role in these theories, but
53  pitch based on resolved harmonics is more salient and easier to discriminate than pitch based on
54  unresolved harmonics (Houtsma and Smurzynski, 1990; Carlyon and Shackleton, 1994; Carlyon,
55  1998; Bernstein and Oxenham, 2003). Finally, detecting the periodicity of a waveform with
56  repetition rate f0 = 30 Hz (the lower limit of pitch (Pressnitzer et al., 2001)) would require
57  delays of about 30 ms, of which there is no clear physiological evidence.

58  In addition, the domain of existence of pitch is complex, which neither type of theory explains:
59  the existence of pitch depends not only on f0 but also on resolvability of harmonics and spectral
60  content (Pressnitzer et al., 2001; Oxenham et al., 2004b, 2011). For example, high frequency
61  complex tones (>4 kHz) with f0 = 120 Hz do not have a clear pitch while a pure tone with the
62  same f0 does (Oxenham et al., 2004b); but high frequency complex tones with f0>400 Hz do
63  have a clear pitch (Oxenham et al., 2011). Finally, while pitch is generally independent of sound
64  intensity (contradicting place theories (Micheyl and Oxenham, 2007)), a few studies suggest a
65  small but significant intensity dependence of pitch for low frequency pure tones (Morgan et al.,
66  1951; Verschuure and Van Meeteren, 1975; Burns, 1982) (contradicting temporal theories).

67  Here we propose to address these issues by reexamining the postulate that pitch is the
68  perceptual correlate of the periodicity of the acoustic waveform. Starting from the observation
69  that the input to the auditory system is not the acoustic waveform but the vibration pattern of

2

70 the BM, we propose instead that pitch is the perceptual correlate of the regularity structure of
71 the BM vibration pattern, across place and time. While this proposition also attributes pitch to
72 periodic sounds, we show that it predicts differences between resolved and unresolved
73 harmonic complexes and a complex domain of existence of pitch. We also present a possible
74 neural mechanism for pitch estimation based on coincidence detection, which does not require
75 long delays.

76

## Materials and Methods

78 **Auditory filters**

79 Auditory filters were modeled as gammatone filters (Slaney, 1993; Fontaine et al., 2011), which
80 approximate reverse correlation filters of cat auditory nerve fibers (Boer and Jongh, 1978;
81 Carney and Yin, 1988) and have been matched to psychophysical measurements in humans
82 (Glasberg and Moore, 1990). Their impulse response defined by: $H(t) =$
83 $t^{n-1}e^{-t/\tau}\cos(2\pi.CF.t)$, where CF is the characteristic frequency, n is the order and the
84 bandwidth is set by $\tau = \left(2\pi \cdot 1.019 \cdot (24.7 + 0.108 \cdot CF)\right)^{-1}$. Filters were spaced uniformly in
85 ERB scale (Glasberg and Moore, 1990) with CF between 100 and 8000 Hz.

86 **Neural model of pitch estimation**

87 The neural model of pitch estimation includes two layers: 1) the input layer (putatively cochlear
88 nucleus) and 2) coincidence detector neurons.

89 *Input layer*

90 Each neuron receives the output x(t) of a gammatone filter, after half-wave rectification and
91 compression with a power law with exponent $\gamma = 0.3$ (Stevens, 1971; Zwislocki, 1973):
92 $y(t) = \kappa([x(t)]^+)^\gamma$ (varying the exponent between 0.2 and 0.5 did not affect the results).

93 We tested different spiking neuron models (Fig. 4), defined by a membrane equation of the
94 following form:

95
$$C\frac{dV}{dt} = g_L(E_L - V) + y(t) + \sigma\xi(t) + I(V), \tag{1}$$

96 where V is the membrane potential, $g_L(E_L - V)$ represent the non-specific leak current, $\sigma$ is the
97 noise level, C is the membrane capacitance and I(V) represents currents from voltage-gated
98 channels.

99 The chopper cell model (T-multipolar) is based on the model of Rothman and Manis (Rothman
100 and Manis, 2003a), with maximal conductances $g_{Na} = 1000$ nS, $g_{KHT} = 150$ nS, and $g_h = 0.5$
101 nS. Octopus cells are also based on the same model but include a low threshold potassium
102 channel (KLT) and model of $I_h$ taken from (Khurana et al., 2011), with $g_{Na} = 1000$ nS,
103 $g_{KHT} = 150$ nS, $g_{KLT} = 600$ nS, and $g_h = 40$ nS. These two models were used only in Fig. 4.

104 We also used a leaky integrate-and-fire model (LIF), a phenomenological model with good
105 predictive value for a broad class of neurons (Jolivet et al., 2004; Gerstner and Naud, 2009). The
106 membrane time constant was $\tau = g_L/C = 1.5$ ms. The model spikes when V(t) reaches the

107 threshold $\theta = -40$ mV, and V(t) is then reset to $V_r = -60$mV and clamped at this value for a
108 refractory time of $\tau_r = 1$ ms. This model was used in all simulations, unless otherwise specified.

109 *Coincidence detectors*

110 The second layer consists of coincidence detectors, which are modeled as integrate-and-fire
111 models (as above) with an adaptive threshold governed by the following equation (Platkiewicz
112 and Brette, 2010, 2011; Fontaine et al., 2014):

$$\tau_\theta \frac{d\theta}{dt} = \theta_0 - \theta + V - E_L, \tag{2}$$

114 where $\theta_0 = -40$mV is the value of threshold at rest and $\tau_\theta = 5$ ms. (note that half-wave
115 rectification can be discarded here because V is always above $E_L$, as there are only excitatory
116 synapses). This equation ensures that the neuron is always in a fluctuation-driven regime where
117 it is sensitive to coincidences (Platkiewicz and Brette, 2011). The response of the coincidence
118 detectors was only considered after 30 ms following note onset.

119 *Synaptic connections*

120 For each possible f0, we build a group of coincidence detectors whose inputs are synchronous
121 when a sound of period 1/f0 is presented. For any sound, the *synchrony partition* is defined as
122 the set of groups of input neurons that fire in synchrony for that particular sound (Brette, 2012)
123 (synchrony is within group, not across groups). One coincidence detector neuron is assigned to
124 each group (synaptic connections from each input neuron to the coincidence detector), so that
125 the synchrony partition corresponds to a set of coincidence detector neurons.

126 To build a group of coincidence detector neurons tuned to periodic sounds with fundamental
127 frequency f0, we consider the synchrony partition of the complex tone made of all harmonics of
128 f0, i.e., tones of frequency k.f0. For each harmonic, we select all pairs of channels in our filter
129 bank that satisfy the following properties (Fig. 2D): 1) the gain at k.f0 is greater than a threshold
130 $G_{min}$ = 0.25 (dashed line in Fig. 2D); 2) the two gains at k.f0 are within $\epsilon = 0.02$ of each other; 3)
131 the gain at neighboring harmonics (order k-1 and k+1) is lower than the threshold $G_{min}$
132 (resolvability criterion). For each selected pair of channels, we connect the corresponding input
133 neurons to a single coincidence detector neuron. The connection from the neuron with higher CF
134 has an axonal delay $\delta = \Delta\phi/kf_0$, where $\Delta\phi$ is the phase difference between the two filters at k.f0
135 (which is known analytically for a gammatone (Zhang et al., 2001)). In addition, for each
136 channel, multiple neurons receiving inputs from the same filter project to a single coincidence
137 detector neuron with axonal delays $\delta = k/f_0$ (as in Licklider's model), where k in the integer
138 varying between 1 and a value determined by the maximum delay $\delta_{max}$.

139 **Sounds**

140 *Musical instruments*

141 To test the neural model in a pitch recognition task, we used recordings of musical instruments
142 and vowels from the RWC Music Database (Musical Instrument Sound), including 762 notes
143 between A2 and A4, 41 instruments (587 notes) and 5 sung vowels (175 notes). Notes were
144 gated by a 10 ms cosine ramp and truncated after 500 ms.

145 *Environmental noises*

4

146     We also used a set of 63 environmental sounds containing stationary noises including: airplanes,
147     thunderstorm, rain, water bubbles, sea waves, fire and street sounds (recordings obtained from
148     [www.freesound.org)](www.freesound.org). We selected 500 ms segments from these sounds, gated by a 10 ms cosine
149     ramp.

150     **Analytical description of the auditory nerve phase response**

151     To analyze the discriminability of cross-channel structure (Fig. 6E-F), we fitted an analytical
152     formula to the phase $\phi(L, f, CF)$ of auditory nerve responses recorded at different levels L and
153     tone frequencies f in fibers with different CF, using a data set from Palmer and Shackleton
154     (2009) (Palmer and Shackleton, 2009), similarly to Carlyon et al. (Carlyon et al., 2012). For each
155     level, we fitted a function corresponding to the phase response of a gammatone filter bank:

$$\phi(L, f, CF) = f\psi(L, CF) + n \arctan\big(2\pi\tau(CF, L)(f - CF)\big)$$

156     where $\psi(L, CF)$ is the initial delay of the travelling wave (a parameterized function of CF,
157     equation (3) in (Zhang et al., 2001)), n in the order of the gammatone filter and $\tau(CF, L) =$
158     $\alpha(L)CF^{\beta(L)}$ is inversely related to the bandwidth of the filter.

159     We also tested another function: $\phi(L, f, CF) = \alpha(L, f) + \beta(L, f) \arctan\big(CF/\gamma(L, f)\big)$ as in
160     Carlyon et al. (Carlyon et al., 2012), where $\alpha, \beta$ and $\gamma$ were second-order polynomial functions of
161     L and f. The fits gave similar results.

162     **Discriminability of cross-channel and within-channel structure**

163     We used signal detection theory (GREEN and SWETS, 1966) to estimate the discriminability of
164     tone frequency based on regularity structure, using only phase information (to simplify). We
165     consider two places on the cochlea tuned to frequencies $f_A$ and $f_B$. A tone of frequency f is
166     detected when the two waveforms at places A and B are in phase after a delay d is introduced in
167     channel B: $\phi(f_B, f) + fd = \phi(f_A, f) + n$ , where $n$ is an integer (phases are expressed in cycles).
168     Note that n is related to related to the maximum delay $\delta_{max}$ (when f < 1/ $\delta_{max}$, there is at most
169     one possible value for n).

170     We note $\Delta\phi_{AB}(f) = \phi(f_B, f) - \phi(f_A, f)$ the phase difference between the two places (before the
171     delay is introduced), so that the equation reads:

172     $$\Delta\phi_{AB}(f) + f\delta = n \qquad\qquad (3)$$

173     That is, the phase difference after the delay is introduced is 0 cycle. When a tone of frequency
174     f+df is presented, the phase difference after the delay is introduced is $\Delta\phi_{AB}(f + df) +$
175     $(f + df)\delta = \Delta\phi_{AB}(f) + f\delta + (\Delta\phi'_{AB}(f) + \delta). df = n + (\Delta\phi'_{AB}(f) + \delta). df$ . Thus, a frequency
176     shift of df induces a phase shift of $(\Delta\phi'_{AB}(f) + \delta). df$ between the two channels, after
177     introduction of the delay.

178     We consider that neurons corresponding to channels A and B fire spikes in a phase-locked
179     manner with precision $\sigma$ (standard deviation of spike phase). Then the discriminability index d'
180     is the mean phase shift divided by the precision:

$$d' = \frac{(\Delta\phi'_{AB}(f) + \delta). df}{\sigma}$$

5

181 The just-noticeable difference (JND) for 75% correct discrimination is then:

$$JND = 1.35 \frac{\sigma}{\Delta\phi'_{AB}(f) + \delta}$$

182 The Weber fraction is JND/f. For two identical channels (within-channel structure), $\delta$ =1/f and
183 the formula simplifies to:

$$JND_{75\%} = 1.35\sigma f$$

184 For distinct channels (cross-channel structure), d is determined by equation (3), and the formula
185 reads:

$$JND_{75\%} = 1.35 \frac{\sigma f}{(f.\Delta\phi'(f) + n - \Delta\phi_{AB}(f))}$$

186 Finally, we relate phase precision with vector strength VS using the following formula, based on
187 the assumption that phases are distributed following a wrapped-normal distribution:

$$\sigma = \sqrt{-\ln(VS^2)}\,/2\pi$$

188

## Results

### The proposition

191 In the cochlea, the BM vibrates in response to sounds. We denote by S(x,t) the displacement of
192 the BM at time t and place x. This displacement is represented in Fig. 1A as the output of a
193 gammatone filterbank with bandwidth based on psychophysical measurements (see Methods).
194 Each auditory nerve fiber transduces the temporal vibration S(x,t) at a specific place into a spike
195 train. In Licklider's delay line model (the classical temporal model (Licklider, 1951)), the
196 periodicity of the mechanical vibration is detected by a coincidence detector neuron receiving
197 synaptic inputs from a single cochlear place x. It fires when it receives coincidences between a
198 spike train produced by a fiber originating from that place and the same spike train delayed by a
199 fixed amount δ (Fig. 1B). Conceptually, this neuron detects the identity S(x,t+δ)=S(x,t) for all t,
200 that is, the fact that S(x,.) is periodic with period T = δ.  This mechanism must be slightly
201 amended to account for the refractory period of fibers, which sets a lower limit to the period that
202 can be detected. This issue can be addressed by postulating that the neuron receives inputs from
203 two different fibers originating from the same place (Fig. 1C).

204 We now consider the possibility that these two fibers may originate from slightly different
205 cochlear places x and y. In this case, the neuron detects the identity S(y,t+δ)=S(x,t), that is,
206 similarity of sensory signals across both place and time (Fig. 1D). We note in this example (a
207 harmonic sound) that the delay δ may now be different from the period T of the vibration.
208 Compared to the detection of periodicity, this does not require any additional anatomical or
209 physiological assumption. Thus we propose to examine the proposition that pitch is the
210 perceptual correlate of the regularity structure of the BM vibration pattern, across both time and
211 place, defined as the set of identities of the form S(x,t)=S(y,t+δ) for all t. A few previous models

212    of pitch also use cross-channel comparisons (Loeb et al., 1983; Shamma, 1985; Carney et al.,
213    2002), and we will relate them to our theory in the discussion.

214    To illustrate our proposition, Fig. 1E-F show the cochleograms obtained by filtering two sounds
215    with a gammatone filterbank. A noise-like sea wave (Fig.1E) produces no regularity structure in
216    the cochleogram, that is, there are no identities S(x,t)=S(y,t+δ) in the signals. A clarinet note, on
217    the other hand, produces a rich regularity structure (Fig. 1F). Because this is a periodic sound,
218    the BM vibrates at the sound's period T at all places (or more generally T/k, where k is an
219    integer), as shown by horizontal arrows: S(x,t+T)=S(x,t) for all t and x. We call this set of
220    identities the *within-channel structure*. More interestingly, we also observe identities across
221    places, as shown by oblique arrows: S(x,t)=S(y,t+δ) for all t. These occur for specific pairs of
222    places x and y, which tend to be in low frequency regions. We note that the time shift δ is
223    different from the sound's period T. We call this set of identities the *cross-channel structure*.

224

225    **Resolvability and regularity structure**

226    We now examine the type of regularity structure produced by sounds. First, if the sound is
227    periodic, then the BM vibrates at the sound's period T at all places, provided there is energy at
228    the corresponding frequency. That is, S(x,t+T)=S(x,t) for all x and t. Conversely, the identity
229    S(x,t+T)=S(x,t) means that the BM vibrates periodically, which can only occur if the sound itself
230    is periodic, at least within the bandwidth of the cochlear filter at place x. Thus, within-channel
231    structure is simply the periodicity structure at each cochlear place.

232    Cross-channel structure is less trivial. What kind of sound produces the same vibration (possibly
233    delayed) at different places of the cochlea? To simplify the argument, we consider that cochlear
234    filters are linear (we come back to this point in the discussion), and we examine the identity
235    S(x,t)=S(y,t+δ) in the frequency domain. If the two signals at place x and y match, then all their
236    frequency components must match, both in phase and amplitude. But these two signals originate
237    from the same sound, filtered in two different ways. Fig. 2A shows the gain (left) and phase
238    (right) of the two filters A and B as a function of frequency. The only way that a frequency
239    component is filtered in the same way by the two filters is that the gains are identical at that
240    frequency, which happens in this case at a single frequency f (illustrated on Fig. 2A, bottom).
241    Additionally, the phases of the two filters must match at frequency f, taking into account the
242    delay δ. That is, the phase difference $\Delta\phi$ must equal f.δ (modulo 1 cycle).

243    In summary, the only type of sound that produces cross-channel structure is a sound with a
244    single frequency component within the bandwidth of the two considered cochlear filters. This is
245    a notion of *resolvability*, and we will say that the frequency component is *resolved* with respect to
246    the pair of filters. Fig. 2B illustrates what happens when a periodic sound with unresolved
247    harmonics is passed through the two filters. Here the output of filter A is a combination of
248    harmonics k and k-1, while that of filter B is a combination of harmonics k and k+1. Therefore,
249    the two resulting signals are different (bottom): there is no cross-channel structure.

250    Thus, the amount of cross-channel structure produced by a harmonic sound depends on the
251    resolvability on its frequency components. Fig. 2C shows the amplitude spectrum of a periodic
252    sound with all harmonics k.f0 (bottom). Because harmonics are linearly spaced but cochlear
253    filter bandwidth increases with frequency (filter amplitude in gray), the excitation pattern of the

7

254  BM as a function of center frequency (top) shows distinct peaks for low-order harmonics (which
255  are thus considered "resolved") but not for high-order harmonics (unresolved). More precisely,
256  low-order harmonics are resolved for many pairs of cochlear filters, meaning that they produce
257  cross-channel structure for many filter pairs (Fig. 2D, left); high-order harmonics produce little
258  or no cross-channel structure (Fig. 2D, right). The amount of cross-channel structure is directly
259  determined by the spacing between frequency components (f0) relative to the cochlear filter
260  bandwidth. With the approximation that filter bandwidth is proportional to center frequency
261  (k.f0 if centered at the $k^{th}$ harmonic), this means that the amount of cross-channel structure is
262  determined by the harmonic number k. Therefore, there is a direct relationship between
263  resolvability defined in a conventional sense and the amount of cross-channel structure
264  produced by the sound.

265  Figure 2E illustrates this point with a resolved harmonic complex consisting of resolved
266  components (left) and with an unresolved harmonic complex (right). Both sounds produce
267  within-channel structure (horizontal arrows), but the resolved complex additionally produces
268  cross-channel structure. Thus, the structural theory attributes a pitch to all periodic sounds, but
269  the amount of regularity structure, and therefore of information about f0, depends on
270  resolvability. It follows in particular that discrimination of f0 based on regularity structure
271  should be more precise for resolved than unresolved sounds (Houtsma and Smurzynski, 1990;
272  Carlyon and Shackleton, 1994; Carlyon, 1998; Bernstein and Oxenham, 2003), since there is
273  more information (the exact quantitative assessment would depend on the specific estimator
274  chosen).

275

276  **The domain of existence of pitch**

277  From the definitions above, the set of sounds that produce regularity structure is exactly the set
278  of periodic sounds. However, perceptually, not all periodic sounds have a melodic pitch. In
279  particular, pitch only exists for f0 between 30 Hz (Pressnitzer et al., 2001) and 5kHz (Semal and
280  Demany, 1990). Within this range, periodic sounds may or may not have a clear pitch, depending
281  on their harmonic content. In the structural theory, the domain of existence of pitch is restricted
282  when we impose constraints on the comparisons between signals (cross- or within-channel)
283  that the auditory system can do. Two physiological constraints seem unavoidable: 1) there is a
284  maximum time shift $\delta_{max}$ (possibly corresponding to a maximum neural conduction delay), 2)
285  temporal precision is limited (possibly corresponding to phase locking precision). We may also
286  consider that there is a maximum distance along the BM across which signals can be compared,
287  but it will not play a role in the discussion below. The temporal precision sets an upper limit to
288  pitch, exactly in the same way as in standard temporal theories. Thus we shall restrict our
289  analysis to the constraint of a maximum delay $\delta_{max}$. We consider the simplest possible
290  assumption, which is a constant maximal delay $\delta_{max}$, independent of frequency.

291  We start by analyzing the domain of existence of within-channel structure (Fig. 3A). Since this is
292  just the periodicity structure, its domain of existence is the same as in standard temporal
293  theories of pitch. When the sound's period exceeds the maximum delay $\delta_{max}$, periodicity cannot
294  be detected anymore. Therefore, the lower limit (minimum f0) is the inverse of the maximum
295  delay: f0 = $1/\delta_{max}$.

296 A different limit is found for cross-channel structure, because the delay δ between signals across
297 channels is not the same as the sound's period (see e.g. Fig. 1F). In fact, this delay can be
298 arbitrary small, if the two places are close enough on the BM. Figure 3B shows an example of a
299 100 Hz pure tone passed through two filters A and B. The gains of the two filters are the same at
300 100 Hz and there is a phase difference of 8/10 cycle, which is equivalent to -2/10 cycle. As a
301 result, the output of the two filters is a pair of tones with identical amplitude and delay δ = 2 ms
302 (2/10 of 10 ms), much smaller than the sound's period. This delay would be even smaller if the
303 center frequencies of the two filters were closer. Thus the lower limit of cross-channel structure
304 is not set by the maximum delay $\delta_{max}$. Instead, it is set by the center frequencies of the filters.
305 Indeed the frequency of the tone (or resolved harmonic) must lie between the two center
306 frequencies of the filters, and therefore the lowest such frequency corresponds to the lowest
307 center frequency of cochlear filters. This minimum frequency is not known in humans, but the
308 lower limit of the hearing range is about 20 Hz, which suggests a lower limit of cross-channel
309 structure slightly above 20 Hz. This is consistent with psychophysical measurements of the
310 lower limit of pitch, around 30 Hz for tones (Pressnitzer et al., 2001).

311 Therefore, the structural theory of pitch predicts different lower limits of pitch depending on
312 whether the sound contains resolved harmonics or not. When it does, the lower limit is
313 determined by cross-channel structure, and thus by the lowest center frequency of cochlear
314 filters, on the order of a few tens of Hz. When it does not, the lower limit of pitch is determined
315 by within-channel structure, and is thus $1/\delta_{max}$. We now compare these theoretical predictions
316 with two recent psychophysical studies. In Oxenham et al. (2004) (Oxenham et al., 2004a),
317 transposed stimuli were created by modulating a high frequency carrier (>4 kHz) with the
318 temporal envelope of a half-wave rectified low frequency tone (<320 Hz) (Fig. 3C, top). Human
319 subjects displayed poor pitch perception for these stimuli, even though the repetition rate $f_0$ was
320 in the range of pitch perception for pure tones. This finding poses a challenge for temporal
321 theories, but is consistent with the structural theory, as is illustrated in Fig. 3C. Indeed, these
322 transposed tones do not contain resolved harmonics, and therefore only produce within-channel
323 structure (horizontal arrows in Fig. 3C). As described above, the lower limit of pitch is $1/\delta_{max}$ in
324 this case. If this maximal delay is $\delta_{max} < 3$ ms, then transposed tones do not produce a pitch when
325 the frequency of the tone is lower than 330 Hz. On the other hand, for pure tones, the lower limit
326 of pitch is much lower than 330 Hz because of the presence of cross-channel structure (oblique
327 arrows in Fig. 3D). In Oxenham et al. (2011) (Oxenham et al., 2011), it was shown that complex
328 tones with f0 between 400 Hz and 2 kHz and all harmonics above 5 kHz elicit a pitch. In the
329 structural theory, all periodic sounds with f0 > $1/\delta_{max}$ produce a pitch, irrespective of their
330 harmonic content. This is shown in Fig. 3E, which shows the cochlear filter responses to a
331 complex tone with $f_0 = 1.2$ kHz and all harmonics above 5 kHz. Therefore, this psychophysical
332 study is consistent with the structural theory if $\delta_{max} > 2.5$ ms. In summary, both psychophysical
333 studies are consistent with the structural theory if $\delta_{max}$ is on the order of 3 ms.

334

335 **A possible neural mechanism**

336 We now propose a possible neural mechanism to estimate f0 based on the vibration structure of
337 the BM. Since the theory is based on similarity between signals, the same mechanism as for
338 temporal models can be suggested. A straightforward generalization of Licklider's model
339 (Licklider, 1951) is illustrated in Fig. 1D: a neuron receives inputs from two presynaptic neurons

340 (X and Y), which encode the BM vibration at two cochlear locations x and y in precisely timed
341 spike trains, and there is a mismatch $\delta$ in their conduction delays. We assume that the
342 postsynaptic neuron responds preferentially when it receives coincident input spikes. Indeed,
343 neurons are highly sensitive to coincidences in their inputs, under broad conditions (Rossant et
344 al., 2011). By acting as a coincidence detector, the postsynaptic neuron signals a particular
345 identity $S(y, t + \delta) = S(x, t)$.

346 Anatomically, neurons X and Y could be auditory nerve fibers and the postsynaptic neuron could
347 be in the cochlear nucleus. Alternatively, neurons X and Y could be primary-like neurons in the
348 cochlear nucleus, for example spherical bushy cells, and the postsynaptic neuron could be in the
349 inferior colliculus or in the medial superior olive. Indeed, as demonstrated in Fig. 4A-B, the
350 synchrony between two neurons depends on the similarity between the signals they encode,
351 rather than on their specific cellular properties. Fig. 4A shows the cochleogram of a trumpet note
352 with f0 = 277 Hz (top). The red and blue boxes highlight the BM vibration at characteristic
353 frequencies 247 Hz and 307 Hz, around the first harmonic. This harmonic produces cross-
354 channel similarity with delay $\delta$, as seen on the red and blue signals shown below (grey shading
355 is the mismatch). As a result, neurons that encode these two signals into spike trains fire in
356 synchrony, as is shown below for three different models: a biophysical model of a type Ic
357 chopper neuron (Rothman and Manis, 2003b), a type II model of an octopus cell, and a leaky
358 integrate-and-fire model. In contrast, when an inharmonic sound is presented, such as a rolling
359 sea wave (Fig. 4B), the inputs do not match and neural responses are not synchronous, for any of
360 the three models.

361 The same mechanism applies for within-channel structure. In Fig. 4C, we consider two high-
362 frequency neurons with the same characteristic frequency CF = 2700 Hz but a delay mismatch $\delta$
363 = 4.5ms. When a periodic sound with repetition rate 220 Hz is presented (here a harpsichord
364 note), their input signals match, which results in synchronous discharges. We note that not all
365 output spikes are coincident. This occurs because the neurons discharge in more complex
366 spiking patterns (Laudanski et al., 2010) and do not fire one spike per cycle: they may miss a
367 cycle or fire several times in one cycle. Nevertheless, coincidences of output spikes occur much
368 less often with an inharmonic sound (Fig. 4D). This mechanism is analog to Licklider's model
369 (Licklider, 1951), in which each neuron signals a particular identity $S(x, t + \delta) = S(x, t)$. Thus
370 the neural mechanism we describe is simply an extension of Licklider's model to cross-channel
371 similarity.

372 As a proof of concept, we now build a simple neural model that estimates f0 by detecting
373 regularity structure. For each f0 between notes A2 and A4 (110 Hz to 440 Hz), we build a group
374 of coincidence detector neurons, one for each similarity identity $S(y, t + \delta) = S(x, t)$ that is
375 present for sounds with that particular f0. To this aim, we examine the BM response (modeled as
376 gammatone filters) to a complex tone with all harmonics n.f0 (Fig. 4E, red comb on the left). On
377 Fig. 4E-F, we represent the BM response using color disks arranged as a function of cochlear
378 location (vertical axis) and delay (horizontal axis): color saturation represents the amplitude of
379 the filter output while hue represents its phase. For low-order harmonics (resolved, bottom), the
380 BM vibrates as a sine wave and therefore disks with the same color correspond to identical
381 signals, and thus to encoding neurons firing in synchrony. For high-order harmonics
382 (unresolved, top), the BM vibrates in a more complex way and there only identically colored
383 disks within the same channel correspond to identical signals. We then set synaptic connections
384 from neurons encoding the same BM signal to a specific coincidence detector neuron (all

385 modeled as integrate-and-fire neurons). Thus we obtain a group of neurons that fire
386 preferentially when the identities $S(y, t + \delta) = S(x, t)$ corresponding to a particular f0 occur
387 (note that we have omitted a number of possible identities for simplicity, e.g. cross-channel
388 identities occurring with high frequency pure tones). In this way, the mean firing rate of the
389 group of neurons is tuned to f0.

390 We iterate this construction for every f0 between A2 and A4 (by semitone steps). As illustrated
391 in Fig. 4F, a different f0 produces a different regularity structure (colored disks), from which we
392 build a different set of synaptic connections to the pitch-tuned group of coincidence neurons
393 (one group per f0). To estimate f0, we then simply look for the pitch-tuned group with the
394 highest mean firing rate.

395 We presented two types of natural sounds to this model (spectrograms shown in Fig. 5A, top):
396 inharmonic sounds (e.g. an airplane, a sea wave and street noise), and harmonic sounds (e.g.
397 clarinet, accordion and viola) with f0 between A2 and G#4. For each sound, we measure the
398 average firing rate of all pitch-tuned neuron groups (Fig. 5A, bottom). Inharmonic sounds
399 generally produce little activation of these neurons, whereas harmonic sounds activate specific
400 groups of neurons (with some octave confusions, see below). In Fig. 5A, musical notes were
401 played in chromatic sequence, which appears in the response of pitch-tuned groups. Fig. 5B
402 shows the distribution of group firing rates, measured in the entire neuron model, for
403 inharmonic (grey) and harmonic sounds (blue), at three different sound levels. Although an
404 increase in sound level produces an overall increase in population firing rate, there is little
405 overlap between the rate distributions for harmonic and inharmonic sounds.

406 From the activity of these neurons, we estimate the pitch of a presented harmonic sound as the
407 pitch associated to the maximally activated group of neurons. This estimation was correct in
408 77% of cases, and was within one semitone of the actual pitch in 88% of cases (Fig. 5C, top).
409 Most errors greater than one semitone correspond to octaves or fifths (octaves: 5.5%, fifth:
410 <2%), which also shows in the distribution of firing rate of pitch-tuned groups (Fig. 5C, bottom).
411 This performance was obtained with 400 frequency channels spanning 50 Hz to 8 kHz, and it
412 degrades if the number of channels is reduced (e.g. 35% score for N = 100, Fig. 5D, top), because
413 the model relies on comparisons between neighboring channels. We then tested how
414 performance was affected by constraints on the maximum delay (Fig. 5D, bottom). We found no
415 difference in performance when maximum delay $\delta_{max}$ was varied between 2 and 15 ms. The
416 highest f0 in our sound database was 440 Hz (A4), which corresponds to a period greater than 2
417 ms. Therefore with $\delta_{max}$ = 2 ms, the model reached the same level of performance with only
418 cross-channel comparisons.

419

### 420 **Pitch discriminability**

421 Finally, we examine the discriminability of pure tones based on regularity structure. To simplify,
422 we ignore amplitude differences and focus on phase differences between channels. We start with
423 within-channel structure and consider two neurons (e.g. auditory nerve fibers) encoding BM
424 vibration from the same place x (i.e., same characteristic frequency) into phase-locked spike
425 trains, with a delay mismatch $\delta$ = 1/f (Fig. 6A). These two neurons fire in synchrony when a pure
426 tone of frequency f is presented. More precisely, given that there is some stochasticity in neural
427 firing, the two neurons produce spikes with the same mean phase relative to the tone, so the

428    difference of phases of spikes ΔΦ(f) is distributed around 0 (Fig. 6A, left). When a tone of
429    frequency f+df is presented, ΔΦ(f) shifts by an amount of δ.df = df/f (Fig. 6A, right).

430    The same analysis applies for cross-channel structure, where the two neurons encode BM
431    vibration at two different places A and B (different CFs, Fig. 6B). Here the delay δ is related to the
432    mismatch in phase response at the places at tone frequency f. When a tone of frequency f+df is
433    presented, ΔΦ(f) shifts because of both the delay and the relative change in response phase at
434    the two places on the BM (see Methods).

435    Thus, discriminating between tones of nearby frequencies corresponds to discriminating
436    between two circular random variables ΔΦ(f) and ΔΦ(f+df) with different means, which can be
437    analyzed with signal detection theory (GREEN and SWETS, 1966). Specifically, the
438    discriminability index d' is the mean phase shift μ divided by the precision σ (standard deviation
439    of phase) (Fig. 6C). The precision of phase locking is often measured by the vector strength (VS),
440    which is relatively independent of frequency below a critical frequency above which it decays
441    rapidly to 0 (Fig. 6D, guinea pig auditory nerve). We estimate the standard deviation σ from VS
442    assuming a wrapped normal distribution (see Methods). To estimate μ, we used spike trains
443    recorded in guinea pig auditory nerve fibers with different CFs in response to tones with various
444    frequencies (Palmer and Shackleton, 2009) and estimated the average spike phase as function of
445    both CF and tone frequency (see Methods) (Fig. 6E).

446    We used these estimates to calculate the just-noticeable difference (JND) for 75% correct
447    discrimination, which is the frequency change df producing a discriminability index d' = 1.35.
448    Figure 6F shows the JND relative to tone frequency (JND(f)/f), called the Weber fraction, as a
449    function of tone frequency, for within-channel structure (black) and for cross-channel structure
450    (colors), for pairs of channels varying by their spacing in CF (1 to 6 semitones). For both types of
451    structure, the Weber fraction increases in high frequency because of the loss of phase locking
452    (VS goes to 0). The two types differ in the low-frequency end: while the Weber fraction is
453    independent of frequency for within-channel structure, it tends to increase with lower frequency
454    for cross-channel structure. We also note that discriminability is better for widely spaced
455    channels (orange) than for neighboring channels (blue), but the former require larger delays.

456

## Discussion

458    We have proposed that pitch is the perceptual correlate of the regularity structure of the BM
459    vibration pattern, defined as the set of identities of the form S(x,t)=S(y,t+δ) for all t, where S(x,t)
460    is the displacement of the BM at time t and place x. The regularity structure generalizes the
461    notion of periodicity. This proposition assigns a pitch to periodic sounds and therefore has many
462    similarities with the standard view that pitch is the perceptual correlate of the periodicity of the
463    acoustic waveform. However, it also predicts that resolved harmonic complexes elicit a stronger
464    pitch than unresolved harmonic complexes (richer structure), and it predicts a complex region
465    of existence of pitch that depends on harmonic content. In particular, it predicts that high
466    frequency complex tones only elicit a clear pitch if f0 is high, in agreement with experiments
467    (Oxenham et al., 2004b, 2011). Finally, it does not rely on the existence of long conduction delays
468    in the auditory system.

469 Previous studies have proposed mechanisms to extract the fundamental frequency of either
470 resolved or unresolved harmonic complexes (see detailed discussion in section "Related theories
471 of pitch" below). Some share common ideas with our proposition: for example, classical
472 temporal models address the extraction of within-channel periodicity (S(x,t) = S(x,t+T)) (de
473 Cheveigné, 2010), which does not distinguish between resolved and unresolved components;
474 other authors have proposed that the frequency of resolved components can be estimated with
475 cross-channel comparisons or operations (Loeb et al., 1983; Shamma, 1985; Carney et al., 2002).
476 These ideas are also present in our proposition. However, instead of proposing a particular
477 mechanism to extract f0, we propose that pitch is not the correlate of the periodicity of the
478 sound waveform but of the regularity structure of the BM vibration pattern (with a limited
479 temporal window). The main implications for pitch perception (as shown in Fig. 3) are to a large
480 extent independent of the particular mechanism that extracts that structure. In particular, this
481 single proposition implies that resolved and unresolved harmonic complexes have different
482 perceptual properties.

483

484 **Neural mechanism**

485 A separate issue is the physiological implementation of this theory, that is, how pitch defined
486 according to the regularity structure of the BM vibration pattern might be estimated by the
487 auditory system. There are different ways in which the auditory system might extract that
488 information. It may also be the case that pitch is not conveyed by the increased firing of pitch-
489 tuned neurons but by temporal relationships in their firing (Cariani, 2001). Here we have simply
490 made a suggestion of a possible mechanism that makes minimal physiological assumptions. But
491 we stress that our core proposition does not rely on a particular mechanism, but on the
492 regularity structure of the BM vibration. The most straightforward implementation is a
493 generalization of Licklider's delay line model (Licklider, 1951), in which a pitch-selective neuron
494 detects coincidences between two inputs with different axonal conduction delays. In the original
495 model, the two inputs originate from the same place in the cochlea. An implementation of the
496 structural theory is obtained simply by allowing the two inputs to originate from slightly
497 different places. If a neural circuit resembling Licklider's model indeed exists in the auditory
498 brainstem, then it is plausible that inputs to these coincidence detector neurons are not exactly
499 identical. Because our proposition relies on the temporal fine structure of sounds, the matching
500 mechanism between the outputs of two channels (whether it is based on coincidence detection
501 or not) should occur early in the auditory periphery. Input neurons could be auditory nerve
502 fibers and the coincidence detector neuron could be in the cochlear nucleus. Alternatively, input
503 neurons could be primary-like neurons in the cochlear nucleus, for example spherical bushy
504 cells, and the coincidence detector neuron could be in the inferior colliculus or in the medial
505 superior olive (MSO). The latter possibility has some appeal because neurons in the MSO are
506 thought to receive few synaptic inputs (Couchman et al., 2010) and are known to act as
507 coincidence detectors (Yin and Chan, 1990), although possibly not monaurally (Agmon-Snir et
508 al., 1998), and there are cases of binaural pitch for sounds that have no monaural structure. In
509 the inferior colliculus, there is some physiological evidence of tuning to pitch (Langner, 1992).
510 Specifically, in a number of mammalian species, IC neurons are tuned in their firing rate to the
511 modulation frequency of amplitude-modulated tones, up to about 1000 Hz, independently of
512 their characteristic frequency, although the best modulating frequency may depend on carrier

513  frequency. There is also some evidence of a topographic organization of periodicity tuning,
514  orthogonal to the tonotopical organization.

515  As a proof of principle, we have shown with a spiking neural model that such a mechanism can
516  indeed estimate the pitch of harmonic sounds, even with short conduction delays. Standard
517  temporal models of pitch have been criticized because they require long delays for low f0, up to
518  30 ms for the lowest pitch (Pressnitzer et al., 2001). There is no experimental evidence of such
519  long axonal delays in the auditory brainstem. In a recent anatomical study of axons of spherical
520  bushy cells in cats (cochlear nucleus projections to the MSO) (Karino et al., 2011), the range of
521  axonal delays was estimated to be just a few hundred μs, far from the required 30 ms (although
522  these were anatomical estimates, not functional measurements). This range could be larger in
523  humans as axons are presumably longer, but it could also be similar if axonal diameter scales in
524  the same way (since conduction speed is approximately proportional to diameter in myelinated
525  axons (Rushton, 1951)). In either case, the range of axonal delays is unlikely to be much greater
526  than a few ms. Another possibility is to consider dendritic propagation delays or intrinsic delays
527  induced by ionic channels. These could contribute additional delays, but the duration of
528  postsynaptic potentials measured at the soma of auditory brainstem neurons tends to be short
529  (Trussell, 1997, 1999), which makes this scenario rather implausible for large delays. We have
530  shown that the structural theory is compatible with psychophysical results when the delays are
531  limited to a few ms, and the neural mechanism based on coincidence detection remains
532  functional even for low f0.

533

534  **Related theories of pitch**

535  Two previous propositions are directly related to the structural theory. Loeb et al. (Loeb et al.,
536  1983) proposed that the frequency of a pure tone can be estimated by comparing signals across
537  the BM: the distance that separates places that vibrate in phase is indeed related to the tone's
538  frequency. This is a special case of the structural theory, when the maximal delay is 0 ms (i.e.,
539  identities of the form $S(x,t) = S(y,t)$ for all t). However, this proposition restricts pitch to resolved
540  harmonic complexes only, and in fact to complexes made of widely separated tones.

541  The phase opponency model (Carney et al., 2002) is a similar proposition, in which a tone of a
542  particular frequency is detected when signals at two different places on the BM are out of phase.
543  This corresponds to detecting identities of the form $S(x,t) = -S(y,t)$ for all t. This model suffers
544  from the same problem as Loeb's model, that is, it applies to a limited subset of pitch-evoking
545  sounds.

546  We may also consider a variation of the structural theory, in which amplitude is discarded (as
547  we did when analyzing frequency discrimination). This variation corresponds to considering
548  identities of the form $S(x,y) = a.S(y,t+\delta)$ for all t. This variation has the same qualitative
549  properties as the original formulation, and is physiologically motivated by the observation that
550  low threshold AN fibers saturate quickly when intensity is increased (Sachs and Abbas, 1974).

551  Place theories of pitch are based on the comparison of internal templates with the spatial
552  pattern of BM vibration encoded in the firing of auditory nerve fibers. A weakness of these
553  theories is that the firing rate of auditory nerve fibers as well as of most neurons in the cochlear
554  nucleus saturate at high levels (Sachs and Young, 1979; Cedolin and Delgutte, 2005). To address

14

555  this problem, it has been proposed that the spatial profile is first sharpened by lateral inhibition,
556  prior to template matching (Shamma, 1985). This preprocessing step enhances the responses at
557  places where the phase changes rapidly, which occurs where the BM is tuned to the sound's
558  frequency. A recent analysis of cat auditory nerve responses has shown that such preprocessing
559  produces spatial profiles from which f0 can indeed be extracted even at high levels (Cedolin and
560  Delgutte, 2010), although a more recent analysis (in guinea pigs and with different methods)
561  suggested that the estimated f0 is very sensitive to level (Carlyon et al., 2012). Because this
562  preprocessing step relies on temporal cues, template-based models of pitch using this stage as
563  input are often described as spatiotemporal models (Cedolin and Delgutte, 2010). However,
564  these are very different from the structural theory we have presented, as they are in fact models
565  based on matching spatial templates where temporal information is discarded, only with an
566  input that is obtained from a spatiotemporal transformation of the auditory nerve response. In
567  contrast, matching in the structural theory as well as in the two related models mentioned above
568  and in standard temporal models is performed on the entire temporal signals.

569  Unlike the structural theory, none of these three models addresses the pitch of unresolved
570  harmonic complexes.

571

572  **The nature of pitch in theories of pitch**

573  In standard temporal theories of pitch, pitch is the perceptual correlate of the periodicity of the
574  acoustical waveform. Independently of how the periodicity is physiologically extracted, this
575  proposition implies for example that: periodic sounds have a pitch, non-periodic sounds do not
576  have pitch, and pitch saliency is related to how close to periodic a sound is. It also implies that
577  two sounds with the same periodicity are similar, and that two sounds with fundamental
578  frequencies differing by an octave are similar, in the sense that they have a periodicity in
579  common. Thus, this characterization of pitch entails a particular region of existence of pitch
580  (what sounds produce pitch) and a particular topology of pitch  (how pitch-evoking sounds
581  relate to each other). These two aspects do not rely on learning, in the sense that they do not
582  depend on the specific sounds the auditory system is exposed to. Instead, they derive from the
583  existence of a general mechanism that identifies periodicity.

584  In a similar way, the structural theory of pitch defines pitch as the perceptual correlate of the
585  regularity structure of the BM vibration pattern. It also entails an existence region of pitch,
586  which is more complex than in temporal theories, and a particular topology of pitch, which is
587  similar to that implied by temporal theories (but see below for the effect of level on pitch). In the
588  same way, these two aspects do not rely on learning.

589  In standard place theories of pitch based on templates, what characterizes pitch-evoking sounds
590  is that they are similar to some internal template (Terhardt, 1974). Thus pitch is the perceptual
591  correlate of a particular category of sounds, which is formed by previous exposure to pitch-
592  evoking sounds. There is an obvious problem of circularity in this characterization, which means
593  that in addition to exposure to the sounds, these sounds must be labeled as having or not having
594  a pitch. That is, pitch is characterized independently of the sounds themselves. An example
595  would be that vocalizations are those special sounds that are considered as producing pitch.
596  Accordingly, a more rigorous characterization of pitch in place theories is the following: pitch is

15

597 the perceptual correlate of spectral similarity to vocalizations (or any other externally defined
598 category of sounds).

599 This characterization is problematic for several reasons. First, it defines an existence region of
600 pitch but not a topology of pitch, unless the spatial activation profiles produced by sounds with
601 the same pitch are similar. This issue might be addressed to some extent by spatial sharpening
602 as previously mentioned (Shamma, 1985), although there is no indication that such an operation
603 occurs in the auditory system. A second problem is that not all pitch-evoking sounds are
604 spectrally similar to vocalizations, for example low-frequency pure tones. Finally, infants have a
605 sense of musical pitch (Montgomery and Clarkson, 1997). The latter two issues have been
606 addressed in a model in which harmonic templates are learned from inharmonic sounds
607 (Shamma and Klein, 2000). Indeed auditory nerve fibers with harmonically related CFs are
608 expected to fire with some degree of correlation in response to noise, because of nonlinearities
609 in their response. Thus a Hebbian mechanism could form harmonic templates by selecting
610 temporally correlated fibers. In this scheme, pitch is then the perceptual correlate of the
611 similarity between the places of activation on the BM and places that are generally expected to
612 be correlated.

613 In addition to the fact that this only addresses the pitch of unresolved harmonic complexes, this
614 proposition is somehow paradoxical. On one hand, the formation of internal templates critically
615 relies on the temporal fine structure of the sounds, and fine correlations between channels.
616 Indeed in Hebbian models, the learning signal is the correlation between input and output (pre-
617 and postsynaptic neurons), and therefore it requires that the output firing is sensitive to input
618 correlations. On the other hand, pitch estimation by template matching assumes that this
619 temporal fine structure is then entirely discarded: only average spectrum is considered, and
620 correlations between channels (relative phases of harmonics in a complex tone) are assumed to
621 have no effect on pitch. To reconcile the two aspects of the model requires either that the
622 neurons are initially sensitive to input correlations and become insensitive to them after a
623 critical period (after learning), or that learning is based on input correlations but not through a
624 Hebbian mechanism (i.e., not involving input-ouput correlations).

625

626 **Experimental predictions**

627 We can formulate two types of predictions, for psychophysical experiments and for physiological
628 experiments. The strongest psychophysical prediction concerns the effect of level on pitch. The
629 phase of the BM response to tones depends on level (Robles and Ruggero, 2001), because of
630 nonlinear effects. Consequently, cross-channel structure should depend on level. However,
631 within-channel structure should not depend on level because such nonlinearities have no effect
632 on periodicity. If we assume that sounds are matched in pitch when they produce some common
633 regularity structure on the BM, then a pitch-matching experiment between sounds with different
634 levels should reveal an effect of level on the pitch of sounds that produce cross-channel structure
635 but not within-channel structure. According to our analysis, these are pure tones of low
636 frequency, i.e., with period larger than the maximum delay. The few studies on such effects
637 support this prediction (Morgan et al., 1951; Verschuure and Van Meeteren, 1975; Burns, 1982),
638 but a more exhaustive and controlled study would be required.

16

639 Predictions for physiological experiments can be made for specific hypotheses about the neural
640 mechanism. For example, low-frequency spherical bushy cells are primary-like neurons of the
641 cochlear nucleus with strong phase locking properties (Joris et al., 1994; Fontaine et al., 2013)
642 (possibly stronger than the auditory nerve), and their pattern of synchrony in response to
643 sounds could then reflect the regularity structure of the BM vibration. The prediction is then that
644 the synchrony receptive field of two such cells, defined as the set of sounds that produce
645 synchronous responses in the two cells (Brette, 2012), should consist of pitch-evoking sounds -
646 in fact of a pure tone of specific frequency. Ideally, such recordings should be done
647 simultaneously, because shared variability (e.g. due to local synaptic connections or shared
648 modulatory input) affects phase locking and reproducibility but not synchrony (Brette, 2012).

649

## References

651 Agmon-Snir H, Carr CE, Rinzel J (1998) The role of dendrites in auditory coincidence
652 detection. Nature 393:268–272.

653 Bernstein JG, Oxenham AJ (2003) Pitch discrimination of diotic and dichotic tone
654 complexes: Harmonic resolvability or harmonic number? The Journal of the Acoustical
655 Society of America 113:3323.

656 Boer E de, Jongh HR de (1978) On cochlear encoding: Potentialities and limitations of
657 the reverse-correlation technique. The Journal of the Acoustical Society of America
658 63:115–135.

659 Brette R (2012) Computing with neural synchrony. PLoS computational biology
660 8:e1002561.

661 Burns EM (1982) Pure-tone pitch anomalies. I. Pitch-intensity effects and diplacusis in
662 normal ears. The Journal of the Acoustical Society of America 72:1394.

663 Cariani PA (2001) Neural timing nets. Neural Netw 14:737.

664 Cariani PA, Delgutte B (1996a) Neural correlates of the pitch of complex tones. I. Pitch
665 and pitch salience. J Neurophysiol 76:1698.

666 Cariani PA, Delgutte B (1996b) Neural correlates of the pitch of complex tones. II. Pitch
667 shift, pitch ambiguity, phase invariance, pitch circularity, rate pitch, and the dominance
668 region for pitch. J neurophysiol 76:1717.

669 Carlyon RP (1998) Comments on "A unitary model of pitch perception" [J. Acoust. Soc.
670 Am. 102, 1811–1820 (1997)]. The Journal of the Acoustical Society of America
671 104:1118–1121.

672 Carlyon RP, Long CJ, Micheyl C (2012) Across-Channel Timing Differences as a Potential
673 Code for the Frequency of Pure Tones. JARO 13:159–171.

674 Carlyon RP, Shackleton TM (1994) Comparing the fundamental frequencies of resolved

675 and unresolved harmonics: Evidence for two pitch mechanisms? The Journal of the
676 Acoustical Society of America 95:3541–3554.

677 Carney LH, Heinzy MG, Evilsizer ME, Gilkeyz RH, Colburn HS (2002) Auditory phase
678 opponency: A temporal model for masked detection at low frequencies. Acta Acustica
679 United with Acustica 88:334–347.

680 Carney LH, Yin TC (1988) Temporal coding of resonances by low-frequency auditory
681 nerve fibers: single-fiber responses and a population model. J Neurophysiol 60:1653–
682 1677.

683 Cedolin L, Delgutte B (2005) Pitch of Complex Tones: Rate-Place and Interspike Interval
684 Representations in the Auditory Nerve. J Neurophysiol 94:347–362.

685 Cedolin L, Delgutte B (2010) Spatiotemporal Representation of the Pitch of Harmonic
686 Complex Tones in the Auditory Nerve. J Neurosci 30:12712–12724.

687 Couchman K, Grothe B, Felmy F (2010) Medial Superior Olivary Neurons Receive
688 Surprisingly Few Excitatory and Inhibitory Inputs with Balanced Strength and Short-
689 Term Dynamics. The Journal of Neuroscience 30:17111 –17121.

690 De Cheveigné A (2010) Pitch perception. The oxford handbook of auditory science:
691 Hearing:71–104.

692 Fontaine B, Benichoux V, Joris PX, Brette R (2013) Predicting spike timing in highly
693 synchronous auditory neurons at different sound levels. J Neurophysiol:jn.00051.2013.

694 Fontaine B, Goodman DFM, Benichoux V, Brette R (2011) Brian Hears: Online Auditory
695 Processing Using Vectorization Over Channels. Front Neuroinform 5 Available at:
696 http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3143729/ [Accessed November 2,
697 2013].

698 Fontaine B, Pena JL, Brette R (2014) Spike-threshold adaptation predicted by membrane
699 potential dynamics in vivo. PLoS Comput Biol 10:e1003560.

700 Gerstner W, Naud R (2009) How Good Are Neuron Models? Science 326:379–380.

701 Glasberg BR, Moore BC. (1990) Derivation of auditory filter shapes from notched-noise
702 data. Hearing Research 47:103–138.

703 GREEN DM, SWETS JA (1966) SIGNAL DETECTION THEORY AND PSYCHOPHYSICS.
704 Oxford, England: John Wiley.

705 Houtsma AJM, Smurzynski J (1990) Pitch identification and discrimination for complex
706 tones with many harmonics. The Journal of the Acoustical Society of America 87:304–
707 310.

708 Jolivet R, Lewis TJ, Gerstner W (2004) Generalized integrate-and-fire models of neuronal

709  activity approximate spike trains of a detailed model to a high degree of accuracy. J
710  Neurophysiol 92:959–976.

711  Joris PX, Carney LH, Smith PH, Yin TC (1994) Enhancement of neural synchronization in
712  the anteroventral cochlear nucleus. I. Responses to tones at the characteristic frequency.
713  J neurophysiol 71:1022.

714  Karino S, Smith PH, Yin TCT, Joris PX (2011) Axonal Branching Patterns as Sources of
715  Delay in the Mammalian Auditory Brainstem: A Re-Examination. J Neurosci 31:3016–
716  3031.

717  Khurana S, Remme MWH, Rinzel J, Golding NL (2011) Dynamic Interaction of Ih and IK-
718  LVA During Trains of Synaptic Potentials in Principal Neurons of the Medial Superior
719  Olive. J Neurosci 31:8936–8947.

720  Langner G (1992) Periodicity coding in the auditory system. Hearing Research 60:115–
721  142.

722  Laudanski J, Coombes S, Palmer AR, Sumner CJ (2010) Mode-Locked Spike Trains in
723  Responses of Ventral Cochlear Nucleus Chopper and Onset Neurons to Periodic Stimuli.
724  Journal of Neurophysiology 103:1226–1237.

725  Licklider JCR (1951) A duplex theory of pitch perception. Experientia 7:128.

726  Loeb GE, White MW, Merzenich MM (1983) Spatial cross-correlation. Biol Cybern
727  47:149–163.

728  Meddis R, O'Mard L (1997) A unitary model of pitch perception. The Journal of the
729  Acoustical Society of America 102:1811.

730  Micheyl C, Oxenham AJ (2007) Across-frequency pitch discrimination interference
731  between complex tones containing resolved harmonics. The Journal of the Acoustical
732  Society of America 121:1621.

733  Montgomery CR, Clarkson MG (1997) Infants' pitch perception: Masking by low- and
734  high-frequency noises. The Journal of the Acoustical Society of America 102:3665–3672.

735  Morgan CT, Garner WR, Galambos R (1951) Pitch and Intensity. The Journal of the
736  Acoustical Society of America 23:658–663.

737  Oxenham AJ (2012) Pitch Perception. J Neurosci 32:13335–13338.

738  Oxenham AJ, Bernstein JG, Penagos H (2004a) Correct tonotopic representation is
739  necessary for complex pitch perception. Proceedings of the National Academy of
740  Sciences of the United States of America 101:1421–1425.

741  Oxenham AJ, Bernstein JGW, Penagos H (2004b) Correct Tonotopic Representation Is
742  Necessary for Complex Pitch Perception. PNAS 101:1421–1425.

743 Oxenham AJ, Micheyl C, Keebler MV, Loper A, Santurette S (2011) Pitch perception
744 beyond the traditional existence region of pitch. Proceedings of the National Academy of
745 Sciences 108:7629–7634.

746 Palmer AR, Shackleton TM (2009) Variation in the Phase of Response to Low-Frequency
747 Pure Tones in the Guinea Pig Auditory Nerve as Functions of Stimulus Level and
748 Frequency. J Assoc Res Otolaryngol 10:233–250.

749 Platkiewicz J, Brette R (2010) A threshold equation for action potential initiation. PLoS
750 Comput Biol 6:e1000850.

751 Platkiewicz J, Brette R (2011) Impact of fast sodium channel inactivation on spike
752 threshold dynamics and synaptic integration. PLoS Comput Biol 7:e1001129.

753 Pressnitzer D, Patterson RD, Krumbholz K (2001) The lower limit of melodic pitch. The
754 Journal of the Acoustical Society of America 109:2074–2084.

755 Ritsma RJ (1962) Existence Region of the Tonal Residue. I. The Journal of the Acoustical
756 Society of America 34:1224–1229.

757 Robles L, Ruggero MA (2001) Mechanics of the mammalian cochlea. Physiol Rev
758 81:1305–1352.

759 Rossant C, Leijon S, Magnusson AK, Brette R (2011) Sensitivity of noisy neurons to
760 coincident inputs. The Journal of Neuroscience 31:17193–17206.

761 Rothman JS, Manis PB (2003a) The Roles Potassium Currents Play in Regulating the
762 Electrical Activity of Ventral Cochlear Nucleus Neurons. J Neurophysiol 89:3097–3113.

763 Rothman JS, Manis PB (2003b) The roles potassium currents play in regulating the
764 electrical activity of ventral cochlear nucleus neurons. Journal of neurophysiology
765 89:3097–3113.

766 Rushton WAH (1951) A theory of the effects of fibre size in medullated nerve. J Physiol
767 115:101–122.

768 Sachs MB, Abbas PJ (1974) Rate versus level functions for auditory-nerve fibers in cats:
769 tone-burst stimuli. The Journal of the Acoustical Society of America 56:1835–1847.

770 Sachs MB, Young ED (1979) Encoding of steady-state vowels in the auditory nerve:
771 Representation in terms of discharge rate. The Journal of the Acoustical Society of
772 America 66:470–479.

773 Semal C, Demany L (1990) The Upper Limit of "Musical" Pitch. Music Perception: An
774 Interdisciplinary Journal 8:165–175.

775 Shamma S, Klein D (2000) The case of the missing pitch templates: How harmonic
776 templates emerge in the early auditory system. The Journal of the Acoustical Society of

777    America 107:2631–2644.

778    Shamma SA (1985) Speech processing in the auditory system II: Lateral inhibition and
779    the central processing of speech evoked activity in the auditory nerve. The Journal of the
780    Acoustical Society of America 78:1622.

781    Slaney M (1993) An efficient implementation of the Patterson-Holdsworth auditory
782    filter bank. Apple Computer, Perception Group, Tech Rep Available at:
783    http://rvl4.ecn.purdue.edu/~malcolm/apple/tr35/PattersonsEar.pdf [Accessed
784    November 2, 2013].

785    Stevens SS (1971) Sensory Power Functions and Neural Events. In: Principles of
786    Receptor Physiology (Loewenstein WR, ed), pp 226–242 Handbook of Sensory
787    Physiology. Springer Berlin Heidelberg. Available at:
788    http://link.springer.com/chapter/10.1007/978-3-642-65063-5_7 [Accessed November
789    2, 2013].

790    Terhardt E (1974) Pitch, consonance, and harmony. J Acoust Soc Am 55:1061–1069.

791    Trussell LO (1997) Cellular mechanisms for preservation of timing in central auditory
792    pathways. Curr Opin Neurobiol 7:487–492.

793    Trussell LO (1999) SYNAPTIC MECHANISMS FOR CODING TIMING IN AUDITORY
794    NEURONS. Annu Rev Physiol 61:477–496.

795    Verschuure J, Van Meeteren AA (1975) The effect of intensity on pitch. Acta Acustica
796    united with Acustica 32:33–44.

797    Yin TC, Chan JC (1990) Interaural time sensitivity in medial superior olive of cat. J
798    Neurophysiol 64:465–488.

799    Zhang X, Heinz MG, Bruce IC, Carney LH (2001) A phenomenological model for the
800    responses of auditory-nerve fibers: I. Nonlinear tuning with compression and
801    suppression. The Journal of the Acoustical Society of America 109:648–670.

802    Zwislocki JJ (1973) On intensity characteristics of sensory receptors: A generalized
803    function. Kybernetik 12:169–183.

804

805

806  **Figures**

807  Figure 1. Regularity structure of the basilar membrane (BM) vibration pattern. (A) Vibration of
808  the basilar membrane produced by a periodic sound S(x,t) (clarinet musical note), at places x
809  tuned to different frequencies (modeled by band-pass filters). (B) The vibration at one place is
810  transformed into spikes produced by an auditory nerve fiber (bottom : post-stimulus time
811  histogram of spikes). In Licklider's model, the fiber projects to a coincidence detector neuron
812  through two axons with conduction delays differing by δ. The neuron fires maximally when the
813  signal's periodicity T equals $\delta$. (C) If the signal's period $T$ is smaller than the neuron's refractory
814  time, then the neuron must detect coincidences between spikes coming from different fibers. (D)
815  If the fibers originate from slightly different places x and y on the cochlea, then the neuron
816  responds to similarities between BM vibrations at different places. (E) Vibration pattern of the
817  BM produced by a non-periodic sound (noise): there is no regularity structure across place and
818  time. (F) Vibration pattern produced by a musical note: there are signal similarities across time
819  (horizontal arrows) and place (oblique arrow).

820  Figure 2. Harmonic resolvability and cross-channel structure. (A) Amplitude and phase
821  spectrum of two gammatone filters. Only a pure tone of frequency f ("Input" waveform) is
822  attenuated in the same way by the two filters (red and blue waveforms: filter outputs). At that
823  frequency, the delay between the outputs of the two filters is $\delta = \Delta\phi/f$. (B) If several harmonic
824  components fall within the bandwidths of the two filters, then the outputs of the two filters
825  differ (no cross-channel similarity). (C) Excitation pattern produced on the cochlea by a
826  harmonic complex. Top: amplitude vs. center frequency of gammatone filters; bottom: spectrum
827  of harmonic complex and of gammatone filters. Harmonic components are "resolved" when they
828  can be separated on the cochlear activation pattern. Higher frequency components are
829  unresolved because cochlear filters are broader. (D) Resolved components produce cross-
830  channel similarity between many pairs of filters (as in A). Unresolved components produce little
831  cross-channel structure (as in B). (E) Thus the vibration pattern produced by resolved
832  components displays both within-channel and cross-channel structure (left), while unresolved
833  components only produce within-channel structure (right).

834  Figure 3. Domain of existence of pitch. (A) Within-channel structure produced by a periodic
835  sound can be decoded if the sound's period is smaller than the maximal neural delay $\delta_{max}$. With
836  $\delta_{max}$ = 4 ms, it occurs for sounds of fundamental frequency greater than 250 Hz. (B) A pure tone
837  or resolved harmonic produces cross-channel structure with arbitrarily small delays between
838  channels, corresponding to the phase difference between the two filters at the sound's
839  frequency: here a 100 Hz tone produces two identical waveforms delayed by δ = 2 ms, while the
840  sound's period is 10 ms. (C) A transposed tone with a high-frequency carrier (>4 kHz)
841  modulated by a low-frequency envelope (<320 Hz) elicits a very weak pitch (Oxenham et al.,
842  2004a) (top: f0 = 120 Hz). Such sounds produce only within-channel structure because they only
843  have high-frequency content (middle). The structural theory of pitch predicts an absence of
844  pitch when the envelope's periodicity is larger than $\delta_{max}$, which is consistent with psychophysics
845  if $\delta_{max}$< 3 ms. (D) A pure tone with the same fundamental frequency (f0 = 120 Hz) produces
846  cross-channel structure with short delays. The structural theory of pitch predicts the existence
847  of pitch in this case, consistently with psychophysical results (Oxenham et al., 2004a). (E)
848  Complex tones with f0 between 400 Hz and 2 kHz and all harmonics above 5 kHz elicit a pitch
849  (Oxenham et al., 2011) (top, spectrum of a complex tone; middle, temporal waveform). Such
850  tones produce only within-channel structure in high-frequency (bottom), and the structural
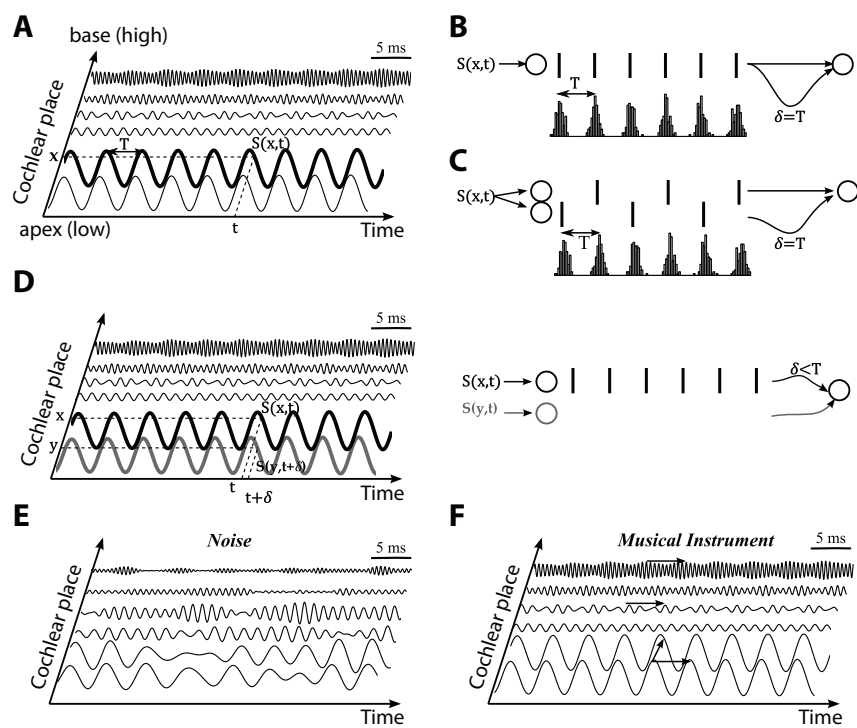
22

851 theory of pitch predicts the existence of pitch if the sound's period is smaller than $\delta_{max}$, which is
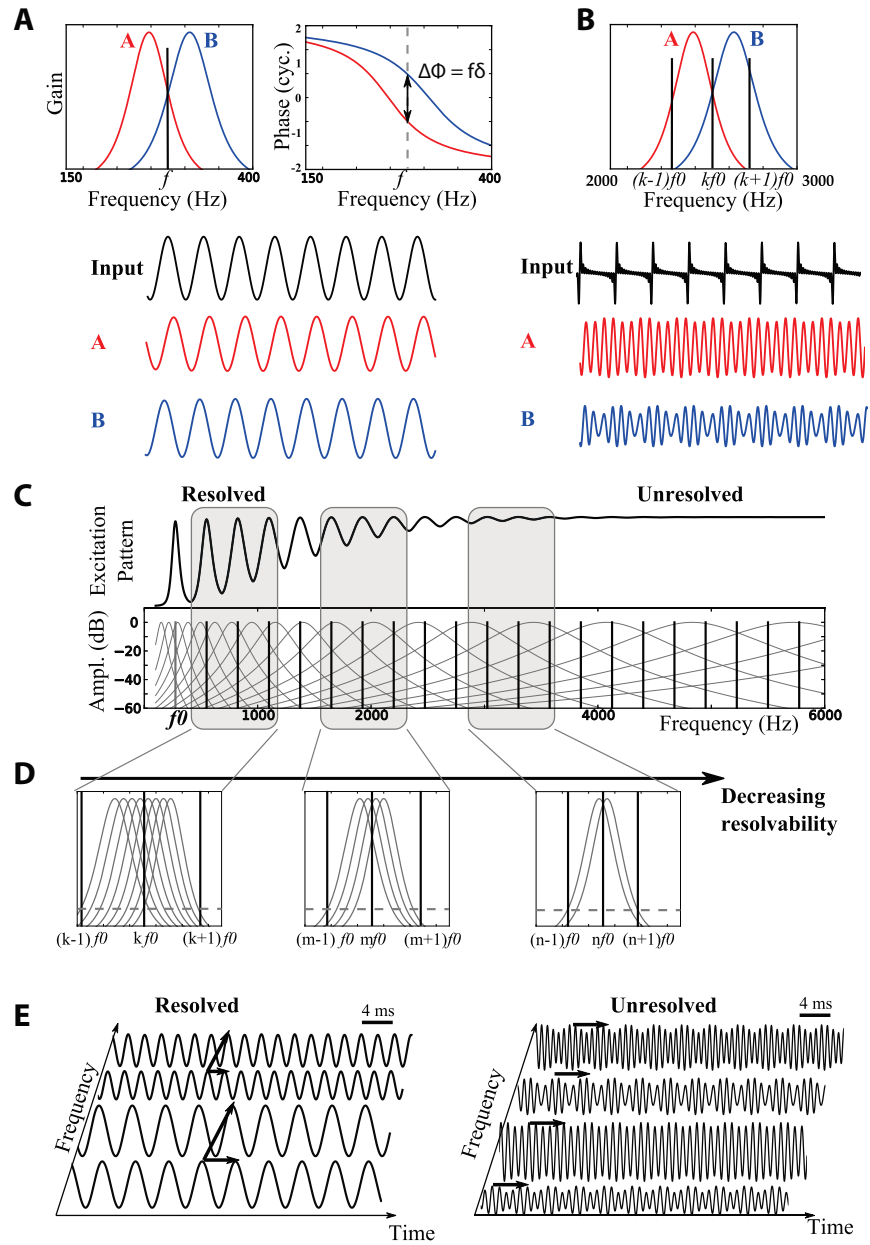852 consistent with psychophysics if $\delta_{max} > 2.5$ ms.

853 Figure 4. Neural network model of pitch estimation using within- and cross-channel structure.
854 (A) Spectrogram of a trumpet sound showing the first two harmonics. Two neurons with CF
855 around the first harmonic and input delay $\delta$ receive the same signal (red and blue rectangles and
856 input signals below). As a result, the two neurons fire synchronously, for all 3 neuron models
857 used: biophysical model of chopper and octopus cells, leaky integrate-and-fire model (voltage
858 traces). (B) Spectrogram of a rolling sea wave sound, which shows no regularity structure. In
859 particular, the two neurons do not receive the same signals (input, shaded area: difference
860 between the two signals) and thus do not fire synchronously. (C) Spectrogram of a harpsichord
861 sound with unresolved harmonics in high frequency. The inset shows the periodicity of the
862 envelope. Two neurons fire synchronously if they receive inputs from the same place delayed by
863 $\delta = 1/f0$. (D) In the same high frequency region, the inharmonic sound of a sea wave does not
864 produce within-channel structure and therefore the two neurons do not fire synchronously. (E)
865 Synaptic connections for a pitch-selective group tuned to f0 = 220 Hz. Harmonics are shown on
866 the left (red comb) superimposed on auditory filters. Resolved harmonics (bottom) produce
867 regularity structure both across and within channels: color saturation represents the amplitude
868 of the filter output while hue represents its phase, for different delays (horizontal axis) and
869 characteristic frequencies (vertical axis). Neurons with the same color fire synchronously and
870 project to a common neuron. Unresolved harmonics (top) produce regularity structure within
871 channels only. Here two identical colors correspond to two identical input signals only when the
872 neurons have identical CF (same row). (F) Same as (E) for a f0 = 261 Hz, producing a different
873 regularity structure, corresponding to a different synchrony pattern in input neurons.
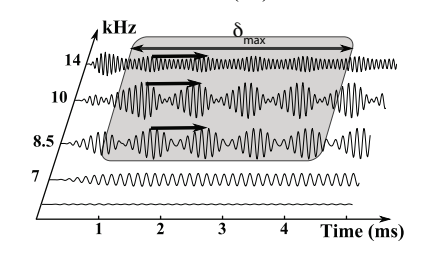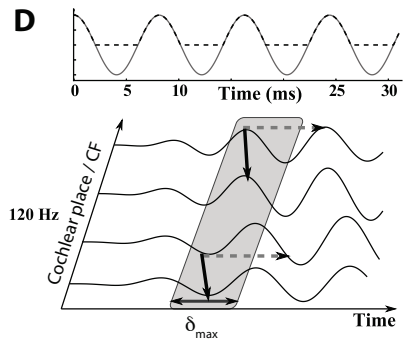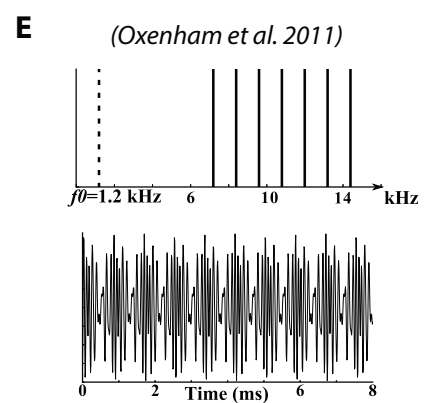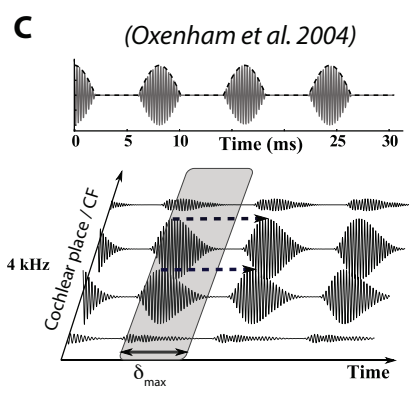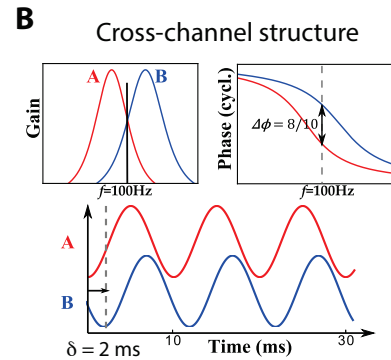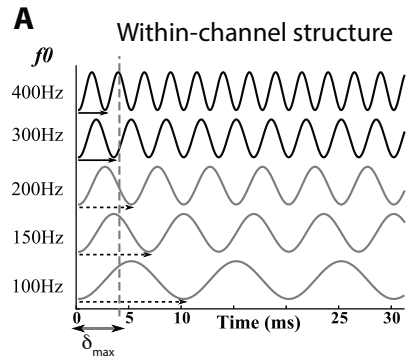874 Synchronous neurons project to another group of neurons, selective for this pitch.

875 Figure 5. Pitch recognition by a neural network model based on the structural theory. (A) Top,
876 Spectrogram of a sequence of sounds, which are either either environmental noises
877 (inharmonic) or musical notes of the chromatic scale (A3-A4) played by different instruments.
878 Bottom, Firing rate of all pitch-specific neural groups responding to these sounds (vertical axis:
879 preferred pitch, A3-A4). (B) Distribution of firing rates of pitch-specific groups for instruments
880 played at the preferred pitch (blue) and for noises (grey), for 3 different sound levels. (C) Top,
881 Pitch recognition scores of the model (horizontal axis: error in semitones) on a set of 762 notes
882 between A2 and A4, including 41 instruments (587 notes) and 5 sung vowels (175 notes).
883 Bottom, Firing rate of all pitch-specific groups as a function of the difference between presented
884 f0 and preferred f0, for all sounds (solid black: average). Peaks appear at octaves (12 semitones)
885 and perfect fifths (7 semitones). (D) Impact of the number of frequency channels (top) and
886 maximal delay $\delta_{max}$ (bottom) on recognition performance.

887 Figure 6. Pitch discriminability. (A) Two neurons tuned to the same frequency (within-channel)
888 but with delay mismatch $\delta = 1/f$ produce phase-locked spikes (red and blue crosses) in response
889 to a tone (sine waves). When the tone frequency is f (left), the two input signals match and the
890 difference of phases of spikes $\Delta\Phi(f)$ between the two neurons is distributed around 0 (shaded
891 curve). When the tone frequency is f+df (right), the two signals are slightly mismatched and the
892 distribution of $\Delta\Phi(f)$ is not centered on 0. (B) Two neurons tuned to different frequencies
893 (cross-channel) respond at different mean phases to tones (red and blue curves). (C) The
894 discriminability index d' is defined as the distance $\mu$ between the centers of be two phase
895 difference distributions ($\Delta\Phi(f)$ and $\Delta\Phi(f+df)$) relative to their standard deviation $\sigma$. (D) The

896    standard deviation of the phase distribution is related to the precision of phase locking,
897    measured by the vector strength (dots: vector strength vs. characteristic frequency for guinea
898    pig auditory fibers; solid curve: fit). (E) Mean phase of spikes produced by auditory nerve fibers
899    of guinea pigs for different tone frequencies (data from Palmer and Shackleton (2009) (Palmer
900    and Shackleton, 2009)), as a function of CF (crosses), with fits (solid lines). (F) Weber fraction
901    ($\Delta f/f$, where $\Delta f$ is the just noticeable difference in frequency) as a function of tone frequency
902    for cross-channel structure (colored curves) and within-channel structure (black curve). Color
903    represent different frequency spacings between the two channels (1-6 semitones). Dotted lines
904    represent the limitations implied by a maximal delay $\delta_{max}$ = 5 ms.

**A**

Gain

A    B

150    $f$    400
Frequency (Hz)

Phase (cyc.)

$\Delta\Phi = f\delta$

150    $f$    400
Frequency (Hz)

**B**

A    B

2000    $(k\text{-}1)f0$    $kf0$    $(k\text{+}1)f0$    3000
Frequency (Hz)

**Input**

**A**

**B**

**Input**

**A**

**B**

**C**

Excitation Pattern

**Resolved**                                        **Unresolved**

Ampl. (dB)

0
−20
−40
−60

$f0$    1000    2000    4000    Frequency (Hz)    6000

**D**

**Decreasing resolvability**

$(k\text{-}1)\,f0$    $k\,f0$    $(k\text{+}1)f0$        $(m\text{-}1)\,f0$    $mf0$    $(m\text{+}1)f0$        $(n\text{-}1)f0$    $nf0$    $(n\text{+}1)f0$

**E**

**Resolved**                    4 ms

Frequency

Time

**Unresolved**                    4 ms

Frequency

Time

**A** Within-channel structure

*f0*

400Hz
300Hz
200Hz
150Hz
100Hz

$\delta_{max}$

Time (ms)
0   5   10        25   30

**B** Cross-channel structure

Gain

A          B

*f*=100Hz

Phase (cycl.)

$\Delta\phi = 8/10$

*f*=100Hz

A

B

$\delta$ = 2 ms        10        Time (ms)        30

**C** *(Oxenham et al. 2004)*

Time (ms)
0   5   10        25   30

Cochlear place / CF

4 kHz

$\delta_{max}$        Time

**D**

Time (ms)
0   5   10        25   30

Cochlear place / CF

120 Hz

$\delta_{max}$        Time

**E** *(Oxenham et al. 2011)*

*f0*=1.2 kHz   6        10        14   kHz

Time (ms)
0        2        6        8

kHz        $\delta_{max}$

14
10
8.5
7

1   2   3   4   Time (ms)

A **Trumpet C#4 (277 Hz)**
B **Sea wave**

$\delta = \Delta\phi/f_0$

Input
Chopper
Octopus
Leaky IF

C **Harpsichord A3 (220 Hz)**
D **Sea wave**

$\delta = 1/f0$

Input
Leaky IF

E $f0 = 220Hz\ (A3)$
F $f0 = 261Hz\ (C4)$

Input neurons
Coincidence detectors

Characteristic Frequency (Hz)

Unresolved harmonics
Resolved harmonics

$10^3$
$10^2$

**Pitch**
$f0$ (Hz)

A3 (220 Hz)
A#3 (233 Hz)
B3 (247 Hz)
C4 (261 Hz)
C#4 (277 Hz)
G#4 (415 Hz)

**A**

### Spectrogram



Frequency

1s

*Airplane*    *Sea waves*    *Street*    *Clarinet*    *Accordion*    *Viola*

### Response of pitch tuned groups

1s

Pitch

Time

Rate (Hz)
270
180
90
0

**B**



45 dB — Instruments — Noises

Probability

0        Firing Rate (Hz)        35

55 dB

Probability

0        Firing Rate (Hz)        80

65 dB

Probability

0        Firing Rate (Hz)        160

**C**



100 %

Recognition score

77 %

5 %

−15  −10  −5    0    5    10   15

Firing Rate (Hz)

−15  −10  −5    0    5    10   15

f0$_{estimate}$ - f0 (semitones)

**D**



Recognition score
80
70
60
50
40
30
20
10

−15 −10 −5 0 5 10 15

# of channels
400
300
200
100

f0$_{estimate}$ - f0 (semitones)

Recognition score
70
60
50
40
30
20
10
0

−15 −10 −5 0 5 10 15

Max delay (ms)
14
12
10
8
6
4
2

f0$_{estimate}$ - f0 (semitones)

**A**  Within-channel

Tone at $f$

$\delta = 1/f$

Tone at $f+df$

$\delta$

$\Delta\phi(f)$

$0$

$\Delta\phi(f+df)$

$0$

**B**  Cross-channel

Phase

$\Delta\phi(f)$

$\Delta\phi(f+df)$

Stimulus freq. (Hz)

$\Delta\phi(f)$

$\Delta\phi(f+df)$

**C**

$$\mathbf{d'} = \mu/\sigma$$

$\mu$

$\sigma$

$\langle\Delta\phi(f)\rangle$  $\langle\Delta\phi(f+df)\rangle$

**D**

Vector strength

0.9

0.6

0.3

0

$10^2$  $10^3$  $10^4$

Frequency (Hz)

**E**

Phase (cycle)

0
-1
-2
-3
-4
-5
-6
-7

0  500  1000  1500  2000  2500  3000

CF (Hz)

— 250 Hz
— 353 Hz
— 500 Hz
— 707 Hz
— 1000 Hz
— 1414 Hz
— 2000 Hz

**F**

Weber Fraction (%)

10

0.1

$10^2$  Frequency (Hz)  $10^3$

— ΔCF: 1 st
— ΔCF: 2 st
— ΔCF: 3 st
— ΔCF: 4 st
— ΔCF: 5 st
— ΔCF: 6 st
— Within channel