# The neural response to the temporal fine structure of continuous musical pieces is not affected by selective attention

**Octave Etard, Rémy Ben Messaoud, Gabriel Gaugain, and Tobias Reichenbach[†]**

Department of Bioengineering and Centre for Neurotechnology, Imperial College London, South Kensington Campus, SW7 2AZ, London, U.K.

[†]To whom correspondence should be addressed (email: reichenbach@imperial.ac.uk)

## Abstract

Speech and music are spectro-temporally complex acoustic signals that a highly relevant for humans. Both contain a temporal fine structure that is encoded in the neural responses of subcortical and cortical processing centres. The subcortical response to the temporal fine structure of speech has recently been shown to be modulated by selective attention to one of two competing voices. Music similarly often consists of several simultaneous melodic lines, and a listener can selectively attend to a particular one at a time. However, the neural mechanisms that enable such selective attention remain largely enigmatic, not least since most investigations to date have focussed on short and simplified musical stimuli. Here we study the neural encoding of classical musical pieces in human volunteers, using scalp electroencephalography (EEG) recordings. We presented volunteers with continuous musical pieces composed of one or two instruments. In the latter case, the participants were asked to selectively attend to one of the two competing instruments and to perform a vibrato identification task. We used linear encoding and decoding models to relate the recorded EEG activity to the stimulus waveform. We show that we can measure neural responses to the temporal fine structure of melodic lines played by one single instrument, at the population level as well as for most individual subjects. The neural response peaks at a latency of 7.6 ms and is not measurable past 15 ms. When analysing the neural responses elicited by competing instruments, we find no evidence of attentional modulation. Our results show that, much like speech, the temporal fine structure of music is tracked by neural activity. In contrast to

25   speech, however, this response appears unaffected by selective attention in the context of our

26   experiment.

27   **Acknowledgement**

33   **Competing financial interests**

34   The authors declare no competing financial interest

# Introduction

36   Music is a fascinatingly complex acoustic stimulus. Listeners can follow multiple melodic lines played

37   by different instruments by separating them on the basis of characteristics such as pitch and timbre

38   (Cross et al., 2008). However, the neural mechanisms that group the sounds in music into distinct

39   melodic lines, forming distinct auditory streams, and allow attention to be directed to one of the lines

40   remain largely unknown (Albert S Bregman, 1994). This is partly due to the difficulty in assessing the

41   neural processing of real-world acoustic signals that have a much richer structure than the simple pure

42   tones and short simplified music patterns that have traditionally dominated research in auditory

43   neuroscience.

44   A better understanding of the neural mechanisms of music processing may emerge from combining

45   statistical models with neuroimaging. Recent studies have indeed shown how these methods can relate

46   key features of a complex sound such as speech to electrophysiological recordings and inform on the

47   neural mechanisms of speech processing (Di Liberto et al., 2015; Ding & Simon, 2012a, 2014;

48   Wöstmann et al., 2017). For example, cortical activity has been found to track slow (< 8Hz) amplitude

49   fluctuations in speech (Ding & Simon, 2012b; Edmund C. Lalor & Foxe, 2010; Nourski et al., 2009;

50   Pasley et al., 2012), while subcortical as well as, presumably to a lesser degree, cortical responses

51    emerge to the higher frequency (> 80 Hz) stimulus structure (Bidelman, 2018; Coffey et al., 2016; Etard

52    et al., 2019; Forte et al., 2017; Maddox & Lee, 2018). The temporal fine structure of speech originates

53    from the periodic opening and closing of the vocal folds at the so-called fundamental frequency. The

54    spectrum of these voiced speech parts is therefore dominated by the fundamental frequency as well as

55    its many higher harmonics, leading to a pitch perception in the listeners.

56    Understanding how the brain can focus on a single instrument amongst others relates to a major

57    challenge in auditory neuroscience, the cocktail party problem. This problem acquired its name from

58    the observation that humans do remarkably well at understanding a target speaker in a noisy

59    environment such as in a busy restaurant or in a loud bar (Cherry, 1953; Haykin & Chen, 2005). A

60    recent study showed that neural responses to the pitch of continuous speech are stronger when the

61    stimulus is attended rather than ignored (Forte et al., 2017). This result suggests that the pitch of a

62    speaker could be used by the brain to perceptually segregate the speech signal from background noise,

63    a finding that agrees with previous psychophysical studies that have found it easier to differentiate two

64    concurrent speech signals if their fundamental frequencies differ (de Cheveigné et al., 1997; Madsen et

65    al., 2017).

66    Musical tones are similarly characterized by a fundamental frequency and higher harmonics, resulting

67    in a characteristic temporal structure that causes a pitch perception. The proximity of fundamental

68    frequencies of subsequent tones has been found to aid the formation of an auditory stream (A. S.

69    Bregman et al., 1990; Oxenham, 2008). Consequently, just as the neural tracking of temporal fine

70    structure could help listeners attend to a voice in background noise, such a neural mechanism may aid

71    with attending to a particular melodic line (Micheyl & Oxenham, 2010).

72    Here we investigated this hypothesis by using linear models to assess neural responses to the temporal

73    fine structure of continuous melodic lines. We first presented volunteers with continuous classical Bach

74    pieces while recording their brain activity through a bipolar EEG montage. We then related the neural

75    activity to the stimulus waveforms using encoding and decoding methods. To assess a putative effect

76    of selective attention on these neural activities, we also presented the volunteers with two competing

77    instruments, a guitar and a piano, that were playing two different melodic lines simultaneously. Subjects

78    were asked to selectively attend to one of the two lines, and we contrasted the neural responses to each

79    instrument when it was attended to when it was ignored.

## Methods

81    Due to multiple nonlinearities in the auditory periphery, both the temporal fine structure and the

82    envelope of the stimuli are represented in neural responses. This encoding has traditionally been

83    investigated in humans by studying time-locked responses to transient or periodic features of repeated

84    short sound tokens such as clicks, pure or complex tones, syllables and words (Skoe & Kraus, 2010).

85    These paradigms typically present a particular stimulus as well as its opposite waveform many times.

86    The neural responses to each polarity are then summed to emphasise responses to the envelope or

87    subtracted to emphasise response to the temporal fine structure (Aiken & Picton, 2008; Krizman &

88    Kraus, 2019). Here we used continuous, long stimuli to derive auditory neural responses to their

89    temporal fine structure using linear convolutive models.

90    **Experimental design and statistical analysis.** Seven of Bach's Two-Part Inventions were used in this

91    study. Each Two-Part Invention is a short keyboard composition that consists of two melodic lines: one

92    played by the left hand, and one by the right. We synthesized the stimuli in GarageBand (Apple, U.S.A)

93    from Musical Instrument Digital Interface (MIDI) files, with the left hand being played by a piano and

94    the right by a guitar. To assess the attention of subjects to a particular melodic line, vibratos were

95    inserted in both lines.

96    Volunteers were presented with two type of stimuli. The first type, "Single Instrument" (SI), consisted

97    of one single instrument, piano or guitar, that played one melodic line. The second type "Competing

98    Instruments" (CI), contained both melodic lines of a Two-Part Invention, one played by the piano and

99    the other by the guitar.

100    The different stimuli were presented in blocks (figure 1). Each block contained one SI stimulus and one

101    subsequent CI stimulus, both of which were obtained from the same Two-Part Invention. During the CI

102    stimulus, the volunteers were asked to selectively listen to the instrument that they heard before in the

103    SI stimulus. They were also asked to identify the vibratos embedded into that melodic line.
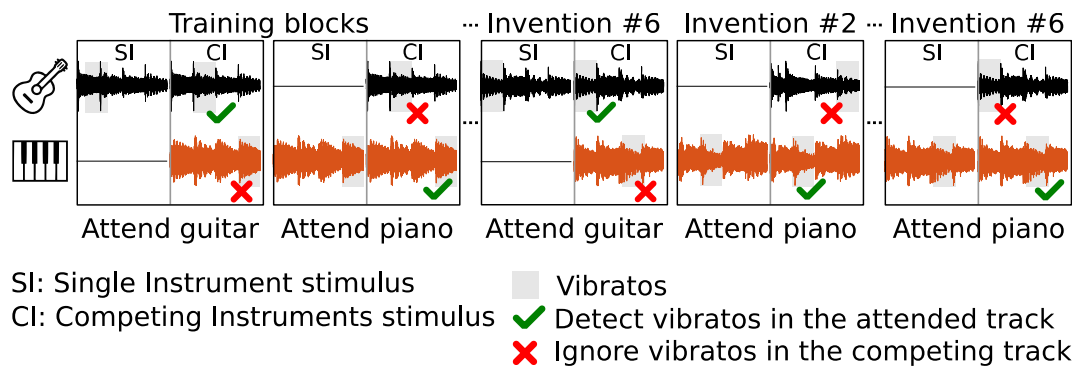
**Figure 1: Schematic representation of the experiment**. Volunteers were presented with continuous classical music pieces, Bach's Two-Part Inventions, that consisted of either a single melodic line (SI) or of two melodic lines (CI). Each melodic line was played either by a guitar or by a piano. In the CI stimuli, each melodic line was played by a different instrument. Vibratos were inserted into the acoustic waveforms of each melody (grey shading). In the CI condition, the subjects had to attend to one of the two instruments and identify the corresponding vibratos (green tick marks) while ignoring the other instrument and its vibratos (red crosses). The stimuli were presented in blocks composed of a SI stimulus followed by a CI stimulus during which the subject was asked to attend to the instrument that they heard before in the SI stimulus. The attended instrument was alternated between blocks, and each block was played twice such that the attended instrument differed in the two presentations. The volunteers' neural responses were recorded throughout the experiment through bipolar two-channel EEG recordings.

104    Blocks with the SI stimulus played by the piano alternated with those played by the guitar. Each of the

105    seven Two-Part inventions was presented twice: once with the SI stimulus played by the guitar, and

106    once with the SI stimulus played by the piano. Each participant therefore heard seven SI stimuli played

107    by the guitar, seven SI stimuli played by the piano, and fourteen CI stimuli.

108    All participants were initially presented with the same two training blocks, one with a SI stimulus played

109    by the guitar and one played by the piano, that corresponded to the same invention. These stimuli

110    presentations served to familiarise the subjects with the task of attending to one melodic line in the CI

111    stimulus and to identify the embedded vibratos. These training blocks were excluded from further

112    analysis, leaving six inventions in each condition.

113  The presentation order of the remaining blocks was pseudo-randomised across participants. In the

114  second presentation of a given CI stimulus, a subject was asked to attend to the instrument they ignored

115  in the first presentation. Two consecutive blocks did not correspond to the same invention. Each

116  participant therefore heard each CI stimulus twice, but attended a different instrument in each

117  presentation. Whether a participant was initially asked to attend to the guitar or the piano was randomly

118  decided. The participant's neural responses were measured through scalp electroencephalography

119  (EEG) with a two-channel bipolar montage (head vertex minus mastoids).

120  We used encoding and decoding approaches (linear forward and backward models) to relate the acoustic

121  stimuli to the recorded neural data. We specifically investigated the neural representation of the

122  temporal fine structure by using the stimulus waveform as a feature. We first established that we could

123  indeed record a significant neural response to this feature, by comparing the neural responses to a null

124  distribution at the level of individual subjects as well as on the population level. We then studied the

125  time course of the response in the region between 0 to 45 ms, using both forward and backward models.

126  Finally, we investigated a putative attentional modulation of this neural response through contrasting

127  the encoding of each instrument in the neural data when attended versus ignored. We used conservative

128  filters to reduce distortions to the neural responses and their latencies, but verified that our results, and

129  in particular the ones related to attention, did not change with stronger filtering (data not shown).

130  **Code and data availability.** The analysis presented in this manuscript was implemented using

131  MATLAB (R2019b, The MathWorks Inc.) with the EEGLAB toolbox (Delorme & Makeig, 2004). The

132  linear     forward     and     backward     models     were     trained     using     the     LMpackage

133  (github.com/octaveEtard/LMpackage). The raw data as well as analysis code and processed data

134  required     to     reproduce     the     results     presented     here     were     made     available

135  (https://github.com/octaveEtard/EEGmusic2020; https://zenodo.org/record/4470135).

136  **Participants.** 17 volunteers (aged $23.8 \pm 2.9$ year, 9 females) participated in this experiment. The

137  number of participants was chosen based on previous studies investigating similar neural responses to

138  continuous speech (Etard et al., 2019; Forte et al., 2017). All participants were right-handed, had no

139    history of auditory or neurological impairments, and provided written informed consent. The

140    experimental procedures were approved by the Imperial College Research Ethics Committee.

141    **Music stimuli.** To generate neural responses to each instrument that were of similar magnitude, the

142    notes of the guitar were lowered by one octave so that their fundamental frequencies fell below 500 Hz.

143    They remained nonetheless somewhat higher than those of the piano notes (figure 2, A, B). The music

144    stimuli were synthesized from MIDI files to generate wav files. These were then processed using

145    MATLAB to apply vibratos to ten segments in each melodic line. Each vibrato was constructed by

146    introducing a sinusoidal warp at a modulation frequency of $f_m$ = 8 Hz on the waveform of a single note.

147    The onset and offset times of the notes were obtained from the MIDI files using the Miditoolbox for

148    Matlab (Eerola & Toiviainen, 2004). The notes were selected such that the onsets of any two vibratos

149    in a given piece, whether both played by the same or different instruments, were separated by at least

150    one second.

151    The waveforms of the CI stimuli, $w_{mixed}$, were constructing by normalising and mixing the waveform

152    $w_g$ of the guitar and the waveform $w_p$ of the piano according to their root-mean-square values (RMS):

153    $w_{mixed} = \frac{w_g}{RMS(w_g)} + 1.25 \cdot \frac{w_p}{RMS(w_p)}$. The mixing parameter of 1.25 for the piano was chosen following

154    a small pilot study to balance the difficulty in attending either the guitar or the piano.

155    The duration of the seven Two-Part Inventions was, taken together, 11.2 minutes. In the SI conditions,

156    only the first half of the corresponding invention was played.

157    **Behavioural task.** In the CI condition, the subjects were instructed to attentively listen to one

158    instrument while ignoring the other. They were also asked to classify the vibratos they heard by pressing

159    a key to indicate the ones that belonged to the attended instrument. A key press within two seconds after

160    the onset of a vibrato in the attended or ignored instrument was classified respectively as "true positive"

161    (TP) or "false positive" (FP). Key presses outside of these ranges were classified as "unprompted" and

162    were not analysed further. Due to a technical error, behavioural data was not recorded for one subject,

163    and only the results for the 16 remaining subjects were analysed. The sensitivity index d-prime was

164    computed for each subject when attending to the guitar and the piano, and it was compared between the

165    two conditions at the population level using a two-tailed paired Wilcoxon signed rank test. Moreover,

166    for each condition the TP rate (TPR) was compared to the FP rate (FPR), and the TPR and FPR were

167    compared between conditions at the population level using two-tailed paired Wilcoxon signed rank tests

168    with FDR correction for multiple comparisons (four tests).

169    **Neural data acquisition and stimulus presentation.** Scalp EEG was recorded through five passive

170    Ag/AgCl electrodes (Multitrode, BrainProducts, Germany). Two electrodes were positioned on the

171    cranial vertex (Cz), and two electrodes were placed on the left and right mastoid processes. A ground

172    electrode was placed on the forehead. The impedance between each electrode and the skin was reduced

173    below 5 kOhm using abrasive electrolyte gel (Abralyt HiCl, Easycap, Germany). One vertex electrode

174    was paired with the left mastoid electrode, and they were connected to, respectively, the non-inverting

175    and inverting ports of a bipolar amplifier (EP-PreAmp, BrainProducts, Germany). The remaining vertex

176    and mastoid electrodes were similarly connected to a second identical amplifier. The output of each

177    bipolar pre-amplifier was fed into an amplifier (actiCHamp, BrainProducts, Germany) and digitized

178    with a sampling frequency of 5 kHz, thus yielding two electrophysiological data channels. The audio

179    stimuli were simultaneously recorded at 5 kHz by the amplifier through an acoustic adapter (Acoustical

180    Stimulator Adapter and StimTrak, BrainProducts, Germany). This channel and independent analogue

181    triggers delivered through an LPT port were used to temporally align the EEG data and stimuli through

182    cross-correlation. The stimuli were delivered diotically at a comfortable loudness level through insert

183    tube earphones (ER-3C, Etymotic, USA) to minimise stimulation artifacts. These earphones introduced

184    a 1 ms delay that was compensated for by shifting the neural data forward in time by 1ms.

185    **EEG data filtering.** To analyse the neural responses to the temporal fine structure of the stimuli, the

186    EEG data was high-pass filtered above 130 Hz (windowed-sinc filters, Kaiser window, one pass forward

187    and compensated for delay; cut-off: 115 Hz, transition bandwidth: 30 Hz, order: 536). These filters

188    rejected lower-frequency neural activity but reduced the temporal precision of the data, as evidenced

189    by the auto-correlation function of the filtered EEG data (figure 2C). Notably, they were non-causal

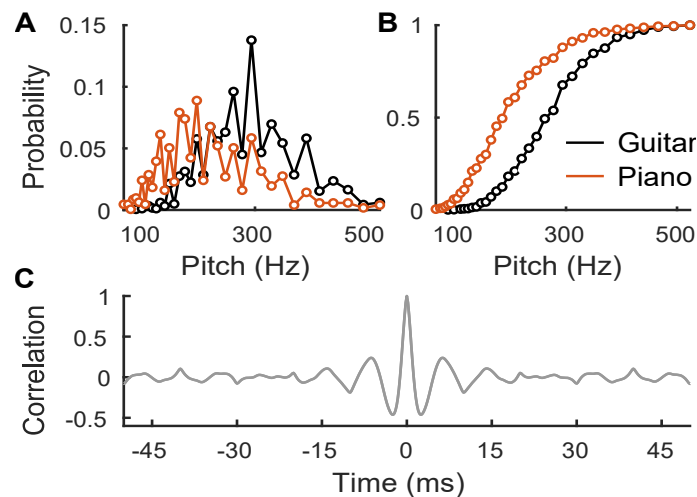190    filters that spread responses in both temporal directions.

**Figure 2: Properties of the acoustic stimuli and of the filtered EEG data.** (**A**), The probability mass function of the fundamental frequency of the notes peaked at around 196 Hz for the piano (red), and at about 294 Hz for the guitar (black). (**B**), The cumulative distribution of the fundamental frequency of the notes showed likewise that most fundamental frequencies lied between 100 Hz and 400 Hz, with the distribution of the guitar notes being shifted to somewhat higher frequencies. (**C**), To eliminate frequencies below the range of the fundamental frequencies, the EEG data was high-pass filtered above 130 Hz. The filtered EEG data consequently displayed some periodicity and correlation in time as evident from its auto-correlation function.

191 **Stimulus representation.** Since the vibratos might lead to neural responses deviating from the ones

192 elicited by the rest of the tracks, the parts of the stimulus waveforms that corresponded to them were

193 replaced with zeros to create the stimulus representations (features) used in the encoding and decoding

194 models. These waveforms were then low pass filtered and resampled from 44.1 kHz to 5 kHz, the

195 sampling frequency of the EEG data, using a linear phase FIR anti-aliasing filter (windowed-sinc filter,

196 Kaiser window, one pass forward and compensated for delay; cut-off: 2,250 Hz, transition bandwidth:

197 500 Hz, order: 14,126).

198 **Encoding models.** We used regularised linear forward models to derive the neural response to the

199 stimulus waveform. In these convolutive encoding models, the measured EEG response $e$ is modelled

200 as $e(t) = (r * s)(t) + n(t)$, where $s$ is the stimulus waveform, $r$ the neural response or Temporal

201 Response Function (TRF), $n$ is noise, and $*$ is the convolution symbol. In practice, assuming a non-zero

202 response in a time interval $(\tau_{min}; \tau_{max})$ only, and with discrete data, the EEG activity $e_i(t_n)$ at channel

203    $i \in \{1; 2\}$ and at time $t_n$ can be estimated as $\hat{e}_i(t_n) = \sum_{k=1}^{N} r(\tau_k) \cdot s(t_n - \tau_k)$, with $\tau_1 = \tau_{min}$ and

204    $\tau_N = \tau_{max}$ . Given the bipolar montage we used, as well as the diotic stimulus presentation, we did not

205    expect any difference between the two EEG channels and assumed the same neural response for both.

206    The model was estimated for time lags spanning $\tau_{min} = -100$ ms to $\tau_{max} = 45$ ms. A population-

207    averaged TRF $r$ was fitted using ridge regression coupled with a leave-one-subject-out and leave-one-

208    data-part out cross-validation (Crosse et al., 2016; Hastie et al., 2009; E. C. Lalor et al., 2009). The

209    model was fitted using the data corresponding to all the stimulus parts bar one and all the subjects bar

210    one, and evaluated on the left-out data part for the left-out subject. The stimulus part and the subject

211    used for testing were hence not seen by the model during training. This constituted one cross-validation

212    fold. The left-out subject and the left-out data parts were then iterated until all combinations were

213    exhausted, for a total of $17 \times 6 = 102$ folds. The validation performance of the model was quantified

214    by dividing the predicted neural response $\hat{e}_i$ and the measured EEG activity $e_i$ from the testing data in

215    each fold into 10-s long segments, and by computing Pearson's correlation coefficients between each

216    segment. The correlation coefficients thus obtained were then averaged over all cross-validation folds

217    as well as over all EEG channels.

218    The performance was assessed for models corresponding to 25 normalised regularisation coefficients

219    $\lambda_n$ that were distributed uniformly on a logarithmic scale between $10^{-6}$ and $10^6$. The regularisation

220    coefficient was thereby $\lambda = \lambda_n \cdot m$, with $m$ the mean eigenvalue of the predictor's auto-correlation

221    matrix (Biesmans et al., 2017). The model yielding the highest reconstruction performance was chosen

222    as representing the neural response. To assess the significance of the obtained TRFs, the negative, non-

223    causal part of the response, -100 ms to 0 ms, was used to construct a null distribution. For each

224    instrument, a Gaussian distribution was fitted to the pooled data points from the negative part of the

225    response. From the distribution we determined the $p$-values of all the points in the positive part of the

226    response (0 ms to 45 ms), and applied an FDR correction for multiple comparison over time points and

227    instruments.

228      To ascertain the relative contributions of the onset and of the sustained parts of the notes to the neural

229      response, we created a new representation of the stimuli in which the note onsets were suppressed. This

230      was achieved by multiplying the original stimulus waveforms by a 60-ms window $w$ centred on each

231      note onsets, with $w(t) = 1 - h(t)$ and $h$ representing a 60-ms Hann window. Forward models were

232      then derived for the original stimuli and their onset-suppressed versions for the two SI conditions taken

233      together, by pooling the data from both instruments. These two models were fitted and their significance

234      was ascertained as described above, that is, by comparing the causal part to the null models, with FDR

235      correction for multiple comparison over time points and over the two models. In the cross-validation

236      procedure, two data parts, one from each SI condition and corresponding to the same invention, were

237      left out at each stage.

238      **Decoding models.** We also used backward models to reconstruct the stimulus waveform $s$ as a linear

239      combination of the neural activity $e_i$ on each channel $i$ at different time lags: $\hat{s}(t_n) =$

240      $\sum_{i=1}^{2} \sum_{k=1}^{N} \beta_i(\tau_k) \cdot e_i(t_n + \tau_k)$, with $\tau_{min} \leq \tau_k \leq \tau_{max}$. The coefficients $\beta$ were trained for each

241      subject independently, using ridge regression with a leave-one-part-out cross-validation and a

242      normalised regularisation coefficient $\lambda_n = 10^{-0.5}$ (Biesmans et al., 2017). As with the forward models,

243      the performances of the backward models were measured through computing the correlation

244      coefficients between the reconstructed stimulus and the actual one on 10-s long segments of the testing

245      data. The set of correlation coefficients pooled from all cross-validation folds for a given participant

246      was used when performing statistical testing at the level of individual subjects, and the corresponding

247      average correlation coefficient was used for each subject when testing at the population level. The

248      performances of the models were thereafter used to quantify the neural encoding of each stimulus for a

249      given reconstruction time window $\tau_{min} - \tau_{max}$.

250      **Significance of the stimulus reconstruction.** The neural encoding of the SI stimuli for each instrument

251      was measured through the backward models using reconstruction time windows of equal duration but

252      centred on different delays. To establish the significance of the stimulus reconstruction procedure at the

253      level of individual subjects, a window of delays between $\tau_{min} = -15$ ms and $\tau_{max} = 0$ ms was used

254      to provide a null distribution for each subject. The neural encoding in the window of interest, from

255 $\tau_{min} = 0$ ms to $\tau_{max} = 15$ ms, was compared to the null distribution for each subject using one-tailed

256 paired Wilcoxon signed rank tests with FDR correction for multiple comparisons over subjects and

257 instruments. Significance was also derived at the population level using the mean correlation

258 coefficients for each subject from the null window of negative delays to create a null population-level

259 distribution. To test the time windows in which a significant response could be detected, the mean

260 reconstruction accuracies from three windows of interest (0 to 15 ms; 15 to 30 ms; 30 to 45 ms) were

261 compared to this null distribution using one-tailed paired Wilcoxon signed rank tests with FDR

262 correction for multiple comparisons over windows and instruments.

263 Since the guitar and piano waveforms formed pairs derived from the same inventions, and although

264 their frequency contents were different, one may wonder whether one instrument could be predicted

265 from the other, and in turn whether the neural responses to one instrument could be predicted or used

266 to decode the other one. To address this question, we trained linear backward models that sought to

267 reconstruct the waveform of one instrument from the neural data that was recorded when the other

268 instrument from the same invention was played in the SI conditions (0 to 15 ms reconstruction window).

269 The model performance was then compared to the null distribution previously described (obtained from

270 a -15 to 0 ms reconstruction window) at the population level, using one-tailed paired Wilcoxon signed

271 rank tests.

272 **Competing conditions, attended and ignored instruments.** In the CI conditions, we trained backward

273 models to reconstruct the waveform of either the attended or the ignored instrument independently,

274 using a window of temporal delays from $\tau_{min} = 0$ ms to $\tau_{max} = 15$ ms as detailed above. We then

275 compared the neural encoding of each instrument, when attended and when ignored, at the population

276 level, using two-tailed paired Wilcoxon signed rank tests with FDR correction for multiple comparisons

277 over instruments.

278 We also used forward models reconstructing the neural activity as the sum of two neural responses, one

279 to the attended instrument and one to the ignored one. In this instance, the EEG response $e$ is modelled

280 as $e(t) = (r_A * s_A)(t) + (r_I * s_I)(t) + n(t)$, where $s_A$ and $s_I$ are the attended and ignored stimulus

281 waveforms, and $r_A$ and $r_I$ the corresponding TRFs. In a similar manner to the procedures previously

282  described, population-averaged TRFs were fitted using ridge regression coupled with a leave-one-

283  subject-out and leave-one-data-part out cross-validation for time lags spanning $\tau_{min} = -100$ ms to

284  $\tau_{max} = 45$ ms on the pooled data from the two CI conditions. To assess the presence of a putative

285  attentional modulation in the obtained TRFs, the distribution of amplitude across subjects was compared

286  between the attended and ignored TRFs for each time point in the 0 ms to15 ms region of interest (two-

287  tailed paired Wilcoxon signed rank tests with FDR correction for multiple comparisons over time

288  points).

## Results

290  We asked volunteers to attend to continuous musical pieces consisting of either one single instrument

291  (SI) or of two competing instruments (CI) while we recorded their neural activity using EEG (figure 1).

292  We first sought to analyse the neural response to the temporal fine structure of a single melodic line.

293  To this end, we computed a linear forward model to derive neural responses to the stimulus waveform

294  at the population level in the SI conditions (figure 3A). The temporal response functions that we

295  obtained for the two instruments were qualitatively similar to each other. They displayed a major

296  significant response at a latency of 7.6 ms, as well as a minor positive peak at 2.2 ms, with sidelobes

297  reminiscent of the EEG auto-correlation function (figure 2C).

298  The neural response to temporal fine structure may be related to the well-established frequency-

299  following response (FFR). Because the latter is known to first exhibit a response to a stimulus onset and

300  to then follow the sustained features, we explored the relative contributions of the note onsets and their

301  sustained oscillations to the neural response. We therefore trained a forward model with stimulus

302  waveforms in which the note onsets were suppressed (figure 3B). The obtained temporal response

303  functions had similar significant regions, and resembled the temporal response functions to the original

304  stimulus waveforms. Moreover, the causal parts of the two temporal response functions, those with

305  positive delays, were highly correlated ($r = 0.96$).

306  As an alternative method to the forward models, we then also used decoding models that reconstructed

307  the stimulus waveforms based on the EEG data. We computed these models for each subject in the SI
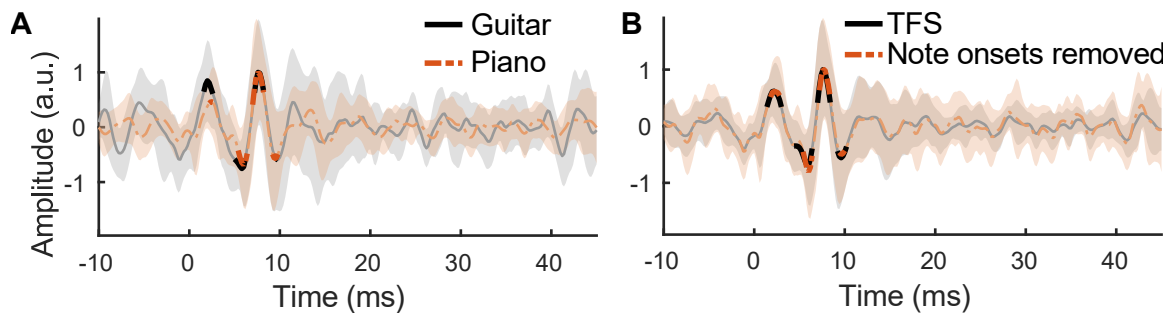
**Figure 3: Temporal Response Functions (TRFs).** (**A**), We obtained TRFs on the population level from forward models that predicted the neural responses from the stimulus temporal fine structure in the SI conditions for the guitar (black) and for the piano (red). Shaded regions denote plus/minus one standard deviation across subjects around the mean TRFs. Significant regions (thick lines) emerged at similar latencies for the guitar and piano, with a first peak at 2.2 ms, followed by a main positive peak at 7.6 ms. (**B**), We also computed TRFs for both instruments taken together from stimulus waveforms in which the note onsets where removed (red). The obtained TRFs exhibited nonetheless the same significant peaks as the TRFs from the original temporal fine structure feature (TFS, black), indicating that the neural response was not influenced by the note onsets.

308   condition. To ascertain the statistical significance of the reconstructions, we used a window from -15

309   ms to 0 ms to provide a null distribution of performance. Compared to this chance level, we found that

310   a significant reconstruction accuracy could be obtained for most subjects when using time lags from 0

311   to 15 ms for both guitar and piano (figure 4A). Indeed, significant reconstructions of the guitar

312   waveforms were obtained in 11 out of 17 subjects ($p \leq 0.05$), in 10 subjects for the piano waveforms,

313   and in 8 subjects for both types of stimuli. The reconstructions of the waveforms for the guitar and for

314   the piano were also significant at the population level (figure 4B; guitar: $p = 1.6 \cdot 10^{-2}$; piano: $p =$

315   $1.54 \cdot 10^{-3}$). Finally, on the population level, when assessing the statistical significance of the stimulus

316   reconstructions using each of three windows of interest (0 to 15 ms, 15 to 30 ms and 30 to 45 ms), we

317   found that only the window from 0 to 15 ms yielded a significant reconstruction accuracy, for either
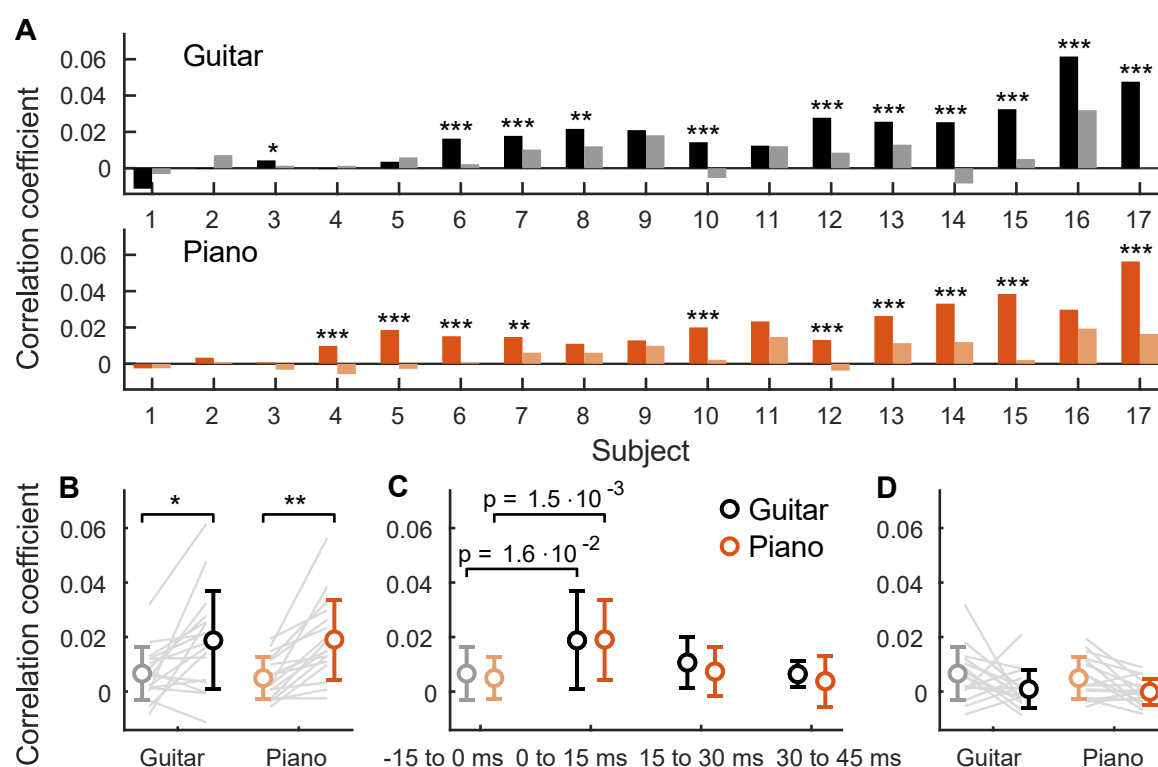
318   instrument (figure 4C).

**Figure 4: Backward models that reconstruct the stimulus waveform from the EEG data in the SI condition.**

(**A**), In most subjects, the backward models gave a stimulus reconstruction that had a significantly larger correlation (dark colour) with the original waveform than a null model (light colour). The volunteers were sorted by mean performance, and asterisks indicate $p$-values (*: $p \leq 0.05$, **: $p \leq 0.01$, ***: $p \leq 0.001$). (**B**), The mean reconstruction accuracy for each subject was used to test the significance of the reconstruction at the population level. Both the guitar and the piano stimuli could be reconstructed significantly better from the EEG recordings than from null models. (**C**), We also assessed the reconstruction of the backward models using three windows of temporal delays: 0 ms to 15 ms, 15 ms to 30 ms, and 30 ms to 45 ms (dark colours), and compared them to a null model obtained from the negative delays of -15 ms to 0 ms (light colours). Only the temporal window of 0 ms to 15 ms allowed for a stimulus reconstruction that was significantly better than that of the null model. (**D**), Reconstructing one instrument waveform using the EEG recorded during the presentation of the other instrument (0 to 15 ms window; dark colours) did not yield significant performances as compared to the null model derived by using negative delays (-15 to 0 ms; light colours).While the two instrument waveforms formed pairs corresponding to an invention, the waveform of one instrument could not be predicted from the neural responses to the other instrument.
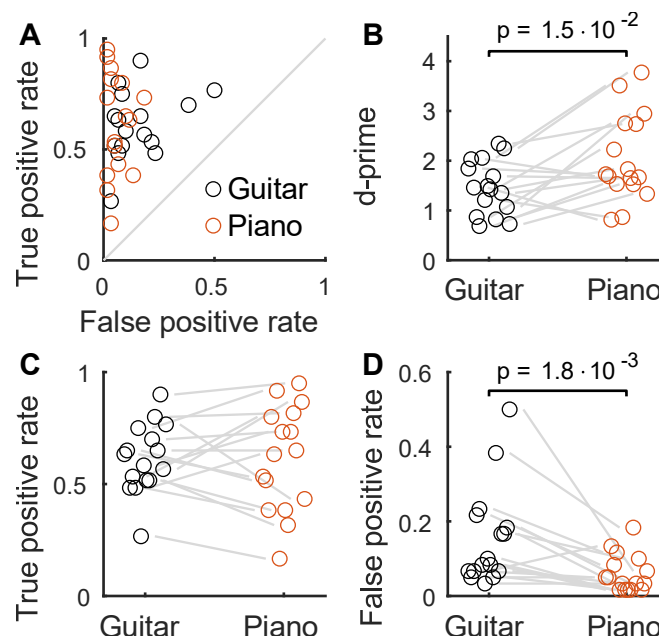
319

**Figure 5: Behavioural results for the vibrato classification, task**. Each circle represents a subject. (**A**), The receiver operating characteristics (ROC) shows that each subject performed above chance level in the CI condition, both when attending to the guitar (black) and when attending to the piano (red). (**B**) The average sensitivity index d' was significantly larger when attending to the piano than to the guitar ($p = 1.5 \cdot 10^{-2}$) with an average value of 2.0 and 1.5, respectively. (**C**), The rate of true positives was similar when attending to the guitar and then attending to the piano. (**D**), Attending to the guitar led to more false positives than attending to the piano ($p = 1.8 \cdot 10^{-3}$).

320      As the stimuli we used were derived from the left and right hands of inventions, one may wonder

321      whether two instrument waveforms derived from the same piece are independent, and whether the

322      neural responses to one instrument could be used to decode the other one. This is particularly relevant

323      in the context of the attention experiment where such an effect could obscure a putative attentional

324      modulation. However, the stimulus reconstruction accuracy when mismatching the EEG – stimuli pairs

325      in such a way (0 to 15 ms reconstruction window) was not significant as compared to the null

326      distribution using matched EEG – stimuli pairs and a -15 to 0 ms reconstruction window (figure 4D;

327      guitar: $p = 0.99$; piano: $p = 0.96$).

328      Armed with the ability to measure neural responses to the temporal fine structure of the notes in a

329      particular melody, we then investigated whether this response was affected by selective attention. To
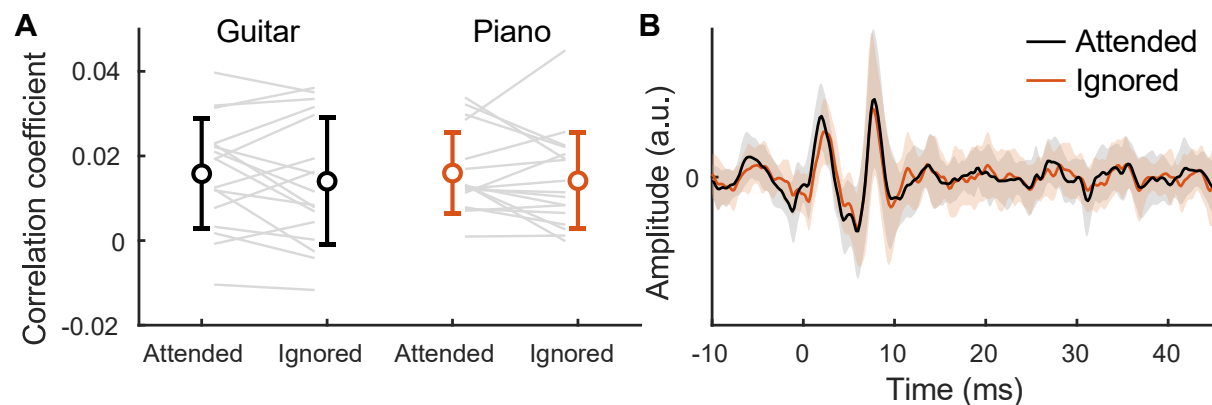
**Figure 6. Absence of attentional modulation of neural responses.** (**A**), Backward models were trained to reconstruct the stimulus waveforms for the guitar (black) and piano (red) in the CI conditions when they were attended or ignored. The reconstruction accuracies, as assessed by the correlation coefficient between the reconstructed and the original signals, did not differ significantly between the attended versus the ignored cases ($p = 0.49$ for guitar and piano). (**B**) Population-average TRFs were derived over the two CI conditions taken together for the attended (black) and ignored (red) instruments. The amplitude of the obtained TRFs did not significantly differ in the 0 ms to 15 ms region of interest. Shaded regions denote plus/minus one standard deviation across subjects around the mean TRFs.

330    this end, we analysed the CI stimuli, in which the participants had to attend selectively to one instrument

331    while ignoring the other. We monitored attention by asking the volunteers to classify vibratos inserted

332    into the melodic line played by the target instrument. The participants exhibited varied performances

333    on this task, however they all had an average performance that was better than that of a random observer,

334    as shown by their receiver operating characteristics, when selectively attending to either of the two

335    instruments (figure 5A). Accordingly, at the population level, the TPR was significantly larger than the

336    FPR when attending to either instrument ($p < 10^{-3}$ for guitar and piano). The sensitivity index d' was

337    significantly larger when attending to the piano than when attending to the guitar ($p = 1.5 \cdot 10^{-2}$), with

338    an average value 2.0 and 1.5, respectively (figure 5B). The TPR did not differ significatively between

339    the two CI conditions ($p = 0.88$; figure 5C), but the FPR was significantly higher when the subjects

340    were attending to the guitar compared to the piano (FPR: $p = 1.79 \cdot 10^{-3}$, figure 5D).

341    In order to test for a putative attentional modulation of the encoding of the stimulus temporal fine

342    structure, we first used backward models with a window from 0 to 15 ms to reconstruct each instrument

343   waveform when it was attended as well as when it was ignored. The reconstruction accuracies did,

344   however, not exhibit a statistically significant difference between the attended and the ignored case

345   (figure 6A; $p = 0.49$ for guitar and piano).

346   We then computed a linear forward model that included two features, the attended and ignored

347   instruments. The linear forward model was trained using the pooled data from the two CI conditions.

348   The model then allowed us to compare the amplitude of the attended and ignored TRFs at each time lag

349   from 0 to 15 ms. No significant difference between the amplitudes emerged at any temporal lag (figure

350   6B).

## Discussion

352   We showed for the first time that neural responses to the temporal fine structure of continuous musical

353   melodies can be obtained from EEG recordings using linear convolutive models. In particular, we

354   demonstrated that the EEG recordings could in part be predicted from the acoustic waveform (forward

355   model, figure 3). *Vice versa,* the temporal fine structure of the musical stimuli could be decoded from

356   the corresponding EEG recordings (backward model, figure 4). Significant responses could be obtained

357   in most individual subjects when they were exposed to about 5 minutes of a single melodic line.

358   The neural response at the population level revealed further information about its origin. Indeed, the

359   significant parts of the response, as obtained from the forward models, emerged most strongly at the

360   latency of 7.6 ms (figure 3A). The responses at the other latencies may have reflected our use of high-

361   pass filters for the EEG data, which spread the response in time in both directions (Widmann et al.,

362   2015). The autocorrelation of the filtered EEG data exhibited sidelobes that are reminiscent of the

363   structure of some of the peaks that we obtained in the neural responses (figure 2C).

364   The backward model showed likewise that only delays between 0 ms and 15 ms allowed for a significant

365   reconstruction of the stimulus waveform. Together with the evidence from the forward model, these

366   delays suggest a sub-cortical origin of the neural response, putatively in the inferior colliculus, although

367   different sub-cortical structures may contribute as well (Bidelman, 2015, 2018; Skoe & Kraus, 2010;

368   Sohmer et al., 1977). Recent MEG work uncovered cortical contributions to the FFR in humans (Coffey

369    et al., 2016; Hartmann & Weisz, 2019; Ross et al., 2020), although they may be limited to frequencies

370    below 150 Hz (Bidelman, 2018). The scalp-recorded FFR may accordingly combine multiple

371    subcortical and cortical sources (Coffey et al., 2019). While the neural response that we have described

372    here is arguably of subcortical origin, our use of only two EEG channels may have obstructed the

373    observation of later cortical sources with different dipole orientations.

374    Neural responses can occur to both transient (e.g. clicks, onsets) and sustained (e.g. temporal fine

375    structure) features of complex stimuli. When investigating the frequency-following response (FFR), for

376    instance, these two aspects can be segregated by time regions (Skoe & Kraus, 2010). However, the

377    continuous nature of the stimuli that we used here did not allow for this type of analysis. Instead, we

378    trained a forward model with stimulus waveforms where note onsets were suppressed, and compared it

379    to a forward model trained using the intact waveforms (figure 3A,B). The two responses were strikingly

380    similar, suggesting that they are primarily driven by the sustained periodic oscillations of individual

381    notes rather than their onsets. This may be expected, as these sustained oscillations constituted most of

382    our music stimuli. In a click train, in contrast, transients dominate the temporal fine structure.

383    When the participants were presented with stimuli consisting of two competing instruments, they had

384    to selectively attend to one of them, and identify vibratos that were inserted in the melodic line of that

385    instrument. We used this task as a marker of selective attention, comparable to the use of comprehension

386    questions in the case of speech stimuli. We found that most subjects were able to identify the target

387    vibratos whilst ignoring the distractors (figure 5). The sensitivity index d' was significantly larger when

388    attending to the piano than the guitar. When attending to either instruments, the true-positive rate did

389    not significantly differ, but the false-positive rate was lower when attending to the piano, indicating that

390    this effect mediated the difference in d' values. We hypothesise that since pianos cannot naturally

391    produce vibratos, the participants may have had a bias leading to a higher propensity to attribute vibratos

392    to the guitar. The two tasks were thus overall balanced, but attending to the piano may have been

393    somewhat easier for the participants.

394    The task of attending to one of two melodic lines allowed us to investigate whether the neural response

395    to the temporal fine structure of a particular melodic line was modulated by selective attention.

396 Following our results on the statistical methods for obtaining this neural response, we employed

397 backward models to reconstruct the stimulus waveform from the EEG recording, using temporal delays

398 between 0 ms and 15 ms. We did not, however, find any significant difference between the resulting

399 reconstruction accuracies of a melodic line when it was being attended or ignored for either instruments

400 (figure 6A). To verify this result using a different methodological approach, we also trained a forward

401 model that used the attended and the ignored instruments as features. Comparing the amplitude of the

402 attended and ignored TRFs between 0 ms and 15 ms did not reveal any significant difference (figure

403 6B).

404 Our negative finding regarding attentional modulation contrasts with previous work on similar neural

405 responses to the temporal fine structure of speech, that were found to be modulated by selective attention

406 (Etard et al., 2019; Forte et al., 2017). It also contrasts with recent MEG work that showed that the

407 cortical components of the FFR can be modulated by intermodal attention (Hartmann & Weisz, 2019).

408 These differences may point to underlying differences between music and speech. First, the two melodic

409 lines that we used in the present work may have been difficult to selectively attend, since they originated

410 from one musical piece, were contrapuntal, and often followed or responded to each other. The resulting

411 interaction between the two melodic lines makes their juxtaposition rather different from that of two

412 independent competing voices that do not interact but merely generate informational and acoustical

413 masking. While two competing speakers may encourage selective attention and neural processing of

414 one of them, our two melodic lines may therefore rather encourage attention, as well as neural

415 processing, of the acoustic mixture.

416 Second, the subjects that participated in the competing speaker experiments effectively had a lifelong

417 training in isolating one speaker from noise, due to the relevance of this task in daily life. As already

418 hinted at above, we speculate that musical stimuli are instead generally perceived as a whole, and that

419 most subjects are unfamiliar with focussing on one of several instruments. Musicians, in contrast, may

420 in general be more familiar and trained at this task. Previous studies have indeed demonstrated that

421 subcortical encoding of the temporal fine structure and FFR responses can exhibit long-term plasticity

422 and that they can be modulated by musical experience (Bidelman, Gandour, et al., 2011; Bidelman,

423    Krishnan, et al., 2011; Kraus & White-Schwoch, 2017). Similarly, musicians might exhibit attentional

424    modulation of the neural response to the temporal fine structure of melodies, although people without

425    musical training might not.

426    Finally, this study design was informed by published work analysing similar neural responses to speech

427    (Etard et al., 2019; Forte et al., 2017; Maddox & Lee, 2018). A combination of the factors listed above

428    may have contributed to produce neural responses differing from the ones previously reported for

429    speech stimuli, and thus yielding no attentional modulation, or one of a much smaller magnitude.

430    Further work is required to disentangle the potential effects of these hypotheses.

431    Music is a rich signal that consists of many transient and sustained features. Here, we focussed on the

432    comparatively high-frequency neural response to the temporal fine structure. Other features, however,

433    could be studied as well from the same stimuli, including notably cortical responses to the onsets of the

434    notes as well as to amplitude fluctuations. Similar cortical responses to continuous speech have received

435    significant attention in the past years, and have been shown to reflect attention (Ding & Simon, 2012a;

436    O'Sullivan et al., 2015; Power et al., 2012) as well as semantic features (Broderick et al., 2018),

437    surprisal (Weissbart et al., 2019) or comprehension (Etard & Reichenbach, 2019; Kösem & van

438    Wassenhove, 2017). It has indeed been found recently that the cortical encoding of sequences of tones

439    in a melody reflects a listener's expectation of the upcoming notes (Di Liberto et al., 2020). Studying

440    the interaction of such cortical responses with the subcortical activity related to the temporal fine

441    structure that we have uncovered here may further clarify the neural mechanisms that allow us to

442    perceive complex musical stimuli in their entirety, while also allowing us to selectively focus on a

443    particular instrument or melodic line.

## 444    References

445    Aiken, S. J., & Picton, T. W. (2008). Human cortical responses to the speech envelope. *Ear and*

446        *Hearing*, *29*(2), 139–157. https://doi.org/10.1097/AUD.0b013e31816453dc

447    Bidelman, G. M. (2015). Multichannel recordings of the human brainstem frequency-following

448        response: Scalp topography, source generators, and distinctions from the transient ABR.

449  *Hearing Research*, *323*, 68–80. https://doi.org/10.1016/j.heares.2015.01.011

450 Bidelman, G. M. (2018). Subcortical sources dominate the neuroelectric auditory frequency-following

451  response to speech. *NeuroImage*, *175*, 56–69. https://doi.org/10.1016/j.neuroimage.2018.03.060

452 Bidelman, G. M., Gandour, J. T., & Krishnan, A. (2011). Cross-domain effects of music and language

453  experience on the representation of pitch in the human auditory brainstem. *Journal of Cognitive*

454  *Neuroscience*. https://doi.org/10.1162/jocn.2009.21362

455 Bidelman, G. M., Krishnan, A., & Gandour, J. T. (2011). Enhanced brainstem encoding predicts

456  musicians' perceptual advantages with pitch. *European Journal of Neuroscience*.

457  https://doi.org/10.1111/j.1460-9568.2010.07527.x

458 Biesmans, W., Das, N., Francart, T., & Bertrand, A. (2017). Auditory-inspired speech envelope

459  extraction methods for improved EEG-based auditory attention detection in a cocktail party

460  scenario. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, *25*(5), 402–

461  412. https://doi.org/10.1109/TNSRE.2016.2571900

462 Bregman, A. S., Liao, C., & Levitan, R. (1990). Auditory grouping based on fundamental frequency

463  and formant peak frequency. *Canadian Journal of Psychology*.

464  https://doi.org/10.1037/h0084255

465 Bregman, Albert S. (1994). *Auditory scene analysis: The perceptual organization of sound*. MIT

466  press.

467 Broderick, M. P., Anderson, A. J., Di Liberto, G. M., Crosse, M. J., & Lalor, E. C. (2018).

468  Electrophysiological Correlates of Semantic Dissimilarity Reflect the Comprehension of

469  Natural, Narrative Speech. *Current Biology*, *28*(5), 803-809.e3.

470  https://doi.org/10.1016/J.CUB.2018.01.080

471 Cherry, E. C. (1953). Some Experiments on the Recognition of Speech, with One and with Two Ears.

472  *The Journal of the Acoustical Society of America*, *25*(5), 975–979.

473  https://doi.org/10.1121/1.1907229

474    Coffey, E. B. J., Herholz, S. C., Chepesiuk, A. M. P., Baillet, S., & Zatorre, R. J. (2016). Cortical

475        contributions to the auditory frequency-following response revealed by MEG. *Nature*

476        *Communications*, *7*, 1–11. https://doi.org/10.1038/ncomms11070

477    Coffey, E. B. J., Nicol, T., White-Schwoch, T., Chandrasekaran, B., Krizman, J., Skoe, E., Zatorre, R.

478        J., & Kraus, N. (2019). Evolving perspectives on the sources of the frequency-following

479        response. *Nature Communications*. https://doi.org/10.1038/s41467-019-13003-w

480    Cross, I., Hallam, S., & Thaut, M. (2008). The Oxford Handbook of Music Psychology. In *The*

481        *Oxford Handbook of Music Psychology*.

482        https://doi.org/10.1093/oxfordhb/9780199298457.001.0001

483    Crosse, M. J., Di Liberto, G. M., Bednar, A., & Lalor, E. C. (2016). The Multivariate Temporal

484        Response Function (mTRF) Toolbox: A MATLAB Toolbox for Relating Neural Signals to

485        Continuous Stimuli. *Frontiers in Human Neuroscience*, *10*, 604.

486        https://doi.org/10.3389/fnhum.2016.00604

487    de Cheveigné, A., Kawahara, H., Tsuzaki, M., & Aikawa, K. (1997). Concurrent vowel identification.

488        I. Effects of relative amplitude and F0 difference. *The Journal of the Acoustical Society of*

489        *America*. https://doi.org/10.1121/1.418517

490    Delorme, A., & Makeig, S. (2004). EEGLAB: An open source toolbox for analysis of single-trial

491        EEG dynamics including independent component analysis. *Journal of Neuroscience Methods*,

492        *134*(1), 9–21. https://doi.org/10.1016/j.jneumeth.2003.10.009

493    Di Liberto, G. M., O'Sullivan, J. A., & Lalor, E. C. (2015). Low-frequency cortical entrainment to

494        speech reflects phoneme-level processing. *Current Biology*, *25*(19), 2457–2465.

495        https://doi.org/10.1016/j.cub.2015.08.030

496    Di Liberto, G. M., Pelofi, C., Bianco, R., Patel, P., Mehta, A. D., Herrero, J. L., de Cheveigné, A.,

497        Shamma, S., & Mesgarani, N. (2020). Cortical encoding of melodic expectations in human

498        temporal cortex. *ELife*. https://doi.org/10.7554/eLife.51784

499    Ding, N., & Simon, J. Z. (2012a). Emergence of neural encoding of auditory objects while listening to

500        competing speakers. *Proceedings of the National Academy of Sciences*, *109*(29), 11854–11859.

501        https://doi.org/10.1073/pnas.1205381109

502    Ding, N., & Simon, J. Z. (2012b). Neural coding of continuous speech in auditory cortex during

503        monaural and dichotic listening. *Journal of Neurophysiology*, *107*(1), 78–89.

504        https://doi.org/10.1152/jn.00297.2011

505    Ding, N., & Simon, J. Z. (2014). Cortical entrainment to continuous speech: functional roles and

506        interpretations. *Frontiers in Human Neuroscience*, *8*, 311.

507        https://doi.org/10.3389/fnhum.2014.00311

508    Eerola, T., & Toiviainen, P. (2004). *MIDI toolbox: Matlab tools for music research*. University of

509        Jyväskylä: Kopijyvä, Jyväskylä, Finland. www.jyu.fi/musica/miditoolbox/

510    Etard, O., Kegler, M., Braiman, C., Forte, A. E., & Reichenbach, T. (2019). Decoding of selective

511        attention to continuous speech from the human auditory brainstem response. *NeuroImage*, *200*,

512        1–11. https://doi.org/10.1016/j.neuroimage.2019.06.029

513    Etard, O., & Reichenbach, T. (2019). Neural speech tracking in the theta and in the delta frequency

514        band differentially encode clarity and comprehension of speech in noise. *The Journal of*

515        *Neuroscience*, *39*(29), 5750–5759. https://doi.org/10.1523/jneurosci.1828-18.2019

516    Forte, A. E., Etard, O., & Reichenbach, T. (2017). The human auditory brainstem response to running

517        speech reveals a subcortical mechanism for selective attention. *ELife*, *6*.

518        https://doi.org/10.7554/eLife.27203

519    Hartmann, T., & Weisz, N. (2019). Auditory cortical generators of the Frequency Following Response

520        are modulated by intermodal attention. *NeuroImage*.

521        https://doi.org/10.1016/j.neuroimage.2019.116185

522    Hastie, T., Tibshirani, R., & Friedman, J. (2009). The Elements of Statistical Learning. *Elements*, *2*.

523        https://doi.org/10.1007/b94608

524     Haykin, S., & Chen, Z. (2005). The Cocktail Party Problem. *Neural Comput.*, *17*(9), 1875–1902.

525         https://doi.org/10.1162/0899766054322964

526     Kösem, A., & van Wassenhove, V. (2017). Distinct contributions of low- and high-frequency neural

527         oscillations to speech comprehension. *Language, Cognition and Neuroscience*, *32*(5), 536–544.

528         https://doi.org/10.1080/23273798.2016.1238495

529     Kraus, N., & White-Schwoch, T. (2017). Neurobiology of Everyday Communication: What Have We

530         Learned from Music? In *Neuroscientist*. https://doi.org/10.1177/1073858416653593

531     Krizman, J., & Kraus, N. (2019). Analyzing the FFR: A tutorial for decoding the richness of auditory

532         function. *Hearing Research*. https://doi.org/10.1016/j.heares.2019.107779

533     Lalor, E. C., Power, A. J., Reilly, R. B., & Foxe, J. J. (2009). Resolving Precise Temporal Processing

534         Properties of the Auditory System Using Continuous Stimuli. *Journal of Neurophysiology*,

535         *102*(1), 349–359. https://doi.org/10.1152/jn.90896.2008

536     Lalor, Edmund C., & Foxe, J. J. (2010). Neural responses to uninterrupted natural speech can be

537         extracted with precise temporal resolution. *European Journal of Neuroscience*, *31*(1), 189–193.

538         https://doi.org/10.1111/j.1460-9568.2009.07055.x

539     Maddox, R. K., & Lee, A. K. C. (2018). Auditory Brainstem Responses to Continuous Natural Speech

540         in Human Listeners. *ENeuro*, *5*(1). https://doi.org/10.1523/ENEURO.0441-17.2018

541     Madsen, S. M. K., Whiteford, K. L., & Oxenham, A. J. (2017). Musicians do not benefit from

542         differences in fundamental frequency when listening to speech in competing speech

543         backgrounds. *Scientific Reports*. https://doi.org/10.1038/s41598-017-12937-9

544     Micheyl, C., & Oxenham, A. J. (2010). Pitch, harmonicity and concurrent sound segregation:

545         Psychoacoustical and neurophysiological findings. In *Hearing Research*.

546         https://doi.org/10.1016/j.heares.2009.09.012

547     Nourski, K. V, Reale, R. A., Oya, H., Kawasaki, H., Kovach, C. K., Chen, H., Howard, M. A., &

548         Brugge, J. F. (2009). Temporal envelope of time-compressed speech represented in the human

549        auditory cortex. *The Journal of Neuroscience*, *29*(49), 15564–15574.

550        https://doi.org/10.1523/JNEUROSCI.3065-09.2009

551    O'Sullivan, J. A., Power, A. J., Mesgarani, N., Rajaram, S., Foxe, J. J., Shinn-Cunningham, B. G.,

552        Slaney, M., Shamma, S. A., & Lalor, E. C. (2015). Attentional Selection in a Cocktail Party

553        Environment Can Be Decoded from Single-Trial EEG. *Cerebral Cortex*, *25*(7), 1697–1706.

554        https://doi.org/10.1093/cercor/bht355

555    Oxenham, A. J. (2008). Pitch Perception and Auditory Stream Segregation: Implications for Hearing

556        Loss and Cochlear Implants. *Trends in Amplification*.

557        https://doi.org/10.1177/1084713808325881

558    Pasley, B. N., David, S. V., Mesgarani, N., Flinker, A., Shamma, S. A., Crone, N. E., Knight, R. T., &

559        Chang, E. F. (2012). Reconstructing speech from human auditory cortex. *PLoS Biology*, *10*(1).

560        https://doi.org/10.1371/journal.pbio.1001251

561    Power, A. J., Foxe, J. J., Forde, E. J., Reilly, R. B., & Lalor, E. C. (2012). At what time is the cocktail

562        party? A late locus of selective attention to natural speech. *European Journal of Neuroscience*,

563        *35*(9), 1497–1503. https://doi.org/10.1111/j.1460-9568.2012.08060.x

564    Ross, B., Tremblay, K. L., & Alain, C. (2020). Simultaneous EEG and MEG recordings reveal vocal

565        pitch elicited cortical gamma oscillations in young and older adults. *NeuroImage*, 116253.

566        https://doi.org/10.1016/j.neuroimage.2019.116253

567    Skoe, E., & Kraus, N. (2010). Auditory brainstem reponse to complex sounds : a tutorial. *Ear Hear*,

568        *31*(3), 302–324. https://doi.org/10.1097/AUD.0b013e3181cdb272.Auditory

569    Sohmer, H., Pratt, H., & Kinarti, R. (1977). Sources of frequency following responses (FFR) in man.

570        *Electroencephalography and Clinical Neurophysiology*, *42*(5), 656–664.

571        https://doi.org/10.1016/0013-4694(77)90282-6

572    Weissbart, H., Kandylaki, K. D., & Reichenbach, T. (2019). Cortical tracking of surprisal during

573        continuous speech comprehension. *Journal of Cognitive Neuroscience*.

574          https://doi.org/10.1162/jocn_a_01467

575    Widmann, A., Schröger, E., & Maess, B. (2015). Digital filter design for electrophysiological data - a

576          practical approach. *Journal of Neuroscience Methods*, *250*, 34–46.

577          https://doi.org/10.1016/j.jneumeth.2014.08.002

578    Wöstmann, M., Fiedler, L., & Obleser, J. (2017). Tracking the signal, cracking the code: speech and

579          speech comprehension in non-invasive human electrophysiology. In *Language, Cognition and*

580          *Neuroscience* (Vol. 32, Issue 7, pp. 855–869). https://doi.org/10.1080/23273798.2016.1262051

581