1    **Title**

2

3    Acoustic timbre recognition

4

5    **Synonyms**

6

7    Sound source identification; Auditory recognition

8

9    **Definition**

10

11    Timbre is what allows a listener to distinguish two sounds that have otherwise the same

12    subjective pitch, loudness, location, and duration. For instance, when orchestral

13    musicians tune at the beginning of a concert, they all play the same note, but one can still

14    tell the difference between instruments. This is largely because of timbre.

15

16    **Detailed Description**

17

18    The standard definition of timbre has several shortcomings. First, it says what timbre is

19    not, rather than what it is. Second, it relates to the comparison between two sound

20    tokens, whereas a more useful function for hearing is to associate a single timbre

21    directly with a sound source (the timbre of the piano, the timbre of the voice of a friend).

22    Perhaps as a consequence, there is still a lively debate about the acoustic features,

23    mental representations, and neural mechanisms underlying timbre recognition. Here,

24    we first outline the basic principles that make timbre such a powerful potential cue for

25    sound source identification. Then we put forward two possible approaches to timbre,

26    which we follow into the fields of acoustics, perception, neural mechanisms, and

27    computational applications.

28

29    *Why do different sound sources produce different timbres?*

30

31    Sound sources are physical objects that come in all shapes and sizes. Sound is produced

32    when some energy makes the object vibrate. The vibrations spread around the source,

33    which then propagate to the air and reach the ear of a listener in the form of pressure

34   waves (Figure 1). Simple physics shows that the wave pattern at the ear can contain a lot

35   of information about what happened at the source (Helmholtz, 1877). For instance, if the

36   energy input was brief, such as a door knock, the chances are that the sound itself will be

37   brief and have most of its energy concentrated around the time of the knock. After the

38   knock, the way the door continues to vibrate is closely related to its geometry, because

39   some wave patterns are consistent with some geometries and some are not. One such

40   rule is that waves with low frequency and thus a long wavelength are not stable within

41   small objects. Thus, the proportions of different frequency components that combine to

42   make the sound of a door knock will be constrained by the size of the door. Other, more

43   complex rules apply, depending on the shape of the object, the nature of the materials

44   involved, and so on.

45

46   Being able to decode the intricate links between wave patterns and sound sources is

47   extremely useful for humans and other animals. It allows the auditory system to serve as

48   a warning sense, for instance to identify sound-producing objects that are out of sight.

49   For people, it is also the very basis of spoken language: vowels and consonants are

50   produced by modulating the shape of the vocal apparatus, resulting in changes in timbre

51   that are the building blocks of oral communication.

52

53

54   *Dimensions versus features*

55

56   There is no consensus on what makes timbre recognition possible for human listeners.

57   To outline current controversies, it is useful to consider two opposite viewpoints

58   (Figure 2). A first view is that timbre is composed of a reasonably small number of

59   perceptual dimensions, which are subjective descriptions of sound just as pitch or

60   loudness. Such dimensions must be metameric, in that several different sounds may

61   project to the same point on the dimension.

62

63   A second view is that timbre recognition relies on the distinctive features of a given

64   sound source, learnt through experience and selected amongst a very large space of

65   potential features. The grain of a friend's voice may be unique, which is what allows us

66   to recognize her instantly. Such features would be conceptually different from

67  dimensions in that a feature does not necessarily apply to all possible sound sources; in
68  fact, it is precisely because it is unique to only a few sources (or even a single source)
69  that it could be efficient for recognition.
70
71  It is likely that a full account of timbre will lie somewhat in between these two simplified
72  hypotheses. However, for clarity, we continue to contrast each approach for different
73  aspects of timbre research.
74
75  *Sound representations*
76
77  To investigate timbre, it is useful to represent sound visually. Classically, this has been
78  done with tools such as the trace of the pressure waveform over time; the spectral
79  analysis of component frequencies through e.g. Fourier analysis; or spectro-temporal
80  transformations such as the short-term Fourier transform or wavelet analyses. More
81  recently, computational models that aim to mimic peripheral or central auditory
82  processing have been suggested (e.g. Patil *et al.*, 2012).
83
84  In the "dimensions" approach, summary statistics are computed on sound
85  representations to define what are referred to as descriptors of timbre. For instance, the
86  center of mass of all frequency components of a sound produces a single number that is
87  correlated with the apparent "brightness" of a sound (McAdams *et al.*, 1995). In the
88  "features" approach, the tendency is rather to maximize the richness of the
89  representation, by including complex spectro-temporal selectivities. Such a feature-
90  based representation need not be orderly. It can be over-complete with thousands of
91  partially overlapping features, or sparse, in the sense that a given sound would only
92  activate a small number of features within that large possible space (Hromadka and
93  Zador, 2009).
94
95  *Perceptual data*
96
97  The basic aim of the dimensions approach is to uncover the nature and number of the
98  perceptual dimensions underlying timbre. To this effect, statistical techniques based on
99  multidimensional scaling have been used: a pair of sounds is presented to the listener,

100 who has to rate how similar to each other the two sounds seem. This is repeated for all

101 possible pairs within a given sound set. Then, the similarity judgments are treated as

102 perceptual distances and used to obtain the dimensionality and geometry of the

103 corresponding mental representation. For musical instruments, classic studies point

104 towards two to three main dimensions: one related to the attack time, one related to the

105 spectral centre of mass, and one additional dimension that is less consistently observed

106 (Grey, 1977; McAdams *et al.*, 1995). More recent investigations, using both

107 multidimensional scaling and verbal descriptions, suggest five main dimensions with

108 more complex interpretations (Elliott *et al.*, 2013).

109

110 In the features approach, the focus is not on similarity but rather on the recognition of

111 the sound source. Again using musical instruments, fast recognition times have been

112 observed (Agus *et al.*, 2012) and recognition was found to be preserved even for

113 severely impoverished signals (Suied *et al.*, 2013). Moreover, recognition was faster and

114 more robust for highly familiar sources such as the human voice, an observation that

115 could not be traced back to simple acoustic dimensions (Agus *et al.*, 2012). These results

116 strongly suggest the existence of diagnostic features that were learnt by listeners,

117 through experience, to recognize e.g. voices in a robust and efficient manner.

118

119 *Neural bases*

120

121 Neural correlates of generic timbre dimensions have been investigated with brain

122 imaging. Using an EEG paradigm to probe sensory memory known as mismatch

123 negativity, it has been found that timbre dimensions such as brightness or onset time

124 could each be represented separately within auditory cortex (Caclin *et al.*, 2006).

125

126 From the features perspective, single-unit recordings have uncovered a rich variety of

127 selectivities, at many levels of the auditory system, often without any obvious ordering

128 principle (other than by frequency). Using linear analysis techniques such as reverse

129 correlation, spectro-temporal receptive fields have been derived. Various spectral and

130 temporal modulation preferences have been observed e.g. in primary auditory cortex

131 (Depireux *et al.*, 2001). Adding a nonlinear component to the analysis adds another layer

132 of complexity (Machens *et al.*, 2004). Furthermore, the neural encoding of timbre may

133    interact with supposedly independent sound characteristics, such as pitch or location

134    (Bizley *et al.*, 2009).

135

136    A further question is whether the identity of a source will be encoded by the activity of a

137    wide network shared by many sound sources, or by the activity of only a small network

138    specifically tuned to that source category. Evidence has been put forward for both

139    models. Using fMRI, the identity of a sound source can be inferred from distributed

140    activity (Staeren *et al.*, 2009). At the same time, there are clear indications of localized

141    brain areas specialized for familiar sound sources such as the human voice (Belin,

142    2006).

143

144    *Timbre recognition by machines*

145

146    There are several applications for acoustic timbre recognition, such as speaker

147    identification or music information retrieval. Even though the techniques used are fast-

148    evolving and a detailed description is beyond the scope of this section, it is interesting to

149    note that the dimensions vs. features contrast can also be seen in the architectures of the

150    computational systems.

151

152    Automatic speech recognition, which can to some extent be viewed as a timbre-decoding

153    exercise, has a long tradition of performing classification on a small number of generic

154    coefficients (e.g. mel-frequency cepstrum coefficients and their variants, Hermansky,

155    1990). For musical instruments, a descriptors-based approach has been directly

156    inspired by the perceptual dimensions of multidimensional studies, with a reasonably

157    small number of explicit descriptors (Peeters *et al.*, 2011). However, other systems exist

158    that are based on feature generation from a huge potential feature space, followed by *ad*

159    *hoc* selection for a given classification task (Coath and Denham, 2005; Pachet and Roy,

160    2009). For musical-instrument classification, machine-learning algorithms applied on a

161    high-dimensional auditory model representation have also been successfully

162    demonstrated (Patil *et al.*, 2012).

163

164    *Perspectives*

165

166 The outstanding issues for timbre research will probably benefit from considering the

167 various strategies available to a listener. For instance, when asked for subjective

168 distance judgments, the most reasonable thing to do may be to abstract common

169 dimensions to a sound set, and then use those for the comparisons. However, when

170 asked to recognize a source as fast as possible, the mere presence of a diagnostic feature

171 may be sufficient. The set of useful timbre dimensions or features can also depend on

172 the task: for a same set of spoken words, different strategies are used if listeners are

173 asked to identify the speaker or report the word content (Formisano *et al.*, 2008).

174 Finally, the very neural representation of timbre may be dynamically tuned to the

175 immediate acoustic context, through rapid plasticity (Fritz *et al.*, 2003). A fundamental

176 reason that makes timbre so elusive may therefore be that timbre recognition is a

177 profoundly adaptive mechanism, able to create and use opportunistic strategies that

178 depend on the sounds and task at hand.

179

180

181 **Cross-References/Related terms (optional)**

182

183 Pulse Resonance Sounds; Auditory Event Related Potentials

184

185 **References**

186

187 Agus, T. R., Suied, C., Thorpe, S. J., and Pressnitzer, D. (**2012**). "Fast recognition of

188     musical sounds based on timbre," J Acoust Soc Am **131**, 4124-4133.

189 Belin, P. (**2006**). "Voice processing in human and non-human primates," Philosophical

190     transactions of the Royal Society of London. Series B, Biological sciences **361**,

191     2091-2107.

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

193     "Interdependent encoding of pitch, timbre, and spatial location in auditory

194     cortex," The Journal of neuroscience : the official journal of the Society for

195     Neuroscience **29**, 2064-2075.

196 Caclin, A., Brattico, E., Tervaniemi, M., Naatanen, R., Morlet, D., Giard, M. H., and

197     McAdams, S. (**2006**). "Separate neural processing of timbre dimensions in

198     auditory sensory memory," Journal of cognitive neuroscience **18**, 1959-1972.

199    Coath, M., and Denham, S. L. (**2005**). "Robust sound classification through the
200        representation of similarity using response fields derived from stimuli during
201        early experience," Biological cybernetics **93**, 22-30.

202    Depireux, D. A., Simon, J. Z., Klein, D. J., and Shamma, S. A. (**2001**). "Spectro-temporal
203        response field characterization with dynamic ripples in ferret primary auditory
204        cortex," Journal of neurophysiology **85**, 1220-1234.

205    Elliott, T. M., Hamilton, L. S., and Theunissen, F. E. (**2013**). "Acoustic structure of the five
206        perceptual dimensions of timbre in orchestral instrument tones," J Acoust Soc Am
207        **133**, 389-404.

208    Formisano, E., De Martino, F., Bonte, M., and Goebel, R. (**2008**). ""Who" is saying "what"?
209        Brain-based decoding of human voice and speech," Science **322**, 970-973.

210    Fritz, J., Shamma, S., Elhilali, M., and Klein, D. (**2003**). "Rapid task-related plasticity of
211        spectrotemporal receptive fields in primary auditory cortex," Nature
212        neuroscience **6**, 1216-1223.

213    Grey, J. M. (**1977**). "Multidimensional perceptual scaling of musical timbres," J Acoust
214        Soc Am **61**, 1270-1277.

215    Helmholtz, H. (**1877**). *On the sensations of tone* (Dover, New York).

216    Hermansky, H. (**1990**). "Perceptual linear predictive (PLP) analysis of speech," J Acoust
217        Soc Am **87**, 1738-1752.

218    Hromadka, T., and Zador, A. M. (**2009**). "Representations in auditory cortex," Current
219        opinion in neurobiology **19**, 430-433.

220    Machens, C. K., Wehr, M. S., and Zador, A. M. (**2004**). "Linearity of cortical receptive fields
221        measured with natural sounds," The Journal of neuroscience : the official journal
222        of the Society for Neuroscience **24**, 1089-1100.

223    McAdams, S., Winsberg, S., Donnadieu, S., De Soete, G., and Krimphoff, J. (**1995**).
224        "Perceptual scaling of synthesized musical timbres: common dimensions,
225        specificities, and latent subject classes," Psychological research **58**, 177-192.

226    Pachet, F., and Roy, P. (**2009**). "Analytical features: a knowledge-based approach to
227        audio feature generation," EURASIP Journal on Audio, Speech, and Music
228        Processing **2009**.

229    Patil, K., Pressnitzer, D., Shamma, S., and Elhilali, M. (**2012**). "Music in our ears: the
230        biological bases of musical timbre perception," PLoS computational biology **8**,
231        e1002759.

232  Peeters, G., Giordano, B. L., Susini, P., Misdariis, N., and McAdams, S. (**2011**). "The Timbre

233      Toolbox: extracting audio descriptors from musical signals," J Acoust Soc Am

234      **130**, 2902-2916.

235  Staeren, N., Renvall, H., De Martino, F., Goebel, R., and Formisano, E. (**2009**). "Sound

236      Categories Are Represented as Distributed Patterns in the Human Auditory

237      Cortex," Current Biology **19**, 498-502.

238  Suied, C., Agus, T. R., Thorpe, S., and Pressnitzer, D. (**2013**). "Processing of short auditory

239      stimuli: The Rapid Audio Sequential Presentation paradigm (RASP)." in *Basic*

240      *Aspects of Hearing: Physiology and Perception*, edited by B. C. J. Moore, R. D.

241      Patterson, I. M. Winter, R. P. Carlyon, and H. E. Gockel (Springer, New York).

242

243
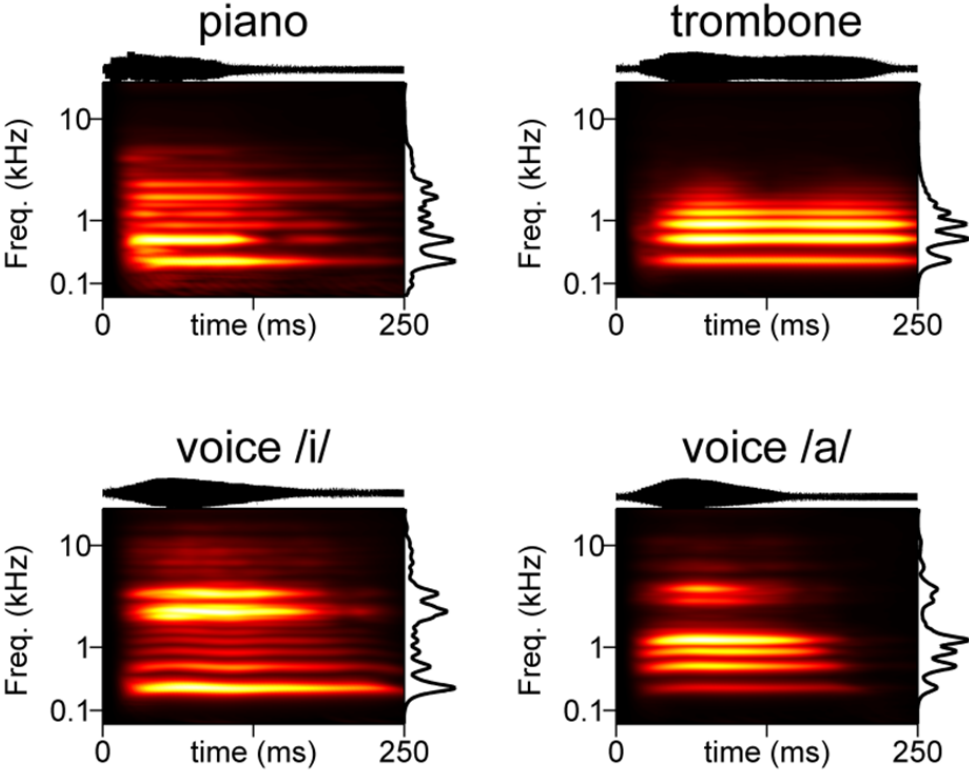
244

245

246  **Figure legends**

247

248

249  **Figure 1.** Visual representations of four sounds with the same duration, loudness and

250  pitch, so only differing by timbre. Each panel displays a time-frequency analysis derived

251  from an auditory model (see Agus *et al.*, 2012 for details). Briefly, color indicates the

252  pattern of energy within frequency channels (*y*-axis) as it evolves over time (*x*-axis).

253  The top trace is the corresponding pressure waveform. The right-hand trace is the

254  average energy over time. The two instruments illustrate classic dimensions of timbre:

255  depending on the sound source and how it is excited, the attack time can be fast (piano)

256  or slow (trombone); the spectral centre of mass can be high (piano) or low (trombone).

257  The two vowels illustrate that other, possibly more complex features may also be used

258  to distinguish e.g. vowels from instruments, or vowels from each other.

259

260  **Figure 2.** Schematic representation of the dimensions approach versus the features

261  approach for timbre. A) For the dimensions approach, all different timbres can be

262  projected in a low-dimensional space of continuous dimensions. B) For the features

263  approach, each timbre is defined by a set of distinctive features among a very large and

264  unordered set of possible features.

**Figure 1**



265

**Figure 2**