

Modeling the onset advantage in musical instrument recognition

Kai Siedenburg,¹ Marc René Schädler,¹ and David Hülsmeyer¹

Dept. of Medical Physics and Acoustics and Cluster of Excellence Hearing4all,

Carl von Ossietzky University of Oldenburg, Oldenburg,

Germany

(Dated: 26 November 2019)

1 Sound onsets provide particularly valuable cues for musical instrument identification
2 by human listeners. It has yet remained unclear whether this *onset advantage* is
3 due to enhanced perceptual encoding or the richness of acoustical information during
4 onsets. Here this issue was approached by modeling a recent study on instrument
5 identification from tone excerpts [Siedenburg 2019, JASA, 145(2), 1078-1087]. A
6 simple Hidden Markov Model classifier with separable Gabor filterbank features sim-
7 ulated human performance and replicated the onset advantage observed previously
8 for human listeners. These results provide evidence that the onset advantage may
9 be driven by the distinct acoustic qualities of onsets.

Copyright (2019) Acoustical Society of America. This article may be downloaded for personal use only. Any other use requires prior permission of the author and the Acoustical Society of America.

The following article appeared in K. Siedenburg, M.R. Schädler, D. Hülsmeyer (in press). Modeling the onset advantage in musical instrument identification The Journal of the Acoustical Society of America, 146(6), EL523–EL529 and may be found at <https://doi.org/10.1121/1.5141369>

10 I. INTRODUCTION

11 The identification of musical instruments is a central task in music perception (e.g., [Rent-](#)
12 [frow and Levitin, 2019](#)). Research on the acoustical underpinnings of instrument identifi-
13 cation still constitutes rough terrain, mainly because the candidate acoustical feature sets
14 are high-dimensional and redundant ([Handel, 1995](#); [McAdams, 2019](#)). A landmark effect
15 concerns sound onsets, which are suspected to provide particularly valuable cues for instru-
16 ment identification: if presented with sound excerpts, human listeners more easily identify
17 instrument sounds from onset portions compared to other portions of the sound ([Saldanha](#)
18 [and Corso, 1964](#); [Schaeffer, 2017](#))—a behavioral effect that we refer to as *onset advantage*.
19 Importantly, this effect does not imply that all instrumental sounds become unidentifiable
20 without onsets, because informative cues may be extracted across the full sound duration
21 and the degree to which this is possible may depend on the specific instrument at hand ([Agus](#)
22 [et al., 2019](#)). However, the psychoacoustical factors playing into the onset advantage largely
23 remain unclear. Auditory modeling approaches are in a position to provide valuable insights
24 into this issue.

25 Previous research has started to characterize the onset advantage, although results have
26 not been unequivocal. Among the more recent studies, [Suied et al. \(2014\)](#) had listeners
27 categorize sounds into broad categories such as sung voices, percussion sounds, or string
28 instruments, using gated excerpts of musical sounds. Categorization performance was above
29 chance for very short gates: 4 ms for voices and 8 ms for instruments whereas identification
30 scores were at ceiling at 64 ms gate duration or more. The authors obtained mixed results

31 regarding the importance of onsets: instrumental sounds, but not vocal sounds benefited
32 from gates being positioned at sound onsets. This means for vocal sounds there was sufficient
33 redundancy of cues across the whole duration for listeners to achieve robust categorization,
34 which may be due to the general robustness of voice recognition (Agus *et al.*, 2012). Ogg *et al.*
35 (2017) measured the durations required for human listeners to discriminate between musical
36 instrument sounds, human speech, and human environmental sounds. Results suggested
37 that listeners required 25 ms for robust discrimination and that the presence of onsets was
38 beneficial, in particular for instrument sounds. In the present study, only musical instrument
39 sounds were considered.

40 For a refined discussion, it is important to distinguish between the notions of onset and
41 the so-called *transient*. Here, transients are understood as short-lived and stochastic sound
42 bursts that are measurable in the sound signal (e.g., the hammer hitting the piano string
43 without the quasi-stationary sound waves that emanate from the harmonically vibrating
44 string). Therefore, transients should not be confused with the onset as a whole: all sounds
45 have onsets but not necessarily pronounced transients (e.g., the smooth onset of a clarinet
46 sound). Siedenburg (2019) then quantified the individual contribution of transient compo-
47 nents to instrument identification. Stationary and transient components were extracted from
48 the audio signal and instrument identification was tested for gated excerpts containing sta-
49 tionary plus transient components, or stationary components alone. Results indicated that
50 the omission of transient components at the onset impaired identification accuracy only by 6
51 percentage points. A much stronger effect was obtained by shifting the position of the gate
52 from the onset to the middle portion of the tone, impairing overall identification accuracy

53 by 25 percentage points. These results portrayed short-lived transient as of relatively minor
54 importance in instrument identification compared to the importance of retaining the onset.

55 Nonetheless, in the experiment by [Siedenburg \(2019\)](#) the important question regarding
56 the origin of the onset advantage was left open: Are listeners focussing on informative
57 acoustic features that are only available during the onset? Or do equally informative acoustic
58 features exist throughout the full duration of instrumental sounds that listeners ignore, either
59 because of their redundancy or because of the particular salience of onsets in auditory neural
60 processing? Based on the analysis of timbre dissimilarity ratings, [Grey \(1977\)](#) suggested that
61 the buildup of sinusoidal components acts as a perceptual dimension of musical instrument
62 sound. Because of the various differences between dissimilarity rating and identification
63 tasks (cf., [Siedenburg and McAdams, 2017](#)), and because of a lack of replication this finding
64 has not been very conclusive. Alternatively, one may suppose that neural coding in the
65 auditory system is tuned to onsets. It is known that already the cochlear nucleus exhibits
66 specialized onset units ([Rhode and Greenberg, 1992](#)) and neurons all along the auditory
67 pathway exhibit particularly strong responses to onsets ([Heil, 1997](#)). In psychophysics,
68 onset dominance in binaural processing has been thoroughly documented ([Houtgast and](#)
69 [Aoki, 1994](#)) and the adaptation mechanisms implemented in models of auditory processing
70 yield a pronounced response overshoot at onsets ([Jepsen et al., 2008](#)), even for simple signals
71 such as ramped sinusoids. All together, these factors suggest a more elaborate encoding
72 of onsets which in turn could imply that acoustic onset features are taken as more reliable
73 for sound identification, whether they are acoustically more informative or not. Thus, the
74 degree to which acoustical and neural factors play into the onset advantage remains unclear.

75 Here, a modeling approach is used to further disentangle these issues. We utilize a
76 Hidden Markov Model (HMM) classifier in conjunction with separable Gabor filterbank
77 (SGBFB) features, an approach that has proven valuable in the domain of speech recognition
78 and psychoacoustic modelling for normal-hearing listeners as part of the Framework for
79 Auditory Discrimination Experiments (FADE) ([Schädler et al., 2016](#)). Originally conceived
80 as speech recognizer, FADE has quite successfully modelled human performance on a variety
81 of psychoacoustic tasks without requiring any internal calibration data beyond training on
82 the specific task at hand. As additional baseline features, we use Mel-frequency cepstral
83 coefficients (MFCCs) and log-Mel spectra. The simulations are set up in an analogous way
84 to the main experiment from [Siedenburg \(2019\)](#), and the results are interpreted in terms of
85 their implications on the role of acoustical factors in the onset advantage.

86 II. METHODS

87 A. Previous experiment used for modeling

88 The main experiment by [Siedenburg \(2019\)](#) serves as the starting point for the present
89 modeling study. That experiment was divided into a training and a test phase. In the
90 training phase, musician participants were presented with sounds from ten test instruments:
91 piano, guitar, harp, vibraphone, marimba, trumpet, clarinet, flute, violin, and cello. Sounds
92 were of 250 ms duration and for each instrument there were sounds with twelve different pitch
93 levels (C4/262Hz to B4/494Hz). Subsequently, listeners were trained in the identification
94 task and obtained feedback on 60 trials (6 per instrument). Notably, before the start of each

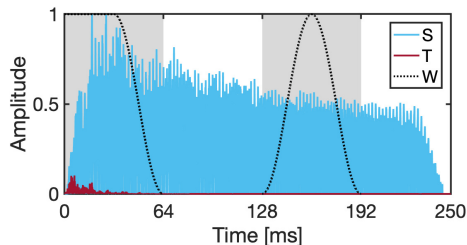


FIG. 1. (Color online). A rectified waveform of a piano tone, decomposed into stationary (S) and transient (T) components. The dashed line corresponds to the gating window (W). In the experimental conditions S+T@0ms and S@0ms, the down-ramped window starting at 0 ms was used. In the condition S+T@128ms, the window extending from 128 to 192 ms was used.

95 experimental block, participants again went through a passive exposure phase, listening
 96 to all the original 250 ms sounds, as in the very first part of the training. This means,
 97 participants had extensive exposure to the full 250 ms sounds before being tested on short
 98 excerpts. In the subsequent test phase, listeners identified instruments from 64 ms segments
 99 of sounds composed of stationary and transient components ($S+T$) or stationary components
 100 alone (S), taken either from the onset ($@0ms$), or from the middle portion of the sounds
 101 ($@128ms$). Stationary and transient extraction was achieved by applying a specifically
 102 developed algorithm (Siedenburg and Doclo, 2017). Figure 1 shows the example of a piano
 103 tone, including the stationary and transient components and the gating window at the $@0ms$
 104 and $@128ms$ positions.

105 B. Modeling rationale

106 The modeling was set up in such a way as to create an analogous scenario compared to
 107 the instrument identification experiment with human listeners: a classifier was trained on

108 the set of full 250 ms sounds and tested on the short 64 ms excerpts. Hence, the scenario
109 required the classifier to generalize to sounds with durations of around one fourth of the
110 training sounds.

111 As back-end, we used a classic HMM classifier with one Gaussian component per state.
112 Several points speak for using HMM for the present modeling task: HMM perform well with
113 small training sets, explicitly encode the temporal dimension of auditory stimuli, and are not
114 overly powerful (in comparison to more recent architectures such as deep neural networks,
115 which may even learn high-level tasks from raw data), hence allowing us to differentiate
116 the quality of the acoustic input features. Recent research has demonstrated that this
117 modeling approach is well suited for various aspects of auditory modeling, including speech
118 intelligibility, elementary psychoacoustics, and hearing loss (Kollmeier *et al.*, 2016; Schädler
119 *et al.*, 2016· 2018). Here, we tested HMM with a variable number of 1-6 states. For every
120 instrument and number of states, separate classifiers were trained for stimuli adjusted to
121 the input levels 65 dB SPL \pm 0, 3, 6, 9, 12, 15 dB. This range was selected to cover short-
122 and long-term level changes that usually occur in music recordings. Each classifier was
123 then tested in the 65 dB condition and we report average recognition performance. A more
124 detailed description of the modeling architecture is provided by Schädler *et al.* (2016).

125 C. Acoustic features

126 Three sets of acoustic features were used in the simulations: log-Mel Spectrograms, Mel-
127 frequency cepstral coefficients (MFCC), and separable Gabor Filterbank features (SGBFB).
128 Log-Mel spectra were computed with a window length of 25 ms and 10 ms hop size, and

129 the linear frequency axis was subsequently warped to 36 bins with Mel spacing, that is,
130 with frequency centers between 64 and 11874 Hz. MFCCs were computed by applying
131 a discrete cosine transform in the spectral dimension, and the first 21 coefficients were
132 used. These were concatenated with the first and second order derivatives along the time
133 axis (the so-called *delta* and *double-delta* coefficients, respectively). SGBFB features were
134 computed by using a temporal and a spectral modulation filterbank operating on the log-
135 Mel spectrogram. These covered spectral modulations from 0.03 to 0.25 cycles per channel,
136 and temporal modulations from 6.2 to 25 Hz. The combination of temporal and spectral
137 filters resulted in a feature set of 570 coefficients (as compared to 63 MFCC and 36 Log-Mel
138 features). The SGBFB and MFCC features were used with mean and variance normalization,
139 the log-Mel spectrum was used without normalization. For more information on the feature
140 sets and details of their implementation, the reader is referred to [Schädler *et al.* \(2012, 2016\)](#)
141 or the model code¹.

142 III. RESULTS

143 A. Performance comparison

144 Accuracies (i.e., proportion correct classifications) for classifiers trained on the full 250 ms
145 sounds are depicted in Fig. 2, together with experimental results from human listeners.
146 The SGBFB features performed best and yielded accuracies of .99, .62, .55, and .41 (aver-
147 aged across the different number of states of the classifier) for the four conditions 0-250ms,
148 S+T@0ms, S@0ms, and S+T@128ms, respectively. MFCCs obtained weaker results with

149 averages of .99, .31, .27, .20 for the four conditions. This means that both feature sets
150 could be easily fitted to the training set 0-250ms, but generalized considerably worse to the
151 shorter test excerpts. With accuracies of .61, .36, .35, and .31 across the four conditions,
152 log-Mel spectra did not yield a similarly good fit to the training set, but slightly better
153 results compared to the MFCC on the short test excerpts. Importantly, both the SGBFB
154 and MFCC features still provided a pattern of results that qualitatively resembled that of
155 human listeners.

156 As further visible in Fig. 2, the number of states of the HMM was a critical factor for
157 the classification performance: SGBFB yields highest performance for an HMM with three
158 states, and notably, this classifier was the only one that reaches human performance in
159 the test conditions up to error tolerance (95% confidence intervals of human performance).
160 With MFCCs, the best classifier contained five states and performed worse than human
161 performance by around 20 percentage points, but showed a very similar decay of performance
162 across the test conditions. Surprisingly, log-Mel spectra did not at all resemble the human
163 pattern of performance and the best performing classifier had only one state, that is, it only
164 encoded static spectral information. Note that log-Mel spectra are the only features that
165 did not explicitly encode spectro-temporal information or modulations (SGBFB are tailored
166 to do so; MFCC do so by virtue of its discrete cosine transform and the delta-coefficients).
167 This supports the view that robust generalization in instrument identification relies on the
168 explicit encoding of spectro-temporal information (cf., [Patil et al., 2012](#)).

169 None of the classifiers performed better than human listeners in the test conditions, which
170 leaves open the possibility that this was due to general difficulties in classifying the short

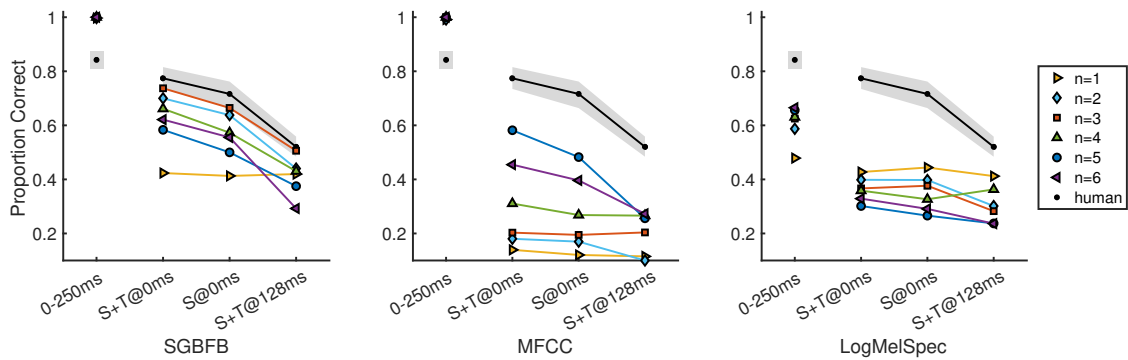


FIG. 2. (Color online). Identification results for the three different feature sets. Test conditions on x-axis: condition 0-250 ms refers to the unaltered sounds with 250 ms duration (the reference condition for human listeners and the model). The other conditions are indexed by whether stationary or transient components were used: for instance, $S + T@0\text{ms}$ indicates that test signals consisted of stationary (S) and transient (T) components with the gating window of 64 ms length starting at 0 ms. Performance for the three different feature sets is provided in the three different panels. The number of states of the HMM is indicated in the legend. The shaded area indicates 95% confidence intervals of human scores.

171 64 ms excerpts. Table I provides accuracies for classifiers that were trained on a merged set
 172 of sounds, containing both the full 250 ms sounds and the 64 ms excerpts from the three
 173 conditions $S+T@0\text{ms}$, $S@0\text{ms}$, and $S+T@128\text{ms}$. With this merged training set, perfor-
 174 mance was much better for the short excerpts with averages of 75% correct classifications
 175 or more for the SGBFB and MFCC models. That is, the recognition performance for this
 176 merged training set was on par with or even exceed human performance. Therefore, this

TABLE I. Proportion correct classification for models trained on a merged set of both 250 ms and 64 ms sounds. Columns index the test sets, rows the feature sets. Table entries correspond to mean and range (square brackets) across the number of states (1-6).

	0-250ms	S+T@0ms	S@0ms	S+T@128ms
SGBFB	.60 [.54, .70]	.83 [.79, .86]	.76 [.72, .79]	.80 [.79, .82]
MFCC	.27 [.21, .34]	.83 [.72, .88]	.75 [.63, .79]	.75 [.69, .78]
Log-Mel Spec	.56 [.45, .68]	.48 [.40, .55]	.49 [.44, .53]	.40 [.37, .43]

177 latter simulation suggests that the primary difficulty for the present classifiers was not to
 178 classify short excerpts per se, but to generalize from 250 ms sounds to 64 ms excerpts.

179 B. Recognition of excerpts over time

180 In order to more detailedly probe the distribution of acoustic information over time, the
 181 classifiers trained on the full 250 ms sounds were tested on 64 ms excerpts that were obtained
 182 through gating with a raised cosine window starting at different temporal positions (0, 32,
 183 64, ..., 192 ms). Note that the last excerpt only had a duration of 58 ms (extending from 192-
 184 250 ms). An additional test condition extending from 0-64 ms was included, which did not
 185 use the full cosine window but only the fade out part, hence preserving the original attack
 186 and closely resembling the experiment condition S+T@0ms (although it also contained the
 187 residual noise that was missing in S+T@0ms). Here, classifiers with three or five states

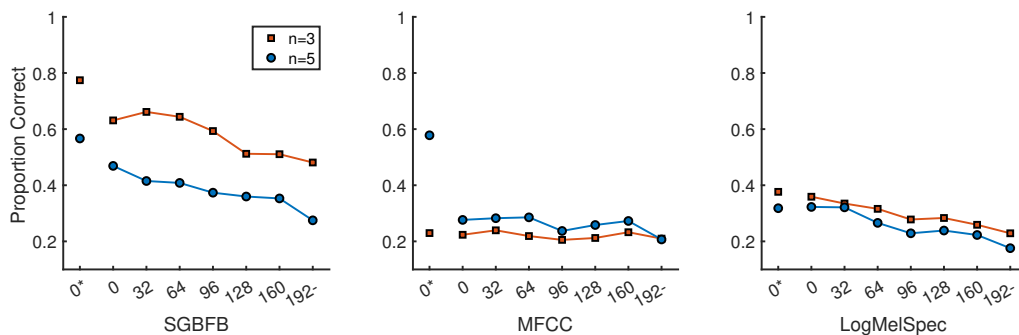


FIG. 3. (Color online). Identification results for models with different feature sets and number of states tested on different excerpts. The x-axis indicates the starting point (in ms) of the 64 ms gating window. The condition 0* denotes the first 64 ms of the sound without fade in, but fade out; all other sounds contain both fade in and out. The condition 192- corresponds to the gate extending from 192-250 ms.

188 were evaluated because these models had performed best on the test set in conjunction with
 189 SGBFB and MFCC features.

190 Figure 3 depicts the results for the three different feature sets. Performance dropped by
 191 more than 10 percentage points from the 0-64ms* to the 0-64ms condition for the classifier
 192 using SGBFB features with three states. The SGBFB features further exhibited a gradual
 193 decline of performance over time after the onset, which suggests that the distinctiveness of
 194 acoustical cues gradually worsens over time for classifiers trained on the full 250 ms sounds.
 195 Even more drastic was the drop of more than 30 percentage points for the classifier using
 196 MFCC features with five states, which continued to perform poorly for the consecutive
 197 excerpts. This finding indicates that classifiers that best performed in these simulations are

198 those that rely on the encoding of precise onset information, providing further evidence for
199 acoustic information to play an important role in the onset advantage.

200 **IV. CONCLUSION**

201 To investigate the importance of acoustic cues for musical instrument identification, we
202 trained HMM classifiers using SGBFB, MFCC, and log-Mel spectrum features on 250 ms
203 sounds and tested generalization to short 64 ms excerpts that contained stationary and
204 transient sound components. Classifiers using SGBFB with three states generalized best
205 to the test excerpts and showed very similar results compared to human listeners. Testing
206 the classifiers on excerpts gated at different time points of the sound indicated that the
207 best performing classifiers rely on precise onset information. More specifically, performance
208 dropped drastically when the initial onset was manipulated through gating and performance
209 degraded gradually the more the excerpts stemmed from later portions of the sound. These
210 results provide converging evidence that the acoustical richness of sound onsets itself could
211 be exploited by listeners as a cue, which then may give rise to the onset advantage.

212 It is important to bear in mind that the current results stem from a scenario which
213 trained the classifiers on the full 250 ms sounds, that is, that the full sounds acted as a
214 reference. Although this seems to be the best possible analogy to the main experiment of
215 [Siedenburg \(2019\)](#) where listeners were trained and heavily (re-)exposed to the full sounds,
216 future research could contextualize this scenario by exposing listeners and machine classifiers
217 to sounds of varying durations in a training phase. These pursuits could also attempt to

218 account for the long-term knowledge about more diverse classes of instrument sounds that
219 the musician participants may have utilized in the experiment.

220 These simulations demonstrate that a simple Hidden Markov Model classifier with sepa-
221 rable Gabor filterbank features replicates the onset advantage in instrument identification.
222 The classifier essentially implemented an elaborate encoding of acoustic information and no
223 component of the classifier architecture was dedicated to onsets per se (e.g. much in contrast
224 to [Newton and Smith, 2012](#)). On the basis of these results, one does not need to assume a
225 specialized neural encoding of onsets to explain the onset advantage in instrument identifi-
226 cation. In the convoluted reality of auditory processing, however, it may well be the case
227 that the acoustical properties and the neural encoding of sounds act in concert and jointly
228 contribute to the onset advantage. In fact, neural sound coding could even have evolved
229 to optimally exploit the acoustic richness of sound onsets ([Młynarski and McDermott, 2018](#);
230 [Theunissen and Elie, 2014](#)).

231 **ACKNOWLEDGMENTS**

232 The authors thank the anonymous reviewers for insightful comments. KS has received
233 funding from the European Unions Framework Programme for Research and Innovation
234 Horizon 2020 (2014-2020) under the Marie Skłodowska-Curie Grant Agreement No. 747124.
235 This work was also funded by the Deutsche Forschungsgemeinschaft (DFG, German Re-
236 search Foundation) Project ID 390895286 EXC 2177/1 and by the Deutsche Forschungs-
237 gemeinschaft (DFG, German Research Foundation) Projektnummer 352015383 SFB 1330
238 A 3.

239 ¹Model code is available at <https://github.com/m-r-s/fade/>

240

241 Agus, T. R., Suied, C., and Pressnitzer, D. (2019). “Timbre recognition and sound source
242 identification,” in *Timbre: Acoustics, Perception, and Cognition*, edited by K. Siedenburg,
243 C. Saitis, S. McAdams, A. N. Popper, and R. R. Fay (Springer), pp. 59–85.

244 Agus, T. R., Suied, C., Thorpe, S. J., and Pressnitzer, D. (2012). “Fast recognition of
245 musical sounds based on timbre,” *The Journal of the Acoustical Society of America* **131**(5),
246 4124–4133.

247 Grey, J. M. (1977). “Multidimensional perceptual scaling of musical timbres,” *The Journal*
248 *of the Acoustical Society of America* **61**(5), 1270–1277.

249 Handel, S. (1995). “Timbre perception and auditory object identification,” in *Hearing*,
250 edited by B. C. Moore, *Handbook of Perception and Cognition* (Academic Press, San
251 Diego, CA), pp. 425–461.

252 Heil, P. (1997). “Auditory cortical onset responses revisited. I. First spike timing,” *Journal*
253 *of Neurophysiology* **77**(5), 2616–2641.

254 Houtgast, T., and Aoki, S. (1994). “Stimulus-onset dominance in the perception of binaural
255 information,” *Hearing research* **72**(1-2), 29–36.

256 Jepsen, M. L., Ewert, S. D., and Dau, T. (2008). “A computational model of human audi-
257 tory signal processing and perception.,” *The Journal of the Acoustical Society of America*
258 **124**(1), 422–438.

259 Kollmeier, B., Schädler, M. R., Warzybok, A., Meyer, B. T., and Brand, T. (2016). “Sen-
260 tence recognition prediction for hearing-impaired listeners in stationary and fluctuation

261 noise with fade: Empowering the attenuation and distortion concept by plomp with a
262 quantitative processing model,” *Trends in Hearing* **20**, doi: 10.1177/2331216516655795.

263 McAdams, S. (2019). “The perceptual representation of timbre,” in *Timbre: Acoustics,*
264 *Perception, and Cognition*, edited by K. Siedenburg, C. Saitis, S. McAdams, A. N. Popper,
265 and R. R. Fay (Springer), pp. 23–57.

266 Młynarski, W., and McDermott, J. H. (2018). “Learning midlevel auditory codes from
267 natural sound statistics,” *Neural Computation* **30**(3), 631–669.

268 Newton, M. J., and Smith, L. S. (2012). “A neurally inspired musical instrument classi-
269 fication system based upon the sound onset,” *The Journal of the Acoustical Society of*
270 *America* **131**(6), 4785–4798.

271 Ogg, M., Slevc, L. R., and Idsardi, W. J. (2017). “The time course of sound category
272 identification: Insights from acoustic features,” *The Journal of the Acoustical Society of*
273 *America* **142**(6), 3459–3473.

274 Patil, K., Pressnitzer, D., Shamma, S. A., and Elhilali, M. (2012). “Music in our ears:
275 The biological bases of musical timbre perception,” *PLOS Computational Biology* **8**(11),
276 e1002759.

277 Rentfrow, P. J., and Levitin, D. J. (2019). *Foundations in music psychology: Theory and*
278 *research* (MIT Press, Cambridge, MA).

279 Rhode, W. S., and Greenberg, S. (1992). “Physiology of the cochlear nuclei,” in *The mam-*
280 *malian auditory pathway: Neurophysiology*, edited by A. N. Popper and R. R. Fay, Springer
281 *Handbook of Auditory Research* (Springer, Heidelberg, Germany), pp. 94–152.

282 Saldanha, E., and Corso, J. F. (1964). “Timbre cues and the identification of musical
283 instruments,” *The Journal of the Acoustical Society of America* **36**(11), 2021–2026.

284 Schädler, M. R., Meyer, B. T., and Kollmeier, B. (2012). “Spectro-temporal modulation
285 subspace-spanning filter bank features for robust automatic speech recognition,” *The Jour-
286 nal of the Acoustical Society of America* **131**(5), 4134–4151.

287 Schädler, M. R., Warzybok, A., Ewert, S. D., and Kollmeier, B. (2016). “A simulation
288 framework for auditory discrimination experiments: Revealing the importance of across-
289 frequency processing in speech perception,” *The Journal of the Acoustical Society of Amer-
290 ica* **139**(5), 2708–2722.

291 Schädler, M. R., Warzybok, A., and Kollmeier, B. (2018). “Objective prediction of
292 hearing aid benefit across listener groups using machine learning: Speech recognition
293 performance with binaural noise-reduction algorithms,” *Trends in Hearing* **22**, doi:
294 10.1177/2331216518768954.

295 Schaeffer, P. (2017). *Treatise on Musical Objects: An Essay Across Disciplines*, **20** (Univ
296 of California Press).

297 Siedenburg, K. (2019). “Specifying the perceptual relevance of onset transients for musical
298 instrument identification,” *The journal of the Acoustical Society of America* **145**(2), 1078–
299 1087.

300 Siedenburg, K., and Doclo, S. (2017). “Iterative structured shrinkage algorithms for station-
301 ary/transient audio separation,” in *Proc. of the 20th Int. Conf. on Digital Audio Effects*
302 (*DAFx-20*), *Edinburgh, Sep 5–8*.

303 Siedenburg, K., and McAdams, S. (2017). “Four distinctions for the auditory “wastebasket”
304 of timbre,” *Frontiers in Psychology* **8**, 1747.

305 Suied, C., Agus, T. R., Thorpe, S. J., Mesgarani, N., and Pressnitzer, D. (2014). “Auditory
306 gist: Recognition of very short sounds from timbre cues,” *Journal of the Acoustical Society*
307 *of America* **135**(3), 1380–1391.

308 Theunissen, F. E., and Elie, J. E. (2014). “Neural processing of natural sounds,” *Nature*
309 *Reviews Neuroscience* **15**(6), 355–366.