# Specifying the perceptual relevance of onset transients for musical instrument identification

Kai Siedenburg<sup>1</sup>

Department of Medical Physics and Acoustics, Carl von Ossietzky University of

Oldenburg, Oldenburg, Germany<sup>a</sup>)

(Dated: 28 February 2019)

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 (2019). Specifying the relevance of onset transients in musical instrument identification. The Journal of the Acoustical Society of America, 145(2), 1078–1087 and may be found at https://asa.scitation.org/doi/ 10.1121/1.5091778.

Sound onsets are commonly considered to play a privileged role in the identification of 1 musical instruments, but the underlying acoustic features remain unclear. By using 2 sounds resynthesized with and without rapidly varying transients (not to be confused 3 with the onset as a whole), this study set out to specify precisely the role of transients 4 and quasi-stationary components in the perception of musical instrument sounds. In 5 experiment 1, listeners were trained to identify ten instruments from 250 ms sounds. 6 In a subsequent test phase, listeners identified instruments from 64 ms segments of 7 sounds presented with or without transient components, either taken from the onset, 8 or from the middle portion of the sounds. The omission of transient components at the 9 onset impaired overall identification accuracy only by 6%, even though experiment 2 10 suggested that their omission was discriminable. Shifting the position of the gate from 11 the onset to the middle portion of the tone impaired overall identification accuracy 12 by 25%. Taken together, these findings confirm the prominent status of onsets in 13 musical instrument identification, but suggest that rapidly varying transients are less 14 indicative of instrument identity compared to the relatively slow build-up of sinusoidal 15 components during onsets. 16

<sup>&</sup>lt;sup>a)</sup>kai.siedenburg@uni-oldenburg.de;

## 17 I. INTRODUCTION

It is a common idea in music psychoacoustics that timbre cues at sound onsets are of 18 central importance for the identification of musical instruments by human listeners. Acous-19 tical explorations of this idea may date back as far as the 1940s, when the advent of tape 20 recording technology allowed sounds to be systematically manipulated by means of cutting 21 and splicing. The radio engineer and musician Pierre Schaeffer pioneered in testing the per-22 ceptual implications of different temporal gatings of sounds (cf., Schaeffer, 2017) and made 23 the observation that sounds such as piano tones lose aspects of their identity if presented 24 bare of onsets. This has led to the idea that onset information is perceptually more valu-25 able compared to other sound components that are present in the so-called *steady state*, the 26 portion of a tone where its waveform (or short-time spectrum) is relatively constant. As of 27 today, however, surprisingly little is known about the specific acoustic ingredients that give 28 rise to this effect. 29

A component of specific importance to onsets is the so-called *transient*. Here, transients 30 are defined as short-lived and chaotic bursts of acoustical energy, such as the sound of the 31 hammer hitting the piano string (without the sound from the harmonically vibrating string). 32 It is important to note that according to this definition, transients should not be confused 33 with the full onset: all sounds have onsets but not necessarily pronounced transients—think 34 of a clarinet tone with a smooth attack. Neither do transients exclusively occur at the 35 onset—think of the return of the hopper of the harpsichord at the release of the key (usually 36 accompanied by sustained harmonic resonance in the soundboard). 37

Regarding the perceptual identification of instruments, rapidly varying onset transients 38 are often claimed to be of prime importance, particular in the audio processing literature 39 (Daudet, 2005; Zaunschirm et al., 2012), although no definitive proof has been provided 40 to date. There yet exist alternative acoustic properties of sound onsets that could bear 41 diagnostic information about sound identity, such as the comparatively slow build-up of 42 sinusoids which could be particularly informative at sound onsets (Grey, 1977). The primary 43 goal of the present study was to better understand the relevance of transients and more slowly 44 varying sinusoidal components for the identification of musical sounds. 45

## 46 A. Previous research

Rigorous empirical research on instrument identification has emerged in the 1960s. Early 47 studies used tape recordings of musical instrument tones that were manipulated by means 48 of cutting and splicing for experimental purposes. In a well-known study, Saldanha and 49 Corso (1964) suggested that several factors contribute to the identification of orchestral 50 instruments: pitch, the presence of vibrato, the experimental session (test/re-test), and 51 the presented excerpt (onset, steady state, offset). Although identification accuracy was 52 generally poor (around 40% correct identifications), offsets did not bear perceptually useful 53 information and shortening the steady state from 9 to 3 seconds did not negatively affect 54 the results. On the contrary, discarding onsets decreased identification accuracy by 15 55 percentage points, although performance remained above chance. Unfortunately, no clear 56 criterion was provided as to how the endings of onsets were determined and hence the 57 durations of the segments that were used as onsets remained unclear. 58

Other research from around that time came to similar conclusions regarding the role of 59 onsets. Clark Jr et al. (1963) presented excerpts from the onset or steady part of recorded 60 instrument tones to listeners with durations varying from 60 to 600 ms. The authors observed 61 that even short portions such as the first 60 ms of tones contained sufficient information for 62 musicians to discriminate instruments. Using recorded tones of 6 s duration, Elliott (1975) 63 observed that discarding the first and last half second from sustained instrument tones with 64 an overall duration of 6 seconds significantly impaired identification performance of several 65 orchestral instruments. 66

Exploring timbre dissimilarity perception, Grey (1977) used musical instrument tones 67 emulated by additive synthesis and observed that the ordering of sounds along one dimension 68 of a timbre space obtained from dissimilarity ratings corresponded to the synchronicity of the 69 onsets of sounds' sinusoidal components. Also studying dissimilarity ratings, Iverson and 70 Krumhansl (1993) tested the role of onsets by using three sets of tones: full tones (duration: 71 2-3.3 s), onsets (first 80 ms), and the remainder (first 80 ms removed). They found strong 72 commonalities between the multidimensional scaling solutions of all three sets, which was 73 interpreted as reflecting a form of acoustical invariance across segments. However, today it is 74 known that an excerpt of 80 ms can be more than enough for instrument identification (Suied 75 et al., 2014), making it likely that listeners also relied on instrument identity or sound source 76 properties in their dissimilarity judgments (cf., Siedenburg et al., 2016). Unfortunately, it 77 thus seems hard to differentiate whether the supposed invariance in Iverson and Krumhansl 78 (1993) arose from invariance of aspects of the sensory representations or from invariance 79

<sup>80</sup> in the inferred sound source mechanism (which in turn may have affected dissimilarity <sup>81</sup> judgments) or a combination of both aspects.

Subsequent research has shown that relatively short durations are necessary to discrim-82 inate instruments. Robinson and Patterson (1995) presented listeners with short sound 83 excerpts, excised from synthetic emulations of brass, flute, harpsichord, and string sounds. 84 For the identification of isolated sounds, it was observed that even for single cycles of peri-85 odic tones (corresponding to 2.9 - 30.5 s depending on pitch), musicians and nonmusicians 86 achieved an impressive performance of around 75% and 50% of correct responses, respec-87 tively. Note that because cycles were presented repeatedly, no temporal cues (onset, offset) 88 were present in the sounds, which highlights the importance of spectral cues for instrument 89 identification. In a similar vein, Suied et al. (2014) tested the minimal duration required 90 for the correct recognition of sound source categories. Listeners heard cosine-shaped gated 91 segments of musical sounds and were required to respond to target categories (sung voices, 92 percussion sounds, string instrument sounds). Categorization performance was above chance 93 for surprisingly short gates, 4 ms for voices, and 8 ms for instruments, and scores were at 94 ceiling at 64 ms gate duration. Mixed results were obtained for the effect of onset infor-95 mation: instrumental, but not vocal sounds benefited from gates being positioned at sound 96 onsets. 97

<sup>98</sup> Most recently, Thoret *et al.* (2016<sup>,</sup> 2017) showed that instrument identification is deter-<sup>99</sup> mined by specific instrument-specific spectrotemporal modulations, although their approach <sup>100</sup> did not allow them to draw specific conclusions about the role of onsets. Ogg *et al.* (2017) <sup>101</sup> studied the minimal duration required to discriminate between musical instrument sounds, <sup>102</sup> human speech, and human environmental sounds. They found that listeners required 25 ms
<sup>103</sup> for robust discrimination and that the presence of onsets was beneficial, even for vocal
<sup>104</sup> sounds.

Two conclusions may be drawn from this review regarding the role of onsets in instrument 105 identification. First, the presence of the onset portion appears to improve sound identifi-106 cation but does not seem to be strictly necessary for correct identification. The relative 107 importance of onsets appears to depend on the specific instrument at hand. Second, and 108 more generally, whether implemented by digital gating or by excised tape, the experimental 109 approach of presenting temporal segments has conceptually remained identical throughout 110 the last 60 years (even though the analog scalpel may be less precise than today's digital 111 means). This approach assumes that sounds can be meaningfully separated into discrete 112 temporal states. However, as it will be demonstrated in Sec. II, short-lived transients and 113 quasi-stationary sinusoidal components cannot be strictly separated in time because both 114 regimes overlap and one dynamically transforms into the other (Levine and Smith, 2007; 115 Reuter, 1995). Therefore, the studies outlined above can only to a limited degree allow for 116 conclusions about the importance of specific acoustical components such as transients—more 117 flexible tools for separating signal components (sharpened acoustical scalpels) are needed. 118

119 B. The present study

The goal of this study was to use a novel transient/stationary separation algorithm to circumvent some of the methodological limitations of the literature. This algorithm is described in the following section Sec. II. In the main experiment described in Sec. III, listeners iden-



FIG. 1. (Color online) Example of a piano sound A4 (440 Hz) of 250 ms separated into stationary and transient components. A) Spectrogram of original sound (window length 25 ms). B) Zoom into first 16 ms of the original sound's spectrogram. C) Waveform of separated stationary components (dark blue) and transients (light red). D) Zoom into first 16 ms of the separated components' waveform. E) Estimated stationary coefficients. F) Estimated transient coefficients. G) Waveform of residual. H) Spectrogram of residual.

tified short segments extracted from the sounds of ten musical instruments. These segments were processed by the separation algorithm and contained both stationary and transient information, or only stationary information. Segments were extracted from the onset or from the middle portion of the sound. The goal of an additional control experiment described in Sec. IV was to assess whether the transient components were generally discriminable.

## 128 II. TRANSIENT SEPARATION

## 129 A. Description of the algorithm

Developments in audio signal processing have made it possible to separate overlapping 130 stationary and transient components from mixtures (for a general review, see Müller, 2015, 131 Chap. 8). A classical approach to this problem was provided by Serra and Smith (1990), 132 approximating transients in a global manner by time-varying filtered noise. Recently, the 133 present author presented a more fine-grained algorithm to estimate transients by using an 134 iterative multi-resolution analysis (Siedenburg and Doclo, 2017). The algorithm exploits 135 the orthogonal orientation of components in the time-frequency plane: Whereas the quasi-136 stationary (S) components are sparse in frequency and persistent over time, rapidly varying 137 transient (T) are sparsely distributed in time and persistent across frequency. Both types 138 of components are extracted iteratively from Short-Term Fourier Transform (STFT) repre-139 sentations, using long window lengths (46 ms) for stationary components, yielding spectral 140 precision, and short window lengths (3 ms) for transient components, yielding temporal 141 precision. In technical terms, the separation process is based on a shrinkage operation of 142 STFT coefficients that specifically extracts coefficients which are part of groups of relatively 143 strong coefficients that extend over time or frequency (so-called *neighborhoods*, see Sieden-144 burg and Doclo, 2017; Siedenburg and Dörfler, 2011). The result is an approximation of the 145 original signal y in terms of three components, y = S + T + e, where e denotes the residual 146 signal. The residual signal usually is of rather low energy and captures reverberation and 147 microphone noise, but also faint phase-distorted versions of the stationary and transient 148

<sup>149</sup> components. The algorithm accurately separates stationary and transient components in
<sup>150</sup> synthetic examples and provides plausible separation results for recorded audio signals from
<sup>151</sup> acoustic musical instruments (although by definition there is no ground truth in this case).
<sup>152</sup> In the following experiment, S and S+T were used to study instrument identification.
<sup>153</sup> Consequently, if there was unintended distortion from the signal processing, it would have
<sup>154</sup> appeared not only in S but also in S+T.

# 155 B. Acoustic analyses

Figure 1 depicts the example of an A4 (440 Hz) piano sound of 250 ms duration. Through-156 out this study, the same settings of the algorithm were used as described in the original 157 publication (Siedenburg and Doclo, 2017). The algorithm separates the impulsive sound of 158 the hammer from the vibrating string (sound examples are provided as part of the supple-159 mentary information<sup>1</sup>). Panel A depicts the spectrogram (using a window length of 25 ms) 160 of the original sound and a zoom into the onset is shown in panel B. The figure illustrates 161 that beyond harmonic components, there is transient energy present in the onset portion of 162 the sound. Moreover, the more detailed visualization in panel B suggests that the partial 163 tones do not all start at the same time, but that lower components precede higher ones. 164 Panels C and D depict the waveform of the separated stationary and transient components. 165 The extracted time-frequency coefficients are shown in Panels E and F. Stationary compo-166 nents are sparse in frequency (although some subharmonic energy seems to be captured by 167 the stationary estimate, because of its relatively long extension in time). Transients have 168 impulsive characteristics. Notably, the extracted transients are short-lived but overlap in 169



FIG. 2. (Color online) Temporal amplitude envelopes (rows 1-2) and spectral envelopes (rows 3-4). Level corresponds to signal intensity raised by 0.3 to approximate loudness according to Steven's law. Original sounds: gray, separated stationary components: blue, transient components: red, dashed-dotted. Lines depict averages across all twelve pitch levels. For temporal amplitude envelopes (rows 1-2), shaded areas correspond to the position of the gating used in experiment 1.

time with the stationary components. This example hence demonstrates the limitations of considering musical sounds as a sequence of discrete states that can be neatly spliced apart in the time domain. To the contrary, components overlap and are continuously transformed over time, and thus transient components should not be confused with onsets as a whole.

The residual signal is depicted in panels G and H. It is visible that the residual contains residual traces of both the harmonic stationary components and the impulsive transient of this piano tone.

In the perceptual experiment reported below, ten instruments at twelve different pitch 177 levels were used (see Sec. III B 2 for details). Analyses indicated that these sounds had tran-178 sients of much lower overall energy compared to the stationary components. Specifically, the 179 stationary-to-transient energy level ratios averaged across pitch was highest for the vibra-180 phone (mean 18 dB), followed by the marimba (24 dB), trumpet (28 dB), guitar (30 dB), 181 piano (31 dB), cello (36 dB), harp (40 dB), violin (42 dB), flute (44 dB), and finally the 182 clarinet (52 dB) with the weakest transient. Somewhat surprisingly, these ratios indicate 183 that it is not generally possible to infer the sound excitation mechanisms of instruments by 184 virtue of the relative transient energies, because the harp (an impulsive instrument) had 185 lower relative transient energy compared to the trumpet (a sustained instrument). 186

The temporal evolution of transient and stationary energy is depicted in Fig. 2 (rows 1–2). The figure shows the average temporal and spectral envelopes of the stationary and transient signal components (for temporal envelopes, gray background indicates the positioning of the gates in experiment 1). Here, temporal envelopes were extracted by computing the magnitude of the analytic signal, filtered with a third-order Butterworth lowpass-filter at a

cutoff frequency of 50 Hz. The levels plotted in Fig. 2 correspond to signal intensities taken 192 to the power of 0.3 (following Steven's law to approximate loudness). Figure 2 shows that 193 the extracted transients do not extent much further than 64 ms into the tone and exhibit 194 exponential decay characteristics for the impulsive instruments. This also holds, albeit to a 195 much smaller degree, for the trumpet, violin, and cello. For the flute and clarinet, however, 196 transients are of very low intensity, potentially more reflecting continuous blowing noise. 197 Regarding the stationary component, the figure further indicates marked differences in en-198 velope slope of impulsively excited instruments (top row) compared to sustained instruments 199 (bottom row), the latter only reaching their energy peak in the middle portion of the tone. 200

The two bottom rows of Fig. 2 shows the average spectral power for the original signal, 201 and the stationary and transient components (as for temporal envelopes raised to the power 202 of 0.3 to reflect loudness). Spectral envelopes were obtained by smoothing the computed 203 magnitude spectra by using a first-order Butterworth lowpass filter with a cutoff frequency 204 of 1000 Hz. The figure illustrates that the extracted transients had energy at relatively high 205 frequencies, with spectral peaks at frequencies around or higher than 1 kHz. The figure also 206 highlights the distinct spectral shapes of the instruments' stationary components compared 207 to the relatively similar spectral shapes of transient components. Experiment 1 tested the 208 perceptual relevance of these components. 209

## 210 III. EXPERIMENT 1: INSTRUMENT IDENTIFICATION

## A. Rationale

The present experiment compared instrument identification for harmonic instrument 212 sounds resynthesized with and without transient components. In order to avoid ceiling 213 performance and to be able to account for the importance of the onset position, sounds were 214 gated with short gates of 64 ms duration. The resulting segments were taken from the onset 215 of the original sounds and presented with and without transient components. In order to 216 obtain an estimate about the general relevance of the onset, a third signal condition was 217 tested that presented segments obtained from the middle portion of sounds (128-196 ms) 218 with stationary and transient components. Note, however, that in the present sound set, 219 the energy of transients was very small for the middle portion (see Fig. 2). Therefore, ex-220 cerpts from the middle portion with only stationary components were not included in the 221 experiment. 222

#### B. Methods

### 224 1. Participants

Eighteen listeners (13 female, 4 male, 1 other) with self-reported normal hearing and a mean age of M = 26.1 years (SD = 6.7, range: 21–48) participated in this experiment. Participants had played their primary musical instrument for an average of M = 9.3 years (SD = 6.6, range: 1–22) and were dedicating M = 10.5h per week to musical activities



FIG. 3. (Color online) Illustration of the windows used to create the experimental signal conditions for the example of a piano tone. The figure shows the amplitude envelopes of stationary components (S), transient components (T), and the gating windows (W) with start positions at 0 ms or at 128 ms.

(SD = 11.0, range: 1-35). Participants were recruited via advertisements at the University of Oldenburg online job board and received a compensation with 10 EUR per hour.

#### 231 2. Stimuli and apparatus

Stimuli were derived from orchestral instrument samples, obtained from the Vienna Symphonic Library (http://vsl.co.at, last accessed June 12, 2018). The following instruments were used in this study: piano, guitar, harp, vibraphone, marimba, trumpet, clarinet, flute, violin, and cello. Guitar samples were obtained from a Yamaha P155 synthesizer. Each instrument was played at twelve pitch levels: C4 (262 Hz) to B4 (494 Hz). From the stereo samples, only the left channels were used. Tones were played at *forte* dynamics and conceived as 8th-notes at a tempo of 120 quarter notes per minute, corresponding to a duration of 250 ms. The actual recordings were longer than this and of varying duration, so a 25 ms
raised cosine function was applied as fade out to obtain a consistent duration of 250 ms.

In the experiments by Suid *et al.* (2014), instrument categorization performance levelled 241 off at a duration of 64 ms. Hence, this gate duration was chosen for the current experiment 242 in order to ensure that participants would be able to perform the task. Furthermore, this 243 gate duration was short enough to meaningfully compare different placings of the gate within 244 sounds. When the gate started at the beginning of the sound (0-64 ms: @0ms), the original 245 onset was preserved and a raised-cosine fade-out was used (cf., Suied et al., 2014). When 246 the gate was positioned in the middle of the sound (128-192 ms: @128ms), both a raised-247 cosine fade-in and fade-out was used. Gated sounds were normalized in root-mean-square 248 energy. The decomposition algorithm described above was used to extract the stationary 249 and transient signal components from the gated sounds. Overall, there were three signal 250 conditions: 1) stationary (S) and transient (T) components gated at the onset (S+T@0ms), 251 2) stationary components at the onset (S@0ms), and 3) stationary and transient components 252 in the middle of the tone (S+T@128ms). Figure 3 shows the gating function and the 253 temporal envelopes of the individual components for an exemplary piano tone. 254

The experiment was run with Matlab and sounds were converted with an RME Fireface audio interface at an audio sampling frequency of 44.1 kHz and 24 bit resolution. Sounds were presented diotically over Sennheiser HDA 200 headphones at an average level of 65 dBA SPL, as calibrated by a Norsonic Nor140 sound-level meter with a G.R.A.S. IEC 60711 artificial ear to which the headphones were coupled. Listeners were tested individually in a sound-proof lab and provided responses on a computer mouse.

#### 261 3. Procedure

The experiment comprised a training and test phase. The training phase was conducted 262 to ensure that participants were familiar with the full range of perceptual features that 263 characterized the test sounds. In the training phase, the original sounds were used. First, 264 participants were exposed to all sounds at twelve pitch levels from each one of the ten instru-265 ments at an inter-onset interval (IOI) of 750 ms. The order of the presentation of individual 266 sounds and instruments was randomized. In order to further provide visual anchors, pic-267 tures of the instruments were presented concurrently. Pictures had been obtained from a 268 web search and depicted standard tokens of the instruments in front of a white background. 269

In the second part of the training, participants were trained to identify sounds presented in isolation, as in the main experiment. The test contained each of the ten instruments at six randomly drawn pitch levels. In every trial, participants listened to a randomly drawn sound and were required to select the corresponding instrument label from a list of alternatives presented on a computer screen. Feedback about the correct response was provided with instrument labels and pictures. Overall, this amounted to 60 trials of training with response feedback and took around 12 minutes.

All participants continued with the main experiment, where sounds from the same ten instruments were presented at twelve pitch levels for the three signal conditions, S+T@0ms, S@0ms, S+T@128ms, described above (Sec. III B 2). The signal conditions were blocked and blocks were presented in random order. There were 120 sounds per block; each block took around 25 minutes to complete and there were obligatory pauses of at least five minutes



FIG. 4. (Color online) Mean identification scores from experiment 1. Individual results are plotted as gray lines and the dotted line indicates chance level. Error bars: 95% CI.

between blocks. Before the start of each experimental block, participants went through a passive exposure phase with the original sounds, as in the first part of the training. This exposure phase was implemented to ensure that potential differences across blocks were due to the signal conditions, and not due to memory loss of the reference that was established or consolidated during the initial training.

To avoid response bias through a fixed order of the instrument labels on the screen, the list order was randomized for each experimental block. Otherwise, the procedure was identical to the second part of the training although no feedback was provided. The experiment was self-paced.



FIG. 5. (Color online) Average confusion matrices from experiment 1 including the training, normalized by the number of presentations of every instrument.

291 C. Results

Fig. 4 shows the average scores for the training and all experimental conditions, together with individual profiles from all participants. In the training, identification performance was high (proportion of correct identifications: M = .84). In the main experiment, average performance in the S+T@0ms signal condition was around seven percentage points below the training score (M = .77) and slightly higher compared to the S@0ms condition (M = .71). In the S+T@128ms signal condition, there was a strong inflation of confusions (M = .52).

Figure 5 depicts average confusion matrices for the training phase and all experimen-298 tal conditions. In the training, it is visible that, surprisingly, the cello and the trumpet 299 were frequently confused at this stage (although this only occurred in the training). In 300 the main experiment, frequent within-family confusions occurred for the S+T@0ms signal 301 condition, in particular for the violin and cello (strings), and the clarinet and flute (winds). 302 The qualitative confusion patterns were very similar for the S@0ms signal condition. In 303 the S+T@128ms condition, the three impulsive instruments piano, guitar, and harp were 304 frequently confused and even attributed to wind instruments such as the clarinet. Among 305 the sustained (i.e., continuously excited) instruments, the flute was particularly poorly iden-306 tified, and often confused with the trumpet. Four instruments were robustly identified for 307 this condition and achieved accuracies above 0.75: the vibraphone, marimba, trumpet, and 308 clarinet. 309

A repeated-measures ANOVA was conducted with the factors signal condition (S+T@0ms, 310 S@0ms, S+T@128ms) and pitch level (the statistical dependency of instrument-wise accura-311 cies does not allow for an ANOVA on an instrument-wise level). The analysis indicated that 312 there were significant differences between signal conditions,  $F(2, 34) = 155.9, p < .001, \eta_p^2 = 0.001$ 313 .90, and of pitch,  $F(11, 187) = 4.52, p < .001, \eta_p^2 = .21$ , but no significant interaction be-314 tween the two,  $F(22, 374) = 1.52, p = .064, \eta_p^2 = .08$ . Post-hoc tests demonstrated that 315 scores from the three signal conditions were significantly different from each other: paired 316 t(17) = 4.3, p = .0013 for S+T@0ms vs. S@0ms, t(17) = 19.7, p < .001 for S+T@0ms 317

vs. S+T@128ms, and t(17) = 10.7, p < .001 for S@0ms vs. S+T@128ms (Bonferroni-318 corrected for multiple comparisons, n = 3). A comparison to the training indicated that 319 training scores were significantly higher compared to all experimental signal conditions, 320 paired t(17) > 5.4, p < .001. Visual inspection of the data did not reveal any systematic 321 relation of identification accuracy and pitch, and scores in none of the three signal condi-322 tions significantly correlated with pitch height, p > .187 (Bonferroni-corrected, n = 3). This 323 suggests that idiosyncratic stimulus features distributed across different pitch levels most 324 likely caused the observed differences of identification scores across pitch levels. 325

## 326 D. Discussion

This experiment compared harmonic musical instrument identification for 64 ms-long 327 sound segments with and without transient components taken from the onset or the middle 328 portion of the original sound. The data indicated that removing the transient at the sound 329 onset impaired identification scores by around 6 percentage points, whereas moving the gate 330 from the onset to the middle portion of the sound impaired identification accuracy by 25 331 percentage points. Surprisingly, this effect did not appear to strictly depend on whether 332 impulsive or sustained instruments were considered. In the signal condition that presented 333 64 ms segments from the middle portion of the tone (S+T@128ms), the vibraphone and 334 marimba were accurately identified (both impulsive) with accuracy scores above 75%, and 335 the same held for the trumpet and clarinet (both sustained). 336

As it can be observed in Fig. 2, the energy levels of the transient components were almost negligible for the tested middle portions of sounds. Furthermore, there was a drastic drop in performance from S@0ms to S+T@128ms ( $\approx$  S@128ms) paired with small differences from S+T@0ms to S@0ms. Therefore, this pattern of results likely reflects the greater diagnosticity of the cues present in quasi-stationary sinusoidal components at the sound onset, and the lack of the transient component appears to be of smaller importance. Notably, even for sustained sounds, the steady state portion probed by the S+T@128ms condition turned out to be less informative than the onset portion.

A potential explanation of the small effect observed for the removal of transients could be that the combined stationary and transient components (S+T) were not clearly discriminable from the stationary parts (S) alone. This question was addressed in a second experiment, which would further help to more comprehensively characterize the perceptual status of transients in musical instrument sounds.

# 350 IV. EXPERIMENT 2: TRANSIENT DISCRIMINATION

# 351 A. Rationale

The second experiment acted as a control experiment in order to test whether listeners would be sensitive to the presence of transients. Specifically, the aim was to test listeners' discrimination abilities of S from S+T, but also to measure discrimination of S+T from the original sound. This would assess the perceptual relevance of the separation algorithm's residual component. To direct listeners attention to transient information, additional foil conditions were included in the experiment, presenting amplified transients together with the stationary part.

#### 359 B. Methods

## 360 1. Participants

Ten listeners (4 female, 5 male, 1 other) with self-reported normal hearing and a mean age of M = 27.8 years (SD = 4.2, range: 23–37) participated. Participants had played their primary musical instrument for an average of M = 14.8 years (SD = 6.8, range: 4–30) and were dedicating M = 13.8 per week to musical activities (SD = 12.5, range: 2-35). Participant recruiting and compensation was identical to experiment 1.

#### 366 2. Stimuli and apparatus

To keep the overall duration of the experiment within limits, only four of the ten instruments from experiment 1 were tested, two of which were impulsive (vibraphone and guitar) and two sustained (cello and trumpet). The corresponding recordings were presented at the full duration of 250 ms. There were five signal conditions, each testing the discrimination of S+T against i) the original signal, ii) S, iii) 5T, iv) S+10T, v) S+15T, where S+xT indicates that the level of T was raised by x dB. The apparatus was identical to the main experiment.

#### 373 3. Procedure

A 3-interval/2-alternative forced-choice task ("odd one out") was used. On every trial, there were three intervals with inter-stimulus intervals of 250 ms and participants were required to detect the odd interval. It was randomly determined whether S+T or the comparison stimulus from signal condition i)-v) served as the odd stimulus. After providing



FIG. 6. (Color online) Discrimination accuracy from experiment 2. S+T was discriminated from the signal type given on the x-axis, where S+xT indicates that the level of T was raised by x dB. Bar color corresponds to instruments as listed in the legend (guitar, vibraphone, trumpet, cello). Triangles correspond to performance of individual participants, the dotted line indicates chance performance. Error bars: 95% CI.

their response by selecting the interval on a computer screen, participants received feedback
about the correct response.

In order to maximize participant's sensitivity to potentially idiosyncratic timbral features, the presentation of instruments was blocked with a random order of the presentation of the signal conditions. The order of blocks was randomized. Every block contained 180 trials (3 intervals x 12 pitch levels x 5 signal conditions). The completion of any one block took around 25 minutes and there were obligatory pauses of at least five minutes between blocks.

#### 385 C. Results

Performance was above chance level for all the five different signal conditions, as confirmed 386 by tailed t-tests against 0.33, t(9) > 5.5, p < .001. Participants robustly discriminated S+T 387 from S, as reflected by 69% of correct identifications in this condition. Participants had 388 greater difficulties to discriminate S+T from the original signal, yielding an average of 389 only 42% correct responses. This result indicates that the omission of the residual from the 390 original signal, leaving S+T, is barely detectable, which validates the general approach to use 391 S+T as a starting point for studying timbre perception. Participants were further sensitive 392 to an amplification of transients, as indicated by the strong effect across foil conditions. 393 Average percentage of correct responses was 54%, 86%, and 96% for discriminating S+T 394 from S+5T, S+10T, and S+15T, respectively. 395

A repeated-measures ANOVA was conducted to analyse differences for individual instruments. The analysis confirmed strong effects of signal condition,  $F(4, 36) = 171.8, p < .001, \eta_p^2 = .95$ , instrument,  $F(3, 27) = 21.4, p < .001, \eta_p^2 = .70$ , and an interaction of signal condition and instrument,  $F(12, 108) = 14.1, p < .001, \eta_p^2 = .61$ .

The mean scores of all five signal conditions were highly different from each other, t(9) > 4.3, p < .002, as visible in Fig. 6 (left panel). Performance for the two impulsive instruments guitar (76%) and vibraphone (74%) was generally better compared to the sustained instruments trumpet (62%) and cello (65%). Pairwise t-tests confirmed no significant differences between instruments of the same excitation type, t(9) < 1.5, p > .16, but all differences across excitation types were highly significant, t(9) > 5.0, p < .001. This means the task was generally easier for the two impulsive instruments, guitar and vibraphone.

The signal conditions elicited differential effects on impulsive instruments compared to continuously excited instruments such that the interaction was due to the high scores for impulsive instruments in the S condition. Specifically, scores for S did not differ significantly from the original signal for the trumpet and the cello, t(9) < 3.0, p > .057 (Bonferronicorrected for multiple comparisons, n = 4). But there were strong differences between the original signal and S signals for the guitar and the vibraphone, t(9) > 7.2, p < .001.

# 413 D. Discussion

This second experiment tested listeners' sensitivity to discriminate signals with manipu-414 lated transient components. Independent of instrument, the original sounds were only poorly 415 discriminated from the signals that were resynthesized without residual (S+T); discrimina-416 tion performance was barely above chance for this signal condition. This result implies that 417 the residual does not appear to be very important in the current separation, which suggests 418 that using the stationary and transients components, S+T, seems to be a good starting 419 point for the current pursuits. More specifically, the above chance performance in both the 420 S+T vs. S and the S+T vs. S+5T (and S+10T, S+15T) conditions indicates that listeners 421 were sensitive to the amplification as well as to omission of transients. Note that this effect 422 was pronounced for impulsive instruments, but, although not as strong (as indicated by the 423 significant interaction of signal condition and instrument), it remained present for sustained 424 instruments. In comparison to the higher performance for the signal condition that omitted 425

the transient (S+T vs. S), this indicates that listeners were much more sensitive to the presence of the transient than to the presence of the residual noise.

Only four instruments could be tested in this experiment and hence the generality of the findings is limited. It is possible that instruments with low energy transients such as the clarinet and flute would have yielded lower scores, in particular for the S vs. S+T signal condition. Nonetheless, the obtained results show that it is generally not the lack of discriminability that is the underlying reason for the small effect between the S+T@0ms and S@0ms signal condition observed throughout instruments in experiment 1.

# 434 V. CONCLUSION

This study revisited the perceptual relevance of onsets in identification and discrimination 435 tasks. Previous studies suggested that the onset plays a privileged role for identification, 436 but the underlying acoustic factors had not been thoroughly tested. Here, a relatively small 437 set of harmonic orchestral instrument sounds was used to test the importance of transient 438 signal components. Using an algorithm to dissect transient from stationary components 439 (Siedenburg and Doclo, 2017), acoustical analysis indicated that rapidly varying transients 440 and quasi-stationary components are generally overlapping in time and that transient com-441 ponents are of relatively low energy. Importantly, these analyses indicate that the transient, 442 defined via its short-livedness and stochastic nature, should not be confused with the onset 443 portion of sounds as a whole—there is no point in time where transients could be neatly 444 separated from sinusoidal components. Instead, the separation of acoustic components must 445 take place in the time-frequency domain.

Two experiments tested the perceptual relevance of transients and quasi-stationary sinu-447 solid components. In experiment 1, it was shown that the omission of transient components 448 at the onset portion of tones had a relatively small detrimental effect on instrument iden-449 tification, even though experiment 2 suggested that a lack of discriminability of signals 450 presented with and without transient components was not the underlying reason for this. 451 Therefore, these results indicate that quasi-stationary components yield the most informa-452 tive cues for instrument identification. Furthermore, shifting the position of the gate from 453 the onset to the middle portion of the tone had a large detrimental effect on identification 454 performance. The latter result confirms that even without the presence of transient com-455 ponents, onsets seem to be much more informative compared to sounds' middle portions, 456 irrespective of the specific instrument or instrument class (impulsive vs. sustained). Taken 457 together, these findings confirm the prominent status of onsets in musical instrument identi-458 fication suggested by the literature, but specify that rapidly varying transients (which often 450 but not exclusively occur at sound onsets) have relatively limited diagnostic value for the 460 identification of harmonic musical instruments. In conclusion, fairly slowly varying signal 461 components during onsets, likely the characteristic build-up of sinusoidal components in 462 particular, provide the most valuable bundle of acoustic features for perceptual instrument 463 identification. 464

A critical reader may object that the great care that musicians, sound designers, and music producers invest in the shaping of transient aspects of sound refutes this argument. This objection may be countered by noting that identification tasks require listeners to rely on informative acoustic cues for sound source identity, but not on every sound feature that

may be integrated into assessments of sound quality (e.g., Pressnitzer et al., 2013; Siedenburg 469 and McAdams, 2017). Coherent with this notion, experiments 1 and 2 collectively suggested 470 that not every class of discriminable sound feature is essential for sound source identification. 471 In effect, sound production may deal in great length with the sculpting of timbral nuances 472 such as high-frequency transients, even if these are only of minor importance for the inference 473 of sound sources. Generally, this view acknowledges the multiplicity of cues available for 474 sound source identification (Giordano et al., 2010; Handel, 1995), all of which may be used 475 opportunistically depending on the perceptual task and context at hand. Furthermore, one 476 should not forget that this study only considered harmonic musical instruments presented 477 in isolation. The situation may be different for non-harmonic percussion instruments and 478 other sound-producing objects, not to speak of sound source identification in polyphonic 479 mixtures. 480

A topic that should be addressed by future acoustical analyses concerns the question whether the utility of the onset (with or without transients) for instrument identification rests on perceptual or acoustical grounds. In other words, are listeners making use of informative features for identification that are only available in the onset, or do there exist equally informative features throughout the sound but listeners prefer to focus on the onset?

From a more general perspective, the current approach is in line with an upsurge of interest in signal analysis/re-synthesis approaches to the study of auditory perception (Mc-Dermott and Simoncelli, 2011; Overath *et al.*, 2015; Ponsot *et al.*, 2018; Thoret *et al.*, 2017). In order to unravel the intricate workings of auditory perception these types of studies develop specific signal processing tools, which allow to work with naturalistic but precisely <sup>491</sup> controlled stimuli. Although this approach is principally related to the early explorations
<sup>492</sup> of cutting and splicing tapes (Schaeffer, 2017), today's digital tools offer an unprecedented
<sup>493</sup> degree of precision and versatility.

#### 494 ACKNOWLEDGMENTS

The author wishes to thank the two reviewers for productive remarks on this manuscript. The author further thanks Henning Schepker, Etienne Thoret, and Trevor Agus for valuable comments on earlier versions of this manuscript, Daniel Pressnitzer and Christoph Reuter for insightful discussions, Saskia Röttges for data collection, and Simon Doclo for general support of this study.

This project has received funding from the European Unions Framework Programme for Research and Innovation Horizon 2020 (2014-2020) under the Marie Skodowska-Curie Grant Agreement No. 747124. This project was also funded by a Carl von Ossietzky Young Researchers' Fellowship from the University of Oldenburg.

<sup>1</sup>See supplementary material at [URL will be inserted by AIP] for sound examples.

505

<sup>509</sup> Daudet, L. (2005). "A review on techniques for the extraction of transients in musical sig-

<sup>510</sup> nals," in International Symposium on Computer Music Modeling and Retrieval, Springer,

<sup>&</sup>lt;sup>506</sup> Clark Jr, M., Luce, D., Abrams, R., Schlossberg, H., and Rome, J. (1963). "Preliminary
<sup>507</sup> experiments on the aural significance of parts of tones of orchestral instruments and on
<sup>508</sup> choral tones," Journal of the Audio Engineering Society 11(1), 45–54.

<sup>511</sup> pp. 219–232.

- Elliott, C. A. (1975). "Attacks and releases as factors in instrument identification," Journal
  of Research in Music Education 23(1), 35–40.
- Giordano, B. L., Rocchesso, D., and McAdams, S. (**2010**). "Integration of acoustical information in the perception of impacted sound sources: the role of information accuracy and exploitability.," Journal of Experimental Psychology: Human Perception and Performance **36**(2), 462–476.
- <sup>518</sup> Grey, J. M. (**1977**). "Multidimensional perceptual scaling of musical timbres," The Journal <sup>519</sup> of the Acoustical Society of America **61**(5), 1270–1277.
- Handel, S. (1995). "Timbre perception and auditory object identification," in *Hearing*,
  edited by B. C. Moore, 2 of *Handbook of Perception and Cognition* (Academic Press, San Diego, CA), pp. 425–461.
- <sup>523</sup> Iverson, P., and Krumhansl, C. L. (**1993**). "Isolating the dynamic attributes of musical <sup>524</sup> timbre," Journal of the Acoustical Society of America **94**(5), 2595–2603.
- Levine, S. N., and Smith, J. O. (**2007**). "A compact and malleable sines+transients+noise model for sound," in *Analysis, Synthesis, and Perception of Musical Sounds*, edited by J. W. Beauchamp (Springer, New York, NY), pp. 145–174.
- <sup>528</sup> McDermott, J., and Simoncelli, E. P. (**2011**). "Sound texture perception via statistics of
- the auditory periphery: Evidence from sound synthesis," Neuron **71**, 926–940.
- <sup>530</sup> Müller, M. (2015). Fundamentals of Music Processing: Audio, Analysis, Algorithms, Appli-
- <sup>531</sup> cations (Springer, Heidelberg, Germany).

- <sup>532</sup> Ogg, M., Slevc, L. R., and Idsardi, W. J. (**2017**). "The time course of sound category <sup>533</sup> identification: Insights from acoustic features," The Journal of the Acoustical Society of <sup>534</sup> America **142**(6), 3459–3473.
- <sup>535</sup> Overath, T., McDermott, J. H., Zarate, J. M., and Poeppel, D. (**2015**). "The cortical anal-<sup>536</sup> ysis of speech-specific temporal structure revealed by responses to sound quilts," Nature <sup>537</sup> neuroscience **18**(6), 903–911.
- <sup>538</sup> Ponsot, E., Arias, P., and Aucouturier, J.-J. (2018). "Uncovering mental representations of
  <sup>539</sup> smiled speech using reverse correlation," The Journal of the Acoustical Society of America
  <sup>540</sup> 143(1), EL19–EL24.
- Pressnitzer, D., Agus, T. R., and Suied, C. (2013). "Acoustic timbre recognition," in *Encyclopedia of Computational Neuroscience: Springer Reference*, edited by D. Jaeger and
  R. Jung (Springer, Heidelberg, Germany), pp. 1–6.
- Reuter, C. (1995). Der Einschwingvorgang nichtperkussiver Musikinstrumente (Peter Lang
  Frankfurt/M).
- <sup>546</sup> Robinson, K., and Patterson, R. D. (1995). "The duration required to identify the instru<sup>547</sup> ment, the octave, or the pitch chroma of a musical note," Music Perception 13(1), 1–15.
- 548 Saldanha, E., and Corso, J. F. (1964). "Timbre cues and the identification of musical
- instruments," The Journal of the Acoustical Society of America **36**(11), 2021–2026.
- Schaeffer, P. (2017). Treatise on Musical Objects: An Essay Across Disciplines, 20 (Univ
   of California Press).
- <sup>552</sup> Serra, X., and Smith, J. O. (**1990**). "Spectral modeling synthesis: A sound analy-<sup>553</sup> sis/synthesis system based on a deterministic plus stochastic decomposition," Computer

- 554 Music Journal 14(4), 12–24.
- <sup>555</sup> Siedenburg, K., and Doclo, S. (2017). "Iterative structured shrinkage algorithms for station-
- ary/transient audio separation," in Proc. of the 20th International Conference on Digital
   Audio Effects (DAFX-20), Edinburgh, Sep 5–8.
- <sup>558</sup> Siedenburg, K., and Dörfler, M. (**2011**). "Structured sparsity for audio signals," in *Proceed-*<sup>559</sup> ings of the 14th Int. Conference on Digital Audio Effects (DAFx-11), Paris.
- <sup>560</sup> Siedenburg, K., Jones-Mollerup, K., and McAdams, S. (2016). "Acoustic and categorical
- dissimilarity of musical timbre: Evidence from asymmetries between acoustic and chimeric
- sounds," Frontiers in Psychology **6**(1977), doi: 10.3389/fpsyg.2015.01977.
- Siedenburg, K., and McAdams, S. (2017). "Four distinctions for the auditory "wastebasket"
  of timbre," Frontiers in Psychology 8, 1747.
- <sup>565</sup> Suied, C., Agus, T. R., Thorpe, S. J., Mesgarani, N., and Pressnitzer, D. (2014). "Auditory
- gist: Recognition of very short sounds from timbre cues," Journal of the Acoustical Society
  of America 135(3), 1380–1391.
- Thoret, E., Depalle, P., and McAdams, S. (**2016**). "Perceptually salient spectrotemporal modulations for recognition of sustained musical instruments," Journal of the Acoustical Society of America **140**(6), EL478.
- Thoret, E., Depalle, P., and McAdams, S. (2017). "Perceptually salient regions of the
  modulation power spectrum for musical instrument identification," Frontiers in Psychology
  8, 587.
- <sup>574</sup> Zaunschirm, M., Reiss, J. D., and Klapuri, A. (**2012**). "A sub-band approach to modification
- of musical transients," Computer Music Journal 36(2), 23–36.