

¹ ARTICLE TEMPLATE

² **How to sketch timbre: investigating sound-shape associations in**
³ **free-form graphical representations of sound.**

⁴ Sebastian Löbbers and György Fazekas

⁵ Centre for Digital Music, Queen Mary University of London, United Kingdom

⁶ ARTICLE HISTORY

⁷ Compiled November 30, 2023

⁸ Wordcount: 10596

⁹ ABSTRACT

¹⁰ This paper investigates the influence of cross-modal associations on visual representations of sound. Compared to established methods that ask study participants to match existing stimuli, this research explores how people represent sound through free-form graphical sketches when focusing on musical timbre. A total of 2320 sound-sketches were collected in two studies that included 28 and 88 participants. High-level sketch categories were established through qualitative analysis of participant interviews and a card-sorting exercise. Inter-participant agreement on representations and correlations between the auditory and visual domains was computed through statistical analyses of quantitative sound and sketch features. The results show that while sound-shape associations play a significant role in sketched representations, humans incorporate other visual aspects like structural complexity or texture or choose figurative representations of emotions or sound-producing objects. Some level of agreement on how to represent sounds could be found between participants, which appears strongest for sounds dominated by a single perceptual attribute. The analysis further suggests that sound-shape associations found in sound-sketches align with findings from perceptual matching tasks. This research is motivated by designing novel perceptually-informed mappings for digital music production and the results presented in this paper lay the groundwork for the development of a sketch-based sound synthesiser.

³⁰ KEYWORDS

³¹ sound-shape associations, sound sketching, cross-modal mapping, timbre
³² perception

³³ 1. Introduction

³⁴ Humans make sense of sound in various ways often connecting different sensory domains. Cross-modal associations describe how a stimulus from one modality can induce a response in another modality. While often reported between sounds and colour, they can occur across various modalities, for example between colours and odours, sounds and tastes or sound and shapes (Spence, 2011). Cross-modal associations are sometimes wrongly referred to as synaesthesia. Synaesthesia is a rare condition with estimates of prevalence in the population ranging from 5% (Cuskley et al., 2019) to 0.5% Ramachandran and Hubbard (2001) and 0.05% Baron-Cohen et al. (1996).

42 Synaesthetes experience cross-modal connections involuntarily and consistently: the
43 same stimulus always induces the same response. Cross-modal associations on the
44 other hand are experienced in some form by most people, but connections tend to
45 be far less consistent and might only occur situationally. Despite this difference, in
46 a study investigating colour-to-sound mappings control participants and synaesthetes
47 employed the same heuristics for linking auditory and visual domains, such as con-
48 necting pitch with lightness (Ward et al., 2006). The authors conclude that this type
49 of synesthesia involves utilising mechanisms similar to those in typical cross-modal
50 perception.

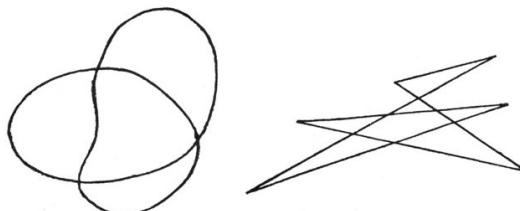


Figure 1.: Visual stimuli used in cross-modality experiments (Köhler, 1929). The left shape is overwhelmingly associated with the made-up word *maluma* or *boubou* and the right one with *takete* or *kiki*.

51 One of the earliest examples of cross-modal research is provided by Wolfgang
52 Köhler, a member of the *Gestaltpsychology* movement in the 1920s, who found
53 that people associate the made-up words *takete* or *kiki* with sharp, jagged shapes
54 and *maluma* or *boubou* with soft, round shapes (Köhler, 1929). The effect was
55 confirmed in multiple studies (Ramachandran & Hubbard, 2001) and generalised to
56 all phonemes (Nielsen & Rendall, 2013). It was observed across cultures (Davis, 1961;
57 Taylor & Taylor, 1962; Bremner et al., 2013), age groups including toddlers (Maurer
58 et al., 2006), to some extent, with the visually impaired (Bottini et al., 2019) and
59 between movement and phonemes (Shinohara et al., 2016). Similar associations were
60 not only found for phonetic sounds but also for musical instruments (Adeli et al.,
61 2014; Gurman et al., 2021) and abstract sonic textures (Grill & Flexer, 2012). While
62 these studies explicitly ask participants to match stimuli from different modalities,
63 a different approach measures how a stimulus in one domain can influence the
64 perception of another stimulus in a different domain. A famous example is the
65 McGurk effect (MacDonald & McGurk, 1978), a multisensory illusion that occurs
66 when merging video and audio recordings of vocalising different consonants. Other
67 examples include the impact of a mug's colour on the taste of coffee (Van Doorn
68 et al., 2014) and the effect of piano music on sexual attractiveness (Marin et al., 2017).

69
70 The work presented in this paper focuses on sound-shape associations, a subset
71 of cross-modality research, that describe how people connect timbre with geometrical
72 shapes (Sidhu & Pexman, 2018). Cross-modal associations are typically investigated
73 by asking participants to match visual and sound stimuli. However as pointed out
74 by Soraghan et al. (2018), this approach can show which of the presented stimuli
75 participants are most likely to match and does not account for whether a participant
76 might favour a different representation altogether. This research asked participants in

77 two studies to create monochromatic sketch representations of their associations with
78 sound. The first study uses an exploratory design that imposes minimal restrictions
79 to gain an understanding of the variety in which participants will represent different
80 types of sound graphically. The second study uses a more controlled design that
81 focuses on simple representations of sound produced by a frequency modulation
82 (FM) synthesiser. Alongside sound-sketches, qualitative feedback was collected
83 from participants through interviews and surveys to find out how the tasks were
84 approached. The analysis categorises the various representational approaches and
85 tests for correlations between sounds and shapes through statistical analyses of
86 quantitative audio and visual features. The results lay the groundwork for the broader
87 context of this research that aims to determine if a graphical sketch input can be
88 used to help improve interaction with digital music production software, specifically
89 for controlling digital synthesisers. The design process of the second study included
90 the development of the digital sketching interface from a simple, generic sketchpad
91 to a specialised sketching interface that could be used for controlling a sketch-based
92 synthesiser.

93

94 This paper first introduces relevant research into visual representation of sound
95 with a focus on sound-shape associations, musical timbre and sketch recognition
96 in Section 2. Methods and material for both studies are described in Section 3.
97 The results for both studies are presented in parallel in multiple sections: analysis
98 of participant feedback in Section 4, sound-sketch categorisation in Section 5 and
99 statistical feature analysis in Section 6. Discussion and conclusion can be found in
100 Sections 7 and 8. Acknowledgement of the research funding body and a disclosure
101 statement including reference to ethical approval for all studies are found in Sections 9
102 and 10. A detailed presentation of extracted audio and visual features is provided in
103 Appendix A.

104 **2. Background**

105 This section first gives an overview of cross-modal representations of sound with a focus
106 on sound-shape associations that build the basis for this research. It then continues to
107 introduce research into timbre and sketch recognition that is relevant for the analysis
108 of the user studies.

109 ***2.1. Cross-modality in visual representations of sound***

110 Cross-modality does not only play a role in human perception of the world but also
111 in executing certain actions. Thoret et al. (2016) found that participants drew circles
112 with a more elliptical skew when listening to sounds that evoked elliptical kinematics.
113 Salgado-Montejo et al. (2016) showed the influence of pitch on the location of free
114 hand movement. They also found movement to be more jagged for higher pitches and
115 rounder for lower pitches. Both examples show that consistencies can be found be-
116 tween people not just in matching or rating tasks but also when given free agency over
117 their response. However, participants' actions were not visualised posing the question
118 of the influence of an action's graphical representation on the cross-modal response.
119 In visual art, connections between the auditory and visual domains are a recurring
120 theme. Notably, Russian artist Wassily Kandinsky developed multiple hypotheses on
121 associations between shapes, colours and music leading to a number of cross-modal

artworks inspired by pieces of composer Arnold Schönberg (Rucsanda et al., 2019). Consistent associations were found between complex stimuli of visual artworks and musical compositions (Albertazzi et al., 2015) and piano music excerpts and simple visual structures (Clemente et al., 2020). M. Küssner (2014) investigated how participants represent pure tones varied in pitch, loudness and tempo through drawing and produced similar findings about the connection of space and pitch as Salgado-Montejo et al. (2016). While timbre is acknowledged to play a role in cross-modal associations, the aforementioned research focuses on the musical context of the sound stimuli. Focusing on music timbre and visual texture, Giannakis and Smith (2000) and Giannakis (2006) provide an overview of audio-visual mappings for computer applications. They critically point out the difference between physical and perceptual representations of sound and arbitrary and sensory mappings. Grill and Flexer (2012) showed how sensory mappings, which are based on cross-modal associations, can help participants retrieve sound samples outside of a specific musical context. Knees and Andersen (2016) further developed this idea by proposing sound retrieval from a graphical sketch input and built a non-functioning prototype of such a system. Compared to Küssner's research, timbre plays a more dominant role here which participants tend to visualise through shapes, contours or textures - described by the researchers as *symbolic* representations. Recent research further investigated how participants represent different timbres in a more controlled way and showed how a functioning sketch-based sound retrieval pipeline could be implemented with the help of deep learning (Engeln & Groh, 2020; Engeln et al., 2021). The results showed that participants chose abstract, geometrical shapes as well as more complex images to represent sound. Work by the authors produced similar results (Löbbers et al., 2021) and further suggests that, to some level, free-form sketches consist of universally recognised patterns that allow participants to extract information about a sound's characteristic (Löbbers & Fazekas, 2022). The research presented in this paper takes a closer look at how to classify different representational approaches and collects a comprehensive sound-sketch dataset.

2.2. Description of musical timbre

Musical timbre is generally described as the qualitative aspects of sound that cannot easily be quantified like loudness or pitch; however, it remains a concept that lacks a unified definition. An often-cited definition from the American National Standards Institute (ANSI) states that timbre is the auditory attribute that enables listeners to distinguish two sounds with the same loudness and pitch (ANSI, 1994). This is often criticised for only describing what timbre is not, but not defining what it actually is (Siedenburg & McAdams, 2017). While humans typically have an intuitive understanding of what constitutes timbre, the language they use to describe it varies individually (Saitis et al., 2020) often borrowing concepts of other sensory domains (Saitis & Weinzierl, 2019). Despite this lack of a common vocabulary, Wallmark and Kendall (2021) argue that some concepts like luminance (bright), texture (rough), and mass (heavy) widely appear across language groups. In the realm of contemporary music production, where digital technologies play a pivotal role, musical timbre occupies a prominent role for many modern music styles that distinguish themselves through their 'specific' sound rather than harmony, melody or musical structure (Provenzano, 2018; Blake, 2012). Finding or crafting a desired sound becomes an integral part of the production process which often involves searching large sample libraries or tweak-

169 ing software parameters. These parameters often relate to the underlying digital signal
170 processing (DSP) rather than perception of sound which can make it difficult to realise
171 sound ideas or explore sonic spaces in an intuitive way (Seago, 2013). The perceptual
172 attributes used to process a sound can be influenced by the source itself or individ-
173 ual factors of the listeners like preference, domain knowledge, cultural upbringing or
174 listening context. This becomes apparent in the descriptors of sound libraries and syn-
175 thesiser presets in digital audio production that can range from describing an acoustic
176 instrument (keys, strings, drums), a compositional function (lead, pad, bass), a play-
177 ing style (staccato, legato, arpeggio), references to genre (rock, hip-hop), hardware
178 (808, clavinet), jargon (wobble bass, acid), mood (happy, sad) or describing a sound's
179 timbre directly often with the help of cross-modal associations (bright, rough, shrill).
180 Without tagging guidelines or a unified approach to describing sound, retrieving spe-
181 cific sounds from sample or preset libraries can be a cumbersome task that impairs
182 the creative process. An increasingly popular approach for the organisation of these
183 libraries is to display samples in a 2-dimensional space where similar samples are dis-
184 played close to each other (Fried et al., 2014; Bruford et al., 2019; Garber et al., 2021).
185 The groundwork for this approach was laid by Grey (1977) and further popularised by
186 Iverson and Krumhansl (1993) and McAdams et al. (1995), who obtained dissimilarity
187 ratings through participant studies and using multi-dimensional scaling (MDS) to low-
188 dimensional timbre space. For applications that are designed to work with a user's own
189 sample library like the *Timbral Explorer*¹ it is not feasible to use samples annotated by
190 humans. Instead, they use computationally extracted audio features. A popular and
191 successful feature is the mel-frequency cepstral coefficient (MFCC) that can be used as
192 input for deep-learning classification of diverse environmental sounds (Cramer et al.,
193 2019). However, MFCCs are not very good at communicating sound characteristics to
194 humans or helping explain the perception of sound (Siedenburg et al., 2016). Other
195 popular computational features that, while still based in acoustic feature analysis,
196 show better alignment with human perception are *spectral flatness* and *zero-crossing*
197 *rate* that describe a sound's noisiness, *spectral centroid* that serves as a measure for
198 brightness and the *root-mean-square* (RMS) describing a sound's loudness (Peeters,
199 2004; McFee et al., 2015). Pearce et al. (2019) aim to bridge the disconnect between
200 computational features and features relevant to human perception by developing a
201 model that, through a mixture of human-annotated sound samples and computed fea-
202 tures, predicts a sound's *hardness*, *depth*, *brightness*, *roughness*, *warmth*, *sharpness*,
203 *booming* and *reverberation*. Further research explores how, prompted by adjective,
204 participants synthesise novel timbres through frequency modulation (FM) (Wallmark
205 et al., 2019). Similarly, Hayes et al. (2022a) implemented a simple FM synthesiser to
206 run in a web browser and collected a dataset of annotated synthesiser sounds through
207 an online participant study. Participants were presented with a prompt to change a
208 reference sound and then asked to annotate the result. The prompts were presented in
209 the form *make the sound less/more: bright/rough/thick*. The attributes were derived
210 from the luminance-texture-mass (LTM) model (Zacharakis & Pasiadis, 2015, 2016)
211 that suggests that timbre descriptions can sufficiently be explained by these three
212 dimensions for sounds with similar amplitude envelopes.

¹<https://www.audiocommons.org/2019/01/30/timbral-explorer.html>

213 **2.3. Sketch recognition**

214 To statistically evaluate cross-modal associations between sound and shapes both do-
215 mains have to be described as quantitative features. The extraction of information from
216 hand-drawn sketches by a computer is called sketch recognition. Typical applications
217 include recognition of handwriting or image retrieval from a sketch input. Sketches
218 can be analysed as rasterised images using a computer vision approach. A benchmark
219 task is the classification of handwritten digits with the MNIST dataset that is typi-
220 cally achieved through machine learning with convolutional neural networks (CNNs)
221 producing the best results (LeCun, 1998; Deng, 2012). Representing sketches sequen-
222 tially either through vectorisation or collecting sketches digitally gives the opportunity
223 for a wider range of analyses. The most notable example is seminal work by Ha and
224 Eck (2017) that introduced the *SketchRNN* architecture that uses a variational au-
225 toencoder (VAE) built on recurrent neural networks (RNNs). *SketchRNN* was trained
226 on *Quick, Draw!*², a large-scale, open-source dataset with over 50 million sketches di-
227 vided into categories ranging from simple geometric shapes to complex representations
228 of animals and objects. *Quick, Draw!* inspired a large number of projects by enabling
229 researchers to experiment with pre-train models and (re-)train them for specific tasks.

230 While *SketchRNN* produces impressive results for sketch classification and gener-
231 ation, algorithmic approaches might be more suitable for describing the shape of a
232 sketch. The *ShortStraw* algorithm (Wolin et al., 2008) provides a simple, effective
233 tool to extract corner points. Xiong and LaViola Jr (2009) extended the algorithm to
234 also recognise curve points. Sezgin (2001) further show that information can not only
235 be extracted from a sketch's shape but also from the sketching speed, for example,
236 corner points can be estimated from a decrease in speed before each point. In the con-
237 text of sound-shape associations, sketch or movement responses are often interpreted
238 qualitatively by humans rather than computationally as it can be seen in the works
239 by Salgado-Montejo et al. (2016) and Knees and Andersen (2016) discussed in Sec-
240 tion 2.1. M. Küssner (2014) proposes a computational approach to extracting features
241 from sound-sketches that, however, expect a representation inside a grid where the
242 x-axis represents time. Engeln et al. (2021) harness advancements in deep learning for
243 sketch recognition to implement a computational method for free-form sketch-based
244 sound retrieval. However, an end-to-end retrieval approach might not reveal which
245 aspects of a sketch are relevant to sound-shape associations.

246 **3. Methods and Material**

247 In order to investigate sound-shape associations in free-form sketch representation,
248 two studies were conducted resulting in two sound-sketch datasets and correspond-
249 ing qualitative and quantitative analyses. In both studies, participants were asked to
250 sketch their associations with sound stimuli using a digital interface. For both stud-
251 ies, quantitative features were extracted from the sound stimuli and sound-sketches
252 to investigate sound-shape connections through statistical correlations. Study 1 was
253 accompanied by qualitative interviews to gain a deeper understanding of how partici-
254 pants chose their representations. This research further provides high-level categories
255 for the broad classification of sound-sketches that were achieved through a card-sorting
256 exercise in Study 1 and an automated, machine-learning approach in Study 2.

²<https://quickdraw.withgoogle.com/>

257 **3.1. Design of perceptual studies**

258 As discussed in Section 2.1, relatively little research has been conducted on how hu-
259 mans represent sounds through graphical sketching. The two studies presented in this
260 paper collect sound-sketches by combining existing research that asks participants to
261 match auditory and visual stimuli or undertake specific hand movements while lis-
262 tening to sound with methods of digital sketch collection introduced in Section 2.3.
263 In both studies, participants are presented with a number of different sound stimuli
264 and are asked to sketch their personal association with them. The first study follows
265 an exploratory design with minimal instructions for participants and different cate-
266 gories of sound from abstract synthesiser pads to acoustic instruments to obtain a wide
267 overview of different representational approaches. Informed by the results of Study 1,
268 the second study followed a more controlled design with a further developed interface
269 and a narrower set of synthesised sound stimuli that encouraged simpler, more ab-
270 stract sketch representations with a stronger link to sound-shape associations. Besides
271 collecting sound-sketches, it is important for this research to understand the reason-
272 ing behind different approaches. Qualitative feedback can give a more detailed insight
273 into the thought processes of a participant and might yield information that cannot
274 be captured through quantitative data. For Study 1, qualitative data was collected
275 through behavioural observation and a semi-structured interview that probed partic-
276 ipants to reflect on their approaches and the interaction with the digital interface.
277 Study 2 gave the option to provide brief written feedback and asked participants to
278 answer a modified System Usability Scale (SUS) questionnaire (Lewis, 2018). All ver-
279 bal feedback was analysed through thematic analysis (Braun & Clarke, 2006; Stuckey,
280 2015). General demographic data including age, gender, occupation and country of
281 origin was collected through survey questions. In addition, the section of the Gold-
282 smiths Music Sophistication Index (Gold MSI) (Müllensiefen et al., 2014) relating to
283 musical training and engagement with music was used to categorise participants by
284 music proficiency.

285 **3.2. Material**

286 Both studies are set entirely in a digital environment and sound stimuli are selected
287 with regard to timbral differences. Participants for Study 1 completed the study on
288 the same laptop and, to keep similar conditions, only laptop or desktop devices were
289 allowed for Study 2.

290 **3.2.1. Study 1**

291 The study was conducted in person using the touchpad on a 15" MacBook Pro and
292 a pair of Beyerdynamic DT 770 headphones in calm, indoor locations. Ten stimuli
293 with distinct timbral characteristics were crafted in reference to two cross-modal ex-
294 periments presented in Section 1, which examine musical sounds (Adeli et al., 2014)
295 and synthesized sounds and textures (Grill & Flexer, 2012). This approach aims to
296 broaden the investigation of cross-modal representations beyond more narrowly de-
297 fined categories to a wider range that could be encountered in a digital music creation
298 environment. The stimuli can be divided into harmonic sounds that include musical
299 instruments (*Piano*, *Strings*, *Electric Guitar*) and synthesised pads (*Telephonic*, *Sub-*
300 *bass*) and inharmonic sounds that include environmental sounds (*Impact*) and abstract
301 textures (*Noise*, *String Grains*, *Crackles*, *Processed Guitar*). *Piano* and *Strings* were

302 created with virtual instruments in the Kontakt 5 plugin by Native Instruments, *Elec-*
303 *tric Guitar* was recorded with an audio interface via line-in. *Impact* was taken from the
304 online audio repository freesound.org and *Telephonic*, *Subbass*, *Noise*, *String Grains*,
305 *Crackles* and *Processed Guitar* were designed with built-in plugins in the digital audio
306 workstation Ableton Live 10. Harmonic sounds were pitched to the MIDI note C3
307 which lies within the frequency range used in the reference experiments. All sound
308 stimuli are monophonic and normalised for equal loudness. They last eight seconds
309 with varying amplitude envelopes and include trailing silence to mark a clear endpoint
310 during looped playback. The perceived base frequency may vary due to prominent har-
311 monics.³

312 **3.2.2. Study 2**

313 The study was conducted entirely online. An initial audio check as well as a mid-study
314 attention test ensured that participants listened to the audio playback either through
315 headphones or speakers. Twenty different FM synthesiser sounds were selected from
316 the dataset by Hayes and Saitis (2020) discussed in Section 2.2. All sounds have a
317 short attack and high sustain to ensure that participants focus on timbre rather than
318 amplitude envelope. The FM synthesiser was implemented with the AudioWorklets
319 interface of the Web Audio API to synthesise sound directly in the browser rather than
320 using audio samples which enabled continuous playback of a sound without looping.
321 All sounds were pitched to the MIDI note A3 and loudness-normalised.⁴

322 **3.3. Interfaces**

323 For both studies, digital sketching interfaces were used that can run in a web browser.
324 This makes it possible to collect sketches directly as sequential data as presented in
325 Section 3.7 and collect data online. As this research investigates sound-shape asso-
326 ciations, the interfaces only allow monochromatic sketches with fixed stroke widths.
327 Sketch analysis of Study 1, described in the following Sections, suggests that sound-
328 shape associations more strongly inform simple sketches. To investigate that connec-
329 tion in greater detail, the interface was developed from a generic to a specialised design
330 through an iterative design process that sought to encourage simple sketches between
331 participants without majorly impacting their perceived sense of expressiveness and
332 autonomy.

333 **3.3.1. Study 1**

334 The requirements for this interface, that is illustrated in Figure 2, were that par-
335 ticipants can only express their ideas through monochromatic strokes, but are not
336 otherwise limited or supported. The black canvas was chosen to clearly distinguish
337 the sketching area from the rest of the website. No option to partially or fully erase
338 a sketch was provided so that participants would not revise their original idea. While
339 the goal of Study 1 was to gain insights into the different sketched representations of
340 sound, it also provided feedback on how to develop the interface for sound-sketching
341 tasks. Feedback evaluation showed that, overall, participants were unsatisfied with the

³All sounds for Study 1 can be accessed online at <https://bit.ly/3ta6crU> together with the sketches pro-
duced by participants.

⁴All sounds for Study 2 can be accessed online at https://sfrl.github.io/study2_gallery/ together with
the sketches produced by participants.



Figure 2.: Interface for Study 1: generic design that allows for fixed-width, white strokes on a black canvas with a fixed size of 750x750 pixels. There are no sketch length limits, stroke simplification and undo or reset options. Sketch points are connected by straight lines.

342 interface, especially in regards to sketching straight lines or smooth transitions. In ad-
 343 dition, the canvas proved to be too small with multiple participants sketching over the
 344 edges. The black background received positive comments, but for some participants, it
 345 seemed to have artistic rather than purely functional meaning. While the design was
 346 successful in forcing participants to stay with their original idea, it did not take into
 347 consideration that participants might want to correct technical mistakes rather than
 348 change their overall approach.

349 3.3.2. Study 2

350 The feedback from Study 1 was addressed in the further interface design. Stroke simpli-
 351 fication and spline interpolation were implemented to make it easier to sketch straight
 352 lines and smooth transitions. The canvas automatically stretches to the size of the
 353 browser window to prevent sketches from extending over the edges. A reset button
 354 ensures that participants can start over, but an erase or undo function was withheld
 355 to prevent participants from focusing on retouching their representations. These de-
 356 sign choices are similar to the *Quick, Draw!* interface introduced in Section 2.3 which
 357 makes collected sketches compatible with their dataset and deep-learning architecture.
 358 A major difference between the two interfaces is the length limit introduced in Study
 359 2. The goal of this research is to collect sketches that encode sound characteristics
 360 through shape which typically results in simple, abstract representations. However,
 361 the categorisation of Study 1 sketches, as described in detail in Section 5, show a sig-
 362 nificant number of figurative representations like scenes or objects. As these sketches
 363 are on average longer, limiting the length permitted by the interface was expected
 364 to reduce those types of representations. This was tested in three small-scale design
 365 studies with 10-15 participants each. Different designs were evaluated on how well
 366 they guided participants towards simple, abstract representations while maintaining
 367 the feeling that the interface allows them to be expressive. This concluded in the setup



MM R

—

Figure 3.: Interface for Study 1: Specialised design that allows for fixed-width, black strokes on a white canvas that scales to a participant’s browser window (a minimal size of 800x600 pixels was enforced). The sketch length is limited to the range of 30 to 150 points. A sketch can only be submitted if it exceeds the lower length limit. If a participant continues sketching after reaching the upper limit, sketch points are erased from the start. A small meter in the top left corner of the canvas indicates the current stroke length. Participants can reset the canvas to start over. The Ramer-Douglas-Peucker algorithm (Douglas & Peucker, 1973) was used for stroke simplification and stroke points were connected with Catmul-Rom splines (Catmull & Rom, 1974).

368 described in Figure 3.

369 *3.4. Participants*

370 Recruitment was open to all adults over the age of 18. The aim of these studies was
 371 to find out about sound-shape associations by sampling from the general population.
 372 As engagement with music was expected to have an influence on representations, for
 373 Study 1 musicians and non-musicians were recruited in equal parts.

374 *3.4.1. Study 1*

375 Twenty-eight participants were recruited through mailing lists and in person at the
 376 School of Electronic Engineering and Computer Science at Queen Mary University
 377 of London. This group was divided equally by gender (14 female, 14 male), 25 were
 378 adults below the age of 34 (three between 34 and 49), 22 had a Western background
 379 (16 from Europe, 4 from North America, 2 from South America) and 5 an Eastern
 380 background (4 from China, 1 from India) with one participant preferring not to disclose
 381 this information. A participant was defined as a musician if they responded to having
 382 at least one year of formal music education and engaged in musical activity (playing
 383 an instrument, producing or composing music) at least once a month. This resulted
 384 in 14 musicians and 14 non-musicians.

385 *3.4.2. Study 2*

386 Eighty-eight participants were recruited through the online platform Prolific.⁵ The
 387 majority of participants were female (48 female, 38 male, 2 other) and 84 were aged
 388 between 18 and 33 ($M = 22.5$, $SD = 4.6$). All participants had a Western background
 389 (56 from North America, 16 from Europe, 13 from South Africa and 3 from South
 390 America) with the majority coming from Mexico (54 participants). Sixty-six were

⁵<https://prolific.co/>

391 students and 28 stated to have had at least one year of formal music education and
392 actively engage in musical activity at least once a month.

393 **3.5. Procedure**

394 A similar procedure was used for both studies. As Study 2 was conducted online a
395 more robust digital environment was deployed to ensure that participants completed
396 the study successfully.

397 **3.5.1. Study 1**

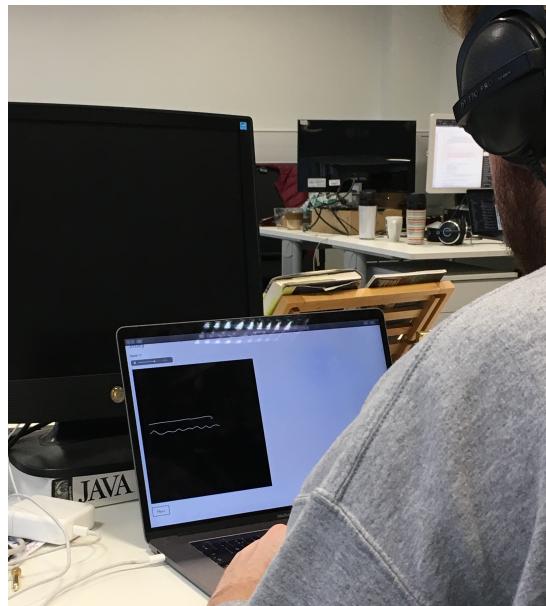


Figure 4.: Participant sketching a sound in Study 1

398 Participants were first asked to adjust playback audio to a comfortable measure
399 using white noise as a reference. They then completed a questionnaire collecting de-
400 mographic data and information about their experience with music. Participants were
401 then asked to familiarise themselves with the sketching interface without audio before
402 they were presented with the sound stimuli. The study intended to encourage a spon-
403 taneous response, therefore no information about the range of sounds was provided
404 and participants were instructed to sketch what they believed to best represent each
405 sound stimulus. Looped playback started automatically with the option to pause and
406 resume. Each sound was played twice in a randomised order resulting in a total of
407 twenty sketches per participant. After completion, a semi-structured interview was
408 conducted asking participants how they approached the task and whether they found
409 it difficult. No time limit was given and the study typically took twenty to thirty
410 minutes to complete.⁶

⁶The setup for Study 1 can be accessed online at <https://bit.ly/3j3FkV0>.

411 **3.5.2. Study 2**

412 Participants were first presented with a set of information ensuring that they use a
413 laptop or desktop device and are able to listen to sound either through headphones
414 or loudspeakers. This was followed by a short introduction of the study that included
415 guidelines on representing sounds in an abstract rather than figurative way, a short
416 explanation of the sketching interface, a guide to adjust playback volume to a com-
417 fortble level and a check that the browser window size was at least 800×600 pixels.
418 Before starting the main task, participants were given the opportunity to familiarise
419 themselves with the interface in two test rounds in which they were asked to sketch
420 their associations with an imagined calm and noisy sound. For the main task, sound
421 stimuli were played back automatically in a randomised order with the instruction *lis-*
422 *ten to the sound and draw your association.* An attention test was played back after
423 completing half of the task asking participants to sketch the number four instead of a
424 sound. The study concluded with a survey asking participants about their experience
425 with the study and the interface and collecting demographic data, experience with
426 music and the hardware they used to complete the study.⁷

427 **3.6. Sketch categorisation**

428 The first step of the analysis categorises different representational approaches that
429 participants deploy and investigates to what extent sound-shape associations inform
430 them. The sketches produced in Study 1 were categorised in an open card-sorting
431 study by six participants (4 female, 3 musicians) who did not take part in the main
432 study. They were asked to sort the collected sketches into three to ten categories
433 including short written descriptions. The study was completed remotely within three
434 hours on participants' devices following detailed instructions that can be accessed
435 online.⁸ Following the methodology of Paea and Baird (2018), the results were
436 encoded and reduced in dimensionality using principal component analysis (PCA).
437 K-means clustering was used together with the silhouette coefficient (Rousseeuw,
438 1987), a measure of cluster goodness, to find the most suitable number of clusters
439 between three and ten. Clusters were named and described qualitatively based on
440 keywords that participants used in their category descriptions.

441

442 Due to the higher participant number, Study 2 produced a considerably larger
443 number of sound-sketches and the categorisation process was automated with the
444 help of machine learning. A variational autoencoder (VAE) using the *SketchRNN*
445 architecture was pre-trained on the *Quick, Draw!* categories *Triangle*, *Square*, *Circle*,
446 *Line*, *Squiggle* and *Zigzag* before feeding in sound-sketches from Study 2. The
447 resulting 128-dimensional latent representation of the dataset was reduced with PCA
448 and the best number of clusters was determined with K-means and the silhouette
449 coefficient as described above. Clusters were named following the findings from Study
450 1.

451

452 Pearson's Chi-squared test was used to test whether sounds are represented
453 equally across sketch categories or if certain types of sound can be connected to a spe-
454 cific category. Similarly, Cochran's Q test was used to assess whether a participant's
455 musical proficiency has an influence on how their sketches are categorised.

⁷The setup for Study 2 can be accessed online at <https://sketching-sounds.web.app/>.

⁸<https://youtu.be/LXTlnaAciWw>

456 **3.7. Quantitative analysis of sketch and sound features**

All sketches are saved digitally as sequential data in nested arrays. Each stroke is described in a separate array that consists of sketch points described by their x and y positions and timestamps. A sketch is rendered into an image by connecting all points within each stroke array. To describe sketches quantitatively, a number of features can be calculated directly from the data structure and through simple arithmetic operations as demonstrated in Equations 1, 2 and 3, where N is the number of strokes in a sketch and \bar{L} , \bar{T} and \bar{S} are their average length, completion time and sketching speed. The number of points in the k^{th} stroke is described by n_k . Each point has a position x_{k_i} and timestamp t_{k_i} . The Euclidean distance between two points is described by $d(p, q)$.

$$\bar{L} = \frac{1}{N} \sum_{k=1}^N \sum_{i=2}^{n_k} d(x_{k_i}, x_{k_{i-1}}) \quad (1)$$

$$\bar{T} = \frac{1}{N} \sum_{k=1}^N t_{k_{n_k}} - t_{k_1} \quad (2)$$

$$\bar{S} = \frac{1}{N} \sum_{k=1}^N \sum_{i=2}^{n_k} d(x_{k_i}, x_{k_{i-1}}) \frac{1}{t_{k_{n_k}} - t_{k_1}} \quad (3)$$

457 Sound-shape associations are usually reported with respect to a shape's contour
 458 focusing on their 'jaggedness' or 'roundness' (Adeli et al., 2014; Grill & Flexer, 2012).
 459 As illustrated in Figure 5, these attributes were quantified by extracting corner points
 460 divided into obtuse, right and acute angles and curve points divided into wide and
 461 narrow shape algorithm (Wolin et al., 2008; Xiong & LaViola Jr, 2009). A qualitative
 462 review suggested that sketches differ by the number of stroke intersections that can be
 463 interpreted as the 'noisiness' of a sketch. The number of intersections was determined
 464 using an adaptation of Bresenham's rasterisation algorithm (Bresenham, 1965). Prior
 465 to extracting features, the sketch data was cleaned by removing consecutive points
 466 with the same position and merging two strokes if a starting point was within a five-
 467 pixel distance to an endpoint. The number of intersections, corner and curve points is
 468 reported relative to the total stroke length of a sketch.

469 In order to investigate sound-shape associations through statistical analysis, the
 470 sound stimuli also have to be described using quantitative features. This was ac-
 471 complished by computing the mean values of *Centroid Frequency*, *Spectral Flatness*,
 472 *Zero Crossing* and *Root Mean Square Power (RMS)* for each sound using the *Librosa*
 473 Python library (McFee et al., 2022) with an FFT window size of 2048 and hop length
 474 of 512. In addition, the *timbral models* by Pearce et al. (2019) provided quantified
 475 measures of *Hardness*, *Depth*, *Brightness*, *Roughness*, *Warmth*, *Sharpness* and *Boomi-
 476 ness* which can more easily be related to human perception of sound. The additional
 477 feature *RMS Slope*, describing how continuous or intersected a sound is, was quantified
 478 by the slope between prominent extrema in the RMS envelope.

479 The sketch categories described in Section 5 provide a qualitative description of rep-
 480 resentational approaches. The quantitative set of features described in this section is
 481 used to investigate sound-shape associations statistically. Spearman's rank correlation
 482 coefficient is used to find out which, if any, visual features correlate with which audio
 483 features. As a linear perceptual relationship between features cannot be assumed and
 484 the underlying distribution is unknown, the non-parametric Spearman test was chosen

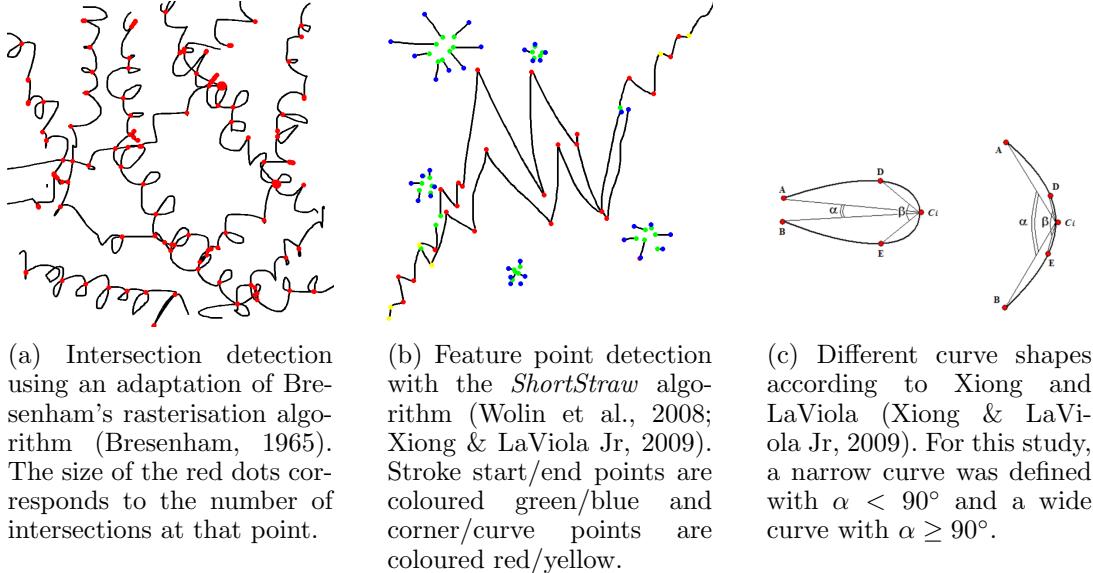


Figure 5.: Sketch feature extraction

as it can find not only linear but monotonic relationships in general. The correlation analysis uses mean values of visual features for each sound. To ensure that averaged features still pose meaningful sketch descriptions the inter-rater reliability was determined using the ICC(2,k) model intraclass correlation coefficient (ICC) (Koo & Li, 2016). The ICC(2,k) measures absolute agreement of average raters by averaging responses of k raters for each subject. In this context, sound stimuli were defined as subjects and sketch features as measurements. Sketch features were first log-transformed to meet the normal distribution assumption of the ICC. This approach was chosen because perceptual data typically includes large variance and patterns emerge more clearly when observing averages. This means that results will provide information about the average sound-sketch representation derived from multiple participants, but might not be applicable for predicting or describing sketches of an individual.

4. Qualitative feedback analysis

Qualitative feedback was given in a semi-structured interview in Study 1 which provided a broad picture of participants' approaches and laid the groundwork for the development of Study 2. A part of these results is illustrated in Figure 6. In Study 2, brief qualitative feedback was given in written form and supported by quantitative survey responses. This section first focuses on feedback about the task itself and then summarises feedback about the interaction with the interface.

4.1. Sketching task

Task difficulty was reported as easy by 50%, neutral by 21% and hard by 29% of participants in Study 1 which changed to 84%, 9% and 7% in Study 2. Despite the positive skew in Study 2, mixed responses were recorded to the question of whether it was easy to think about sound in a visual way with 50% agreeing, 27% disagree-

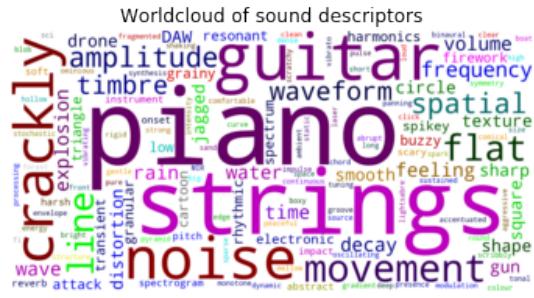


Figure 6.: Reoccurring concepts when describing the sounds in Study 1 extracted from participant interviews.

ing and 23% giving a neutral response. This is well captured in the response ‘I have never had to interpret sounds visually, therefore I found it to be kind of a difficult but interesting task.’ (P2.70). This was echoed in interviews from Study 1 where some participants found it difficult to ‘think of sound in a very visual way’ (P1.8) (P1.16). However, for some, it became easier to visualise a sound once the task started, as one participant stated ‘I wasn’t really expecting to be able to visualize sound, but some of those frequencies were extremely clear to me, as far as how they looked in my brain.’ (P2.53). Participants who felt that the task was easy thought that ‘there was no right or wrong’ (P1.4), ‘it was just about being creative’ (P1.15), they did not have to ‘achieve something’ (P1.10), the setup ‘allowed the listener space to interpret all sorts of sound visually’ (P2.91) or simply found the task ‘interesting’ (P2.39, P2.41, P2.70, P2.83, P2.91) and ‘fun’ (P2.20, P2.33, P2.45). For Study 2, some criticized that the ‘sounds were very alike’ P(2.13) which made it ‘[...] hard to find an specifically draw[ing] to each sound’ P(2.34). On the contrary, in Study 1 some participants struggled with the ‘great variety in the sounds’ (P1.2) which made it difficult to find a consistent approach. While some participants approached the task as an intuitive, creative activity, others were concerned with establishing a consistent visual language, difficulties arose while deciding which sound characteristics to follow because ‘there are too many things to consider’ like ‘brightness or aggressiveness or how it [timbre] develops over time’ (P1.6). Deploying a more systematic rather than intuitive approach appeared more difficult with one participant who did not deviate from their initial concept finding themselves ‘going round in circles, and question how valid the whole approach is’ (P1.9). Some participants reported that ‘complicated ones [sounds] sounded like pictures, and then the simple ones [...] like piano notes were a lot harder to draw’ (P1.8) possibly because they ‘hear [them] all the time’ (P1.1), while other participants thought that ‘it’s pretty straightforward because I know a piano note more than others’ (P1.5). Most participants approached the task by listening to ‘the actual sonic qualities of them [the sounds]’ (P1.1) and representing them with ‘[...]abstract patterns, along with patterns that would come from what I thought the instruments were, and it was kind of a mixture of going back and forth between the two.’ P(1.5). Familiar sounds like the piano can influence participants to choose figurative representations that include the sound-producing source. However, some participants adopted a figurative approach that does not pay attention to specific sound characteristics but rather extracts general information like an emotion from a sound and then depicts a scene or an object that fits this information. One participant explained their sketch which is shown in Figure 7 as follows: ‘I wanted to show the emotion in the sound like

545 the sound was scary so I drew the forest with only one person inside it. If I think the
546 sound is peaceful, I will draw the sea and the sun and the water.' (P1.21). These two
547 different approaches suggest that abstract representations are more strongly informed
548 by sound-shape associations and figurative representations draw more strongly from
549 an emotional response or personal memory.

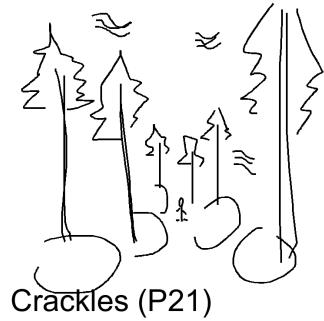


Figure 7.: A participant in Study 1 represented a sound that they perceived as 'scary' with a scene that represented that emotion.

550 **4.2. Sketching interface**

551 Analysing responses for feedback about the interface interaction showed that for Study
552 1, overall, participants were unsatisfied with the interface with twelve mentioning
553 that they would prefer a pen (either digital or analogue) to be able to sketch more
554 accurately. Six stated that they would have liked to utilise additional visual tools
555 like colour, different strokes and textures. However, responses suggested that while
556 the interface 'could be more expressive [...] for the purpose it was expressive enough'
557 (P1.5). The feedback led to the overhaul of the interface design for Study 2 described in
558 Section 3.3.2. For Study 2, 70 and 75 participants responded with agree or completely
559 agree to the questions *I thought the drawing interface allowed me to be expressive* and
560 *I thought the drawing interface was easy to use*. Three participants still mentioned
561 that they would like to add colours to the interface and one mentioned that the task
562 'would be easier to do on my phone' (P2.85). A further three participants mentioned
563 that they would like to increase or discard the stroke length limit.

564 **5. Sketch Categorisation**

565 From the interview analysis of Study 1 two broad representational approaches can be
566 defined: an abstract approach that is guided at least in part by sound-shape associa-
567 tions and a figurative approach that is strongly influenced by imagery. Sketch categori-
568 sation for Study 1 focuses on further formalising these findings and investigating how
569 participants could be guided towards abstract representations. Study 2 focuses on the
570 influence of prominent sound characteristics on abstract representational categories.

571 **5.1. Study 1: manual categorisation through card sorting**

572 As described in Section 3.6, sketches were categorised in an open card sorting study.
573 Analysis of the responses returned an optimal number of five categories that were
574 named: *Chaotic/Jagged* (172 sketches), *Radiating/Round* (126), *Lines* (120), *Object-*
575 *Scenes* (86) and *Grains* (56). The results are visualised in Figure 8. Descriptive
576 keywords and sketch examples for each category can be found in Table 1. A maxi-
577 mal silhouette coefficient of 0.49 suggests that categories are distinguishable, but not
578 clearly separated which is also reflected by occasionally overlapping keywords. The *Ob-*
579 *ject/Scenes* category consists mainly of figurative representations while the remaining
580 categories include mainly abstract representations. Chi-squared test suggests that non-
581 musicians produce *Objects/Scenes* sketches more often ($\chi^2(1,N=28)=22.51$ $p<.0001$)
582 while musicians produce *Lines* sketches at a higher rate ($\chi^2(1,N=28)=7.5$ $p<.01$)
583 possibly because this category contains sketches that appear to reference audio visu-
584 alisations like envelopes or waveforms. Category counts for *Objects/Scenes* sketches
585 significantly differ between sounds ($\chi^2(9)=67.07$ $p<.0001$) with post-hoc analysis re-
586 vealing that *Piano* and *Impact* show significantly higher counts than *Noise*, *String*
587 *Grains* and *Processed Guitar* ($p<.01$ for each pair). A possible explanation is that
588 *Piano* and *Impact* have an easily identifiable source that participants attempted to
589 sketch rather than capturing sound characteristics directly. *Noise* and *String Grains*
590 that show high values for the *roughness* audio feature (81 and 56) as displayed in Ta-
591 ble A1 also have the largest share of sketches in the *Chaotic/Jagged* category. On the
592 other hand, *Subbass*, *Telephonic* and *Impact* with high values for the *warmth* audio
593 feature (65, 54 and 51) have the highest share of the *Radiating/Round* category. Inter-
594 interestingly, despite high values for *warmth* (54), *Piano* and *Strings* are more frequently
595 represented with *Lines* sketches, possibly because they are comparably spectrally sim-
596 ple sounds which is reflected by the low values for spectral flatness. However, it could
597 also relate to auditory aspects such as pitch and loudness, as it is not feasible to eval-
598 uate timbre independently from these attributes. Table 2 shows that *Object/Scenes*
599 has the highest average number of sketch points and the second highest average num-
600 ber of strokes. The highest number of strokes can be found in the *Grains* category
601 which is most prevalent in the sounds *Crackles* and *String Grains*. However, with
602 a low average number of points, this category appears to represent the intersecte-
603 ness of sound quantified by the audio feature *RMS slope* with multiple short, simple
604 structures. *Chaotic/Jagged* and *Radiating* have a similar average number of points
605 as *Object/Scenes*, but a significantly lower number of average strokes. This analysis
606 indicates that *Object/Scenes* sketches are more likely to be long and complex with
607 multiple components which informed the interface development for Study 2 described
608 in Section 3.3.2.

609 **5.2. Study 2: automated categorisation using machine learning**

610 As described in Section 3.6 the *SketchRNN* deep learning architecture was used to
611 create a latent representation of the collected sound-sketches. Figure 10 shows a visu-
612 alisation of the latent space and categorisation through cluster analysis. In contrast to
613 Study 1, three clusters were determined to separate the data best, however, a silhouette
614 score of 0.36 shows that large overlaps exist between them. Three major categories
615 informed by the card sorting annotations from Study 1 were defined: *Lines* (610),
616 *Chaotic/Jagged* (455 sketches), and *Radiating/Round* (694). The categories *Obje-*
617 *cet/Scenes* and *Grains* were no longer present in the dataset. This was expected as the

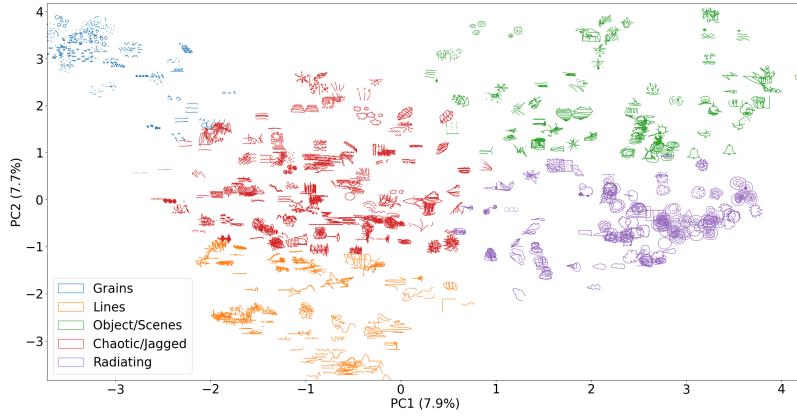


Figure 8.: Study 1 sound sketches organised by the K-Means clusters calculated from the card-sorting study. Colours indicate the different clusters.

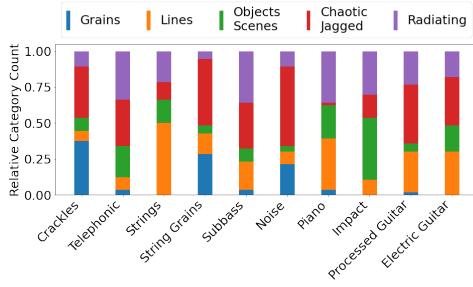


Figure 9.: Categories by sound stimulus in Study 1.

618 interface design discouraged *Objects/Scenes* sketches and *Grains* sketches were largely
 619 found in intersected sounds that were not included in the audio stimuli for Study 2.
 620 To compare differences in category distribution, six sound groups were created from
 621 sounds from the annotated attributes *bright*, *rough* and *thick* derived from the FM syn-
 622thesiser sound dataset by Hayes and Saitis (2020) discussed in Section 2.2. Each group
 623 contains three sounds with either the highest or lowest value for an attribute. Further
 624 analysis showed that category counts for *Chaotic/Jagged* and *Lines* differ significantly
 625 between sound groups ($\chi^2(5)=20.38$ $p<.01$ and $\chi^2(5)=34.29$ $p<.00001$), but no signif-
 626icant differences could be found for *Radiating/Round* ($\chi^2(5)=9.59$ $p>.05$). Post-hoc
 627 analysis showed prominent differences in category distribution between the most and
 628 least rough sounds as illustrated in Figure 11 supporting findings from Study 1 which
 629 showed that rough sounds were more frequently represented with chaotic, complex
 630 sketches and calmer, less rough sounds with simpler lines. Measurements for rough-
 631ness which were extracted automatically with *timbral model* by Pearce et al. (2019)
 632 and presented in Table A1 were high for the *rough* sound group (70, 63 and 59 for
 633 Synth 9,13,19) and low for the *not rough* sound group (29,0,45 for Synth 1,11,14).
 634 This further supports the validity of these features as measurements relevant to hu-
 635 man perception.

Grains	Lines	Object/ Scenes	Chaotic/ Jagged	Radiating/ Round
<i>Small, repeated, grainy, spots, multiple components, layers, abstract, distinct</i>	<i>round, soft, continuous, jagged, irregular, simple, single, lines</i>	<i>real-life objects, environment, actions or feelings, abstract structures</i>	<i>chaotic, intense, jagged, multiple layers, single objects</i>	<i>round, circular, spiral, sharp, shaking, distinct objects, radiating, natural</i>
 Nose (P11)	 Crackles (P4)	 Telephonic (P1)	 Piano (P2)	 Impact (P12)
 Processed Guitar (P20)	 String Grains (P8)	 Processed Guitar (P22)	 String (P2)	 Piano (P16)

Table 1.: Sketch categories with examples. Category names and keywords were obtained through thematic analysis as described in Section 3.6. *Objects/Scenes* mainly refers to real-world associations while other categories highlight different abstract approaches, but category clusters might overlap with a number of sketches showing characteristics of more than one category. Colours were inverted for better visibility.

	Grains		Lines		Object/ Scenes		Chaotic/ Jagged		Radiating/ Round	
	Mean	Std	Mean	Std	Mean	Std	Mean	Std	Mean	Std
Points	560	532	534	338	1068	557	1060	696	914	722
Strokes	14	12	2	2	13	9	7	10	4	4

Table 2.: Average number of points and strokes for each sketch category rounded to the closest integer. *Object/Scenes* shows the highest number of points and second highest number of strokes which implies that this category could be reduced when limiting the size of a sketch.

6. Quantitative feature analysis

This analysis was conducted with the aim of statistically investigating to what extent sound-shape associations emerge from sound-sketches. First, inter-rater reliability was determined to confirm sketch features used in this research can measure agreement between participants. This was followed by calculating correlations between individual audio and visual features. As discussed in Section 3.7, the interpretation of the results needs to consider that this analysis uses averaged values

6.1. Inter-rater reliability

The results of the ICC(2,k) inter-rater reliability measures are illustrated in Figure 14. For Study 1, reliability measures were good to excellent for *Intersections* and *Acute Angles*, poor to good for *Average Speed* and moderate to good for all remaining features within the 95% confidence interval (CI). For Study 2, measures were excellent for *Acute Angles* and good to excellent for *Intersections*, *Wide Curves* and *Right Angles* solidifying findings from Study 1 with higher reliability scores and narrower confidence intervals. In contrast to Study 1, *Number of Strokes* and *Average Time* only returned

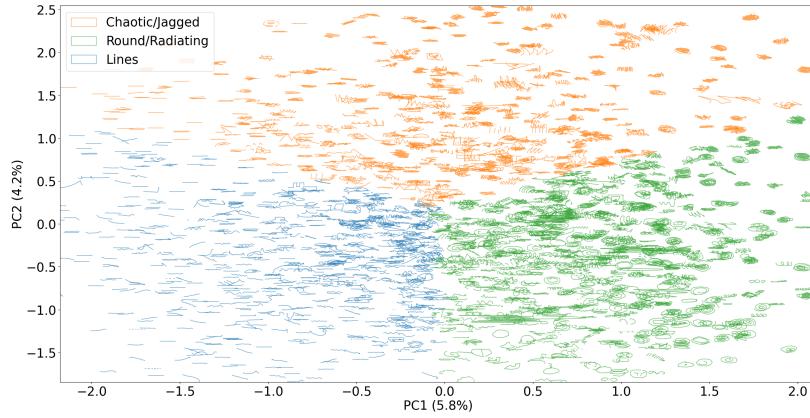


Figure 10.: Study 2 sound sketches organised in the latent space that was reduced to two dimensions with PCA. Colours indicate the different clusters found through K-Means.

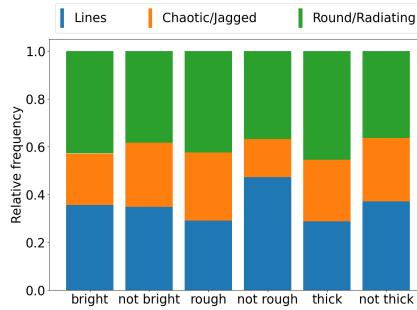


Figure 11.: Categories by sounds grouped by attributes in Study 2.

651 poor to moderate reliability which can be accredited to the sketch length limit of the
 652 interface. *Average Speed* however showed good to excellent reliability compared to poor
 653 to good in Study 1 implying that participants might have expressed characteristics
 654 through their sketching speed that would have been expressed through longer, more
 655 complex sketches with an unrestricted interface. This is supported by a larger variance
 656 in average sketching speed compared to Study 1 ($SD = 0.23$ compared to $SD = 0.12$).
 657 These results suggest that some level of agreement exists between averaged participants
 658 on how to represent sounds visually and that it can be measured with the extracted
 659 sketch features.

660 6.2. Feature correlation

661 Several significant correlations were found between sketch and audio features in both
 662 studies illustrated in Figure 13. For Study 1, *Acute Angles* (11), *Intersections* (9) and
 663 *Number of Strokes* (8) show the highest statistically significant ($p < .05$) number of
 664 strong ($r > .6$) and very strong ($r > .8$) correlations with audio features. The strongest
 665 correlation overall was found between *RMS Mean* and *Average Time* ($r = .95$, $p < .001$).
 666 Opposing audio features like *Warmth* and *Sharpness* showed similar absolute corre-

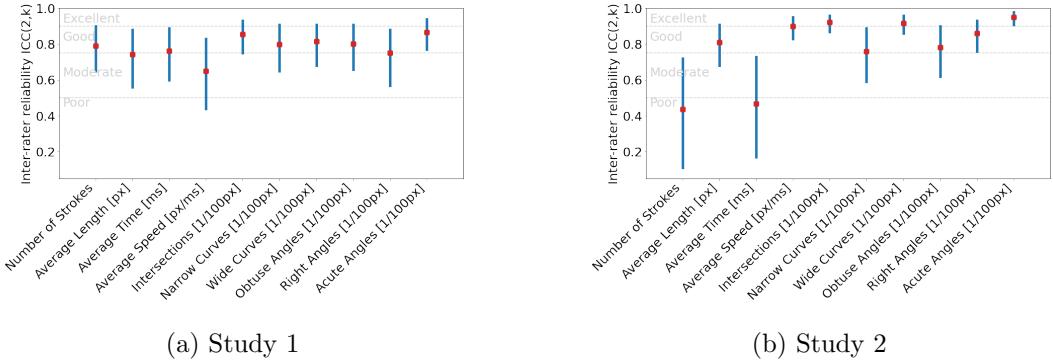


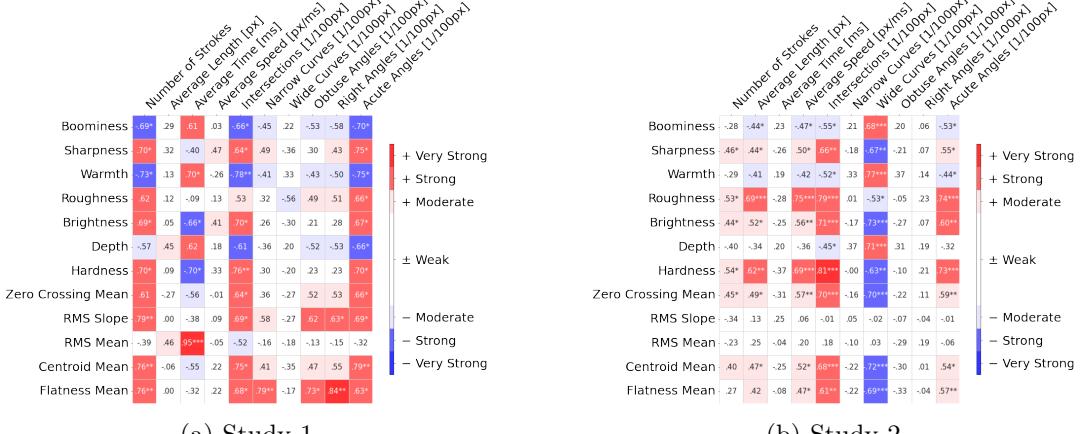
Figure 12.: Mean values and 95% CI of $ICC(2,k)$ inter-rater reliabilities for each sketch feature with evaluation guidelines proposed by Koo and Li (2016) ($df_1=19$, $df_2=513$, $p<.01$ for all features).

667 lation values but opposite directions for *Number of Strokes*, *Intersections* and *Acute*
 668 *Angles*. For Study 2, *Wide Curves* (9), *Intersections* (7) and *Acute Angles* (3) show the
 669 highest statistically significant ($p<.05$) number of strong ($r>.6$) and very strong ($r>.8$)
 670 correlations with audio features. The strongest correlation overall was found between
 671 *Intersections* and *Hardness* ($r=.81$, $p<.001$). While only 12 strong and very strong
 672 correlations were found compared to 19 in Study 1, more results were significant at
 673 $p<.05$ level (49 compared to 37). Similar trends can be observed between both studies
 674 with *Acute Angles* and *Intersections* showing negative correlations with audio features
 675 *Boominess*, *Warmth* and *Depth* and positive correlations with *Roughness*, *Brightness*
 676 and *Hardness*. *Average Speed* shows similar correlations to *Number of Strokes* in Study
 677 1 for example with *Hardness* and *Boominess* further supporting the hypothesis that
 678 sketching speed compensated for the sketch length restrictions. In contrast to Study 1,
 679 *Wide Curves* shows a large number of significant correlations with audio features that
 680 are mirroring *Acute Angles* correlations as expected from known sound-shape associa-
 681 tions. As all sound stimuli were created with the same amplitude envelope in Study
 682 2 the features *RMS Slope* and *RMS Mean* did not provide distinguishing descriptions
 683 and consequently did not show any significant correlations with sketch features.

684 7. Discussion

685 7.1. Abstract and figurative representations of sound stimuli

686 The sound-sketches collected in this research were categorised through a rigorous,
 687 human-centred process consisting of participant interviews and a card-sorting study.
 688 On the highest level, sound-sketches were divided into two groups: *abstract* and
 689 *figurative*. Abstract representations appear to be strongly informed by cross-modal
 690 associations and show many similarities to visual stimuli used in matching tasks.
 691 Figurative representations, on the other hand, depict objects or scenes associated with
 692 a sound that may depict the sound source directly, for example in form of a musical
 693 instrument, or may be informed by an emotional response or memory, for example, a
 694 sound perceived as scary might be represented with a scene similar to the sketch in
 695 Figure 7 that is informed by the memory of a movie scene that corresponds to that
 696 emotion. Which representational approach a participant took was influenced by the



(a) Study 1

(b) Study 2

Figure 13.: Spearman’s rank correlation coefficients between sketch and audio features with annotated p-values: $p < .05$ (*), $.01$ (**), $.001$ (***)

sound type or a participant’s experience. Figurative representations were found to be more prevalent among non-musicians and for sounds with an easily identifiable sound source like a musical instrument. This could be interpreted as a difference in listening modes: coined by Schaeffer (2017) and further explored by Chion (2019), the *reduced listening* mode describes a focus on the sound itself as opposed to *semantic listening* which focuses on the source or meaning of a sound. Described as hardly natural by Chion (2019), *reduced listening* can be more demanding for the untrained ear, especially for sounds from a familiar source. This analysis suggests that two different approaches would be needed when computationally mapping sketches to sound. For abstract sketches that encode information about sound characteristics directly in their form, extracting quantitative sketch features like the ones described in Section 3.7 appears to be a viable approach. For figurative sketches, a more suitable approach might have to include object recognition for sketches of musical instruments and other sound-producing objects or sentiment analysis for emotionally informed scenes.

Section 6 shows that for Study 1 multiple significant correlations were found between audio and sketch features despite including abstract and figurative sketches. With 84% of sketches classified as abstract, figurative sketches did not appear to impact these correlations. Because of the small sample size, differences in correlations between abstract and figurative sketches could not be investigated in a meaningful way. Figurative representations might still be indirectly influenced by cross-modal associations, for example, an uncomfortable noisy and dissonant sound might be represented with a scene similar to the design *Landscape of Thorns* by Moisey (2017) that aims to communicate discomfort and danger through sharp and jagged structures. On the other hand, abstract representations might include references to symbolic representations of sound that are not based on cross-modal associations. The feedback analysis in Section 4 reveals that some sketches, in particular in the *Lines* category, are informed by audio representations like waveforms, spectrograms or amplitude envelopes.

As described in Section 2.1, cross-modal associations are not fixed; they emerge in different forms depending on the stimulus and situation and sketch representations

729 might incorporate multiple associations or references to objects, scenes or symbols.
 730 Without explicitly asking a participant what they had in mind, it may not be
 731 possible to exactly determine their underlying motivation. Generally, only a few
 732 participants strictly adhered to one representational approach. While dominated by
 733 a primary approach, most participants switched between or mixed representational
 734 styles throughout the sketching task.

735 **7.2. Cross-modal associations in free-form sketches**

736 Compared to the figurative category *Object/Scenes* the abstract categories *Chaotic/-*
 737 *Jagged, Grains, Lines* and *Radiating/Rounds* appear to be primarily influenced by cross-
 738 modal associations. This research aimed to investigate sound-shape associations specific-
 739 ally, but it might be useful to define the terminology to better interpret the results.
 740 The visual stimuli used in sound-shape matching tasks like Adeli et al. (2014) focus
 741 very strictly on shape meaning that stimuli only show the outlining *form* through a
 742 non-intersecting connected series of lines. Looking at Figures 8 and 10 it is obvious
 743 that participants did not only use *form* in their abstract sketches. Knees and Ander-
 744 sen (2016) made similar findings and, in their sound-sketch prototype, included tools
 745 to create outlines of simple forms like circles and triangles that can be given more
 746 complexity through free-form lines and filled with *textures*.

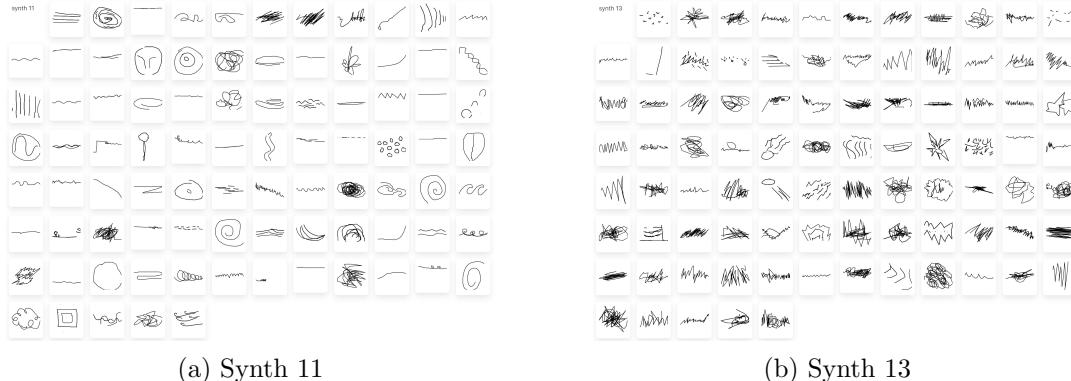


Figure 14.: Sound-sketches taken from Study 2 for the stimuli Synth 11, characterized by high values for *Warmth* and low values for *Roughness* and *Spectral Flatness*, and Synth 13, which exhibits contrasting values for these attributes. There is a difference in form, with Synth 11 more frequently depicted as *Radiating/Round* and Synth 13 as *Chaotic/Jagged* in sketches. In addition, Synth 11 tends to be represented by less complex sketches that fall into the *Lines* category.

747 In this research, the abstract categories *Chaotic/Jagged* and *Radiating/Round*
 748 appear to encode sound characteristics through their *form* with the latter found
 749 more frequently for warm sounds in Study 1 and the former for rough sounds in
 750 both studies which aligns with existing sound-shape research (Köhler, 1929; Adeli et
 751 al., 2014). However, *Chaotic/Jagged* appears to be also contrasted by *Lines* which
 752 is more frequently found in the least rough sounds. Sketches in the *Lines* category
 753 have considerably fewer sharp angles than sketches in *Chaotic/Jagged* which could
 754 be interpreted as a difference in *form*, but *Lines* also exhibits considerably fewer
 755 stroke intersections than *Chaotic/Jagged* which can be interpreted as a difference in
 756 *complexity*. The remaining abstract category *Grains* was particularly prevalent for the

757 intersected sounds *Crackles* and *String Grains* in Study 1 and is described to consist
758 of small, repetitive components which can be interpreted as *texture* and was quantified
759 by the average number of strokes and the average stroke length. Strictly speaking,
760 sound-shape associations are not the only cross-modal associations that influence
761 abstract representations, however, in this context, the definition of *shape* could be
762 extended to describe any 2-dimensional structure composed of straight or curved
763 lines. For quantitative analysis, this broader definition of shape might be sufficient as
764 the sketch features described in Section 3.7 capture more than just the *form* of a sketch.

765

766 The results of the feature correlation analysis in Section 6.2 support that typi-
767 cal sound-shape associations influenced sketch representations with acute angles
768 positively correlating with attributes like *Sharpness*, *Roughness* and *Brightness* and
769 negatively with *Warmth* and *Boominess*. The sketch feature *Wide Curves* which was
770 expected to quantify roundness did not show any significant correlations in Study 1,
771 but multiple significant correlations for Study 2 that point in the opposite direction
772 of *acute angles*. This might be due to noise introduced by higher variance in sketch
773 approaches in Study 1 but could suggest that the extracted features are not a good
774 measure for the overall roundness of a sketch. *Intersections* proved to be meaningful
775 in describing cross-modal associations showing multiple significant correlations with
776 audio features across both studies. The direction of these correlations is the same
777 for acute angles, leading to the belief that the opposing pairs *rough* and *soft* for
778 sound might not only be visualised through jaggedness or roundness but also through
779 complexity and simplicity. Overall, the roughness of a sound appears to have a strong
780 influence on how sketches are represented which is clearly visualised in Figure 11.
781 The figure also suggests that the *Lines* category is more prevalent for *thin* sounds,
782 but the differences found in this study did not prove to be significant. As the sketch
783 and audio features do not directly quantify thinness, feature correlation could not
784 provide any additional information. Further work could focus on this attribute of
785 sound which would provide an additional timbral dimension encoded in sound-sketch
786 representations.

787

788 A general problem that emerged in both studies is that it cannot be clearly
789 determined which sound characteristic participants focused on. As a participant
790 in Study 1 stated, there are a lot of different aspects of sound and it is difficult
791 to represent them all in a simple sketch. This might be made even more difficult
792 with the sketch length limit that was introduced in Study 2. While it succeeded
793 in guiding participants to represent sounds in more abstract ways, it does prevent
794 them from elaborating or overlaying multiple approaches that might capture more
795 aspects of a sound. It might be possible that thinness or thickness can be encoded
796 through sketching, but participants mainly focused on roughness when representing
797 sound. For future research, it could be considered to also ask participants to rate
798 which attributes of the sound they focused on with a design similar to Hayes et
799 al. (2021). An alternative method might involve prompting participants to envision
800 a sound possessing specific attributes instead of exposing them to an actual au-
801 ditory stimulus. This approach was employed during the test session in Study 2,
802 where participants were tasked with sketching a noisy and calm sound. However,
803 understanding how participants mentally conceptualize a sound remains somewhat
804 ambiguous, particularly considering that certain common auditory descriptors, such
805 as "sharpness" or "roundness" may inherently include visual elements. In addition to
806 verbal feedback, allowing participants to select a sound that most accurately captures

807 their mental representation of a descriptor could offer valuable insights into their
808 attribute assignments.

809 **7.3. Implications for development of a sketch-based sound synthesiser**

810 Study 2 was designed with a future implementation of a sketch-based sound synthesiser
811 in mind. The results can shed light on which mapping architecture is more appropriate
812 for such a system:

- 813 • a regression approach where a change in timbre is induced by an incremental
814 change in the corresponding sketch feature. For example, a sound would become
815 noisier with an increase in sharp angles.
- 816 • a classification approach where the overall timbre category is determined. In this
817 scenario, a sketch could represent multiple categories for example either *rough*
818 or *soft* in combination with either *thick* or *thin*.

819 The feature correlations in Section 6.2 describe a monotonic relationship between
820 audio and sketch features which can lead to the conclusion that a gradual increase in
821 an audio feature like roughness coincides with a gradual increase in a sketch feature
822 like *acute angles*. However, it has to be remembered that these sketch features rep-
823 resent the average sketch of the participants meaning that a single participant might
824 not follow a change of timbre in such an incremental fashion. Rather participants
825 could think of sounds in categories and with rising roughness, an increasing number of
826 participants switch from round to jagged shapes. Similarly, statistical agreement for
827 sketch representations was only determined for an average participant in Section 6.1
828 and, in this research, significant results could not be obtained when looking at individ-
829 ual participants. In addition, if a specific, computed audio feature increases linearly
830 for a series of sounds, it does not mean that humans would perceive this as a linear
831 change in timbre. When thinking about a sound that is neither particularly rough
832 nor soft, it is likely that another more prominent feature, for example, *thick* or *thin*
833 would inform the perception of this sound and hence the sketch representation of it. In
834 fact, Hayes et al. (2022b) found that relatively small, localised clusters emerge when
835 participants were asked to define sound within a timbre space according to descriptors
836 like *rough* or *soft*. A linear interpolation between the parameters of two FM syn-
837 thesiser sounds might not be perceived as a linear transition in sound characteristics, but
838 rather different sound environments that emerge along this path. Given this analysis,
839 a categorisation approach appears to be more viable when designing a sketch-based
840 sound synthesiser. It should also be added that sketch to synthesiser mappings do not
841 have to be limited to timbre. An advantage of this design is that multiple parameters
842 can be represented with a single input. Revisiting the research of Salgado-Montejo
843 et al. (2016) and M. B. Küssner et al. (2014) shows that sketch representations can
844 also serve to represent pitch or amplitude envelopes. Future work can explore how a
845 sketch input could manipulate these parameters at the same time, potentially using a
846 combination of mappings that are established through a data-driven machine learning
847 approach and hard-coded mappings that are based on findings from empirical stud-
848 ies. Given that colour and visual texture have been identified as common associations
849 with timbre in previous studies (Ward et al., 2006; Adeli et al., 2014; Gurman et al.,
850 2021; Grill & Flexer, 2012), the exploration of polychromatic colour palettes, different
851 stroke widths and brush styles could be considered in further development to facilitate
852 multi-dimensional mappings.

853 **8. Conclusion**

854 In two participant studies, 2320 free-form sound-sketches (560 from Study 1 and 1760
855 from Study 2) were collected with simple, monochromatic digital sketching interfaces.
856 Through a rigorous human-centred analysis sketches were categorised by their primary
857 representational approach. Several significant correlations between quantitative audio
858 and sketch features were found that align with findings from cross-modal matching
859 tasks. The results show that while sound-shape associations play a significant role
860 in sketched representations, humans incorporate other visual aspects like structural
861 complexity or texture as well or choose figurative representations of emotions or sound-
862 producing objects. Some level of agreement on how to represent sounds could be
863 found between participants which appears strongest for sounds that are dominated
864 by one sound characteristic. This research provides useful insights and suggestions for
865 designing a sketch-based sound synthesiser.

866 **9. Acknowledgments**

867 EPSRC and AHRC Centre for Doctoral Training in Media and Arts Technology
868 (EP/L01632X/1).

869 **10. Disclosure statement**

870 All studies that contributed to this research were approved by the Queen Mary Uni-
871 versity of London Ethics Committee (Reference numbers QMREC2341 for Study 1
872 and QMERC2517a the card sorting exercise, QMERC20.478 for participant evalua-
873 tion of the interface development and QMERC20.009 for Study 2). All participants
874 gave informed consent. No monetary compensation was offered except for Study 2
875 which was conducted via Prolific rewarding £2.50 for 20-minute participation. The
876 authors report there are no competing interests to declare.

877 **References**

- 878 Adeli, M., Rouat, J., & Molotchnikoff, S. (2014). Audiovisual correspondence between
879 musical timbre and visual shapes. *Frontiers in human neuroscience*, 8, 352.
880 doi:
881 Albertazzi, L., Canal, L., & Micciolo, R. (2015). Cross-modal associations between
882 materic painting and classical spanish music. *Frontiers in Psychology*, 6, 424.
883 doi:
884 ANSI. (1994). Timbre. *American National Standards Institute. Psychoacoustic ter-
885 minology*.
886 Baron-Cohen, S., Burt, L., Smith-Laittan, F., Harrison, J., & Bolton, P. (1996).
887 Synaesthesia: prevalence and familiality. *Perception*, 25(9), 1073–1079.
888 Blake, D. K. (2012, June). Timbre as Differentiation in Indie Music. *Music Theory
889 Online*, 18(2). doi:
890 Bottini, R., Barilari, M., & Collignon, O. (2019). Sound symbolism in sighted and
891 blind. the role of vision and orthography in sound-shape correspondences. *Cog-
892 nition*, 185, 62–70. doi:

- 893 Braun, V., & Clarke, V. (2006). Using thematic analysis in psychology. *Qualitative
894 Research in Psychology*, 3(2), 77–101. doi:
- 895 Bremner, A. J., Caparos, S., Davidoff, J., de Fockert, J., Linnell, K. J., & Spence, C.
896 (2013). “Bouba” and “Kiki” in Namibia? a remote culture make similar shape–
897 sound matches, but different shape–taste matches to westerners. *Cognition*,
898 126(2), 165–172. doi:
- 899 Bresenham, J. E. (1965). Algorithm for computer control of a digital plotter. *IBM
900 Systems journal*, 4(1), 25–30.
- 901 Bruford, F., Barthet, M., McDonald, S., & Sandler, M. B. (2019). Groove explorer:
902 An intelligent visual interface for drum loop library navigation. In *Proceedings
903 of the acm iui workshops*. Los Angeles, USA: CEUR-WS.org.
- 904 Catmull, E., & Rom, R. (1974). A class of local interpolating splines. In *Computer
905 aided geometric design* (pp. 317–326). Elsevier. doi:
- 906 Chion, M. (2019). 2. the three listening modes. In *Audio-vision: Sound on screen* (pp.
907 22–34). New York Chichester, West Sussex: Columbia University Press. doi:
- 908 Clemente, A., Vila-Vidal, M., Pearce, M. T., Aguiló, G., Corradi, G., & Nadal, M.
909 (2020). A set of 200 musical stimuli varying in balance, contour, symmetry,
910 and complexity: Behavioral and computational assessments. *Behavior Research
911 Methods*, 52(4), 1491–1509.
- 912 Cramer, J., Wu, H.-H., Salamon, J., & Bello, J. P. (2019, May). Look, listen and learn
913 more: Design choices for deep audio embeddings. In *Ieee int. conf. on acoustics,
914 speech and signal processing (icassp)* (pp. 3852–3856). Brighton, UK.
- 915 Cuskley, C., Dingemanse, M., Kirby, S., & Van Leeuwen, T. M. (2019). Cross-modal
916 associations and synesthesia: Categorical perception and structure in vowel–color
917 mappings in a large online sample. *Behavior Research Methods*, 51, 1651–1675.
918 doi:
- 919 Davis, R. (1961). The Fitness of Names to Drawings. a Cross-Cultural Study in
920 Tanganyika. *British Journal of Psychology*, 52(3), 259–268. doi:
- 921 Deng, L. (2012). The mnist database of handwritten digit images for machine learning
922 research [best of the web]. *IEEE signal processing magazine*, 29(6), 141–142.
- 923 Douglas, D. H., & Peucker, T. K. (1973). Algorithms for the reduction of the number
924 of points required to represent a digitized line or its caricature. *Cartographica:
925 the international journal for geographic information and geovisualization*, 10(2),
926 112–122. doi:
- 927 Engeln, L., & Groh, R. (2020). Coherence of audible shapes—a qualitative user
928 study for coherent visual audio design with resynthesized shapes. *Personal and
929 Ubiquitous Computing*, 1–11. doi:
- 930 Engeln, L., Le, N. L., McGinity, M., & Groh, R. (2021). Similarity analysis of visual
931 sketch-based search for sounds. In *Proceedings of audio mostly 2021* (pp. 101–
932 108). Trento, Italy: Association for Computing Machinery. doi:
- 933 Fried, O., Jin, Z., Finkelstein, A., & Oda, R. (2014). *AudioQuilt: 2D Arrangements
934 of Audio Samples using Metric Learning and Kernelized Sorting*.
- 935 Garber, L., y Ciencia, M. A., Ciccola, T., & Amusategui, J. C. (2021). Audiostellar, an
936 open source corpus-based musical instrument for latent sound structure discovery
937 and sonic experimentation. In *Proceedings of international computer music con-
938 ference* (pp. 86–91). Santiago, Chile. Retrieved from <https://hdl.handle.net/2027/fulcrum.t435gg568>
- 939 Giannakis, K. (2006). A comparative evaluation of auditory-visual mappings for sound
940 visualisation. *Organised Sound*, 11(3), 297–307.

- 942 Giannakis, K., & Smith, M. (2000). Towards a theoretical framework for sound
 943 synthesis based on auditoryvisual associations.
- 944 Grey, J. M. (1977). Multidimensional perceptual scaling of musical timbres. *the
 945 Journal of the Acoustical Society of America*, 61(5), 1270–1277.
- 946 Grill, T., & Flexer, A. (2012). Visualization of Perceptual Qualities in Textural
 947 Sounds. In *Proceedings of international computer music conference* (pp. 589–
 948 596). Ljubljana, Slovenia: Michigan Publishing Services. Retrieved from <http://hdl.handle.net/2027/spo.bbp2372.2012.110>
- 949 Gurman, D., McCormick, C. R., & Klein, R. M. (2021). Crossmodal correspondence
 950 between auditory timbre and visual shape. *Multisensory Research*, 35(3), 221–
 951 241.
- 952 Ha, D., & Eck, D. (2017). A neural representation of sketch drawings. *arXiv preprint
 953 arXiv:1704.03477*.
- 954 Hayes, B., & Saitis, C. (2020). There's more to timbre than musical instruments:
 955 semantic dimensions of FM sounds. In *Proceedings of international conference
 956 on timbre*. Thessaloniki, Greece: Timbre 2020.
- 957 Hayes, B., Saitis, C., & Fazekas, G. (2021). Perceptual and semantic scaling of
 958 fm synthesis timbres: Common dimensions and the role of expertise. *ICMPC-
 959 ESCOM*.
- 960 Hayes, B., Saitis, C., & Fazekas, G. (2022a). Disembodied timbres: A study on
 961 semantically prompted fm synthesis. *Journal of the Audio Engineering Society*,
 962 70(5), 373–391.
- 963 Hayes, B., Saitis, C., & Fazekas, G. (2022b). timbre.fun: A gamified interactive
 964 system for crowdsourcing a timbre semantic vocabulary. In *Proceedings of the
 965 24th international congress on acoustics* (p. 10).
- 966 Iverson, P., & Krumhansl, C. (1993). Isolating the dynamic attributes of musical
 967 timbres. *The Journal of the acoustical society of America*, 94(5), 2595–2603.
 968 doi:
- 969 Knees, P., & Andersen, K. (2016). Searching for audio by sketching mental images
 970 of sound: A brave new idea for audio retrieval in creative music production. In
 971 *Proceedings of international conference on multimedia retrieval* (pp. 95–102).
 972 New York, USA: Association for Computing Machinery. doi:
- 973 Köhler, W. (1929). Gestalt psychology. *Liveright*.
- 974 Koo, T. K., & Li, M. Y. (2016). A guideline of selecting and reporting intraclass
 975 correlation coefficients for reliability research. *Journal of chiropractic medicine*,
 976 15(2), 155–163.
- 977 Küssner, M. (2014). *Shape, drawing and gesture: Cross-modal mappings of sound and
 978 music* (Unpublished doctoral dissertation). King's College London (University
 979 of London).
- 980 Küssner, M. B., Tidhar, D., Prior, H. M., & Leech-Wilkinson, D. (2014). Musicians are
 981 more consistent: Gestural cross-modal mappings of pitch, loudness and tempo
 982 in real-time. *Frontiers in Psychology*, 5. doi:
- 983 LeCun, Y. (1998). The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>.
- 984 Lewis, J. R. (2018, July). The System Usability Scale: Past, Present, and Future.
 985 *International Journal of Human-Computer Interaction*, 34(7), 577–590. doi:
- 986 Löbbers, S., Barthet, M., & Fazekas, G. (2021). Sketching sounds: an exploratory
 987 study on sound-shape associations. In *Proceedings of international computer
 988 music conference* (pp. 299–304). Santiago, Chile: Michigan Publishing Services.
 989 Retrieved from <https://hdl.handle.net/2027/fulcrum.t435gg568>

- 992 Löbbers, S., & Fazekas, G. (2022). Seeing sounds, hearing shapes: a gamified study to
 993 evaluate sound-sketches. In *Proceedings international computer music conference*
 994 (pp. 174–179). Limerick, Ireland: Michigan Publishing Services. Retrieved from
 995 <https://hdl.handle.net/2027/fulcrum.nk322g689>
- 996 MacDonald, J., & McGurk, H. (1978). Visual influences on speech perception pro-
 997 cesses. *Perception & psychophysics*, 24(3), 253–257. doi:
 998 Marin, M. M., Schober, R., Gingras, B., & Leder, H. (2017). Misattribution of musical
 999 arousal increases sexual attraction towards opposite-sex faces in females. *PLoS
 1000 One*, 12(9), e0183531. doi:
 1001 Maurer, D., Pathman, T., & Mondloch, C. J. (2006). The shape of boubas:
 1002 Sound–shape correspondences in toddlers and adults. *Developmental Science*,
 1003 9(3), 316–322. doi:
 1004 McAdams, S., Winsberg, S., Donnadieu, S., De Soete, G., & Krimphoff, J. (1995).
 1005 Perceptual scaling of synthesized musical timbres: Common dimensions, speci-
 1006 ficities, and latent subject classes. *Psychological research*, 58, 177–192.
 1007 McFee, B., Metsai, A., McVicar, M., Balke, S., Thomé, C., Raffel, C., ... Thassilo
 1008 (2022, February). *librosa/librosa: 0.9.1*. Zenodo. Retrieved from <https://doi.org/10.5281/zenodo.6097378> doi:
 1009 1010 McFee, B., Raffel, C., Liang, D., Ellis, D. P., McVicar, M., Battenberg, E., & Nieto,
 1011 O. (2015). librosa: Audio and music signal analysis in python. In *Proceedings of
 1012 the 14th python in science conference* (Vol. 8, pp. 18–25).
 1013 Moisey, A. (2017). Permanent negative value: The waste isolation pilot plant. *Critical
 1014 Inquiry*, 43(4), 861–892. doi:
 1015 Müllensiefen, D., Gingras, B., Musil, J., & Stewart, L. (2014). The musicality of non-
 1016 musicians: an index for assessing musical sophistication in the general population.
 1017 *PloS one*, 9(2), e89642.
 1018 Nielsen, A. K. S., & Rendall, D. (2013). Parsing the role of consonants versus vowels
 1019 in the classic Takete-Maluma phenomenon. *Canadian Journal of Experimen-
 1020 tal Psychology/Revue canadienne de psychologie expérimentale*, 67(2), 153–163.
 1021 doi:
 1022 Paea, S., & Baird, R. (2018). Information architecture (ia): Using multidimensional
 1023 scaling (mds) and k-means clustering algorithm for analysis of card sorting data.
 1024 *Journal of Usability Studies*, 13(3).
 1025 Pearce, A., Brookes, T., & Mason, R. (2019). Modelling Timbral Hardness. *Applied
 1026 Sciences*, 9(3), 466. doi:
 1027 Peeters, G. (2004). A large set of audio features for sound description (similarity
 1028 and classification) in the cuidado project. *CUTDADO 1st Project Report*, 54(0),
 1029 1–25.
 1030 Provenzano, C. (2018). *Auto-Tune, Labor, and the Pop-Music Voice*. Oxford Univer-
 1031 sity Press.
 1032 Ramachandran, V. S., & Hubbard, E. M. (2001). Synesthesia—a window into percep-
 1033 tion, thought and language. *Journal of consciousness studies*, 8(12), 3–34.
 1034 Rousseeuw, P. J. (1987). Silhouettes: a graphical aid to the interpretation and val-
 1035 idation of cluster analysis. *Journal of computational and applied mathematics*,
 1036 20, 53–65.
 1037 Rucsanda, M. D., et al. (2019). Aspects of the relationship between music and painting
 1038 and their influence on schoenberg and kandinsky. *Bulletin of the Transilvania
 1039 University of Brașov, Series VIII: Performing Arts*, 12(2), 91–100.
 1040 Saitis, C., & Weinzierl, S. (2019). The semantics of timbre. In K. Siedenburg, C. Saitis,
 1041 S. McAdams, A. N. Popper, & R. R. Fay (Eds.), *Timbre: Acoustics, perception,*

- 1042 and cognition (pp. 119–149). Cham: Springer International Publishing. Retrieved
1043 from https://doi.org/10.1007/978-3-030-14832-4_5 doi:
1044 Saitis, C., Weinzierl, S., von Kriegstein, K., Ystad, S., & Cuskley, C. (2020). Timbre
1045 semantics through the lens of crossmodal correspondences: A new way of asking
1046 old questions. *Acoustical Science and Technology*, 41(1), 365–368.
1047 Salgado-Montejo, A., Marmolejo-Ramos, F., Alvarado, J. A., Arboleda, J. C., Suarez,
1048 D. R., & Spence, C. (2016). Drawing sounds: representing tones and chords
1049 spatially. *Experimental Brain Research*, 234, 3509–3522.
1050 Schaeffer, P. (2017). *Treatise on musical objects: An essay across disciplines* (Vol. 20).
1051 Univ of California Press.
1052 Seago, A. (2013). A new interaction strategy for musical timbre design. In
1053 S. Holland, K. Wilkie, P. Mulholland, & A. Seago (Eds.), *Music and human-*
1054 *computer interaction* (pp. 153–169). London: Springer London. Retrieved from
1055 https://doi.org/10.1007/978-1-4471-2990-5_9 doi:
1056 Sezgin, T. M. (2001). *Feature point detection and curve approximation for early pro-*
1057 *cessing of free-hand sketches* (Unpublished doctoral dissertation). Massachusetts
1058 Institute of Technology.
1059 Shinohara, K., Yamauchi, N., Kawahara, S., & Tanaka, H. (2016, September). Takete
1060 and Maluma in Action: A Cross-Modal Relationship between Gestures and
1061 Sounds. *PLOS ONE*, 11(9), e0163525. doi:
1062 Sidhu, D. M., & Pexman, P. M. (2018). Five mechanisms of sound symbolic associa-
1063 tion. *Psychonomic bulletin & review*, 25, 1619–1643.
1064 Siedenburg, K., Fujinaga, I., & McAdams, S. (2016). A comparison of approaches to
1065 timbre descriptors in music information retrieval and music psychology. *Journal*
1066 *of New Music Research*, 45(1), 27–41.
1067 Siedenburg, K., & McAdams, S. (2017, October). Four Distinctions for the Auditory
1068 “Wastebasket” of Timbre1. *Frontiers in Psychology*, 8, 1747. doi:
1069 Soraghan, S., Faire, F., Renaud, A., & Supper, B. (2018). A new timbre visualization
1070 technique based on semantic descriptors. *Computer Music Journal*, 42(1), 23–
1071 36.
1072 Spence, C. (2011). Crossmodal correspondences: A tutorial review. *Attention, Per-*
1073 *ception, & Psychophysics*, 73, 971–995.
1074 Stuckey, H. L. (2015, June). The second step in data analysis: Coding qualitative
1075 research data. *Journal of Social Health and Diabetes*, 03(1), 7–10. doi:
1076 Taylor, I. K., & Taylor, M. M. (1962). Phonetic symbolism in four unrelated languages.
1077 *Canadian Journal of Psychology*, 16(4), 344–356. doi:
1078 Thoret, E., Aramaki, M., Bringoux, L., Ystad, S., & Kronland-Martinet, R. (2016,
1079 April). Seeing Circles and Drawing Ellipses: When Sound Biases Reproduction
1080 of Visual Motion. *PLOS ONE*, 11(4), e0154475. doi:
1081 Van Doorn, G. H., Wuillemin, D., & Spence, C. (2014). Does the colour of the mug
1082 influence the taste of the coffee? *Flavour*, 3(1), 1–7.
1083 Wallmark, Z., Frank, R. J., & Nghiêm, L. (2019). Creating novel tones from adjectives:
1084 An exploratory study using fm synthesis. *Psychomusicology: Music, Mind, and*
1085 *Brain*, 29(4), 188.
1086 Wallmark, Z., & Kendall, R. A. (2021, 10). 578579C23Describing Sound: The Cognitive
1087 Linguistics of Timbre. In *The Oxford Handbook of Timbre*. Oxford University
1088 Press. Retrieved from <https://doi.org/10.1093/oxfordhb/9780190637224.013.14> doi:
1089 Ward, J., Huckstep, B., & Tsakanikos, E. (2006). Sound-colour synesthesia: to what
1090 extent does it use cross-modal mechanisms common to us all? *Cortex*, 42(2),
1091

- 1092 264–280.
- 1093 Wolin, A., Eoff, B., & Hammond, T. (2008). ShortStraw: A Simple and Effective
1094 Corner Finder for Polylines. In *Proceedings of eurographics workshop on sketch-*
1095 *based interfaces and modeling*. Annecy, France: The Eurographics Association.
1096 doi:
- 1097 Xiong, Y., & LaViola Jr, J. J. (2009). Revisiting shortstraw: improving corner finding
1098 in sketch-based interfaces. In *Proceedings of eurographics symposium on sketch-*
1099 *based interfaces and modeling* (pp. 101–108). New Orleans USA: Association for
1100 Computing Machinery. doi:
- 1101 Zacharakis, A., & Pastiadiis, K. (2015). A confirmatory approach of the luminance-
1102 texture-mass model for musical timbre semantics. In *Proceedings of the au-*
1103 *dio mostly 2015 on interaction with sound*. New York, NY, USA: Associa-
1104 tion for Computing Machinery. Retrieved from <https://doi.org/10.1145/2814895.2814898> doi:
- 1105 1106 Zacharakis, A., & Pastiadiis, K. (2016). Revisiting the luminance-texture-mass model
1107 for musical timbre semantics: A confirmatory approach and perspectives of ex-
1108 tension. *Journal of the Audio Engineering Society*, 64(9), 636–645.

1109 **Appendix A. Sound and sketch feature extraction**

- 1110 Table A1 shows the values of extracted audio features for Study 1 and Study 2.
1111 Table A2 show the values of extracted sketch features for Study 1 and Study 2.

	Flatness Mean	Centroid Mean	RMS Mean	RMS Slope	Zero Cross. Mean	Hard- ness	Depth	Bright- ness	Rough- ness	Warmth	Sharp- ness	Boomi- ness
Crackles	4.1 * 10^{-2}	1728	0.032	21	0.035	55	28	60	52	41	48	17
Telephonic	1.0 * 10^{-3}	543	0.157	5	0.013	34	67	43	50	54	28	41
Strings	3.9 * 10^{-5}	1042	0.151	10	0.020	42	57	52	52	54	35	33
String Grains	1.9 * 10^{-3}	1177	0.107	22	0.031	52	40	56	56	49	37	25
Subbass	1.5 * 10^{-8}	206	0.416	1	0.002	34	77	36	47	65	34	47
Noise	2.9 * 10^{-1}	7915	0.086	34	0.144	81	44	81	81	29	68	22
Piano	4.5 * 10^{-7}	542	0.051	1	0.013	42	65	49	40	54	31	42
Impact	2.8 * 10^{-4}	1047	0.097	8	0.011	65	77	60	49	51	49	39
Processed Guitar	1.7 * 10^{-3}	229	0.344	3	0.002	30	80	39	40	63	37	46
Electric Guitar	4.5 * 10^{-6}	1295	0.305	5	0.022	47	56	56	56	50	42	34

(a) Study 1

	Flatness Mean	Centroid Mean	RMS Mean	RMS Slope	Zero Cross. Mean	Hard- ness	Depth	Bright- ness	Rough- ness	Warmth	Sharp- ness	Boomi- ness
Synth 1	4.9 * 10^{-8}	562	0.0429	0	0.0199	34	42	39	29	39	29	21
Synth 2	1.4 * 10^{-7}	794	0.0614	0	0.0239	52	65	54	53	50	37	34
Synth 3	3.8 * 10^{-5}	2889	0.0266	0	0.113	67	42	75	72	30	61	9
Synth 4	5.5 * 10^{-3}	9670	0.0529	0	0.437	73	31	86	83	17	77	-6
Synth 5	1.6 * 10^{-8}	247	0.0731	0	0.01	32	68	30	41	53	25	43
Synth 6	5.9 * 10^{-3}	8607	0.0374	0	0.355	76	34	84	77	20	76	12
Synth 7	6.0 * 10^{-6}	2671	0.116	0	0.11	64	43	75	72	31	59	13
Synth 8	2.7 * 10^{-6}	327	0.0893	0	0.0116	11	62	30	42	48	23	40
Synth 9	8.4 * 10^{-7}	1693	0.104	0	0.0583	63	46	67	70	39	48	25
Synth 10	5.9 * 10^{-7}	1259	0.057	0	0.0511	48	45	65	68	40	40	23
Synth 11	1.5 * 10^{-7}	259	0.1	0	0.00995	3	64	27	0	56	26	44
Synth 12	3.5 * 10^{-5}	5493	0.0577	0	0.192	50	33	79	57	20	72	4
Synth 13	7.8 * 10^{-4}	8715	0.0539	0	0.357	61	30	84	63	31	80	25
Synth 14	4.7 * 10^{-8}	432	0.0466	0	0.0157	33	55	36	45	46	25	32
Synth 15	1.8 * 10^{-7}	515	0.0578	0	0.0194	39	61	46	56	52	29	33
Synth 16	1.5 * 10^{-5}	2768	0.0521	0	0.121	55	37	73	56	30	62	4
Synth 17	9.6 * 10^{-3}	9683	0.124	0	0.416	77	38	86	82	17	76	-1
Synth 18	4.1 * 10^{-8}	613	0.0758	0	0.0187	35	56	48	51	45	30	33
Synth 19	8.2 * 10^{-7}	569	0.0574	0	0.02	54	60	48	59	52	29	33
Synth 20	2.6 * 10^{-3}	6691	0.0564	0	0.259	69	25	82	65	25	74	15

(b) Study 2

Table A1.: Sound features extracted from sound stimuli of both studies. For Librosa features, the mean values of all windows are reported.

	Number of Strokes	Average Length [px]	Average Time [ms]	Average Speed [px/ms]	Inter- sections [1/100px]	Narrow Curves [1/100px]	Wide Curves [1/100px]	Obtuse Angles [1/100px]	Right Angles [1/100px]	Acute Angles [1/100px]
Crackles	10.5	471	2796	0.22	1.78	0.55	0.60	0.94	0.22	0.67
Telephonic	6.0	1091	4613	0.27	1.05	0.16	0.37	0.55	0.08	0.26
Strings	3.9	900	4379	0.26	0.31	0.10	0.23	0.28	0.06	0.09
String Grains	10.1	914	3340	0.28	1.48	0.24	0.21	0.58	0.15	0.54
Subbass	4.5	1328	6154	0.23	0.92	0.17	0.43	0.51	0.08	0.30
Noise	12.8	1816	3540	0.57	2.37	0.32	0.22	0.61	0.16	0.56
Piano	3.8	586	3020	0.29	0.61	0.06	0.31	0.20	0.02	0.06
Impact	7.2	1239	3313	0.51	1.12	0.14	0.28	0.24	0.05	0.31
Processed Guitar	4.3	1124	4707	0.33	0.43	0.17	0.23	0.49	0.09	0.22
Electric Guitar	5.3	1121	4561	0.33	0.93	0.12	0.12	0.35	0.08	0.40

(a) Study 1

	Number of Strokes	Average Length [px]	Average Time [ms]	Average Speed [px/ms]	Inter- sections [1/100px]	Narrow Curves [1/100px]	Wide Curves [1/100px]	Obtuse Angles [1/100px]	Right Angles [1/100px]	Acute Angles [1/100px]
Synth 1	3.4	1002	3021	0.44	1.8	0.3	0.84	1.1	0.13	0.31
Synth 2	3.1	876	2923	0.4	2.6	0.55	0.96	1.1	0.24	0.53
Synth 3	2.9	1186	3274	0.46	2.0	0.42	0.88	1.1	0.079	1.2
Synth 4	3.9	2125	3178	0.74	5.4	0.32	0.39	0.86	0.23	1.4
Synth 5	3.0	1636	3239	0.53	1.2	0.9	1.2	1.7	0.53	0.62
Synth 6	4.3	1444	3344	0.49	3.7	0.75	0.67	1.4	0.19	1.6
Synth 7	3.9	3644	2725	1.2	5.8	0.5	0.54	1.0	0.19	1.5
Synth 8	2.03	936	3666	0.31	0.74	0.37	0.84	0.58	0.096	0.21
Synth 9	5.2	1815	2752	0.79	5.7	0.42	0.98	0.95	0.19	0.7
Synth 10	3.9	1555	3008	0.56	2.2	0.45	1.0	1.6	0.18	0.48
Synth 11	2.1	834	3682	0.27	0.6	0.19	1.0	0.92	0.082	0.12
Synth 12	2.6	1352	3533	0.46	2.9	0.34	0.61	0.68	0.12	0.63
Synth 13	3.8	1239	3125	0.48	3.2	0.41	0.68	0.81	0.14	0.6
Synth 14	2.5	995	3181	0.37	0.97	0.5	0.9	0.78	0.045	0.2
Synth 15	2.8	1543	2984	0.55	3.3	0.62	0.74	1.5	0.31	1.4
Synth 16	2.7	1104	2816	0.48	3.1	0.63	0.78	1.1	0.18	0.72
Synth 17	2.3	3055	3171	1.1	6.4	0.34	0.24	0.58	0.15	1.2
Synth 18	2.6	1269	3555	0.39	2.0	0.35	0.68	0.72	0.11	0.34
Synth 19	3.3	1854	3279	0.54	3.8	0.83	0.89	1.6	0.31	1.4
Synth 20	4.1	1993	2961	0.67	3.2	0.26	0.54	0.71	0.12	0.59

(b) Study 2

Table A2.: Sketch features for both studies. The mean value from all participants is presented for each sound stimulus.