

Manuscript Number:

Title: Evaluation of live human-computer music-making: quantitative and qualitative approaches

Article Type: SI: Sonic Interaction Design (Serafin)

Keywords: Music; qualitative; quantitative

Corresponding Author: Mr Dan Stowell,

Corresponding Author's Institution: Queen Mary University of London

First Author: Dan Stowell

Order of Authors: Dan Stowell; Andrew Robertson; Nick Bryan-Kinns; Mark D Plumbley

Abstract: Live music-making using interactive systems is not completely amenable to traditional HCI evaluation metrics such as task-completion rates. In this paper we discuss quantitative and qualitative approaches which provide opportunities to evaluate the music-making interaction, accounting for aspects which cannot be directly measured or expressed numerically, yet which may be important for participants. We present case studies in the application of a qualitative method based on Discourse Analysis, and a quantitative method based on the Turing Test. We compare and contrast these methods with each other, and with other evaluation approaches used in the literature, and discuss factors affecting which evaluation methods are appropriate in a given context.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Evaluation of live human-computer music-making: quantitative and qualitative approaches

D. Stowell*, A. Robertson, N. Bryan-Kinns, M. D. Plumbley

*Centre for Digital Music, School of Electronic Engineering and Computer Science,
Queen Mary University of London, UK*

Abstract

Live music-making using interactive systems is not completely amenable to traditional HCI evaluation metrics such as task-completion rates. In this paper we discuss quantitative and qualitative approaches which provide opportunities to evaluate the music-making interaction, accounting for aspects which cannot be directly measured or expressed numerically, yet which may be important for participants. We present case studies in the application of a qualitative method based on Discourse Analysis, and a quantitative method based on the Turing Test. We compare and contrast these methods with each other, and with other evaluation approaches used in the literature, and discuss factors affecting which evaluation methods are appropriate in a given context.

Key words: Music, qualitative, quantitative

1 Introduction

Live human-computer music-making, with reactive or interactive systems, is a topic of recent artistic and engineering research (d'Escrivan and Collins, 2007, esp. chapters 3, 5, 8). However, the formal evaluation of such systems is relatively little-studied (Fels, 2004). As one indicator, a survey of recent research papers presented at the conference on New Interfaces for Musical Expression (NIME – a conference about user interfaces for music-making)

* Corresponding author. Centre for Digital Music, School of Electronic Engineering and Computer Science, Queen Mary University of London, London E1 4NS, UK.
Email address: dan.stowell@elec.qmul.ac.uk

Evaluation type	NIME conference year		
	2006	2007	2008
<i>Not applicable</i>	8	9	7
None	18	14	15
Informal	12	8	6
Formal qualit.	1	2	3
Formal quant.	2	3	3
Total formal	3 (9%)	5 (19%)	6 (22%)

Table 1

Survey of oral papers presented at the conference on New Interfaces for Musical Expression (NIME), indicating the type of evaluation described. The last line indicates the total number of formal evaluations presented, also given as a percentage of the papers (excluding those for which evaluation was not applicable).

shows a consistently low proportion of papers containing formal evaluations (Table 1).

A formal evaluation is one presented in rigorous fashion, which presents a structured route from data collection to results (e.g. by specifying analysis techniques). It therefore establishes the degree of generality and repeatability of its results. Formal evaluations, whether quantitative or qualitative, are important because they provide a basis for generalising the outcomes of user tests, and therefore allow researchers to build on one another’s work.

Live human-computer music making poses challenges for many common HCI evaluation techniques. Musical interactions have creative and affective aspects, which means they cannot be described as tasks for which e.g. completion rates can reliably be measured. They also have dependencies on timing (rhythm, tempo, etc.), and feedback interactions (e.g. between performers, between performer and audience), which further problematise the issue of developing valid and reliable experimental procedures.

“Talk-aloud” protocols (Ericsson and Simon, 1996, section 2.3) are used in many HCI evaluations. However, in some musical performances (such as singing or playing a wind instrument) the use of the speech apparatus for music-making precludes concurrent talking. More generally, speaking may interfere with the process of rhythmic/melodic performance: speech and music cognition can demonstrably interfere with each other (Salamé and Baddeley, 1989), and the brain resources used in speech and music processing partially overlap (Peretz and Zatorre, 2005), suggesting issues of cognitive “competition” if subjects are asked to produce music and speech simultaneously.

Other observational approaches may be applicable, although in many cases

1 observing a participant’s reactions may be difficult: because of the lack of
2 objectively observable indications of “success” in musical expression, but also
3 because of the participant’s physical involvement in the music-making process
4 (e.g. the whole-body interaction of a drummer with a drum-kit).
5

6 Another challenging aspect of musical interface evaluation is that the partic-
7 ipant populations are often small (Wanderley and Orio, 2002). For example,
8 it may be difficult to recruit many virtuoso violinists, human beatboxers, or
9 jazz trumpeters, for a given experiment. Therefore evaluation methods should
10 be applicable to relatively small study sizes.
11

12
13 In this paper we present two methods developed specifically for evaluation of
14 live musical systems, and which accommodate the issues described above.
15

16 17 18 19 *1.1 Previous work* 20

21
22 Some prior work has looked at HCI issues in “offline” musical systems, i.e.
23 tools for composers (e.g. Buxton and Sniderman (1980); Polfreman (2001)).
24 Others have used theoretical considerations to produce recommendations and
25 heuristics for designing musical performance interfaces (Hunt and Wander-
26 ley, 2002; Levitin et al., 2003; Fels, 2004; de Poli, 2004), although with-
27 out explicit empirical validation. Note that in some such considerations, a
28 “Composer→Performer→Audience” model is adopted, in which musical ex-
29 pression is defined to consist of timing and other variations applied to the
30 composed musical score (Goebel, 2004; de Poli, 2004). In this work we wish
31 to consider musical interaction more generally, encompassing improvised and
32 interactive performance situations. Wanderley and Orio (2002) provide a par-
33 ticularly useful contribution to our topic. They discuss pertinent HCI meth-
34 ods, before proposing a task-based approach to musical interface evaluation
35 using “maximally simple” musical tasks such as the production of glissandi
36 or triggered sequences. The authors propose a user-focused evaluation, using
37 Likert-scale feedback (Grant et al., 1999) as opposed to an objective measure
38 of gesture accuracy, since such objective measures may not be a good represen-
39 tation of the musical qualities of the gestures produced. The authors suggest
40 by analogy with Fitts’ law (Card et al., 1978) that their task-based approach
41 may allow for quantitative comparisons of musical interfaces.
42
43
44
45
46
47
48

49 Wanderley & Orio’s framework is interesting but may have some drawbacks.
50 The reduction of musical interaction to maximally simple tasks risks compro-
51 mising the authenticity of the interaction, creating situations in which the
52 affective and creative aspects of music-making are abstracted away. In other
53 words, the reduction conflates *controllability* of a musical interface with *ex-*
54 *pressiveness* of that interface (Dobrian and Koppelman, 2006). The use of
55
56
57
58

1 Likert-scale metrics also may have some difficulties. They are susceptible to
2 cultural differences (Lee et al., 2002) and psychological biases (Nicholls et al.,
3 2006), and may require large sample sizes to achieve sufficient statistical power
4 (Göb et al., 2007).
5

6 Acknowledging the relative scarcity of research on the topic of live human-
7 computer music-making, we may look to other areas which may provide useful
8 analogies. The field of computer games is notable here, since it carries some of
9 the features of live music-making: it can involve complex multimodal interac-
10 tions, with elements of goal-oriented and affective involvement, and a degree of
11 learning. For example, Barendregt et al. (2006) investigates the usability and
12 affective aspects of a computer game for children, during first use and after
13 some practice. Mandryk and Atkins (2007) use a combination of physiological
14 measures to produce a continuous estimate of the emotional state (arousal
15 and valence) of subjects playing a computer game.
16
17
18
19
20

21 In summary, although there have been some useful forays into the field of
22 expressive musical interface evaluation, and some work in related disciplines
23 such as that of computer games evaluation, the field could certainly benefit
24 from further development. In this paper we hope to contribute to this area
25 by presenting work on two different evaluation approaches which we have
26 developed.
27
28
29
30
31
32

33 *1.2 Outline of paper* 34 35 36

37 We first present two methods for evaluation of live musical systems. We de-
38 scribe each method, along with a case study of its application. The methods
39 are
40
41

- 42 (1) A qualitative method using Discourse Analysis (Section 2), to evaluate a
43 system by illuminating how users conceptually integrate the system into
44 the context of use.
- 45 (2) A Turing-Test method, designed for the case when the system is intended
46 to respond in a human-like manner (Section 3).
47
48
49

50 Then in Section 4 we compare and contrast the methods with each other,
51 and with other evaluation approaches described in the literature, and discuss
52 factors affecting which approaches are appropriate in a given context. Section
53 4.2 aims to distil the discussion down to recommendations which may be used
54 by a researcher wishing to evaluate an interactive musical system.
55
56
57
58

2 A qualitative approach: Discourse Analysis

Interviews and free-text comments are sometimes reported in studies on musical interfaces. However, often they are conducted in a relatively informal context, and only quotes or summaries are reported rather than any structured analysis, therefore providing little analytic reliability. Good qualitative methods penetrate deeper than simple summaries, offering insight into text data (Antaki et al., 2004). *Discourse Analysis* (DA) is one such approach, used in disciplines such as psychology and social sciences (Banister et al. (1994); Silverman (2006), chapter 6).

DA's strength comes from using a *structured method* which can take apart the language used in discourses (e.g. interviews, written works) and elucidate the connections and implications contained within, while remaining faithful to the content of the original text (Antaki et al., 2004). DA is designed to go beyond the specific sequence of phrases used in a conversation, and produce a structured analysis of the conversational resources used, the relations between entities, and the “work” that the discourse is doing.

Uszkoreit (1996) summarises the aim of DA very compactly:

“The problems addressed in discourse research aim to answer two general kinds of questions:

- (1) what information is contained in extended sequences of utterances that goes beyond the meaning of the individual utterances themselves?
- (2) how does the context in which an utterance is used affect the meaning of the individual utterances, or parts of them?”

DA is not a single method, rather an analytical tradition in which various methods have been developed. Our DA method is based on that of Banister et al. (1994, p. 95–102) and is elaborated in section 2.2.

2.1 Method

We wish to use the power of DA as part of a qualitative and formal method which can explore issues such as expressivity and affordances for users of interactive musical systems. Longitudinal studies may also be useful, but imply a high cost in time and resources. Therefore our design aims to provide users with a brief but useful period of exploration of a new musical interface, including interviews and discussion which we can then analyse.

In any evaluation of a musical interface one must decide the context of the evaluation. Is the interface being evaluated as a successor or alternative to

1 some other interface (e.g. an electric cello vs an acoustic cello)? Who is ex-
2 pected to use the interface (e.g. virtuosi, amateurs, children)? Such factors
3 will affect not only the recruitment of participants but also some aspects of
4 the experimental setup.

5
6 Our method is designed either to trial a single interface with no explicit com-
7 parison system, or to compare two similar systems (as is done below in our
8 case study). The method consists of two types of user session, solo sessions
9 followed by group session(s), plus the Discourse Analysis of data collected.

12 13 *2.1.1 Solo sessions*

14
15
16 In order to explore individuals' personal responses to the interface(s), we first
17 perform solo sessions in which a participant is invited to try out the interface(s)
18 for the first time. If there is more than one interface to be used, the order of
19 presentation is randomised in each session.

20
21
22 The solo session consists of three phases for each interface:

23
24 **Free exploration.** The participant is encouraged to try out the interface for
25 a while and explore it in their own way.

26
27 **Guided exploration.** The participant is presented with audio examples of
28 recordings created using the interface, in order to indicate the range of pos-
29 sibilities, and encouraged to create recordings inspired by those examples.
30 This is not a precision-of-reproduction task; precision-of-reproduction is ex-
31 plicitly not evaluated, and participants are told that they need not replicate
32 the examples.

33
34 **Semi-structured interview.** The interview's main aim is to encourage the
35 participant to discuss their experiences of using the interface in the free and
36 guided exploration phases, both in relation to prior experience and to the
37 other interfaces presented if applicable. Both the free and guided phases are
38 video recorded, and the interviewer may play back segments of the recording
39 and ask the participant about them, in order to stimulate discussion.

40
41
42
43 The raw data to be analysed is the interview transcript. Our aim is for the
44 participant to construct their own descriptions and categories, which means
45 it is very important that the interviewer is experienced in neutral interview
46 technique, and can avoid (as far as possible) introducing labels and concepts
47 that do not come from the participant's own language patterns.

48 49 50 51 52 *2.1.2 Group session*

53
54
55 To complement the solo sessions we also conduct a group session. Peer group
56 discussion can produce more and different discussion around a topic, and can

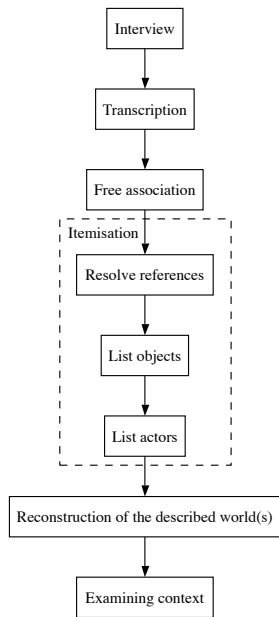


Fig. 1. Outline of our Discourse Analysis procedure.

demonstrate the group negotiation of categories, labels, comparisons, and so on. The focus-group tradition provides a well-studied approach to such group discussion (Stewart, 2007). Our group session has a lot in common with a typical focus group in terms of the facilitation and semi-structured group discussion format. In addition we make available the interface(s) under consideration and encourage the participants to experiment with them during the session.

As in the solo sessions, the transcribed conversation is the data to be analysed. A neutral facilitation technique is also important here, to encourage all participants to speak, to allow opposing points of view to emerge in a non-threatening environment, and to allow the group to negotiate the use of language with minimal interference.

2.2 Data analysis

Our DA approach to analysing the data is based on that of Banister et al. (1994, p. 95–102), adapted to our study context. The DA of text is a relatively intensive and time-consuming method. It can be automated to some extent, but not completely, because of the close linguistic attention required. Our approach is summarised in Figure 1 and consists of the following five steps:

- (a) **Transcription.** The speech data is transcribed, using a standard style of notation which includes all speech events (including repetitions, speech fragments, pauses). This is to ensure that the analysis can remain close to

Transcription	Object (referent)	Description	Is a subject?
...to see how the other person did it,	the other person ((recording the examples))	Participant was trying to work out what this person was ((doing))	Y
Because it was more fun Yeah, I just think the noises were a bit more, bit different, you know	((Interface)) the noises ((made by Q))	Participant preferred this ((to X)) because it was more fun were a bit more, bit different	
you could come up with some slightly more funky noises.	((general person))	could come up with some slightly more funky noises ((in Q, than X))	Y
	noises	((general person)) could come up with some slightly more funky ones ((in Q, than X))	

Fig. 2. Excerpt from a spreadsheet used during the itemisation of interview data, for step (c) of the Discourse Analysis.

what is actually said, and avoid adding a gloss which can add some distortion to the data. For purposes of analytical transparency, the transcripts (suitably anonymised) should be published alongside the analysis results.

(b) Free association. Having transcribed the speech data, the analyst reads it through and notes down surface impressions and free associations. These can later be compared against the output from the later stages.

(c) Itemisation of transcribed data. The transcript is then broken down by itemising every single object in the discourse (i.e. all the entities referred to). Pronouns such as “it” or “he” are resolved, using the participant’s own terminology as far as possible. For every object an accompanying description of the object is extracted from that speech instance – again using the participant’s own language, essentially by rewriting the sentence/phrase in which the instance is found.

The list of objects is scanned to determine if different ways of speaking can be identified at this point. Also, those objects which are also “actors” (or “subjects”) are identified – i.e. those which act with agency in the speech instance; they need not be human.

It is helpful at this point to identify the most commonly-occurring objects and actors in the discourse, as they will form the basis of the later reconstruction.

Figure 2 shows an excerpt from a spreadsheet used during our DA process, showing the itemisation of objects and subjects, and the descriptions extracted.

(d) Reconstruction of the described world. Starting with the list of most commonly-occurring objects and actors in the discourse, the analyst reconstructs the depictions of the world that they produce. This could for example be achieved using concept maps to depict the interrelations between the actors and objects. If different ways of speaking have been identified, there will typically be one reconstructed “world” per way of speaking. Overlaps and contrasts between these worlds can be identified. Figure 3 shows an excerpt of a concept map representing a “world” distilled in this way.

The “worlds” we produce are very strongly tied to the participant’s own discourse. The actors, objects, descriptions, relationships, and relative importances, are all derived from a close reading of the text. These worlds are essentially just a methodically reorganised version of the participant’s own language.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

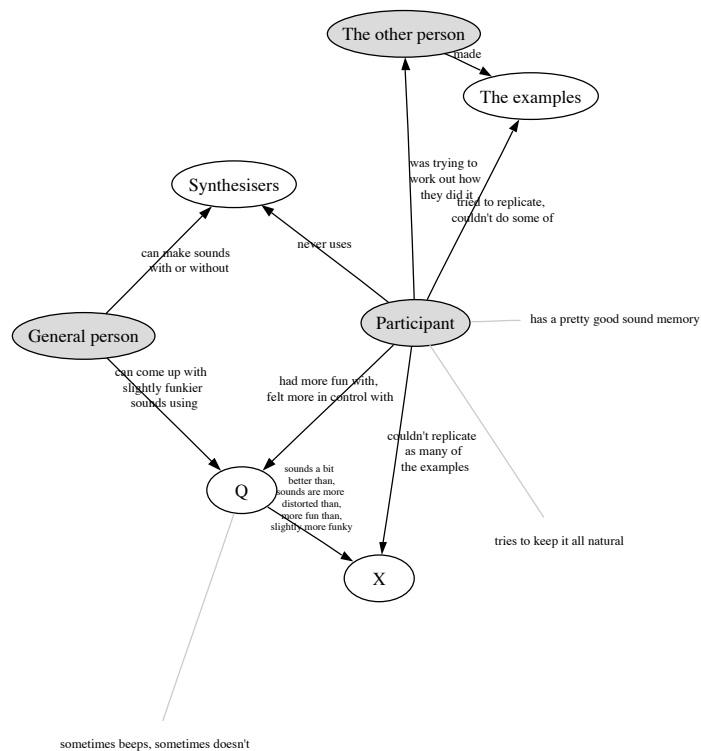


Fig. 3. An example of a reconstructed set of relations between objects in the described world. This is a simplified excerpt of the reconstruction for User 2 in our study. Objects are displayed in ovals, with the shaded ovals representing actors.

In our particular context, we may be interested in the user’s conceptualisation of musical interfaces. It is particularly interesting to look at how these are situated in the described world, and particularly important to avoid preconceptions about how users may describe an interface: for example, a given interface could be: an instrument; an extension of a computer; two or more separate items (e.g. a box and a screen); an extension of the individual self; or it could be absent from the discourse.

- (e) **Examining context.** The relevant context of the discourse typically depends on the field of study, for example whether it is political or psychological. Here we have created an explicit context of other participants. After running the previous steps of DA on each individual transcript, we compare and contrast the described worlds produced from each transcript, first comparing those in the same experimental condition (i.e. same order of presentation, if relevant), then across all participants. We also compare the DA of the focus group session(s) against that of the solo sessions.



Fig. 4. Timbre remapping maps the timbral space of a voice source onto that of a target synthesiser.

2.3 The method in action: evaluating voice timbre remapping

The present case-study was conducted in the context of a project to develop voice-based interfaces for controlling musical systems. Our interface uses a process we call *timbre remapping* to allow the timbral variation in a voice to control the timbral variation of an arbitrary synthesiser (Figure 4). The procedure involves analysing vocal timbre in real-time to produce a multi-dimensional “timbre space”, then retrieving the synthesis parameters that correspond best to that location in the timbre space. The method is described further by Stowell and Plumbley (2007).

In our study we wished to evaluate the timbre remapping system with beatboxers (vocal percussion musicians), for two reasons: they are one target audience for the technology in development; and they have a familiarity and level of comfort with manipulation of vocal timbre that should facilitate the study sessions.

We recruited by advertising online (a beatboxing website) and around London for amateur or professional beatboxers. Participants were paid £10 per session plus travel expenses to attend sessions in our (acoustically-isolated) studio. We recruited five participants from the small community, all male and aged 18–21. One took part in a solo session; one in the group session; and three took part in both. Their beatboxing experience ranged from a few months to four years. Their use of technology for music ranged from minimal to a keen use of recording and effects technology (e.g. Cubase).

1 In our study we wished to investigate any effect of providing the timbre remap-
2 ping feature. To this end we presented two similar interfaces: both tracked the
3 pitch and volume of the microphone input, and used these to control a syn-
4 thesiser, but one also used the timbre remapping procedure to control the
5 synthesiser’s timbral settings. The synthesiser used was an emulated General
6 Instrument AY-3-8910 (General Instrument, early 1980s), which was selected
7 because of its wide timbral range (from pure tone to pure noise) with a well-
8 defined control space of a few integer-valued variables. Participants spent a
9 total of around 30–60 minutes using the interfaces, and 15–20 minutes in in-
10 terview. Analysis of the interview transcripts using the procedure of section
11 2.1 took approximately 9 hours per participant (around 2000 words each).
12
13

14 We do not report a detailed analysis of the group session transcript here: the
15 group session generated information which is useful in the development of our
16 system, but little which bears directly upon the presence or absence of timbral
17 control. We discuss this outcome further in section 4.
18
19
20

21 In the following, we describe the main findings from analysis of the solo ses-
22 sions, taking each user one by one before drawing comparisons and contrasts.
23 We emphasise that although the discussion here is a narrative supported by
24 quotes, it reflects the structures elucidated by the DA process – the full tran-
25 scriptions and Discourse Analysis tables are available online¹. In the study, con-
26 dition “Q” was used to refer to the system with timbre remapping active, “X”
27 for the system with timbre remapping inactive.
28
29
30
31

32 *2.3.1 Reconstruction of the described world*

33

34 **User 1** expressed positive sentiments about both Q (with timbre remapping)
35 and X (without timbre remapping), but preferred Q in terms of sound qual-
36 ity, ease of use and being “more controllable”. In both cases the system was
37 construed as a reactive system, making noises in response to noises made into
38 the microphone; there was no conceptual difference between Q and X – for
39 example in terms of affordances or relation to other objects.
40
41
42

43 The “guided exploration” tasks were treated as reproduction tasks, despite
44 our intention to avoid this. User 1 described the task as difficult for X, and
45 easier for Q, and situated this as being due to a difference in “randomness”
46 (of X) vs. “controllable” (of Q).
47
48

49 **User 2** found the the system (in both modes) “didn’t sound very pleasing
50 to the ear”. His discussion conveyed a pervasive structured approach to the
51 guided exploration tasks, in trying to infer what “the original person” had
52
53

54 ¹ [http://www.elec.qmul.ac.uk/digitalmusic/papers/2008/
55 Stowell08ijhcs-data/](http://www.elec.qmul.ac.uk/digitalmusic/papers/2008/Stowell08ijhcs-data/)
56
57
58

1 done to create the examples and to reproduce that. In both Q and X the
2 approach and experience was the same.

3
4 Again, User 2 expressed preference for Q over X, both in terms of sound
5 quality and in terms of control. Q was described as more fun and “slightly
6 more funky”. Interestingly, the issues that might bear upon such preferences
7 are arranged differently: issues of unpredictability were raised for Q (but not
8 X), and the guided exploration task for Q was felt to be more difficult, in part
9 because it was harder to infer what “the original person” had done to create
10 the examples.
11

12
13 **User 3**’s discourse placed the system in a different context compared to others.
14 It was construed as an “effect plugin” rather than a reactive system, which
15 implies different affordances: for example, as with audio effects it could be
16 applied to a recorded sound, not just used in real-time; and the description
17 of what produced the audio examples is cast in terms of an original sound
18 recording rather than some other person. This user had the most computer
19 music experience of the group, using recording software and effects plugins
20 more than the others, which may explain this difference in contextualisation.
21

22
23
24 User 3 found no difference in sound or sound quality between Q and X, but
25 found the guided exploration of X more difficult, which he attributed to the
26 input sounds being more varied.
27

28
29
30 **User 4** situated the interface as a reactive system, similar to Users 1 and 2.
31 However, the sounds produced seemed to be segregated into two streams rather
32 than a single sound – a “synth machine” which follows the user’s humming,
33 plus “voice-activated sound effects”. No other users used such separation in
34 their discourse.
35

36
37 “Randomness” was an issue for User 4 as it was for some others. Both Q and X
38 exhibited randomness, although X was much more random. This randomness
39 meant that User 4 found Q easier to control. The pitch-following sound was
40 felt to be accurate in both cases; the other (sound effects / percussive) stream
41 was the source of the randomness.
42

43
44 In terms of the output sound, User 4 suggested some small differences but
45 found it difficult to pin down any particular difference, but felt that Q sounded
46 better.
47

48 49 50 51 *2.3.2 Examining context*

52
53 Users 1 and 2 were presented with the conditions in the order XQ; Users 3 and
54 4 in the order QX. Order-of-presentation may have some small influence on the
55 outcomes: Users 3 and 4 identified little or no difference in the output sound
56
57

1 between the conditions (User 4 preferred Q but found the difference relatively
2 subtle), while Users 1 and 2 felt more strongly that they were different and
3 preferred the sound of Q. It would require a larger study to be confident that
4 this difference really was being affected by order-of-presentation.
5

6 In our study we are not directly concerned with which condition sounds better
7 (both use the same synthesiser in the same basic configuration), but this is an
8 interesting aspect to come from the study. We might speculate that differences
9 in perceived sound quality are caused by the different way the timbral changes
10 of the synthesiser are used. However, participants made no conscious connec-
11 tion between sound quality and issues such as controllability or randomness.
12
13

14 Taking the four participant interviews together, no strong systematic differ-
15 ences between Q and X are seen. All participants situate Q and X similarly,
16 albeit with some nuanced differences between the two. Activating/deactivating
17 the timbre remapping facet of the system does not make a strong enough dif-
18 ference to force a reinterpretation of the system.
19
20
21

22 A notable aspect of the four participants' analyses is the differing ways the
23 system is situated (both Q and X). As designers of the system we may have
24 one view of what the system "is", perhaps strongly connected with technical
25 aspects of its implementation, but the analyses presented here illustrate the
26 interesting way that users situate a new technology alongside existing tech-
27 nologies and processes. The four participants situated the interface in differing
28 ways: either as an audio effects plugin, or a reactive system; as a single out-
29 put stream or as two. We emphasise that none of these is the "correct" way
30 to conceptualise the interface. These different approaches highlight different
31 facets of the interface and its affordances.
32
33
34
35

36 During the analyses we noted that all participants maintained a conceptual
37 distance between themselves and the system, and analogously between their
38 voice and the output sound. There was very little use of the "cyborg" discourse
39 in which the user and system are treated as a single unit, a discourse which
40 hints at mastery or "unconscious competence". This fact is certainly under-
41 standable given that the participants each had less than an hour's experience
42 with the interface. It demonstrates that even for beatboxers with strong expe-
43 rience in manipulation of vocal timbre, controlling the vocal interface requires
44 learning – an observation confirmed by the participant interviews.
45
46
47
48

49 The issue of "randomness" arose quite commonly among the participants.
50 However, randomness emerges as a nuanced phenomenon: although two of
51 the participants described X as being more random than Q, and placed ran-
52 domness in opposition to controllability (as well as preference), User 2 was
53 happy to describe Q as being more random and also more controllable (and
54 preferable).
55
56
57
58
59

1 A uniform outcome from all participants was the conscious interpretation of
2 the guided exploration tasks as precision-of-reproduction tasks. This was evi-
3 dent during the study sessions as well as from the discourse around the tasks.
4 As one participant put it, “If you’re not going to replicate the examples, what
5 are you gonna do?”
6

7 A notable absence from the discourses, given our research context, was dis-
8 cussion which might bear on expressivity, for example the expressive range of
9 the interfaces. Towards the end of each interview we asked explicitly whether
10 either of the interfaces was more expressive, and responses were generally non-
11 committal. We propose that this was because our tasks had failed to engage the
12 participants in creative or expressive activities: the (understandable) reduction
13 of the guided exploration task to a precision-of-reproduction task must have
14 contributed to this. We also noticed that our study design failed to encourage
15 much iterative use of record-and-playback to develop ideas. In section 4 we
16 suggest some possible implications of these findings on future study design.
17
18
19
20

21 We have seen the Discourse Analysis method in action and the information it
22 can yield, about how users situate a system in relation to themselves and other
23 objects. In the next section we will turn to consider an alternative evaluation
24 approach based on the Turing Test, before comparing and contrasting the
25 methods.
26
27
28
29
30

31 **3 A quantitative approach: musical Turing Test**

32
33

34 Turing’s seminal paper (Turing, 1950) proposes replacing the question “can
35 a computer think?”, by an “Imitation Game”, now commonly known as the
36 *Turing Test*, in which the computer is required to imitate a human being in an
37 interrogation. If the computer is able to fool a human interrogator a substantial
38 amount of the time, then the computer can be credited with “intelligence”.
39
40
41

42 There has been considerable debate around the legitimacy of this approach
43 as a measure of artificial intelligence (e.g. Searle (1980)). However, without
44 making any claims about the intelligence of musical systems, we can say that
45 often they are designed with the aim of reacting or interacting in a human-like
46 fashion. Therefore the degree of observer confusion between human and auto-
47 mated response is an appropriate route for evaluating systems which perform
48 human-like tasks, such as score-based accompaniment or musical improvisa-
49 tion. Analysing this degree of confusion could allow us to make numerical
50 comparisons between systems, each of which aim to emulate some human
51 skill, and evaluate their relative success at this emulation.
52
53
54
55

56 Our example concerns the task of real-time beat tracking with a live drummer.
57
58
59
60
61
62
63
64
65

1 We have developed a beat tracker specifically for such live use, named “B-
2 Keeper” (Robertson and Plumbley, 2007), which is event-based and uses a
3 method related to the oscillator models used by Large (1995) and Toivainen
4 (1998). Whilst the B-Keeper does not generate musical parts, it does control
5 the tempo of the accompaniment so that it responds to subtle tempo changes
6 being made by the drummer.
7

8 We wished to develop a test suitable for assessing this real-time interaction.
9 Established beat tracking evaluations exist, typically comparing annotated
10 beat positions against ground-truths provided by human annotators (McKin-
11 ney et al., 2007). However, these neglect the component of interaction, and
12 do not attempt to judge the degree of “naturalness” or “musicality” of any
13 variation in beat annotations.
14
15

16 Qualitative approaches such as that described above could be appropriate.
17 However, in this case we are interested specifically in evaluating the beat-
18 tracker’s designed ability to interact in a human-like manner, which the mu-
19 sical Turing Test allows us to quantify.
20
21
22
23
24

25 *3.1 Method*

26
27
28

29 In our application of the musical Turing Test to evaluate the B-Keeper system,
30 we decided to perform a three-way comparison, incorporating human, machine,
31 and a third “control” condition using a steady accompaniment which remains
32 at a fixed tempo dictated by the drummer. Our experiment is depicted in Fig-
33 ure 5. For each test, the drummer gives four steady beats of the kick drum to
34 set the tempo and start, then plays along to an accompaniment track. This is
35 performed three times. Each time, a human tapper (one of the authors, AR)
36 taps the tempo on the keyboard, keeping time with the drummer, but only for
37 one of the three times will this be altering the tempo of the accompaniment.
38 For these trials, controlled by the human tapper, we applied a Gaussian win-
39 dows to the intervals between taps in order to smooth the tempo fluctuation,
40 so that it would still be musical in character. Of the other two performances,
41 one uses accompaniment controlled by the B-Keeper system and the other
42 the same accompaniment but at a fixed tempo. The sequence in which these
43 three trials happen is randomly chosen by the computer and only revealed to
44 the participants after the test so that the experiment is *double-blind*, i.e. nei-
45 ther the researchers nor the drummer know which accompaniment is which.
46 Hence, the quantitative results gained by asking for opinion measures and
47 performance ratings should be free from any bias.
48
49
50
51
52
53

54 We are interested in the interaction between the drummer and the accom-
55 paniment which takes place through the machine. In particular, we wish to
56
57
58

know how this differs from the interaction that takes place with the human beat tracker. We might expect that, if our beat tracker is functioning well, the B-Keeper trials would be ‘better’ or ‘reasonably like’ those controlled by the human tapper. We would also expect them to be ‘not like a metronome’ and hence, distinguishable from the Steady Tempo trials.

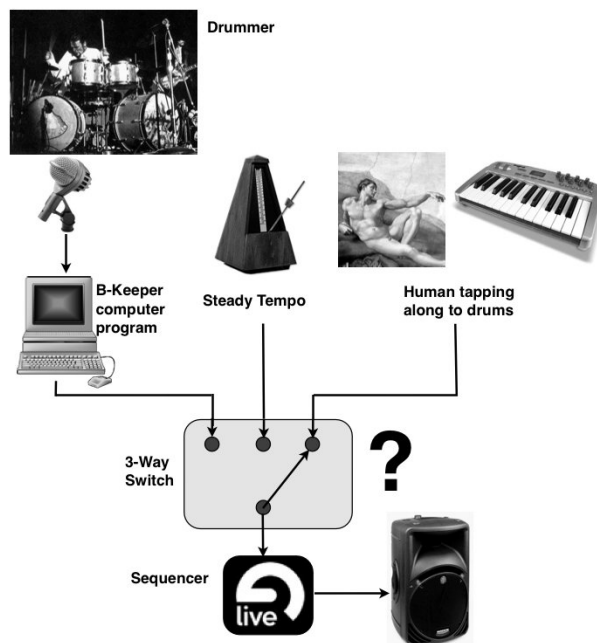


Fig. 5. Design set-up for the experiment. Three possibilities: (a) Computer controls tempo from drum input; (b) Steady Tempo; (c) Human controls tempo by tapping beat on keyboard

We carried out the experiment with eleven professional and semi-professional drummers. All tests took place in an acoustically isolated studio space. Each drummer took the test (consisting of the three randomly-selected trials) twice, playing to two different accompaniments. The first was based on a dance-rock piece first performed at Live Algorithms for Music Conference, 2006, which can be viewed on the internet². The second piece was a simple chord progression on a software version of a Fender Rhodes keyboard with some additional percussive sounds. The sequencer used was Ableton Live³, chosen for its time-stretching capabilities.

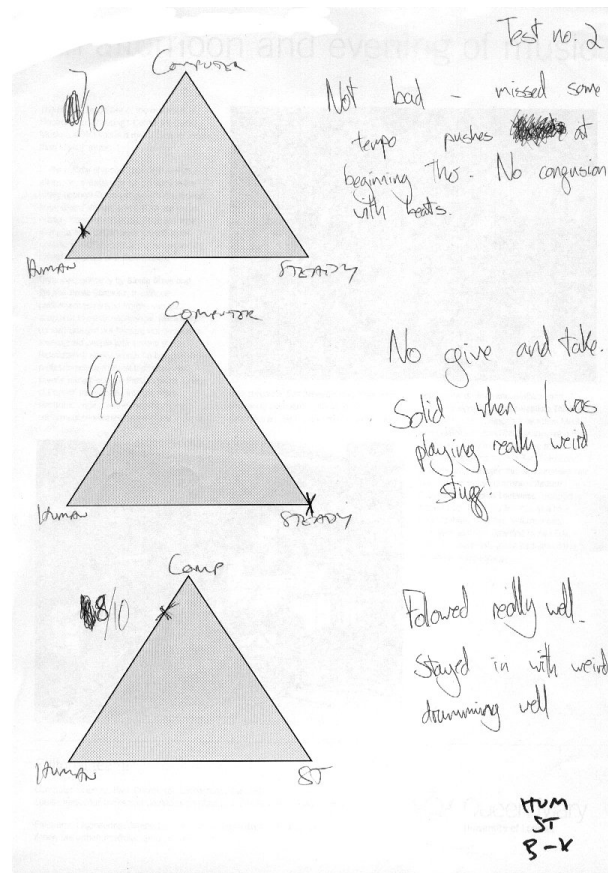
In the classic Turing Test, there would only be two possibilities: the human or the machine. However, since we wish to also contrast the beat tracker against a metronome as a control, we required a three-way choice. After each trial, we asked each drummer to mark an ‘X’ on an equilateral triangle which would indicate the strength of their belief as to which of the three systems was responsible. The three corners corresponded to the three choices and the

² <http://www.elec.qmul.ac.uk/digitalmusic/b-keeper>

³ <http://www.ableton.com>

1 nearer to a particular corner they placed the 'X', the stronger their belief
 2 that that was the tempo-controller for that particular trial. Hence, if an 'X'
 3 was placed on a corner, it would indicate certainty that that was the scenario
 4 responsible. An 'X' on an edge would indicate confusion between the two
 5 nearest corners, whilst an 'X' in the middle indicates confusion between all
 6 three. This allowed us to quantify an opinion measure for identification over all
 7 the trials. The human tapper (AR) and an independent observer also marked
 8 their interpretation of the trial in the same manner.
 9

10
 11 In addition, each participant marked the trial on a scale of one to ten as
 12 an indication of how well they believed that test worked as 'an interactive
 13 system'. They were also asked to make comments and give reasons for their
 14 choice. A sample sheet from one of the drummers is shown in Figure 6.
 15



16
 17
 18
 19
 20
 21
 22
 23
 24
 25
 26
 27
 28
 29
 30
 31
 32
 33
 34
 35
 36
 37
 38
 39
 40
 41
 42
 43
 44
 45
 46
 47
 48
 49
 50
 51
 52
 53
 54
 55
 56
 57
 58
 59
 60
 61
 62
 63
 64
 65

Fig. 6. Sample sheet filled in by a drummer.

3.2 Results

The participants' difficulty in distinguishing between controllers was a common feature of many tests and, whilst the test had been designed expecting that this might be the case, the results often surprised the participants

1 when revealed, with both drummers and the researchers being mistaken in
2 their identification of the controller. We shall contrast the results between
3 all three tests, particularly with regard to establishing the difference between
4 the B-Keeper trials and the Human Tapper trials and comparing this to the
5 difference between the Steady Tempo and Human Tapper trials. In Figure 7,
6 we can see the opinion measures for all drummers placed together on a single
7 triangle. The corners represent the three possible scenarios: B-Keeper, Human
8 Tapper and Steady Tempo with their respective symbols. Each 'X' has been
9 replaced with a symbol corresponding to the actual scenario in that trial. In
10 the diagram we can clearly observe two things:
11

- 12 • There is more visual separation between the Steady Tempo trials than the
13 other two. With the exception of a relatively small number of outliers, many
14 of the steady tempo trials were correctly placed near the appropriate corner.
15 Hence, if the trial is actually steady then it will probably be identified as
16 such.
17
- 18 • The B-Keeper and Human Tapper trials tend to be spread over an area
19 centered around the edge between their respective corners. At best, approx-
20 imately half of these trials have been correctly identified. The distribution
21 does not seem to have the kind of separation seen for the Steady Tempo
22 trials, suggesting that they have difficulty telling the two controllers apart,
23 but could tell that the tempo had varied.
24
25
26
27
28
29
30

31 *3.2.1 Analysis and Interpretation*

32
33 The mean scores recorded by all drummers are given in the first rows of Table
34 2. They show similar measures for correctly identifying the B-Keeper and Hu-
35 man Tapper trials: both have mean scores of 44%, with the confusion being
36 predominantly between which of the two variable tempo controllers is operat-
37 ing. The Steady Tempo trials have a higher tendency to be correctly identified,
38 with a score of 64% on the triangle.
39
40

41
42 Each participant in the experiment had a higher score for identifying the
43 Steady Tempo trials than the other two. It appears that the Human Tapper
44 trials are the least identifiable of the three and the confusion tends to be
45 between the B-Keeper and the Human Tapper.
46
47

48 For analysis purposes, we can express the opinion measures from Figure 7
49 as polarised decisions, by taking the nearest corner to be the participant's
50 decision for that trial. In the case of points equidistant from corners, we split
51 the decision equally. Table 3 shows the polarised decisions made by drummers
52 over the trials. There is confusion between the B-Keeper and Human Tapper
53 trials, whereas the Steady Tempo trials were identified over 70% of the time.
54 The B-Keeper and Human Tapper trials were identified 43% and 45% of the
55
56
57
58
59

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

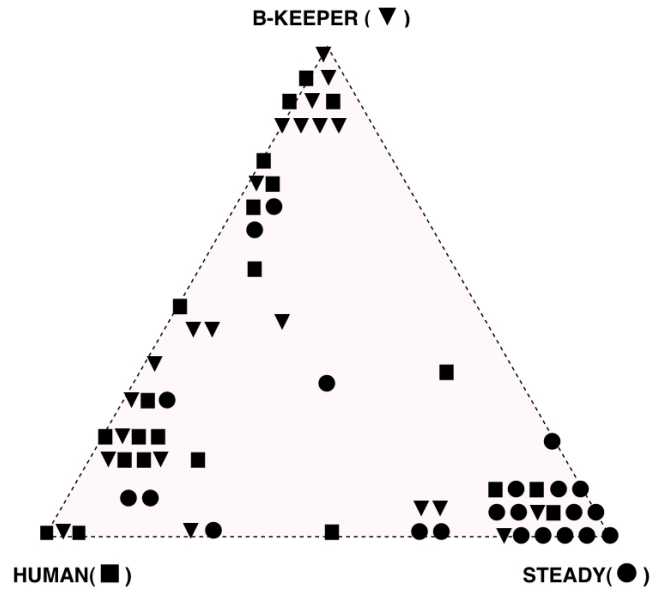


Fig. 7. Results illustrating where the eleven different drummers judged the three different accompaniments (B-Keeper, Human Tapper and Steady Tempo) in the test. The symbol used indicates which accompaniment it actually was (see corners). Where the participants have marked many trials in the same spot, as happens in the corners corresponding to Steady Tempo and B-Keeper, we have moved the symbols slightly for clarity. Hence, a small number of symbols are not exactly where they were placed. The raw data is available in co-ordinate form online (see footnote 1).

Judge	Accompn.t	Judged as:		
		B-Keeper	Human	Steady
Drummer	B-Keeper	44 %	37 %	18 %
	Human	38 %	44 %	17 %
	Steady	12 %	23 %	64 %
Human Tapper	B-Keeper	59 %	31 %	13 %
	Human	36 %	45 %	23 %
	Steady	15 %	17 %	68 %
Observer	B-Keeper	55 %	39 %	6 %
	Human	33 %	42 %	24 %
	Steady	17 %	11 %	73 %

Table 2
Mean Identification measure results for all judges involved in the experiment. Bold percentages correspond to the correct identification

time respectively – little better than the 33% we would expect by random choice.

Controller	Judged as:		
	B-Keeper	Human	Steady
B-Keeper	9.5	8.5	4
Human Tapper	8	10	4
Steady Tempo	2	4	16

Table 3

Polarised decisions made by the drummer for the different trials.

Controller	Judged as:	
	Human Tapper	Steady Tempo
Human Tapper	12	4
Steady Tempo	5	14

Table 4

Polarised decisions made by the drummer over the Steady Tempo and Human Tapper trials.

3.2.2 Comparative Tests

In order to test the distinguishability of one controller from the other, we performed a Chi-Square Test, calculated over all trials with either of the two controllers. If there is a difference in scores so that one controller is preferred to the other (above a suitable low threshold), then that controller is considered to be chosen for that trial. Where no clear preference was evident, such as in the case of a tie or neither controller having a high score, we discard the trial for the purposes of the test.

Thus, for any two controllers, we can construct a table of which decisions were correct. The table for comparisons between the Steady Tempo and the Human Tapper trials is shown in Table 4. We test against the null hypothesis that the distribution is the same for either controller, corresponding to the premise that the controllers are indistinguishable.

The separation between Steady Tempo and Human Tapper trials is significant ($\chi^2(3, 22) = 8.24, p < 0.05$), meaning participants could reliably distinguish them. Partly this might be explained from the fact that drummers could vary the tempo with the Human Tapper controller but the Steady Tempo trials had the characteristic of being metronomic.

Comparing the B-Keeper trials and the Human Tapper trials, we get the results shown in table 5. No significant difference is found in the drummers' identification of the controller for either trial ($\chi^2(3, 22) = 0.03, p > 0.5$). Whilst B-Keeper shares the characteristic of having variable tempo and thus is not identifiable simply by trying to detect a tempo change, we would expect

Controller	Judged as:	
	Human Tapper	B-Keeper
Human Tapper	9	8
B-Keeper	8	8

Table 5

Table contrasting decisions made by the drummer over the B-Keeper and Human Tapper trials.

that if there was a *machine-like* characteristic to the B-Keeper’s response, such as an unnatural response or unreliability in following tempo fluctuation, syncopation and drum fills, then the drummer would be able to identify the machine. It appeared that, generally, there was no such characteristic and drummers had difficulty deciding between the two controllers.

From the above, we feel able to conclude that the B-Keeper performs in a satisfactorily human-like manner in this situation.

3.2.3 Ratings

In addition to the identification of the controller for each trial, we also also asked each participant to rate each trial with respect to how well it had worked as an interactive accompaniment to the drums. The frequencies of ratings aggregated over all participants (drummers, human tapper and independent observer) are shown in Figure 8. The Steady Tempo accompaniment was consistently rated worse than the other two. The median values for each accompaniment are shown in Table 6. The B-Keeper system has generally been rated higher than both the Steady Tempo and the Human Tapper accompaniment.

The differences between the B-Keeper ratings and the others were analysed using the Wilcoxon signed-rank test (Mendenhall et al., 1989, section 15.4). These were found to be significant ($W = 198$ (Human Tapper) and $W = 218$ (Steady Tempo), $N = 22$, $p < 0.05$).

It is encouraging that not only did the beat tracker generally receive a high rating whether judged by the drummer or by an independent observer, but that its performance was sufficiently human-like to confuse participants as to which was the beat tracker and which the human tapper (section 3.2.2). This suggests that musically the beat tracker is performing its task well.

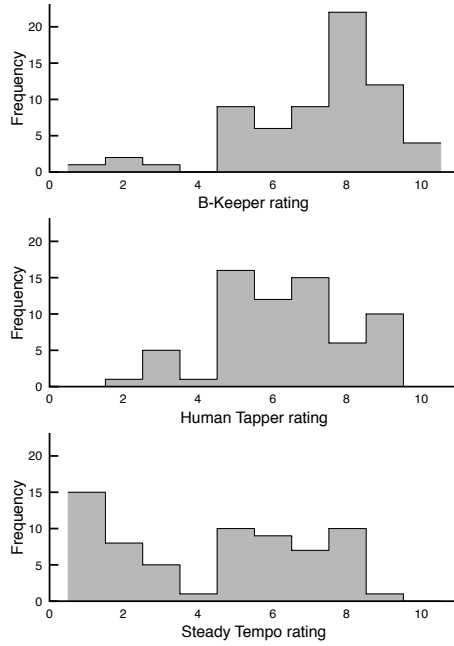


Fig. 8. Frequencies of ratings for the three scenarios: B-Keeper (upper), Human Tapper (middle) and Steady Tempo (lower).

Judge	Median Rating		
	B-Keeper	Human Tapper	Steady Tempo
Drummer	7.5	5.5	5
Human	8	6.5	4
Observer	8	7	5
Combined	8	6	5

Table 6
Median ratings given by all participants for the different scenarios.

4 Discussion

The two evaluation methods described above are designed to evaluate live interactive musical systems, without reducing the musical interaction to unrealistically simple tasks. In the case studies described, we have seen some of the possibilities afforded by the methods.

Firstly, the Discourse Analysis (DA) method can extract a detailed reconstruction of users' conceptualisation of a system. Our investigation of a voice-controlled interface provides us with interesting detail on the interaction between such concepts as controllability and randomness in the use of the interface, and the different ways of construing the interface itself. These findings would be difficult to obtain by other methods such as observation or question-

naire.

1
2
3 However, we see evidence that the discourses obtained are influenced by the
4 experimental context: the solo sessions, structured with tasks in using both
5 variants of our interface, produced discourse directly related to the interface;
6 while the group session, less structured, produced wider-ranging discourse with
7 less content bearing directly on the interface. The order of presentation also
8 may have made a difference to the participants. It is clear that the design
9 of such studies requires a careful balance: experimental contexts should be
10 designed to encourage exploration of the interface itself, while taking care not
11 to “lead” participants in unduly influencing the categories and concepts they
12 might use to conceptualise a system.
13
14
15

16
17 Secondly, the musical Turing Test method can produce a quantitative result
18 on whether a system provides an interactive experience similar to that pro-
19 vided by a human – despite the fact that we cannot evaluate such similarity
20 directly. Our case study found that both participants and observers exhibited
21 significant confusion between the B-Keeper and the Human Tapper, but not
22 between the B-Keeper and the Steady Tempo. Preference ratings alone tell us
23 that the B-Keeper provides a satisfactory experience, but the confusion results
24 go further: they tell us that B-Keeper achieves its aim of synchronising a piece
25 of music with the tempo variations of a live drummer, in a manner similar to
26 that obtained if a human performs the synchronisation.
27
28
29
30
31

32 The musical Turing Test approach is of course limited to situations in which
33 a system is intended to emulate a human musician, or perhaps to emulate
34 some other system. It cannot be applied to the vocal timbre-mapping system
35 of section 2.3, since for that there is no reference against which to compare.
36 However, emulation of human abilities is not uncommon in the literature: for
37 example the Continuator (Pachet, 2003), designed to provide a naturalistic
38 “call and response” interaction with a keyboard player; or BBCut (Collins,
39 2006), designed to produce real-time “beat-slicing” effects like a Drill’n’bass
40 producer. A more general method such as our DA method could be used on
41 these systems, and could produce useful information about users’ cognitive
42 approach to the systems, perhaps even illuminating the extent of human-like
43 affordances. However, the musical Turing Test gives us a more precise analysis
44 of this specific facet of musical human-computer interaction, and for example
45 enables numerical comparison between two systems.
46
47
48
49
50
51

52 Having explored our two methods, we are in a position to compare and con-
53 trast them with approaches used by other investigators, and then to work
54 towards recommendations on the applicability of different methods to differ-
55 ent contexts.
56
57
58

4.1 Comparison with other approaches

1
2
3
4 A useful point of comparison is the approach due to Wanderley and Orio
5 (2002), involving user trials on “maximally simple” tasks followed by Likert-
6 scale feedback. As previously discussed, this approach raises issues of task
7 authenticity, and of the suitability of the Likert-style questionnaire. Indeed,
8 Kiefer et al. (2008) investigate the Wanderley & Orio approach, and find qual-
9 itative analysis of interview data to be more useful than quantitative data
10 about task accuracy. The Wanderley & Orio method may therefore only be
11 appropriate to cases in which the test population is large enough to draw
12 conclusions from Likert-scale data, and in which the musical interaction can
13 reasonably be reduced or separated into atomic tasks. We suggest the crossfad-
14 ing of records by a DJ as one specific example: it is a relatively simple musical
15 task that may be operationalised in this way, and has a large user-base. (We
16 do not wish to diminish the DJ’s art: there are creative and sophisticated
17 aspects to the use of turntables, which may not be reducible to atomic tasks.)
18
19
20
21

22
23 One advantage of the Wanderley & Orio method is that Likert-scale ques-
24 tionnaires are very quick to administer and analyse. In our case study of the
25 Discourse Analysis approach, the ratio of interview time to analysis time was
26 approximately 1:30 or 1:33, a ratio slightly higher than the ratio of 1:25–1:29
27 reported for observation analysis of video data (Barendregt et al., 2006). This
28 long analysis time implies practical limitations for large groups.
29
30

31
32 Our approaches (as well as that of Wanderley & Orio) are “retrospective”
33 methods, based on users’ self-reporting after the musical act. We have ar-
34 gued that concurrent verbal protocols and observation protocols are problem-
35 atic for experiments involving live musicianship. A third alternative, which
36 is worthy of further exploration, is to gather data via physiological measure-
37 ments. Mandryk and Atkins (2007) present an approach which aims to eval-
38 uate computer-game-playing contexts, by continuously monitoring four phys-
39 iological measures on computer-game players, and using fuzzy logic to infer
40 the players’ emotional state. Analogies between the computer-gaming context
41 and the music-making context suggest that this method could be adopted for
42 evaluating interactive music systems. However, there are some issues which
43 would need to be addressed:
44
45
46

- 47
48 • Most importantly, the inference from continuous physiological variables to
49 continuous emotional state requires more validation work before it can be
50 relied on for evaluation.
- 51
52 • The evaluative role of the inferred emotional state also needs clarification:
53 the mean of the *valence* (the emotional dimension running from happiness
54 to sadness) suggests one simple figure for evaluation, but this is unlikely to
55 be the whole story.
56
57
58

- Musical contexts may preclude certain measurements: the facial movements involved in singing or beatboxing would affect facial electromyography (Mandryk and Atkins, 2007), and the exertion involved in drumming will have a large effect on heart-rate. In such situations, the inference from measurement to emotional state will be completely obscured by the other factors affecting the measured values.

Another consideration regarding these physiological approaches is that the finely-calibrated equipment required may in some cases be costly.

We note that the literature, the present work included, is predominantly concerned with evaluating musical interactive systems from a performer-centred perspective. Other perspectives are possible: a composer-centred perspective (for composed works), or an audience-centred perspective. But the performer is the primary user of such systems, and unlike the audience, has access to both the intention and the act. In some situations it may be appropriate to perform e.g. audience-centred evaluation. Our methods can be adapted for use with audiences – indeed, the independent observer in our musical Turing Test case study takes the role of audience. However, for audience-centred evaluations it may be the case that other methods are appropriate, such as voting or questionnaire approaches for larger audiences.

A further aspect of evaluation focus is the difference between solo and group music-making. Wanderley & Orio’s set of simple musical tasks is only applicable for solo experiments. Our evaluation methods can apply in both solo and group situations, with the appropriate experimental tasks for participants. The physiological approach may also apply equally well in group situations.

4.2 Recommendations

From our studies, we suggest that an investigator wishing to formally evaluate an interactive music system, or live music interface, should consider the following:

- (1) **Is the system primarily designed to emulate the interaction provided by a human, or by some other known system?** If so, the musical Turing Test method can be recommended.
- (2) **Is the performer’s perspective sufficient for evaluation?** In many cases the answer to this is “yes”, although there may be cases in which it is considered important to design an experiment involving audience evaluation. Many of the same methods (interviews, questionnaires, Turing-Test choices) are applicable to audience members – and because audiences can often result in large sample sizes compared against performer populations, survey methods such as Likert scales are more likely to be appropriate.

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
- (3) **Is the system designed for complex musical interactions, or for simple/separable musical tasks?** If the latter, then Wanderley & Orio’s approach using simplified tasks may hold some attraction. If the former, then we recommend a more situated evaluation such as our Discourse Analysis approach, which avoids the need to reduce the musical interaction down to atomic tasks.
 - (4) **Is the system intended for solo interaction, or is a group interaction a better representation of its expected use pattern?** The experimental design should reflect this, using either solo or group sessions.
 - (5) **How large is the population of participants on which we can draw for evaluation?** Often the population will be fairly small, which raises issues for the statistical power of quantitative approaches. Qualitative approaches should then be considered.

19
20
21
22
23
24
25
26
27
28
29
30
31
32

One of the key themes in our recommendations is that the design of an evaluation experiment should aim as far as possible to reflect an authentic context for the use of the system. Experimental design should include a phase which encourages use and exploration of the system. Approaches such as our Discourse Analysis of interview data can then be applied in a wide variety of cases to probe the participants’ cognitive constructs produced during the experiment. Discourse Analysis is not the only way to analyse interview data (Silverman, 2006), and others may be worth pursuing; we have argued for Discourse Analysis as a principled approach which extracts a structured picture of the described world from a relatively small amount of interview data.

33
34
35
36
37
38
39
40
41
42
43

In any design using interview data, it is important that the facilitator is experienced in neutral interview technique, able to avoid “leading” participants in their choice of concepts and language. It is also important that the reporting of the experiment demonstrates the difference between formal and informal qualitative analysis: a formal qualitative analysis makes clear the route from data to conclusions, by describing the methodological basis and the steps taken to process the data, and ideally by publishing transcripts etc.

44
45
46
47
48
49
50

Approaches based on continuous physiological measures (Mandryk and Atkins, 2007) may become viable for evaluating interactive systems, although there are at present some issues to be resolved, discussed above. We consider this a topic for future research, rather than an approach to be generally recommended at present, although we look forward to developments in this area.

51
52
53
54
55
56
57
58

Finally, from our experience we repeat the advice given by others (Kiefer et al., 2008) that the importance of piloting should not be underestimated, as it can reveal issues with an experimental design that do not otherwise become apparent beforehand.

5 Conclusions

We have introduced two methods for evaluating interactive musical systems which allow for evaluation in the context of realistic musical interactions. Both our methods avoid problems with other methods proposed in the literature, in particular the oversimplification of musical tasks and the small populations typical of specific musical performer types.

The Discourse Analysis method is a general method which aims to characterise the conceptual structures created by participants while interacting with a system. In our case study with a voice timbre remapping interface, we found that the timbral aspect of the system was unproblematic for users, and we highlighted the nuanced interaction of control and randomness issues in the use of such a system.

The musical Turing Test method is applicable to the case where a system aims to emulate some aspect of human musical performance. We applied it to the evaluation of a real-time beat tracker (the B-Keeper) and found that although participants could reliably distinguish the steady-state system from the B-Keeper, they could not reliably distinguish the B-Keeper from a human tapper – in other words, the B-Keeper “passes” the musical Turing Test.

Furthermore, we have placed our methods in context with other methods in the literature, and derived recommendations (Section 4.2) for researchers wishing to evaluate an interactive musical system.

6 Acknowledgements

Dan Stowell and Andrew Robertson are supported by Doctoral Training Account research studentships from the EPSRC.

References

- Antaki, C., Billig, M., Edwards, D., Potter, J., 2004. Discourse analysis means doing analysis: A critique of six analytic shortcomings. *Discourse Analysis Online* 1 (1).
URL <http://extra.shu.ac.uk/daol/articles/v1/n1/a1/antaki2002002-paper.html>
- Banister, P., Burman, E., Parker, I., Taylor, M., Tindall, C., 1994. *Qualitative Methods in Psychology: A Research Guide*. Open University Press, Buckingham.

- 1 Barendregt, W., Bekker, M. M., Bouwhuis, D., Baauw, E., 2006. Identifying
2 usability and fun problems in a computer game during first use and after
3 some practice. *International Journal of Human-Computer Studies* 64 (9),
4 830–846, doi:10.1016/j.ijhcs.2006.03.004.
- 5 Buxton, W., Sniderman, R., 1980. Iteration in the design of the human-
6 computer interface. In: *Proceedings of the 13th Annual Meeting, Human*
7 *Factors Association of Canada*. pp. 72–81.
- 8 Card, S. K., English, W. K., Burr, B. J., 1978. Evaluation of mouse, rate-
9 controlled isometric joystick, step keys, and text keys for text selection on
10 a CRT. *Ergonomics* 21 (8), 601–613, doi:10.1080/00140137808931762.
- 11 Collins, N., 2006. BBCut2: Integrating beat tracking and on-the-
12 fly event analysis. *Journal of New Music Research* 35 (1), 63–70,
13 doi:10.1080/09298210600696600.
- 14 de Poli, G., 2004. Methodologies for expressiveness modelling of and for
15 music performance. *Journal of New Music Research* 33 (3), 189–202,
16 doi:10.1080/0929821042000317796.
- 17 d’Escrivan, J., Collins, N. (Eds.), 2007. *The Cambridge Companion to Elec-*
18 *tronic Music*. Cambridge University Press.
- 19 Dobrian, C., Koppelman, D., 2006. The ‘E’ in NIME: Musical expression with
20 new computer interfaces. In: *Proceedings of New Interfaces for Musical Ex-*
21 *pression (NIME)*. IRCAM, Centre Pompidou Paris, France, pp. 277–282.
22 URL http://www.nime.org/2006/proc/nime2006_277.pdf
- 23 Ericsson, K. A., Simon, H. A., 1996. *Protocol analysis: Verbal reports as data*
24 (Revised edition). Massachusetts Institute of Technology, Cambridge, MA.
- 25 Fels, S., 2004. Designing for intimacy: Creating new interfaces for
26 musical expression. *Proceedings of the IEEE* 92 (4), 672–685,
27 doi:10.1109/JPROC.2004.825887.
- 28 General Instrument, early 1980s. GI AY-3-8910 Programmable Sound Gener-
29 ator datasheet.
- 30 Göb, R., McCollin, C., Ramalhoto, M. F., 2007. Ordinal methodology in
31 the analysis of Likert scales. *Quality and Quantity* 41 (5), 601–626,
32 doi:10.1007/s11135-007-9089-z.
- 33 Goebel, W., September 2004. Computational models of expressive music perfor-
34 mance: The state of the art. *Journal of New Music Research* 33, 203–216(14),
35 doi:doi:10.1080/0929821042000317804.
- 36 Grant, S., Aitchison, T., Henderson, E., Christie, J., Zare, S., McMurray,
37 J., Dargie, H., 1999. A comparison of the reproducibility and the sensi-
38 tivity to change of visual analogue scales, Borg scales, and Likert scales
39 in normal subjects during submaximal exercise. *Chest* 116 (5), 1208–1217,
40 doi:10.1378/chest.116.5.1208.
- 41 Hunt, A., Wanderley, M. M., 2002. Mapping performer parameters to synthesis
42 engines. *Organised Sound* 7 (2), 97–108, doi:10.1017/S1355771802002030.
- 43 Kiefer, C., Collins, N., Fitzpatrick, G., 2008. HCI methodology for evaluating
44 musical controllers: A case study. In: *Proc. International Conference on New*
45 *Interfaces for Musical Expression (NIME)*. pp. 87–90.
- 46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 1 Large, E. W., 1995. Beat tracking with a nonlinear oscillator. In: Working
2 Notes of the IJCAI-95 Workshop on Artificial Intelligence and Music, Mon-
3 treal. pp. 24–31.
- 4 Lee, J. W., Jones, P. S., Mineyama, Y., Zhang, X. E., 2002. Cultural differences
5 in responses to a Likert scale. *Research in Nursing & Health* 25 (4), 295–306,
6 doi:10.1002/nur.10041.
- 7 Levitin, D. J., McAdams, S., Adams, R. L., 2003. Control parameters for
8 musical instruments: A foundation for new mappings of gesture to sound.
9 *Organised Sound* 7 (02), 171–189, doi:10.1017/S135577180200208X.
- 10 Mandryk, R. L., Atkins, M. S., 2007. A fuzzy physiological approach for
11 continuously modeling emotion during interaction with play technolo-
12 gies. *International Journal of Human-Computer Studies* 65 (4), 329–347,
13 doi:10.1016/j.ijhcs.2006.11.011.
- 14 McKinney, M. F., Moelants, D., Davies, M. E. P., Klapuri, A., 2007. Evaluation
15 of audio beat tracking and music tempo extraction algorithms. *Journal of*
16 *New Music Research* 36 (1), 1–16, doi:10.1080/09298210701653252.
- 17 Mendenhall, W., Wackerly, D. D., Scheaffer, R. L., 1989. *Mathematical statis-*
18 *tistics with applications*, 4th Edition. PWS-Kent.
- 19 Nicholls, M. E. R., Orr, C. A., Okubo, M., Loftus, A., 2006. Satisfaction guar-
20 anteed: The effect of spatial biases on responses to Likert scales. *Psycholog-*
21 *ical Science* 17 (12), 1027–1028, doi:10.1111/j.1467-9280.2006.01822.x.
- 22 Pachet, F., 2003. The Continuator: Musical interaction with style. *Journal of*
23 *New Music Research* 32 (3), 333–341, doi:10.1076/jnmr.32.3.333.16861.
- 24 Peretz, I., Zatorre, R. J., 2005. Brain organization for mu-
25 sic processing. *Annual Review of Psychology* 56 (1), 89–114,
26 doi:10.1146/annurev.psych.56.091103.070225.
- 27 Polfreman, R., 2001. A task analysis of music composition and its applica-
28 tion to the development of Modalyser. *Organised Sound* 4 (01), 31–43,
29 doi:10.1017/S1355771899001053.
- 30 Robertson, A., Plumbley, M. D., 2007. B-Keeper: A beat-tracker for live per-
31 formance. In: *Proc. International Conference on New Interfaces for Musical*
32 *Expression (NIME)*, New York, USA. pp. 234–237.
- 33 Salamé, P., Baddeley, A., 1989. Effects of background music on phonologi-
34 cal short-term memory. *The Quarterly Journal of Experimental Psychology*
35 *Section A* 41 (1), 107–122, doi:10.1080/14640748908402355.
- 36 Searle, J., 1980. Minds, brains and programs. *Behavioural and Brain Sciences*
37 3, 417–457.
- 38 Silverman, D., 2006. *Interpreting Qualitative Data: Methods for Analysing*
39 *Talk, Text and Interaction*, 2nd Edition. Sage Publications Inc.
- 40 Stewart, D. W., 2007. *Focus groups: Theory and practice*. SAGE Publications.
- 41 Stowell, D., Plumbley, M. D., July 2007. Pitch-aware real-time timbral remap-
42 ping. In: *Proceedings of the Digital Music Research Network (DMRN)*
43 *Summer Conference*.
- 44 URL [http://www.elec.qmul.ac.uk/digitalmusic/papers/2007/
45 StowellPlumbley07-dmrn.pdf](http://www.elec.qmul.ac.uk/digitalmusic/papers/2007/StowellPlumbley07-dmrn.pdf)
- 46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 Toiviainen, P., 1998. An interactive MIDI accompanist. *Computer Music Journal* 22 (4), 63–75.
2
3 Turing, A., 1950. Computing machinery and intelligence. *Mind* 59, 433–460,
4 doi:10.1007/978-1-4020-6710-5.
5 Uszkoreit, H., 1996. Discourse and dialogue. In: *Survey of the State of the*
6 *Art in Human Language Technology*. R. A. Cole, J. Mariani, H. Uszkoreit,
7 A. Zaenen, and V. Zue, eds., Center for Spoken Language Understanding,
8 Oregon Health and Science University, Ch. 6.
9 URL <http://cslu.cse.ogi.edu/HLTsurvey/ch6node2.html>
10
11 Wanderley, M. M., Orio, N., 2002. Evaluation of input devices for musical
12 expression: Borrowing tools from HCI. *Computer Music Journal* 26 (3),
13 62–76, doi:10.1162/014892602320582981.
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65