



Deposited via The University of Sheffield.

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/id/eprint/187303/>

Version: Published Version

---

**Article:**

Holzappel, A., Benetos, E., Killick, A. et al. (2021) Humanities and engineering perspectives on music transcription. *Digital Scholarship in the Humanities*, 37 (3). pp. 747-764. ISSN: 2055-7671

<https://doi.org/10.1093/lc/fqab074>

---

**Reuse**

This article is distributed under the terms of the Creative Commons Attribution-NonCommercial (CC BY-NC) licence. This licence allows you to remix, tweak, and build upon this work non-commercially, and any new works must also acknowledge the authors and be non-commercial. You don't have to license any derivative works on the same terms. More information and the full terms of the licence here:

<https://creativecommons.org/licenses/>

**Takedown**

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing [eprints@whiterose.ac.uk](mailto:eprints@whiterose.ac.uk) including the URL of the record and the reason for the withdrawal request.

# Humanities and engineering perspectives on music transcription

---

Andre Holzapfel 

Division of Media Technology and Interaction Design, KTH Royal Institute of Technology, Stockholm, Sweden

Emmanouil Benetos 

School of Electronic Engineering and Computer Science, Queen Mary University of London, London, UK

Andrew Killick

Department of Music, University of Sheffield, Sheffield, UK

Richard Widdess

Department of Music, SOAS University of London, UK

---

## Abstract

Music transcription is a process of creating a notation of musical sounds. It has been used as a basis for the analysis of music from a wide variety of cultures. Recent decades have seen an increasing amount of engineering research within the field of Music Information Retrieval that aims at automatically obtaining music transcriptions in Western staff notation. However, such approaches are not widely applied in research in ethnomusicology. This article aims to bridge interdisciplinary gaps by identifying aspects of proximity and divergence between the two fields. As part of our study, we collected manual transcriptions of traditional dance tune recordings by eighteen transcribers. Our method employs a combination of expert and computational evaluation of these transcriptions. This enables us to investigate the limitations of automatic music transcription (AMT) methods and computational transcription metrics that have been proposed for their evaluation. Based on these findings, we discuss promising avenues to make AMT more useful for studies in the Humanities. These are, first, assessing the quality of a transcription based on an analytic purpose; secondly, developing AMT approaches that are able to learn conventions concerning the transcription of a specific style; thirdly, a focus on novice transcribers as users of AMT systems; and, finally, considering target notation systems different from Western staff notation.

### Correspondence:

Andre Holzapfel, KTH Royal Institute of Technology, Stockholm, Sweden.

### E-mail:

holzap@kth.se

---

## 1 Introduction

Music transcription is a process of creating a notation of musical sounds, with music notation being the

representation of musical sound through some other medium. This can take many forms, but this article, and the field of Music Information Retrieval (MIR) in general, is concerned with static, visual representations

of sound, and specifically with Western staff notation. When performers use notation, it serves them as a set of instructions for producing certain sound patterns, but in music transcription, the notation is produced from sound (usually from a recorded performance) rather than the reverse. Transcriptions may be made either by human transcribers or by machines such as computers, and may be human readable (e.g. staff notation, guitar tab), and/or machine readable (e.g. piano roll, MIDI (Musical Instrument Digital Interface) file). If human readable, they may be used as instructions for performers, but in ethnomusicology they are more often used to support analysis of music that was not previously notated (or not notated in a way that is useful for the analysis).

In ethnomusicology and its parent discipline, comparative musicology, transcription long played a central role (Ellingson, 1992). Initially, it was the only means of communicating the sound and style of the music to readers who had never heard it. Over time, recordings took over this function, and transcription became more restricted to its other role, that of supporting and illustrating analyses (Nettl, 2015, p. 75). As the comparative musicologists' agenda of arranging all music into global evolutionary schemes gave way to the more locally situated, fieldwork-based inquiries of ethnomusicologists, the role and focus of music analysis became more variable. For some, the whole project of transcription and analysis, and especially transcription into Western notation, felt uncomfortably close to the colonialist legacy that ethnomusicology was trying to leave behind (see Marian-Bălașa, 2005). For others, transcription and analysis still had a part to play in characterizing particular forms of music and thus helping to answer ethnomusicology's big question of why people make and use the particular forms of music that they do. But those who still practised transcription and analysis used it to address specific analytical questions relating to the particular music they were concerned with, and developed approaches and solutions that were tailored to these questions and often quite personal to the transcriber. Few were content with staff notation as used by classical composers: even the early comparative musicologists had recognized that modifications were necessary for transcribing music from outside that tradition (Abraham and von Hornbostel, 1994). Some turned to technologies for producing automatic transcriptions

that transcended the limitations of staff notation and human listeners (Metfessel, 1928; Seeger, 1958).

Research on computational—as opposed to more broadly technological—methodologies for music transcription emerged in the 1970s (e.g. Moorer, 1975). Research on the topic has intensified during the last two decades in the context of MIR research, resulting in development and increasing refinement of methods for automatic music transcription (AMT). Most of these methods assume Western staff notation as the final goal of the transcription process, and the related problems of this notation format have not been discussed to a large extent within MIR. Besides AMT approaches, various metrics have been proposed within MIR that aim at the evaluation of the quality of AMT-produced transcriptions based on comparison. So far, ethnomusicologists did not make much use of AMT methods or evaluation metrics. Thus, a central question for this article is whether, and how, AMT might be made more useful to ethnomusicologists. The potential may lie within, for instance, the support of inexperienced transcribers, the discovery of melodic motives in large corpora, or the visualization of longer recordings in the form of notation. We believe that an increased awareness of MIR research among ethnomusicologists, and conversely of problems long discussed in ethnomusicology among MIR researchers, will help to provide answers to our central question and lead to an improved and more widely applicable AMT technology.

We therefore approach our central question through a combination of perspectives. In the following section, we contrast perspectives on transcription in MIR and in ethnomusicology, in order to identify aspects of proximity and divergence between the fields. In Section 3, we present our method that, first, extends a previous user study (Holzapfel and Benetos, 2019), which collected a large number of manual transcriptions for a collection of traditional dance tune recordings. Secondly, our method employs a combination of expert and computational evaluation of these transcriptions. This enables us to investigate the limitations of computational transcription metrics and AMT methods in Section 4. In Section 5, we discuss promising avenues to make automatic transcription more useful for music studies in the Humanities.





evaluated using the benchmark metrics proposed as part of the MIR Evaluation eXchange public evaluation campaigns (Bay *et al.*, 2009). Informal observations have been made on how the evaluation of systems with respect to producing lists of notes (typically referred to as note-based evaluation) is more perceptually relevant compared with the evaluation of groups of pitches over small time segments (typically referred to as frame-based evaluation). However, community efforts towards proposing evaluation metrics for AMT which are linked with how humans would judge transcriptions are fairly limited and have not reached a broad consensus. An early perceptual study by Daniel *et al.* (2008) proposed evaluation metrics that take some common local errors related to automatic transcription into account (such as note insertions and octave errors), but do not take into account aspects related to metre or tonality. Recently, for the more relevant task of complete transcription, Nakamura *et al.* (2018) proposed a set of evaluation metrics that address local errors in a musical score (e.g. insertions, deletions). Higher-level evaluation metrics have also been proposed which draw knowledge from music theory and focus on typesetting (Cogliati and Duan, 2017; McLeod and Steedman, 2018), although their links with human assessment of music transcription are unclear. Finally, Ycart *et al.* (2020) proposed a single evaluation metric for AMT systems that produce outputs in physical time following perceptual evaluations; the focus of the study was however only on piano music, and the metric's ability to generalize to other instruments is as yet unclear.

An important issue of AMT methods refers to their implicit or explicit algorithmic biases. The development of AMT methods using supervised machine learning methods (which are the most commonly used methodologies nowadays) assumes the presence of a 'target' or a 'reference' transcription that the transcription system should try to approximate. This, however, can bias any resulting systems to music recordings for which notation exists—most commonly pieces that have been composed in written form. A second point that is also linked to the previous one refers to the availability of data to train AMT systems. The availability of music recordings with corresponding annotations of notes and musical events over physical time is scarce since the process of producing such annotations is an extremely laborious task (Su and Yang, 2015). This has

led to the creation of datasets for AMT research that have been created using acoustic musical instruments that can automatically create such annotations—most commonly acoustic pianos with sensors that can capture music performance characteristics, such as Yamaha or Bösendorfer Disklavier pianos. This has led to a large imbalance in the diversity of AMT datasets towards the piano, e.g. through the commonly used MAPS (Emiya *et al.*, 2010) and MAESTRO (Hawthorne *et al.*, 2019) datasets, and therefore led to the creation of AMT methods that are focused on piano transcription. Combined with the greater availability of reference scores for Western art music compared with other musical cultures or styles, this has led to the vast majority of AMT systems being exclusively focused on transcribing Western art music performed on the piano. This stylistic focus leads to the fact that the analytic purpose of a transcription has largely not been taken into account in MIR until now.

On the value of AMT methods for studies in musicology, in recent years, studies have been attempted in the intersection of digital and computational musicology and MIR, mostly aiding musical study for large corpora. For example, Tidhar *et al.* (2014) used AMT methods as a first step towards a large-scale analysis of temperament profiles and temperament trends over time in harpsichord recordings. Similarly, AMT methods have been used towards estimating trends in tuning frequencies in the context of archival Western art music recordings (Abdallah *et al.*, 2017). In the context of jazz studies, melody estimation algorithms have been used to group melodic patterns in jazz improvisation (Höger *et al.*, 2019). In the field of computational ethnomusicology, AMT methods have been used in the context of Turkish makam music (Benetos and Holzapfel, 2015), where a discrepancy between music theory and practice was observed with respect to the pitch values implied by notation. Finally, melody estimation methods have been used to infer music similarity and dissimilarity in folk and traditional music recordings in a global sample (Panteli *et al.*, 2017).

### 2.3 Music transcription in between the fields

In ethnomusicology, mechanical devices for automatic transcription and analysis have been used since before the advent of computer technology (Ellingson,

















**Fig. 3** Example transcription of segment 2 (bottom staff) with a large divergence between rating by experts (A.K.: 3, R.W.: 4) and computational metric (9.1). Transcriptions by A.K. and R.W. are depicted in the upper two staves. The pick-up measure was added to the transcriptions of A.K. and R.W. for alignment purposes. Dashed boxes denote mistakes that the experts specified as motivation for their low rating



**Fig. 4** Example transcription of segment 6 (bottom staff) with a large divergence between rating by experts (5) and computational metric (9.7). Transcriptions by A.K. and R.W. are depicted in the upper two staves. Dashed boxes denote mistakes that the experts specified as motivation for their low rating

rhythm and pitch contours more complex; better results might be produced by providing tighter constraints regarding rhythm and pitches as an additional learning stage for the algorithms based on a corpus of example transcriptions.

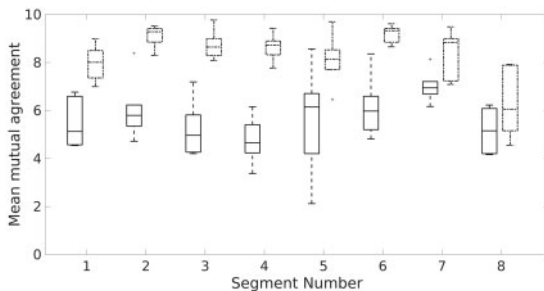
#### 4.5 Agreement between transcriptions

The analysis in this article is able to profit from the availability of expert ratings, and obtains additional insight by specific expert assessment of individual transcriptions. In the absence of such expert evaluation and established reference transcriptions, it may be of advantage to estimate the quality of a group of

transcriptions automatically. The question is whether we can automatically identify a set of good transcriptions, based on their mutual agreement. As a first step in this direction, we investigate how the mutual agreement among the  $N$  best transcriptions compares with the mutual agreement among the  $N$  lowest rated. In order to compute the mutual agreement in a group of transcriptions, we employ the combined metric depicted in Table 2 between all transcriptions in a group, and compute the mean of the obtained values. Our comparison demonstrates that the  $N$  best transcriptions agree mutually more than  $N$  lowest rated transcriptions (Fig. 5). This implies that high-quality transcriptions tend to agree more in their basic pitch,

note onset, and duration characteristics than low-quality transcriptions. The difference is significant over the whole dataset, and only for segment 8, an overlap can be observed. This is the segment with the generally lowest quality ratings by the experts, and it was rated as the most difficult to transcribe by the eighteen participants in our user study. Stylistically, it is highly idiosyncratic compared with the other segments, characterizing it as a special case among our eight segments. Figure 6 depicts the two transcriptions with the highest average rating for this segment, in which relatively large differences in the interpretation of both pitch and rhythm are apparent.

In a real-world scenario, quality ratings on the level of individual transcriptions will not be available, and shaping a reference committee of a group of transcribers that are assumed to have high expertise may be a viable alternative. We evaluated such an alternative and were able to confirm that the mutual agreement among a group of transcribers is strongly correlated



**Fig. 5** Mutual agreement among the highest rated (dashed-dotted line boxes) and the lowest rated (solid line boxes). The combination of two metrics found to correlate most with the expert rating was used to mutually compare each group of transcriptions



**Fig. 6** Transcriptions of segment 8 that received the highest average ratings by the experts. The tempo of the transcribed segment is about 120 bpm

with the average expert rating of that group of transcribers. This implies that the mean mutual agreement among a group of transcribers or transcriptions can be used as an indicator for the choice of reference transcriptions in a corpus. It remains to be explored if the findings in our case study generalize to other musical styles and analytic tasks.

## 5 Conclusion

By comparing ratings of human experts with computational metrics through corpus and close analysis, we documented differences in how the quality of a transcription is assessed in ethnomusicology and in MIR. We revealed several aspects that the metrics seem to be ‘missing’ in Section 4.3. Computational metrics are only partially correlated with human ratings. Specifically, the highest correlation between metrics and human ratings can be found for metrics that focus on onset times and pitch detection errors, and computational metrics are less sensitive to rhythmic problems compared with experts. An important methodological aspect that is shared is the assessment procedure, which is based on comparison with a reference transcription, which indicates that the MIR procedure is not substantially wrong. A conceptual aspect that is missing is the consideration of the analytic purpose, which importantly guides the shape of the reference transcription for evaluation. In applications where such purpose is not clearly stated—such as the development of generic transcription tools in MIR—we recommend to use more than one reference transcription to cover a range of such purposes. To identify such a group of references in a larger corpus, the mean mutual agreement among a group of







