

Artificial Intelligence in the Concertgebouw

Andreas Arzt^(1,2), Harald Frostel⁽¹⁾, Thassilo Gadermaier⁽²⁾

Martin Gasser⁽²⁾, Maarten Grachten⁽²⁾, Gerhard Widmer^(1,2)

⁽¹⁾Johannes Kepler University, Linz, Austria

⁽²⁾Austrian Research Institute for Artificial Intelligence, Vienna, Austria

andreas.arzt@jku.at

Abstract

In this paper we present a real-world application (the first of its kind) of machine listening in the context of a live concert in a world-famous concert hall – the Concertgebouw in Amsterdam. A real-time music tracking algorithm listens to the Royal Concertgebouw Orchestra performing Richard Strauss’ *Alpensinfonie* and follows the progress in the sheet music, i.e., continuously tracks the most likely position of the live music in the printed score. This information, in turn, is used to enrich the concert experience for members of the audience by streaming synchronised visual content (the sheet music, explanatory text and videos) onto tablet computers in the concert hall. The main focus of this paper is on the challenges involved in tracking live orchestral music, i.e., how to deal with heavily polyphonic music, how to prepare the data needed, and how to achieve the necessary robustness and precision.

1 Introduction

Real-time music listening is a big challenge for machine perception and AI. While ‘listening’ is a broad concept, involving various aspects of structure comprehension and abstraction (e.g., perceiving and tracking beat and tempo, identifying the melody, recognising voices, instruments, style, genre, etc.) – all of this is the domain of the Music Information Retrieval (MIR) research field –, even the more specialised task of listening to a live performance and synchronously reading along in the printed music score (e.g., for the purpose of page turning) is a big challenge [Arzt *et al.*, 2008]. This task is generally called *score following* or *music tracking*. What it involves is listening to a live incoming audio stream, extracting higher-level features from the raw audio that somehow capture aspects of the ‘sound’ of the current moment, and tracking the most likely position in the score that the sound seems to correspond to – regardless of the specific tempo chosen by the musicians on stage, of continuous or abrupt tempo changes due to expressive timing, and robust to varying sound quality and instrument sounds.

Real-time music tracking, which started in the 1980s (see [Dannenberg, 1984; Vercoe, 1984]), has attracted quite some research in recent years [Raphael, 2010; Cont, 2009;

Arzt *et al.*, 2008; Korzeniowski *et al.*, 2013]. While there still are many open research questions (such as on-line learning of predictive tempo models during a performance), real-time score following is already beginning to be used in real-world applications. Examples include Antescofo¹, which is actively used by professional musicians to synchronise a performance (mostly solo instruments or small ensembles) with computer realised elements, and Tonara², a music tracking application focusing on the amateur pianist and running on the iPad.

In this paper, we lift the problem to a new level of complexity: we wish to track an entire orchestra playing complex polyphonic music. This presents specific challenges to a music tracking algorithm. First and foremost, it has to deal with heavily polyphonic music, with many different instruments – and sometimes very unusual ones. Furthermore, a long piece like a symphony challenges a music tracking algorithm with many different situations (e.g. very soft and quiet sections, immensely powerful, dynamic and fast parts, solo sections for different instruments). The music tracking algorithm has to cope with all these situations and has to be able to track all the instruments, be it a trumpet, a violin or an organ. All of this is done live, in contrast to studies like e.g. [Miron *et al.*, 2014], where orchestral recordings are aligned to a score off-line in a non-causal way. Tracking can only work if a good representation of the score of the underlying piece is provided. While this process is relatively straightforward for e.g. a solo piano piece, it is far from trivial for a complicated classical symphony.

In addition to describing how we solved these challenges, we report on the first public live demonstration of our system in a regular concert in a famous concert hall, and provide a quantitative evaluation of the precision and robustness of our algorithm in solving this task.

The paper is organised as follows. Section 2 explains the general idea, context, and setup of the experiment. In Section 3 the requirements on the internal score representation and the audio features are discussed. Sections 4 and 5 give a description of the tracking algorithm. In Section 6 we present a qualitative analysis of the tracking during the live concert, and Section 7 gives a detailed quantitative evaluation.

¹repmus.ircam.fr/antescofo

²tonara.com

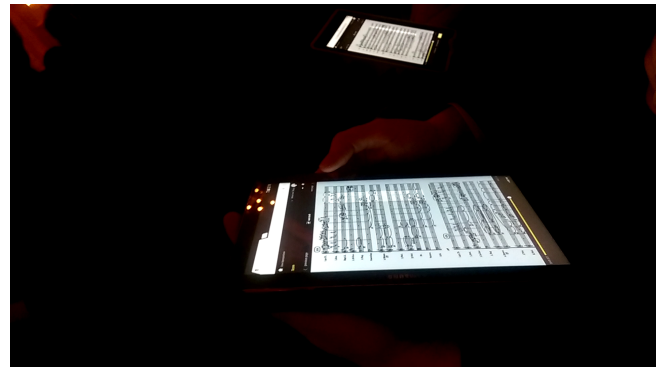
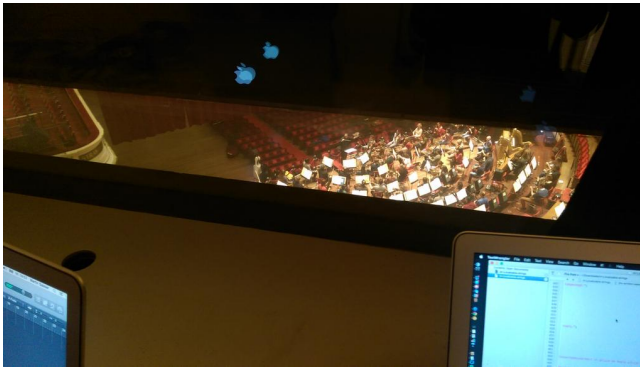


Figure 1: Left: View from the control room onto the stage (during orchestra rehearsal); right: synchronised score display in the audience during the concert.

2 The Challenge: Tracking the Concertgebouw Orchestra

The multi-national European research project PHENICX³ provided us with the unique opportunity (and challenge) to demonstrate our score following technology in the context of a big, real-life symphonic concert. The general goal of the project is to develop technologies that enrich the experience of classical music concerts. In the experiment to be described, this was done by using the live performance tracker to control, in real time and via WiFi, the transmission and display of additional visual and textual information, synchronised to the live performance on stage. The user interface and the visualisations were provided by our project partner Videodock⁴.

The event took place on February 7th, 2015, in the Concertgebouw in Amsterdam. The Royal Concertgebouw Orchestra, conducted by Semyon Bychkov, performed the *Alpensinfonie* (Alpine Symphony) by Richard Strauss. This concert was part of a series called ‘Essentials’, during which technology developed within the project can be tested in a real-life concert environment. All the tests during this concert series have to be as non-invasive as possible. For the demonstration during the concert in question, a test audience of about 30 people was provided with tablet computers and placed in the rear part of the concert hall.

The setup was as follows. Two microphones were placed a few meters above the conductor, in an AB-setup, picking up the music, but also a lot of noise, e.g. coughing in the audience and noise made by orchestra members, and a fair amount of reverberation from the hall. In a control room behind the scenes a regular consumer laptop was receiving the audio signal and feeding it to a music tracking algorithm, computing at any point during the performance the current position in the score. This information was sent to the tablets of the test audience and triggered pre-prepared visualisations at the appropriate times. The audience could choose between 3 different kinds of synchronised visualisations: the sheet music (with synchronised highlighting of the current bar, and automatic page turning), textual information and explanations, and an

artistic video, visualising the story of the symphony (which is ‘Program Music’ *par excellence*). Two pictures with impressions from the live setup are shown in Figure 1.

This specific application of music tracking poses some unique challenges. Most importantly, so far the focus of music tracking has mostly been on solo or small ensemble music, like solo violin or flute, solo piano or string quartets, but nothing comparable to a full sized orchestra (according to Strauss’ notes the optimal size of the orchestra for the *Alpensinfonie* is 129 or more musicians!). This level of polyphony and of variety of instruments has to be considered when choosing the internal representation of the score and the features used during the on-line alignment process.

Furthermore, this piece challenges the music tracking algorithm with a vast range of musical situations: very quiet and slow parts without a clear melodic line (only a sound texture), very sparse parts with long pauses, energetic, loud and fast parts and even solo sections. Ideally, the tracker has to do equally well in all these situations, or at least well enough to not get lost completely. Thus, the main focus of our music tracking algorithm is placed on robustness. It actually does not matter much if an event is detected with a short delay, but it is very important that the algorithm does not get lost during this long piece (a typical performance takes about 50 minutes and contains no breaks).

3 The Score: Data Representation

To make the live tracking possible some internal representation of the musical score is needed. In most cases in music tracking the score is provided in symbolic form (e.g. MIDI or MusicXML). For the tracking, some kind of features are computed from the score representation that can be compared to the audio signal of the live performance.

Furthermore, for our task the content to be visualised has to be linked to the symbolic score, ideally via bar and beat numbers. For each video and every text snippet timing information is needed. Additionally, for the score visualisation also the area to be highlighted in the sheet music for each point in time needs to be known. We decided to provide this information at the level of musical bars.

The most natural approach (and most common in music

³<http://phenicx.upf.edu>

⁴<http://videodock.com>

The image shows a musical score excerpt for various instruments. The instruments listed on the left are: I. II. (Flutes), 3 Flg. (Flute), III. (Flute), Contrafag. (Contrabassoon), (F) I. (Horn), 2 Hörner (Horns), (B) IV. (Horn), 4 Pos. (Trumpets), 1. Baßtuba (Tuba), I. Viol. (vierfach) (Violins), II. Viol. (vierfach) (Violas), Br. (vierfach) (Trumpets), Violone. (vierfach) (Violone), and C.-B. (vierfach) (Contrabass). The score features sustained notes in the strings, horns, and contrabassoon, with dynamic markings of *p* (piano) and *pp* (pianissimo). A circled '2' is visible above the first flute staff. The bottom section of the score shows triplet figures in the bass section.

Figure 3: Excerpt from the score. This part is played very slowly and softly (note the *p* and *pp* dynamic markings), without a distinct melody (sustained notes in the strings, horns and the contrabassoon). The triplet figures in the bass section are so soft that they don't stand out but add to the overall sound texture.

sequences. We decided on using the features presented in [Arzt *et al.*, 2012] – alternatives approaches can be found in [Joder *et al.*, 2010]. These are well tested and reliable for audio to audio alignment tasks. Originally they were developed for tracking piano music, but as they try to combine two important properties of musical tones (the attack phase and the harmonic content), they are also well suited for a wider range of instruments.

Specifically, two different types of features are combined. 1) onset-emphasised features (an onset is the start time of a tone), which work particularly well for instruments and playing styles that produce sudden increases in the amplitude at onset times, and 2) harmonic features, which model the spectral content. As both features map the spectral data to the semitone scale, they are relatively robust to differences in timbre.

4 Following Live Orchestral Music: Tracking Algorithm

The music tracking algorithm we used is based on an on-line version of the dynamic time warping (DTW) algorithm, making it suitable for real-time music tracking (see [Dixon, 2005]). As input it takes two time series consisting of feature vectors – one known completely beforehand (the score) and one coming in real-time (the live performance) –, computes an on-line alignment, and at any time returns the current position in the score. In contrast to the standard dynamic time warping algorithm the alignment is computed incrementally, and it has linear time and space complexity due to the fact that instead of the whole alignment matrix, in each step only a small window centered at the current position is considered for computation.

Subsequently, improvements to this basic algorithm were proposed. This includes an extension called the ‘backward-forward strategy’ [Arzt *et al.*, 2008], which reconsiders past decisions and tries to improve the precision of the current score position hypothesis, and relatively simple tempo models [Arzt and Widmer, 2010] which are used to stretch or compress the score representation dynamically and therefore reduce differences in absolute tempo between the score representation and the live performance.

5 Adding a Safety Net: Multi-agent Tracking

While in preliminary experiments the presented combination of score representation, features and tracking algorithm generally gave robust and accurate results, we also witnessed some problems in specific circumstances. In the Alpensinfonie there is a part early in the piece (see Figure 3), starting around bar 14, that is played very softly and slowly. There is no distinct melody, only a relatively monotonic sound texture. Given sub-optimal sound quality (some noise, some distracting sounds), the tracker sometimes got lost or recovered very slowly.

Furthermore, a brief comparison of tracking results shows that, given multiple performances of a piece of music, some pairs of performances are ‘easier’ (i.e. more accurately and robustly) to align to each other, than others. Generally the features and the tracking algorithm are very robust to differences in tempo, timbre, tuning and even instrumentation, but the tracking process works best when the selected performances are similar to each other in these respects. In case an unsuitable pair of performances is selected, more inaccuracies will occur, and in very rare cases the tracking process might even fail at some point. To alleviate this problem, in [Wang *et al.*, 2014] an off-line music alignment algorithm was presented that improved the quality of pairwise alignments by also taking into account the additional information provided by alignments to all the other performances. Inspired by this we came up with a simple way to increase the robustness in the on-line case.

Instead of using one single instance of the tracking algorithm aligning the live performance to a ‘score performance’, multiple instances are initialised as independent agents, each using a different performance as its score representation (see Figure 4). The performances used for the Alpensinfonie are

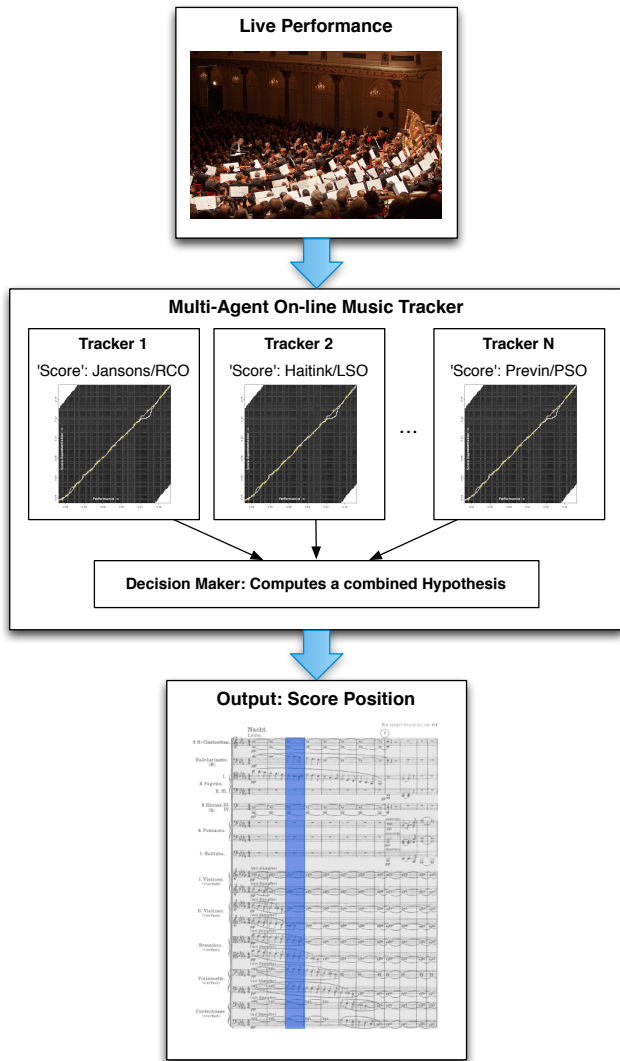


Figure 4: The Multi-agent Tracker. The live input is fed to N independent instances of the tracking algorithm. Each aligns the input to its own score representation, based on different performances of the same piece. Then, the individual hypotheses are combined and the estimate of the current position in the score is returned.

given in Table 1 – thus we had 7 trackers in total. As can be seen, only 1 performance was actually annotated manually while the other 6 were then aligned to it with the help of an off-line audio alignment algorithm (see [Grachten *et al.*, 2013]; other approaches to this problem include [Hu *et al.*, 2003] and [Ewert *et al.*, 2009]) to produce the information about the location of the downbeats. This means that these 6 additional ‘score performances’ were produced without any additional manual effort. Off-line alignment generally is much more accurate than on-line tracking, although it will still lead to some (but for our case acceptably small) inaccuracies.

During the concert the trackers run in parallel and each

tracker tries to align the incoming live performance to its score representation, each producing its own, independent hypothesis of the current position in the score. In the end the hypotheses are combined to form one collective hypothesis of the music tracking system.

Generally, many different ways of combining the hypotheses would be possible, e.g. based on voting or on the current alignment error of the individual trackers. As we had a clear goal in mind – robustness to single trackers getting lost – we decided on a very simple method: taking the median of the positions that are returned by the individual trackers. In our case this effectively allows us to come up with a good estimate of the score position even when 3 out of 7 trackers get lost.

6 The Event: Live Tracking in the Concertgebouw

The event on February 7, 2015 in the Concertgebouw was a big success. The tracking went smoothly and there were no glitches, only some minor inaccuracies. An obvious mistake happened at the quiet section in the beginning that was already discussed above. The sound texture here essentially consists of a very soft and repeating pattern. In cases like this the trackers sometimes tend to ‘wait’, because they try to align newly incoming instances of the pattern to past positions in the score (that also represent the same pattern). This resulted in a perceived delay of roughly 1 bar, for a period of about 5 bars. As soon as the texture changed and more distinct sounds could be recognised, the trackers recovered quickly. There were no further noticeable problems and in the end all of the trackers could follow the whole concert, and there was never any concern that the system might fail.

The general opinion amongst the project staff and the test audience was that the tracking worked very well and the accuracy was more than sufficient to trigger the visualisation in time. Only a few inaccuracies were noticed.

A formal in-depth evaluation, based on questionnaires the test audience had to fill out after the concert, will be published at a later point. This does not directly cover the tracking part as a technical process, but focuses on the user experience and on the value added by the provided additional information.

7 Evaluation

To be able to perform quantitative experiments the concert was recorded and annotated in the same way as the score representation above. Thus, the correct times of the down beats in the performance are known, and can be compared to the output of the music tracker. For the evaluation the error for each aligned downbeat is computed as the absolute value of the difference between the point in time given by the algorithm and the actual time according to the ground truth.

The results are presented in Tables 2 and 3. As can be seen, there are only slight differences in the results for the single tracker and the multi-agent approach. Keeping in mind that the goal of the multi-agent approach was to increase the robustness – the tracker would still produce similar results even when 3 out of 7 trackers fail –, this is a good result: extra robustness and a slight increase in accuracy were achieved

Err. (sec)	Single	Multi-agent
≤ 0.25	78.25%	81.80%
≤ 0.50	92.20%	93.24%
≤ 0.75	95.57%	96.44%
≤ 1.00	97.49%	98.01%

Table 2: Real-time alignment results for the single tracker and the multi-agent tracker, shown as cumulative frequencies of errors of matching pairs of downbeats. For instance, the first number in the first row means that the single tracker aligned 78.25% of the downbeats with an error smaller than or equal to 0.25 seconds.

	Single	Multi-agent
Average Error	0.20 sec.	0.19 sec.
Standard Dev.	0.35 sec.	0.36 sec.
First Quartile	0.06 sec.	0.05 sec.
Median Error	0.11 sec.	0.10 sec.
Third Quartile	0.22 sec.	0.19 sec.
Maximum Error	5.33 sec.	5.38 sec.

Table 3: Real-time alignment results for the single tracker and the multi-agent tracker.

without any extra manual efforts as the additional data was prepared by automatic methods.

Generally the results were more than sufficient for the task in question. The median error for the multi-tracking approach is about 0.1 seconds. Only in very rare cases did the tracker make major mistakes. Specifically the section already discussed above (see Figure 3) still causes problems, culminating in a maximum error of 5.38 seconds at bar 24 (which translates to about 1.5 bars, as this part has a relatively slow tempo). Actually the extent of the problem was not as apparent during the concert itself, also because even for humans it is very hard to follow the orchestra during this part.

8 Conclusions and Future Work

In this paper we presented a real-world application of machine-listening in the context of an actual concert in a world-famous concert hall. A music tracking algorithm was listening to the on-going live performance and was used to synchronise additional content (the sheet music, textual information and an artistic video), provided to the audience on tablet computers, to the live music.

The general impression during the concert was that the live tracking worked very well. This was confirmed later by a detailed quantitative evaluation.

As discussed above, our algorithm still runs into problems during soft and slow passages with very little structure or information. We are planning to solve this problem by both looking at additional features and by making stronger use of the tempo model at these parts.

A common problem of real-time music tracking and audio to score alignment are structural differences between the score and the performance. For example, if a piece has some

repeated sections, the performers might decide to play the repeat or to leave it out. For the Alpensinfonie this was not an issue, but in the future we will try to cope with this fully automatically – in the preparation phase via the technique used in [Grachten *et al.*, 2013], and in the live tracking phase with the approach presented in [Arzt *et al.*, 2014], extended to orchestral music.

We will also further investigate the multi-agent approach and will evaluate its merits in two scenarios: 1) in more noisy surroundings, and 2) for music of different genres, e.g. romantic piano music, where extreme differences in performing one and the same piece exist. Ultimately, we would like to use the multi-agent approach not only to increase the robustness, but also the accuracy of the tracking process.

In the future we wish to continue using the system in concert halls, possibly also on other genres. Specifically, we are considering opera – which would also lead directly to an additional use case: most opera houses provide subtitles for the audience, which so far are synchronised manually. Tracking an opera is a very challenging task, due to it being a combination of music (including singing voice), spoken language and acting and thus the best approach might be to combine multiple modalities (audio and video) within the tracking process.

Acknowledgments

This research is supported by the Austrian Science Fund (FWF) under project number Z159 and the EU FP7 Project PHENICX (grant no. 601166).

References

- [Arzt and Widmer, 2010] Andreas Arzt and Gerhard Widmer. Simple tempo models for real-time music tracking. In *Proceedings of the 7th Sound and Music Computing Conference (SMC 2010)*, Barcelona, Spain, 2010.
- [Arzt *et al.*, 2008] Andreas Arzt, Gerhard Widmer, and Simon Dixon. Automatic page turning for musicians via real-time machine listening. In *Proceedings of the 18th European Conference on Artificial Intelligence (ECAI 2008)*, Patras, Greece, 2008.
- [Arzt *et al.*, 2012] Andreas Arzt, Gerhard Widmer, and Simon Dixon. Adaptive distance normalization for real-time music tracking. In *Proceedings of the 20th European Signal Processing Conference (Eusipco 2012)*, Bucharest, Romania, 2012.
- [Arzt *et al.*, 2014] Andreas Arzt, Sebastian Böck, Sebastian Flossmann, Harald Frostel, Martin Gasser, and Gerhard Widmer. The complete classical music companion v0. 9. In *Proceedings of the 53rd AES Conference on Semantic Audio*, London, England, 2014.
- [Cont, 2009] Arshia Cont. A coupled duration-focused architecture for realtime music to score alignment. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(6):837–846, 2009.
- [Dannenberg, 1984] Roger Dannenberg. An on-line algorithm for real-time accompaniment. In *Proceedings of the International Computer Music Conference (ICMC 1984)*, Paris, France, 1984.

- [Dixon, 2005] Simon Dixon. An on-line time warping algorithm for tracking musical performances. In *Proceedings of the 19th International Joint Conference on Artificial Intelligence (IJCAI 2005)*, Edinburgh, Scotland, 2005.
- [Ewert *et al.*, 2009] Sebastian Ewert, Meinard Müller, and Peter Grosche. High resolution audio synchronization using chroma onset features. In *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2009)*, Taipei, Taiwan, 2009.
- [Grachten *et al.*, 2013] Maarten Grachten, Martin Gasser, Andreas Arzt, and Gerhard Widmer. Automatic alignment of music performances with structural differences. In *Proceedings of the 14th International Society for Music Information Retrieval Conference (ISMIR 2013)*, Curitiba, Brazil, 2013.
- [Hu *et al.*, 2003] Ning Hu, Roger Dannenberg, and George Tzanetakis. Polyphonic audio matching and alignment for music retrieval. In *IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*, New Paltz, USA, 2003.
- [Joder *et al.*, 2010] Cyril Joder, Slim Essid, and Gael Richard. A comparative study of tonal acoustic features for a symbolic level music-to-score alignment. In *Proceedings of the IEEE International Conference on Acoustics Speech and Signal Processing (ICASSP 2010)*, Dallas, USA, 2010.
- [Korzeniowski *et al.*, 2013] Filip Korzeniowski, Florian Krebs, Andreas Arzt, and Gerhard Widmer. Tracking rests and tempo changes: Improved score following with particle filters. In *Proceedings of the International Computer Music Conference (ICMC 2013)*, Perth, Australia, 2013.
- [Miron *et al.*, 2014] Marius Miron, Julio José Carabias-Orti, and Jordi Janer. Audio-to-score alignment at note level for orchestral recordings. In *Proceedings of the 15th International Society for Music Information Retrieval Conference (ISMIR 2014)*, Taipei, Taiwan, 2014.
- [Raphael, 2010] Christopher Raphael. Music Plus One and machine learning. In *Proceedings of the 27th International Conference on Machine Learning (ICML 2010)*, Haifa, Israel, 2010.
- [Vercoe, 1984] Barry Vercoe. The synthetic performer in the context of live performance. In *Proceedings of the International Computer Music Conference (ICMC 1984)*, Paris, France, 1984.
- [Wang *et al.*, 2014] Siying Wang, Sebastian Ewert, and Simon Dixon. Robust joint alignment of multiple versions of a piece of music. In *Proceedings of the 15th International Society for Music Information Retrieval Conference (ISMIR 2014)*, Taipei, Taiwan, 2014.