

A Lightweight Multi-Agent Musical Beat Tracking System

Simon Dixon

Austrian Research Institute for Artificial Intelligence,
Schottengasse 3, A-1010 Vienna, Austria.
simon@ai.univie.ac.at

Abstract

Beat tracking is what people do when they tap their feet in time to music. We present a software system which performs this task, processing music in a standard digital audio format and estimating the locations of musical beats. A time-domain algorithm detects salient acoustic events, and then a clustering algorithm groups the time intervals between events to obtain hypotheses about the current tempo. Multiple competing agents track these hypotheses throughout the music, with further agents being created at decision points. The output for each agent is a sequence of beat locations, which is evaluated for its closeness of fit to the data. This approach to beat tracking assumes no previous knowledge of the music such as the style, time signature or approximate tempo; all required information is derived from the audio data. The system has been tested with various styles of music (popular, jazz, and classical) and performs robustly, rarely making errors in popular music, and recovering quickly from errors in more complex styles of music, despite the fact that no high level musical knowledge is encoded in the system. We describe several applications, including musical score extraction and an automatic disc jockey that performs beat mixing in real time.

Introduction

Although most people can tap their feet in time with music, equivalent performance on a computer has proved remarkably difficult to achieve. One reason for this is that these systems have been based only on the codification of high level metrical knowledge. We show that such knowledge is secondary to the beat tracking task. Just as people can follow the beat of music without musical training and without previous knowledge of the particular piece of music, so can a computer program.

We do not attempt to model or describe the cognitive mechanisms involved in human rhythm perception. However we do note certain features of perception which motivate an ambitious unsupervised approach to the beat tracking problem. Firstly, human rhythm perception sets its own parameters; the tempo and the metrical structure are not specified to a listener at the beginning of a piece, and if they change suddenly during

the piece, the perceptual system adjusts within seconds to the new listening framework. Secondly, it copes well with “noise” in the input. That is, deviations from precise timing and variations in tempo do not destroy the overall perception of the beat. Thirdly, human perception is able to cope with syncopation, that is, sections of music where the more salient events are occurring between the beats rather than on the beat.

In contrast with these capabilities, computer music software does not normally cope well in these situations. Commercial sequencing and transcription programs usually require the beat to be declared explicitly before music can be processed, so that all data can then be indexed relative to this given beat. Even many research systems are limited by the fact that once they get out of synchronization with the music, it is very difficult for them to recover and resume correct interpretation of the rhythmic structure (Dannenberg 1991). The robustness of human perception is one feature which has proved difficult to reproduce in a computer system.

In this paper, we present a system which processes musical audio signals, estimating the tempo and determining the locations of musical beats. No specific assumptions are made about the music being analyzed, so the system performs robustly, recovering quickly from tracking errors. The system has been tested on various types of popular, jazz and classical music.

The subsequent sections of the paper contain a review of related work, then a description of the underlying model of musical timing and the assumptions on which the system is based. The next sections present the algorithm for onset detection from raw audio data, followed by the algorithm for *tempo induction*, defined as the estimation of the time interval between successive occurrences of the main rhythmic pulse of the music. We then describe the multi-agent approach to *beat tracking*, which is the determination of beat locations (and therefore tempo fluctuations) in the light of the previous tempo estimations. The results from testing the system with various types of music are then presented and discussed, and the paper concludes with a description of applications of the beat tracking system.

Related Work

A substantial amount of research has been performed in the area of rhythm recognition by computer, including a demonstration of various beat tracking methods using a computer to control a shoe which tapped in time with the calculated beat of the music (Desain & Honing 1994). These systems are difficult to compare directly, as they make different assumptions about the input format, style and complexity of the music.

Much of the work in machine perception of rhythm has used MIDI files as input, which contain control information for a synthesizer rather than audio data. MIDI files consist of sequences of events, usually corresponding to pressing and releasing keys on a piano-style keyboard, plus an encoding of the time duration between successive events. Structural information such as the time signature and tempo can also be stored in MIDI files, but it is usually assumed that such information is not available to rhythm recognition programs.

Using MIDI files, the input is usually interpreted as a series of event times, ignoring the event duration, pitch, amplitude and chosen synthesizer voice. That is, each note is treated purely as an uninterpreted event. It is assumed that the other parameters do not provide *essential* rhythmic information, which in many circumstances is true. However, there is no doubt that these factors provide useful rhythmic cues; for example, more salient events tend to occur on stronger beats.

Notable work using MIDI file input is an emulation of human rhythm perception (Rosenthal 1992), which produces multiple hypotheses of possible hierarchical structures in the timing, assigning a score to each hypothesis, corresponding to the likelihood that a human listener would choose that interpretation of the rhythm. This technique gives the system the ability to adjust to changes in tempo and meter, as well as avoiding many of the implausible rhythmic interpretations produced by commercial systems.

Allen and Dannenburg (1990) describe a beat tracking system that uses beam search to consider multiple hypotheses of beat timing and placement. A heuristic evaluation function directs the search, preferring interpretations that have a “simple” musical structure and make “musical sense”, although they do not define what they mean by these terms. They also do not describe the input format or any specific results.

One early project on rhythm using audio input was the percussion transcription system of Schloss (1985). Onsets were detected as peaks in the slope of the amplitude envelope, where the envelope was defined to be equal to the maximum amplitude in each period of the high-pass filtered signal, and the period defined as the inverse of the lowest frequency expected to be present in the signal. The system was limited in that it required parameters to be set interactively, and it was evaluated only by resynthesis of the signal.

A more complete approach to beat tracking of acoustic signals was developed by Goto and Muraoka (1995; 1997b; 1998). They developed two systems for follow-

ing the beat of popular music in real time. The earlier system (BTS) used frequency histograms to find significant peaks in the low frequency regions, corresponding to the frequencies of the bass and snare drums, and then tracked these low frequency signals by matching patterns of onset times to a set of pre-stored drum beat patterns. This method was successful in tracking the beat of most of the popular songs on which it was tested. A later system allowed music without drums to be tracked by recognizing chord changes, assuming that significant harmonic changes occur at strong rhythmic positions. These systems required a powerful parallel computer in order to run in real time.

Commercial transcription and sequencing programs do not address the issues covered by these research systems. They generally require that the tempo and time signature are specified before the music is played, and the system then aligns each note with the nearest position on a metrical grid. Recent systems allow parameterization of this grid in terms of its resolution (the shortest allowed note length), and adjustment of restrictions on the complexity of rhythm that can be produced by the system. These systems often produce implausible rhythmic interpretations, and cannot be used in an unsupervised manner for anything but simple rhythms.

Musical Timing

Despite the large amount of research in time and rhythm in music, the beat tracking problem remains poorly defined. The reason is that the beat is a subjective property of performed music. Formal musical models tend to be based on the notational representation rather than performance (Lerdahl & Jackendoff 1983; Longuet-Higgins & Lee 1982), and those which address performance timing do so from the point of view of generation and/or transformation of timing rather than extraction or explanation of performance data (Desain & Honing 1991). We follow (Goto & Muraoka 1997a) in evaluating the correctness of the beat tracking system relative to a subjective labelling of beat positions.

A theoretical definition of *beat* is a perceived pulse marking off equal durational units; in practice, the durational units marked off by the onsets of notes on successive beats are only approximately equal. Performance studies have shown that asynchronies of events (with respect to notation and other events) are often in the range of 20-50ms, in both ensemble situations (Keil 1995) and solo performances (Palmer 1996). In such situations there can be more than one “correct” beat location. However, for a large amount of music, the subjective differences in perceived beat are minor, otherwise human activities such as ensemble playing and dancing would not be possible. In this study we restrict our attention to music which has such an agreed beat, and rely on the onset detection algorithm to choose the more salient events as possible beat locations. Aural testing confirms that this method is sufficient for the types of music we are examining.

Audio Processing

In this and the following sections, we describe the stages of processing performed by the beat tracking system. All of the software is written in C++ and runs on a Unix platform (Linux or Solaris). The complete processing of a song takes about 10 seconds on a current PC, making it viable for use in real time audio applications, although the software is not currently designed for real time use. The input to the system is a digitally sampled acoustic signal, such as is found on audio compact discs. The stereo compact disc data is converted to a single channel format by averaging the left and right channels, resulting in a single channel 16 bit linear pulse code modulated (PCM) format, with a sampling rate of 44.1kHz.

The aim of the initial signal processing stage is to detect *events* in the audio data, from which rhythmic information can be derived. For the purposes of this work, events correspond to note onsets, that is, the beginnings of musical notes, including percussive sounds. By ignoring note durations and offset times, we discard valuable information, but our results justify the present assumption, that there is sufficient information in note onsets to perform beat tracking.

A time-domain method similar to (Schloss 1985) is employed for onset detection. This method involves passing the signal through a simple high-pass filter, calculating the absolute sum of small overlapping windows of the signal, and then finding peaks in the slope of these window sums using a 4 point linear regression. Only the more salient event onsets are detected with the method, which is ideal for the subsequent task of tempo induction.

Tempo Induction

The tempo induction section of the system determines a set of hypotheses about the tempo of a given section of music, which may be expressed in beats per minute (BPM) or in seconds, as the inter-beat interval (IBI). The algorithm, described further in (Dixon 1997; 1999), is based on clustering of *inter-onset intervals* (IOI's). In the literature, an IOI is defined as the time between the onsets of two successive events, but we extend the definition to include times between onsets of pairs of events that are separated by intervening event onsets. All possible pairs of onsets that occur within 2.5 seconds of each other are grouped by the clustering algorithm. Figure 1 shows clustering for five events (A, B, C, D, E) into intervals of similar size. For example, cluster C1 contains the intervals AB, BC and DE, while cluster C2 contains AC and CD. Each cluster is identified by its average interval size.

After grouping IOI's into clusters, a score is calculated for each cluster, based on the number of IOI's in the cluster. The highly ranked clusters usually correspond to the beat or small integer multiples or fractions of the beat. For example, supposing that C2 represents the IBI, then C1 represents half of the IBI and C4 rep-

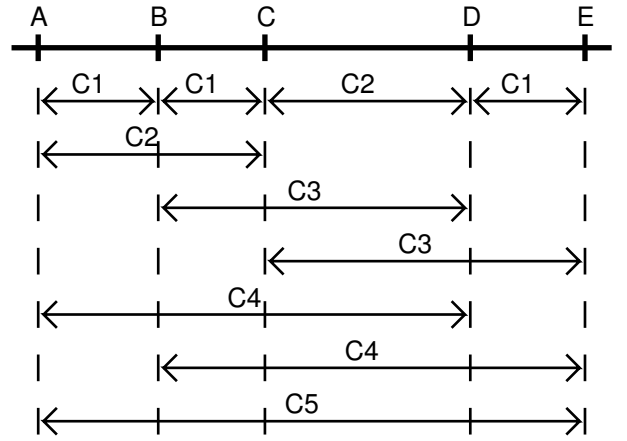


Figure 1: Clustering of inter-onset intervals

resents double the IBI. Each cluster's score is increased for each other cluster to which it is related by a small integer ratio, and a final ranking of the inter-beat interval hypotheses is determined.

In previous work, it was found that the correct tempo can be induced from a 5-10 second excerpt of the music with 90% reliability, and by using multiple (or longer) excerpts, the reliability quickly approaches 100%. In this work, it does not matter if the initial estimate is correct, as multiple hypotheses are checked in the beat tracking stage, so that an error in tempo induction can be corrected at a later time.

Beat Tracking Agents

The tempo induction algorithm computes the inter-beat interval, that is, the time between successive beats, but does not calculate the location of the beat. In Figure 1, the clustering might determine that C2 represents the inter-beat interval, but it does not reveal whether events A, C and D are beat locations or whether B and E are beat locations. By analogy with wave theory, we could say that it calculates the *frequency* but not the *phase* of the beat.

The techniques used in previous work did not produce a reliable estimate of phase. The main difficulty with phase calculations is that they are extremely sensitive to errors in the inter-beat interval, because they are measured in fractions of a beat, so any tempo error is multiplied by the number of beats between events being examined. Also, it is not possible to average phase values, as the metrical positions of events are unknown (in the absence of a musical score), and it would only be meaningful to average the phases of events if they were known to occur in the same relative position within the beat.

The phase calculation problem was solved by employing an agent-based architecture to examine multiple hypotheses simultaneously throughout the music. The agents are characterized by their *state* and *history*.

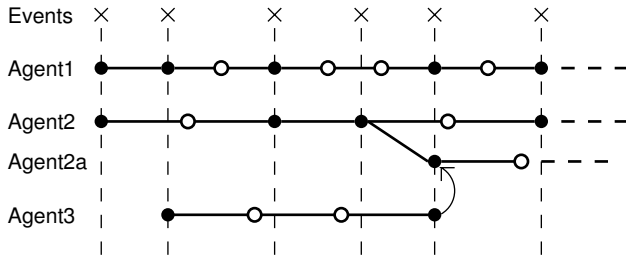


Figure 2: Beat tracking agents (see text for details)

The state is the agent’s current hypothesis of the beat frequency and phase, and the history is the sequence of beat locations selected so far by the agent. Each agent is evaluated on the basis of its history, with higher scores being awarded for greater regularity in the spacing between events, greater salience of chosen events, and fewer gaps in the sequence.

Initially, a number of agents are created for each of the tempo hypotheses from the tempo induction stage; for each tempo, one agent is created for each of the first few events in the piece, with its phase set to 0 at the time of the event onset. A simplified example is shown in Figure 2, where there are 2 tempo hypotheses, and two starting locations. Agents 2 and 3 start with the same tempo, but different phase, while Agent 1 has a different tempo, but same initial phase as Agent 2. Note there is no need to start an agent with the tempo of Agent 1 and phase of Agent 3, since Agent 1 covers the second event itself.

The main loop of the beat tracking section passes each event to each agent, which compares the event’s onset time with the predicted beat location. The agents have two windows of tolerance, an inner window, within which the agent is sure that the event corresponds to the predicted beat location, and an outer window, within which the agent is unsure if the event should be accepted as a beat location. If the event falls in the inner window, it is added to the agent’s history, and the agent’s tempo and phase hypotheses are updated. In the case that the event falls in the outer window around the predicted location, the agent creates a clone which accepts the event as a beat location, while the current agent rejects the event. In Figure 2, Agent 2 creates Agent 2a when the 5th event falls in its outer window of expected beat locations. This guards against the current beat location being lost due to a rogue event, whilst also allowing for moderate deviations in tempo and phase to occur. When an accepted event is more than one beat from the previous beat location determined by the agent, the missing beats are filled in by interpolation (shown by hollow circles in the figure).

As the agents track the beats, a confidence value is maintained for each agent. This value is increased each time an event is accepted as a beat location. The amount of increase depends on the salience of the event and its proximity to the predicted beat location. The

salience is measured in terms of the amplitude of the event onset, calculated on a logarithmic scale. The value is then reduced according to the difference between the predicted and actual beat locations, and then added to the agent’s confidence value. The final confidence value for an agent is calculated by reducing the confidence for each beat which had to be interpolated, and normalizing the result so that the agents with a faster tempo are not advantaged by their greater number of beats.

It often occurs that two or more agents come to the same conclusion about the current tempo and beat phase. Since these are the only variables that determine the agent’s state (and therefore its future behavior), it is computationally advantageous to remove all but one of these agreeing agents, retaining only the agent with the highest confidence (that is, with the best scoring history). In Figure 2, Agent 3 is terminated when its state coincides with that of Agent 2a, as indicated by the arrow at the 5th event. Removal of duplicate agents is performed after each event is processed. Agents are also checked for currentness, and are removed if they are unable to find any event corresponding to a predicted beat for some fixed length of time.

After the last event is processed, the highest scoring current agent is selected, and its history is output as the beat tracking “solution”. It is also possible to view a trace of the agents and their scores at each event during processing. For aural testing (and demonstrations) of the system, the music can also be played back or saved to file with a click track added to it, that is, a percussion track indicating the positions of beats as detected by the system. In the following section, testing methodology is discussed, and the results of beat tracking with various styles of music are presented.

Results and Discussion

Informal testing was performed by listening to the music with synthesized percussion strokes (e.g. cow-bell) played at the beat locations computed by the system. With this method, it is easy to check that the tempo estimation and tracking are approximately correct, but it is not a very precise form of testing. It is also very time-consuming if used repeatedly to test the effects of small adjustments to the system. However, aural testing provides intuition about the situations in which the beat tracker fails, which is useful for determining which aspects of the system require further work.

More precise testing was performed by comparison of results with manually calculated beat positions. The audio files were examined with standard digital audio editing software (GoldWave) which permits viewing the data at arbitrary resolution, and playing selections at arbitrary speeds. Beat locations were determined for a number of beats on which clear percussive events occurred. These beat boundaries were then used to interpolate the locations of the intervening beats. This avoided the problem of determining beat locations in

Song Title	Artist	CD Code	Style	Date	Tempo range	Time sign.	Results Tempo Phase	
I Don't Remember A Thing	Paul Kelly and the Coloured Girls	Mushroom CD 53248	Pop/rock	1987	139-142	4/4	Yes	Yes
Dumb Things	" "	Mushroom CD 53248	Pop/rock	1987	151-154	4/4	Yes	Yes
Untouchable	" "	Mushroom CD 53248	Pop/rock	1987	145-146	4/4	Yes	Yes
Superstition	Stevie Wonder	Motown37463-03192-9	Motown	1972	96-104	4/4	Yes	Yes
You Are The Sunshine of My Life	Stevie Wonder	Motown37463-03192-9	Motown	1972	127-136	4/4	Yes	Yes
On A Night Like This	Bob Dylan	Columbia CD 32154	Country	1974	136-140	4/4	Yes	Part
Rosa Moreña	João Gilberto Trio	Jazz Roots CD 56046	Bossa nova	1964	128-134	4/4	Yes	Part
Michelle	Béla Fleck and the Flecktones	Warner 7599-26562-2	Jazz swing	1991	180-193	3/4	Yes	Part
Jitterbug Waltz	James Morrison	WEA 9031-71211-1	Jazz waltz	1990	155-175	3/4	Yes	Part

unclear passages, as discussed in the previous section on musical timing.

We now discuss the results shown above. The two results columns on the right indicate whether the highest scoring agent had the correct tempo, and whether its beat locations agreed with those calculated manually for all (Yes), part (Part) or none (No) of the song.

The first 3 songs are standard modern pop/rock, characterized by very steady tempo, which is clearly defined by simple and salient drum patterns, similar to the data used in early audio beat tracking work of (Goto & Muraoka 1995). In the production of this style of music, it is common practice for each instrument to be recorded separately, using a metronome track for synchronization. In this case one expects the performed beat to be very regular, with only small deviations from mechanical regularity. The beat tracking system made no errors on these songs.

The next style examined was Motown/Soul, characterized by more syncopation, greater tempo fluctuations (5-10% in these examples), and more freedom to anticipate or lag behind the beat. Despite the greater difficulty in beat tracking, the complete songs were tracked correctly.

The Bob Dylan song was more difficult to track, because of his idiosyncratic style of singing and playing against the rhythmic context. Although the beat is reasonably clear to a human listener, the drums are not prominent, and there is a much lower correlation between the conceptual beat and the actual musical events than in the other styles. The beat tracking system tracked correctly up to the instrumental section after the final verse, in which it lost synchronization and tracked the off-beats (i.e., it continued at the correct tempo but half a beat out of phase).

The next test involved a live bossa nova performance with syncopated guitar and very little percussion to indicate beat positions. The song was tracked correctly except in one passage where it went out of phase, but the error was corrected within about 10 beats.

The two jazz pieces were chosen for their particularly complex, syncopated rhythms, which are difficult for humans to follow. These pieces also provided exam-

ples of a different time signature, swing eighth notes, and greater tempo variation. In both cases, the highest scoring agent was able to track the majority of the piece correctly, but encountered phase errors in some parts. (This was not the first time that phase errors were encountered. With the salience calculation removed, the system tracks the whole of *I Don't Remember A Thing* at half a beat out of phase.) For rock music, the salience of events differentiates the beat from the offbeat at most points in the music. This is not true in jazz, where the offbeat is often accentuated for long periods of time, so the system requires an alternative way of choosing the correct path through the data.

Finally, a classical piece was tested, the third movement of Mozart's Piano Sonata in C major (KV279). The system lost synchronization several times, tracking the off-beats rather than the beats, due to large tempo variations and the system's lack of musical knowledge for distinguishing between beats and off-beats. Note that the beat tracking system is not equipped with musical knowledge — no notion of off-beats or expected rhythmic patterns has been programmed into it. Its apparent musical intelligence comes from patterns *in the data*, without any high-level knowledge or reasoning (Brooks 1991). Apart from the simplicity of this approach, a great advantage is that the system is quite robust, and generalizes well to different styles of music, as long as there is a salient beat. In order to disambiguate complex or ambiguous rhythmic patterns, the system will need sources of musical knowledge other than timing of events; these are not presently available to it. In current work, we are examining a specialization of the system for solo piano music which will incorporate a level of musical and stylistic knowledge with the aim of extracting the score from performance data.

Conclusion

We have described a beat tracking system which analyses acoustic data, detects the salient note onsets, determines possible inter-beat intervals and then employs multiple agents to find a sequence of events which represents the beat of the music. The system successfully

tracks the beat in most popular music, but makes some phase errors when presented with extremely complex rhythms or music with large tempo deviations. Even in these situations, the performance is quite robust, with the system recovering from its errors and resuming correct tracking after a short period.

Unlike previous audio beat tracking systems which required a large parallel computer (Goto & Muraoka 1998), our system has modest requirements, processing a song in under 10 seconds on a current personal computer, leaving sufficient resources for real time applications using the beat tracking system as one component.

One such application is an automatic disc jockey (DJ), which plays a list of songs, cross-fading between the songs so that the beats of successive songs are synchronized (beat-mixing). Another application which is currently being pursued is that of a score extraction system. This application uses MIDI input rather than audio, and the system's job is to make "musical sense" of the performed rhythm. The nature of this problem is different, in that we seek a musical explanation for every event, whereas the current system ignores events which are determined not to lie on the beat. MIDI input also facilitates the use of other knowledge from the data, such as duration, pitch, repeated melodic patterns and musical voice, as well as external musical knowledge concerning, for example, the use of ornaments. It is still an open problem how such details can be extracted reliably directly from audio data. A further application of beat tracking, and one which requires reliable recognition of pitch and duration of notes, is an automatic music transcription system, that is, a system which produces musical scores directly from audio data.

The use of manual beat tracking for evaluation of the system limits the amount of testing that can be performed, but is necessary when analyzing performed music. It would also be useful to perform a study of beat tracking in synthetically generated music, where the variations in tempo and onset times could be controlled precisely, and performance could be evaluated automatically.

Acknowledgements

This research is part of the project Y99-INF, sponsored by the Austrian Federal Ministry of Education, Science and Culture in the form of a START Research Prize. The Austrian Research Institute for Artificial Intelligence is supported by the Austrian Federal Ministry of Education, Science and Culture. The author also wishes to thank Emiliós Cambouropoulos and Gerhard Widmer for comments on earlier drafts of this paper.

References

- Allen, P., and Dannenburg, R. 1990. Tracking musical beats in real time. In *Proceedings of the International Computer Music Conference*, 140–143. International Computer Music Association, San Francisco CA.
- Brooks, R. 1991. Intelligence without reason. In *Proceedings of the International Joint Conference on Artificial Intelligence*.
- Dannenberg, R. 1991. Recent work in real-time music understanding by computer. *Proceedings of the International Symposium on Music, Language, Speech and Brain*.
- Desain, P., and Honing, H. 1991. Towards a calculus for expressive timing in music. *Computers in Music Research* 3:43–120.
- Desain, P., and Honing, H. 1994. Foot-tapping: a brief introduction to beat induction. In *Proceedings of the International Computer Music Conference*, 78–79. Computer Music Association, San Francisco CA.
- Dixon, S. 1997. Beat induction and rhythm recognition. In *Proceedings of the Australian Joint Conference on Artificial Intelligence*, 311–320.
- Dixon, S. 1999. A beat tracking system for audio signals. In *Proceedings of the Diderot Forum on Mathematics and Music*, 101–110. Austrian Computer Society.
- Goto, M., and Muraoka, Y. 1995. A real-time beat tracking system for audio signals. In *Proceedings of the International Computer Music Conference*. Computer Music Association, San Francisco CA.
- Goto, M., and Muraoka, Y. 1997a. Issues in evaluating beat tracking systems. In *Issues in AI and Music – Evaluation and Assessment: Proceedings of the IJCAI'97 Workshop on AI and Music*. International Joint Conference on Artificial Intelligence.
- Goto, M., and Muraoka, Y. 1997b. Real-time rhythm tracking for drumless audio signals – chord change detection for musical decisions. In *Proceedings of the IJCAI'97 Workshop on Computational Auditory Scene Analysis*. International Joint Conference on Artificial Intelligence.
- Goto, M., and Muraoka, Y. 1998. An audio-based real-time beat tracking system and its applications. In *Proceedings of the International Computer Music Conference*. Computer Music Association, San Francisco CA.
- Keil, C. 1995. The theory of participatory discrepancies: a progress report. *Ethnomusicology* 39(1):1–19.
- Lerdahl, F., and Jackendoff, R. 1983. *A Generative Theory of Tonal Music*. MIT Press.
- Longuet-Higgins, H., and Lee, C. 1982. The perception of musical rhythms. *Perception* 11:115–128.
- Palmer, C. 1996. Anatomy of a performance: Sources of musical expression. *Music Perception* 13(3):433–453.
- Rosenthal, D. 1992. Emulation of human rhythm perception. *Computer Music Journal* 16(1):64–76.
- Schloss, W. 1985. *On the Automatic Transcription of Percussive Music — From Acoustic Signal to High Level Analysis*. PhD thesis, CCRMA, Stanford University.