

A Survey of Computer Systems for Expressive Music Performance

ALEXIS KIRKE

Interdisciplinary Centre for Computer Music Research (ICCMR),
University of Plymouth, Plymouth, UK, Alexis.Kirke@Plymouth.ac.uk

EDUARDO RECK MIRANDA

Interdisciplinary Centre for Computer Music Research (ICCMR),
University of Plymouth, Plymouth, UK, Eduardo.Miranda@Plymouth.ac.uk

We present a survey of research into automated and semi-automated computer systems for expressive performance of music. We will examine the motivation for such systems and then examine the majority of the systems developed over the last 25 years. To highlight some of the possible future directions for new research, the review uses primary terms of reference based on 4 elements: Testing Status, Expressive Perception, Non-monophonic Ability, and Performance Creativity.

Categories and Subject Descriptors: J.5 [**Computer Applications**]: Arts and Humanities – Performing Arts

Additional Key Words and Phrases: Music Performance, Generative Performance, Computer Music, Machine Learning

1. INTRODUCTION

In the early 1980s the seeds of a problem were sown as a result of synthesizers being developed and sold with built-in sequencers. The introduction of MIDI into this equation led to an explosion in the use of sequencers and computers, thanks to the new potential for connection and synchronisation. These computers and sequencers performed their stored tunes in perfect metronomic time, a performance which sounded inhuman. They sounded inhuman because human performers normally perform expressively – for example speeding up and slowing down while playing, and changing how loudly they play. The performer’s changes in tempo and dynamics allow them to express a fixed score – hence the term expressive performance [Widmer and Goebel 2004]. However rather than looking for ways to give the music performances more human-like expression, pop performers developed new types of music, such as synth-pop and dance music, that actually utilized this metronomic perfection to generate “robotic” performances.

Outside of pop, the uptake of sequencers for performance (as opposed to for composition) was less enthusiastic, except for occasional novelties like *Switched on Bach* by Wendy Carlos, computer performance of classical music was a rarity. Computer composition of classical music had been around since 1957 when *The Illiac Suite* for String Quartet - the first published composition by a computer - was published by Lejaren Hiller [Hiller and Isaacson 1959]. Since then there has been a large body of such music and research published, with many successful systems produced for automated and semi-automated computer composition [Buxton 1977][Roads 1996][Miranda 2001]. But publications on computer expressive performance of music lagged behind composition by almost quarter of a century. During the period when MIDI and computer use exploded amongst pop performers, and up to 1987 - when Yamaha had released their first Disklavier MIDI piano - there were only 2 or 3 researchers publishing on algorithms for expressive performance of music [Todd 1985; Sundberg et al. 1983]. However, from the end of the 1980s onwards there was an increasing interest in automated and semi-automated Computer Systems for Expressive Music Performance (CSEMP). A CSEMP is a computer system able to generate expressive performances of music. For example software for music typesetting

will often be used to write a piece of music, but some packages play back the music in a relatively robotic way – the addition of a CSEMP enables a more realistic playback. Or an MP3 player could include a CSEMP which would allow performances of music to be adjusted to different performance styles.

In this paper we will review the majority of research on automated and semi-automated CSEMPs. By automated, we refer to the ability of the system – once set-up or trained - to generate a performance of a new piece, not seen before by the system, without manual intervention. Some automated systems may require manual set-up, but then can be presented with multiple pieces which will be played autonomously. A semi-automated system is one which requires some manual input from the user (for example a musicological analysis) to deal with a new piece.

1.1 Human Expressive Performance

How do humans make their performances sound so different to the so-called “perfect” performance a computer would give? In this paper the strategies and changes which are not marking in a score but which performers apply to the music will be referred to as expressive Performance Actions. Two of the most common performance actions are changing the Tempo and the Loudness of the piece as it is played. These should not be confused with the tempo or loudness changes marked in the score, like *accelerando* or *mezzo-forte*, but to additional tempo and loudness changes not marked in the score. For example, a common expressive performance strategy is for the performer to slow down as they approach the end of the piece [Friberg and Sundberg 1999]. Another performance action is the use of expressive articulation – when a performer chooses to play notes in a more staccato (short and pronounced) or legato (smooth) way. Those playing instruments with continuous tuning, for example string players, may also use expressive intonation, making notes slightly sharper or flatter; and such instruments also allow for expressive vibrato. Many instruments provide the ability to expressively change timbre as well.

Why do humans add these expressive performance actions when playing music? We will set the context for answering this question using a historical perspective. Pianist and musicologist Ian Pace offers up the following as a familiar historical model for the development of notation (though suggests that overall it constitutes an over-simplification) [Pace 2007]:

In the Middle Ages and to a lesser extent to the Renaissance, musical scores provided only a bare outline of the music, with much to be filled in by the performer or performers, freely improvising within conventions which were essentially communicated verbally within a region or locality. By the Baroque Era, composers began to be more specific in terms of requirements for pitch, rhythm and articulation, though it was still common for performers to apply embellishments and diminutions to the notated scores, and during the Classical Period a greater range of specificity was introduced for dynamics and accentuation. All of this reflected a gradual increase in the internationalism of music, with composers and performers travelling more widely and thus rendering the necessity for greater notational clarity as knowledge of local performance conventions could no longer be taken for granted. From Beethoven onwards, the composer took on a new role, less a servant composing to occasion at the behest of his or her feudal masters, more a freelance entrepreneur who followed his own desires, wishes and convictions, and wrote for posterity, hence bequeathing the notion of the master-work which had a more palpable autonomous existence over and above its various manifestations in performance. This required an even greater degree of notational exactitude; for example in the realms of tempo, where generic Italianate conventions were both rendered in the composer’s native language and finely nuanced by qualifying clauses and adjectives. Through the course

of the nineteenth century, tempo modifications were also entered more frequently into scores, and with the advent of a greater emphasis on timbre, scores gradually became more specific in terms of the indication of instrumentation. Performers phased out the processes of embellishment and ornamentation as the score came to attain more of the status of a sacred object. In the twentieth century, this process was extended much further, with the finest nuances of inflection, rubato, rhythmic modification coming to be indicated in the score. By the time of the music of Brian Ferneyhough, to take the most extreme example, all minutest details of every parameter are etched into the score, and the performer's task is simply to try and execute these as precisely as he or she can.

So in pre-20th century music there has been a tradition of performers making additions to a performance which were not marked in the score (though the reason Pace calls this history an oversimplification is that modern music does have the capacity for expressive performance, as we will discuss later).

A number of studies have been done into this pre-20th Century (specifically Baroque, Classical and Romantic) music performance. The earliest studies began with Seashore [1938], and good overviews include Palmer [1997] and Gabrielsson [2003]. One element of these studies has been to discover what aspects of a piece of music – what Musical Features - are related to a performer's use of expressive performance actions. One of these musical features expressed is the performer's structural interpretation of the piece [Palmer 1997]. A piece of music has a number of levels of meaning – a hierarchy. Notes make up motifs, motifs make up phrases, phrases make up sections, sections make up a piece (in more continuous instruments there are intranote elements as well). Each element - note, motif, etc - plays a role in other higher elements. Human performers have been shown to express this hierarchical structure in their performances. Performers have a tendency to slow down at boundaries in the hierarchy – with the amount of slowing being correlated to the importance of the boundary [Clarke 1998]. Thus a performer would tend to slow more at a boundary between sections than between phrases. There are also regularities relating to other musical features in performers' expressive strategies. For example in some cases the musical feature of higher pitched notes causes a performance action of the notes being played more loudly; also note which introduce melodic tension relative to the key may be played more loudly. However for every rule there will always be exceptions.

Another factor influencing expressive performance actions is Performance Context. Performers may wish to express a certain mood or emotion (e.g. sadness, happiness) through a piece of music. Performers have been shown to change the tempo and dynamics of a piece when asked to express an emotion as they play it [Gabrielsson and Juslin 1996]. For a discussion of other factors involved in human expressive performance, we refer the reader to [Juslin 2003].

1.2 Computer Expressive Performance

Having examined human expressive performance, the question now becomes why should we want computers to perform music expressively? There are at least five answers to this question:

1. Investigating human expressive performance by developing computational models – Expressive performance is a fertile area for investigating musicology and human psychology [Seashore 1938; Palmer 1997; Gabrielsson 2003]. As an alternative to experimentation with human performers, models can be built which attempt to simulate elements of human expressive performance. As in

all mathematical and computational modelling, the model itself can give the researcher greater insight into the mechanisms inherent in that which is being modelled.

2. Realistic playback on a music typesetting or composing tool – There are many computer tools available now for music typesetting and for composing. If these tools play back the compositions with expression on the computer, the composer will have a better idea of what their final piece will sound like. For example, Sibelius, Notion and Finale have some ability for expressive playback.
3. Playing computer-generated music expressively – There are a number of algorithmic composition systems that output music without expressive performance but which audiences would normally expect to hear played expressively. These compositions in their raw form will play on a computer in a robotic way. A CSEMP would allow the output of an algorithmic composition system to be played directly on the computer which composed it (for example in a computer game which generates mood music based on what is happening in the game).
4. Playing data files - a large number of non-expressive data files in formats like MIDI and MusicXML [Hirata et al 2003] are available on the internet, and they are used by many musicians as a standard communication tool for ideas and pieces. Without CSEMPs most of these files will playback on a computer in an unattractive way, whereas the use of a CSEMP would make such files much more useful.
5. Computer accompaniment tasks - it can be costly for a musician to play in ensemble. Musicians can practice by playing along to recordings with their solo part stripped out. But some may find it too restrictive since such recordings cannot dynamically follow the expressiveness in the soloist's performance. These soloists may prefer to play along with an interactive accompaniment system that not only tracks their expression but also generates its own expression.

2. A GENERIC FRAME FOR PREVIOUS RESEARCH IN COMPUTER EXPRESSIVE PERFORMANCE

Figure 1 shows a generic model for the framework that most (but not all) previous research into automated and semi-automated CSEMPs tends to have followed. The modules of this diagram are described beneath Figure 1.

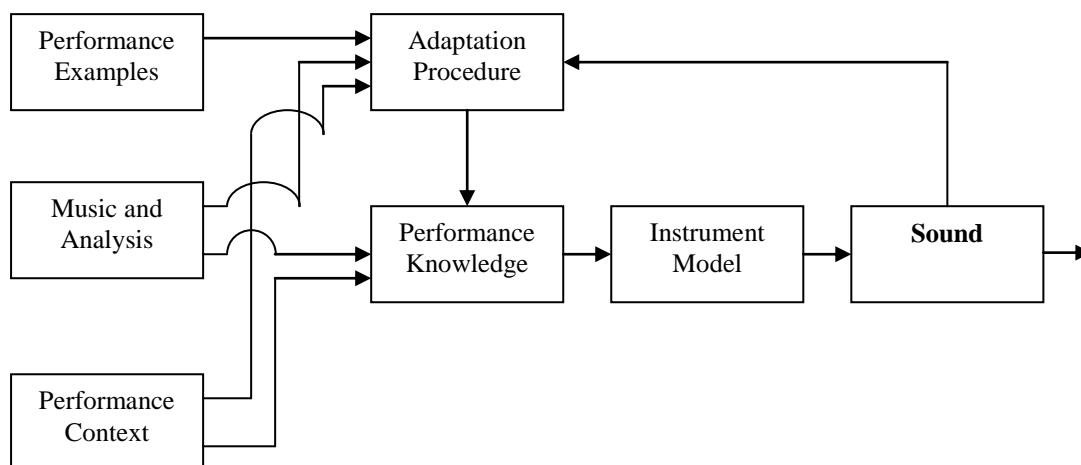


Figure 1: Generic model for most current CSEMPs

Performance Knowledge - This is the core of any performance system. It is the set of rules or associations that controls the performance action. It is the “expertise” of the system which contains the ability, implicit or explicit, to generate an expressive performance. This may be in the form of an Artificial Neural Network, a set of cases in a Case-based Reasoning system, or a set of linear equation with coefficients. To produce performance actions, this module uses its programmed knowledge together with any inputs concerning the particular performance. Its main input is the Music / Analysis module.

Music / Analysis - The Music / Analysis module has two functions. First of all, in all systems, it has the function of inputting the music to be played expressively (whether in paper score, MIDI, MusicXML, audio or other form) into the system. The input process can be quite complex, for example paper score or audio input will require some form of analytical recognition of musical events. This module is the only input to the Performance Knowledge module that defines the particular piece of music to be played. In some systems, it also has a second function – to provide an analysis of the musical structure. This analysis provides information about the Music Features of the music – for example metrical, melodic or harmonic structure. (It was mentioned earlier how it has been shown that such structures have a large influence on expressive performance in humans.) This analysis can then be used by the Performance Knowledge system to decide how the piece should be performed. Analysis methods used in some of the systems include Lerdahl and Jackendoff’s Generative Theory of Tonal Music [Lerdahl and Jackendoff 1983], Narmour’s Implication Realisation [Narmour 1990], and various bespoke musical measurements. The analysis may be automated, manual or a combination of the two.

Performance Context - Another element which will effect how a piece of music is played is the performance context. This includes such things as how the performer decides to play a piece for example happy, perky, sad, or lovelorn. It can also include whether the piece is played in a particular style, e.g. baroque or romantic.

Adaptation Process - The adaptation process is the method used to develop the Performance Knowledge. Like the Analysis module this can be automated, manual or a combination of the two. In some systems, a human expert listens to actual musical output of the performance system and decides if it is appropriate. If not then the Performance Knowledge can be adjusted to try to improve the musical output performance. This is the reason that in Figure 1 there is a line going from the Sound module back to the Adaptation Procedure module. The Adaptation Procedure also has inputs from Performance Context, Music / Analysis, Instrument Model, and Performance Examples. All 4 of these elements can influence the way that a human performs a piece of music, though the most commonly used is Music / Analysis and Performance Examples.

Performance Examples - One important element that can be incorporated in the Performance Knowledge building is the experience of past human performances. These examples can be used by the Adaptation procedure to analyse when and how performance actions are added to a piece of music by human performers. The examples may be a database of marked-up audio recordings, MIDI files together with their source scores, or (in the manual case) a person’s experience of music performance.

Instrument Model - By far the most common instrument used in computer generated performance research is the piano. This is because it allows experiments with many aspects of expression, but requires only a very simple instrument model. In fact the instrument model used for piano is often just the MIDI/media player and soundcard in a PC. Alternatively it may something more complex but still not part of the simulation system, for example a Yamaha Disklavier. However a few simulation systems use non-keyboard instruments, for example Saxophone, and Trumpet. In these cases the issue of a performance is more than just expressiveness. Just simulating a human-like performance, even if it is non-expressive, on these instruments is non-trivial. So systems simulating expressive performance on such instruments may require a relatively complex instrument model in addition to expressive performance elements.

3. A SURVEY OF PREVIOUS RESEARCH IN COMPUTER EXPRESSIVE PERFORMANCE

The review presented here is meant to be representative rather than exhaustive but will cover the majority of published automated and semi-automated CSEMP systems to date. Table I lists the systems reviewed, together with information about their modules. This information will be explained in the detailed part of the review. A number of abbreviations are used in Table I and throughout the paper. Table II lists these abbreviations and their meaning.

CSEMP	Performance Knowledge	Input	Music Feature Analysis	Performance Context	Adaptation Procedure	Performance Examples	Instrument	Performance Actions
Director Musices	Rules	MIDI, Score	Custom	Mood space	-	MIDI Performances	All (Piano)	T/D/A/P
Hierarchical Parabola Model	Parabola equation	MIDI, Score	GTTM TSR	-	-	-	Piano	T/D
Composer Pulse	Multiplier set	MIDI	-	-	Manual	Tapping	All	T/D
Bach Fugue	Rules	Score	Custom	-	Manual	Books, Experts	Keyboard	T/A
Rubato	Operators	MIDI, Score	Custom	-	Manual	-	All (Piano)	T/D
Trumpet Synthesis	Linear model	Audio, Score	Custom	-	Manual	Audio Performances	Trumpet	A/P/O
Linear Regression	Linear model	MIDI, Score	GTTM / Meyer	-	Regression with ANDs	Audio Performances	Piano	T/D/A
ANN Piano	ANN	MIDI	Custom	-	ANN training	MIDI Performances	Piano	T/D
Music Plus One	BBN	Audio / MIDI, Score	Custom	Soloist tempo	BBN training	Audio soloist performances	All	T/D
SaxEx	Fuzzy Rules	Audio, Score	Narmour / IR / GTTM, Custom	Mood space	CBR training	Audio Performances by mood	Saxophone	T/D/V/K
CARO	Linear model	Audio	Custom	Mood space	PCA, Linear Regression	Audio Performances by mood	All	T/K/D/A

Emotional Flute	ANN and rules	Audio, Score	Custom	Mood space	ANN training	Audio Performances by mood	Flute	T/D/V
Kagurame	Rules	MIDI, Score	Custom	Performance conditions	CBR training	MIDI Performances by context	Piano	T/D/A
Ha-Hi-Hun	Rules	MIDI, Score	GTTM TSR	Language Performance conditions	CBR training	MIDI Performances by condition	Piano	T/D
PLCG	Learned rules	MIDI, Score	Custom	-	Meta-Sequential Learning	MIDI Performances	Piano	T/D/A
Phrase-decomposition/ PLCG	Learned rules	MIDI, Score	Custom, Harmonic by musicologist	-	CBR training, Meta-Sequential Learning	MIDI Performances	Piano	T/D/A
DISTALL	CBR	MIDI, Score	Custom, Harmonic by musicologist	-	CBR training	MIDI Performances	Piano	T/D/A
Pop-E	Ruled-based	MIDI, Score	GTTM, Custom	-	-	MIDI Performances	Piano	T/D
ESP Piano	HMM	MIDI, Score	Custom	-	HMM training	MIDI Performances	Piano	T/D/A
Drumming	Non-linear mapping	Audio	Custom	-	KRR, GPR, kNN	Audio Performances	Drums	T
Non-linear Piano System	Non-linear mapping	Worm, Score	Custom	-	KCCA	Performance Worm	Piano	T/D
Genetic Programming	Regression trees	Audio	IR, Custom	-	Genetic programming	Audio performances	Saxophone	T/D/A/N
Sequential Covering GAs	Rule-based	Audio	IR, Custom	-	Sequential covering by GA	Audio performances	Saxophone	T/D/A
Generative Performance GAs	Pulse set	MIDI, Score	LBDM, Kruhman, Melisma	-	GA	None	Piano	T/D
MAS with Imitation	Pulse set	MIDI, Score	LBDM, Kruhman, Melisma	-	Imitation	None	Piano	T/D
Ossia	Fitness rules	MIDI	-	-	-	None	Piano	T/D/P
Music Emotionality	Linear model	MIDI, Score	Custom	Mood space	-	Performances and Experiments	Piano	T/D/P/O
pMIMACS	Rule-based (2)	MIDI	Performance skill	Excitability state	Imitation	None	Piano	T/D

Table I. Systems reviewed.

A	Articulation
ANN	Artificial Neural Network
BBN	Bayesian Belief Network
CBR	Case-based Reasoning
CSEMP	Computer System for Expressive Music Performance
D	Dynamics
DM	Director Musices (KTH System)
EC	Evolutionary Computing
GA	Genetic Algorithm
GP	Genetic Programming
GPR	Gaussian Process Regression
GTTM	Lerdahl and Jackendoff's Generative Theory of Tonal Music
HMM	Hidden Markov Model
IBL	Instance-based Learning
IR	Narmour's Implication/Realisation Theory of Melody
K	Attack
KCCA	Kernel Canonical Correlation Analysis
kNN	k-Nearest Neighbour
KRR	Kernel Ridge Regression
LBDM	Local Boundary Detection Model of Cambouropoulos
MAS	Multi-agent System
MES	Music Emotionality System by Livingstone et al
MIDI	Musical Instrument Digital Interface
MIMACS	Mimetics-inspired Multi-agent Composition System
MusicXML	Music Extended Markup Language
MIS	Music Intepretation System by Katayose et al
N	Note addition/consolidation
P	Pitch
PCA	Principal Component Analysis
T	Tempo
TSR	Time Span Reduction Technique (from GTTM)
V	Vibrato

Table II. Abbreviations

Before discussing our primary terms of reference for this review, we first need to observe that the issue of evaluation of CSEMPs is an open problem. How does one evaluate what is essentially a subjective process? If the CSEMP is trying to simulate a particular performance, then correlation tests can be done. But even if the correlations are low for a generated performance, it is logically possible

for the generated performance to be more preferable to some people than the original performance. [Papadopoulos and Wiggins 1999] discuss the evaluation issue in a different but closely related area - computer algorithmic composition systems. They list four points that they see as problematic in relation to such composition systems:

1. The lack of evaluation by experts, for example professional musicians.
2. Evaluation is a relatively small part of the research with respect to the length of the research paper.
3. Many systems only generate melodies. How do we evaluate the music without a harmonic context? Most melodies will sound acceptable in some context or other.
4. Most of the systems deal with computer composition as a problem solving task rather as a creative and meaningful process.

All of these four points are issues in the context of computer systems for expressive performance as well. So from these observations we will extract three of our primary terms of reference for this review: Performance Testing Status (points 1 and 2), Non-monophonic Ability (point 3) and Performance Creativity (point 4). We will now examine these first three dimensions.

Performance Testing Status refers to how and to what extent the system has been tested. It is important to emphasise that Testing Status is not a measure of how successful the testing was, but how extensive it was. There are 3 main approaches to CSEMP testing: (a) trying to simulate a particular human performance or an average of human performances, (b) trying to create a performance which does not sound machine-like, (c) trying to create as aesthetically pleasing a performance as possible. For the first of these, a correlation can be done between the computer performance and the desired target performance/performances. (However this will not be an “perceptually-weighted” correlation; errors may have a greater aesthetic/perceptual effect at some points than at others.) For approaches (b) and (c) we have listening tests by experts and non-experts. A wide variety of listening tests are used in CSEMPS, from the totally informal and hardly reported, to the results of formal competitions.

Each year since 2002 a formal competition, that has been described as a “musical Turing Test”, called the RenCon (Contest for Performance Rendering Systems) Workshop, has been held [Hiraga et al 2004]. About a third of the systems we will review have been entered into a RenCon competition - see Table III for the results. RenCon is a primarily piano-based competition for Baroque, Classical and Romantic music, and includes manual as well as automated systems (though the placings in Table III are displayed relative to automated and semi-automated CSEMPS, ignoring the manual submissions to RenCon.) Performances are graded and voted on by a jury of attendees from the sponsoring conference. Scores are given for “humanness”, and “expressiveness”, giving an overall “preference” score. It is the preference score we will focus on in the review. The size of RenCon juries, and their criteria have varied over time. In previous years (apart from its first year - 2002) RenCon did not have a separate autonomous section – it had two sections: Compulsory and Open, where compulsory was limited to a fixed piano piece for all contestants, and open was open to all instruments and pieces. In these competitions, automated CSEMPs went up against human pianists and renditions which were carefully

crafted by human hand. Thus many past RenCon results are not the ideal evaluation for automated and semi-automated CSEMPs. However they are the only published common forum available, so in the spirit of points (1) and (2) from Papadopoulos and Wiggins, they will be referred to where possible in our review.

CSEMP	RenCon Placings
Director Musices	2002 (4 th), 2004 compulsory (1 st), 2005 compulsory (2 nd)
SuperConductor (includes Composer Pulse)	2004 open (1 st), 2006 open (1 st and 4 th)
Rubato	2004 open (4 th)
MIS	2002 (2 nd)
Music Plus One	2003 (1 st)
Kagurame	2002 (6 th), 2004 compulsory (3 rd), 2006 compulsory (2 nd)
Ha-Hi-Hun	2002(5 th), 2004 compulsory (4 th), 2006 compulsory (3 rd)
DISTALL	2002 (1 st)
Pop-E	2005 compulsory (1 st), 2006 compulsory (1 st), 2006 open (2 nd and 3 rd)

Table III. RenCon placings of CSEMPS in this Paper

From 2008 onwards the competition has three sections: an “autonomous” section, a “type-in” section and the open section. The autonomous section aims to only evaluate performances rendered by automated CSEMPs. Performances are graded by the composer of the test pieces as well as by a jury. For the autonomous section, the 2008 RenCon contestants are presented with two one-minute pieces of unseen music: one in the style of Chopin and one in the style of Mozart. An award will be presented for the highest scored performance and for the performance most preferred by the composer of the two test pieces. The type-in section is for computer systems for manually generating expressive performance.

The third term of reference in this review, Performance Creativity, refers to the ability of the system to generate novel and original performances, as opposed to simulating previous human strategies. For example the Artificial Neural Network Piano system [Bresin et al 1990][Bresin 1998] is designed to simulate human performances (an important research goal), but not to create novel performances; whereas a system like Director Musices [Friberg et al 2006], although also designed to capture human performance strategies, has a parameterisation ability which can be creatively manipulated to generate entirely novel performances. There is an important proviso here – a system which is totally manual would seem at first glance to have a high creativity potential, since the user could entirely shape every element of the performance. However this potential may never be realised due to the manual effort required to implement the performance. Not all systems are able to act in a novel and practically controllable way. Many of the systems generate a model of performance which is basically a vector or matrix of coefficients. Changing this matrix by hand (“hacking it”) would allow the technically knowledgeable to creatively generate novel performances. However the changes could require too much effort, or the results of such changes could be too unpredictable (thus requiring too many iterations or “try outs”). So performance creativity includes the ability of a system to produce novel performances with a reasonable amount of effort. Having said that, simple controllability is not

the whole of Performance Creativity; for example there could be a CSEMP which has only 3 basic performances rules which can be switched on and off with a mouse click and the new performance played immediately. However the results of switching off and on the rules would in all likelihood generate a very uninteresting performance.

So for performance creativity, a balance needs to exist between automation and creative flexibility, since in this review we are only concerned with automated and semi-automated CSEMPs. An example of such a balance would be an almost totally automated CSEMP, but with a manageable number of parameters that can be user-adjusted before activating the CSEMP for performance. After activating the CSEMP, a performance is autonomously generated but is only partially constrained by attempting to match past human performances. Such creative and novel performance is often applauded in human performers. For example Glenn Gould has created highly novel expressive performances of pieces of music and has been described as having a vivid musical imagination [Church 2004]. Expressive computer performance provides possibilities for even more imaginative experimentation with performance strategies.

We will now add a fourth and final dimension to the primary terms of reference which – like the other three - also has parallels in algorithmic composition. Different algorithmic composition systems generate music with different levels of structural sophistication – for example some may just work on the note-to-note level, like [Kirke and Miranda 2008]. Whereas some may be able to plan at the higher structure level, generating forms like ABCBA, for example [Anders 2007]. There is an equivalent function in computer systems for expressive performance: Expressive Perception. Expressive Perception is the level of sophistication in the CSEMP's perception of the score. We have already mentioned the importance of the music's structure to a human expressive performance. A piece of music can be analysed with greater and greater levels of complexity. At its simplest it can be viewed a few notes at a time; or from the point of view of melody only. At its most complex the harmonic and hierarchical structure of the score can be analysed – as is done in Widmer's DISTALL system [Widmer and Tobudic 2003a; 2003b]. The greater the expressive perception of a CSEMP, the more of the music features it can potentially express.

So to summarise, our 4 primary terms of reference will be:

- Testing Status
- Expressive Perception
- Non-monophonic Ability
- Performance Creativity

At some point in our description of each system, these points will be implicitly or explicitly addressed, and are summarised at the end of the paper (see Table V). It is worth noting that these are not an attempt measure of how successful the system is overall, but an attempt to highlight some key issues which will help to show potential directions for future research.

What now follows is the actual descriptions of the CSEMPs, divided into a number of groups. Each CSEMP is grouped according to how their performance model is built – i.e. by learning method.

This provides a manageable division of the field, shows which learning methods are most popular, and shows where there is room for development in the building of performance models. The grouping will be:

1. Non-learning (10 systems)
2. Rule/Case-based learning (6 systems)
3. Linear regression (2 systems)
4. Non-linear regression (2 systems)
5. Artificial Neural Networks (2 systems)
6. Statistical Graphical Models (2 systems)
7. Evolutionary computation (6 systems)

3.1 Non-learning Systems

3.1.1 Director Musices. Director Musices (DM) [Sundberg et al. 1983; Friberg et al. 2006] has been an ongoing project since 1982. Researchers including violinist Lars Fryden developed and tested performance rules using an analysis-by-synthesis method (later using analysis-by-measurement and studying actual performances). Currently there are around 30 rules which are written as relatively simple equations that take as input Music Features such as height of the current note pitch, the pitch of the current note relative to the key of the piece, or whether the current note is the first or last note of the phrase. The output of the equations defines the Performance Actions. For example the higher the pitch the louder the note is played, or during an upward run of notes, play the piece faster. Another DM rule is the Phrase Arch which defines a “rainbow” shape of tempo and dynamics over a phrase. The performance speeds up and gets louder towards the centre of a phrase and then tails off again in tempo and dynamics towards the end of the phrase. Some manual score analysis is required – for example harmonic analysis and marking up of phrase start and ends. DM’s ability for expressive perception is at the note and phrase level – it does not use information at higher levels of the musical structure hierarchy.

Each equation has a numeric “k-value” - the higher the k-value the more effect the rule will have and a k-value of 0 switches the rule off. The results of the equations are added together linearly to get the final performance. Thanks to the adjustable k-value system, DM has much potential for performance creativity. Little work has been reported on an active search for novel performances, though it is reported that negative k-values reverse rule effects and cause unusual performances. DM’s ability as a semi-automated system comes from the fact it has a “default” set of k-values, allowing the same rule settings to be applied automatically to different pieces of music (though not necessarily with the same success).

Rules are also included for dealing with non-monophonic music [Friberg et al. 2006]. To understand more deeply the issues raised by polyphonic expression, consider that each voice of an ensemble has its own melodic structure. Many monophonic methods described in our survey would lead to a number of voices each being given their own expressive timing deviations, causing desynchronization and cacophony. In DM, the “Melodic-sync” rule generates a new voice consisting of

all timings in all other voices (if two voices have simultaneous notes, then the note with the greatest melodic tension is selected.) Then all rules are applied to this synchronisation voice, and resulting durations are mapped back onto the original voices. The “Bar-sync” rule can also be applied to make all voices re-synchronise at each bar end.

DM is also able to deal with some Performance Contexts, specifically emotional expression [Bresin and Friberg 2000], drawing on work by Gabrielsson and Juslin [1996]. Listening experiments were used to define the k-value settings on the DM rules for expressing emotions. The music used was a Swedish nursery rhyme and a computer-generated piece in a minor mode written using Cope’s [1992] algorithmic composition system in the musical style of Chopin. Six rules were used from DM to generate multiple performances of each piece. Subjects were asked to identify a performance emotion from the list: fear, anger, happiness, sadness, solemnity, tenderness or no-expression. As a result parameters were found for each of the 6 rules which mould the emotional expression of a piece. For example for “tenderness”: inter-onset interval is lengthened by 30%, sound level reduced by 6dB, and two other rules are used: the Final Ritardando rule (slowing down at the end of a piece) and the Duration Contrast rule (if two adjacent notes have contrasting durations, increase this contrast).

Director Musices has a good test status, having been evaluated in a number of experiments. In [Friberg 1995] k-values were adjusted by a search algorithm, based on 28 human performances of 9 bars of of Schumann’s *Träumerei*. A good correlation was found between the human performances and the resulting DM performance. Another experiment involved manually fitting to one human performance the first 20 bars of the *Adagio* in Mozart’s sonata K.332 [Sundberg et al 2003]. The correlations were found to be low, unless the k-values were allowed to change dynamically when the piece was performed. An attempt was made to fit k-values using to a larger corpus of piano music using Genetic algorithms in [Kroiss 2000], and the results were found to give a low correlation as well. In an attempt to overcome this [Zanon and De Poli 2003] allowed k-values to vary in a controlled way over a piece of music. This was tested on Beethoven’s *Sonatine in G Major* and Mozart’s K.332 piano sonata (the slow movement) – but the results were found to be poor for the Beethoven. In the first RenCon in 2002, the second prize went to a DM rendering, however the first placed system (a manually rendered performance) was voted for by 80% of the jury. In RenCon 2005, a Director Musices default-settings (i.e. automated) performance of Mozart’s *Minuette KV 1(1e)* came a very close 2nd in the competition, behind Pop-E (Section 3.1.7). However 3 of the other 4 systems competing were versions of the DM-system.

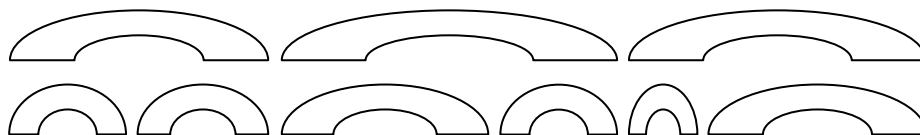
So the results have been mixed, perhaps because DM is a “building-block” approach to music performance, developed more out of a desire to learn about the building blocks of music performance, rather than finding an efficient way of putting them together. The DM model has been influential, and as will be seen in the later systems, DM-type rules appear again and again.

3.1.2 Hierarchical Parabola Model. One of the first CSEMPs with a hierarchical expressive perception was Todd’s Hierarchical Parabola Model [Todd 1985; Todd 1989; Todd 1992; Todd 1995]. Todd argues it was consistent with a kinematic model of expressive performance, where tempo changes are viewed as being due to accelerations and decelerations in some internal process in the human

mind/body, for example the auditory system. For tempo the hierarchical parabola model uses a rainbow shape like DM’s phrase arch, which is consistent with Newtonian kinematics. For loudness the model uses a “the faster the louder” rule, creating a dynamics rainbow as well.

The key difference between DM and this hierarchical model is that the hierarchical model has greater expressive perception and wider performance action. Multiple levels of the hierarchy are analysed using Lerdahl and Jackendoff’s Generative Theory of Tonal Music (GTTM). GTTM Time Span Reduction (TSR) examines each note’s musicological place in all hierarchical levels. The rainbows/parabolas are generated at each level, from the note-group level upwards (Figure 2) and added to get the performance. This generation is done by a parametrized parabolic equation which takes as input the result of the GTTM TSR analysis.

Tempo Parabolas



Note Groupings

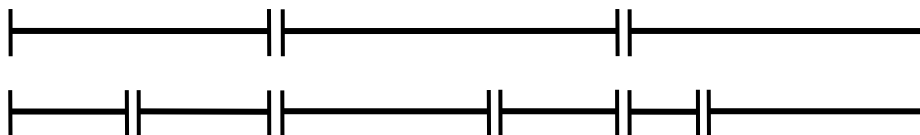


Figure 2: Todd Parabola Model

The performance was shown to correlate well by eye with a short human performance but no correlation figures were reported. [Clarke and Windsor 2000] tested the first four bars of Mozart’s K.331; comparing two human performers with two performances by the Hierarchical Parabola model. Human listeners found the Parabola version unsatisfactory compared to the human ones. In the same experiment however, the Parabola model was found to work well on another short melody. The testing also showed that the idea of “the louder the faster” did not always hold. [Desain and Honing 1993] claim through informal listening tests that in general the performances do not sound convincing.

Like DM, the Hierarchical Parabola Model is not purely designed to produce pleasant performances. It attempts to investigate a fundamental part of the expressive performance – hierarchy expression. The constraint of utilising the hierarchy and the GTTM TSR approach limits the Performance Creativity. Note groupings will be limited to those generated by a GTTM TSR analysis, and the parabolas generated will be constrained by the model’s equation. Any adjustments to a performance will be constrained to working within this framework.

3.1.3 Composer Pulse and Predictive Amplitude Shaping. Manfred Clynes’ Composer Pulse [Clynes 1986] also acts on multiple levels of the hierarchy. Clynes hypothesises each composer has a unique pattern of amplitude and tempo variations running through performances – a pulse. This is captured as a set of numbers multiplying tempo and dynamics values in the score. It is hierarchical with separate

values for within the beat, the phrase and at multiple bar level. Table III shows the values of pulses for phrase level for some composers. The pulses were measured using a sentograph to generate pressure curves from musicians tapping their finger whilst thinking of or listening to a specific composer. Figure 3 shows the structure of a pulse set in three-time (each composer has a three-time and a four-time pulse set defined). This pulse set is repeatedly applied to a score end on end. So if the pulse is 12 beats long and the score is 528 beats, the pulse will repeat $528/12 = 44$ times end on end.

Level 2 Composers' Pulses - 4 Pulse					
Beethoven	Duration	106	89	96	111
	Amplitude	1.00	0.39	0.83	0.81
Mozart	Duration	105	95	105	95
	Amplitude	1.00	0.21	0.53	0.23
Schubert	Duration	97	114	98	90
	Amplitude	1.00	0.65	0.40	0.75
Haydn	Duration	108	94	97	102
	Amplitude	1.00	0.42	0.68	1.02
Schumann	Duration	96	116	86	102
	Amplitude	0.60	0.95	0.50	1.17
Mendelssohn	Duration	118	81	95	104
	Amplitude	1.00	0.63	0.79	1.12

Table III. Level 2 Composers' Pulses

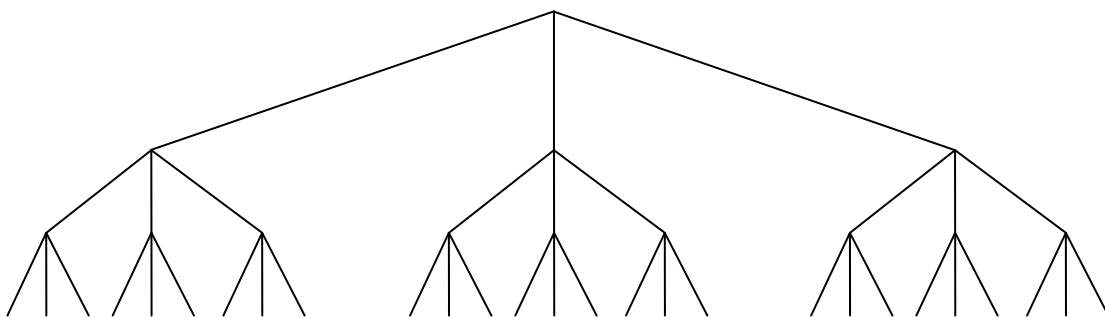


Figure 3. Structure of a pulse set in three-time

Another key element of Clyne's approach is Predictive Amplitude Shaping. This adjusts a note's dynamic based on the next note simulating "a musician's unconscious ability to sculpt notes in this way" that "makes his performance flow beautifully through time, and gives it meaningful coherence even as the shape and duration of each individual note is unique." A fixed envelope shape model is used (some constants are manually defined by the user), the main inputs being distance to the

next note and duration of the current note. So the Pulse/Amplitude system has only note level expressive perception.

Clynes' test of his own model [Clynes 1995] showed that a number of expert and non-expert listeners preferred music with a composers pulse than with a different pulse. However not all tests on Clynes' approach have supported a universal pulse for each composer [Thompson 1989; Repp 1990], suggest instead that the pulse may be effective for a subset of a composer's work. Clynes' pulses and amplitude shaping have been combined with other performance tools (e.g. vibrato generation) as part of his commercial software SuperConductor. Two SuperConductor generated performances were submitted to RenCon 2006 open section: Beethoven's Eroica Symphony, Op.55, Mvt.4 and Brahms' Violin Concerto, Op.77, Mvt.1. The Beethoven piece scored low, but the Brahms piece came 1st in the open section (beating two pieces submitted by Pop-E - Section 3.1.7). The generation of this piece could have involved significant amounts of manual work were required. Also because it was the open section, the pieces submitted by Pop-E were not the same as submitted by SuperConductor – hence like was not compared to like. SuperConductor also won the open section in RenCon 2004 with J. S. Bach, Brandenburg Concerto No.5, D Major, 3rd Movement. The only competitor included from this review was Rubato (Section 3.1.6) performing a Bach piece. It should be re-emphasised that these results were for SuperConductor and not solely for the Pulse and Amplitude tools.

In the context of SuperConductor, Clynes approach allows for significant Performance Creativity. The software is designed to allow a user to control the expressive shaping of a MIDI performance, giving significant amounts of control. However, outside of the context of SuperConductor, the pulse has little scope for performance creativity – though the amplitude shaping does. The pulse and amplitude shaping do not explicitly address non-monophonic music, though SuperConductor can be used to generate polyphonic performances.

3.1.4 Bach Fugue System. In the Bach Fugue System [Johnson 1991] Expert System methods are used to generate performance actions. Johnson generated the knowledge base through interviews with two musical expert performers, and through a performance practice manual and an annotated edition of the Well-Tempered Clavier; so this system is not designed for performance creativity. Twenty eight conditions for tempo and articulation are so-generated for the knowledge base. For example: “If there is any group of 16th notes following a tied note, then slur the group of 16th notes following the long note.” Expressive perception is focused on the note to phrase level. The CSEMP does not perform itself, but generates instructions for 4/4 fugues. So testing was limited to examining the instructions. It gave the same instructions as human experts 85 to 90% of the time, though it is not said how many tests were run. The system is working in the context of polyphony.

3.1.5 Trumpet Synthesis. Three out of the last four CSEMPs reviewed have focused on keyboard, and this pattern will continue through the paper - most CSEMPs focus on the piano because it is easier to collect and analyze data for the piano than for other instruments. One of the first non-piano systems was Dannenberg and Derenyi's Trumpet Synthesis [Dannenberg and Derenyi 1998; Dannenberg et al. 1998]. The authors' primary interest here was to generate realistic trumpet synthesis, and adding

performance factors improves this synthesis. It is not designed for performance creativity but for simulation. This trumpet system synthesizes the whole trumpet performance, without needing any MIDI or audio building blocks as the basis of its audio output. The performance actions are amplitude and frequency, and these are controlled by envelope models which were developed using a semi-manual statistical analysis-by-synthesis method. A 10-parameter model was built for amplitude, based on elements such as articulation, direction and magnitude of pitch intervals, and duration of notes. This system works by expressively transforming one note at a time, based on the pattern of the surrounding two notes. In terms of expressive perception the system works on a 3 note width. The pitch expression is based on envelopes which were derived and stored during the analysis-by-synthesis.

No test results are reported. Dannenberg and Derenyi placed two accompanied examples online: parts of a Haydn Trumpet Concerto and of a Handel Minuet. The start of the trumpet on the Concerto without accompaniment is also online, together with a human playing the same phrase. The non-accompanied synthesis sounds quite impressive, only being let down by a synthetic feel towards the end of the phrase – though the note-to-note expression (as opposed to the synthesis) consistently avoids sounding machine-like. In both accompanied examples it became clear as the performances went on that a machine was playing, particularly in faster passages. But once again note-to-note expression did not sound too mechanical. Despite the reasonably positive nature of these examples, there is no attempt to objectively qualify how good the trumpet system is.

3.1.6 Rubato. Mazzola, a mathematician and recognised Jazz pianist, developed a mathematical theory of music [Mazzola 1994][Mazzola 2002]. Music is represented in an abstract geometrical space whose co-ordinates include onset time, pitch, and duration. A score will exist in this space, and expressive performances are generated by performing transformations on the space. The basis of these transformations are a series of “Operators” which can be viewed as a very generalised version of the rule-based approach taken in Director Musices. For example, the Tempo operator and the Split operator allow the generation of tempo hierarchies. These give Rubato a good expressive perception. However the definition of the hierarchy here differs somewhat from that found in the Hierarchical Parabola Model or DISTALL (Section 3.2.6). A tempo hierarchy, for a piano performance, may mean that the tempo of the left hand is the dominant tempo, at the top of a hierarchy, and the right hand tempo is always relative to the left hand tempo – and so is viewed as being lower in the hierarchy. Mazzola also discusses the use of tempo hierarchies to generate tempo for grace notes and arpeggios – the tempo of these is relative to some global tempo higher in the hierarchy. Ideas from this theory have been implemented in a piece of software called Rubato, which is available online. The expressive performance module in Rubato is the “Performance Rubette”. A MIDI file can be loaded in Rubato and pre-defined operators used to generate expressive performances. The user can also manually manipulate tempo curves using a mouse and GUI, giving Rubato good scope for performance creativity.

Test reports are limited. In RenCon 2004, a performance of Bach’s Contrapunctus III modelled using Rubato was submitted, and came 4th in the open section (SuperConductor came 1st in the section with a different piece). It is not clear how automated the generation of the performance was. Listening to the submission it can be heard that although the individual voices are quite expressive and

pleasant (except for the fastest parts), the combination sounds relatively unrealistic. An online MIDI example is available of Schumann's *Kindersezenen* op. 15 Nr. 2, "Kuriose Geschichte" which evidences both tempo and dynamics expression and is quite impressive, though once again it is not clear how automated the production of the music was.

3.1.7 Pop-E. Pop-E [Hashida et al. 2006], a Polyphrase Ensemble system, was developed by some of the team involved in MIS (Section 3.3.1). It applies expression features separately to each voice in a MIDI file, through a synchronisation algorithm. The music analysis uses GTTM local level rules, and utilizes beams and slurs in the score to generate note groupings. So the expressive perception is up to phrase level. Expressive actions are applied to these groupings through rules reminiscent of Director Musices. The 5 performance rules have a total of 9 manual parameters between them. These parameters can be adjusted, providing scope for performance creativity. In particular jPop-E [Hashida et al 2007], a java implementation of the system provides such tools for shaping new performances.

To deal with polyphony, Synchronisation Points are defined at the note grouping start and end points in the attentive part. The attentive part is that voice which is most perceptually prominent to a listener. The positions of notes in all other non-attentive parts are linearly interpolated relative to the synchronisation points (defined manually). This means that all parts will start and end at the same time at the start and end of groupings of the main attentive part.

Pop-E was evaluated in the laboratory to see how well it could reconstruct specific human performances. After setting parameters manually, performances by three pianists were reconstructed. The average correlation values between Pop-E and a performer were 0.59 for tempo and 0.76 for dynamics. This has to be viewed in the context that the average correlations between the human performers were 0.4 and 0.55 respectively. Also the upper piano part was more accurate on average. (It is interesting to note that for piano pieces whose attentive part is the right hand, the Pop-E synchronisation system is similar to the methods in the DISTALL system for dealing with polyphony – see Section 3.2.6.) Pop-E won the RenCon 2005 compulsory section, beating Director Musices. In RenCon 2006 Pop-E won the compulsory section beating Kagurame (Section 3.2.2) and Ha-Hi-Hun (Section 3.2.3). In the open section in 2006 SuperConductor beat Pop-E with one performance, and lost to Pop-E with another.

3.1.8 Hermode Tuning. In the next two sub-sections we will describe commercial CSEMPs. Despite the lack of details available on these proprietary systems we felt they should be included here since they are practical CSEMPs that people are paying money for, and illustrate some of the commercial potential of CSEMPs for the music business. However because of the lack of some details we will not attempt to apply the four review terms of reference. The first system is Hermode Tuning [Sethares 2004]. Most systems in this review focus on dynamics and timing. However intonation is another significant area of expression for many instruments – for example many string instruments. (In fact, three intonation rules were added to Director Musices in its later incarnations; for example, the higher the pitch, the sharper the note.) Hermode Tuning is a dedicated expressive intonation system which can work in real time, its purpose being to “imitate the living intonation of well educated instrumentalists in

orchestras and chamber music ensembles.” Instrumentalists do not perform in perfect intonation – in fact if an orchestra performed music in perfect tuning all the time the sound would be less pleasant than one that optimised its tuning through performance experience. A series of algorithms are using in Hermode tuning not just to avoid perfect intonation but to attempt to achieve optimal intonation. The algorithms have settings for different types of music, for example Baroque and Jazz/Pop. Examples are available on the website, and the system has been successful enough to be embedded in a number of commercial products – for example Apple Logic Pro 7.

3.1.9 Sibelius. As mentioned in the introduction of this paper, the music typesetting software package Sibelius has built-in algorithms for expressive performance. These use a rule-based approach. Precise details are not available for these commercial algorithms but some information is available [Finn 2007]. For dynamics, beat groups such as barlines, sub-bar groups and beams are used to add varying degrees of stress. Also the higher the note is the louder it is played, though volume resets at rests and dynamic expression is constrained to not be excessive. Some random fluctuation is added to dynamics to make it more human-sounding as well. Tempo expression is achieved using a simple rubato system; however this rubato does not include reliable phrase analysis. The manufacturer reports that “phrasing need only be appropriate perhaps 70% of the time – the ear overlooks the rest” and that “the ear is largely fooled into thinking it’s a human performance.” Notion and Finale also have expressive performance systems built into them, which are reportedly more advanced than Sibelius’, but even fewer details are available for the proprietary methodologies in these systems.

3.1.10. Music Emotionality System. In relation to the philosophy behind the Music Emotionality System (MES) [Livingstone et al 2007], Livingstone observes that “the separation of musical rules into structural and performative is largely an ontological one, and cedes nothing to the final audio experienced by the listener.” The Music Emotionality System has a rule-set of 19 rules developed through analysis-by-synthesis. The rules have an expressive perception up to the phrase level, some requiring manual mark-up of the score. These rules are designed not only to inject microfeature deviations into the score to generate human-like performances, but also to use microfeature and macrofeature deviations to express emotions to the listener. To this end MES is able to change the score itself, recomposing it.

MES has a 2-D model of human emotion space with four quadrants going from very active and negative to very active and positive, to very passive and positive through to very passive and negative. These four elements combine to give such emotions as angry, bright, contented and despairing. The quadrants were constructed from a review of 20 studies of music and emotion. The rules for expressing emotions include: moving between major and minor modes, changing note pitch classes, as well as DM-type rules for small changes in dynamics and tempo. It was found that the addition of the microfeature humanisation rules improved the accuracy of the emotional expression (as opposed to solely using macrofeature “recomposition” rules). The rules for humanising the performance include some rules which are similar to Director Musices, such as Phrase Arch and emphasising metrically important beats. Creative Performance is possible in MES by adjusting the

parameters of the rule set, and the emotional specification would allow a user to specify different emotions for different parts of a performance.

A significant number of formal listening tests have been done by Livingstone and they support the hypothesis that MES is more successful than DM at expressing emotions. MES is one of the better tested systems in this review – one reason being that its aim is more measurable than a purely aesthetic goal. Examples of MES are available on the author's webpage.

3.2 Case and Instance-based Systems

Case-based Reasoning is the first learning method we will address. Learning CSEMPs can incorporate more knowledge more quickly than non-learning systems. However such methods do not always provide tools for creative performance because they are strongly rooted in past performances. Before continuing we should explain that any CSEMP that learns expressive deviations needs to have a non-expressive reference point, some sort of representation of the music played robotically/neutrally. The CSEMP can then compare this to the score played expressively by a human, and learn the deviations.

3.2.1 SaxEx. Arcos and Lopez de Mantaras' SaxEx [Arcos et al. 1997; Arcos et al. 1998; Arcos and Mantaras 2001a; Arcos and Mantaras 2001b; Mantaras and Arcos 2002] was one of the first systems to learn performances based on the performance context of mood. Like the trumpet system described earlier (Section 3.1.5), SaxEx includes algorithms for extracting notes from audio files, and generating expressive audio files from note data. SaxEx also looks at intranote features like vibrato and attack. Unlike the trumpet system, SaxEx needs a non-expressive audio file to perform transformations upon. Narmour's IR theory is used to analyse the music. IR considers features of the previous two notes in the melody, and postulates that a human will expect the melody to move in a certain direction and distance; thus it can classify each note as being part of a certain expectation structure. Other elements used to analyse the music are ideas from Jazz theory, as well as GTTM TSR. This system's expressive perception is up to phrase level and is automated.

SaxEx was trained on cases from monophonic recordings of a tenor sax playing 4 Jazz standards with different moods (as well as a non-expressive performance). The moods are designed around three dimensions: tender-aggressive, sad-joyful, and calm-restless. The mood and local IR, GTTM and Jazz structures around a note are linked to the expressive deviations in the performance of that note. These links are stored as performance cases. SaxEx can then be given a non-expressive audio file and told to play it with a certain mood. A further AI method is used then to combine cases: Fuzzy Logic. For example - if two cases are returned for a particular note in the score and one says play with low vibrato, and the other says play with medium vibrato, then fuzzy logic combines them into a low-medium vibrato. The learning of new CBR solutions can be done automatically or manually through a GUI, which affords some performance creativity giving the user a stronger input to the generation of performances. Though this is limited by SaxEx's focus on being a simulation system. There is, like the Music Emotionality System, the potential for the user to generate a performance with certain moods at different points in the music.

There is no formal testing reported, but SaxEx examples are available online. The authors report “dozens” of positive comments about the realism of the music from informal listening tests, but no formal testing is reported or details given. The two short examples online (Sad and Joyful) sound realistic to us, more so than - for example - the trumpet system examples. But the accuracy of the emotional expression was difficult for us to gauge.

3.2.2 Kagurame. Kagurame [Suzuki et al. 1999; Suzuki 2003] is another case-based reasoning system which – in theory - also allows expressiveness to be generated from moods, this time for piano. However it is designed to incorporate a wider degree of performance conditions than solely mood, for example playing in a Baroque or Romantic style. Rather than GTTM and IR, Kagurame uses its own custom hierarchical note structures to develop and retrieve cases for expressive performance. This hierarchical approach gives good expressive perception. Score analysis automatically divides the score into segments recursively with the restriction that the divided segment must be shorter than one measure. Hence manual input is required for boundary information for segments longer than one measure. The score patterns are derived automatically after this, as is the learning of expressive actions associated with each pattern. Kagurame acts on timing, articulation, and dynamics. There is also a polyphony action called Chord Time Lag –notes in the same chord can be played at slightly different times. It is very much a simulation system with little scope for creative performance.

Results are reported for monophonic Classical and Romantic styles. Tests were based on learning 20 short Czerny etudes played in each style. Then a 21st piece was performed by Kagurame. Listeners said it “sounded almost human like, and expression was acceptable” and that the “generated performance tended to be similar to human, particularly at characteristic points.” A high percentage of listeners guessed correctly whether the computer piece was Romantic or Classical style. In RenCon 2004 Kagurame came 4th in one half of the compulsory section, one ahead of Director Musices, but was beaten by DM in the second half, coming 5th. At RenCon 2006 a polyphonic performance of Chopin's piano Etude in E major came 2nd (with Pop-E taking 1st place).

3.2.3 Ha-Hi-Hun. Ha-Hi-Hun [Hirata and Hiraga 2002] utilizes data structures designed to allow natural language statements to shape performance conditions (these include data structures to deal with non-monophonic music). The paper focuses on instructions of the form “render performance of piece X in the style of an expressive performance of piece Y”. As a result, there are significant opportunities for performance creativity through rendering a piece in the style of a very different second piece; or perhaps performing the second piece bearing in mind that it will be used to generate creative performances of the first piece. The music analysis of Ha-Hi-Hun uses GTTM TSR to highlight the main notes that shape the melody. TSR gives Ha-Hi-Hun an expressive perception above note level. The deviations of the main notes in the piece Y relative to the score of Y are calculated, and can then be applied to the main notes in the piece X to be performed by Ha-Hi-Hun. After this the new deviations in X's main notes are propagated linearly to surrounding notes like “expressive ripples” moving outwards. The ability of Ha-Hi-Hun to automatically generate expressive performances comes from its ability to generate a new performance X based on a previous human performance Y.

In terms of testing, two pieces were rendered, each in the style of another piece and formal listening results were reported as positive, but few experimental details are given. In RenCon 2002 Ha-Hi-Hun learned to play Chopin Etude Op.10, No. 3 through learning the style of a human performance of Chopin's Nocturne Op. 32, No. 2. The performance came 9th out of 10 submitted performances by other CSEMPs (many of which were manually produced). In RenCon 2004, Ha-Hi-Hun came last in the compulsory section, beaten by both Director Musices and Kagurame (Section 3.2.2). In RenCon 2006, a performance by Ha-Hi-Hun also came third out of six in the compulsory section, beaten by Pop-E (Section 3.1.7) and Kagurame.

3.2.4 PLCG System. Gerhard Widmer has applied various versions of a rule-based learning approach, attempting to utilise a larger database of music than previous CSEMPs. The PLCG system [Widmer 2000; Widmer 2002; Widmer 2003] uses data mining to find large numbers of possible performance rules and cluster each set of similar rules into an average rule. This is a system for musicology and simulation rather than one for creative performance. PLCG is Widmer's own meta-learning algorithm – the underlying algorithm being Sequential Covering [Mitchell 1997]. PLCG runs a series of sequential covering algorithms in parallel on the same monophonic musical data, and gathers the resulting rules into clusters, generating a single rule from each cluster. The data set was thirteen Mozart Piano sonatas performed by Roland Batik in MIDI form (only melodies were used - giving 41,116 notes). A note-level structure analysis learns to generate tempo, dynamics and articulation deviations based on the local context – e.g. size and direction of intervals, durations of surrounding notes, and scale degree. So this CSEMP has a note level expressive perception. As a result of the PLCG algorithm, 383 performance rules were turned into just 18 rules. Interestingly, some of the generated rules had similarities to some of the Director Musices rule set.

Detailed testing has been done on the PLCG, including its “Generalisation” ability (the ability for the system to perform music or composers that weren't explicitly included in its learning). Widmer's systems are the only CSEMPs reviewed here that have had any significant generalisation testing. The testing methods were based on correlation approaches. Seven pieces of the learning scores were regenerated using the rule set, and their tempo/dynamics profiles compared to the original performances very favourably. Regenerations were compared to performances by a different human performer Phillipe Entremont and showed no degradation relative to the original performer comparison. The rules were also applied to some music in a romantic style (two Chopin pieces), giving encouraging results. There are no reports of formal listening tests.

3.2.5 Combined Phrase-decomposition/PLCG. The above approach was extended by Widmer and Tobudic into a monophonic system whose expressive perception extends into higher levels of the score hierarchy. This was the combined Phrase-decomposition/PLCG system [Widmer and Tobudic 2003]. (Once again this is a simulation system rather than one for creative performance.) When learning, this CSEMP takes as input scores that have had their hierarchical phrase structure defined to three levels by a musicologist (who also provides some harmonic analysis), together with an expressive MIDI performance by a professional pianist. Tempo and Dynamics curves are calculated from the MIDI

performance, and then the system does a multi-level decomposition of these expression curves. This is done by fitting quadratic polynomials to the tempo and dynamics curves (like the curves in Todd's Parabola Model in Section 3.1.2).

Once the lowest level fitting has been done, there is still a "residual" expression. This is hypothesised as being due to note-level expression, and the PLCG algorithm is run on the residuals to learn the note-level rules which generate this residual expression. The learning of the non-PLCG tempo and dynamics is done using a case-based learning type method - by a mapping from multiple-level features to the parabola/quadratic curves. An extensive set of music features are used including: length of the note group, melodic intervals between start and end notes, where the pitch apex of the note group is, whether the note group ends with a cadence, and the progression of harmony between start, apex and end. This CSEMP has the most sophisticated expressive perception of all the systems described in our review.

To generate an expressive performance of a new score, the system moves through the score and in each part runs through all its stored features vectors learned from the training; it finds the closest one using a simple distance measure. It then applies the curve stored in this case to the current section of the score. Data for curves at different levels, and results of the PLCG, are added together to give the expression performance actions.

A battery of correlation tests were performed. Sixteen Mozart sonatas were used to test the system – training on 15 of them and then testing against the remaining one. This process was repeated independently selecting a new 1 of the 16 and then re-training on the other 15. This gave a set of 16 results which the authors described as "mixed". Dynamics generated by the system correlated better with the human performance than a non-expressive performance curve (straight line) did in 11 out of 16, this was only case for 6 out of 16 for the timing curves. There are no reports of formal listening tests.

3.2.6 DISTALL system. Widmer and Tobudic did further work to improve the results of the Combined Phrase-decomposition/PLCG, developing the DISTALL system [Widmer and Tobudic 2003a; 2003b] for simulation. The learned performance cases in the DISTALL system are hierarchically linked, in the same way as the note groupings they represent. So when the system is learning sets of expressive cases, it links together the feature sets for a level 3 grouping with all the level 2 and level 1 note groupings it contains. When a new piece is presented for performance, and the system is looking at a particular level 3 grouping of the new piece, say X – and X contains a number of level 2 and level 1 subgroupings - then not only are the score features of X compared to all level 3 cases in the memory, but the subgroupings of X are compared to the subgroupings of the compared level 3 cases as well. There have been measures available which can do such a comparison in case-based learning before DISTALL (e.g. RIBL [Emde and Wettschereck 1996]). However DISTALL does it in a way more appropriate to expressive performance – giving a more equal weighting to subgroupings within a grouping, and giving this system a high expressive perception.

Once again correlation testing was done with a similar set of experiments to Section 3.2.5. All 16 generated performances had smaller dynamics errors relative to the originals than a robotic/neutral

performance had. For tempo, 11 of the 16 generated performances were better than a robotic/neutral performance. Correlations varied from 0.89 for dynamics in Mozart K283 to 0.23 for tempo in Mozart K332. The mean correlation for dynamics was 0.7 and for tempo was 0.52. A performance generated by this DISTALL system was entered into RenCon 2002. The competition CSEMP included a simple accompaniment system where dynamics and timing changes calculated for the melody notes were interpolated to allow their application to the accompaniment notes as well. Another addition was a simple heuristic for rendering grace notes: the sum of durations of all grace notes for a main note is set equal to 5% of the main note's duration, and the 5% of duration is divided equally amongst the grace notes. The performance was the top scored automated performance at RenCon 2002 - ahead of Kagurame, MIS (Section 3.3.1), and Ha-Hi-Hun - and it beat one non-automated system.

3.3 Linear Regression

Linear Regression models assume a basically linear relationship between the Music Features and the Expressive Actions. The advantage of such models is their simplicity, the disadvantage is that assuming music expressive performance a linear process is almost certainly an oversimplification.

3.3.1 Music Interpretation System. The Music Interpretation Systems (MIS) [Aono et al 1997][Ishikawa et al. 2000] generates expressive performances in MIDI format, but learns expressive rules from audio recordings. This is done using a spectral analysis system with dynamic programming for note detection. The system is a simulatory CSEMP and uses a set of linear equations which map score features on to performance deviation actions. Its expressive perception is on the note and phrase level. MIS has methods to include some non-linearities using logical ANDs between music features in the score, and a way of reducing redundant music features from its equations. This redundancy reduction improves generalisation ability. MIS learns links between music features and performance actions of tempo, dynamics, and articulation. The music features used include score expression marks, and aspects of GTTM and two other forms of musicological analysis: Leonard Meyer's Theory of Musical Meaning [Meyer 1957] and Narmour's IR Theory. Meyer's Theory - like Narmour's - is an expectation-based approach, but coming from the perspective of Game Theory.

For testing, MIS was trained on the first half of a Chopin Waltz and then used to synthesize the second half. Correlations (accuracies when compared to a human performance of the second half) were: for velocity 0.87, for tempo 0.75, and for duration 0.92. A polyphonic MIS interpretation of Chopin, Op. 64, No. 2 was submitted to RenCon 2002. It came 3rd behind DISTALL beating 3 of the other 4 automated systems (DM, Kagurame, Ha-Hi-Hun).

3.3.2 CARO. CARO [Canazza et al. 2000a; Canazza et al. 2000b; Canazza et al. 2003; Canazza et al. 2004; de Poli 2004] is a monophonic CSEMP designed to generate audio files which - like SaxEx and MES - express certain moods/emotions. It does not require a score to work from, but works on audio files which are mood-neutral. The files are however assumed to include the performer's expression of the music's hierarchical structure. Its expressive perception is at the local note level. CARO's performance actions at the note and intra-note level include changes to inter-onset interval, brightness,

and loudness-envelope centroid. A linear model is used to learn actions - every action has an equation characterised by parameters called Shift and Range Expansion. Every piece of music in a particular mood has its own set of Shift and Range Expansion values. This limits the generalisation potential.

CARO also learns “how musical performances are organised in the listener’s mind” in terms of moods: hard, heavy, dark, bright, light and soft. To do this, a set of listening experiments analysed by Principal Component Analysis (PCA) generates a two dimensional space that captures 75% of the variability present in the listening results; this space is used to represent listeners’ experience of the moods. A further linear model is learned for each piece of music which maps the mood space onto Shift and Range Expansion values. The user can select any point in the mood space, and CARO generates an expressive version of the piece. A line can be drawn through mood space, and following that line in time CARO can generate a performance morphing through different moods. Apart from the ability to adjust Shift and Range expansion parameters manually, CARO’s potential for creative performance is extended by its ability to have a line drawn through the mood space. Users can draw trajectories through this space to create entirely novel performances.

For testing, 20 second clips each from 3 piano pieces by different composers were used. A panel of 30 listeners evaluated CARO’s ability to generate pieces with different expressive moods. Results showed that the system gave a good modelling of expressive mood performances as realised by human performers.

3.4 Non-linear Regression

3.4.1 Drumming System. So far, we have surveyed pitched instruments – piano, saxophone and trumpet. We will now look at a system for non-pitched (drumming) expression. In the introduction we talked about how pop music had enthusiastically utilised the “robotic” aspects of MIDI sequencers. However eventually pop musicians wanted a more realistic sound to their electronic music, and Humanization systems were developed for drum machines that added random tempo deviations to beats. Later systems also incorporated what are known as Grooves - a fixed pattern of tempo deviations which are applied to a drum beat or any part of a MIDI sequence (comparable to a one level Clynes pulse set, see 3.1.3). Such groove systems have been very useful commercially in systems like Propellorhead Reason, where it is possible to generate Groove Templates from a drum track track and apply it to any other MIDI track [Carlson et al. 2003]. However just as some research has suggested limitations in the application of Clynes’ composer pulses, so Wright and Berdahl’s [2006] research shows the limits of groove templates. They did an analysis of multi-voiced Brazillian drumming recordings and found that groove templates could only account for 30% of expressive timing.

They investigated other methods to capture the expressive timing using a system that learned from audio files. The audio features examined were a note’s timbre, metric position and rhythmic context (i.e. timbres and relative temporal position of notes within 1 beat of the input notes – this gives a low expressive perception). The system learned to map these audio features onto the timing deviations of each non-pitched note; it is not designed to generate creative performances, but to simulate. The mapping model was based on non-linear regression between audio features and timing deviation (versus a quantized version of the beat). They tried three different methods of learning the

mapping model: Kernel Ridge Regression, Gaussian Process Regression and kNN methods. This learning approach was found to track the expressive timing of the drums much better than the groove templates, clearly demonstrated in their graphs showing the error over the drum patterns. All 3 learning methods were found to give approximately equal accuracy, though it was believed that Gaussian Process Regression had the greatest room for improvement. Examples are provided online. Note that the system is not only limited to Brazilian drumming, Wright and Berdahl also tested it on some Reggae rhythms with similar success.

3.4.2 Non-linear Piano System. The most recent application of a non-linear regression methods to expressive performance is the system by [Dorard et al 2007]. Their main aim is simulatory, to imitate the style of a particular performer and allow new pieces to be automatically performed using the learned characteristics of the performer. A performer is defined based on the “Worm” representation of expressive performance [Dixon et al 2002]. The worm is a useful visualisation tool for the dynamics and tempo aspects of expressive performance. It uses a 2D representation with tempo on the x-axis and loudness on the y-axis. Then, as the piece plays, at fixed periods in the score (e.g. once per bar) an average is calculated for each period and a filled circle plotted on the graph at the average. Past circles remain on the graph, but their colour fades and size decreases as time passes – thus creating the illusion of a wriggling worm whose tail fades off into the distance in time. If the computer played an expressionless MIDI file then its worm would stand still, not wriggling at all.

The basis of Dorard’s approach is to assume that the score and the human performances of the score are two views of the musical semantic content, thus enabling a correlation to be drawn between the worm and the score. The system focuses on homophonic piano music – a continuous upper melody part and an accompaniment – and divides the score into a series of chord and melody pairs. Kernel Canonical Correlation Analysis (KCCA) is then used, a method which looks for a common semantic representation between two views. Its expressive perception is based on the note group level, since KCCA is looking to find correlations between short groups of notes and the performance worm position. An addition needed to be made to the learning algorithm to prevent extreme expressive changes in tempo and dynamics. This issue is a recurring problem in a number of CSEMPs (see the Artificial Neural Network Models in Section 3.5, Sibelius in Section 3.1.9, and also Section 3.7.3).

Testing was performed on Chopin’s Etude 3 Opus 10 – the system was trained on the worm of the first 8 bars, and then tried to complete the worm for bars 9 to 12. The correlation between the original human performance worm for 9 to 12 and the reconstructed worm was measured to be 0.95 (whereas the correlation with a random worm was 0.51). However the resulting performances were reported - through presumably informal listening tests - to not be very pleasant to listen to.

3.5 Artificial Neural Networks

Although many Artificial Neural Networks (ANNs) are just another form of non-linear regression, the connectionist visualisation has supported the development of new “architectures”. At least two CSEMPs have been produced using dynamic time-based ANN architectures.

3.5.1 Artificial Neural Network Piano System. The earliest ANN approach is the Artificial Neural Network Piano System [Bresin et al 1990][Bresin 1998]. It has two incarnations. The first did not learn from human performers: a set of 7 monophonic Director Musice rules were selected, and two (loudness and timing) feedforward ANNs learned these rules through being trained on them. By learning a fixed model of Director Musices, the ANN loses the performance creativity of the k-values. When monophonic listening tests were done with 20 subjects, using Mozart’s Piano Sonatas K331 and K281, the Director Musices performance was rated above the non-expressive computer performance, but the Neural Network performance rated highest of all. One explanation for the dominance of the ANN over the original DM rules was that the ANN generalised in a more pleasant way than the rules. The ANN system by Bresin was a simulation CSEMP which also used a separate loudness and timing feedback ANN. The ANNs were trained using actual pianist performances from MIDI, rather than on DM rules; but some of the independently learned rules turned out to be similar to some DM rules. Informal listening tests judged the ANNs as musically acceptable. The network looked at a context of four notes (loudness) and five notes (timing), so had note to phrase-level expressive perception, though it required the notes to be manually grouped into phrases before being input.

3.5.2 Emotional Flute. The Camurri et al. [2000] Emotional Flute system uses explicit Music Features and Artificial Neural Networks, thus allowing greater generalisation than the related CARO system (Section 3.3.2). The music features are similar to those used in Director Musices. This CSEMP is strongly related to Bresin’s second ANN, extending it into the non-piano realm and adding mood space modelling. Expressive actions include inter-onset interval, loudness, and vibrato. Pieces need to be segmented into phrases before being input - this segmentation is performed automatically by another ANN. There are separate nets for timing and for loudness – net designs are similar to Bresin’s, and similar levels of expressive perception. There is also a third net for the duration of crescendo and decrescendo at the single note level. However the nets could not be successfully trained on vibrato, so a pair of rules were generated to handle it. A flautist performed the first part of Telemann’s Fantasia no.2 in nine different moods: cold, natural, gentle, bright, witty, serious, restless, passionate and dark. Like CARO a 2-D mood space was generated and mapped on to the performances by the ANNs, and this mood space can be utilised to give greater Performance Creativity.

To generate new performances the network drives a physical model of a flute. Listening tests gave an accuracy of approximately 77% when subjects attempted to assign emotions to synthetic performances. To put this in perspective, even when listening to the original human performances, human recognition levels were not always higher than 77%; the description of emotional moods in music is a fairly subjective process.

3.6 Statistical Graphical Models

3.6.1 Music Plus One. The Music Plus One system [Raphael 2001a; Raphael 2001b; Raphael 2003] is designed to deal with non-piano and polyphonic performances. It has the ability to adjust performances of polyphonic sound files (e.g. orchestral works) to fit as accompaniment for solo performers. This CSEMP contains two modules: Listen and Play modules. Listen uses a Hidden Markov Model (HMM) to track live audio and find the soloist's place in the score in real time. Play uses a Bayesian Belief Network (BBN) which, at any point in a soloist performance and based on the performance so far, tries to predict the timing of the next note the soloist will play. Music Plus One's BBN is trained by listening to the soloist. As well as timing, the system learns the loudness for each phrase of notes. However loudness learning is deterministic - it performs the same for each accompaniment of the piece once trained, not changing based on the soloist changing their own loudness. Expressive perception is at the note level for timing and phrase level for loudness.

The BBN assumes a smooth changing in tempo, so for any large changes in tempo (e.g. a new section of a piece) need to be manually marked up. For playing MIDI files for accompaniment, the score needs to be divided up manually into phrases for dynamics; for using audio files for accompaniment such a division is not needed. When the system plays back the accompaniment it can play it back in multiple expressive interpretations dependent on how the soloist plays. So it has learned a flexible (almost tempo-independent) concept of the soloist's expressive intentions for the piece.

There is no test reported for this system – the authors state their impression that the level of musicality obtained by the system is surprisingly good, and ask readers to evaluate the performance themselves by going to the website and listening. Music Plus One is actually being used by composers, and for teaching music students. It came first at RenCon 2003 in the compulsory section with a performance of Chopin's Prelude No. 15 "Raindrop", beating Ha-Hi-Hun, Kagurame, and Widmer's system. To train Music Plus One for this, several performances were recorded played by a human, using a MIDI keyboard. These were used to train the BBN. The model was extended to include velocities for each note, as well as times, with the assumption that the velocity varies smoothly (like a random walk) except at hand-identified phrase boundaries. Then a mean performance was generated from the trained model.

As far as performance creativity goes, the focus on this system is not so much to generate expressive performances, as to learn the soloist's expressive behaviour and react accordingly in real time. However the system has an "implicit" method of creating new performances of the accompaniment – the soloist can change their performance during playback. There is another creative application of this system: multiple pieces have been composed for use specifically with the Music Plus One system - pieces which could not be properly performed without the system. One example contains multiple sections where a musician plays 7 notes while the other plays 11. Humans would find it difficult to do this accurately, whereas a soloist and the system can work together properly on this complicated set of polyrhythms.

3.6.2 ESP Piano System. Grindlay's [2005] ESP Piano system is a polyphonic CSEMP designed – like Dorard's system (section 3.4.2) - to simulate expressive playing of pieces of piano music which consist of a largely monophonic melody, with a set of accompanying chords. A Hidden Markov Model learns expressive performance using music features such as whether the note is the first or last of the piece, the position of the note in its phrase, and the notes duration relative to its start and the next note's start (called its "articulation" here). The expressive perception is up to the phrase level. Phrase division is done manually, though automated methods are discussed. The accompaniment is analysed for a separate set of music features some of which are like the melody music features, and some of which are unique to chords – for example the level of consonance/dissonance of the code (based on a method called Euler's Solence). Music features are then mapped on to a number of expressive actions such as (for melody) the duration deviation, and the velocity of the note compared to the average velocity. For the accompaniment similar actions are used as well as some chord-only actions, like the relative onset of chord notes (similar to the Kagurame Chord Time Lag, in section 3.2.2). These chordal values are based on the average of the values for the individual notes in the chord.

Despite the focus of this system on homophony, tests were only reported for monophonic melodies, training the HMM on 10 graduate performances of Schumann's *Träumerei*. 10 out of 14 listeners ranked the expressive ESP output over the inexpressive version. 10 out of 14 ranked the ESP output above that of an undergraduate performance. 4 out of 7 preferred the ESP output to a graduate student performance.

3.7 Evolutionary Computation

A number of more recent CSEMPs have used evolutionary computation methods. In general (but not always) such systems have interesting opportunities for performance creativity. They often have a parameterization that is simple to change – for example a fitness function. They also have an emergent result which can sometimes produce unexpected but coherent results.

3.7.1 Genetic Programming Jazz Sax. Some of the first researchers to use EC in computer systems for expressive performance were Ramirez and Hazan. They did not start out using EC, beginning with a Regression Tree system for Jazz Saxophone first [Ramirez and Hazan 2005]. We will describe this before moving on to the Genetic Programming (GP) approach, as it is the basis of their later GP work. A performance Decision Tree was first built using C4.5 [Quinlan 1993]. This was built for musicological purposes – to see what kinds of rules were generated – not to generate any performances. The Decision tree system had a 3-note-level expressive perception, and music features used to characterise a note included metrical position and some Narmour IR analysis. These features were mapped on to a number of performance actions from the training performances, such as lengthen/shorten note, play note early/late and play note louder/softer. Monophonic audio was used to build this decision tree using the authors' own spectral analysis techniques and five Jazz standards at 11 different tempos. The actual performing system was built as a Regression rather than Decision Tree, thus allowing continuous expressive actions. The continuous performance features simulated were duration, onset, and energy variation (i.e. loudness). The learning algorithm used to build the tree was

M5Rules[Witten and Frank 2000] and performances could be generated via MIDI and via audio thanks to the synthesis algorithms. In tests, the resulting correlations with the original performances were 0.72, 0.44 and 0.67 for duration, onset and loudness respectively. Other modelling methods were tried (linear regression and 4 different forms of Support Vector Machines) but didn't fare as well correlation-wise.

Ramirez and Hazan's next system [Hazan and Ramirez 2006] was also based on Regression Trees, but these trees were generated using Genetic Programming (GP), which is ideal for building a population of "if-then" Regression Trees. GP was used to search for Regression Trees which best emulated a set of human audio performance actions. The Regression Tree models were basically the same as in their previous paper, but in this case a whole series of trees was generated, they were tested for fitness and then the fittest were used to produce the next generation of trees/programs (with some random mutations added). Fitness was judged based on a distance calculated from a human performance. Creativity and expressive perception are enhanced because, in addition to modelling timing and dynamics, the trees modelled the expressive combining of multiple score notes into a single performance note (consolidation), and the expressive insertion of one or several short notes to anticipate another performance note (ornamentation). These elements are fairly common in Jazz Saxophone. It was possible to examine these deviations because the fitness function was implemented using an Edit Distance [Levenshtein 1966] to measure score edits.

This evolution was continued until average fitness across the population of trees ceased to increase. The use of GP techniques was deliberately applied to give a range of options for the final performance since, as the authors say - "performance is an inexact phenomenon". Also because of the mutation element in Genetic Programming, there is the possibility of unusual performances being generated. So this CSEMP has quite a good potential for performance creativity. No evaluation was reported of the resulting trees' performances - but average fitness stopped increasing after 20 generations.

3.7.2 Sequential Covering Algorithm GAs. The Sequential Covering Algorithm Genetic Algorithm (GA) [Ramirez and Hazan 2007] uses the Sequential Covering to learn performance. Each covering rule is learned using a GA; and a series of such rules are built up covering the whole problem space. In this paper the authors return to their first (non-EC) paper's level of expressive perception - looking at note level deviations without ornamentation or consolidation. However they make significant improvements over their original non-EC paper. The correlation coefficients for onset, duration, and energy/loudness in the original system were 0.72, 0.44 and 0.67 - but in this new system they were 0.75, 0.84 and 0.86 - significantly higher. And this system also has the advantage of slightly greater creativity due to its GA approach.

3.7.3 Generative Performance GAs. A parallel thread of EC research is the Zhang and Miranda's [2006a; 2006b] monophonic Generative Performance GAs which evolve Pulse Sets (see Section 3.1.3). Rather than comparing the generated performances to actual performances, the fitness function here expresses constraints inspired by the generative performance work of Eric Clarke [1998]. When a score is presented to the GA system for performance, the system constructs a theoretical timing and

dynamics deviation curve for the melody (one advantage of this CSEMP being that this music analysis is automatic). However this curve is not used directly to generate the actual performance, but to influence the evolution. This, together with the GA approach, increases the performance creativity of the system. The timing deviation curve comes from an algorithm based on Cambouropoulos' [2001] Local Boundary Detection Model (LBDM) – the inputs to this model are score note timing, pitch and harmonic intervals. The resulting curve is higher for greater boundary strengths. The approximate dynamics curve is calculated from a number of components – the harmonic distance between two notes (based on a measure by Krumhans [1991]), the metrical strength of a note (based on the Melisma software [Temperley and Sleator 1999]), and the pitch height. These values are multiplied for each note to generate a dynamics curve. The expressive perception of this system is the same as the expressive perception of the methodologies used to generate the theoretical curves.

A fitness function is constructed referring to the score perception curves. It has 3 main elements – fitness is awarded if: (1) the pulse set dynamics and timing deviations follow the same direction as the generated dynamics and timing curves; (2) timing deviations are increased at boundaries; (3) timing deviations are not too extreme. Part (1) does not mean that the pulse sets are the same as the dynamics and timing curves, but – all else being equal - that if the dynamic curve moves up between two notes, and the pulse set moves up between those two notes, then that pulse set will get a higher fitness than one that moves down there. Regarding point (3) – this is reminiscent of the restriction of expression used in the ANN models and the non-Linear Piano model described earlier. It is designed to preventing the deviations from becoming too extreme.

There has been no formal testing of this GA work, though the authors demonstrate - using an example of part of Schumann's *Träumerei* and a graphical plot - that the evolved pulse sets are consistent in at least one example with the theoretical timing and dynamics deviations curves. They claim that “when listening to pieces performed with the evolved pulse sets, we can perceive the expressive dynamics of the piece.” However, more evaluation would be helpful because repeating pulse sets as the expressive action format have been shown to not be universally applicable, and post-Clynes CSEMPs have shown more success using non-cyclic expression.

3.7.4 Multi-Agent System with Imitation. Zhang and Miranda [2007] developed the above into a Multi-Agent System (MAS) with Imitation – influenced by Miranda's [2003] evolution of music MAS study, and inspired by the hypothesis that expressive music performance strategies emerge through interaction and evolution in the performers' society. In this model each agent listens to other agents' monophonic performances, evaluates them, and learns from those whose performances are better than their own. Every agent's evaluation equation is the same as the fitness function used in the previous GA paper, and performance deviations are modelled as a hierarchical pulse set. So it has the same expressive perception. When an agent hears a performance it evaluates as being better than its own, it attempts to imitate it but ends up giving a slightly mutated performance. This imitating agent then generates a pulse set internally that matches its mutated performance. So pulse sets are not exchanged, performances are (just as humans imitate behaviour, not neural and underlying biological processes). The performances of the system, which were generated through 100 generations of imitation, were

once again not evaluated by the researchers. Though it was found that the imitation approach generated performances more slowly than the previous GA approach.

This CSEMP has significant performance creativity, one reason being that the pulse sets generated may have no similarity to the hierarchical constraints of human pulse sets. They are generated mathematically and abstractly from agent imitation performances. So entirely novel pulse set types could be produced by agents that a human would never generate. Another element that contributes to creativity is that although a global evaluation function approach was used, a diversity of performances was found to be produced in the population of agents.

This Zhang/Miranda MAS system has much potential for future work. For example, an approach could be investigated that allowed each agent to have its own evaluation function; with some agents actually basing their evaluation on comparing with human performances. Furthermore, some agents could be equipped with entirely novel non-human-fitness evaluation functions. Thus by altering the proportion of agents with different types of evaluation functions, different types of performances could be generated in the population. Another possibility would be to allow more direct user intervention in shaping performances by letting the user become an agent in the system themselves.

3.7.5 Ossia. Like the Music Emotionality System (Section 3.1.10), Dahlstedt's [2007] Ossia is a CSEMP which incorporates both compositional and performance aspects. However, whereas MES was designed to operate on a composition, Ossia is able to generate entirely new and expressively performed compositions. Although we have grouped it as an EC learning system, technically Ossia is not a learning system. It is not using EC to learn how to perform like a human, but to generate novel compositions and performances. However we include it in this section because its issues relate more closely to EC and learning systems than to any of the non-learning systems (the same reason applies for the system described in the next subsection 3.7.6). Ossia generates music through a novel representational structure that encompasses both composition and performance – Recursive Trees (generated by GAs). These are “upside down trees” containing both performance and composition information. The bottom leaves of the tree going from left to right represent actual notes (each with their own pitch, duration and loudness value) in the order they are played. The branches above the notes represent transformations on those notes. To generate music the tree is flattened – the “leaves” higher up act upon the leaves lower down when being flattened to produce a performance/composition. So going from left to right in the tree represents music in time. The trees are generated recursively – this means that the lower branches of the tree are transformed copies of higher parts of the tree. Here we have an element we argue is the key to combined performance and composition systems - a common representation – in this case transformations.

This issue of music representation is not something we have addressed explicitly in this review, being in itself an issue worthy of its own review, for examples see [Dannenberg 1993][Anders 2007]. However we will take a moment now to briefly discuss it. The representation chosen for a musical system has a significant impact on the functionality – Ossia's representation is what leads to its combined composition and performance generation abilities. The most common music representation mentioned in this review has been MIDI, which is not able to encode musical structure directly. As a

result some MIDI-based CSEMPs have to supply multiple files to the CSEMP, a MIDI file together with files describing musical structure. More flexible representations than MIDI include MusicXML, ENP-score-notation [Laurson and Kuuskankare 2003], WEDEMUSIC XML [Hirata et al 2003], MusicXML4R [Hashida et al 2006], and the proprietary representations used by commercial software such as Sibelius, Finale, Notion, and Zenph High-Resolution MIDI [LaVerne 2006] (which was recently used on a released CD of automated Disklavier re-performances of Glenn Gould).

Many of the performance systems we have described so far transform an expressionless MIDI or audio file into an expressive version. Composition is often done in a similar way – motifs are transformed into new motifs, and themes are transformed into new expositions. Ossia uses a novel transformation-based music representation. In Ossia, transformations of note, loudness and duration are possible – the inclusion of note transformations here emphasising the composition aspect of the Ossia. The embedding of these transformations into recursive trees leads to the generation of gradual crescendos, decrescendos and duration curves – which sound like performance strategies to a listener. Because of this Ossia has a good level of performance creativity. The trees also create a structure of themes and expositions. Ossia uses a GA to generate a population of trees, and judges for fitness using such rules as number of notes per second, repetivity, amount of silence, pitch variation, and level of recursion. These fitness rules were developed heuristically by Dahlstedt through analysis-by-synthesis methods.

Ossia's expressive perception is equal to its compositional perception. Dahlstedt observes "The general concept of recapitulation is not possible, as in the common ABA form. This does not matter so much in short compositions, but may be limiting." So Ossia's perception would seem to be within the A's and B's, giving it a note to section level expressive perception. In terms of testing – the system has not been formally evaluated; though it was exhibited as an installation at Gaudeamus Music Week in Amsterdam. Examples are also available on the website, including a composed suite. The author claims that the sound examples "show that the Ossia system has the potential to generate and perform piano pieces that could be taken for human contemporary compositions." Having listened to examples on the website, we were impressed by their natural quality. The question of how to test a combined performance and composition, when that system is not designed to simulate but to create, is a sophisticated problem which we will not try to address here. Certainly listening tests are a possibility but these may be biased by the preferences of the listener (e.g. preferring pre-1940s classical music, or pop music). Another approach is musicological analysis but the problem then becomes that musicological tools are not available for all genres and all periods – for example musicology is more developed for pre-1940 than post-1940 art music.

An example score from Ossia is described which contains detailed dynamics and articulations, and subtle tempo fluctuations and rubato. This subtlety raises another issue – scores generated by Ossia in common music notation had to be simplified to be simply readable by humans. The specification of exact microfeatures in a score can lead to it being unplayable except by computer or the most skilled concert performer. This has a parallel in a compositional movement which emerged in the 1970s "The New Complexity", involving composers such as Brian Ferneyhough and Richard Barret [Toop 1988] In "The New Complexity" elements of the score are often specified down to the microfeature level, and

some scores are described as almost unplayable. Compositions such as this, whether by human or computer, bring into question the whole composition/performance dichotomy. (These issues also recall the end of Ian Pace's quote in the first section of this review.) However, technical skill limitations and common music notation scores are not necessary for performance if the piece is being written on and performed by a computer. Microfeatures can be generated as part of the computer (or computer-aided) composition process if desired. In systems such as Ossia and MES (Section 3.1.10), as in *The New Complexity*, the composition/performance dichotomy starts to break down - the dichotomy is really between macrofeatures and microfeatures of the music.

3.7.6 pMIMACS. Before discussing the final system in this survey, we will highlight another motivation for bringing composition and performance closer in CSEMPs. A significant amount of CSEMP effort is in analysing the musical structure of the score/audio. However, many computer composition systems generate a piece based on some structure which can often be made explicitly available. So in computer music it is often inefficient to have separate composition and expressive performance systems – where a score is generated and CSEMP sees the score as a black box and performs a structure analysis. Greater efficiency and accuracy would require a protocol allowing the computer composition system to communicate structure information directly to the CSEMP, or – like Ossia - simply combine the systems using for example a common representation (where microtiming and microdynamics are seen as an actual part of the composition process). A system which was designed to utilise this combination of performance and composition is pMIMACS, developed by the authors of this survey. It is based on a previous system MIMACS (Mimetics-Inspired Multi-agent Composition System), which was developed to solve a specific compositional problem (generating a multi-speaker spatial composition).

pMIMACS combines composition and expressive performance – the aim being to generate contemporary compositions on a computer which when played back on a computer do not sound too machine like. In Zhang/Miranda's system (Section 3.7.4) the agents imitate each others' expressive performances, whereas in pMIMACS agents can be performing entirely different pieces of music. The agent cycle is a process of singing and assimilation. Initially all agents are given their own tune – these may be random or chosen by the composer. An agent (A) is chosen to start. A performs its tune, based on its “performance skill” (explained below). All other agents listen to A and the agent with the most similar tune, say agent B, adds its interpretation of A's tune to the start or end of B's current tune. There may be pitch and timing errors due to its “mishearing”. Then the cycle begins again, but with B performing its extended tune in the place of A.

An agent's initial performing skills are defined by the average pitch and standard deviation of their initial tune - this could be interpreted as the tune they are familiar with performing, or as the range they are comfortable performing in. The further away a note's pitch is from the agent's average learned pitch, the slower the tempo at which the agent will perform. Also, further away pitches will be played more quietly. An agent updates its skill/range as it plays. Every time it plays a note, that note changes the agent's average and standard deviation pitch value. So when an agent adds an interpretation of another agent's tune to its own, then as the agent performs the new extended tune its average and

standard deviation (skill/range) will update accordingly – shifting and perhaps widening - changed by the new notes as it plays them. In pMIMACS an agent also has a form of performance context, called an Excitability State. An “excited” agent will play its tune with twice the tempo of an “unexcited” agent, making macro-errors in pitch and rhythm as a result.

The listening agent has no way of knowing whether the pitches, timings and amplitude that it is hearing are due to the performance skills of the performing agent, or part of the “original” composition. So the listening agent attempts to memorize the tune as it hears it, including any performance errors or changes. As the agents perform to each other, they store internally and exponentially growing piece of transforming music. The significant and often smooth deviations in tempo generated by the performance interaction will create a far less robotic-sounding performance than rhythms generated by a quantized palette would do. On the downside, the large-scale rhythmic texture has the potential to become repetitive because of the simplicity of the agents’ statistical model of performance skill. Furthermore the system can generate rests that are so long that the composition effectively comes to a halt for the listener. But overall the MAS is expressing its experience of what it is like to perform the tune, by changing the tempo and dynamics of the tune; and at the same time this contributes to the composition of the music. No formal listening tests have been completed yet, but examples of an agent’s tune memory after a number of cycles can be listened to at: <http://cmr.soc.plymouth.ac.uk/pmimacs>.

There is also a more subtle form of expression going on relating to the hierarchical structure of the music. The hierarchy develops as agents pass around an ever growing string of phrases. Suppose an agent performs a phrase P and passes it on. Later on it may receive back a “super-phrase” containing two other phrases Q and R – in the form QPR. In this case A will perform P faster than Q and R (since it knows P). Now suppose in future A is passed back a super-super-phrase of, say, SQPRTQPRTS, then potentially it will play P fastest, QPR second fastest (since it has played QPR before) and the S and T phrases slowest. So the tempo and amplitude at which an agent performs the parts SQPRTQPRTS is affected by how that phrase was built up hierarchically in the composition/imitation process. Thus there is an influence on the performance from the hierarchical structure of the music. This effect is only approximate because of the first order statistical model of performance skill.

Despite the lack of formal listening tests, we report pMIMACS here as a system designed from the ground up to combine expressive performance and composition.

4. SUMMARY

We began our review of automated and semi-automated computer systems for expressive performance, by developing four terms of reference which were inspired from research into computer composition systems, and we have brought it to a close with a pair of systems that question the division between expressive performance and composition. So there are clearly opportunities for to these two areas to learn from each other in the context of computer music. Before we summarise further, another viewing of Table I at the start of the review may be helpful to the reader.

Expressive Performance is a complex behaviour with many causative conditions – so it is no surprise that in this review that more than half the systems produced have been learning CSEMPS,

usually learning to map music features on to expressive actions. Expressive performance actions most commonly included timing and dynamics adjustments, with some articulation, and the most common non-custom method for analysis of music features was GTTM, followed by IR. Dues to its simplicity in modelling performance, the most common instrument simulated was piano – but interestingly this was followed closely by saxophone – possibly because of the popularity of the instrument in the Jazz genre. Despite, and probably because, of its simplicity - MIDI is still the most popular representation.

To help structure the review, four primary terms of reference were selected: Testing Status, Expressive Perception, Non-monophonic Ability, and Performance Creativity. Having applied these, it can be seen that only a subset of the systems have had any formal testing, and for some of them designing formal tests is a challenge in itself. This is not that unexpected – since testing a creative computer system is an unsolved problem. Also about half of the systems were only been tested on monophonic tunes. Polyphony and Homophony introduce problems both in terms of synchronisation and in terms of music feature analysis. Further to music feature analysis, most of the CSEMPs had an expressive perception up to one bar/phrase, and over half did not look at the musical hierarchy. However avoiding musical hierarchy analysis can have the advantage of increasing automation. We have also seen that most CSEMPs are designed for simulation of human expressive performances, general or specific – a valuable research goal, and one which has possibly been influenced by the philosophy of human simulation in machine intelligence research.

The results for the primary terms of reference are summarized in Table V. The numerical measures in columns 1,2 and 4 are an attempt to quantify observations, scaled from 1 to 10. The more sophisticated the expressive perception of music features (levels of the hierarchy, harmony, etc), the higher the number in column 1. The more extensive the testing (including informal listening, RenCon submission, formal listening, correlation and/or successful utilisation in the field) the higher the number in column 2. The greater we perceived the potential of a system to enable the creative generation of novel performances, the higher the number in column 4. Obviously such measures contain some degree of subjectivity but should be a useful indicator for anyone wanting an overview of the field, based on the 4 elements discussed at the start of this paper. Figure 4 shows a 3D plot of summary table V.

5. CONCLUSIONS

There have been significant achievements in the field of simulating human musical performance in the last 25 years, and there many opportunities ahead for future improvements in simulation. In fact one aim of the RenCon competitions is for a computer to win the Chopin competition by 2050. Such an aim begs some philosophical and historical questions, but nonetheless captures the level of progress being made in performance simulation. The areas of expressive perception and non-monophonic performance appear to be moving forwards. However the issue of testing and evaluation still requires more work and would be a fruitful area for future CSEMP research.

CSEMP	Expressive Perception	Testing Status	Non-monophonic	Performance Creativity
Director Musices	6	10	Y	8
Hierarchical Parabola Model	9	7		3
Composer Pulse	6	9	Y	5
Bach Fugue	6	4	Y	3
Rubato	6	4	Y	8
Trumpet Synthesis	4	1		3
MIS	6	6	Y	3
ANN Piano	4	6		3
Music Plus One	4	7	Y	4
SaxEx	6	1		8
CARO	3	6		7
Emotional Flute	4	6		6
Kagurame	9	8	Y	3
Ha-Hi-Hun	6	6	Y	8
PLCG	4	6		3
Phrase-decomposition/ PLCG	10	6		3
DISTALL	10	9	Y	3
Pop-E	6	10	Y	8
ESP Piano	6	6	Y (untested)	3
Non-linear Piano	4	6	Y	3
Drumming	3	6		3
Genetic Programming	7	6		6
Sequential Covering GAs	6	6		6
Generative Performance GAs	9	4		8
MAS with Imitation	9	1		9
Ossia	6	4	Y	10
Music Emotionality	6	6	Y	10
pMIMACS	9	1		10

Table V. Summary of the 4 primary terms of reference

Another fruitful area for research is around the issue of the automation of the music analysis. Of the eight CSEMPs with the highest expressive perception, almost half of them require some manual input to perform the music analysis. Also manual marking of the score into phrases is a common requirement. There has been some research into automating musical analysis such as GTTM [Masatoshi et al 2005]. The usage of such techniques, and the investigation of further automation analysis methods specific to expressive performance, would be a useful contribution to the field.

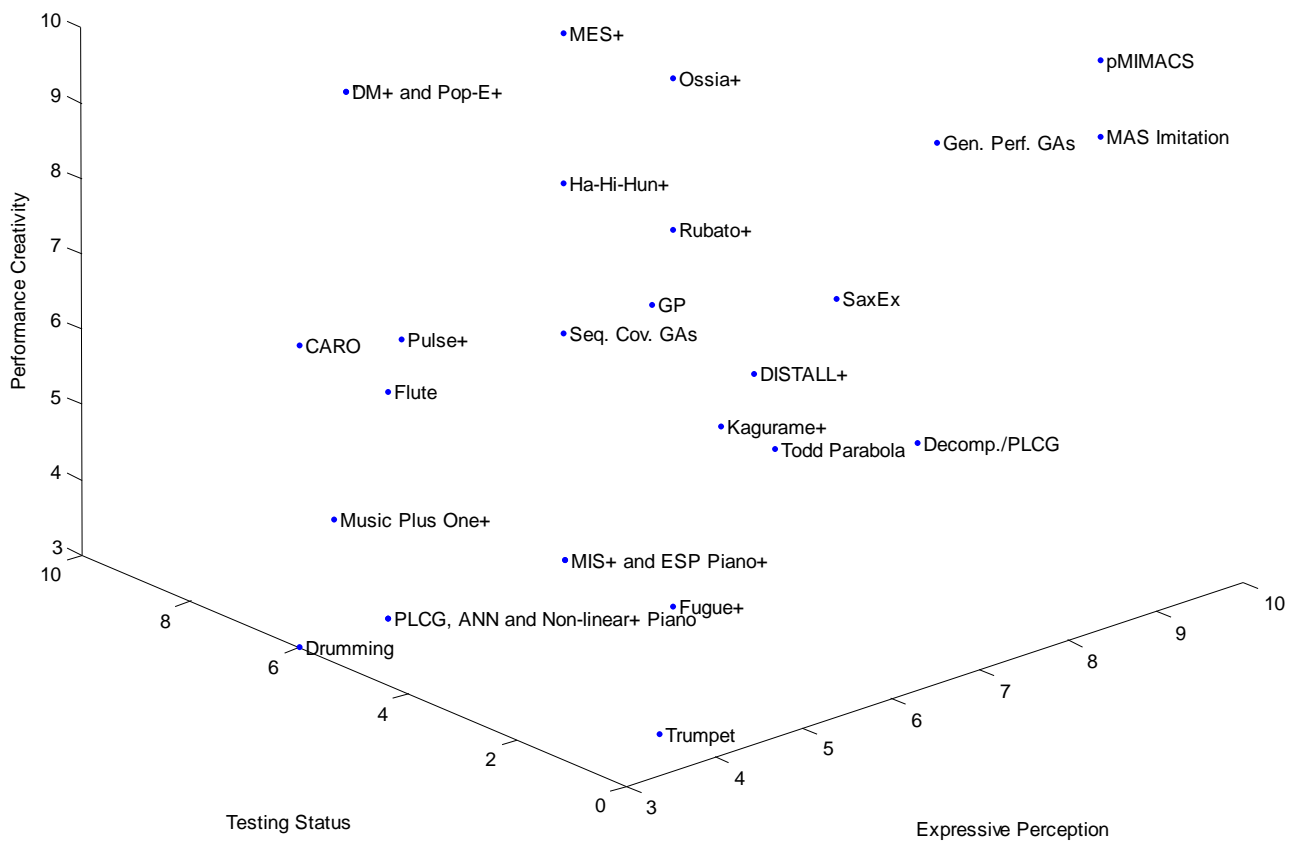


Figure 4. Evaluation of Reviewed systems relative to terms of reference ('+' means Non-monophonic)

The field could also benefit from a wider understanding of the convergence of performance and composition elements in computer music. For reasons of efficiency, controllability and creativity which have already been detailed, research into combined performance/composition systems would be a fertile area for further computer music research. Here computer composition and computer performance research can cross-fertilise: performance algorithms for expressing structure and emotion/mood can help composition, as well as composition providing more creative and controlled computer performance. The question is also open as to what forms of non-human expression can be developed and provide whole new vistas of the meaning of the phrase “expressive performance”, perhaps even for human players.

One final observation regards the lack of neurological and physical models of performance simulation. The issue was not included as a part of our terms of reference, since it is hard to objectively quantify. But we would like to address this in closing. Neurological and physical modelling of performance should go beyond ANNs and instrument physical modelling. The human/instrument performance process is a complex dynamical system about which there has been some deeper psychological and physical study. However attempts to use these hypotheses to develop computer performance systems have been rare. More is being learned about the neural correlates of music and

emotion [Koelsch and Siebel 2005; Britton et al. 2006][Durrant et al 2007], and Eric Clarke [1993] has written on the importance of physical embeddedness of human performance. But although researchers such as Parncutt [1997] (in his virtual pianist approach) and Widmer [Widmer and Goebel 2004] have highlighted the opportunity for deeper models, there has been little published progress in this area of the CSEMP field.

So to conclude - the overarching focus so-far means that there are opportunities for some better tested, more creative, and neurological and physical models of human performance. These will be systems which not only help to win the Chopin contest, but also to utilise the innate efficiencies and power of computer music techniques. Music psychologists (and musicologists) will be provided with richer models; composers will be able to work more creatively in the micro-specification domain, and more easily and accurately generate expressive performances of their work. And the music industry will be able to expand the power of humanisation tools creating new efficiencies in recording and performance.

REFERENCES

- ARCOS, J.L., DE MANTARAS, R.L., AND SERRA, X. 1997. SaxEx: A Case-based Reasoning System for Generating Expressive Musical Performances. In Proceedings of 1997 International Computer Music Conference, Thessalonikia, Greece, September 1997, COOK P.R., Eds. ICMA, San Francisco, CA, 329-336.
- ARCOS, J.L., LOPEZ DE MANTARAS, R., AND SERRA, X. 1998. Saxex: A Case-Based Reasoning System for Generating Expressive Musical Performance. *Journal of New Music Research* 27, 194-210.
- ARCOS, J.L. AND DE MANTARAS, R.L. 2001. The SaxEx System for Expressive Music Synthesis: A Progress Report. In Proceedings of the Workshop on Current Research Directions in Computer Music, Barcelona, Spain, November 2001, LOMELI C. and R. LOUREIRO, Eds. 17-22.
- ARCOS, J.L. AND LOPEZ DE MANTARAS, R. 2001. An Interactive Case-Based Reasoning Approach for Generating Expressive Music. *Journal of Applied Intelligence* 14, 115-129.
- BRESIN, R. 1998. Artificial Neural Networks Based Models For Automatic Performance of Musical Scores. *Journal of New Music Research* 27, 239-270.
- BRESIN, R. AND FRIBERG, A. 2000. Emotional Coloring of Computer-Controlled Music Performances. *Computer Music Journal* 24, 44-63.
- BRITTON, J.C., PHAN, K.L., TAYLOR, S.F., WELSCH, R.C., BERRIDGE, K.C., AND LIBERZON, I. 2006. Neural correlates of social and nonsocial emotions: An fMRI study. *NeuroImage* 31, 397-409.
- BUXTON, W.A.S. 1977. A Composers Introduction to Computer Music. *Interface* 6, 57-72.
- CAMBOUROPOULOS, E. 2001. The Local Boundary Detection Model (LBDM) and its Application in the Study of Expressive Timing. In Proceedings of the 2001 International Computer Music Conference, Havana, Cuba, September 2001, SCHLOSS R. and R. DANNENBERG, Eds. International Computer Music Association, San Fransisco, CA,
- CAMURRI, A., DILLON, R., AND SARON, A. 2000. An Experiment on Analysis and Synthesis of Musical Expressivity. In Proceedings of 13th Colloquium on Musical Informatics, L'Aquila, Italy, September 2000
- CANAZZA, S., DRIOLI, C., DE POLI, G., RODÀ, A., AND VIDOLIN, A. 2000. Audio Morphing Different Expressive Intentions for Multimedia Systems. *IEEE Multimedia* 7, 79-83.
- CANAZZA, S., DE POLI, G., DRIOLI, C., RODÀ, A., AND VIDOLIN, A. 2001. Expressive Morphing for Interactive Performance of Musical Scores. In Proceedings of First International Conference on WEB Delivering of Music, Florence, Italy, November 2001, IEEE, 116-122.

- CANAZZA, S., DE POLI, G., RODÀ, A., AND VIDOLIN, A. 2003. An abstract control space for communication of sensory expressive intentions in music performance. *Journal of New Music Research* 32, 281-294.
- CANAZZA, S., DE POLI, G., DRIOLI, C., RODÀ, A., AND VIDOLIN, A. 2004. Modeling and control of expressiveness in music performance. *The Proceedings of the IEEE* 92, 686-701.
- CARLSON, L., NORDMARK, A., AND WIKLANDER, R. 2003. Reason Version 2.5 - Getting Started. Propellorhead Software,
- CLARKE, E.F. 1993. Generativity, mimesis and the human body in music performance. *Contemporary Music Review* 9, 207-219.
- CLARKE, E.F. 1998 Generative Principles in Music Performance. In *Generative Processes in Music: The Psychology of Performance, Improvisation, and Composition*, J.A. SLOBODA, Eds. Clarendon Press, Oxford, 1-26.
- CLYNES, M. 1986. Generative principles of musical thought: Integration of microstructure with structure. *Communication and Cognition* 3, 185-223.
- COPE, D. 2005. *Computer Models of Musical Creativity*. Cambridge, MA, USA.
- DAHLSTEDT, P. 2007. Autonomous Evolution of Complete Piano Pieces and Performances. In *Proceedings of ECAL 2007 Workshop on Music and Artificial Life (MusicAL 2007)*, Lisbon, Portugal, September 2007.
- DANNENBERG, R.B. 1993. A Brief Survey of Music Representation Issues, Techniques, and Systems. *Computer Music Journal* 17, 20-30.
- DANNENBERG, R.B., PELLERIN, H., AND DERENYI, I. 1998. A Study of Trumpet Envelopes. In *Proceedings of the 1998 International Computer Music Conference*, Ann Arbor, Michigan, October 1998, International Computer Music Association, San Francisco, 57-61.
- DANNENBERG, R.B. AND DERENYI, I. 1998. Combining Instrument and Performance Models for High-Quality Music Synthesis. *Journal of New Music Research* 27, 211-238.
- DE POLI, G. 2004. Methodologies for expressiveness modeling of and for music performance. *Journal of New Music Research* 33, 189-202.
- DERENYI, I. AND DANNENBERG, R.B. 1998. Synthesizing Trumpet Performances. In *Proceedings of the 1998 International Computer Music Conference*, Ann Arbor, Michigan, October 1998, International Computer Music Association, San Francisco, CA, 490-496.
- DESAIN, P. AND HONING, H. 1993. Tempo Curves Considered Harmful. *Contemporary Music Review* 7, 123-138.
- EMDE, W. AND WETTSCHERECK, D. 1996. Relational Instance Based Learning. In *Proceedings of 13th International Conference on Machine Learning*, Bari, Italy, July 1996, SAITTA L., Eds. Morgan Kaufmann, 122-130.
- FINN, B. 2007. Personal Communication. August 2007
- FRIBERG, A., BRESIN, R., FRYDÉN, L., AND SUNDBERG, J. 1998. Musical Punctuation on the Microlevel: Automatic Identification and Performance of Small Melodic Units. *Journal of New Music Research* 27, 271-292.
- FRIBERG, A., BRESIN, R., AND SUNDBERG, J. 2006. Overview of the KTH rule system for musical performance. *Advances in Cognitive Psychology* 2, 145-161.
- GABRIELSSON, A. AND JUSLIN, P. 1996. Emotional expression in music performance: between the performer's intention and the listener's experience. *Psychology of Music* 24, 68-91.
- GABRIELSSON, A. 2003. Music performance research at the millenium. *Psychology of Music* 31, 221-272.
- GRINDLAY, G.C. 2005. Modelling expressive musical performance with Hidden Markov Models. PhD Thesis, University of Santa Cruz, CA
- HASHIDA, M., NAGATA, N., AND KATAYOSE, H. 2006. Pop-E: a performance rendering system for the ensemble music that considered group expression. In *Proceedings of 9th International Conference on Music Perception and Cognition*, Bologna, Spain, August 2006, BARONI M., R. ADDESSI, R. CATERINA, and M. COSTA, Eds. ICMP, 526-534.
- HASHIDA, M., NAGATA, N., AND KATAYOSE, H. 2008. jPop-E: An Assistant System for Performance Rendering of Ensemble Music. In *Proceedings of 2007 Conference on New Interfaces for Musical Expression (NIME07)*, CRAWFORD L., Eds. 313-316.

- HAZAN, A. AND RAMIREZ, R. 2006. Modelling expressive performance using consistent evolutionary regression trees. In Proceedings of 17th European Conference on Artificial Intelligence (Workshop on Evolutionary Computation), Riva del Garda, Italy, August 2006, BREWKA G., S. CORADESCHI, A. PERINI, and P. TRAVERSO, Eds. IOS Press,
- HILLER, L. AND ISAACSON, L. 1959. *Experimental Music. Composition with an Electronic Computer*. McGraw Hill, New York.
- HIRAGA, R., HASHIDA, M., HIRATA, K., AND KATAYOSE, H. 2002. RENCON: Toward a New Evaluation Method for Performance Rendering Systems. In Proceedings of the 2002 International Computer Music Conference, Gothenborg, Sweden, September 2002, ICMA, San Francisco, CA
- HIRATA, K., NOIKE, K., AND KATAYOSE, H. 2003. Proposal for a Performance Data Format. In Working Notes of IJCAI-03 Workshop on methods for automatic music performance and their applications in a public rendering contest,
- HIRAGA, R., BRESIN, R., AND HIRATA, K. & K.H. 2004. Rencon 2004: Turing Test for Musical Expression Proceedings of International Conference on New Interfaces for Musical Expression. In Proceedings of 2004 New Interfaces for Musical Expression Conference, Hamatsu, Japan, June 2004, NAGASHIMA Y. and M. LYONS, Eds. Shizuoka University of Art and Culture, 120-123.
- HIRATA, K. AND HIRAGA, R. 2002. Ha-Hi-Hun: Performance Rendering System of High Controllability. In Proceedings of the ICAD 2002 Rencon Workshop on Performance Rendering Systems, Kyoto, Japan, July 2002, 40-46.
- ISHIKAWA, O., AONO, Y., KATAYOSE, H., AND INOKUCHI, S. 2000. Extraction of musical performance rule using a modified algorithm of multiple regression analysis. In Proceedings of the International Computer Music Conference, Berlin, Germany, August 2000, International Computer Music Association, 348-351.
- JOHNSON, M.L. 1991. Toward an Expert System for Expressive Musical Performance. *Computer* 24, 30-34.
- JUSLIN, P. 2003. Five Facets of Musical Expression: A Psychologist's Perspective on Music Performance. *Psychology of Music* 31, 273-302.
- KATAYOSE, H., FUKUOKA, T., TAKAMI, K., AND INOKUCHI, S. 1990. Expression extraction in virtuoso music performances. In Proceedings of the 10th International Conference on Pattern Recognition, Atlantic City, New Jersey, USA, June 1990, IEEE Press, 780-784.
- KOELSCH, S. AND SIEBEL, W.A. 2005. Towards a neural basis of music perception. *Trends in Cognitive Sciences* 9, 579-584.
- KRUMHANS, C. 1991. *Cognitive Foundations of Musical Pitch*. Oxford University Press, Oxford.
- LERDAHL, F. AND JACKENDOFF, R. 1983. *A Generative Theory of Tonal Music*. The MIT Press, Cambridge.
- LEVENSHTAIN, V.I. 1966. Binary codes capable of correcting deletions, insertions, and reversals. *Soviet Physics Doklady* 10, 707-710.
- LIVINGSTONE, S.R., MUHLBERGER, R., BROWN, A.R., AND LOCH, A. 2007. Controlling Musical Emotionality: An Affective Computational Architecture for Influencing Musical Emotions. *Digital Creativity* 18
- LOPEZ DE MANTARAS, R. AND ARCOS, J.L. 2002. AI and music: From composition to expressive performances. *AI Magazine* 23, 43-57.
- MAZZOLA, G. 2002. *The Topos of Music - Geometric Logic of Concepts, Theory, and Performance*. Birkhäuser, Basel/Boston.
- MEYER, L.B. 1957. Meaning in Music and Information Theory. *Journal of Aesthetics and Art Criticism* 15, 412-424.
- MIRANDA, E.R. 2001. *Composing Music With Computers*. Focal Press, Oxford, UK.
- MIRANDA, E.R. 2003. On the evolution of music in a society of self-taught digital creatures. *Digital Creativity* 14, 29-42.
- MITCHELL, T. 1997. *Machine Learning*. McGraw-Hill, New York.
- NARMOUR, E. 1990. *The Analysis and Cognition of Basic Melodic Structures: The Implication-Realization Model*. The University of Chicago Press, Chicago.
- PACE, I. 2007. Complexity as Imaginative Stimulant: Issues of Rubato, Barring, Grouping, Accentuation and Articulation in Contemporary Music, with Examples from Boulez, Carter, Feldman, Kagel, Sciarrino, Finnissey. In Proceedings of the 5th International Orpheus Academy for Music & Theory, Gent, Belgium, April 2007

- PALMER, C. 1997. Music performance. *Annual Review of Psychology* 48, 115-138.
- PAPADOPOULOS, G. AND WIGGINS, G.A. 1999. AI Methods for Algorithmic Composition: A Survey, A Critical View, and Future Prospects. In *Proceedings of the AISB'99 Symposium on Musical Creativity*, AISB
- PARNCUTT, R. 1997. Modeling piano performance: Physics and cognition of a virtual pianist. In *Proceedings of 1997 International Computer Music Conference*, Thessalonika, Greece, September 1997, COOK P.R., Eds. ICMA, San Francisco, CA, 15-18.
- PENNYCOOK, B.W. 1985. Computer-Music Interfaces: A Survey. *Computing Surveys* 17, 267-289.
- QUINLAN, J.R. 1993. *C4.5: Programs for Machine Learning*. Morgan Kaufmann
- RAMIREZ, R. AND HAZAN, A. 2007. Inducing a Generative Expressive Performance Model using a Sequential-Covering Genetic Algorithm. In *Proceedings of 2007 Genetic and Evolutionary Computation Conference*, London, UK, July 2007, ACM Press
- RAMIREZ, R. AND HAZAN, A. 2005. Modeling Expressive Performance in Jazz. In *Proceedings of 18th International Florida Artificial Intelligence Research Society Conference (AI in Music and Art)*, Clearwater Beach, FL, USA, May 2005, AAAI Press, 86-91.
- RAPHAEL, C. 2001a. Can the Computer Learn to Play Music Expressively? In *Proceedings of Eighth International Workshop on Artificial Intelligence and Statistics*, 2001, JAAKKOLA T. and T. RICHARDSON, Eds. Morgan Kaufmann, San Francisco, CA, 113-120.
- RAPHAEL, C. 2001b. A Bayesian Network for Real-Time Musical Accompaniment. *Neural Information Processing Systems* 14
- RAPHAEL, C. 2003. Orchestra in a Box: A System for Real-Time Musical Accompaniment. In *Proceedings of 2003 International Joint Conference on Artificial Intelligence (Working Notes of RenCon Workshop)*, Acapulco, Mexico, August 2003, GOTTLÖB G. and T. WALSH, Eds. Morgan Kaufmann, 5-10.
- REPP, B.H. 1990. Composer's pulses: Science or art. *Music Perception* 7, 423-434.
- SEASHORE, C.E. 1938. *Psychology of Music*. McGraw-Hill, New York.
- SETHARES, W. 2004. *Tuning, Timbre, Spectrum, Scale*. Springer, London
- SUNDBERG, J., ASKENFELT, A., AND FRYDÉN, L. 1983. Musical performance. A synthesis-by-rule approach. *Computer Music Journal* 7, 37-43.
- SUNDBERG, J., FRIBERG, A., AND BRESIN, R. 2003. Attempts to reproduce a pianist's expressive timing with Director Musices performance rules. *Journal of New Music Research* 32, 317-325.
- SUZUKI, T., TOKUNAGA, T., AND TANAKA, H. 1999. A Case Based Approach to the Generation of Musical Expression. In *Proceedings of the 16th International Joint Conference on Artificial Intelligence*, Stockholm, Sweden, August 1999, Morgan Kaufmann, San Francisco, CA, USA, 642-648.
- SUZUKI, T. 2003. Kagurame Phase-II. In *Proceedings of 2003 International Joint Conference on Artificial Intelligence (Working Notes of RenCon Workshop)*, Acapulco, Mexico, August 2003, GOTTLÖB G. and T. WALSH, Eds. Morgan Kauffman, Los Altos, CA,
- TEMPERLEY, D. AND SLEATOR, D. 1999. Modeling meter and harmony: a preference rule approach. *Computer Music Journal* 23, 10-27.
- THOMPSON, W.F. 1989. Composer-specific aspects of musical performance: An evaluation of Clynes's theory of pulse for performances of Mozart and Beethoven. *Music Perception* 7, 15-42.
- TOBUDIC, A. AND WIDMER, G. 2003a. Relational ibl in music with a new structural similarity measure. In *Proceedings of the 13th International Conference on Inductive Logic Programming*, Szeged, Hungary, September 2003, HORVÁTH T. and A. YAMAMOTO, Eds. Springer Verlag, Berlin, 365-382.
- TOBUDIC, A. AND WIDMER, G. 2003b. Learning to Play Mozart: Recent Improvements. In *Proceedings of the IJCAI'03 Workshop on Methods for Automatic Music Performance and their Applications in a Public Rendering Contest (RenCon)*, Acapulco, Mexico, August 2003, HIRATA K., Eds.
- TOBUDIC, A. AND WIDMER, G. 2003c. Technical Notes for Musical Contest Category. In *Proceedings of 2003 International Joint Conference on Artificial Intelligence (Working Notes of RenCon Workshop)*, Acapulco, Mexico, August 2003, GOTTLÖB G. and T. WALSH, Eds. Morgan Kauffman, Los Altos, CA,

- TOBUDIC, A. AND WIDMER, G. 2005. Learning to play like the great pianists. In Proceedings of the 19th International Joint Conference on Artificial Intelligence, Edinburgh, UK, August 2005, KAEHLING P. and A. SAFFIOTTI, Eds. Professional Book Center, USA, 871-876.
- TODD, N.P. 1985. A model of expressive timing in tonal music. *Music Perception* 3, 33-58.
- TODD, N.P. 1989. A computational model of Rubato. *Contemporary Music Review* 3, 69-88.
- TODD, N.P. 1992. The Dynamics of Dynamics: A Model of Musical Expression. *Journal of the Acoustical Society of America* 91, 3540-3550.
- TODD, N.P. 1995. The kinematics of musical expression. *Journal of Acoustical Society of America* 97, 1940-1949.
- TOOP, R. 1988. Four Facets of the New Complexity. *CONTACT* 32, 4-50.
- WIDMER, G. 2000. Large-scale induction of expressive performance rules: first quantitative results. In Proceedings of the 2000 International Computer Music Conference, Berlin, Germany, September 2000, ZANNOS I., Eds. International Computer Music Association, San Francisco, CA, 344-347.
- WIDMER, G. 2002. Machine discoveries: A few simple, robust local expression principles. *Journal of New Music Research* 31, 37-50.
- WIDMER, G. 2003. Discovering simple rules in complex data: A meta-learning algorithm and some surprising musical discoveries. *Artificial Intelligence* 146, 129-148.
- WIDMER, G. AND TOBUDIC, A. 2003. Playing Mozart by analogy: Learning multi-level timing and dynamics strategies. *Journal of New Music Research* 32, 259-268.
- WIDMER, G. 2004. Computational Models of Expressive Music Performance: The State of the Art. *Journal of New Music Research* 33, 203-216.
- WIDMER, G. AND GOEBL, W. 2004. Computational Models of Expressive Music Performance: The State of the Art. *Journal of New Music Research* 33, 203-216.
- WITTEN, I.H. AND FRANK, E. 2000. *Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations*. Morgan Kaufmann, San Mateo, CA.
- WRIGHT, M. AND BERDAHL, E. 2006. Towards Machine Learning of Expressive Microtiming in Brazilian Drumming. In Proceedings of the 2006 International Computer Music Conference, New Orleans, USA, November 2006, ZANNOS I., Eds. ICMA, San Francisco, CA, 572-575.
- ZHANG, Q. AND MIRANDA, E.R. 2006a. Evolving Musical Performance Profiles using Genetic Algorithms with Structural Fitness. In Proceedings of the 8th annual conference on Genetic and evolutionary computation, Seattle, Washington, July 2006, DIGGELEN J.V., M.A. WIERING, and E.D.D. JONG, Eds. ACM Press, New York, USA, 1833-1840.
- ZHANG, Q. AND MIRANDA, E.R. 2006b. Towards an Interaction and Evolution Model of Expressive Music Performance. In Proceedings of the 6th International conference on Intelligent Systems Design and Applications, Jinan, China, October 2006, CHEN Y. and A. ABRAHAM, Eds. IEEE Computer Society, Washington, DC, USA, 1189-1194.
- ZHANG, Q. AND MIRANDA, E.R. 2007. Evolving Expressive Music Performance through Interaction of Artificial Agent Performers. In Proceedings of ECAL 2007 Workshop on Music and Artificial Life (MusicAL 2007), Lisbon, Portugal, September 2007