



HAL
open science

No title

Olivier Lartillot, Shlomo Dubnov, Gérard Assayag, Gill Bejerano

► **To cite this version:**

Olivier Lartillot, Shlomo Dubnov, Gérard Assayag, Gill Bejerano. No title. 8èmes Journées d'Informatique Musicale, Jun 2001, Bourges, France. pp.113-119. hal-01105887

HAL Id: hal-01105887

<https://hal.science/hal-01105887>

Submitted on 2 Oct 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Automatic Modeling of Musical Style

O. Lartillot¹, S. Dubnov², G. Assayag¹, G. Bejerano³

¹Ircam (Institut de Recherche et Coordination Acoustique/Musique), Paris, France

²Ben Gurion University, Israel

³Institute of Computer Science, Hebrew University, Jerusalem, Israel

email: lartillo@ircam.fr, dubnov@bgumail.bgu.ac.il, assayag@ircam.fr, jill@cs.huji.ac.il

Abstract

In this paper, we describe and compare two methods for unsupervised learning of musical style, both of which perform analyses of musical sequences and then compute a model from which new interpretations / improvisations close to the original's style can be generated. In both cases, an important part of the musical structure is captured, including rhythm, melodic contour, and polyphonic relationships. The first method is a drastic improvement of the Incremental Parsing (IP) method, a method derived from compression theory and proven useful in the musical domain. The second one is an application to music of Prediction Suffix Trees (PST), a learning technique initially developed for statistical modeling of complex sequences with applications in linguistics and biology.

1 Style Modeling

By Style Modeling, we imply building a computational representation of the musical surface that captures important stylistic features hidden in the way patterns of rhythm, melody, harmony and polyphonic relationships are interleaved and recombined in a redundant fashion. Such a model makes it possible to generate new instances of musical sequences that respect this explicit style (Assayag & al 1999a). It is therefore an analysis by synthesis scheme, where the closeness of the synthetic to the original may be evaluated and may validate the analysis. Our approach is unsupervised, that is we want an automatic learning process that may be run on huge quantities of musical data. Interesting applications include style characterization tools for the musicologist (Dubnov & al 1998), generation of stylistic meta-data for intelligent retrieval in musical data

bases, convincing music generation for web and game applications, machine improvisation with human performers, computer assisted composition.

The interesting aspect of our method can be seen if the whole process is considered as a sort of statistical learning algorithm. By statistical learning we mean a method of acquiring certain statistical properties of a data source so that new sequences can be created, having the same properties as the source. One of the main purposes of learning is creating a capability to sensibly generalize. Composers are interested in finding out the possibilities of certain musical material, without necessarily "explaining" it. So, let us consider here the possibilities that a piece (or a set of pieces in a given style) offer. Statistical analysis of a corpus reveals some of the recombination possibilities that comply to constraints or redundancies typical of the particular style. The concept of redundancy is closely related to information or entropy.

2 The IP Method

The Incremental Parsing algorithm is inspired by the analysis part of compression techniques of the Lempel-Ziv family. To understand how an idea derived from the compression field might be useful for our purpose, it is important to see that, as has been stated by several authors, compressing is equivalent to understanding, because in order to encode efficiently incoming information one has to perform a fine analysis of the way redundancy is organized. Compression algorithm can be split into two phases: first, it reads the input sequence and constructs a model that captures redundancy, and then it generates the compressed code of this sequence with respect to the model. In our case, the second phase is replaced by a stochastic navigation through the model in order to generate new sequences.

2.1 Description

During the generation process, a context-inference scheme is applied. The sequence formed by recently generated objects (a particular suffix of this sequence) is the context, from which a prediction on the next object to come is made with regard to a contextual probability distribution. So, the analysis part must provide a dictionary of such possible contexts along with their possible continuations.

First, the dictionary is reduced to the empty pattern and IP incrementally reads the sequence. At each cycle, it selects a pattern, from the current position to a further position, such that this pattern is the shortest one which is not already in the dictionary. Every left prefix of this pattern may become a context, and every object that follows this prefix may become a continuation. It is easy to see that an optimal representation (in space) for such a dictionary is a prefix tree where a branch descending from the root to any node is a context, the childs of any nodes are the continuation objects to the context going from root to this node, and where the cardinalities of the subtrees at a certain node encodes the probability distribution for the objects at the the root of each subtree. An optimal representation in generation time is a suffix tree where the context are reversed (i.e. from a leaf to the root) and continuations stored as pointers associated to every node.

The probability of each continuation is derived according to the length of the branches of the subtree at every node. We associate to each node in the dictionary tree a weight which is the number of nodes¹ that belong to the subtree which this node is a root of. Thus all leaves in the tree get a number 1, the root gets the total number of nodes² and the weight associated with every node is the sum of the weight associated with all its descendents. The probability of continuing in one of the deciples is the ratio between the weight of the child note and the weight associated to the current node.

Asymptotically it has been shown that IP predictor outperforms a Markov predictor of any fixed finite order. This surprising property of the IP scheme derives from the counting interpretation of the IP procedure. In this interpretation, the IP predictor is viewed as a set of sequential predictors operating on separate bins, where each phrase derived in the process of IP parsing is referred as the

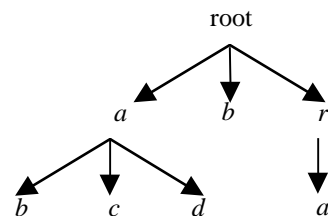
bin label. Since IP is unbounded in its length, it can be shown that for any finite k-th order Markov predictor, the long contexts of IP serve as refinements of the k-th order Markov predictors for sufficiently long sequences. Thus, the total number of errors due to a long context prediction is smaller for IP than the error in the case of a limited memory Markov predictor. When calculating the overall error regime, it turns out that asymptotically the longer terms dominate (since the length of the string grows) and eventually the IP scheme outperforms any finite order Markov predictor.

2.2 Example

Here is how IP analyzes the very simple sequence “*abracadabra*”. First, as we said, the dictionary contains only the empty pattern. Therefore, the shortest pattern starting from the beginning of the sequence that is not already in the dictionary is simply “*a*”. This pattern is added to the dictionary and IP goes on its analysis on the suffix of the pattern, that is “*bracadabra*”. Once again, the shortest pattern consists only of the first letter of the sequence, that is “*b*”. Idem for the next pattern “*r*”.

At this point, the current remaining sequence is “*acadabra*” and the dictionary contains “”, “*a*”, “*b*”, “*r*”. Now the shortest pattern is no more “*a*”, because it already belongs to the dictionary, but “*ac*”. The next patterns are “*ad*”, “*ab*” and “*ra*”.

The corresponding prefix tree is therefore :



Now for each pattern in the dictionary, the last letter is considered as the continuation of its prefix. In this way, for the pattern “*a*”, “*b*” and “*r*”, we say that *a*, *b*, and *r* are three possible continuations of the context “”; for “*ab*”, “*ac*”, and “*ad*”: *b*, *c* and *d* are three possible continuations of the context “*a*”. The probabilities associated to these continuations are computed as below. We have in this very simple example three contexts : “”, “*a*” and “*r*”. The corresponding suffix tree is:

¹ A variant of the algorithm considers the number of leaves.

² The variant considers the total number of leaves.

context : ""
 continuations : a (4/7), b (1/7), r (2/7).

context : "a" context : "r"

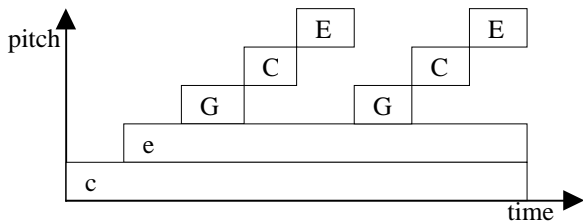
continuations : b (1/3), c (1/3), d (1/3). continuation : a (1/1).

The generation process is incremental too. The first letter is one continuation of the context "". The next letters are one continuation of the longest context that is a suffix of what has been already generated. In our simple example, if the process generates a *a*, then the next letter is one continuation associated with the context "a" ; idem for *r* ; else the context is "".

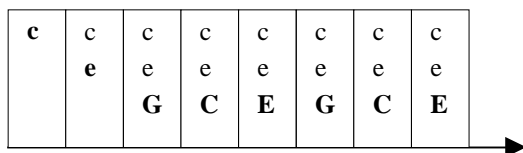
2.3 Music as a Sequence

We have loosely used the term sequence as an ordered list of objects. In order to capture a significant amount of musical substance, we shall, in a pre-analytic phase, cut the musical data (generally in Midi format) into slices which beginning and end are determined by the appearance of new events and the extinction of past events. Every slice has a duration information, and contains a series of channels, each of which contains pitch and velocity information and whatever available musical parameters. These slices are serialized in sequences submitted to analysis (these slices will be referred to by the word object or symbol, and the set of possible symbols as alphabet).

Here is for example the piano roll representation of the beginning of Bach's prelude in C, where lower case letters represent the third octave, and upper case letters the fourth one.



And here is how this sequence is sliced into symbols. The bold letters represent beginning of notes.



2.4 Improvements

We have described here the basic IP algorithm. This algorithm had already been tested with interesting results, but had certain drawbacks that made it quite impractical in specific situations. The improvements presented here are divided into four sections: pre-analytic simplification, generative constraints, loop escape, analysis-synthesis parameter distribution.

Pre-analytic Simplification. Real musical sequences - for example MIDI files of a piece interpreted by a musician - feature fluctuations of note onsets, durations and velocities, inducing a complexity which fools the analysis: the alphabet size tends to grow in an intractable way, leading to unexpected failures, and poor generalisation power. We have thus developed a toolkit containing five simplification filters:

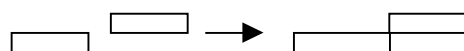
- the *arpeggio filter* vertically aligns notes which are attacked nearly at the same time,



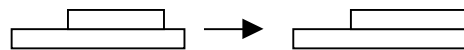
- the *legato filter* removes overlap between successive notes



- the *staccato filter* ignores silence between successive notes,



- the *release filter* vertically aligns note releases,



- the *duration filter* statistically quantizes the durations in order to reduce the duration alphabet.



These features can be tweaked by the user, using thresholds (e.g. a threshold of 50ms separates a struck chord on the piano from an intended arpeggiated chord). Using the simplification toolkit, Midifiles containing real performances that were intractable with the basic IP now become manageable, opening new perspectives, because this particular kind of musical data, full of idiosyncrasy, is of great value as a model for synthetic improvisation.

Generative Constraints. It is now possible to specify constraints during the synthetic phase. At each synthetic cycle, if the constraint is not respected by the new generated symbol, this generation is canceled and a new symbol is tried. If no symbol is satisfying, the algorithm backtracks to the previous cycle, and more if necessary. One interesting constraint is called the continuity constraint: at any cycle of the synthetic phase, it is possible that no context sequence be a suffix of the already generated sequence. The maximum context is thus the empty context. In these cases, the algorithm generates a continuation symbol of the empty context, that is to say, any symbol, with a stochastic model corresponding to its occurrence in the original sequence. The obtained result is a musical discontinuity. To avoid this, it is possible to specify a continuity constraint which imposes a minimum size of context anytime during the synthetic phase.

Loop Escape. The synthetic phase may easily enter into an infinite loop state. Here is one example: at one cycle, the maximal context A proposes only one continuation symbol; once this symbol is generated the new maximal context B features only one continuation symbol, and the new maximal context is A, again. The next contexts will be B, A, B, A, ... This tends to happen when the input data contains contiguous repetitions, which is often the case in music. We describe a mechanism to detect and escape these loops.

The previous example is very simple, because this loop was totally deterministic, and had only two states. At one state of a loop, it may be possible to consider several possible continuation symbols, but this choice leads, sooner or later, to return back to this present context or a previous one. We introduce the concept of context-generated subtree, which consists of the exhaustive set of all the possible contexts that may be met after the present one. Practically, we just need the size of this subtree, which is computed only once before the generation phase. Its computation consists of a marking, directly in the original tree, of all the possible future contexts and a counting of these markings. When the size of the context-generated subtree is below a user-specified threshold N , an N -order-loop is detected. The loop phenomenon is principally due to the fact that the synthetic phase searches for the maximal context, which is unique and proposes few alternative continuation symbols. In the case of a loop, we loosen this constraint and examine not only

the maximal context, but also smaller ones, which escape the loop in most cases.

Analysis-Synthesis Parameter Distribution. The analytic phase consists of finding the redundancy inside an original sequence of symbols. The trouble is there may be so much complexity and diversity in musical sequences that the size of the alphabet may be of the same order of the length of the sequence. Therefore, little redundancy would be observable. Moreover, each symbol consists, as said before, of several channels, each channel consists of notes, and each note consists of different parameters. So a symbol is in fact a Cartesian product of several musical parameters.

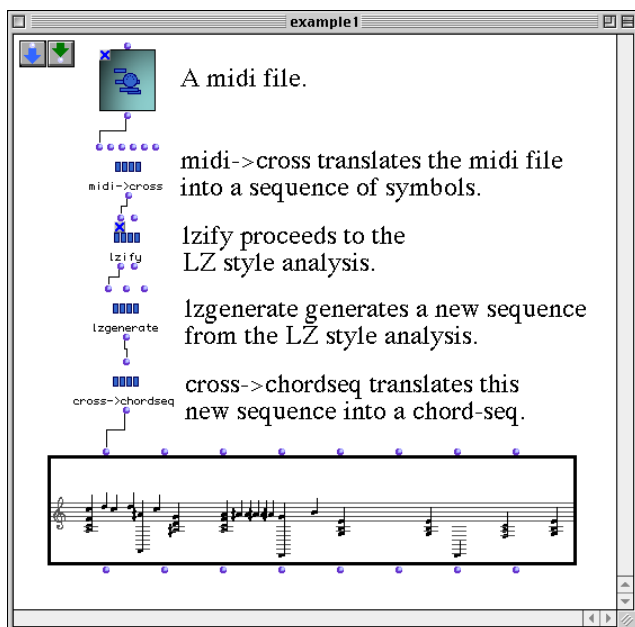
In order to increase abstraction and power in the analysis, we allow the system to discard some parameters, for example the velocity. The retained parameters will be called analysis information. In this way, it is obviously possible to increase redundancy, because it implicitly organizes the alphabet into equivalence classes: a chord struck two times with different velocities is nevertheless the same chord with the same harmonic function, and it will be detected as such. The problem is, in the synthetic phase, all the discarded information cannot be retrieved. For example, if we choose to discard note durations and dynamics, we finally obtain isorhythmic and dynamically flat musical sequences, which sounds like musical box production. A better solution is to store in the model the excluded information. This information is thus called synthetic information. The analytic phase is now performed on classes inside the initial symbol set, and during the synthetic phase, synthetic information, e.g. expressivity, may be reconstructed.

This solution has significant advantages. First, generated music regains much diversity, spirit and human appearance of the original one. Moreover, since it is possible to restrict analytic information and therefore find more redundancy in the original sequence, the synthetic phase becomes less constrained. At each synthetic cycle, every context features many more possible continuations than before.

3.2 Implementation

The software is implemented as a user library in an Open Source visual programming language developed at IRCAM, called OpenMusic (Assayag & al 1999b). Each step of the algorithm - pre-analysis, analysis, synthesis, and post-synthesis, is a function, which, in the musical representation

software, is represented by a box featuring inlets and outlets. All parameters may be tuned up by the user. Indeed, synthetic constraints, the distribution of analytic and synthetic information, and the reconstruction of the synthetic information may be explicitly formulated through visual expressions.



2.6 Musical Experiments

A lot of musical experiments have been carried in order to test the new IP algorithm. Midifiles gathered from several sources, including polyphonic music, piano music, and which style ranges from early music to hard-bop jazz have been submitted to the learning process. Experiments show that the combination of the simplification tool box and the new analysis-synthesis distribution scheme improves dramatically the results in the case of ‘live’ music, and in cases where the overall complexity leads to a huge alphabet. We show convincing examples, including, a set of piano improvisations in the style of Chick Corea, another one in the style of Chopin Etudes, polyphonic Bach counterpoint, 19th century symphonic music, and modern jazz style improvisations derived from a training set fed by several performers asked to play for the learning system.

3 PST

We have seen that IP predictor asymptotically outperforms a Markov predictor of any fixed finite order. In

practice the strings (music sequences) are of a finite order, and moreover, the size of the IP tree is bounded by a small finite size. Due to these limitations it is desirable to consider modeling schemes that might be more optimal for shorter-term situations. Another important feature of IP that seems redundant for our needs is the sequential nature of its operation. In our application, the goal is generation of new sequences that maintain similar statistical properties to the reference source. We use the prediction probabilities as the statistics generator, but we are not bounded by a requirement to rely on the past only. Allowing one to use the whole sequence for estimation of the statistics might help improve the performance.

Ron & al. (1996) developed a variable length Markov model termed Prediction Suffix Tree (PST). It has been shown that PST is a subclass of Probabilistic Finite Automata called PSAs. PSA is a variant order Markov chain in which the memory is variable and in principle unbounded. Given a finite size PSA, there exists an equivalent PST of a slightly larger size that produces the same probability output. Moreover, an efficient learning algorithm exists that allows one to construct a PST from samples generated by PSA. Now, if we consider the PSA as a common ground relating it to IP, we can consider the similarities and differences between the two approaches. One can see that both methods are similar in terms of their use of a variable length context for determining the probability for next symbol. The basic estimation procedure of the PST though is significantly different from that of IP.

3.1 Description

(Bejerano & al 1999) First, we define L to be the memory length of the PST, i.e. the maximal length of a possible string in the tree. We work out gradually through the space of all possible subsequences of length 1 through L , starting at single letter subsequences, and abstaining from further extending a subsequence whenever its empirical probability has gone below a certain threshold (P_{min}), or on having reached the maximal L length boundary. The P_{min} cutoff avoids an exponentially large (in L) search space. At the beginning of the search we hold a PST consisting of a single root node. Then, for each subsequence we decide to examine, we check whether there is some symbol in the alphabet for which the empirical probability of observing that symbol right after the given subsequence is non negligible, and is also significantly different (i.e. the

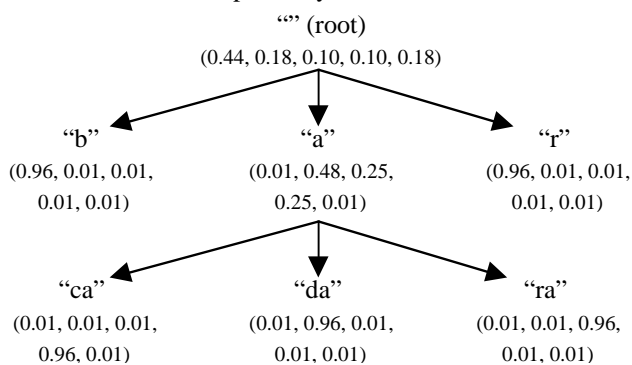
quotient exceeds a certain threshold r) from the empirical probability of observing that same symbol right after the string obtained from deleting the leftmost letter from our subsequence. This string corresponds to the label of the direct father of the node we are currently examining (note that the father node has not necessarily been added itself to the PST at this time). Whenever these two conditions hold, the subsequence, and all necessary nodes on its path, are added to our PST.

The reason for the two step pruning (first defining all nodes to be examined, then going over each and every one of them) stems from the nature of PSTs. A leaf in a PST is deemed useless if its prediction function is identical (or almost identical) to its parent node. However, this in itself is no reason not to examine its sons further while searching for significant patterns. Therefore, it may, and does happen that consecutive inner PST nodes are almost identical.

Finally, the node prediction functions are added to the resulting PST skeleton, using the appropriate conditional empirical probability, and then these probabilities are smoothed using a standard technique so that no single symbol is absolutely impossible right after any given subsequence (even though the empirical counts may attest differently).

3.2 Example

Here is the PST analysis of “*abracadabra*”, with $P_{min} = 0.1$, $r = 2$, $L = 10$ and a minimum smoothed probability of 0.01. For each node is associated the list of probabilities that the continuation be, respectively, a , b , c , d and r .



3.3 Implementation

This algorithm has just been adapted to Open Music, and integrated in exactly the same framework than IP. Therefore, all the features presented in the paragraphs 2.3 to 2.6 are also disponible for PST.

4 IP/PST comparison

4.1 Batch vs. On-line

PST parsing is not on-line in nature, as IP is — the training text is viewed as a whole unit and symbol frequencies are observed over the whole text. Thus there is no arbitrary parsing as in LZ where a single symbol change in the sequence may have deep influence on the dictionary structure. So PST may prove more powerful for short sequences but is not practicable for real time improvisation situations, where the on-line nature of IP will be adapted.

4.2 Selectivity

While IP goes over the training text and parses all of it, the PST learning algorithm looks at the text as a whole and picks for its dictionary only patterns considered relevant. This selective parsing can lead to more discontinuities and a smaller repertoire to improvise on, but the context/inference rules stores in the trees may be considered as more motivated than in the case of IP.

From these two remarks, one could say that IP is more adapted in generative applications, especially improvisation and real-time, and PST is more precise and complete, and is a good start for musicology applications where one wants to achieve the finest description.

4.3 Example

Thanks to the integration of both algorithms inside a unified framework, it is now possible to compare their respective analysis and generation for a specific musical example. Below is a score generated by Open Music of an interpretation of the beginning of a *Gnossienne* by Erik Satie followed by an improvisation of each algorithm. As the original example does not feature a lot of complexity, one can easily see that the generative process of both algorithms consists of a kind of sampling of parts of the original piece, so that the result be the most similar to the original sequence as possible. Nevertheless, these algorithms may sometime decide some transitions between the very little samples (which may consist of even one note) that were not expected, but that follow their modelling of the style of the original piece. As a result, the original style can be recognized, but in the same time, these generated sequences feature a certain amount of creativity. We invite you to listen to more complex improvisations, included in

the CDROM, in order to appreciate the creative skills of these algorithms, that are limited here by the simplicity of this didactic example.

5 Conclusion and New Directions

In the present state of researches, we have to consider music as a sequence of symbols. It could be more interesting to analyze directly the bi-dimensional score structure. This would enable an intelligent analysis of a fugue for example.

Although the context-inference approach may be compared with the implication/realisation view of Meyer, we have to acknowledge that his view is more powerful since his expectation is in long term, and not a systemic first order expectation like ours. But we suspect hard computability problems behind the long term approach.

In our work, learning is indeed unsupervised since we do not feed our system with musical knowledge, rather it analyzes music with constrained algorithms and does not proceed to any inductive inference. Another approach would be trying to induce automatically some conclusions about the presence of musical structures, with the help of minimal cognitive mechanisms.

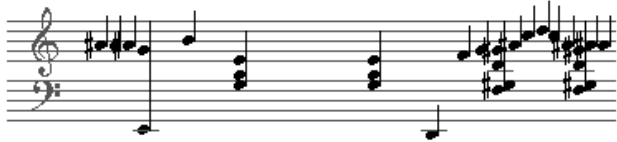
Other statistical techniques have to be experimented in music and compared to the two presented here. PPM for instance, can be thought of as going half-way between IP and PSTs. They are on-line (like IP) but they model k-terms of growing k with respect to data wealth (like PSTs).

Finally, one may try instead of modeling sequence $x_1 x_2 \dots x_n$ based on its own redundancies, to try and model $x_1 x_2 \dots x_n$ based on the redundancies of a second, correlated (eg 2nd voice) sequence $y_1 y_2 \dots y_n$ (transducer). Then, the generation of a X-type sequence could be constrained by an incoming Y-type sequence. This could lead to a synchronous improvisation scheme, where the computer, instead of providing 'answers' to a human player, as usual, would play synchronously with him, keeping an overall polyphonic consistency.

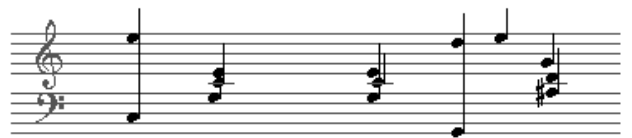
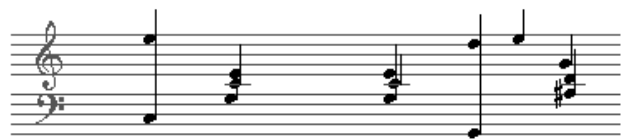
References

- Assayag, G., S. Dubnov and O. Delerue. 1999a. "Guessing the Composer's Mind : Applying Universal Prediction to Musical Style." *Proceedings of the International Computer Music Conference*. International Computer Music Association, pp. 496-499.
- Assayag, G., C. Rueda, M. Laurson, C. Agon and O. Delerue. 1999b. "Computer Assisted Composition at Ircam : PatchWork & OpenMusic.", *Computer Music Journal*, 23:3.
- Bejerano, G. and G. Yona. 1998. "Modeling protein families using probabilistic suffix trees." *Proceedings of RECOMB*.
- Dubnov, S., G. Assayag and R. El-Yaniv. 1998. "Universal Classification Applied to Musical Sequences." *Proceedings of the International Computer Music Conference*. International Computer Music Association, pp. 332-340.
- Ron, D., Y. Singer and N. Tishby. 1996. "The power of amnesia: Learning probabilistic automata with variable memory length." *Machine Learning* 25(2-3):117-149.

Gnossienne, by Erik Satie (beginning).



LZ improvisation.



PST improvisation.