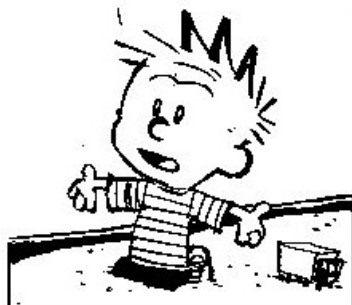


## Mood - A music recognizer

Auditory recognition and classification based on changes in attentional state.

YOU CAN'T JUST TURN ON  
CREATIVITY LIKE A FAUCET.  
YOU HAVE TO BE IN THE  
RIGHT MOOD.



[Click to see entire strip...](#)

### The team

**M**anolis Kalvin-selis - [manoli@mit.edu](mailto:manoli@mit.edu)

**O**vidiu Marina - [strider@mit.edu](mailto:strider@mit.edu)

**H**Obbes-man Vassef - [hvassef@mit.edu](mailto:hvassef@mit.edu)

**D**edric Carter - [dedric@mit.edu](mailto:dedric@mit.edu)

## Introduction: What does Calvin really mean?

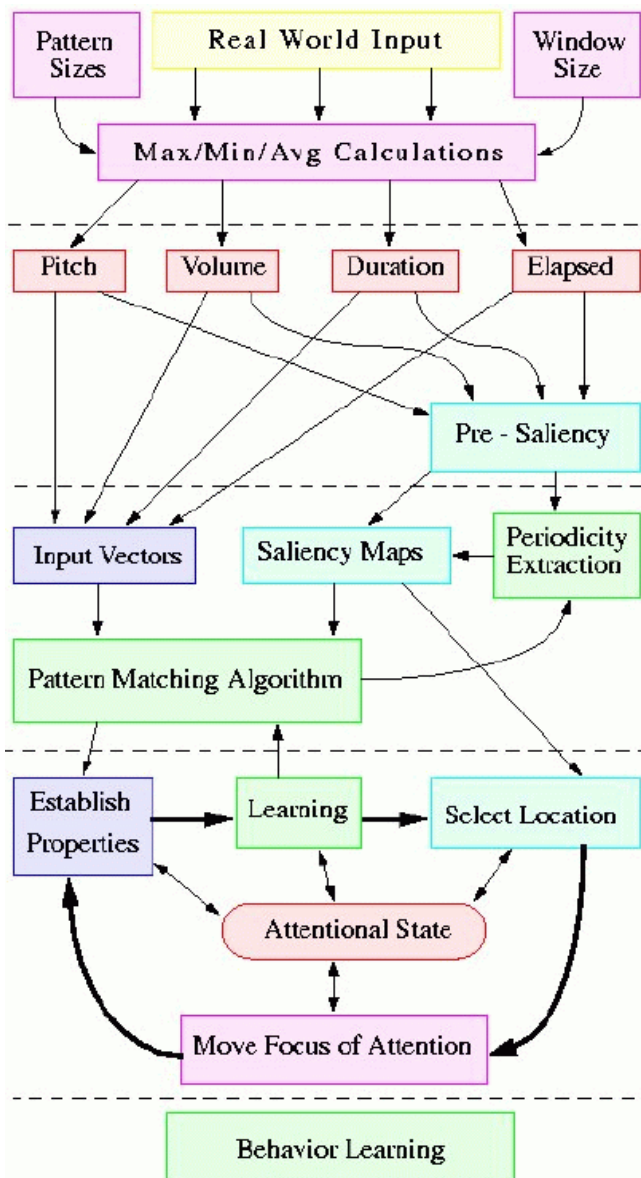
What Calvin describes here is the main idea motivating our project. To understand a pattern, you have to be in the right mood, the right attentional state. We constructed a pattern recognition system that modelled the state transitions of the human brain, as described by Seymour Ullman and implemented by Sajit Rao for the visual cortex. We used the same states that Sajit used for vision, namely:

- Select Location
- Focus attention on new location
- Establish Properties at the new attentional state
- Learn low-level patterns and patterns of state changes.

What makes the system able to generalize from one music piece to an entire category is the one thing that remains constant within a class of music compositions: namely, the pattern of changes in attentional state.

---

## Mood - Architecture Overview



### Input Layer

The input provided by the MIDI file is preprocessed to calculate Pitch, Volume, and Duration values, as well as a voice index. The voice index need is replaced by the time elapsed value, which makes the use of voices general enough to allow for a dynamic reordering of voices. Maximum, minimum, and average values within a window are calculated for each of them in a pre-processing step. [Input Details...](#)

### Saliency Layer

Based on these values and their maximums, minimums and averages, saliency values are calculated for each piece of data and their derivatives. This results in saliency maps which make different vectors stand out, depending on what feature one is paying attention to. [Saliency Details...](#)

### Pattern Matching

The pattern matching algorithm will be biased in matching the values which are salient on the vectors which are salient. Patterns of varying lengths will be stored for salient vectors, and matched as new inputs come in. [Pattern Details...](#)

### Prediction

When a pattern is matched partially with the current input, the system will predict future notes as the continuation of the stored pattern. Since more than one patterns will be matched (at least one for each pattern length - one for every pattern table), a weighted average will be calculated, depending on certainty and pattern length. [Prediction Details...](#)

## 1. Vision

### Input-Output Processing is at the heart of Human Intelligence

Some of the research in Artificial Intelligence research has been oriented mainly towards building a general computational theory of intelligence, of which human intelligence would be only one instance. However, some of the recent research in Artificial Intelligence has been primarily directed towards understanding human intelligence, taking the results of neuroscience and psychology into large consideration.

Our research fits in the second category, as a part of the effort to understand general intelligence through a study of animal and particularly human intelligence, the only form of intelligence we know. We use methods of Computer Science, and especially Computer Engineering, as tools to exploit the facts revealed in biology, psychology and neuroscience. Most importantly, **we believe that perception and interaction with the real world is at the heart of intelligence.**

In this perspective, our research is aimed at expressing, with computational tools, a model of human perceptual information processing. In particular, this project is a proposal to explore an **attention-based model for auditory processing**, inspired by the model proposed by Rao [1] and Ullman [2] for visual spatial perception.

### Why Audition, Why Attention?

A lot of work has been done in fields related to visual information processing, for example computer vision, character recognition, etc. However, less work has been done in auditory information processing, in areas other than language processing. Furthermore, only a small fraction of the work on vision is directly aimed at understanding practically how the human brain does it. Similarly, little work has been done in understanding how we listen to and react to sounds other than conversations, for instance, music. How do we understand, interpret, recognize music? Moreover, why does it affect our emotions, our mood?

For the human visual system, Ullman [1] proposed a theory of cognitive routines to account for the flexibility and polyvalence of our visual system in performing various complex information extraction tasks. Rao [2] proposed a model on how the routines can be performed and learned using a generic architecture around attention.

Our task in this project was to try out a model similar to Rao's model for musical perception. We chose music/audition instead of vision, to see if the generic model for spatial information processing is also a generic model for other senses, as Rao speculates. The concept of an attention-based model is an attractive one, for two reasons, which are related to each other:

1. It is what makes Rao's model generic for various visual tasks. It reduces their complexity by organizing them into serial routines, sequences of attentional states. It would make the model generic to all senses if we can apply it here.
2. Attention is at the center of the mystery of the human mind. It is goes together with the unresolved issue of how a machine could have consciousness, and be self-aware. Here, our engineer's approach is the following: if we include manifestations of attention as a computational tool in our model, maybe that will give us insights on the consciousness issue.

## Our Goals, for this project and beyond

- **Understand the human auditory system**
- **Unify the workings of the brain**  
: explain audition and vision under the same architecture, understand the hierarchy of the brain
- **Extend Rao's language of attention to auditory perception**
- **Extend the notion of implicit memory (routines) to understand explicit memory**
- **Attempt to recreate music in a learned style**

## Our Focus for this project

- **Perception of music**  
: to narrow down the extend of this first project, we abstracted out the lower-level auditory mechanisms, and made reasonable assumptions about them. In vision, the evidence from neuroscience and psychology shows that the brain separates the task in two: object recognition and spatial information processing. It allowed Rao to do a correct abstraction of object recognition while he focused on spatial information processing. Similarly, we abstract out the sound recognition process (which among other things, lets us recognize a particular instrument, or spoken language), and focus on the perception of music. To this end we made a set of reasonable assumptions: that the lower-level mechanisms yield to our system a notion of pitch, volume and duration.
- **Learning**  
: we will focus on learning musical routines. Understanding music composition will be left for later work.

## Our Contributions

- **Auditory World in terms of Attentional Patterns**  
: we proposed a new understanding of the auditory world in terms of changes in attentional patterns. Applying Rao's work to audio was no easy task, since at first sight, the two worlds have nothing in common.
- **Periodicity is Learned, not hard coded**  
: New theory of how periodicity emerges from matching successively larger patterns, and understanding emergent patterns.
- **Language of Attention**  
: Provided a language of attention for the understanding of musical pieces. Adding a new category to Rao's system, and filing in the auditory equivalents of Rao's visual routines.
- **Unifying visual and auditory world**  
: Same architecture of attention can be used to explain both visual and auditory worlds. Implications on evolution of those behaviors and on the separation and specialization of our brains to the different tasks it accomplishes today.

## 2. MIDI

### Real World Input

Although we abstract away the issues of note harmonics and playing instrument types, we use a keyboard as input. This allows for a more rapid creation of new tunes, as well as for future real-time input. The same mechanisms used in loading the MIDI file output from the

electric keyboard can also be used in loading pre-made MIDI files.

### Keyboard

Dedric playing puts in the mood. ;o)

The input is limited to one instrument at a time, although multiple notes can be played simultaneously. The notes can be loaded by our utilities from either the MIDI0 or MIDI1 file format. The notes are separated after loading into our system into different voices.

### Time ON / OFF to Durations

Matching notes can be tricky, because the MIDI format stores only time ON/OFF, with a code format which can be confusing. Dedric's program successfully decodes this information and returns a Pitch/Volume/Duration/Time on encoding for each of the notes played. This encoding includes a cleaning up of slight note overlaps.

### Distinguishing Voices

Ovidiu's code, after updating the internal MOOD note store, takes the new note and tags it with a voice. The voices denote notes which overlap, or notes which follow a pattern which is distinct from a separate set of notes being played in the same tune. This is akin to two objects which move about in space: each has a trajectory which is separate from the other object's trajectory. When doing vision, one of the first tasks is to separate potential objects based on their motion. The equivalent is being attempted here in the auditory realm.

As Manolis observes, this is yet another great MOOD success.

On to [Augmenting the Input...](#)

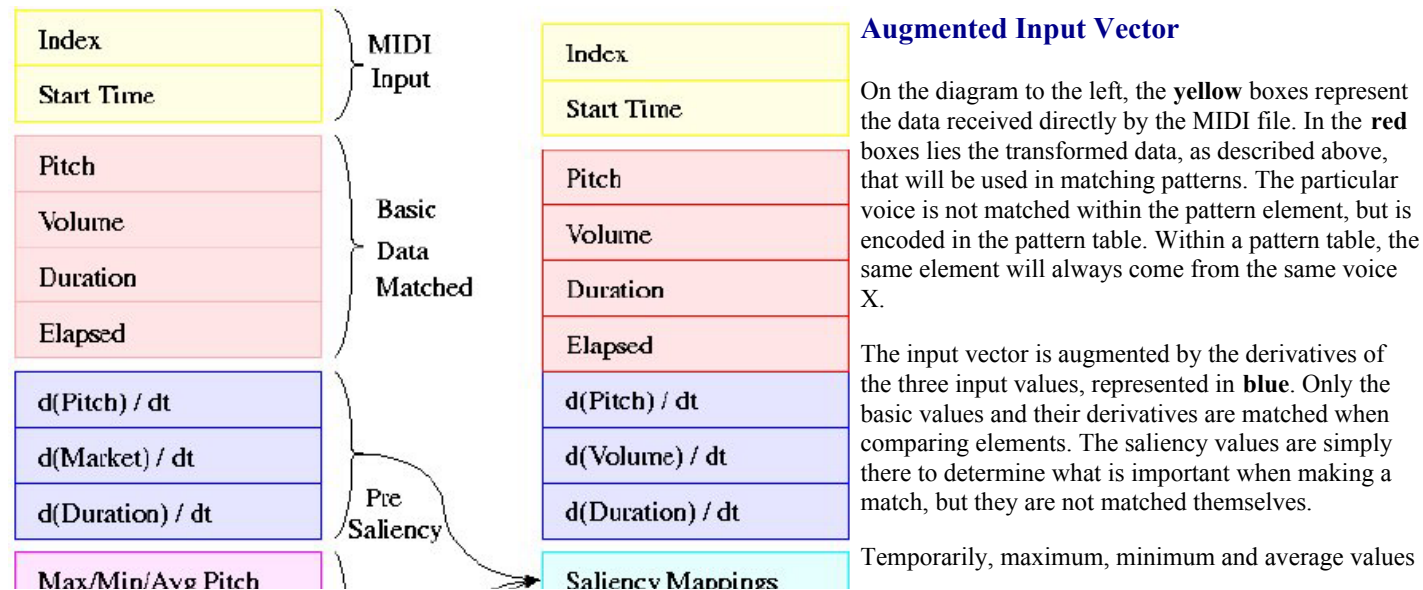
## 3. Input

### Duration Elapsed

The duration elapsed is the time difference between the current time and the time of the note. When a pattern of 5 consecutive notes is constructed, the duration fields will be [4, 3, 2, 1, 0], and they will match exactly any other pattern of consecutive notes. However, there is no constraint on which data will enter a pattern.

For example, a pattern table can have the first element come from voice X, the second from voice Y, and the third from voice Z. A pattern in that table could be "voice X falls 2 notes ago, voice Y rises 10 notes ago, and voice Z has maximum pitch now". In that case the duration fields will be [2, 10, 1] for voices [X Y Z]. The pattern will also match [2, 8, 1], where Y rose 8 notes ago instead. This allows the system to match a larger class of patterns independent of the exact time of occurrence, and which voice the notes belong to.

If such a duration field is not included, one has to train the system to recognize only cases where a pattern depends on exact time relationships. This varying duration field allows for more dynamic relationships between voices. Since pattern tables are created dynamically, depending on what data is salient, the relationships to look for need not be hard-coded. A different pattern table will be created for different combinations of notes, but that's another story, which belongs to the [pattern matching](#) domain.



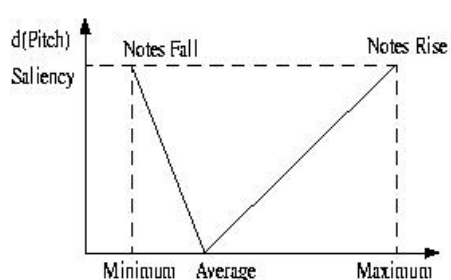
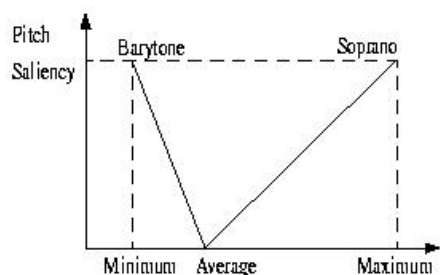
of each of the six fields are calculated for the window-size (**magenta**), and they are used to calculate the **saliency** maps (**cyan**). The window size represents how far back in the past we should look for normalizing values when comparatively analyzing the current data. It is relative to those min max and avg values that values are **matched** within patterns, as will be described later on. They also serve as a normalization step, so that data given in different units can be compared in an unbiased way.

Along with the other saliencies, a periodicity saliency is calculated, which makes salient notes based on a periodically salient history (**green**).

Thus we have reached from a 3-value input (pitch, volume, time) a 7-value vector (pitch, volume, duration, their derivatives, and time elapsed) to match upon. Adding periodicity, each vector also has an 8-field saliency, which determines what values should be matched more carefully.

On to [Saliency](#) or back to [overview](#).

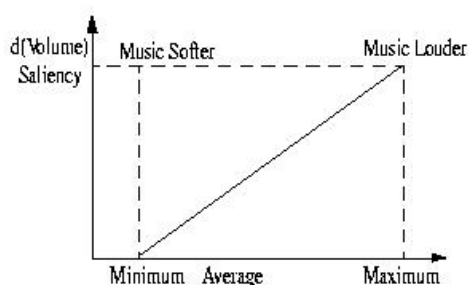
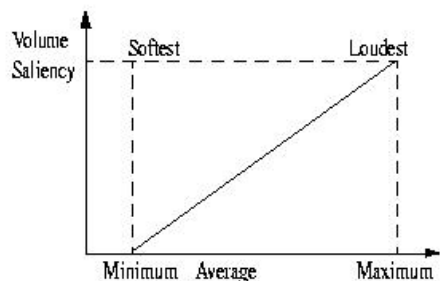
## 4. Saliency Maps



### Vector Saliency

Saliency is a value between 0 and 1 that indicates how important a vector or a field inside a vector is. The most important ones have a value, the least important, a saliency of 0.

The saliency value for a vector is computed as the average of individual saliencies calculated for each of the components.



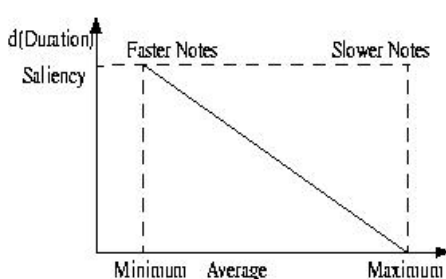
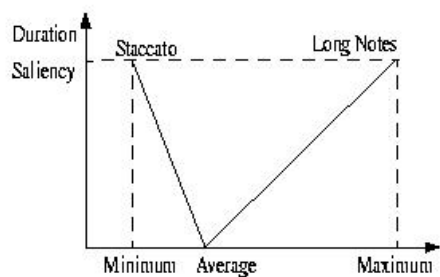
### Individual Saliencies

Individual saliencies for each value are computed relative to the window of attention. The highest value in the window will generally have a saliency value of 1, as well as the lowest one. An average value will have a saliency value of 0.

The highest volume in the current window yields a saliency of 1. The lowest has a saliency of 0.

The longest duration and the shortest one have saliencies of 1 (either long = salient or staccato = salient), and the average duration will have a saliency of 0.

Similarly for pitch - highest and lowest notes have saliency of 1, and average pitch has saliency of 0.



Derivatives follow the saliency rules as the components they come from.

### Pre-Saliency

Based on the average of these saliencies, each note gets a *pre-saliency* value, which is used to compute the periodicity of the input. This can lead the computer to expect salient notes at specific times if it has been seeing them periodically. Hence, if falling on a periodically salient time, even a silence could appear important.

### Saliency Maps

We have thus constructed many saliency maps. One can pay attention to notes salient in volume, or in pitch, or in duration, or in their derivatives. The pattern matcher also computes a periodicity saliency map, which gives a value of 1 to notes expected to be salient, and a 0 to the non-salient ones. Together with the pre-saliency values, this final addition constitutes the saliency part of the code.

On to [Pattern Matching](#)...

## 5. Pattern Matching

Learning Patterns of Size 2



### Pattern Matching

Pattern matching, along with learning the periodicity of the input, does the basic computations the code is based upon. It is able to learn patterns of different length from the input, and match partially completed patterns to either predict future inputs, or turn the focus of attention to different patterns interesting to a specific behavior.

The pattern matcher receives commands from the attention loop, which directs it on what to match, what to record (learn), and what to pay attention to.

	Pattern 0	Value	Certainty	Value	Certainty
	Pitch	1	1		
	Volume	1	1		
	Duration	2	1		
	d(Pitch) / dt	2	1		



$d(\text{Volume}) / dt$	1	1		
$d(\text{Duration}) / dt$	1	1		



Pattern 1	Value	Certainty	Value	Certainty
Pitch	1	1	1	1
Volume	1	1	1	1
Duration	2	1	2	1
$d(\text{Pitch}) / dt$	3	1	3	1
$d(\text{Volume}) / dt$	4	1	4	1
$d(\text{Duration}) / dt$	4	1	4	1

### Calculating Certainties

Every time a pattern is matched, individual elements get averaged, and their certainties updated. The certainty of a match between two values is one minus the difference between the two values over the difference between the maximum value and the minimum value at the current window. The certainty of a match between two pattern elements is the average of the certainties of their values. The certainty of a match between two patterns, is the average of the certainties of their elements, weighted by the saliency of each element. Intuitively, one will try to match within a sequence of notes the important notes, and will care less about less salient notes which do not match.

---

Table Size 2												
Pitch	1	1	1	1	1	1	1	1	1	1		
Volume	1	1	1	1	1	1	1	1	1	1		
Duration	2	2	1	1	2	2	1	1	2	2	1	1
d(Pitch) / dt	2	2	1	1	2	2	1	1	2	2	1	1
d(Volume) / dt	1	1	1	1	1	1	1	1	1	1	1	1
d(Duration) / dt	1	1	1	1	1	1	1	1	1	1	1	1

Table of Size 3																	
Pitch	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
Volume	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
Duration	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
d(Pitch) / dt	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
d(Volume) / dt	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
d(Duration) / dt	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1

Table of Size 4																
Pitch	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
Volume	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
Duration	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
d(Pitch) / dt	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
d(Volume) / dt	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
d(Duration) / dt	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1

### Pattern Tables

The figure shows an individual pattern of two elements. A pattern table will hold many such patterns. A different pattern table is created for every pattern length that we are matching. Moreover, for the same length (the same number of notes that we match upon), there will be different pattern tables for different voice combinations. For patterns of 3 notes for example, we can have a pattern table storing all patterns found on voices [a b c] another one storing all patterns found on [a b d], and so on.

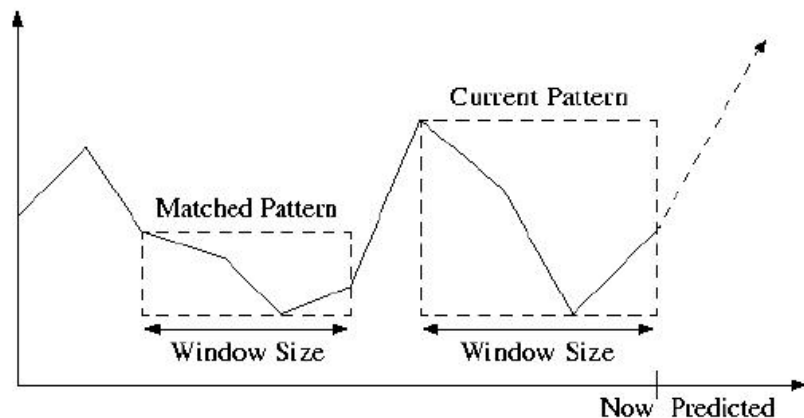
### Updating Values

When the certainty with which two patterns are matched is too low, then a new pattern is added to the pattern table.

When the certainty is high enough, then the matched pattern is extended to include the current pattern with which it matched.

## 6. Making a match





### Making a Match

In the figure, we would like that the second pattern enclosed in the box matches the first one, previously recorded. Clearly, the values do not match, nor do their derivatives, since everything seems to have doubled (for example in volume).

Therefore, when a match is made, what is compared is not absolute values, but instead values relative to the current maximum and minimum of the window-size considered.

### Making a Prediction

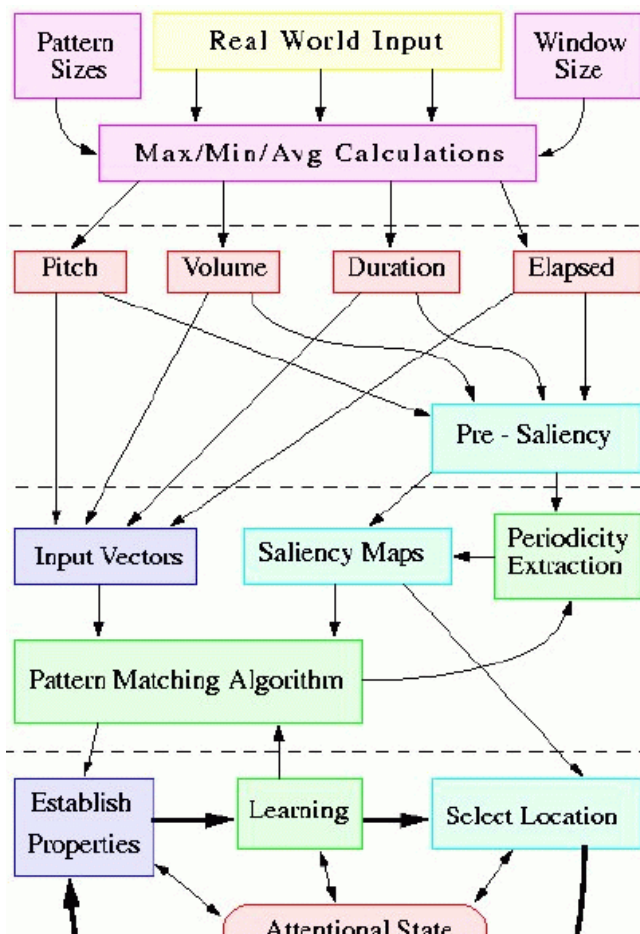
Predictions are made accordingly, relative to the current window maximum and minimum. If the current pattern has length 4, then the system will match the first 4 elements of a stored pattern of length 5, and predict the next note as the final element of the matched stored pattern.

Patterns of different sizes are matched in making the predictions. Their predictions are weighted then averaged, to obtain the final prediction. For each length, let best[length] be the pattern of that length that best matched the current pattern (of length-1). Each best[length] will have a certainty value associated with it, and with each element. This certainty value, along with the pattern length, are used in weighting the different size predictions.

The longer the length of the pattern is, the higher its weight should be. The reason for this bias of trusting longer patterns is that they are less likely to match. For almost anything, a two-note pattern will match with high certainty.

On to [Attention...](#) Mood - Report

## 7. Attentional Loop



The attentional loop is composed of four types of procedures, called serially always in the same order, establishing properties at specific locations, learning patterns, selecting locations, and moving to new locations. A behavior is specified by the commands that are going to be called from each family.

### Establish Properties

This is the largest of the procedure categories. It works on the specified focus of attention to establish local properties.

### Establish relations between properties

**High duration followed by rise in volume then silence/applauding**  
This could be hard coded in our brain. When we hear sssssssssSSSSSSSSHHH-BOOM!

### Determine movement type at current location

A movement could be either sad, or angry, or romantic - determine based on low-level + patterns

### Learning

The procedures in this category are responsible for the learning that gets done in the system. Different types of learning are possible, and they're basically controlling the pattern matching code into doing useful things.

### Learn patterns of length i-j ending at current note

Records the patterns either learning them as new entries in the pattern table or changing old patterns to be able to match the new input.

**Do not learn anything.**

If the input is not salient, or not interesting, or simply that we're sure that we're already matching something well known, we can simply not record anything.

**Learn pattern of patterns.**

One can simply abstract at a higher level and learn patterns of patterns instead of patterns of notes. Useful for more complex behaviors.

**Select Location**

There are different ways of selecting a location, and different criteria to base this selection upon. A location is either a note or a window size, or a voice to listen to.

**Select most salient overall**

One can simply pay attention to the note that is the most salient in the current window, or to the voice that's most salient

**Select most salient on a particular feature**

The feature can be volume, duration, pitch or their derivatives

**Select based on pattern**

From the pattern matching code, we can select the location that matches optimally the current pattern, or a particular pattern we are interested in

**Select on periodicity**

Location either where the next salient note should appear, or where a periodic pattern should start appear again.

**Move Focus of Attention**

As mentioned above, there are more than one type of locations to move the focus of attention to, and there are more than one ways to move the focus of attention there.

**Move to selected voice**

Changes the current voice, where patterns are recorded from

**Move to selected time**

Changes the current time to be something different than the next time

**Move to selected window size**

Changes the window size based on which saliencies are calculated

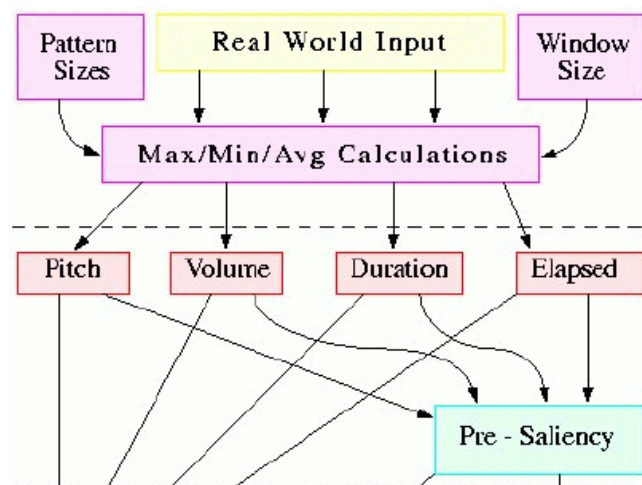
**Move to selected pattern**

Changes the current pattern to which we're matching

**Move completely or partially**

Finally, one can move there completely or still pay attention to many other things

## 8. Behavior



Behaviors are a specification of which procedures are called within each type in the attentional state loop.

**Attentional State**

The attentional state is the shared state which coordinates the different routines of the different categories. They are the means of communication among them. The properties established at the focus of attention will be saved there. The learning routines will keep there the patterns that the pattern matching algorithm matches upon and modifies. The location to move to (along with the type of location it is) will be saved for the MFOA procedure to read it.

It's basically a database of shared state, and provides no functionality by itself, but simply a place to scribble on at every iteration of the attentional loop.

## Choosing Procedures

The choice of procedures is thus crucial in understanding the music that is presented to the system. This choice can be either hard coded, to recognize specific types of music, according to how we think we understand a piece we listen to. A behavior will simply be a series of routines alternating types.

### Behavior Example

To recognize a fugue, the system will

- EP: Number of voices playing
- LN: nothing (no input long enough yet)
- SL: Select location of most salient notes overall
- MF: Move FOA to that voice
  
- EP: Pitch - Duration - Volume - Period to expect for next note
- LN: Learn patterns at selected pitch/duration/volume/period
- SL: Select voice where pattern repeats
- MF: move to voice (or adjust window-size = period when not found)
  
- EP: Changes between theme and its repetition
- LN: Learn changes as new patterns of length = length\_of(theme)
- SL: Select voice of pattern
- MF: move to voice

### Learning Procedures

Alternatively, those procedure sequencies can be learned, by selecting the most salient at each time, and seeing which one it was. When in the future it sees that the first routines emerging by following movements and most salient properties, it will match patterns of changes in attentional state, and follow those emerging patterns as behaviors.

Our hope is that similar behaviors correspond to similar types of music, and thus our program would be able to recognize and categorize types of music simply by the way it reacts to each of them.

## 9. Contributions

### Auditory World in terms of Attentional Patterns

We proposed a new understanding of the auditory world in terms of changes in attentional patterns. Applying Rao's work to audio was no easy task, since at first sight, the two worlds have nothing in common.

### Periodicity is Learned, not hard coded

New theory of how periodicity emerges from matching successively larger patterns, and understanding emergent patterns.

### Language of Attention

Provided a language of attention for the understanding of musical pieces. Adding a new category to Rao's system, and filling in the auditory equivalents of Rao's visual routines.

### Unifying visual and auditory world

Same architecture of attention can be used to explain both visual and auditory worlds. Implications on evolution of those behaviors and on the separation and specialization of our brains to the different tasks it accomplishes today.

## 10. Future Work

### Use more Complex and Accurate Musical Primitives

We reduced the input to a 3-dimensional vector for pitch, duration and volume. This is an oversimplification of our perception. One direction for future work is to explore the lower-level sound wave perception with the correct level of complexity. We would understand

how the many harmonics of an instrument enriches our perception of music.

## **Extend the model to the whole spectrum of Auditory Perception**

Our project expressed a language of attention for audition, within the limits of music perception, given one instrument. With a more accurate model for lower-level auditory perception, we can extend this language to all of the auditory world. We could understand how we perceive several instruments coherently together, or how we can select one instrument, one conversation among others, or one pattern of sounds in the middle of nature.

## **Understand Memory**

Visual and Auditory Routines are nothing else but elements of our implicit memory. With this insight, we can explore how implicit memory relates to explicit memory and understand how it all works in the brain.

## **Levels of Attention - Distributed Attention**

When we first learn how to play an instrument, we must pay very close attention to the movement of our fingers. But once the learning is done, we can play the piece like a routine, without even paying attention to what our fingers are doing, and focus on something else, for example a conversation with a friend. How does attention move horizontally between senses? How does attention move vertically between levels of processing, when it can follow every movement of the fingers or follow the global flow of a routine? Furthermore, how do we distribute attention, in other words pay attention to several things at the same time?

## **Understand the Hardware**

Having a unified model for auditory and visual perception can help us understand the functions of the neural structure of the human brain. How is attention implemented in our brain?

## **Attention and Consciousness**

If we can put attention at the center of all our perception, will that explain the notion of consciousness and self-awareness?

## **Music and Emotions**

What is beautiful about music? How can music affect our mood, make us sad, happy, angry, energetic, or even empowered or enraged? Music has the potential to create a strong emotional response in people, sometimes pushing it to an extreme. For example, Muslims give to some forms of music the same forbidden status as alcohol or gambling, because such music is believed to alienate people. If our model shows how we perceive music, then by studying what patterns of music perception affect our emotional state, we might get insights about how emotions work.