

# SoundThimble: A High Resolution Gesture Sonification Framework

Ben Trovato  
Institute for Clarity in  
Documentation  
1932 Wallamaloo Lane  
Wallamaloo, New Zealand  
trovato@corporation.com

Grigore Burloiu  
CINETic  
UNATC "I.L. Caragiale"  
Bucharest, Romania  
gburloiu@gmail.com

G.K.M. Tobin  
Institute for Clarity in  
Documentation  
P.O. Box 1212  
Dublin, Ohio 43017-6221  
webmaster@marysville-  
ohio.com

Bogdan Golumbeanu  
CINETic  
UNATC "I.L. Caragiale"  
Bucharest, Romania  
bogdangolumbeanu@yahoo.com

## ABSTRACT

The installation is based on a high-resolution, state-of-the-art motion capture system consisting of eight high-resolution, infra-red cameras. The system is used alongside MaxMSP to track, interpret and sonify the movement and gestures of a performer in 3d space.

## Author Keywords

sonification, motion tracking, gesture spotting, interactive installation, synthesis

## ACM Classification

H.5.5 [Information Interfaces and Presentation] Sound and Music Computing, H.5.2 [Information Interfaces and Presentation] User Interfaces—Auditory (non-speech) feedback

## 1. INTRODUCTION

- motivation
- challenges
- the Vicon system

The SoundThimble project started with the aim of exploring the technical possibilities offered by the Vicon technology and implementing them into the realm of Sonic Arts and interactive sound installations. Three-dimensional motion capture systems — including Vicon — have been extensively used in animation, movies and games where actors’/performers’ movement in 3d space is being tracked with a high degree of both spatial and temporal resolution (scieti voi pls valorile exacteă?mm discriminare / ??fps) and assigned to animate various models such as characters, animals, aliens and other creatures. Although motion capture devices such as video cameras and Microsoft Kinect devices have provided researchers and sound artist multiple expressive possibilities, the idea of using advanced

motion capture systems like Vicon is relatively new and un-exploited. (ref despre Kinect, proiecte misto, proiecte cu Vicon)

Vicon Vantage system used in this project contains 8 5-megapixel cameras which transmit movement in realtime to Nexus. Having accelerometers and temperature sensors, each camera detects any decalibration of the system. <sup>1</sup>.

In addition to a high spatial resolution and discrimination of only (cati mm?) the ability of placing a large number of markers makes it possible to track head, hands, elbows, shoulders, torso, feet, thighs and even fingers. By assigning numerous markers to groups, there is possible to reduce the chance of the system temporarily losing markers, which can happen when using only one marker per limb because cameras’ blind spots. These characteristics result in a very accurate and complex motion tracking.

==== FIGURE : Nexus / Bogdan ====

## 2. STATE OF THE ART

- Vicon & related projects
  - interactive / movement sonification examples[2].
  - *Sound in space* represents another innovative element in this project because introduces the idea of *controlling sound* by movement. This means that the sound is an entity, gets a materialization and it becomes switchable [1]. It is not about the **localization** of sound in space (detecting the direction of sound source), it is about its **position** in space, the coordinates of the sound object in space, like an actual *object* in a room. *Sound objects* concept represents an innovative tool for multimedia arts such as sketches, imaginary games and realtime interactions.

Interaction between sound control and human gesture has constantly increased over the last years [3]. Probabilistic models for analysing motion and sound relationships became a necessity and a forthcoming tool [4].

## 3. PROJECT DESCRIPTION

### 3.1 Concept

==== FIG: schema logica cu etape & gesturi recording/recall  
====

- Gestures, virtual objects, dynamic mapping



Licensed under a Creative Commons Attribution 4.0 International License (CC BY 4.0). Copyright remains with the author(s).

NIME'17, May 15-19, 2017, Aalborg University Copenhagen, Denmark.

<sup>1</sup>See <https://www.vicon.com>.

SoundThimble can be described as either being an interactive sound installation or an auditory game. - 3 etape: finding, manipulation, arrangement.

The narrative of SoundThimble implies a performer which tries to find a virtual object, randomly placed in 3d space, by analyzing cues that are constantly shifting in the sonic fabric in relation to his/her movement. Similarly to the traditional game, the distance between the performer and the virtual object is correlated to sound synthesis and modulation parameters. The closer one gets to the object the more organized and constant the sound and vice-versa.

Once the object is found, the performer has full freedom to record a number of gestures that can be re-performed later, re-called, and used to trigger or manipulate various sound events/sonic shifts.

We can think of SoundThimble as having two different narratives: one before the object is found, where only distance in relation to the object is used to control sound and the second where gestures are being employed, after the object has been found. Generally, the level of interaction is higher after the object is found, however both situations can overlap (e.g: the performer tries to find the second object while recalling gestures that were assigned to the first object). This scenario keeps repeating: other objects are randomly generated, the performer finds them, defines gestures and interacts sonically with them.

- Performance aesthetic

## 3.2 Implementation

=== FIG: Framework Diagram: Camere, Nexus, C++, Max, Speakers ===

One of the first challenges was to create a method of sending real-time from Vicon to MaxMSP where data would be processed and used to generate and manipulate sound. By default, the Vicon system does not support the Open Sound Control (OSC) protocol (ref), needed to communicate with MaxMSP. To overcome this limitation, some code modifications and additions have been implemented inside Vicon's Blade sdk.

### 3.2.1 Character design

Implementation of the concept presented above requires Nexus, Vicon SDK and MAX software. Every character involved in the scene is defined by a limited number of markers. In this case, two markers are positioned on the head, one marker is positioned on elbow and the others 2 on the hand (thumb and index finger). Every marker has associated a name in Nexus and between them 6 segments are drawn. It is very important in realtime capture motion, that the marker to have assigned correct coordinates.

- **Vicon extensions (SDK plugin)**

Vicon's SDK is a versatile and simple tool for users to gain easy access to Vicon DataFlow created in Nexus, Blade or Tracker applications. The Vicon DataStream Software Development Kit (SDK) provides intuitive programable access to data with custom functions created in C++. With the help of some functions, Vicon's SDK forwards the Vicon DataStream to other constructive softwares and plug-ins to create custom applications<sup>2</sup>. In combination with Open Sound Control protocol, Vicon's SDK forwards data to any software compatible with this communication protocol (eg. MAX). OSC is a protocol for communication among computers, sound synthesizers and other multimedia devices<sup>3</sup>.

<sup>2</sup>See <https://www.vicon.com/products/software/datastream-sdk/>.

<sup>3</sup>See <http://opensoundcontrol.org>.

Hence, any marker can be routed in MAX using its parameters and coordinates. Also, the Vicon DataStream Software Development Kit (SDK) admits inside changes such as labeling markers, timecode generation and framerate.

### 3.2.2 Objects generation & performance mechanics

Manipulating objects algorithm consists of some big steps: object generation, finding the object, picking up the object, throwing the object on the floor.

Object generation is realized by random generators with the help of *drunk* object, but with certain limits. These limitations are influenced by the dimensions of the room in which the Vicon system is installed. Finding the object supposes continuous mathematic operations between the coordinates of the object and coordinates of the left hand's marker. This process comes with an audio feedback. When these coordinates are close enough one to another, the object is retrieved and manipulated by performer (eg. define gesture). After all these processes, a simple comparison between the coordinates of the floor and the value of the z axes of the marker is done in order to put down the object. According to this, a performer can handle as many objects as he wants.

### 3.2.3 Gesture recognition

*Mubu* containers provided by Ircam laboratories in MAX software represent a handy tool to record and analyze gesture, captured with Vicon system [5]. Our gesture recognition algorithm is based on Hierarchical Hidden Markov Models (HHMM) implemented in *mubu.hhmm* object of MAX/MSP. HHMMs are a generalization of HMM where each state is considered to be a self-contained probabilistic model [6]. The system is trained by captured data which is essentially a gesture. This process requires a predefined indicator in order to delimitate gestures from all data flow. The algorithm analyzes all input data and generates a probability of similarity between data and saved gestures. In order to control every generated object, there are associated 2 or 3 gestures saved by the performer, but there is a limited time for the gestures to be executed. Predefined gestures offer the possibility to delete the gesture just saved and also indicate the moment the gesture is recorded.

### 3.2.4 Sound design

Although there is a wide range of sonification alternatives available, including triggering of pre-stored sounds, modulating/LFO-ing stored sounds etc. we started experimenting with a few synthesis patches with parameters that could be correlated to XYZ coordinates in a reactive and expressive manner. In order to exploit the increased spatial resolution, marker positions, grouping of markers and other possibilities of the Vicon system, the development of synthesis algorithms in MaxMSP seems to be very flexible. In this way, the whole soundscape can be generated in a continuous, organic manner by correlating markers' positions with synthesis parameters.

The interactive experience can be described as having two main paradigms: object finding and object interaction. For object finding we have been experimenting with two straightforward patches in MaxMSP: the first is based on a sawtooth wave that is sent to 8 delays. These output a random signal with a settable range which continuously alter the phase and frequency of the sawtooth iterations. This chorus-inspired algorithm has three controllable parameters: main frequency, range and speed of frequency variation. The farthest someone is from the virtual object the more detuning and phase shifting occurs on each of the

eight iterations, while approaching it the effect becomes less pronounced to the point where only slight variations of the signal occur. Also, main frequency is associated with movement in the vertical plane. Further modulation can be used to make the sound more complex. The second patch is somewhat based on the same principle of decorrelation: six sine waves with different frequencies are modulated in amplitude by a random object which generates control waves with a variable degree of complexity (more or less noise-like shapes). As in this case, distance between the performer and the virtual object is associated with the level of randomness and decorrelation: longer distance translates to a higher degree of decorrelation and randomness. What is interesting about this patch is the effect of amplitude modulation (AM) (referinta ce e AM) when the carrier frequency (?) goes beyond 20hz and sidebands occur. By using noise-like carriers, complex sonorities occur with a variable harmonic content. In both cases, the performer tries to find the object by listening to these variations. By correlating small and large variations to its position in the 3d field the performer receives meaningful clues about where the object might be as well as an interesting and engaging soundscape.

An additional granular synthesis (ref) patch that was initially created for object interaction also seems to be effective for object finding. The patch reads short grains from a pre-loaded sound and scatters into "clouds" in either a controlled or random manner. Among the parameters that can be mapped are: grain position, grain size, envelope shape, level of scattering, pitch, stereo width. Not only this, but by using the [patrrstorage] and [patrr] objects we can group multiple parameters's state, save them as presets and interpolate between them. By doing this, we can control an undefined number of parameters by linking only one marker/one gesture to the interpolation amount box. This patch also seems to be effective for finding the object, differentiating the two paradigms by the level of control: less control for object finding and more control for object interaction.

Real-scenario testing and simulations using the Vicon cameras is quite limited because two main reasons: we don't always have access to the space and we would always need a performer that could cope with long hours of code debugging, errors, MaxMSP programming and so on. This is the reason why in order to simulate interactions inside the sound design patches, we also created a basic interface in Jitter that receives data from Vicon Blade via OSC and represents each marker in the 3d space. By using this interface it is possible to move a particular marker anywhere along the XYZ planes, by only using the mouse, and get instant auditory feedback. The downside being that experiments done in virtual space do not always translate well to the real space: calibration by value scaling, experimenting with different function shapes etc. is tedious but absolutely necessary.

So far, all sound-design is based on two channels that can either be routed to multiple pairs of speakers or downmixed to mono and diffused on an arbitrary number of speakers, however multichannel sound is taken into consideration for future improvements.

- Visualisation (jitter)

## 4. CASE STUDIES

### 4.1 Interactive Installation

- performance analysis

### 4.2 Performance

## 5. CONCLUSIONS AND FUTURE WORK

- Areas of improvement

- Eye tracking?

Although we have been working on the project for only a few months it is fairly safe to state that the expressive opportunities of the Vicon system are superior to others such as cameras and Microsoft Kinects. Although a lot more challenging (no native OSC support, difficulties in conducting test at any time, arduous task of integrating Blade to the Jitter interface to the Gesture Recognition patch, and then to the Sound Synthesis patch, great deal of calibration and scaling), new other possibilities are being offered by the high temporal and spatial resolution.

Future work on SoundThimble could/will include: multiple performers, eye tracking, generative visuals, spatial sound, more powerful synthesis and sound manipulation algorithms, different rules added to the narrative of the auditory game.

## 6. ACKNOWLEDGMENTS

This section is optional; it is a location for you to acknowledge grants, funding, editing assistance and what have you.

## 7. REFERENCES

- [1] Brian Kane. *Sound Unseen - Acousmatic Sound in Theory and Practice*. Oxford University Press, New York.
- [2] T. Hermann, A. Hunt, and J. G. Neuhoff. *The sonification handbook*. Logos Verlag Berlin, 2011.
- [3] K. N. Jorge Solis. *Musical Robots and Interactive Multimodal Systems*. Springer-Verlag Berlin Heidelberg, Berlin, 2011.
- [4] R. B. F. B. Jules Francoise, Norbert Schnell. Probabilistic models for designing motion and sound relationships. *International Conference on New Interfaces for Musical Expression*, pages 287–292, June 2014.
- [5] D. S. G. P. R. B. Norbert Schnell, Axel Robel. Mubu and friends assembling tools. *International Computer Music Association*, pages 423–426, August 2009.
- [6] N. T. Shai Fine, Yoram Singer. The hierarchical hidden markov model: Analysis and applications. *Machine Learning*.