

# SoundThimble: A High Resolution Real-Time Gesture Sonification Framework

Ben Trovato

Authors hidden for review  
1932 Wallamaloo Lane  
Wallamaloo, New Zealand  
trovato@corporation.com

G.K.M. Tobin

Authors hidden for review  
P.O. Box 1212  
Dublin, Ohio 43017-6221  
webmaster@marysville-ohio.com

## ABSTRACT

We introduce *SoundThimble*, a platform for sonic interaction based on the relationship between human motion and virtual objects in 3D space.

A Vicon motion capture system and custom software are used to track, interpret and sonify the movement and gestures of a performer.

We identify three possible interaction scenarios, centred around object searching, manipulation and arrangement. We illustrate the resulting possibilities of extended perception and expression via two case studies: a participative installation and a dance performance.

The software developed is open source and portable to similar hardware systems, leaving room for further extension of the interaction mechanics.

## Author Keywords

Sonification, motion capture, gesture spotting, interactive installation, synthesis

## ACM Classification

- Applied computing → Sound and music computing
- Computing methodologies → Motion capture
- Human-centered computing → Gestural input
- Human-centered computing → Auditory feedback

## 1. INTRODUCTION

High resolution three-dimensional motion tracking is traditionally used for animation in film and games, as well as for life sciences research and engineering applications [18]. This technology has long been mined by the NIME community, although in many early cases, technical limitations meant that the motion data transmission and the sound generation processes were not simultaneous [3, 10].

The *SoundThimble* project harnesses current motion capture technology and gesture detection algorithms to enable new modes of real-time sound exploration. Our aim is to push beyond the standard paradigms of isolated body motion audification [3, 10] or sound control interfaces [4, 12], towards deeper narrative structures coupled with layered arrangement of music patterns.

Our implementation uses a state-of-the-art Vicon motion capture system<sup>1</sup> containing eight Vantage 5-megapixel infrared cameras and two Bonita video cameras. Since the open-source software developed in this project<sup>2</sup> is built around Vicon's Datastream SDK,<sup>3</sup> the platform can be ported to both older and future Vicon-based systems.

In the remainder of the paper, we review relevant literature and technology (section 2), we describe the *SoundThimble* concept and implementation (section 3), we propose two applied case studies (section 4), and we conclude with a survey of remaining challenges and future perspectives (section 5).

## 2. STATE OF THE ART

Sonification, as the auditory representation of a datastream, is a rich tool for interpreting human movement [7]. From the musical perspective, infrared motion capture systems have been revealed as a technically superior means for expressive interaction in a controlled environment, with many features being translatable to more portable technologies [15, 17].

In particular, Vicon motion capture systems have been used for over a decade for music applications [3, 10, 4, 17]. A software bridge for streaming OSC data from Vicon exists<sup>4</sup> [4] as part of a concluded project, which proved to be incompatible with our current setup.

The current decade has seen qualitative advances in the interaction between human gesture and sound behaviour, made possible by real-time gesture recognition and following tools [1, 2, 6]. These allow for more complex scenarios, where movement is used both for direct sonification and for multi-level control of system behaviour—features which our project channels into a coherent framework.

## 3. PROJECT DESCRIPTION

### 3.1 Concept

The “sound-thimble”, as the basic building block of our framework, is based on the concept of *sound object* in the Schaefferian sense, as a clearly delimited sounding unit, open to manipulation, arrangement and composition [13].

Such an entity, once instantiated, can retain an ambiguous nature (spatially and acoustically) or can switch to a more material state (positioned in space and tied to a causal source) [9]. The duality between the latent positioning of the object (which can be inferred from phenomena other than spatial sound reproduction), and the active sound spatialisation and transformation, becomes an innovative tool

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

<sup>2</sup>Available at <https://github.com/RVirmoors/viconOSC>.

<sup>3</sup>See <https://www.vicon.com/products/software/datastream-sdk>.

<sup>4</sup>See <http://sonenvir.at/downloads/qvicon2osc/>.



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.

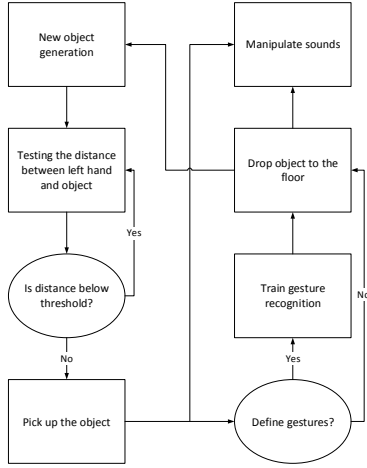


Figure 1: *SoundThimble* interaction workflow.

for the sonic arts through sound sketching, auditory games and other real-time interaction setups.

### 3.1.1 Interaction scenario

We present the initial application of our framework, in the form of an auditory game comprising three phases: search, manipulation, arrangement.

The game’s narrative starts with a human player attempting to find a sound-thimble (a stationary virtual object, randomly positioned in 3D space), by analysing cues that are constantly shifting in the sonic fabric based on the hand’s movement relative to the object. Analogously to the traditional game of *Hunt the Thimble* (a.k.a *Hot or Cold*), the space between the human and the virtual object is correlated to sound synthesis and modulation parameters. Briefly, the closer one comes to the object, the more coherent the sound and vice-versa.

Once the object is found and attached to the hand, its sonic manifestation gains a richer causal relationship: the player becomes a performer, and is now able to explore the object’s sonic palette, and record a number of gestures that can be re-performed later, re-called, and used to trigger or manipulate sonic shifts and events.

Finally, the performer can drop the virtual object to the floor, or discard it by “pushing” it outside of the installation boundaries. This triggers a new object to be randomly generated, while the player retains a degree of control over the initial object via the recorded gestures. Both objects are now in a latent state, with the new one guiding the player’s search, and the previous one responding to the learned set of gestures.

This repeating scenario is outlined in Figure 1: objects are randomly generated, the performer finds them, defines gestures and interacts sonically, before arranging them in a pleasing configuration. With each spawning of an object or assignment of a gesture, the game becomes more challenging and complex, but also more flexible and rewarding.

### 3.1.2 Performance aesthetic

The human-object dynamic at the core of our framework results in certain interaction features which circumscribe the aesthetics of any application of the platform.

Our approach is informed by Worrall’s study [19], which reveals a necessity for the mapping of minute gestural inflections to alter sonic material with a view to certain modes of

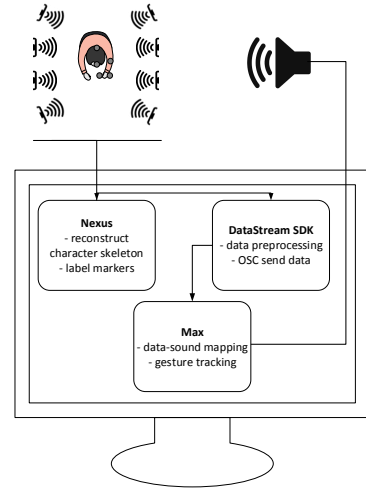


Figure 2: *SoundThimble* framework architecture.

listening. The aim of *SoundThimble* is to fluctuate between: reflexive, kinaesthetic, connotative, empathetic, reduced.

The cross-modality between different kinds of sense perception guides the performer’s attention to the various sonic responses to physical actions. The multi-modal information is processed in real time, continuously redefining the affordances enabled by the system.

Moreover, since the system’s responsiveness is reliant on marker visibility, the performer becomes, to a degree, existentially dependent on the camera eye. This, coupled with the coexistence of virtual, responsive objects in the same scene, can lead to novel mechanics at the limits between presence and absence, real and virtual.

## 3.2 Implementation

The framework architecture diagram is laid out in Figure 2. Three-dimensional sensor data is streamed into the Vicon Nexus software, which is able to reconstruct and label the underlying character skeleton. The gesture recognition and sonification algorithms are programmed in Max<sup>5</sup>, which generally receives control data via the OSC<sup>6</sup> protocol. Since Vicon systems do not support OSC out of the box, we used the *oscpack*<sup>7</sup> library to extend the DataStream C++ SDK and send OSC bundles to Max.

The following description is tailored to our installation application, but any *SoundThimble*-based project involves similar conditions.

### 3.2.1 Character design

Figure 3 shows a skeletal reconstruction in the Nexus environment. We pursued a minimal amount of markers, for ease of setup and prototyping. The resulting configuration—sufficient for tracking hand gestures, while ensuring redundancy for when one marker is obscured from the cameras—consists of 5 markers: two positioned on the head, one on the forearm, and two on the hand (thumb and index finger).

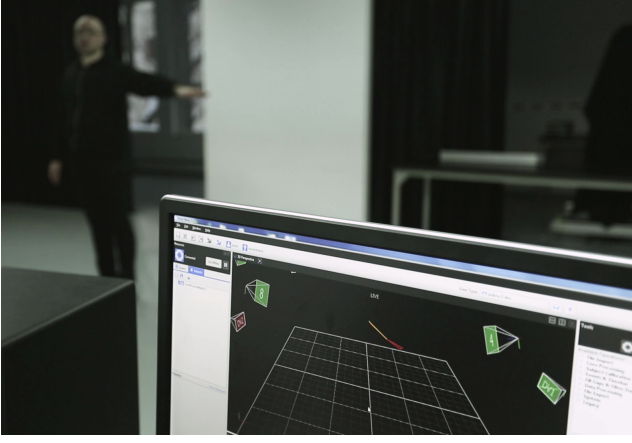
Each OSC bundle sent through the SDK consists of the following:

- 3D coordinates for the head (averaged from the two head markers);

<sup>5</sup>Max is a state-of-the-art programming environment for real-time multimedia: <http://cycling74.com/>.

<sup>6</sup>OpenSoundControl is a multimedia communication protocol: <http://opensoundcontrol.org/>.

<sup>7</sup>See <http://www.rossbencina.com/code/oscpack>.



**Figure 3: A performer tracked in Vicon Nexus. Two visible segments: head-forearm, forearm-hand.**

- 3D coordinates for the hand (averaged from the two hand markers);
- distance between thumb and index finger.

Each of the three items is sent only if non-zero, i.e. for the head and hand coordinates at least one of the respective markers is active, while for the distance computation, both markers need to be visible and correctly labelled. The forearm marker only serves for skeleton reconstruction, and is not sent via OSC.

To maximise responsiveness and minimise data loss, we raised the frame rate to 500Hz, taking into account the maximum movement speed and the minimum spacing between markers [16]. This configuration produces highly stable and responsive inputs into the Max system, with a spatial resolution of 1mm and a latency of around 5ms.

### 3.2.2 Object generation & interaction mechanics

The following object-related mechanics are implemented as basic algorithms in Max: generation, detection, dropping.

Object generation is executed randomly within the boundaries of the motion capture field, as defined in the Vicon calibration phase. Detection of the sound-thimble occurs when the distance between hand and object falls below a set threshold. By default this threshold is set at 200mm radial distance; lowering it can make the game considerably more difficult. Once the object is detected, it becomes mobile, its coordinates tracking those of the hand's.

Finally, a simple thresholding of the  $z$ -axis (height) value of the hand position serves to put down the object. If the velocity computed on the  $x$  or  $y$  axes (horizontal plane) is high enough, then the object is pushed in the respective direction, and is able to leave the area of the installation, essentially being removed from the game. At this point, a new object is introduced. The system keeps track of all object coordinates, as they appear and disappear over time.

### 3.2.3 Gesture recognition

We use the thumb-index finger distance value to enable gesture recording while the two fingers are kept close together. The input features captured into *MuBu* multi-buffer containers [14] consist of pairs of values:

- $\sqrt{\frac{(\Delta x)^2 + (\Delta y)^2}{2}}$ ;
- $\Delta z$ ,

where  $\Delta x$ ,  $\Delta y$ ,  $\Delta z$  are the respective differences between head and hand coordinates. This feature preprocessing serves two purposes:

Firstly, gestures are recorded based on the position of the hand relative to the head, thus becoming invariable to the performer's absolute *position* within the space. Secondly, by composing the  $x$  and  $y$  values into a single feature, gestures become invariable to the performer's *orientation* on the horizontal plane. Thus, gestures can be recorded and recalled anywhere within the space, irrespective of the direction the performer is facing.

The input features are fed to one of two gesture recognition algorithms, both part of the *MuBu* package. The first, based on Hierarchical Hidden Markov Models (HHMM), is implemented in the *mubu.hhmm* Max object. HHMMs are a generalization of HMM where each state is considered to be a self-contained probabilistic model [6, 5]. The second is the *Gesture Follower gf* Max object, based on a Sequential Monte Carlo inference engine [2].

The first method allows for *gesture spotting*, i.e. it constantly produces likelihood values of a certain gesture being active, together with an approximation of its completion rate. If these exceed a certain threshold, the respective gesture is triggered. The second method is more flexible and precise: the detected gesture can be followed at a variable rate or scale, even backwards. The only drawback is that it requires a start trigger, which we send by quickly attaching and separating the thumb and index finger.

In the game, each generated object has a number of gestures associated to it. When an object is dropped to the floor, the classifier is (re)trained with the new data, and consequently gestures can act as on/off switches for a particular sonic behaviour (if they are performed/spotted once at a time), or as continuous controllers (if they are repeated). When several objects exist on the floor, a specific movement might act on one or more objects, depending on which detection likelihoods exceed the threshold.

### 3.2.4 Sound design

Each sound-thimble has a corresponding sound design patch, differing in (a) the source sound material used, (b) the synthesis techniques applied, and/or (c) the control mapping schema to the object search and manipulation variables. The various combinations of (a), (b) and (c) give rise to a growing library of objects, each with its own character. By designing various interaction rules for each object, segments are linked to different synthesis parameters or groups of parameters resulting in a continuously evolving, organic soundscape.

We differentiate between the three phases of the installation in terms of mapping technique and level of sonic interactivity. The search mode employs straightforward parameter mapping, where human-object distance measures are linked to synthesis parameters, while the manipulation phase relies on a model-based mapping approach where different gestures and actions reach deeper levels of control [7]. Finally, variations on both these techniques occur in the arrangement phase.

The synthesis patches used for search are built around a process of decorrelation [11]: the farther the human's hand from the object the more decorrelation occurs, up to the point where each instance of the signal becomes a distinct sonic entity. This is done by continuously modulating each of the copies in terms of pitch (FM) and amplitude (AM) with low frequency oscillators (LFO's). Changes in frequency, amplitude and wave shape of the LFO's lead to complex sonorities ranging from coupled streams of sound to distinct iterations with a high degree of randomness. This

mechanism is subtly mixed with a granular engine in a latent state, which gains more prominence in the next mode.

In the manipulation phase, the 3D space is split into chunks, each one acting as a zone with its own mappings and interaction laws. The synthesis patches are based on the segmentation of a source sound into short grains. We built a granular synth with these controllable parameters: grain size, grain position, envelope shape, level of scattering, pitch, timbre and stereo width. Another patch implements a concatenative synth that traverses the grains guided by movement velocity and trajectory. Certain gestures trigger sonic events while interpolating between sets of parameter values via a convergent-mapping schema [8].

In the search phase, 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. In the manipulation phase, the soundscape is spatialised to track the position of the performer. In the arrangement phase, the dropped object retains a “root” source location, which can relate dynamically to the performer position when a gesture is executed: for instance, the granular patch sends spatialised grains back and forth between the dropped object and the performer.

In developing the sound design algorithms, we used two input data sources. The first one is motion capture data recorded in Nexus and played back through the SDK. For more flexibility and immediate control, we also developed a basic visual interface to monitor the input data and to manipulate it in virtual 3D space using the mouse, for instant auditory feedback. This feature comes in useful when specific motion data is not available and needs to be roughly simulated.

## 4. CASE STUDIES

We present two applications of our platform: a participative installation, and a performance piece. They represent two manifestations of the *SoundThimble* concept and infrastructure, revealing a pair of experiences: (a) the direct interaction with sound in space, and (b) the mediated interaction with sound and space of an audience member. In both cases, the range of sense and perception is extended through new experiences.

### 4.1 Interactive installation

The first form conceived for the *SoundThimble* platform is that of an open installation, where visitors become participants and observers in the situation outlined in section 3.1.1.

Beyond the considerations in section 3.1.2, the major feature of the auditory game is the layered sequence of search-manipulation-arrangement which translates to a layering of awareness of accumulated experiences. The meshing of sound objects (supplied by the game) and gestures (defined by the player) leads from a state of uncertainty and potential, to phases of exploration, play, composition. Snapshots of interest can be stored for later recall—currently this is done manually, but might be automated in the future.

### 4.2 Dance performance

The following is a sketch for a work currently under development. A performer enters the motion capture space, and commences a series of improvised motions, silently establishing the action space.

Once a threshold is reached, a virtual object is generated and the movement is no longer completely free, its range of action being directed by the relationship to the object position. Searching and finding the sound-thimble is a process

of correlating the sonic characteristics to one’s own movement patterns.

By “resonating” with the object’s manifestation and locating it, the performer enters the manipulation phase, extending her auditive and tactile perception through embodied listening. The multimodal information processed by the performer is both cochlear and kinaesthetic/proprioceptive, tactile, vestibular and visual. Thus, the information she forms about the shape of the movement, its trajectory and spatial dynamics, is intrinsically linked to the sonic connections to these parameters. This dynamic knowledge influences the timing of new actions and the anticipation of sonic feedback. The resulting action schemas are composed as both spatial and sonic shapes.

Gradually, as the performer’s control patterns become crystallised, the sonic nature of the source object moves between the abstract and the concrete. In the latter phase, the performer is able to control higher-level parameters such as tempo and orchestration of a coherent musical structure which was partly a result of the generative interaction in the more abstract/exploratory phase.

## 5. CONCLUSIONS AND FUTURE WORK

This paper introduced *SoundThimble*, a multi-layered platform for real-time motion-music interaction. All software developed for the project (including the C++ code for data preprocessing and transmission, and the Max patches for gesture tracking and sound design) is open source and publicly available.

The team is currently pursuing several directions for improvement and extension. We are working to support more than one participant at once, implementing mechanics for sharing control of the virtual object. Video projection support is planned, e.g. each object having a corresponding reactive video, with the arrangement phase also producing a background visual collage. We are also considering eye tracking technology to enable an audience member’s active involvement in shaping the sound. Meanwhile, work continues on further interaction mechanics and sound design.

Finally, we are commencing the outreach to composers, artists and creative programmers, to apply our platform to new innovative projects and engage in practice-led research.

## 6. ACKNOWLEDGMENTS

Hidden for review.

## 7. REFERENCES

- [1] B. Caramiaux, J. Françoise, N. Schnell, and F. Bevilacqua. Mapping through listening. *Computer Music Journal*, 38(3):34–48, 2014.
- [2] B. Caramiaux, N. Montecchio, A. Tanaka, and F. Bevilacqua. Adaptive gesture recognition with variation estimation for interactive systems. *ACM Transactions on Interactive Intelligent Systems (TiiS)*, 4(4):18, 2015.
- [3] C. Dobrian and F. Bevilacqua. Gestural control of music: using the vicon 8 motion capture system. In *Proceedings of the international conference on New interfaces for musical expression*, pages 161–163. National University of Singapore, 2003.
- [4] G. Eckel, D. Pirro, and G. K. Sharma. Motion-enabled live electronics. In *Proceedings of the 6th Sound and Music Computing Conference, Porto, Portugal*, 2009.
- [5] S. Fine, Y. Singer, and N. Tishby. The hierarchical hidden markov model: Analysis and applications.

- Machine Learning*, pages 41–62, 1998.
- [6] J. Françoise, N. Schnell, R. Borghesi, and F. Bevilacqua. Probabilistic models for designing motion and sound relationships. In *Proceedings of the 2014 international conference on new interfaces for musical expression*, pages 287–292, 2014.
  - [7] T. Hermann, A. Hunt, and J. G. Neuhoff. *The sonification handbook*. Logos Verlag Berlin, 2011.
  - [8] A. Hunt and R. Kirk. Mapping strategies for musical performance. *Trends in Gestural Control of Music*, 21:231–258, 2000.
  - [9] B. Kane. *Sound Unseen - Acousmatic Sound in Theory and Practice*. Oxford University Press, New York.
  - [10] A. Kapur, G. Tzanetakis, N. Virji-Babul, G. Wang, and P. R. Cook. A framework for sonification of vicon motion capture data. In *Conference on Digital Audio Effects*, pages 47–52, 2005.
  - [11] G. S. Kendall. The decorrelation of audio signals and its impact on spatial imagery. *Computer Music Journal*, 19(4):71–87, 1995.
  - [12] K. Nymoen, S. A. v. D. Skogstad, and A. R. Jensenius. Soundsaber-a motion capture instrument. In *Proceedings of the international conference on New interfaces for musical expression*, 2011.
  - [13] P. Schaeffer, G. Reibel, B. Ferreyra, H. Chiarucci, F. Bayle, A. Tanguy, J.-L. Ducarme, J.-F. Pontefract, and J. Schwarz. *Solfège de l’objet sonore*. INA GRM, 1998.
  - [14] N. Schnell, A. Röbel, D. Schwarz, G. Peeters, R. Borghesi, et al. Mubu and friends-assembling tools for content based real-time interactive audio processing in max/msp. In *International Computer Music Conference*, 2009.
  - [15] S. A. v. D. Skogstad, A. R. Jensenius, and K. Nymoen. Using ir optical marker based motion capture for exploring musical interaction. 2010.
  - [16] M.-H. Song and R. I. Godøy. How fast is your body motion? determining a sufficient frame rate for an optical motion tracking system using passive markers. *PloS one*, 11(3):e0150993, 2016.
  - [17] G. Vigliensoni and M. M. Wanderley. A quantitative comparison of position trackers for the development of a touch-less musical interface. In *NIME*, 2012.
  - [18] G. Welch and E. Foxlin. Motion tracking: No silver bullet, but a respectable arsenal. *IEEE Computer graphics and Applications*, 22(6):24–38, 2002.
  - [19] D. Worrall. Understanding the need for micro-gestural inflections in parameter-mapping sonification. In *International Conference on Auditory Display*, 2013.