

TH KÖLN

PROJECT REPORT

# Realizing spatial listening tests in Virtual Reality using Unity

*Alexander Müller*

Supervised by  
M.Sc. Melissa Andrea Ramirez Caro  
&  
Prof. Dr. Christoph Pörschmann

November 12, 2021

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Motivation</b>	<b>1</b>
<b>3</b>	<b>Fundamentals</b>	<b>2</b>
3.1	Auditive localization . . . . .	2
3.2	Auralization . . . . .	2
3.3	Hearing perception (energy and similarity) . . . . .	3
3.4	Training abilities . . . . .	3
<b>4</b>	<b>Approach</b>	<b>3</b>
4.1	Unity overview . . . . .	3
4.2	Orientation on LiSN project . . . . .	4
4.3	Audio assets . . . . .	4
4.4	Visuals . . . . .	4
4.5	User interface . . . . .	5
4.5.1	Audio Sources . . . . .	5
4.6	Training-Game . . . . .	5
4.6.1	Audio files . . . . .	5
4.6.2	Game loop . . . . .	5
<b>5</b>	<b>Conclusion</b>	<b>5</b>
<b>6</b>	<b>Outlook</b>	<b>6</b>
	<b>List of Figures</b>	<b>8</b>
	<b>Sources</b>	<b>8</b>

# 1 Introduction

Spatial hearing describes the ability to use binaural clues to localize the source of a noise. This ability relies upon three main features: the difference in level (1) and phase (2) of a signal between the left and the right ear and the spectral coloring of a known noise, caused by the geometry of the listener's head and upper body.

Apart from the obvious advantages of spatial hearing, like quickly identifying incoming dangers (e.g. a car), it provides an additional aid to the way we perceive sound. In noisy environments, binaural clues can be used to focus on a particular sound source. This situation is most commonly described in the "cocktail-party-effect". A listener can use the information of position of a given audio source to reduce noises from other directions within their perception.

Since noisy environments with a lot of similar audio signals (in particular speech) are very common within a lot of every day occasions, the advantage of using binaural clues has a great importance. In most cases the lack of this ability is found on people with hearing disabilities. Even if hearing aids have advanced throughout the last decades, they still can't offer their users the same subtle binaural clues as a person with intact hearing abilities.

But apart from physiological issues causing problems with spatial hearing it has been found, that children with perfect hearing still might not be able to correctly interpret binaural clues. This effect is described as Spatial Processing Disorder (SPD). Since the actual hearing abilities of the affected children is not impaired, they are - in contrast to persons with typical hearing disabilities - provided with all the information required to localize a given sound source. Based on this knowledge the idea was formed, to treat SPD by *training* affected children and therefore teaching them, how to correctly interpret binaural clues.

Within the *Development and Evaluation of the LiSN & Learn Auditory Training Software for Deficit-Specific Remediation of Binaural Processing Deficits in Children: Preliminary Findings* [2] paper by Sharon Cameron and Harvey Dillon it is described how such a training process could look like and already gave some hints <sup>1</sup>.

Building upon this approach the training concept described by Cameron and Dillon shall be transferred into a virtual reality application. This shall bring the additional advantage of combining auditory and visual clues and also allows to further extend the process.

# 2 Motivation

Realizing this listening test as a Unity-based VR application offers several advantages. The first one being the open nature of this project. Since only free-to-use assets have

---

<sup>1</sup> The original experiments have only been done with a small sample size. Within the conclusion of the paper the requirement of a clinical trial is mentioned, to validate the efficacy of the training process.

been included, both the compiled application and the source files can be made publicly available. Combined with simple setup required to perform the listening experiments, this will hopefully allow a variety of interested groups to perform these tests to collect further data and offer affected children access to work on the condition.

Furthermore the Unity framework in conjunction with the prefabs and scripts, that have been created for this project, it would be relatively easy to extend upon the original training concept. For example the addition of multiple distracters within a given scene could very easily be realized.

Apart from the previously mentioned aspects the usage VR peripherals for this project also offers a lot more possibilities. Some of them will be discussed in 6. But especially adaptive 3D audio in combination with head tracking would allow for a much more realistic scenario. Also the novelty of the VR headset itself will most likely already spark a lot of interest in the participants in comparison to a regular learning/training game.

### 3 Fundamentals

Before discussing the approach towards this project, some basic principles of modern day audio processing shall be explained.

#### 3.1 Auditive localization

The human ability of locating audio sources is based upon two main principles. The first being the perception of the time and level differences between an acoustic event on the left and right ear. The different time points at which the signal is perceived at the two ears, translates to a different phase at which the audio wave is registered. On the other side the level differences mainly derive from the shadowing caused by the head. Both of these effects are not frequency independent, for lower tones the localization via phase differences works better, while at higher frequencies the level variations offer a better indication.

The second way of locating a sound source is based upon the spectral differences that occur through the different paths a sound wave can take around the listeners head and upper body. This effect is strongly influenced by the shape of the outer ear, but also the general geometry of the head makes a difference. Basically the head and the ear can be described as directional filters. This kind of localization is mostly based on experiences. The listener can learn to associate the spectral differences in a known noise to the location from which the noise is originated (e.g. by looking for the sound source).

#### 3.2 Auralization

Auralization describes the process of recreating not only a given audio signal but also the spatial information. With correct auralization a recorded scene (e.g. a classical concert in a concert hall) can be rec

This requires advanced recording processes to capture the spatial information of a scene (e.g. a concert hall)

Auralization describes the process of adding spatial features to audio signals. The simplest example for this would be the move from mono to stereo audio. However, since these initial steps the options of simulating and re-creating spatial features in digital audio processing have advanced quite a lot. This can be easily seen in many consumer electronic products such as AV-Receivers, Soundbars or Headphones supporting 5.1, 7.1 etc. or in particular VR peripherals coming with their own 3D audio frameworks for developers.

The basic principle remains more or less the same for all these products. Within a given *listening space* the original auditory environment shall be recreated.

### 3.3 Hearing perception (energy and similarity)

### 3.4 Training abilities

## 4 Approach

Before evaluating possible improvements of virtual reality support in listening tests, it has to be made sure that there are no major other influences alternating the results of this test. Therefore the first iteration of the new test environment will be developed as close as possible to the original LiSN software. This allows to compare test and training results from this test environment with the data already available from previous experiments. If the new software allows to recreate similar test results, it can be assumed that the unenivitable difference don't affect the results in a too drastic manner. Afterwards the data from these initial tests can also be used to validate whether VR exclusive features like head tracking add further benefitis to the training process.

**Development tools** As already mentioned, the major part of the development will be done in Unity with the addition of Visual Studio as IDE. The other major tool is the Oculus VR SDK, which will be used for audio spatialization and of course rendering the visuals to be properly displayed on the given VR peripherals.

### 4.1 Unity overview

Unity is a development framework, which is mainly used as a engine for video games. However due to the requirements to realism in modern video games, the amount of audio processing included, is already pretty elaborate. In conjunction with the large user base and support by hardware manufacturer (e.g. Oculus), Unity offers a great starting point for this project.

## 4.2 Orientation on LiSN project

### 4.3 Audio assets

As basis for this project the word lists given in the appendix of the LiSN paper [?] are used. Since due to the huge variety of options, the only feasible way to construct the sentences is by assembling a sentence through separate audio files for each word. This however comes with a problem. Simply recording single words and then playing the files one after another strongly alter the speech flow and accentuation, which would normally be present when the sentence is spoken.

To handle this problem, the words won't be recorded separately but instead a subset of the possible sentences from each list will be recorded. Within this subset all words from the list are included within the recordings. Afterwards the recorded sentences will be cut into the individual words and then will be used to assemble randomly generated sequences based on these assets. This should work fine, since the structure of all sentences within any list remains the same. So the speech flow and pronunciation between the sentences should not be too larger. GET SOME REFERENCES FOR THIS ASSUMPTION!

### 4.4 Visuals

Since this project shall be used to perform research with children as participants the graphical presentation is of great importance. Especially since Virtual Reality (as of today) is in general a quite unfamiliar setting, it is very important to give the user a sense of space within the setting.

**Matching audio parameters to visuals** The audio scenario shall be based upon free-field sound propagation. To avoid confusion of the listener the environment presented within the VR application should represent this setting. So a closed room would not be a good option, since we inherently expect reflections and reverberation within such a setting. A better approach would be an open field, where the lack of reflexions feels more natural.

**Avatars** Another important part is the graphical representation of the audio sources. Since we have the visual component given through the VR headset, the sound should not simply come from an invisible sources, but have a origin which the user can identify through the graphics. This however includes another challenge. Of course it would be possible to create humanoid avatars with complex animations - including lip syncing - to convey that the object is indeed active and not just a passive talking rock. However this would take a huge amount of effort to pull off. Instead a different path will be used here, which is also pretty common in budget oriented game design (e.g. indie games). Through abstract avatars and simple animations/movements it is possible to convince the user, that the object is alive and active, while at the same time requiring a lot less effort to pull off. SOURCE...REFERENCES

## 4.5 User interface

### 4.5.1 Audio Sources

Within the level both interactive and regular audio sources are implemented. All of them use 3D rendering, so that the perceived sound changes based on the position of the player to the object, emitting the sound. The advantage of adding interactive sound sources is that,...

## 4.6 Training-Game

### 4.6.1 Audio files

Target sentences are constructed in a way that the co-articulation in between the sentences of a list is similar enough the allow randomly scrambling up sentences will maintaining a natural sounding result.

### 4.6.2 Game loop

Continuous distracter stories creating a noise with similar characteristics as the target sentences.

Both the target sentences and the distracter stories shall be omitted by actual objects within the game world. This association between audio and visual representation shall create further immersion (FIND REFERENCE!).

Whenever a target sentence was played, four word options will be displayed. (IMAGES OR WORDS?). When the correct result is picked, the selected option turns green and a celebratory sound is played. Otherwise if an incorrect option is selected, it is marked red and a failure sound is played. Alternatively an "uncertain" option is presented.

SNR: For every correct guess, the SNR is decreased by 2.5 dB. For every incorrect guess, SNR is increased by 1.5 dB. At uncertain at first the same sentence is repeated with 1.5 dB higher SNR. At a second consecutive *uncertain* a new sentence will be played, with again increased SNR.

## 5 Conclusion

Within the scope of this project a simple binaural recreation plug-in has been implemented using the Juce framework and the FFTW library. Unfortunately the focus of this project has been moved heavily onto software development and typical coding issues, instead of digging deeper into the actual subject of audio auralization. But nonetheless the current state of the plug-in might be a good starting point for further projects, at cutting some of the very time consuming setup procedures. Of course this would require to encapsulate the currently implemented processing further into individual functions, or possibly even better an additional class. With a bit added abstraction, a generic processing chain could be implemented (similar to the one already implemented by Juce), which would allow to

add and move individual auralization processes much more easily. Another working point would be to amp up the current convolution, to also feature e.g. elevation. This not too much would have to be changed in the elemental parts of the source code, this should be a rather simple addition.

## **6 Outlook**



## Abbreviations

<b>API</b>	Application Programming Interface
<b>DAW</b>	Digital Audio Workstation
<b>DSP</b>	Digital Signal Processing
<b>GUI</b>	Graphical User Interface
<b>HRTF</b>	Head Related Transfer Function
<b>IDE</b>	Integrated Development Environment
<b>LISN</b>	Listening in Spatialized Noise
<b>SDK</b>	Software Development Kit
<b>SNR</b>	Signal to Noise Ratio
<b>SPD</b>	Spatial Processing Disorder
<b>UI</b>	User Interface

## List of Figures

## Sources

- [1] Oculus Spatializer: <https://developer.oculus.com/documentation/unity/audio-osp-unity-spatialize/>
- [2] Development and Evaluation of the LiSN & Learn Auditory Training Software for Deficit-Specific Remediation of Binaural Processing Deficits in Children: Preliminary Findings, Sharon Camerion, Harvey Dillon, Date: 2011
- [3] Correlating performance on the Listening in Spatialized Noise Sentences Test (LiSN-S) with the Listening in Spatialized Noise Universal Test (LiSN-U), Kiri Mealingsa, Sharon Camerona and Harvey Dillon, Published: April 2020, Date: 2019