

An efficient 3-D sound experience for mobile applications

Group 3 Aalborg University
Copenhagen, Denmark

Abstract—We developed an application for mobile devices that uses an efficient HRTF model based on to create 3D soundscapes. The HRTF model is based on a combination of filters and delays provided on theoretical basis by Brown et. al. [1]. The device orientation is used to set the angle between the sound and the user. This model has many possibilities on the limited hardware capabilities of the mobile devices in comparison to HRTF database based models. The prototype will present the user several sounds scattered in a limited area where the user can interact with

Index Terms—HRTF, mobile devices, soundscape, Pure Data, OpenFrameworks

I. INTRODUCTION

A lot of research on modeling 3-D sound has been done with fairly good solutions (references). The best results can be obtained by those models implementing 3-D sound using personalized HRTF functions which needs tedious measuring of impulse responses etc (reference). Since mobile applications as well as multimedia productions are usually aimed at a big number of users, implementing personalized HRTF functions would require measurements of the individual user which is almost impossible. Additionally, using huge databases with HRTF's for subjects having similar head diameters and pinnae structures (dunno, how those HRTF's work) are not feasible for mobile applications as memory space is limited. Therefore, the aim of this project was to develop a computationally efficient and general HRTF model, yet, keeping the quality of the 3-D sound as good as possible. We therefore used the model provided on a theoretical basis in [?] by implementing a combination of filters and delays in Pure Data and C. This audio engine was then embedded in a mobile application. Sensor data provided by the mobile device was used to compute from which direction the sound should appear to be coming from. With the application users should be able to discover a virtual sound space situated in Copenhagen. This means that our HRTF model should satisfy two major requirements.

1. It should provide direction cues that are good enough to guide the user to the position where a specific sound is situated.
 2. The implemented HRTF model should provide an intuitive way of locating a sound source from an qualitative aspect.
- The first aim is bound to accuracy of the GPS signal and orientation data provided by the individual mobile phone. We hypothesized that using a general HRTF model should be of sufficient quality to satisfy the first requirement since sensory data can be rather noisy anyways which would be overfitted

by a very precise HRTF model.

===Something about the second hypothesis===

Taken both hypotheses together the user should experience a soundscape in at least 2-dimensions, azimuth and distance.

The purpose of this paper is...*(it should be short, specific and immediate)*

First, some state of the art research on audio augmented reality will be provided. That section will be followed by the presentation of our implementation of the audio engine as well as the main application. In the last part we will present the results obtained by testing our application and discuss those in light of the two hypotheses stated above.

A. State of the art

The concept of audio augmented reality (AAR) deals with techniques where a real sound environment is extended with virtual auditory environments. Nowadays, the progress in audio technology and computing predicts the introduction of completely new type of interactive audio applications []. Advances in mobile technologies have made it possible to create audio augmented spaces almost anywhere. For instance, spatial auditory displays that can provide the user with landmarks and are capable to attract the user's attention have been tested and introduced []. Many experiments, qualitative and quantitative researches have been designed so far to better understand the way in which people usually perceive multiple simultaneous sources differently placed and to increase the level of immersion in the experience.

1) *Audio reality vs Virtual reality*: First of all, the possibility to hear the natural acoustic environment around a user differentiates the concept of audio augmented reality from the traditional concept of a virtual reality audio environment []. In virtual reality, generally participants are abstracted from the natural environment and are surrounded only by a completely synthetic one (acoustic and/or visual). On the opposite, in augmented reality a virtual environment is *superimposed* on a real one. To be more specific, in a mobile audio augmented environment participants are able to interact with the virtual audio mixed with real vision and soundscape.

2) *Augmented audio reality*: The aim of AAR is to combine both real and virtual sound scenes so that virtual sounds are perceived as an extension to the natural ones. Virtual sound, as expressed in [], could be originated from another environment or artificially created. Two extreme possible augmentations are possible. The former is a perfect augmentation of the listener's auditory environment and is achieved when the listener is

unable to predict whether a sound source is part of the real or the virtual audio environment. The latter is a high-quality virtual auditory scene with some characteristic that are not possible in acoustic environments.

3) *Previous work*: Several work has been done in this field although this area is relatively new. *Audio Aura* (Mynatt 1995) was one of the first project to deal exclusively with audio in augmented reality system. It basically consisted in providing information to users as they travelled through their workspace. These information was triggered by particular locations in the workspace. The use of audio in *Audio Aura* is particularly interesting because most of its cues were associated to sounds from nature rather than recorded vocal, speech or synthetic sounds. Similar in approach to *Audio Aura* was the *Automated Tour Guide* (Bederson 1995). In both cases, triggers are readily identifiable, still and rarely changing. Those augmented sounds were associated to pre-determined locations in space and there was no need to determine the precise location of an individual. A successive work, *Hear&There*, was able to determine the location and head position of the user using the information from GPS and a digital compass []. A user could listen to these ‘audio imprints’ by walking into the area that a specific imprint occupied, which was triggered by proximity. The essential premise of *Hear&There* was that a physical environment has been augmented with audio. All the sounds and the data were gathered inside that system. Since a ‘Field Use’ has been developed, in which the user wear a hardware portable system, *Hear&There* has undoubtedly contributed to an improved definition of mobile augmented reality environment.

4) *Human sound localization*: The auditory system is used by the humans for several purposes in the daily life. One of this purposes is to provide the necessary information to localize sound sources in various dimensions, (width, height, depth) and it is even possible to guess the size of the source.

When a sound event occurs, the waves travel in all directions and, when they reach us, our brain compares the signals received by the left and right ears to localize it in the horizontal plane. The spectrum of the signal reaching each ear is different, since the amplitude and phase information differs. These binaural cues are called *interaural intensity difference* (IID) and *interaural time difference* (ITD). However, these cues are not enough to localize accurately the source since with this information the listener can not determine if the sound is in front, above or behind. This region of positions where all sounds yield the same ITD and IID is called *cone of confusion*. This ambiguity can be solved with the information provided by the filter effect caused by the pinnae, head, shoulders and torso, which modify the spectrum of the sound that reach the listener’s ears. The sum of all these features are characterized by the Head Related Transfer Functions (HTRF) which are not only frequency and direction-dependent but also differ from person to person. That dependency makes therefore hard to generalize the spectral features among individuals. It is well known that using a HTRF from one person in another can significantly impair the perception due to the individual differences in the anatomy, but it has also been shown that some people localize better the sounds than others, and their

HTRFs are suitable for a large group of listeners.

5) *Localization/Lateralization*: Headphones are much related to the concept of mobile and wearable application and have been used successfully in many virtual reality applications. Headphones, and earphones as well, make the ear receive the sound separately. That is advantageous because the whole scale of differences between the signals coming to ears, i.e. *interaural time differences* (ITD) and *interaural intensity differences* (IID), can be manipulated separately and individually [].

Generally speaking, when listening to sounds coming from headphones, the sound sources seem to be located inside the subject’s head []. This is usually called *lateralization*, or intracranial, or ‘inside-head-localization’ (IHL). On the opposite, there is the effect of having the sound outside the head, according to specific direction and distance, i.e. *localization* or ‘outside-head-localization’ (OHL). It has been necessary distinguish localization inside and outside the head. This difference in terminology has reinforced the assumption that there is a basic difference between perception of sound conveyed by headphones and sound that reaches the ear from an external and real source []. It has also demonstrated that a listener can make a clear distinction in headphones listening between localized (that is, sound inside the head) and lateralized sounds sources and that both these type can coexist in the listener’s experience [].

6) *Issues in headphones-conveyed sound*: Headphones auralization often produces an incorrect localization of virtual sound sources [] and many issues could be experienced.

a) *Externalization errors*: externalization is related to the perception of auditory distance such that there is a continuum in perceived locations of sources from inside the listener’s head to any external position. One of the most severe problem in AAR is a perceived effect of *lateralization* even in sounds that should be located in a real environment. To avoid such effect several techniques can be used. For instance, as expressed in [], the effect of a lateralized sound in headphone listening can be produced using amplitude and delay differences in two headphone channels corresponding to each source. The main goal of virtual acoustic synthesis should be to produce sounds that seem *externalized*, that is, outside the listener’s body []. In order to make a sound source externalized and let the user be capable of a correct judgement of the distance of the sound more sophisticated binaural techniques are needed. In particular, spectrum differences in the two ear signal due to *head-related transfer functions*, HRTF’s, play an important role. Moreover, acoustic cues such as the amount of reverberation and control of signal level are necessary for a successful auralization of a virtual sound source. These are explained into details in further sections.

b) *Localization errors*: *localization* error refers to the deviation of the reported position of a sound source from a measured ‘target’ location, i.e. the listener fails in matching the correct location of the sound. Localization errors can be divided in *azimuth* (deviations along the horizontal plane) and *elevation* errors (deviation from eye-level elevation) []. Dynamic cues related to head turning and other movements

of either a listener or a sound source should be taken into consideration []. These errors might come from some location accuracy. Determining the accurate position of the user of one of the most important task in an AR system owing to that system always produces output to the user based on his or her location in space []. That said, any location inaccuracy should be avoided using GPS receiver with high sensitivity and reliability. Here, the implementation design takes a crucial role as noted in [].

c) *Reversal errors*: sometimes called front-back or back-front ‘confusion’, that error refers to the judgement of a sound source as located on the opposite side of the interaural axis than the target position []. An informal proposal has been made, which tries to help the user in front-back discrimination on the basis of the familiarity of the effects on timbre cues, e.g. unique patterns in *early reflections* depending on the virtual sound location. Indeed, this has not yet been verified experimentally.

7) *Conclusion*: Although it is a recent field of research, several works and experiments in the context of audio augmented reality has been made so far. Most of research has been done mainly on the use of visual information with very little attention paid to the audio aspects.

The use of audio and its refinement in an augmented reality environment is a big deal. There are many issues and challenges that an audio environment presents []. For instance, providing a way for the user to orient himself or herself in an only-audio virtual space is even more difficult than providing visual cues. Another interesting aspects of audio, as also explained in [], is the fact that audio does not persist and it starts disappearing as soon as it has heard. That deals with a more attention due to the audio to capture every detail that may help the user retrieve a specific sound in a specific location.

Finally, with the advent of portable systems such as sensors-equipped smartphone and laptop a further development of *geolocated AR* has been made possible. That systems use locational sensing and user tracking technologies (e.g. Global Positioning System, or GPS) to let the user navigate through a audio-augmented space and to get a more defined listener’s position, which is essential for a correct externalization. More information will be given in further sections.

II. DESIGN AND IMPLEMENTATION CHOICES

also some argumentation why we used a HRTF model based on filtering rather than using a database (citations why database not efficient and why quality better with databases but good enough with filtering)

III. THE APPLICATION

A. openFrameworks

We decided to develop our application for Android platform. Nevertheless, we did not use the facilities and the built-in code provided by Android Studio, but instead we decided to use openFrameworks.

The reasons why we decided to develop using this programming environment was the need of a fast communication between the mobile platform and our

program code. It works using C++ language, which implies that the code is compiled directly into assembly language and hence works very fast [] and does not need to be interpreted by a virtual machine as it is the case for a java-built environment such as Android Studio [].

Discussing all the steps involved in the building process is beyond the scope of this paper (and more information can be found in [] and []), but a short explanation of what is openFrameworks and the main operations used to achieve our application will be given.

1) *Description*: OpenFrameworks, by definition given on its website, is:

an open source C++ toolkit for creative coding []

Since it is entirely written in C++, distributed under MIT license and actually runs on five operative systems and four IDEs, it undoubtedly presents some kind of advantages. It gives the opportunity to several users to deal with code designed to be minimal and easy to grasp and use it on the favourite IDE. In that case, we use both OS X and Windows platform, and Xcode and Visual Studio as integrated development environment. That *simple and intuitive framework for experimentation* [] is designed to work as a general purpose glue and wraps together several commonly used libraries, such as *OpenGL*, *OpenCv*, *PortAudio* and many more. Nowadays, this is a popular platform for experiments in generative and sound art and creating interactive installations and audiovisual performances []. The current operative version is 0.9.0.

2) *Addons*: its design philosophy claims for a collaborative environment. It thrives on the contributions of many people, and collaborate mainly on *addons* and projects. An *addons* is made of several snippets of code put together in order to extend openFrameworks functionality, bring some external framework and allow it to be integrated into openFrameworks project or make specific and complicated tasks easier and reusable in other contexts []. It also generally contains the library itself in a form that is ready to be linked to project binaries. In our case, as depicted in the next section, we used several *third-party* addons and built our own library called *ofxOrientation* to access sensors information about localization and orientation and provide the HTRF model with right values.

3) *Main issues*: openFrameworks requires such lots efforts to be mastered and its learning rate is quite steep. Its weakness is in fact due to the lack of documentation and people who don’t have at least some programming background could feel quite uncomfortable. Hopefully, its community is continuously providing discussion and projects, while addons are released with some example to compile and have to be created using a bunch of accurate rules.

4) *The openFrameworks project*: all openFrameworks project have a similar structure of folders and files. The most important folder among them is for sure the *src* folder. It contains all the source codes and consists at least of *main.cpp* (containing the *main()* function to let the operating system start the application), *ofApp.h* (containing declaration of the specific class) and *ofApp.cpp*, which contains definition of all functions declared in the previous file. All the methods in

that class are *event-handling* methods, that is they are triggered in response to events that happens inside the application such as mouse scrolling and program quitting.

To create a new project, we used the Project Generator wizard located in the same directory and directly provided by the environment. Such a way is simple and it is especially useful when dealing with several addons, which are automatically linked. Once we got the structure, we could access every source file in the *src* folder. At its simplest, working with an openFrameworks project is adding new code to the appropriate method, or just create a new one and declare it in the *ofApp.cpp*. In [] a further explanation of the main methods and their workflow is provided.

5) *Compiling: soon...*

B. Sensor data retrieving

Our application required two kind of sensors data: *GPS* and *orientation*.

- GPS is a way of determining location that has become amazingly common with the increasing of portable technologies. Any device that receives a GPS signal is called a *GPS receiver*. As explained in [] the receiver calculates its position by precisely timing the signals sent by the GPS satellites. We needed the GPS user position to calculate the distance to the sounds that were placed on fixed locations in the real environment.
- orientation usually is derived using a combination of field sensors, such as accelerometer and magnetometer. In this case, we needed orientation to define the angle the user is heading relative to the sound sources.

To access the GPS data we used *ofxMaps* and *ofxGeo* addons (reference to webpage/gitrepository). Reading the GPS position on an Android platform requires a listening mechanism and we implemented it in openFrameworks by calling the *ofRegisterGPSEvent()* in the *setup()* method of our main code. We also added a method to our application, in the *ofApp.h* precisely, to handle the updates from the Android OS system calls. It is named *locationChanged()* and adds an event handler for the dispatched event []. Moreover, the *startGPS()* method is called at the beginning and *stopGPS()* should call when quitting the application to avoid an overconsumption of phone battery.

(Something about smoothing if we used it, and maybe we have to look into the code how they used smoothing).

Since we are dealing with some kind of filtering depending on the sound sources position (e.g. the sound could appear in front of the user as well as behind), it is necessary to take into consideration the user's orientation. That is basically the angle created between the heading of the user and the geographic north pole. Initially, we assumed the user always faced to the north, i.e. the angle between the user and the geographic north pole¹ is equal to 0. If we can discard the deviation of the user's orientation, it is possible to define the angle between the user and the sound source as the *true bearing*². This is, by definition

¹There is a difference between *geographic* and *magnetic* north pole.

²It is possible refer to it simply as *bearing*.

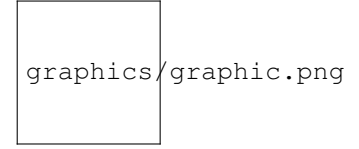


Fig. 1. Components of the model provided on a theoretical basis in [?]

the angle measured in degrees in a clockwise direction³ from the north line [].

If we also assume a sound source as a point that radiates in every directions, we could easily calculate the bearing of the user from the sound source, as shown in (add figure). The bearing could be accessed via the *ofxGeo* addon, which contains a specific method called *GeoUtils::bearingHaversine* that uses the ‘*haversine*’ formula:

$$\theta = \text{atan2}(\sin(\delta\gamma)\cos(\varphi_2))\cos(\phi_1)\sin(\varphi_2) - \sin(\varphi_1)\cos(\varphi_2)\cos(\delta\gamma))$$

where φ is the *latitude*, δ is the *longitude*.

Once we got the bearing, we can start supposing that the angle between the heading of the user and the sound source should be the sum of the bearing and the direction the user is facing. Doing so, we can retrieve it using a compass. In order to access the compass orientation we built our own addon accessing the ‘SENSOR TYPE ORIENTATION’ provided by the android API. Having this information we now had to differentiate between four cases: (some images and explanation about how we calculated the angle between our heading orientation and the sound). In order to smooth the sensor data (especially for the orientation) we filtered the sensory information by computing the median of a chunk of *mediansize* instances of the sensor output.

C. Audio engine

The audio engine was the most crucial part of our applications. As shortly mentioned in the introduction we decided to use a filter based HRTF model which claimed to be both ‘efficient’ and of reasonable ‘quality’ [?] in order to achieve a satisfying user experience and an efficient real time application. We basically followed the workflow (order of filter and delay lines) given in [?] and depicted in figure 1.

First, the monoaural sound which has been processed to add a range effect, goes to a head model and a room model. The binaural output of the former is then processed by the pinna model, and finally each one is added to the output of the room model to provide an externalization effect.

In the following formulas, the angles are measured in radians in the interaural-polar system, where θ is the azimuth angle and ϕ is the elevation angle. However, in our application the latter is no considered since all the sound sources are placed in the horizontal plane without elevation.

³All bearings are measured in a horizontal plane.

1) The head model:

The path length difference from the sound source to the two ears and the shadowing effect produced by the head at the far ear, lead to a delay and a intensity difference in the sound arriving at the left and right ears. In order to estimate the time delay, the Woodworth's formulas are used [?]:

$$ITD = (a/c)[\theta + \sin(\theta)] \quad [0 \leq \theta \leq \pi/2] \quad (1)$$

$$ITD = (a/c)[\pi - \theta + \sin(\theta)] \quad [\pi/2 \leq \theta \leq \pi] \quad (2)$$

where a is the approximated head radius, θ is the azimuth angle in radians and c is the head radius. The time differences between the audio signal reaches the head and the ears are therefore

$$T_L(\theta) = \frac{a + a\theta}{c} \quad (3)$$

$$T_R(\theta) = \frac{a - a\sin(\theta)}{c} \quad (4)$$

These formulas refers to a source in front of the head and on the right, with azimuths $0 \leq \theta \leq \pi/2$. If the source is placed on the right ($-\pi/2 \leq \theta \leq 0$), the expressions are reversed.

The head shadow effect is characterized in the following analog transfer function:

$$H(s, \theta) = \frac{\alpha(\theta)s + \beta}{s + \beta}, \text{ where } \beta = \frac{2c}{a} \quad (5)$$

====Maybe this part in the appendix...?====

Since this is an analog transfer function we had to derive the digital version by applying a bilinear transform applying the following substitution:

$$s = \frac{2}{T} \frac{z - 1}{z + 1}, \text{ where } T \text{ is the sampling interval in seconds} \quad (6)$$

Applying the substitution in equation 6 to equation 5, we get the following filter function in the digital domain:

$$H(z, \theta) = \frac{\alpha(\theta)(\frac{2}{T} \frac{z-1}{z+1}) + \beta}{(\frac{2}{T} \frac{z-1}{z+1}) + \beta} \quad (7)$$

To identify the filter coefficients we transformed Equation 7 and obtained the following frequency response⁴:

$$H(z, \theta) = \frac{2\alpha(\theta) + T\beta + z^{-1}(-2\alpha(\theta) + T\beta)}{2 + T\beta + z^{-1}(-2 + T\beta)} = \frac{Y(z)}{X(z)} \quad (8)$$

=====

This gives us the filter coefficients $a_0 = 2\alpha(\theta) + T\beta$ and $a_1 = -2\alpha(\theta) + T\beta$ as well as $b_0 = 2 + T\beta$ and $b_1 = -2 + T\beta$. Since we could not use the filter function in Matlab which has as input values the filter coefficients we had to isolate the output $Y(z)$ in Equation 8 as well as going from the frequency to the time domain. The result is given in Equation 9⁵:

$$Y[n] = \frac{a_0 X[n] + a_1 X[n-1] - b_1 Y[n-1]}{b_0} \quad (9)$$

⁴For the derivation of Equation 8 see appendix A

⁵For the derivation of Equation 9 see appendix B

2) The pinna model:

High frequency components that arrive to the listener's ears are influenced by the pinna, which provide some azimuth information, but it has been studied mainly because of its contribution to the estimation of the elevation [?].

Our model has the following form:

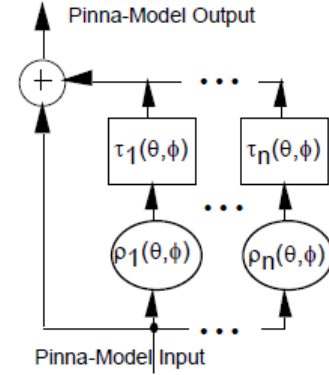


Fig. 2. Pinna model. Image taken from [?]

where the ρ_k are the reflection coefficients and the τ_k are the time delays of the k th event of a total of n . Informal listening tests showed that 5 events were enough to represent the pinna response and that it was convenient to use constant values for the amplitudes ρ_k , independent of azimuth, elevation and the subject [?]. The time delays seems to be properly approximated by the following formula:

$$\tau_k(\theta, \phi) = A_k \cos(\theta/2) \sin(D_k(90^\circ - \phi)) + B_k \quad (10)$$

In this equation, dependent on the azimuth and elevation, the A_k is an amplitude, B_k an offset and D_k is a scaling factor that has to be adapted to the individual listener. In the following table one can see the values for the parameters used in the pinna model. Only one set of values for D_k in our application have been used.

TABLE I
PINNA MODEL COEFFICIENTS

k	ρ	A_k	B_k	D_k
1	0.5	1	2	1
2	-1	5	4	0.5
3	0.5	5	7	0.5
4	-0.25	5	11	0.5
5	0.25	5	13	0.5

3) The room model: Finally, a room model is added to the two outputs of the pinna model...

D. User interface

include pictures of user interface etc.

IV. THE EXPERIMENT

One of the aims of our experiment was to investigate how fast users were able to locate and navigate to a single

sound source using our 3-D audio engine compared to a simple panning and 2/d distance calculation which served as a baseline performance. Therefore the time needed from the beginning of a trial until the sound source was found was measured.

Additionally we assessed a qualitative analysis of how intuitive it was for subjects to locate the sound as well as investigating whether the sound seemed to be embedded in the real scenery. These results were again compared to the panning model. The data was obtained with a 10 point likertscale questionnaire going from "don't agree at all" to "completely agree".

(Here we can write a one paragraph preview of our results.)

A. Materials and methods

1) *Participants*: One participant had to be discarded because of Qualitative data were collected from 11 participants (3 women and 8 men; mean age: 25, ranging from 22 to 30 years). For the quantitative data 6 more participants were measured summing up to in total 17 participants (4 women and 13 men; mean age: 25, ranging from 22 to 30 years). All participants were students coming from different backgrounds.⁶ Six students were already familiar with 3-D audio sound and four of them had experience in audio navigation before conduction of the experiment. The other participants were naive. All participants reported normal or corrected to normal vision and no hearing deficits. Participants received no form of compensation other than gratitude.

2) *Apparatus and stimuli*: The experiment took place in a free field (grass) in "Valby Parken" in Copenhagen, Denmark (Coordinates: 55°38'22.3"N, 12°31'27.4"E). The mobile device used was a Motorola Moto G running with Android version 5.0.2. The device was connected to Sony MDR ZX600 stereo headphones. Only one person at a time was tested.

To create the auditory stimulus, digital vocal recordings of a male voice (age: ...) were collected using a Recordings were done with ... bits of resolution for amplitude at a sampling rate of ... kHz. Auditory stimuli were normalized in peak amplitude so that at the minimal distance to the soundsource (7 meter) using the maximal level of volume on the mobile device, the soundlevel of each soundsource was at ... dB.

3) *Design*: The experiment comprised one within subject factor at 2 levels (sound engine). The two levels of the sound engine were our 3-D audio engine (3-D) and a simple panning and distance amplitude modulation (panning).

The starting point was fixed and the distance of the sound-source from that point was always 48 meters.⁷ The sound could appear at any angle from the starting point.

For all trials the same soundfile was played. Each participant performed a total of 6 trials meaning 3 trials in each of the two conditions. Participants were encouraged to take a short break whenever they felt fatigued.

⁶However, 14 subjects were master students in Sound and Music Computing at Aalborg University Copenhagen.

⁷However, because of poor GPS accuracies, the location of the soundsource was calculated with the momentary GPS coordinates provided by the phone sensor instead of the actual fixed position in the free field. This ensured that the initial distance to the sound source was always the same.

4) *Procedure*: The position (angle from the starting point) of the sound sources in the different trials were completely randomized within and between subjects independently from the sound engine. This ensured that participants could not predict the location of the sound in the next trial. The order of trials with the one or the other sound engine was not randomized within subjects. Nevertheless, half of the participants started the experiment with the 3-D audio engine while the other half started with the audio panning sound engine to counterbalance for effects related to training. After each block of trials (one block were 3 trials with one of the sound engines) participants had to fill out the qualitative questionnaire.

Each trial was initiated by the experimenter with a button press after which the sound source was played. No visual feedback about the location of the sound source or the already passed time to locate it was given. For localization, subjects had to rely solely on the auditory cues given via the headphones. /footnoteAdmittedly participants could have used the position of the starting point as some kind of fixpoint. However, participants did know nothing about the possible distances of the sound source. Additionally, questions after the experiment about applied strategies did not reveal that participants used visual cues for localizing the sound source. As soon as the participants reached a radius of 7 meter around the sound source, an earcon was played indicating the success of the localization. The participants went back to the initial starting point and the next trial was started.

Before the experiment, participants were familiarized with the target sound as well as the earcons that were played when having reached the location of the virtual sound source. Additionally, all participants performed at least two practice trials with the sound engine they were tested with first before the actual measurements with that sound engine began. The location in the practice trials were fixed at the same positions for all participants to allow us to guide the subject when they got lost.

Participants were instructed to find the sound source as fast as possible. No information about the possible locations (distances and angles) was given to the subjects. Trials which lasted longer than 5 minutes were aborted and labeled as "not found".

5) *Analyses*: For the quantitative analysis 26 trials were excluded because either analyses of GPS data for those trials revealed very poor precision (sudden jumps or no changes in the GPS signal for a period of time) or the orientation sensor got stuck which either increased the time needed to find the sound source to a great extent or made it even impossible for the subjects to locate the sound source within 5 minutes. One subject had to be discarded completely from the analyses because there was no single trial with reliable GPS for the panning engine. The other 16 subjects performed at least one trial in each condition. The remaining trials to be analysed summed up to 41 trials in the 3-D condition and 42 trials in the panning condition.

The reasons for choosing a 10 point likertscale instead of the often referenced 5 or 7 point scales [] were that a greater number of responses should equalize the distances between the possible answer possibilities supporting treating likertscale

data also as interval data []. This gives us the grounds to use paired samples t-test for the qualitative analysis which is more powerful than the Wilcoxon signed rank test enabling us to detect smaller differences in ratings []. Since we only had few questions, having so many answer possibilities should not have fatigued the subject too much so that accuracy of the responses should not have suffered due to too many answer possibilities []. Though, it has to be noted that subject were outside when filling out the questionnaire which due to cold temperatures might have driven the subjects to respond in a fast pace leading to less accurate responses.

The threshold chosen to correctly reject the null-hypothesis was always 5% ($\alpha = .05$).

B. Results

1) *Quantitative:* On average subjects needed 104.6 seconds (SD = 44.8) to find the sounds for the 3-D conditions and 109.8 seconds (SD = 44.1) for the panning condition.

We computed a paired samples t-test to compare the results for both sound engines which revealed no significant difference between the means for both sound engines ($t(15) = -.48$, $p = .64$). However, checking whether our data is normally distributed using the Shapiro-Wilk test showed that for the 3-D data the assumption of normality was violated ($Z = .81$, $p = .004$) and close to significance in the panning condition ($Z = .89$, $p = .054$). Therefore, we additionally computed the non-parametric Wilcoxon signed rank test which does not assume normal distributed data. Nevertheless, we could not find a significant difference in scores for time needed to find the soundsource between both conditions (median: 3-D=85.3, panning=95.9; $Z = -.83$, $p = 0.41$).

2) *Qualitative:* The Shapiro-Wilk test showed that the data was normally distributed both for the 3D model ($Z = 0.933$, $p = 0.44$) and the Panning model ($Z = 0.95$, $p = 0.67$).

Evaluation of the questionnaires, Table II, showed that the mean and median values for questions 1, 2, 3 and 4 for the 3D model were only numerically different from the panning. Neither the paired samples t-test nor the Wilcoxon signed rank test showed significant results comparing the two conditions (3-D vs. panning). However, performing a one sample t-test on means of the four questions for each model comparing it with a neutral response (5.5) showed that users rated the connection between their position and what was presented via the headphones better than neutral for the 3-D sound engine ($t(10) = 4.32$, $p = .002$; $Z = 2.74$, $p = .006$) but not for the panning ($t(10) = 1.07$, $p = .31$; $Z = .90$, $p = .367$). For all other questions there were no significant differences.⁸

Summing the results for the first four questions to form an overall category of quality of the two sound engines lead to the mean and median values depicted in Table III. Again, the difference between both sound engines were only

⁸It is worth noting that choosing mean and median as 5 instead of 5.5, showed that also for question two, a paired samples t-test and the wilcoxon signed rank test were significant for the 3-D ($t(10) = 2.52$, $p = .03$; $Z = 2.19$, $p = .028$) but not for the panning sound engine ($t(10) = 1.19$, $p = .26$; $Z = 1.2$, $p = .23$).

⁹It should be noted that for this question lower values represent higher quality of a 3-D sound because of the direction of the question.

Question	Mean (s)		Std (s)		Median (s)	
	3D model	Panning	3D model	Panning	3D model	Panning
1	6.36	5.91	1.75	1.70	6	6
2	6.91	5.82	2.51	2.27	8	6
3	7.09	6.00	1.22	1.55	7	6
4 ⁹	6.64	6.09	2.58	1.97	7	5
5	6.09	6.55	2.21	2.54	7	7

TABLE II
QUESTIONNAIRE RESULTS

numerically visible since performing a paired samples t-test and a Wilcoxon signed rank test showed no significant results ($t(10) = 1.05$, $p = 0.32$; $Z = -1.07$, $p = .283$). Furthermore, running a one sample t-test and wilcoxon signed rank test to compare the results of both sound engines to a mean (5.5) that represents being neutral towards a question showed no significant results for both conditions (HRTF: $t(10) = 1.54$, $p = .154$; $Z = 1.54$, $p = .13$; Panning: $t(10) = -.22$, $p = .829$; $Z = -0.5$, $p = .96$). However, summing only the first three questions, showed that the 3-D model ratings were significantly higher than the neutral mean (5.5) ($t(10) = 2.83$, $p = .018$; $Z = 2.09$, $p = .036$) which was not the case for the panning condition ($t(10) = .86$, $p = .41$; $Z = 0.76$, $p = .447$).

Question	Mean (s)		Std (s)		Median (s)	
	3D model	Panning	3D model	Panning	3D model	Panning
1 - 4	23.73	21.64	3.72	5.45	24	21

TABLE III
OVERALL RESULTS

V. DISCUSSION

There was no performance difference between the two sound engines. This could have 2 reasons. first it is actually equally good. Probably the difference if there is a difference at all is so small that with such few data we could not detect the small difference. However the mean value was lower although not significantly. If our results are not significant this could be because of low sample size -> "Numerous methodologists have cautioned that a small sample size implies low statistical power, that is, a high probability of Type II error" [?]

As for the qualitative analyses the explanation why our 3-D was not rated significantly better then the other could be due to a too high reverberation level for a free field. since a lot of participants had background in audio knowledge about a proper reverberation level for a free field can be assumed which even increased the negative effect on the ratings of the 3-D audio engine. This might have nullified the difference in left to right ear panning vs. 3-D sound. We had even worse numerical values for the fourth question for the HRTF. Also the questions might not have been posed in an optimal way...

VI. CONCLUSION

ACKNOWLEDGMENTS

LIST OF FIGURES

1	Components of the model provided on a theoretical basis in [?]	4
2	Pinna model. Image taken from [?]	5

APPENDIX A
DERIVATIONS

$$\begin{aligned}
H(z, \theta) &= \frac{\alpha(\theta)(\frac{2}{T}\frac{z-1}{z+1}) + \beta}{(\frac{2}{T}\frac{z-1}{z+1}) + \beta} \\
&= \frac{\frac{2\alpha(\theta)(z-1)}{T(z+1)} + \frac{T\beta(z+1)}{T(z+1)}}{\frac{2(z-1)}{T(z+1)} + \frac{T\beta(z+1)}{T(z+1)}} \\
&= \frac{\frac{2\alpha(\theta)(z-1) + T\beta(z+1)}{T(z+1)}}{\frac{2(z-1) + T\beta(z+1)}{T(z+1)}} \\
&= \frac{2\alpha(\theta)(z-1) + T\beta(z+1)}{2(z-1) + T\beta(z+1)} \\
&= \frac{z2\alpha(\theta)(1-z^{-1}) + zT\beta(1+z^{-1})}{z2(1-z^{-1}) + zT\beta(1+z^{-1})} \\
&= \frac{2\alpha(\theta)(1-z^{-1}) + T\beta(1+z^{-1})}{2(1-z^{-1}) + T\beta(1+z^{-1})} \\
&= \frac{2\alpha(\theta) - 2\alpha(\theta)z^{-1} + T\beta + T\beta z^{-1}}{2 - 2z^{-1} + T\beta + T\beta z^{-1}} \\
&= \frac{(2\alpha(\theta) + T\beta) - (2\alpha(\theta) + T\beta)z^{-1}}{(2 + T\beta) - (2 + T\beta)z^{-1}} \\
&= \frac{(2\alpha(\theta) + T\beta) + (-2\alpha(\theta) + T\beta)z^{-1}}{(2 + T\beta) + (-2 + T\beta)z^{-1}}
\end{aligned}$$

LIST OF TABLES

I	Pinna model coefficients	5
II	Questionnaire results	7
III	Overall results	7

APPENDIX B
DERIVATIONS

$$\begin{aligned}
H(z, \theta) &= \frac{Y(z)}{X(z)} = \frac{a_0 + a_1 z^{-1}}{b_0 + b_1 z^{-1}} \\
&\iff Y(z)(b_0 + b_1 z^{-1}) = X(z)(a_0 + a_1 z^{-1}) \\
&\iff Y(z)b_0 + b_1 Y(z)z^{-1} = a_0 X(z) + a_1 X(z)z^{-1} \\
&\iff Y(z) = \frac{a_0 X(z) + a_1 X(z)z^{-1} - b_1 Y(z)z^{-1}}{b_0} \\
&\rightarrow Y[n] = \frac{a_0 X[n] + a_1 X[n-1] - b_1 Y[n-1]}{b_0}
\end{aligned}$$

REFERENCES