

Head Tracking Binaural Localization System for Horizontal Sound Source Detection *

Nolan Lem¹ and Jens Ahrens²

¹Center for Computer Research in Music and Acoustics (CCRMA) Stanford University

²Technische Universität Berlin, Germany

Abstract—This project attempted to reconstruct an interactive spatialization auditory scene by using a head tracking sensor in conjunction with a binaural panning system implemented in the MaxMSP programming environment. This study was motivated by the need to be able to demonstrate in more informal settings—such as classroom demonstrations, or paper presentations—the general experimental environment that many researchers have used to examine the localization of point-like sound sources on the horizontal plane.

I. INTRODUCTION

Sound localization tasks involving point-sources, are well-studied in the psychoacoustics literature. The human auditory system exploits several auditory features and percepts to infer source location from an incoming sound. These include ITDs, ILDs, spectral cues, head movement cues, intensity/loudness cues, familiarity to the sources, Direct-to-reverb ratios (DRR), and visual and other non-auditory cues (cite). Among these, the last four are associated with distance cues while the others allow us to align our direct localizations. It is well known that have much more spatial accuracy in identifying sound sources on the horizontal plane as compared to our vertical (azimuth) source detection. Head-related Transfer Functions (HRTFs) describe the frequency responses of our ears. They can be characterized by a filter that describes how a sound from a specific source will arrive at each of the ears. Because of the localization techniques we use to detect sound in space, we can simulate "virtual" sources by taking advantage of these transfer functions typically by deriving head-related impulse responses (HRIRs) which is the Fourier transform of the HRTF. The HRTF reflects the signal differences each ear receives given its orientation in space and the shape of a person's head and ear canal.

When multiple sources are involved, we often create sound-object representations that facilitates our ability to understand and parse out information from an auditory scene. This fusion of synchronous auditory streams (or otherwise) is usually a function of how well temporally correlated the signals are which results in summing localizations from which point-like virtual sources emerge. In relatively reflection-free acoustic spaces, there are several factors that influence the perception of sound sources.

These include the signal level or duration which tends to increase the perceived width of the distribution, and the frequency of the source stimuli which tends to decrease the width percept. Typically, studies have used the "concurrent minimum audible angle (CMAA)" to discriminate between listeners' abilities to discern angular separation among sources presented as point-like sources projected from multiple speaker channels arranged in different spatial orientations. Because of the effects of stimulus frequency, researchers have studied a plethora of stimuli containing various spectral content to examine its effects on listeners' ability to detect sound sources in spatial configurations. Additionally, researchers have often employed the use of broadband noise in the context of point-source identification tasks.

II. CONVOLUTION WITH HRIRs

A. Description

This system can be formulated by the following description. Let a group of left and right ear HRIRs, $h_i^l(n)$ and $h_i^r(n)$, each of size N and separated by 1° of angular horizontal separation be convolved with a source signal, $x_{source}(n)$. (1).

$$\begin{aligned} y_{i,t}^l(n) &= (h_i^l * x)(n) \\ y_{i,t}^r(n) &= (h_i^r * x)(n) \end{aligned} \quad (1)$$

where i indexes HRIR associated with each degree of separation and $y_i^l(n)$ and $y_i^r(n)$ are the outputs of source signal convolved with the respective HRIRs for the left (l) and right (r) ears respectively. Let $W(n)$ be an equal-power window function also of size N that "crossfades" between two $h_i^{l,r}(n)$ impulse responses when we traverse (an auditor moves their head) along two sequential angles, y_i from the y_{i-1} .

$$y_{i,t}^{l,r}(n) = \sum_{t=1}^T W(n - \frac{2Nt-1}{2}) y_{i-1}^{l,r}(n - Nt) \quad (2)$$

where $y_{i,t}^{l,r}$ is the output signal of length NT that accumulates to characterize the auditor's head movements from time t to T .

- [2] L. Ludovico, D.A. Mauro, and D. Pizzamiglio, HEAD IN SPACE: A HEAD-TRACKING BASED BINAURAL SPATIALIZATION SYSTEM , 2010
- [3] J. Blauert, and W. Lindemann, Spatial mapping of intracranial auditory events for various degrees of interaural coherence , in Journal of Acoustical Society of America., vol. 79 (3), Mar. 1986
- [4] K. Hiyama, S. Komiyama, and K. Hamasaki, The minimum number of loudspeakers and its arrangement for reproducing the spatial impression of diffuse sound field , Audio Engineering Society Convention., Oct. 5-8, 2002
- [5] O. Santala, and V. Pulkki, Directional perception of distributed sound sources , Journal of Acoustical Society of America., vol. 129 (3), Mar. 2011
- [6] V. Pulkki, Coloration of Amplitude-Panned Virtual Sources , Audio Engineering Society Convention., (110) May 2001
- [7] A. Lindau, HRTF Datensatz für die Horizontalebene Deutsche Telekom Laboratories, 2008