

Chapter 4

Sound in space

Federico Avanzini

Copyright © 2007 by Federico Avanzini.

All rights reserved except for paragraphs labeled as *adapted from* <reference>.

4.1 Introduction

In the previous chapters, *Sound modeling: signal based approaches*, and *Sound modeling: source based approaches*, we have studied models for sounds and sound sources. We now move a step further and examine the effects of the medium in which sound propagates: one of the most frequently encountered effect is reverberation. We will see the physical and perceptual background of reverberation, as well as some of the most known reverberation algorithms.

Then we turn to the receiver, and specifically examine a human receiver with two ears. We will study how and to what extent such a receiver can gain information about the incoming direction and distance of an emitted sound, and we will review some *3-D sound* processing techniques by which a virtual sound source can be positioned in some point of the space around a listener.

WARNING: this chapter is at a draft stage.

4.2 Reverberation: physical and perceptual background

4.2.1 Basics of room acoustics

4.2.1.1 Sound waves in a closed space

We have analyzed in Chapter *Sound modeling: source based approaches* the D'Alembert equation which describes sound propagation within a perfectly elastic medium. While the 1-D D'Alembert equation can be used to model strings or acoustic tubes, the 3-D equation describes sound propagation in space:

$$\nabla^2 p(\mathbf{x}, t) = \frac{1}{c^2} \frac{\partial^2 p}{\partial t^2}(\mathbf{x}, t), \quad (4.1)$$

where \mathbf{x} represents Euclidean coordinates in space and p is the acoustic pressure. The symbol $\nabla^2 = \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} + \frac{\partial^2}{\partial z^2}$ stands for the 3-dimensional Laplacian operator. As opposed to mechanical vibrations in a string or membrane, acoustic vibrations are *longitudinal* rather than transversal, i.e. the air particles are displaced in the same direction of the wave propagation. The constant c has the dimensions m/s of a velocity and indeed is sound velocity in air.

By adding suitable boundary conditions we can gain a description of waves of particle velocity within a three-dimensional enclosure. Let us start with the simplest possible 3-D enclosure, a rectangular room with perfectly smooth and rigid walls. More precisely, we define the domain \mathcal{D} of the problem to be a parallelepiped with edges of length L_x, L_y, L_z :

$$\mathcal{D} = \{\mathbf{x} = (x, y, z); 0 \leq x \leq L_x, 0 \leq y \leq L_y, 0 \leq z \leq L_z\} \quad (4.2)$$

The boundary conditions on the rigid walls of this enclosure require the air velocity perpendicular to each wall to be zero at the wall. Then one can provide analytically a solution of Eq. (4.1) in terms of stationary modes:

$$p(\mathbf{x}, t) = q(t)f(\mathbf{x}) \quad (4.3)$$

If we are working with acoustic pressure p then the conditions on the boundary \mathcal{B} of the parallelepiped are $\partial p / \partial \mathbf{x}(\mathcal{B}) = 0$. Then

$$f_{n,m,l}(\mathbf{x}) = \sqrt{\frac{2}{L_x}} \sqrt{\frac{2}{L_y}} \sqrt{\frac{2}{L_z}} \cos(k_n^{(x)} x) \cos(k_m^{(y)} y) \cos(k_l^{(z)} z) \quad (4.4)$$

where we can define the *wavenumbers* $\mathbf{k}_{n,m,l}$ as

$$\mathbf{k}_{n,m,l} = \left(k_n^{(x)}, k_m^{(y)}, k_l^{(z)} \right) \quad \text{with } k_n^{(x)} = \frac{n\pi}{L_x}, k_m^{(y)} = \frac{m\pi}{L_y}, k_l^{(z)} = \frac{l\pi}{L_z}. \quad (4.5)$$

Analogously to the 1-D case discussed for modal synthesis in Chapter *Sound modeling: source based approaches*, these functions are a orthonormal basis for the space $\mathcal{L}^2(\mathcal{D})$.

Then the temporal part is given by

$$q_{n,m,l}(t) = \cos(\omega_{n,m,l} t + \phi_{n,m,l}), \quad \text{with } \omega_{n,m,l} = c\pi \sqrt{\left(\frac{n}{L_x}\right)^2 + \left(\frac{m}{L_y}\right)^2 + \left(\frac{l}{L_z}\right)^2} \quad (4.6)$$

The frequencies $\omega_{n,m,l}$ clearly are a non-harmonic series. However each of the three spatial directions (where only one of the three indexes (n, m, l) is varying) is associated to a harmonic subseries.

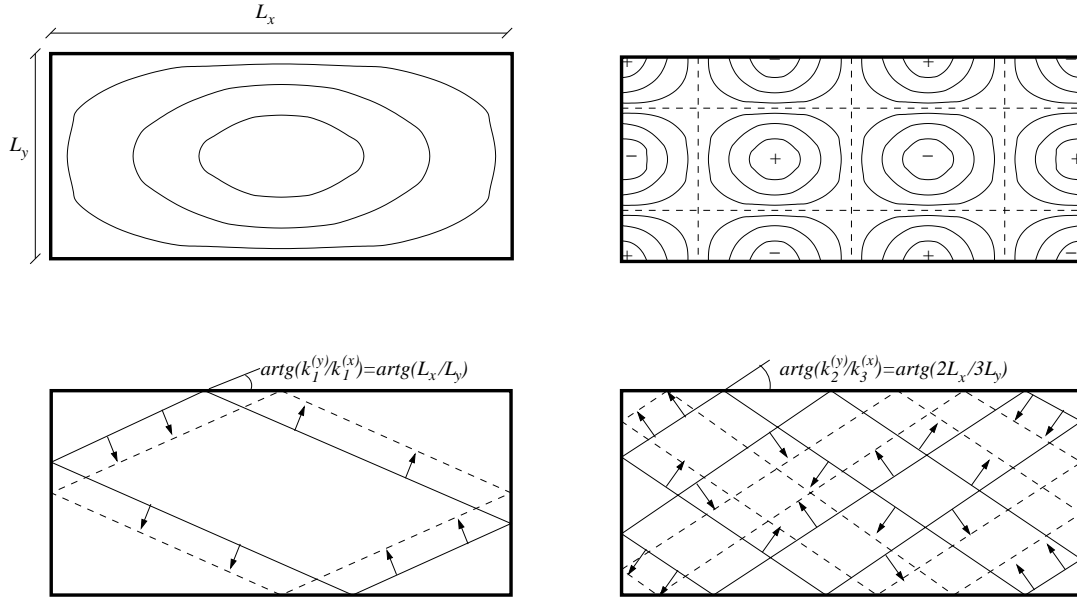


Figure 4.1: Plane wave loops $(1, 1, 0)$ and $(3, 2, 0)$, as seen on the (x, y) plane.

Analogously to the 1-D case a mode (n, m, l) has nodal surfaces, which corresponds to the regions where $f_{n,m,l}(\mathbf{x}) = 0$. It is easy to see that these are planes parallel to the walls of the rectangular room.

A normal mode of the form (4.3) can be written as a superposition of waves traveling in different directions. This can be easily seen through multiple application of Werner formulas¹, which yields

$$p_{n,m,l} = \sqrt{\frac{2}{L_x}} \sqrt{\frac{2}{L_y}} \sqrt{\frac{2}{L_z}} \sum \cos \left[\mathbf{k}_{n,m,l}^{\pm\pm\pm} \cdot \mathbf{x} \pm (\omega_{n,m,l} t + \phi_{n,m,l}) \right], \quad (4.7)$$

where we have defined $\mathbf{k}_{n,m,l}^{\pm\pm\pm} = (\pm k_n^{(x)}, \pm k_m^{(y)}, \pm k_l^{(z)})$, and where the summation has to be extended over the sixteen possible combinations of signs in the argument. This means that for each mode there are eight directions of wave propagation, each one associated to one $\mathbf{k}_{n,m,l}^{\pm\pm\pm}$ vector.

Figure 4.1 visualizes the wavefronts for the modes $(1, 1, 0)$ and $(3, 2, 0)$: these result in plane wave loops having constant length.

We now want to derive an estimate of the *modal density*, i.e. the average density of eigenfrequencies on the frequency axis. From Eq. (4.5) one observes that the allowed values for the wave numbers k are distributed on a regular point lattice in the 3-D space depicted in Fig. 4.2(a). The number N_f of eigenfrequencies in the frequency range from 0 to f equals the number of lattice points contained in the sphere octant of radius k depicted in Fig. 4.2(b). In other words, $N_f = V_f/V_0$, where V_f is the volume of the sphere octant of radius k and V_0 is the average volume per lattice point. The former is one octave of the sphere volume, $V_f = \pi k^3/6$, while the latter can be estimated as the volume of the cube depicted in Fig. 4.2(b), whose edges have lengths $\pi/L_x, \pi/L_y, \pi/L_z$, respectively. Therefore

¹ $2 \cos \alpha \cos \beta = \cos(\alpha - \beta) + \cos(\alpha + \beta)$.

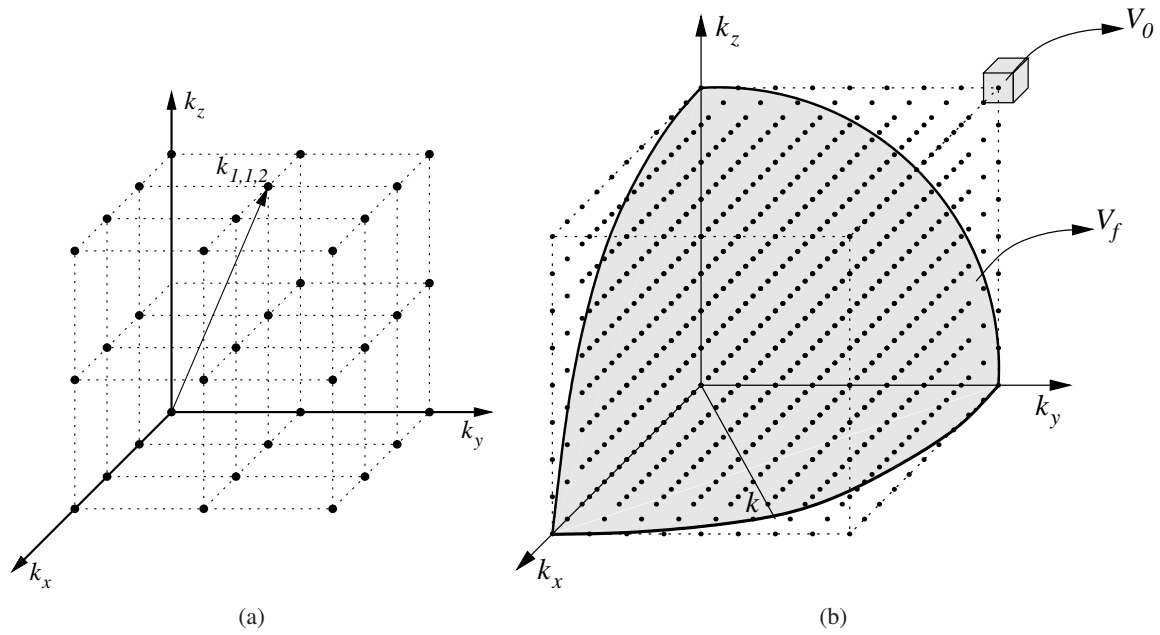


Figure 4.2: Estimation of modal density; (a) distribution of wavenumbers on a regular point lattice, (b) estimation of the amount of wavenumbers contained in a spherical octant of radius k .

$V_0 = \pi^3/V$, where $V = L_x L_y L_z$ is the room volume. One finally obtains

$$N_f = \frac{\pi k^3/6}{\pi^3/V} = \frac{4\pi}{3} V \left(\frac{f}{c}\right)^3. \quad (4.8)$$

. The modal density is estimated as the derivative of N_f with respect to frequency:

$$D_f(f) = \frac{dN_f}{df} = \frac{4\pi V}{c^3} f^2 \quad (4.9)$$

4.2.1.2 Sound sources and room impulse responses

More realistic situation: within the parallelepiped we have a source, moreover we consider complex, non-ideal, boundary conditions (wall absorption)

We want to find the solution of the wave equation in the parallelepiped, in the presence of a sound source: the distribution in space of the source is described by a continuous density function $\bar{f}(\mathbf{x})$, while the time-domain signal emitted by the source is described by a function $\bar{q}(t)$: this means that $\bar{q}(t) \cdot \bar{f}(\mathbf{x}) dV$ is the volume velocity of a volume element dV at time t .

Then

$$\nabla^2 p(\mathbf{x}, t) = \frac{1}{c^2} \frac{\partial^2 p}{\partial t^2}(\mathbf{x}, t) - \rho_{air} \bar{f}(\mathbf{x}) \frac{d\bar{q}}{dt}(t). \quad (4.10)$$

Since the $f_{n,m,l}$ of Eq. (4.4) are a basis for $\mathcal{L}^2(\mathcal{D})$, we can project both the source density function \bar{f} and the solution p of Eq. (4.10) on this basis:

$$\bar{f}(\mathbf{x}) = \sum_{n,m,l} \bar{F}_{n,m,l} f_{n,m,l}(\mathbf{x}), \quad P(\mathbf{x}, s) = \sum_{n,m,l} P_{n,m,l}(s) f_{n,m,l}(\mathbf{x}). \quad (4.11)$$

Note that in the second of the above equations we have implicitly assumed to work in the Laplace domain instead of the time domain. If we can find the unknown coefficients $P_{n,m,l}(s)$ as functions of the known coefficients $\bar{F}_{n,m,l}$, then we have the solution $P(\mathbf{x}, s)$ or equivalently $p(\mathbf{x}, t)$. If one inserts both series into Eq. (4.10) the result is

$$P_{n,m,l}(s) = s\rho_{air}c^2Q(s)\frac{\bar{F}_{n,m,l}}{s^2 + c^2k_{n,m,l}^2}. \quad (4.12)$$

Now we consider the special case of a point source located at a certain point \mathbf{x}_0 of the room and emitting an impulsive sound signal. Under this assumption one has $\bar{f}(\mathbf{x}) = \delta(\mathbf{x} - \mathbf{x}_0)$, where the function $\delta(\cdot)$ is the Dirac delta. This implies that the coefficients $\bar{F}_{n,m,l}$ are in this case $\bar{F}_{n,m,l} = f_{n,m,l}(\mathbf{x}_0)$. Moreover, if the sound source is emitting an impulse $\bar{q}(t) = \delta(t)$, then the corresponding spectrum is $Q(s) = 1$. If one substitutes the coefficients (4.12) into the second of Eqs. (4.11), the result is

$$P(\mathbf{x}, s) := H_{\mathbf{x}_0, \mathbf{x}}(s) = s\rho_{air}c^2 \sum_{n,m,l} \frac{f_{n,m,l}(\mathbf{x})f_{n,m,l}(\mathbf{x}_0)}{s^2 + c^2k_{n,m,l}^2}. \quad (4.13)$$

This is the acoustic pressure generated in \mathbf{x} by a point source located at \mathbf{x}_0 and emitting an impulse. If we take the inverse Laplace transform, $h_{\mathbf{x}_0, \mathbf{x}}(t) = \mathcal{L}^{-1}\{H_{\mathbf{x}_0, \mathbf{x}}\}(t)$, this is what we call a *Room Impulse Response (RIR)*, measured at point \mathbf{x} after an impulse emitted in \mathbf{x}_0 .

Wall absorption: the normal modes have now complex eigenvalues $k_{n,m,l}$:

$$k_{n,m,l} = \omega_{n,m,l}/c + j\delta_{n,m,l}/c, \quad \delta_{n,m,l} \ll \omega_{n,m,l}. \quad (4.14)$$

Then Eq. (4.13) is telling us that the RIR is a superposition of numerous second-order resonant systems, each with center frequency very close to $\omega_{n,m,l}$ and damping constant very close to $\delta_{n,m,l}$.

$$h_{\mathbf{x}_0, \mathbf{x}}(t) = \begin{cases} 0 & t < 0 \\ \sum_{n,m,l} A_{n,m,l}(\mathbf{x}_0, \mathbf{x}) e^{-\delta'_{n,m,l}t} \cos(\omega'_{n,m,l}t + \phi_{n,m,l}) & t \geq 0 \end{cases} \quad (4.15)$$

The function $h_{\mathbf{x}_0, \mathbf{x}}(t)$ completely describes the room response for a source in \mathbf{x}_0 and a receiver in \mathbf{x} : if the emitted sound is not an impulse but a generic signal $\bar{q}(t)$, then the response will be –as usual– the convolution of the signal with the impulse response: $s(t) = [\bar{q} * h_{\mathbf{x}_0, \mathbf{x}}](t)$.

Now in normal rooms damping constants typically lie between 1 and 20 Hz: this justifies the assumption of very small δ coefficients in Eq. (4.14), and moreover tells that half-widths of these resonant systems are of the order of 1 Hz. If we compare this finding to the modal density estimate given in Eq. (4.9), we see that the average spacing of eigenfrequencies is smaller by several orders of magnitude than half-widths. Therefore each single resonant peak always covers many others and it is practically impossible to excite a single room resonance e.g. with a sinusoidal signal.

4.2.1.3 Reverberation time

From Eq. (4.15) we see that room reverberation add a decaying tail to a source signal. One of the most important parameters derived from this equation is the *reverberation time* T_r , i.e. the time required for the sound pressure to decay 60 dB. Clearly T_r is related to the absorption coefficients $\delta_{n,m,l}$. An approximate description of T_r can be derived as follows.

Given a source signal $\bar{q}(t)$, the resulting room response $s(t)$ is

$$s(t) = [\bar{q} * h_{\mathbf{x}_0, \mathbf{x}}](t) = \dots = \sum_{n,m,l} c_{n,m,l} e^{-\delta'_{n,m,l}t} \cos(\omega'_{n,m,l}t + \psi_{n,m,l}) = \sum_{n,m,l} c_{n,m,l} s_{n,m,l}(t) \quad (4.16)$$

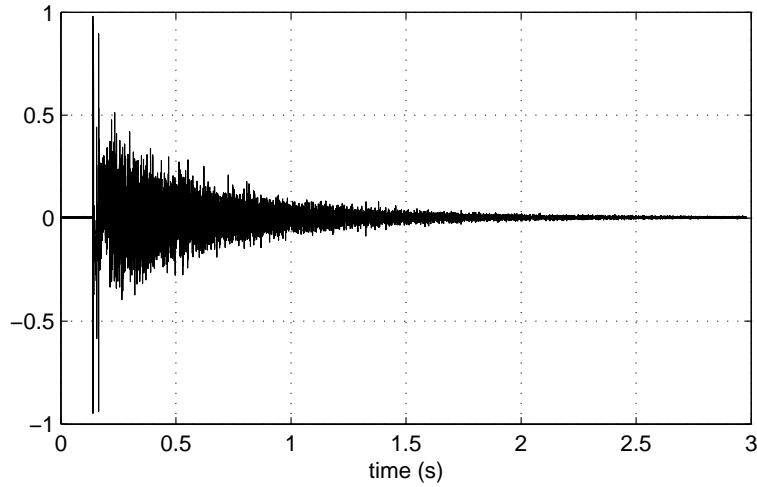


Figure 4.3: Impulse response of a very reverberant environment (a cathedral).

where the $c_{n,m,l}$'s and the $\psi_{n,m,l}$'s will vary depending on the signal \bar{q} , and where we have introduced the notation $s_{n,m,l}(t) = e^{-\delta'_{n,m,l}t} \cos(\omega'_{n,m,l}t + \psi_{n,m,l})$ for brevity. An expression proportional to the energy density is obtained by squaring $s(t)$:

$$w(t) = [s(t)]^2 = \sum_{n,m,l} \sum_{n',m',l'} s_{n,m,l}(t) s_{n',m',l'}(t). \quad (4.17)$$

A simpler expression can be found by averaging $w(t)$ over time and exploiting the circumstance that the exponential terms vary slowly (as the δ 's are small). By averaging the cosine products only, the products with $(n, m, l) \neq (n', m', l')$ cancel on the average, and the products with $(n, m, l) = (n', m', l')$ give a value $1/2$. If one makes the further assumption of nearly uniform damping, i.e. $\delta_{n,m,l} \sim \delta_0$, then we obtain the following result:

$$\langle w(t) \rangle = \sum_{n,m,l} c_{n,m,l}^2 e^{-2\delta_{n,m,l}t} \sim e^{-2\delta_0 t} \sum_{n,m,l} c_{n,m,l}^2. \quad (4.18)$$

This equation tells that for uniform damping the energy of the reverberation tail decays exponentially. In particular the reverberation time T_r is in this case derived as

$$-60 = 10 \log \left(e^{-2\delta_0 T_r} \right), \quad \Rightarrow T_r = \frac{6.91}{\delta_0}. \quad (4.19)$$

In general one cannot assume uniform damping, and as a consequence T_r is a function of frequency. Fortunately, however, the reverberation level falls in many practical cases in a fairly exponential fashion and therefore an overall reverberation time T_r can be defined and measured. We will return on this concept in Section 4.2.2, where we will examine typical reverberation time for various environments.

M-4.1

Write a function that computes the reverberation time T_r given a signal representing a RIR.



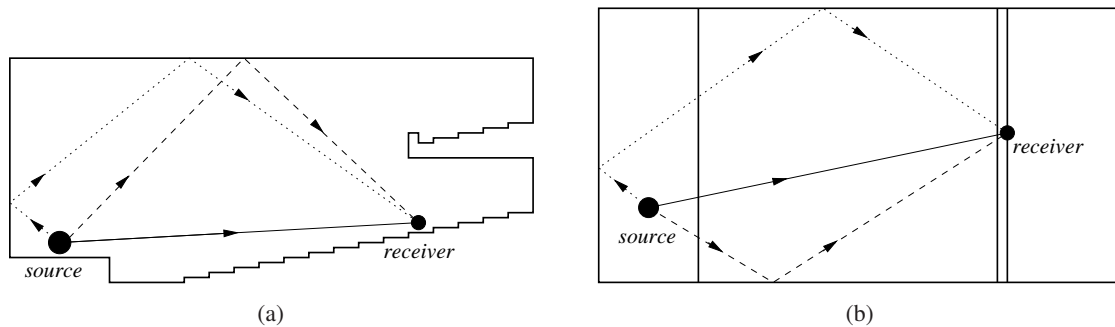


Figure 4.4: Acoustic rays from a source to a receiver (a) in a vertical room section and (b) in a horizontal room section. Solid lines represent the direct sound, dashed lines represent first-order reflections, dotted lines represent second-order reflections.

Figure 4.3 shows an impulse response measured in a very reverberant environment, precisely a chatedral. Note that, apart from the initial spikes, the overall decay is fairly exponential. One could measure that the reverberation time is in this case $T_r \sim 3.82$ s, which is a quite large value as one would expect in a cathedral.

4.2.1.4 Geometrical room acoustics

Few results of practical use are obtained from manipulation of the D'Alembert equation as in the previous sections, especially when we consider rooms of arbitrary shapes instead of parallelepipeds: in that case even the computation of a single normal mode can become extremely difficult

An alternative description of the acoustical properties of a room can be employed if we consider extremely high acoustic frequencies. In this limit situation, the concept of sound waves can be replaced by the concept of *acoustic rays*. By sound ray, we mean a vanishingly small portion of a spherical wave emitted by a point source in a room. This ray has well-defined direction and velocity of propagation, and conveys a total energy which remains constant (provided that it propagates within an ideal medium with no losses).

This simplified description based on acoustic rays takes the name of *geometrical acoustics* and has strict similarities with geometrical optics, although typical wavelengths and propagation velocities are very different in the two cases. Note that the assumption of extremely high frequencies is practically met in many cases of interest in room acoustics: a frequency of 1 kHz corresponds to a wavelength of approximately 34 cm, which is one or two orders of magnitude smaller than typical linear dimensions of rooms, as well as typical distances traveled by sound waves in a room.

Similarly to an optic ray, an acoustic ray that strikes a plane surface is reflected according to the following principles: (a) the reflected ray remains in the plane identified by the incident ray and the normal to the surface, and (b) the angles of the incident and reflected rays with the normal are equal. Figure 4.4 shows a room with a non-trivial shape (something like an auditorium), in which we have positioned a sound source and a receiver. There are many paths from the source to the receiver, that can be characterized according to the number of reflections involved. The single source-receiver path with 0 reflections is the *direct sound*, and is followed by a small number of *first-order reflections* that involve one reflection on the room boundary, a larger number of *second-order reflections* that involve two reflections on the boundary, and so on. In Fig. 4.4 we have drawn two examples of first- and

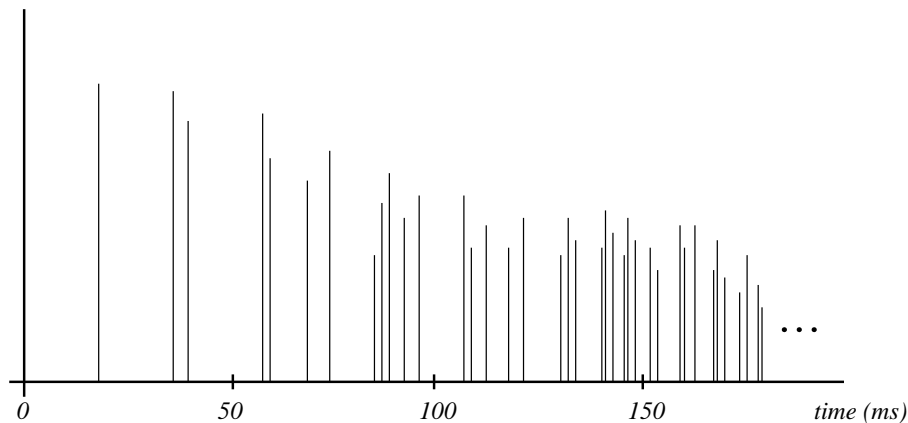


Figure 4.5: Schematical room response to an ideal impulse: the time axis is relative to the direct sound, which reaches the receiver at $t = 0$.

second-order reflections.

We now want to gain a qualitative description of a RIR using geometrical room acoustics. Assume that an ideal impulse shot from a point source reaches a receiver at time $t = 0$. Each reflected ray will then arrive with a certain time delay and also with a certain attenuation, which depends on the path length (absorption in the medium) and on the number of reflections (wall absorption). The first reflections are strong and sporadic, but the temporal density of reflections increases rapidly while the average reflection energy decays accordingly. A schematical reflection diagram is given in Fig. 4.5. Except for the first few isolated reflections, the weaker and denser reflections arriving at later times merge into what is perceived as reverberation. This description of the reverberation of a room as the temporal sum of reflected rays is complementary to the view of reverberation as the sum of free decaying normal modes.

We now want to derive an estimate of the temporal structure of reflections. To this end we employ the usual prototype room, i.e. the parallelepiped, and we introduce the concept of *image sources*.

If the reflecting surface is a plane the reflection of a sound ray can be simulated by constructing an *image source*. This process is illustrated in Fig. 4.6(a). Given a sound source A and a receiver B , the path of a reflected ray r from the wall to B is the same path of the direct ray r' emitted by the image source A' . The process can be iterated in order to take into account higher-order reflections, and results in the construction of a grid of image sources that replace the wall altogether.

Now suppose that at time $t = 0$ all the sources emit generate an impulse. During the time interval from t to $t + dt$, the impulses that reach a receiver located in the center of the room are those emitted by image sources whose distance from the receiver lies between cdt and $c(t + dt)$. These sources are located within the spherical shell with radius ct , thickness cdt , and volume $4\pi c^3 t^2 dt$ illustrated in Fig. 4.6(b). Therefore the volume V of an image room is contained in the spherical shell is $4\pi c^3 t^2 dt/V$ times, and if t is large enough (i.e. the reflection density is high enough) we can assume that this number coincides with the number dN_r of image sources contained in the shell. In conclusion, the temporal density of reflections arriving at time t is

$$D_r(t) = \frac{dN_r}{dt}(t) = 4\pi \frac{c^3 t^2}{V}. \quad (4.20)$$

One could show that this result applies not only to a parallelepiped but to rooms of arbitrary shapes.

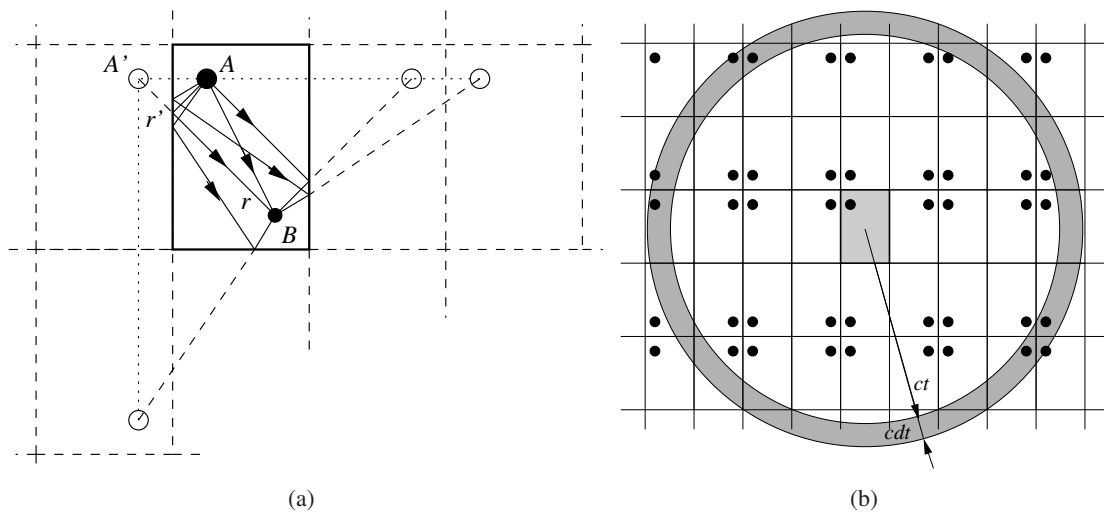


Figure 4.6: Estimation of temporal reflection density through the image source method; (a) construction of two first-order and two second-order reflections, and (b) estimation of acoustic rays reaching a receiver within the time interval $(t, t + dt)$.

4.2.2 Perceptual reverberation parameters

In the previous section we have analyzed reverberation from a purely physical point of view. However in many applications it is important to correlate physical measurements to subjective judgements of acoustical quality, obtained from psychophysical experimentation. This is especially true in the domain of concert hall acoustics, where researchers have tried to isolate the objective parameters that are most relevant in determining the perception of acoustical quality of a hall. Subjective attributes are typically derived from perceptual experiments with musicians and listeners, who answer to detailed interviews, and subsequent comparison of the results with measured objective parameters.

In this section we enter, for the time in this book, the domain of psychoacoustics, and review some of the subjective attributes and objective measures most commonly used in establishing the acoustical quality of reverberant environments. The literature on this topic is vast and the terminology is not always fully consistent, thus we try to cluster together similar or equivalent concepts wherever possible.

Clearly the perceptual attributes of reverberation are of great importance also for the design of reverberation algorithms. The ultimate goal is to determine an orthogonal set of subjective attributes, using e.g. multidimensional scaling techniques, and then providing a reverberation algorithm with a set of knobs each of which controls a different perceptual attribute. A fundamental problem with this kind of approach is that the number of perceptual dimensions is not known *a priori*, and moreover it is hard to assign relevance to dimensions that are added.

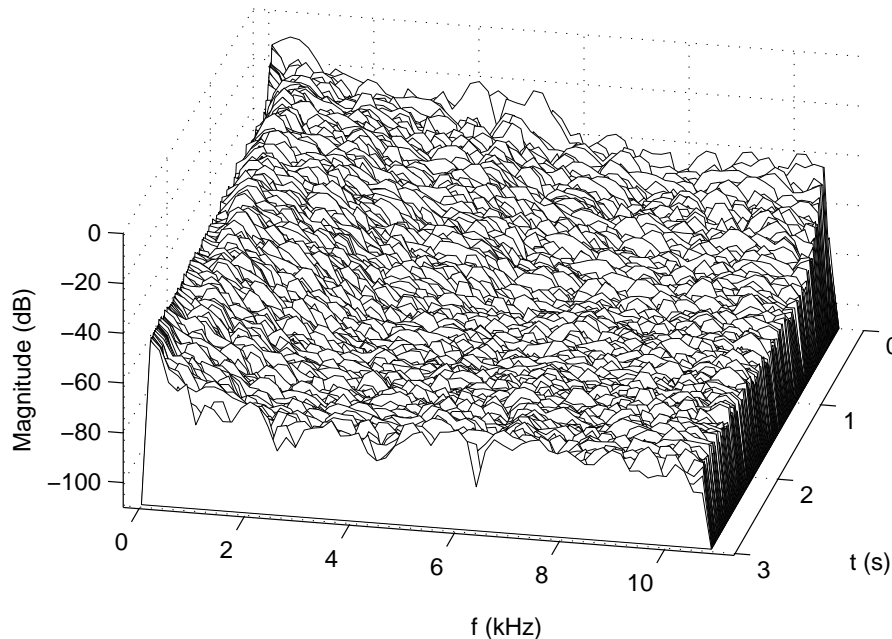


Figure 4.7: Waterfall representation of the impulse response of Fig. 4.3.

4.2.2.1 Reverberance

We have introduced the concept of reverberation time T_r in Sec. 4.2.1: it is the time required for the sound pressure to decay 60 dB.² The reverberation time is one of the most important parameters for the perception of the *reverberance* of an environment, i.e. the property of the environment of adding fullness and loudness to a dry sound, and of giving the listener a sense of being enveloped by the sound. Some researchers and musicians use the term *liveness* to refer to a similar concept, and by contrast call an environment that is not reverberant dry, or dead.

The time T_r defined in Eq. (4.19) is the overall decay time of a RIR, estimated assuming uniform damping. However we have already mentioned that the reverberation time is in general a function of frequency, because absorptive properties of materials vary with frequency and specifically absorption is typically higher at higher frequencies. A confirmation of this is given in Figure 4.7, which shows a waterfall representation of a RIR: one can see that each frequency bin decays with a different rate. This dependence of T_r on frequency is also important perceptually. In general the mid-frequency reverberation time can be considered to be the best measure of the overall reverberant characteristics of a room.

Clearly the audibility of reverberation depends greatly on the sound source. If one thinks at music or speech, the early portion of the reverberant decay contributes more to the perception of reverberance than does late reverberation, because it is audible during pauses and gaps between notes, syllables, and words. For this reason an *early decay time (EDT)* parameter is also used as a complementary measure of reverberance. The EDT is defined as the time required for the sound pressure to decay from 0 to -10 dB, multiplied by a factor of 6. The multiplicative factor merely serves to facilitate comparison with T_r .

²An alternative and more precise definition commonly used in the domain of concert hall acoustics is the following: T_r is the time required for the sound pressure to decay from -5 to -35 dB, multiplied by a factor of 2.

One might wonder what is the “optimal” reverberation time for a reverberant environment. The answer to this question depends first of all on the source signal: in the case of speech a relatively short T_r is generally preferred, while longer values are suitable for music. This can be expected, since when listening to speech we generally want to understand what the speaker is saying and thus we need to perceive each element of the sound signal. This is not the case for music, on the contrary reverberation can make a musical signal more pleasant by masking small imperfections and blending musical sounds. Given this remark, it is not surprising that reverberation times in (good) concert halls are usually in the range 1.8 to 2.2 s, while in opera houses values are usually in the range 0.9 to 1.5 s because the listener has to be able to enjoy the music as well as to understand the text. Note however that reverberation times of renowned opera theaters are more scattered than those of equally renowned concert halls.

4.2.2.2 Early reflections and spatial impression

The subjective attribute of *spatial impression* refers to the sense of a listener of being in close communication with the sound source and surrounded by the sound. Other terms that are often found in the literature and refer to a similar concept are spaciousness, envelopment, ambience, apparent source width. Subjective judgements about this property appear to be strongly correlated to the structure of the early reflections of the environment, with two elements being of specific importance.

A first commonly accepted result is that the degree of spatial impression depends on the difference in arrival times between the direct sound and the first reflection, which is called *initial time-delay gap* and is often indicated as t_I . A lack of early reflections (i.e. a long t_I) has the effect of making the sound source perceived as remote and disconnected from the listener, while a short t_I provides the desired sense of envelopment. Some studies suggest that a parameter t_I defined as above becomes useless if the first reflection is much weaker than the following ones.

A second physical property that correlates to spatial impression is the fraction of lateral energy to the total energy within the early reverberation: a significant amount of *lateral* early reflections, i.e. reflections coming from the sidewalls, provides the listener with the impression of being enveloped by the sound. A rough quantitative estimate of this property is the so-called *lateral energy fraction* LF_t . defined as

$$LF_t = \frac{\int_0^t h_{lat}^2(\tau) d\tau}{\int_0^t h^2(\tau) d\tau}, \quad (4.21)$$

where $h(t)$ is the room impulse response measured with an omnidirectional microphone while $h_{lat}(t)$ is the one measured with a dipole microphone (with null axis facing forward this captures lateral energy in the $\pm 20^\circ \pm 90^\circ$ range). A typical integration time is $t = 80$ ms.

The LF_t measure has been superseded by another parameter, called the *early interaural cross-correlation coefficient* $IACC_E$. Let us first define the interaural cross-correlation function $IACF(t)$ as

$$IACF(t) = \frac{\int_{t_1}^{t_2} h_L(\tau) h_R(\tau + t) d\tau}{\sqrt{\int_{t_1}^{t_2} h_L(\tau) d\tau \int_{t_1}^{t_2} h_R(\tau) d\tau}}, \quad (4.22)$$

where $h_{L,R}(t)$ are the impulse responses measured at the entrance of the left and right ear canals, respectively, with the listener facing the sound source. Such a measurement can be done using e.g. a so-called “dummy-head” (such as those described later on in Sec. 4.6.1). Therefore the $IACF(t)$ function is a *binaural* attribute of reverberation, unlike all of the parameters previously examined in this section which are monaural attributes.

The interaural cross-correlation coefficient $IACC$ is the maximum of this function over a range ± 1 ms:

$$IACC = \max_{t \in (-1,1) \cdot 10^{-3}} IACF(t). \quad (4.23)$$

In particular, if the integration times $t_1 = 0$, $t_2 = 80$ ms are used then the above equations provide the early interaural cross-correlation coefficient $IACC_E$. This is a measure of the similarity of the sound signals arriving at the two ears during the first 80 ms. If the sounds are equal then $IACC_E = 1$, while if they are two independent random signals then $IACC_E = 0$. The $IACC_E$ parameter is a measure of spatial impression because it scales with the fraction of lateral early reflections arriving at the ears: as the number of reflections from outside the median plane increases, the $IACF(t)$ function broadens and consequently $IACC_E$ takes smaller values.

In concert halls initial time-delay gap t_I and the amount of lateral energy are correlated parameters. Measures of t_I in real concert halls show a high correlation of this parameter with the hall width: in a narrow hall it can be shorter than 30 ms, while in a wide hall it can be longer than 50 ms. On the other hand, the hall width is clearly correlated with the fraction of lateral energy arriving at the listener, which will increase as the hall narrows. It is a common finding in the literature of concert hall acoustics that subjective rankings of the acoustic quality of halls scale with their width.

As a final remark, it has to be noted that the subjective attribute of spatial impression is largely independent of the reverberation time: halls with similar T_r values but different t_I and $IACC_E$ values will be perceived to be very different from each other. This finding supports the commonly accepted assumption that early reflections and late reverberation play rather separate roles in the perception of reverberant properties of an environment.

4.2.2.3 Clarity

The subjective attribute of *clarity* refers to the “transparency” of a reverberant environment. If the source signal is music, then clarity is associated to the ability of a listener to perceive musical details, while if the source signal is speech then clarity correlates to speech intelligibility. An alternative term which is sometimes found in the literature is that of distinctness.

Single reflections of a reverberant environment are not perceived as individual events, except for exceptional (and generally undesirable) cases. Roughly speaking, early reflections have the effect of making the sound source appear more extended and to increase the apparent loudness of the direct sound. On the contrary, reflections arriving with longer delays are considered to be detrimental for the transmission of information, since they cause different portions of the direct sound signal to merge.

A quantitative measure of clarity is the *clarity index*, or early-to-reverberant energy ratio C_t :

$$C_t = 10 \log_{10} \left(\frac{\int_0^t h^2(\tau) d\tau}{\int_t^\infty h^2(\tau) d\tau} \right), \quad (4.24)$$

measured in dB. The integration time t is ideally the time instant where late reverberation starts, and is typically selected to be $t = 80$ ms. Thus C_t is a measure of early to late energy ratio.

It is sometimes recommended that $C_t|_{t=0.008}$ for concert halls takes values in the range of -2 to $+1$ dB. Note however that this parameter is not an independent measurable quality, since it correlates to the initial time-delay gap t_I and also on the early decay time EDT . Therefore, subjectively “good” values of the clarity index will also depend on t_I and EDT values. In other words, the subjective attribute of clarity is not orthogonal to reverberance and spaciousness.

Note also that C_t is strongly dependent on the distance between source and listener: the direct sound falls off 6 dBs for each distance doubling, whereas the reverberant level remains approximately

constant throughout the room. For this reason, the ratio of direct to reverberated energy is one of the most important cues for the perception of distance, as we will see in Sec. 4.5.

A second objective parameter that relates to the subjective attribute of clarity is the *center time* t_s , defined as the center of gravity time of the sound field:

$$t_s = \frac{\int_0^\infty \tau \cdot h^2(\tau) d\tau}{\int_0^\infty h^2(\tau) d\tau}. \quad (4.25)$$

Obviously a single reflection with a given strength will contribute the more to t_s the longer it is delayed with respect to the direct sound. Therefore high clarity is associated to low values of t_s . It has to be noted however that many studies report a high correlation of t_s with C_t , in the range $50 < t < 80$ ms. Therefore this parameter does not add new information with respect to the clarity index.

4.2.2.4 Other perceptually relevant parameters

The physical concept of *diffusion*, which we have examined previously, has a direct perceptual counterpart. If one listen to music in a rectangular hall with perfectly flat sidewalls, the sound takes on an undesirable harsh character. In order to produce the effect of a mellower sound and to increase spaciousness during late reverberation, diffusion should be physically realized at fine and large scales. A commonly accepted measure of diffusion is the *late interaural cross-correlation coefficient* $IACC_L$. This is defined from Eqs. (4.22, 4.23) using integration times $t_1 = 80$ ms and $t_2 = 3$ s, i.e. by estimating the IACF function in the late reverberation portion. Similarly to the $IACC_E$ parameter, $IACC_L$ is a *binaural* attribute of reverberation. It provides a measure of the correlation of the signals at the two ears during late reverberation.

Loudness (or *strength*) is often mentioned as a relevant subjective attribute. Of course the overall loudness depends on the power output of the sound source and not only on the reverberation of the environment. Nonetheless it is useful to introduce a measure of loudness of the environment, which is normalized with respect to the the source power. Such a measure can be used e.g. as a complementary parameter to the clarity index (see Eq. (4.24) above), since high clarity is of no use if the sound cannot be heard at proper loudness.

A normalized measure of the environment loudness is achieved by the following quantity, sometimes called *strength index* G :

$$G = 10 \log_{10} \left(\frac{\int_0^\infty h^2(\tau) d\tau}{\int_0^\infty h_0^2(\tau) d\tau} \right), \quad (4.26)$$

where $h(t)$ is as usual the room impulse response and $h_0(t)$ is the response to the same non-directional impulse measured in an anechoic environment at a distance of 10 m. Note however that subjective loudness increases with reverberation time and is affected by the structure of early reflections. Therefore G is not an independent correlate of loudness.

Finally, the most elusive subjective attributes are those related to timbral qualities of a reverberant environment. Roughly speaking, many of the attributes in this family are related to the frequency-dependent shape of the reverberation time. One such attribute is *warmth*, or sometimes *timbre*, which characterizes the musicians' judgement of "richness in bass". This attribute correlates with the variation of the reverberation time in the low- and mid-frequency range: as an example, a quantitative measure of warmth can be the ratio of the average T_r in the range 250 – 500 Hz to that in the range 500 – 1000 Hz, or alternatively the slope of a linear interpolation of the EDT function in the range 125 – 2000 Hz. Other timbre-related attributes are *heaviness* and *liveness*, which roughly relate to low-frequency and high-frequency variations of the reverberation time, respectively.

A compact representation of the perceptually relevant features of a room impulse response is the so-called *Energy Decay Relief (EDR)* function, which is a time-frequency representation of the reverberation energy. Let us first define a new function called, *Energy Decay Curve (EDC)*, as follows:

$$EDC(t) = \int_t^{\infty} h^2(\tau) d\tau, \quad (4.27)$$

where $h(t)$, is a RIR. This integral is often called the Schroeder integral. The value $EDC(t)$ provides a measure of the reverberation energy that is left in the RIR at time t . Using this function we can introduce the Energy Decay Relief function as follows: given a RIR $h(t)$, this is bandpass filtered into a number N of frequency bands, and the Schroeder integral of each of the bandpassed responses $h_i(t)$ ($i = 1 \dots N$) is computed. The resulting function $EDR(t, \omega)$ can be displayed as a surface in the 3-D space. The section $EDR(0, \omega)$ provides the power gain as a function of frequency. A section $EDR(t, \omega_0)$ shows the energy decay curve for a given frequency ω_0 .

M-4.2

Write a function that computes the EDR given a RIR.

The time-frequency EDR function can be parametrized through two functions of frequency only. The first one is $T_r(\omega)$, the frequency-dependent reverberation time. The second one is the *frequency response envelope*, $G(\omega)$. This latter function is constructed by backward interpolating up to $t = 0$ the exponential decay time. For an ideally diffuse reverberation that decays exponentially, one has the equality $G(\omega) = EDR(0, \omega)$ and G coincides with the power gain of the room. In non-ideal cases, $G(\omega)$ only represent a “conceptual” $EDR(0, \omega)$ of the late reverberation, and the parametrization through $T_r(\omega)$ and $G(\omega)$ is only valid for the late portion of $EDR(t, \omega)$.

The EDR function is sometimes regarded as a perceptual “signature” of a room impulse responses, meaning with this that a large number of objective measures of independent perceptual factors can be categorized as energy ratios or energy decay slopes computed within different time-frequency regions of the EDR function.

4.3 Algorithms for synthetic reverberation: the perceptual approach

Even if we have enough computer power to compute convolutions by long impulse responses in real time, there are still serious reasons to prefer reverberation algorithms based on feedback delay networks in many practical contexts. It is not easy to modify a room impulse response to reflect some of the room attributes. If the impulse response has been synthesized with some spatial rendering algorithm, such as ray tracing, these manipulations can be operated at the level of room description, and the coefficients of the room impulse response transmitted to the real-time convolver. However, continuous variations of the room impulse response are rendered more easily using a model of reverberation operating on a sample-by-sample basis. In the second half of the twentieth century, several engineers and acousticians tried to invent electronic devices capable of simulating the long-term effects of sound propagation in enclosures.

Sticking to the terminology introduced in the previous chapters, we can say that these are *signal-based* models, since they aim at reproducing realistic impulse response signals, with no attempt to model the underlying physical phenomena. Often used terminology in this context: *perceptual* approach (add “ambience” in a dry recording) vs. *physical* approach (faithfully simulate the acoustics of a room or hall).

In the remainder of this section we address the perceptual approach.



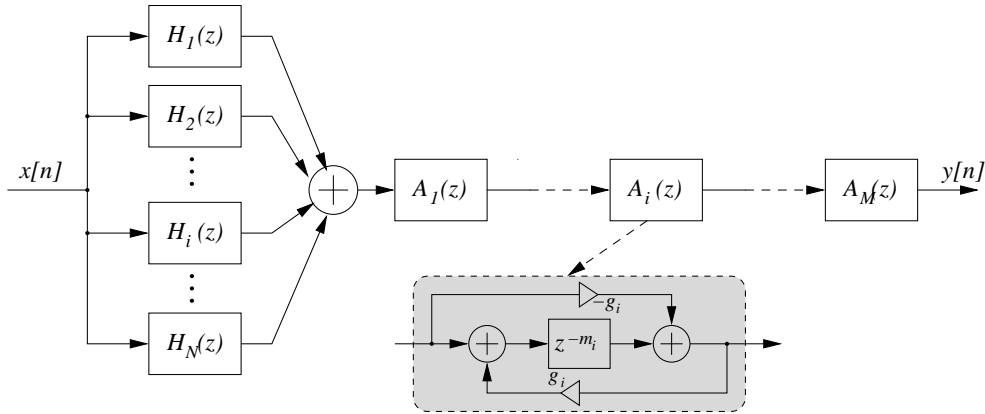


Figure 4.8: Block scheme of a reverberator based on comb filters (the H_i blocks) and all-pass comb filters (the A_i blocks). The internal structure of the A_i filters is shown in the grey box.

4.3.1 Approximating late reverberation

We have previously seen that a RIR can be seen as made of two components, early reflections and late reverberation. In this section we address the modeling of late reverberation, and we postpone early reflection modeling to a Sec. 4.3.2.

4.3.1.1 Recirculating delays

The two main computational structures that can be used for the inexpensive simulation of complex patterns of echoes associated to late reverberation are the recursive comb filter $H(z)$ (see Karplus-Strong in Ch. *Sound modeling: source based approaches*) and the so-called *all-pass comb filter* $A(z)$

$$H(z) = \frac{z^{-m}}{1 - gz^{-m}}, \quad A(z) = \frac{z^{-m} - g}{1 - gz^{-m}}. \quad (4.28)$$

It is easily seen that $A(z)$ is an all-pass structure, since each of the m poles is the reciprocal of one of the m zeros and the amplitude response $|A(z)|$ is therefore flat. For $m = 1$ the structure reduces to the first-order all-pass filter examined in Ch. *Sound modeling: source based approaches*. The (positive) gain g in $A(z)$ has to be less than unity in order to ensure stability.

Figure 4.8 depicts a reverberator constructed using comb-filters and all-pass comb filters, together with a realization of the all-pass comb (see the grey box). The general idea behind this structure is the following. First, the parallel combination of comb filters generates a frequency response that contains peaks contributed by each comb. In theory we can obtain an arbitrary modal density by using a sufficiently large number N of comb filters. Second, the series combination of all-pass combs that receives the output of the parallel combination of combs has the effect of dramatically increasing the temporal density of reflections, because each echo generated by $A_i(z)$ will create a set of echoes in $A_{i+1}(z)$. Again, an arbitrarily high reflection density can be in principle obtained by using a sufficiently large number M of all-pass combs.

4.3.1.2 Tuning the parameters

The choice of a proper set of parameter values is critical in order to obtain convincing results. In the remainder of this section we provide a list of commonly accepted guidelines. The sample delays m_i of

the combs should be mutually coprime (or incommensurate), in order to reduce the superimposition of echoes in the impulse response, thus maximizing the modal density and reducing the so-called flutter echoes. The gains g_i of the combs can be chosen as functions of the sample delays m_i , given a desired reverberation time T_r . It is easy to prove that the following equation holds for the reverberation time of a single comb:

$$\frac{F_s \cdot 20 \log_{10}(g_i)}{m_i} = -\frac{60}{T_r} \Rightarrow g_i = 10^{-3 \frac{m_i}{F_s T_r}}. \quad (4.29)$$

Note that this choice ensures that the pole moduli $m_i \sqrt{g_i} = 10^{-3 \frac{1}{F_s T_r}}$ have the same value for all the combs. If this condition was not verified, then the poles with largest moduli would resonate longer and would add an undesired tonal coloration in the late decay.

A quantitative estimate of the modal density provided by the parallel comb structure can be easily obtained. If the m_i 's of the combs are mutually coprime, then the modal density D_f (which is number of frequency peaks per Hz) can be estimated as

$$D_f = \sum_{i=1}^N \frac{m_i}{F_s} = \frac{N \bar{m}}{F_s}, \quad (4.30)$$

where \bar{m} is the mean sample delay length. Note that this modal density is constant for all frequencies, unlike in real rooms (see Eq. (4.9)). A too low D_f can introduce audible beating between two neighboring modes, especially in response to narrowband signals. In order to avoid this effect, a good rule of thumb is to choose the m_i 's such that $D_f \geq T_r$: this ensures that the average beat period is at least equal to the reverberation time.

In a similar way we can estimate quantitatively the temporal reflection density provided by the parallel combination of combs: each filter outputs one echo every m_i/F_s seconds, therefore the combined reflection density (number of reflections per second) is

$$D_r = \sum_{i=1}^N \frac{F_s}{m_i} \approx \frac{N F_s}{\bar{m}}, \quad (4.31)$$

where the last approximation only holds when the m_i are similar. Again, the reflection density is constant as a function of time, unlike real rooms (see Eq. (4.20)). A value $D_r = 10^3$ is sometimes considered to be sufficient to sound indistinguishable from diffuse reverberation, although higher values (e.g. $D_r = 10^4$) are preferable.

From the two estimates (4.30) and (4.31) provide an estimate of the number of comb filters needed in order to achieve desired modal and reflection densities:

$$N = \sqrt{D_f D_r}. \quad (4.32)$$

Note however that this estimate does not consider the effect of the cascaded series of all-pass comb filters A_i : as already mentioned, the A_i provide a dramatic increase of the reflection density and allow to a number N of comb filters that is smaller than the one estimated from Eq. (4.32).

M-4.3

Realize the reverberant structure of Fig. 4.8. The reverberator can be tried e.g. with $N = 4$, $M = 2$, and with the following settings: time delays m_i/F_s ($i = 1 \dots 4$) of the comb filters distributed 30 and 45 ms, time delays m_i/F_s ($i = 5, 6$) of the all-pass combs between 1.7 and 5 ms, modal density $D_f = 1000$, gains of the all-pass combs $g_i = 0.7$ ($i = 5, 6$). With these settings the structure is known as Schroeder reverberator (see the bibliography).

4.3.2 Improved structures

The reverberator discussed in the previous section sounds reasonably well especially for short reverberation times and low reverberation levels. For different settings however it suffers from a number of problems. First, the reverberation is not dense enough at the beginning, resulting in a “grainy” sound quality (especially in response to impulsive sounds). Second, the late reverberation tends to exhibit an already mentioned “fluttering” effect. Third, especially for long reverberation times a “ringing” effect can be heard, which gives an undesired metallic quality to the reverberation. Fourth, the modal density is not sufficiently large and, as already mentioned, does not increase with frequency. Fifth, the reverberation time T_r does not depend on frequency, unlike in real rooms (see Sec. 4.2.2 and the EDR function there discussed). Finally, no modeling of early reflection is provided.

4.3.2.1 Low-pass combs

A first obvious way of improving the modal density is to increase the number of comb filters in parallel, especially when long reverberation times need to be simulated. A second more substantial improvement amounts to employ, in place of comb filters, a *low-pass comb* filter, where a low-pass filter H_{lp} is inserted in the feedback loop of the comb filter instead of a scalar gain. The purpose of this modification is to simulate the attenuation effects of higher frequencies, due to air viscosity, heat conduction, and energy losses at reflection. As a result, the reverberation time decreases at higher frequencies and makes the reverberation sound more realistic. In addition, the response to impulsive sounds is also improved, due to the smoothing effect of the low-pass filtering.

If a simple one-pole low-pass filter H_{lp} is used, then the low-pass comb filter is given as

$$H(z) = \frac{z^{-m}}{1 - H_{lp}(z)z^{-m}}, \quad \text{with } H_{lp}(z) = \frac{g_1}{1 - g_2z^{-1}}. \quad (4.33)$$

One can easily verify that in order for $H(z)$ to be stable the condition $g_1/(1 - g_2) < 1$ must hold. Note that we have already introduced the low-pass comb filter for the Karplus-Strong algorithm in Ch. *Sound modeling: source based approaches*, although here we are using a different low-pass filter H_{lp} .

Coefficients of the low-pass combs: g_1 can be determined as a function of the delay length and the desired T_r , as explained in the previous section. The g_2 coefficient can also be related with decay time at a specific frequency or fine tuned by direct experimentation.

M-4.4

Realize the reverberant structure of Fig. 4.8 using low-pass comb filters of the form (4.33). The reverberator can be tried e.g. with $N = 6$, $M = 1$, and with the following settings: time delays m_i/F_s ($i = 1 \dots 6$) of the comb filters distributed between and ms, time delay of the all-pass comb $m_7/F_s = 6$ ms, modal density $D_f = \dots\dots\dots$, gain of the all-pass comb $g_7 = 0.7$.

4.3.2.2 Nested all-pass filters

Despite the improvements provided by this latter reverberator, some problems remain. First, it is not possible to tune the reverberator to a desired $T_r(\omega)$ function. Second, the modal density is still constant with respect to frequency and the ringing quality and the fluttering effect in the reverberation tail remain, although reduced to some extent. In order to overcome this problems some researchers have proposed reverberators with entirely different structures than the one shown in Fig. 4.8. One such structure is shown in Fig. 4.9.

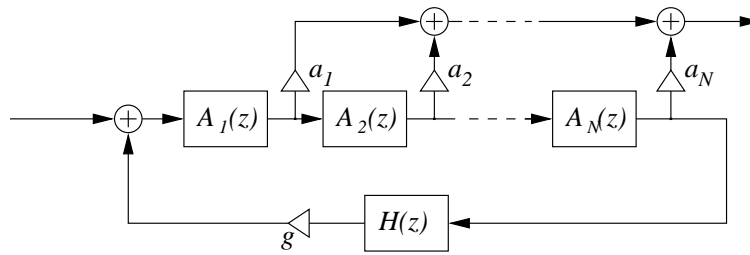


Figure 4.9: A reverberator constructed with a series connection of all-pass filters and a low-pass filter in feedback.

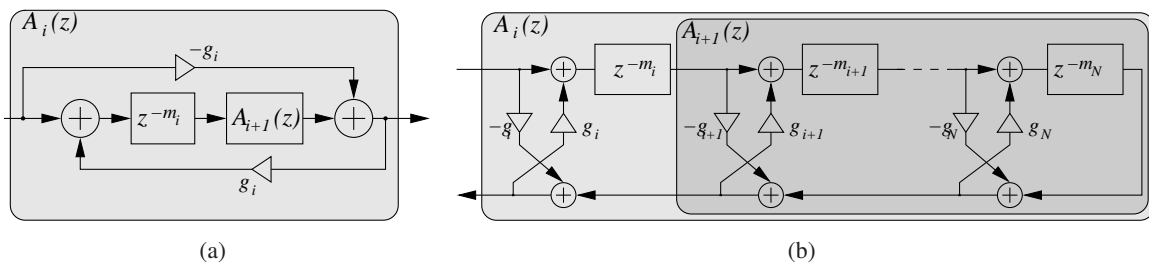


Figure 4.10: Nested all-pass filters; (a) generalization of an all-pass structure (see Fig. 4.8), and (b) realization by means of a lattice structure.

As before, the cascaded all-pass filters $A_i(z)$ have the effect provide a high temporal density of reflections, because each echo generated by a filter will create a set of echoes in the following one. In this case however, the output of the last all-pass filter is recirculated to the series connection through a low-pass filter $H(z)$ and an attenuating gain g . The resulting system is stable, if the condition $|gH(e^{j\omega})| < 1 \forall \omega$ is verified.

The low-pass filter $H(z)$ can be interpreted as simulating frequency-dependent absorptive losses, and the gain g provides control over the reverberation time. An important effect of this outer feedback loop is that the characteristic metallic sound of the series all-pass is drastically reduced. Another peculiarity of this structure is that the output is constructed as a linear combination of the all-pass outputs. Since each each tap outputs a different response shape, the coefficients a_i can be adjusted in order to shape the amplitude envelope of the reverberant decay.

A final remark concerns the possibility of generating a reflection density that increases with time, as in real rooms. A structure that achieves this goal is a *nested all-pass filter* $A_1(z)$, which can be defined recursively as follows:

$$\begin{aligned} A_{N+1}(z) &= 1, \\ A_i(z) &= \frac{z^{-m_i} A_{i+1}(z) - g}{1 - g z^{-m_i} A_{i+1}(z)}, \text{ for } i = 1 \dots N. \end{aligned} \quad (4.34)$$

Figure 4.10(a) shows that this structure can be seen as a generalization of the all-pass comb, in which part of the delay line has been substituted by an all-pass filter. Figure 4.10(b) explodes this structure into a realization based on a lattice structure. It is easy to verify that each of the nested filters $A_i(z)$ are all-pass. Moreover, Fig. 4.10(a) shows that each echo generated by the inner all-pass $A_{i+1}(z)$

is recirculated to itself through the outer feedback path of $A_i(z)$: this intuitively explains why this structure provides a reflection density that increases with time.

M-4.5

Realize the reverberant structure of Fig. 4.9 using nested all-pass filters of the form (4.34).

4.3.2.3 Adding early reflections

So far we have only examined algorithms for the simulation of the late, diffuse reverberation. We have not paid any attention to the simulation of early reflections, which have great importance in the perception of the acoustic space. In this section we address this point.

As previously discussed, the early response of a room is sparsely populated with attenuated impulses. These can be straightforwardly simulated using a direct-form FIR filter that reproduces these impulses explicitly and accurately. For the determination of the filter parameters, a good rule of thumb is to apply to the early reflections delays the same criterion of “mutually-primeness” used before for the comb delays. A better strategy is to derive the parameters from some geometric modeling technique, e.g. the source image method discussed in Sec. 4.2.1.

Figure 4.11 shows an example of early reflection modeling, in which the FIR filter has been realized using a direct form structure. The early reflection filter has to be connected to a late reverberation block: Fig. 4.11(a) and 4.11(b) show two possible connections. In Fig. 4.11(a) the late reverberator receives the delayed input signal, and therefore the FIR response will always occur before the late response in the final output. Figure 4.11(b) shows a more complex coupling between the two blocks. In this case the late reverberator is driven by the output of the FIR filter, with the result of increasing the reflection density in the late reverberation. Moreover, additional control parameters are available: the gain g can be adjusted in order to balance the early/late reverberation ratio, while the delays D_1 , D_2 can be tuned so that the start of the late reverberator output coincides with the last pulse output from the FIR filter, thus avoiding undesired gaps in the overall response.

M-4.6

Realize the reverberator depicted in Fig. 4.11(a), where the early reflection FIR filter has to be coupled to one of the late reverberation structures discussed in the previous sections.

M-4.7

Realize the reverberator depicted in Fig. 4.11(b), where the early reflection FIR filter has to be coupled to one of the late reverberation structures discussed in the previous sections. Compare the resulting impulse responses with the ones obtained from M-4.6.

In order to improve the quality of the FIR structure described above, one has to include some form of low-pass filtering that models frequency dependent losses. One possibility is to substitute each of the gains a_i with a low-pass filter, composed by considering the history of reflections for each echo. Early reflections are not perceived as individual events however, therefore it is not necessary to model accurately the spectral content of each single reflection. A cheaper, and often satisfactory, choice is to sum sets of reflections together and to filter them through the same low-pass.

As a final remark to this section, one should note that early reflections are extremely important for the formation of spatial impression: an echo reaches the two ears with different intensities and at different times, because of the shadowing effect of the head and because of the different distance traveled. For this reason early reverberation is most effective if it is presented *binaurally*, i.e. by taking into account these effects and presenting different echoes to the two ears (e.g. via headphones). This

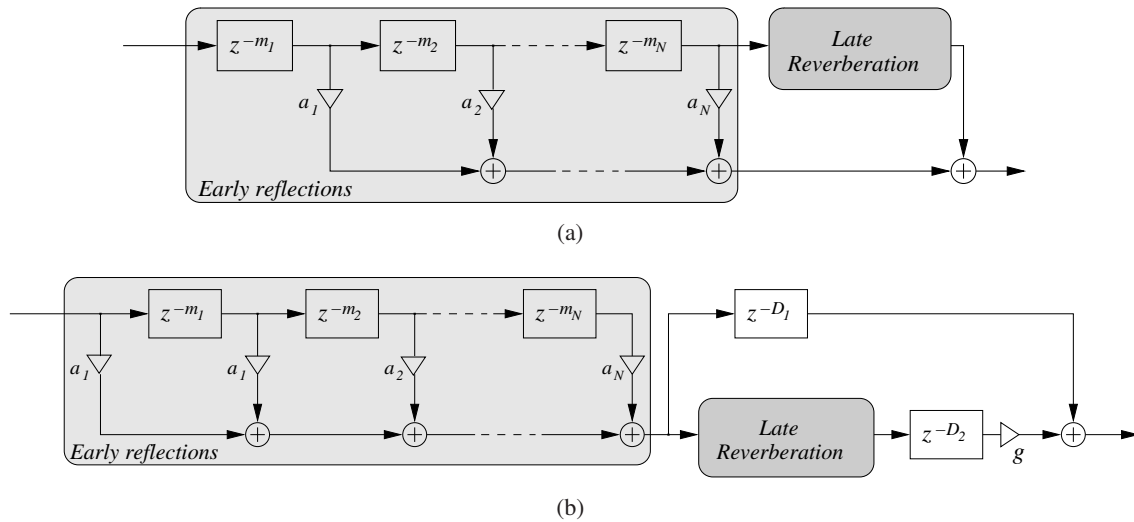


Figure 4.11: Two realizations of a reverberator with early reflections; (a) late reverberation block receiving the delayed input signal, and (b) late reverberation block receiving the output of the early reverberation FIR filter, with additional control parameters D_1 , D_2 , g . The late reverberation block can be one of the structures examined in the previous sections.

consideration anticipates the subject of Sec. 4.6, where we will address the topic of rendering the location in space of a sound source. At this point it worth mentioning that if no binaural processing is done the addition of early reflections can in certain cases deteriorate the quality of a reverberator, as they cause tonal coloration of the sound without producing spatial impression.

4.4 Multidimensional reverberation structures

4.4.1 Feedback delay networks

4.4.1.1 A n-D generalization of the recursive comb filter

In the previous section we have seen that the recursive comb filter of Eq. (4.28) has been extensively used as the main building block of perceptual reverberators, as an inexpensive way to generate patterns of resonances. Now the question is: can we generalize the comb structure in order to achieve higher modal densities? The filter structure depicted in Fig. 4.12 provides a first answer. First, it is easily seen to be a vector generalization of the recursive comb filter, as it reduces to a parallel combination of ordinary comb filters when the feedback matrix $\mathbf{A} = [a_{ij}]$ is diagonal. Second, and more interesting, it recirculates the output of the i th delay line to the input of the j th delay line, for every non-null element a_{ij} . This observation gives the intuition that when \mathbf{A} is non-diagonal this structure is capable of much higher modal densities than a simple parallel of comb filters.

The generalization extend also to stability conditions. While the comb filter of Eq. (4.28) is stable if $|g| < 1$, the multidimensional structure of Fig. 4.12 is stable if $\|\mathbf{A}\|_2 < 1$, where $\|\cdot\|_2$ is the spectral norm of a matrix³ This can be easily verified by applying the conditions for Lyapunov

³The matrix norm corresponding to any vector norm $\|\cdot\|$ may be defined for any matrix \mathbf{A} as $\|\mathbf{A}\| = \max_{\mathbf{x} \neq 0} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|}$.

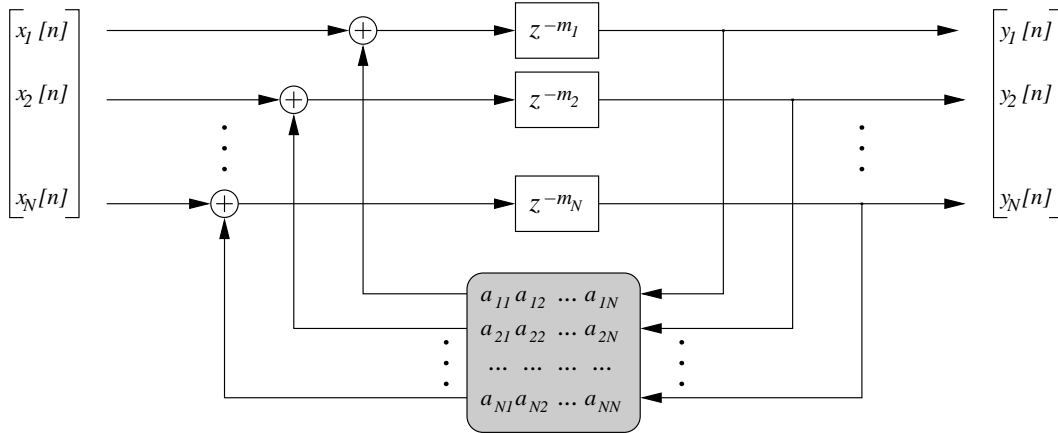


Figure 4.12: ...

stability, i.e. that the output $\mathbf{y}[n]$ decreases in time when the input signal \mathbf{x} is zero:

$$\|\mathbf{y}[n-1]\|_2 > \|\mathbf{y}[n]\|_2 = \left\| \mathbf{A} \begin{bmatrix} y_1[n-M_1] \\ \vdots \\ y_N[n-M_N] \end{bmatrix} \right\|_2 \quad (4.35)$$

where the first inequality is the Lyapunov stability condition, and the second equality holds for the block scheme of Fig. 4.12. Therefore stability is guaranteed whenever the feedback matrix satisfies

$$\|\mathbf{A}\mathbf{y}\|_2 < \|\mathbf{y}\|_2 \quad \forall \mathbf{y}. \quad (4.36)$$

In other words, a sufficient condition for stability is that the feedback matrix decreases the \mathcal{L}^2 norm of its input vector. Since in general $\|\mathbf{A}\mathbf{y}\|_2 < \|\mathbf{A}\|_2 \cdot \|\mathbf{y}\|_2$, we conclude that stability is guaranteed for $\|\mathbf{A}\|_2 < 1$.

A class of matrices that satisfy the stability condition is

$$\mathbf{A} = \mathbf{\Gamma}\mathbf{Q}, \quad \text{where} \quad \mathbf{\Gamma} = \begin{bmatrix} g_1 & 0 & \cdots & 0 \\ 0 & g_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & 0 \\ 0 & 0 & \cdots & g_N \end{bmatrix}, \quad |g_i| < 1, \quad (4.37)$$

and where \mathbf{Q} is an orthogonal matrix. Recall that (1) the spectral norm $\|\mathbf{A}\|_2$ is the square root of the largest eigenvalue of $\mathbf{A}\mathbf{A}^T$, and that (2) by definition \mathbf{Q} is orthogonal if and only if $\mathbf{Q}\mathbf{Q}^T = \mathbf{I}$. Then $\|\mathbf{A}\|_2 = \|\mathbf{\Gamma}\mathbf{Q}\| = \max_i |g_i|$.

The above analysis justify the use of the structure of Fig. 4.12 as a multichannel reverberator in which N input signals $\mathbf{x}[n]$ (or N replicas of a single input signal $x[n]$) produce N outputs $\mathbf{y}[n]$ that are approximately mutually incoherent and thus can be used in a N -channel loudspeaker system to render a diffuse soundfield. A possible choice for the matrix \mathbf{A} is

$$\mathbf{A} = g \frac{1}{\sqrt{2}} \begin{bmatrix} 0 & 1 & 1 & 0 \\ -1 & 0 & 0 & -1 \\ 1 & 0 & 0 & -1 \\ 0 & 1 & -1 & 0 \end{bmatrix}, \quad |g| < 1, \quad (4.38)$$

The spectral norm $\|\cdot\|_2$ is the matrix norm induced by the \mathcal{L}^2 vector norm.

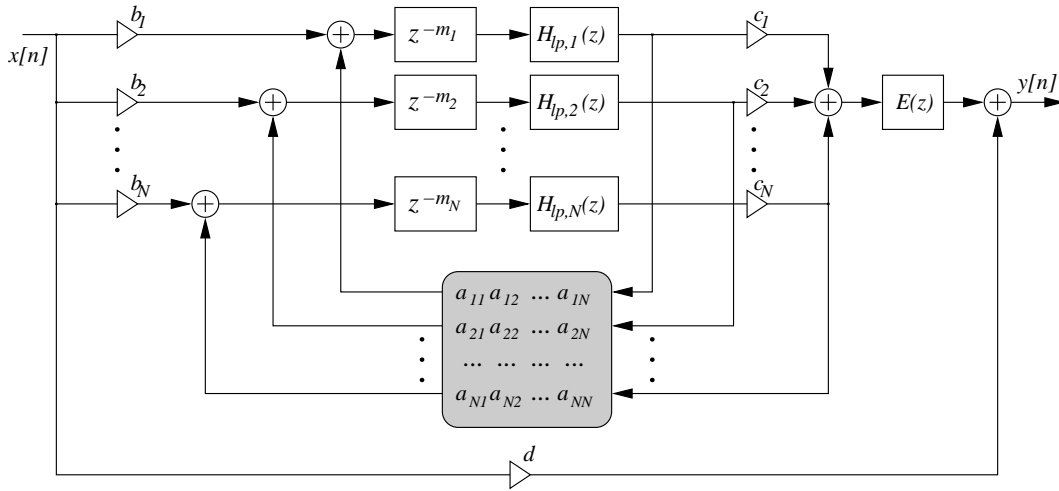


Figure 4.13: A Feedback Delay Network structure for artificial reverberation.

which is immediately seen to belong to the class (4.37).

4.4.1.2 A general FDN reverberators

The “vector comb filter” that we have analyzed in the previous section is an example of a class of filter networks, known as *Feedback Delay Networks (FDNs)*. Figure 4.13 shows a more general FDN structure for artificial reverberation, that extends in many ways the one depicted in Fig. 4.12. First, it is a Single-Input, Single-Output structure which uses two $N \times 1$ vectors $\mathbf{b} = [b_i]$ and $\mathbf{c} = [c_i]$ to split the input into N channels and to combine the N outputs in one channel. Second, low-pass filters $H_{lp,i}(z)$ are cascaded to the delay lines. Third, the final output y is corrected with an additional filter $E(z)$ plus an additive term dx . The transfer function of the system is almost immediately found to be:

$$\frac{Y(z)}{X(z)} = \mathbf{c}^T \left\{ [\mathbf{I} - \mathbf{D}(z)\mathbf{A}]^{-1} \mathbf{D}(z) \right\} \mathbf{b} \cdot E(z) + d = \mathbf{c}^T [\mathbf{D}(z^{-1}) - \mathbf{A}]^{-1} \mathbf{b} \cdot E(z) + d, \quad (4.39)$$

where $\mathbf{A} = [a_{ij}]$ is the *feedback matrix* of the system, and

$$\mathbf{D}(z) = \begin{bmatrix} z^{-m_1} H_{lp,1}(z) & 0 & \cdots & 0 \\ 0 & z^{-m_2} H_{lp,2}(z) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & z^{-m_N} H_{lp,N}(z) \end{bmatrix}$$

is the *delay matrix* of the system. We shall see that this structure allows to orthogonalize to a great extent the reverberation parameters, as the various blocks can be independently tuned to fit desired values of different reverberation parameters.

M-4.8

Realize the reverberant structure of Fig. 4.13. With $N = 4$, and with the matrix given in Eq. (4.38), the structure of Fig. 4.12 is a special case of this.

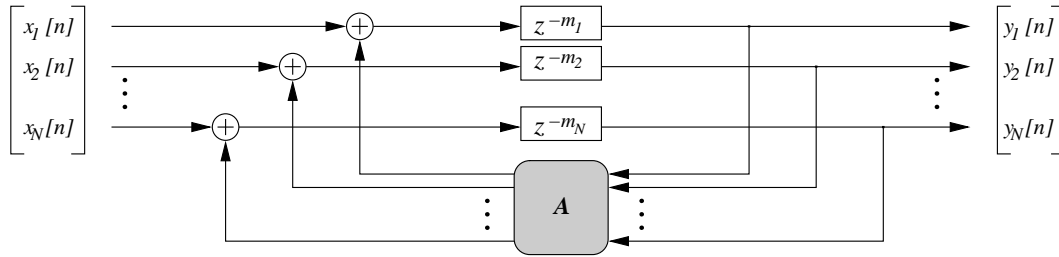


Figure 4.14: Lossless prototype network associated to the Feedback Delay Network of Fig. 4.13.

Since an “ideal” late reverberation impulse response should resemble exponentially decaying noise, it is useful to start designing a lossless reverberator (with infinite reverberation time) and work on making it a good noise generator. Once this *lossless prototype* has been designed, one can work on obtaining the desired reverberation time in each frequency band. We associate to the FDN of Fig. 4.13 the lossless prototype of Fig. 4.14.

What does the losslessness requirement imply to the feedback matrix \mathbf{A} ? We know that by definition of losslessness the equality $\int_{\omega} \left\{ \sum_{i=1}^n |Y_i(e^{j\omega})|^2 \right\} d\omega = \int_{\omega} \left\{ \sum_{i=1}^n |X_i(e^{j\omega})|^2 \right\} d\omega$ must hold. Moreover it is a general result that a multidimensional filter is lossless if and only if its frequency response matrix $\mathbf{H}(e^{j\omega})$ is unitary, i.e. $\mathbf{H}(e^{j\omega})\mathbf{H}^*(e^{j\omega}) = \mathbf{I}$ (where $*$ denotes the complex-conjugate transpose as usual). In our case, it is quite straightforward to prove that \mathbf{A} being unitary is a sufficient condition for the overall frequency response matrix to be unitary. Moreover the entries a_{ij} have to be real in order for the system to output a real signal $y[n]$, and a unitary matrix with real entries is an orthogonal matrix.

In conclusion, if \mathbf{A} is orthogonal then the network of Fig. 4.14 is lossless. Note however that this condition is sufficient but not necessary, thus the system may be lossless even with a non-orthogonal feedback matrix. We will return to this point in Sec. 4.4.2.

4.4.1.3 Designing the lossless prototype

Designing the lossless prototype means choosing the dimension N , the m_i 's, and the feedback matrix \mathbf{A} . Let us start with the dimension N and the delay lengths m_i . Together with the feedback matrix these parameters determine the buildup of reflection density. The criteria that we have examined in Sec. 4.3 (see in particular Eqs. (4.30, 4.31)) can be applied also in this case with satisfactory results. Note however that Eqs. (4.30, 4.31) are no longer valid here, since, a non-diagonal feedback matrix increases the modal and reflection densities. Therefore in general the parameters have to be chosen on the basis of empirical observations. It is generally noted that $N = 8$ to 16 lines with a total delay $\sum_i m_i/F_s$ of 1 to 2 seconds already produce a response perceptually undistinguishable from white noise.

Let us now consider the lossless feedback matrix \mathbf{A} . The simplest orthogonal matrix is a diagonal matrix whose diagonal elements (which are the eigenvalues) have unit modulus: as already seen this choice corresponds to a parallel of ordinary comb filters. A more interesting family of orthonormal matrices are *Householder reflection matrices*. A specific Householder matrix is defined given the

reference vector $\mathbf{u} = [1, \dots, 1]^T$:

$$\mathbf{A} = \mathbf{I} - \frac{2}{N} \mathbf{u} \mathbf{u}^T, \quad \text{then } \mathbf{A} \mathbf{x} = \begin{bmatrix} x_1 - \frac{2}{N} \sum_i x_i \\ \vdots \\ x_N - \frac{2}{N} \sum_i x_i \end{bmatrix}, \quad (4.40)$$

for any input vector \mathbf{x} . We will see in Sec. 4.4.2 that \mathbf{u} can be interpreted as the specific vector about which an input vector is reflected by the matrix \mathbf{A} in an N -dimensional space. A more general formulation may be obtained by replacing the identity matrix in Eq. (4.40) with any $N \times N$ permutation matrix.

The explicit expression for $\mathbf{A} \mathbf{x}$ in Eq. (4.40) shows that applying a Householder matrix to a vector requires $N - 1$ additions and one multiplication to obtain the term $\frac{2}{N} \sum_i x_i$, plus N additions to subtract this term from \mathbf{x} . Therefore the matrix-times-vector operation is only $\mathcal{O}(N)$ as opposed to the usual $\mathcal{O}(N^2)$.

Another interesting feature of the Householder feedback matrix is that for $N \neq 2$ \mathbf{A} does not have null entries. This is a desirable property since it implies that every delay line feeds back to every other delay line, reinforcing the build-up of reflection density. The case $N = 4$ is especially nice, since the matrix entries all have the same magnitude and \mathbf{A} is therefore “balanced”. For larger N the diagonal becomes larger than the off-diagonal elements, and \mathbf{A} approaches a diagonal matrix as $N \rightarrow \infty$. Due to the elegant balance of the $N = 4$ case, a larger ($N = 16$) feedback matrix can be constructed as follows:

$$\mathbf{A} = \frac{1}{2} \begin{bmatrix} \mathbf{A}_4 & -\mathbf{A}_4 & -\mathbf{A}_4 & -\mathbf{A}_4 \\ -\mathbf{A}_4 & \mathbf{A}_4 & -\mathbf{A}_4 & -\mathbf{A}_4 \\ -\mathbf{A}_4 & -\mathbf{A}_4 & \mathbf{A}_4 & -\mathbf{A}_4 \\ -\mathbf{A}_4 & -\mathbf{A}_4 & -\mathbf{A}_4 & \mathbf{A}_4 \end{bmatrix}, \quad \text{where } \mathbf{A}_4 := \frac{1}{2} \begin{bmatrix} 1 & -1 & -1 & -1 \\ -1 & 1 & -1 & -1 \\ -1 & -1 & 1 & -1 \\ -1 & -1 & -1 & 1 \end{bmatrix}. \quad (4.41)$$

is the 4×4 Householder matrix.

Other types of unitary matrices may be used. In particular, unitary feedback matrices can be derived from Hadamard matrices. A Hadamard matrix \mathbf{H} is defined as an $N \times N$, $(-1, 1)$ -matrix (i.e. a matrix whose elements consist only of the numbers -1 or 1) with the additional property that $\mathbf{H} \mathbf{H}^T = N \mathbf{I}$. This means that $\mathbf{A} = \mathbf{H} / \sqrt{N}$ is an orthogonal matrix whose entries all have the same magnitude $1/\sqrt{N}$. In Sec. 4.4.2 we discuss other classes of feedback matrices.

4.4.1.4 Designing lossy components

So far we have designed the lossless prototype. Now we have to correct it by inserting the low-pass filters $H_{lp,i}$ and the correction filter E . The $H_{lp,i}$'s set the reverberation time from infinity to a finite value, by moving the poles slightly inside the unit circle. More precisely, they can be chosen to tune the reverberator to a desired, frequency-dependent reverberation time $T_r(\omega)$.

The following analysis assumes that the filters $H_{lp,i}$ are defined as $H_{lp,i} = [G(z)]_i^m$: this is conceptually equivalent to substituting each delay z^{-1} in the lines with a “damped delay” $G(z)z^{-1}$, where the factor $G(z)$ represents a *filtering per sample* in the propagation medium. We also make the simplifying hypotheses that (1) the response $G(e^{j\omega})$ is zero-phase and that (2) the magnitude $|G(e^{j\omega})|$ is close to 1. Now assume that the lossless prototype has poles $e^{j\omega_i/F_s}$, $i = 1, \dots, N$, then the insertion of the low-pass filters moves the poles to

$$p_i \approx R_i e^{j\omega_i/F_s}, \quad \text{with } R_i = G\left(R_i e^{j\omega_i/F_s}\right) \approx G\left(e^{j\omega_i/F_s}\right), \quad (4.42)$$



where we have exploited our first simplifying hypothesis in assuming that the filters affect the radius of the poles and not their angles, and we have exploited our second simplifying hypothesis in the last approximation for R_i .

We know that the component of the impulse response arising from the i th pole of the system decays as R_i^n , as a function of time n . Therefore the time needed for this response to decay by 60 dB (i.e. $T_r(\omega_i)$) satisfies the relation $20 \log_{10} \left(R_i^{T_r(\omega_i) F_s} \right) = -60$ dB. From Eq. (4.42), and recalling that $H_{lp,i} = G^{m_i}$, we conclude that the ideal low-pass filter satisfies the relation

$$20 \log_{10} \left| H_{lp,i} \left(e^{j\omega_i/F_s} \right) \right| = -60 \frac{m_i}{F_s T_r(\omega_i)}. \quad (4.43)$$

Having been derived in the assumption of zero-phase, this expression disregards the phase response of the $H_{lp,i}$'s, which has the effect of slightly modifying the effective length of the delay m_i . It is usually assumed that in practice this correction has no perceivable effect and can therefore be ignored.

A consequence of incorporating the filters $H_{lp,i}(z)$ into the delay lines is that the envelope of the frequency response of the system will no longer be flat. In particular, for exponentially decaying reverberation the envelope is proportional to the reverberation time at all frequencies. The role of the filter $E(z)$ (often referred to as the *tonal correction filter*) is to compensate for this effect: a flat frequency response envelope is restored if the magnitude response of $E(z)$ is inversely proportional to the reverberation time:

$$\left| E \left(e^{j\omega/F_s} \right) \right| \sim \frac{1}{\sqrt{T_r(\omega)}}. \quad (4.44)$$

Having specified ideal filter responses for the $H_{lp,i}$'s and for E , any number of filter-design methods can be used to find low-order filters that reasonably approximate Eqs. (4.43, 4.44). Note that this design effectively decouples the control over reverberation time from the overall reverberator gain.

M-4.9

Write a function that computes filter coefficients for $H_{lp,i}(z)$ and $E(z)$, given a function $T_r(\omega)$ specified on a set of points $\{\omega_k\}$, and given the filter order k . The native functions `invfreqz` and `stmcb` may help.

Since the function $T_r(\omega)$ is typically very smooth and slowly varying with respect to ω , the filters $H_{lp,i}(z)$ can be chosen to have low order. In particular, first-order filters of the form (4.33) can be used:

$$H_{lp,i}(z) = \frac{g_{1,i}}{1 - g_{2,i}z^{-1}}. \quad (4.45)$$

In this case one can use Eq. (4.43) to find the gains (we only report results):

$$g_{2,i} = \frac{\ln(10)}{4} \log_{10}(a_i) \left(1 - \frac{T_r(0)^2}{T_r(\pi F_s)} \right), \quad g_{1,i} = a_i(1 - g_{2,i}) \quad (4.46)$$

where $a_i = 10^{-3 \frac{m_i}{F_s T_r(0)}}$ is determined from the desired reverberation time at $\omega = 0$, while $g_{2,i}$ sets the reverberation time at high frequencies.

If first-order low-pass filters of the form (4.45) are used, then one can use a correction filter which is also first-order and is determined as follows (we only report results):

$$E(z) = \frac{1 - bz^{-1}}{1 - b}, \quad \text{with } b = \frac{1 - \frac{T_r(\pi F_s)}{T_r(0)}}{1 + \frac{T_r(\pi F_s)}{T_r(0)}}. \quad (4.47)$$

M-4.10

Write a function that computes filter coefficients for $H_{lp,i}(z)$ and $E(z)$ in the first-order case described above, given a function $T_r(\omega)$ specified on a set of points $\{\omega_k\}$.

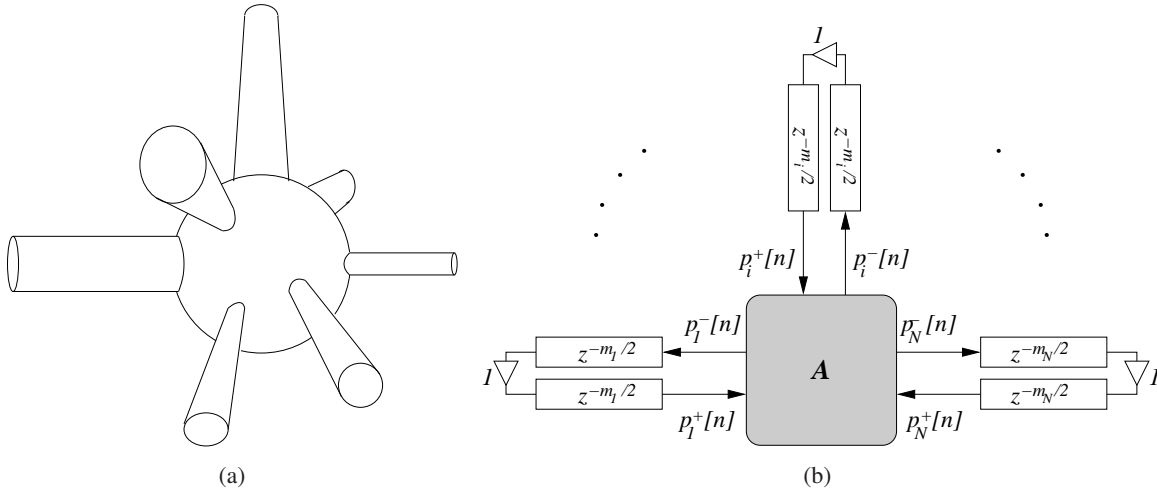


Figure 4.15: DWN reverberator

4.4.2 Digital waveguide networks

4.4.2.1 The link between FDNs and DWNs

In Eq. (4.40) we have introduced a specific Householder reflection matrix, constructed from the reference vector $\mathbf{u} = [1, \dots, 1]^T$. In fact a Householder matrix can be constructed given any reference vector \mathbf{u} . We now want to provide a geometric interpretation of this family of matrices.

Consider the *projection matrix* \mathbf{P}_u , which orthogonally projects any vector \mathbf{x} onto the vector \mathbf{u} :

$$\mathbf{P}_u = \frac{\mathbf{u} \mathbf{u}^T}{\mathbf{u}^T \mathbf{u}} = \frac{\mathbf{u} \mathbf{u}^T}{\|\mathbf{u}\|^2}, \quad \text{then} \quad \mathbf{x}_u := \mathbf{P}_u \mathbf{x} = \mathbf{u} \frac{\langle \mathbf{u}, \mathbf{x} \rangle}{\|\mathbf{u}\|^2} \quad (4.48)$$

is the orthogonal projection of \mathbf{x} onto \mathbf{u} . Now consider the vector $\mathbf{x}_u^\perp := (\mathbf{I} - \mathbf{P}_u)\mathbf{x}$: this is the projection of \mathbf{x} onto the hyperplane orthogonal to \mathbf{u} , since it is easily verified that $\mathbf{x}_u^\perp \perp \mathbf{x}_u$ and that $\mathbf{x}_u^\perp + \mathbf{x}_u = \mathbf{x}$.

Finally consider the vector \mathbf{y} obtained by *reflecting* \mathbf{x} about \mathbf{u} . Elementary geometrical considerations allow to conclude that this vector is the difference between \mathbf{x}_u and \mathbf{x}_u^\perp :

$$\mathbf{y} = \mathbf{x}_u - \mathbf{x}_u^\perp = \mathbf{P}_u \mathbf{x} - (\mathbf{I} - \mathbf{P}_u)\mathbf{x} = (2\mathbf{P}_u - \mathbf{I})\mathbf{x}. \quad (4.49)$$

The matrix $(2\mathbf{P}_u - \mathbf{I})$ is a Householder matrix as defined in Eq. (4.40), except for a sign. Therefore we conclude that given a reference vector \mathbf{u} the corresponding Householder matrix reflects any vector \mathbf{x} about \mathbf{u} .

Having understood the meaning of Householder matrices, we now construct a digital waveguide network (DWN) that is equivalent to the FDN lossless prototypes considered in the previous section. We start by considering the physical resonator depicted in Fig. 4.15(a). It is composed by N acoustic bores connected in parallel. In Chapter *Sound modeling: source based approaches* we have derived the $N \times N$ *scattering matrix* \mathbf{A} that relates the incoming pressure waves \mathbf{p}^+ to the outgoing pressure waves \mathbf{p}^- . In this section we reconsider that matrix when the pressure waves in the i th bore are defined as

$$p_i^+ = \sqrt{\Gamma_i} \frac{p_i + Z_i u_i}{2}, \quad p_i^- = \sqrt{\Gamma_i} \frac{p_i - Z_i u_i}{2}, \quad (4.50)$$

where Z_i and $\Gamma_i = 1/Z_i$ are the wave impedance and admittance of the i th bore. These are often referred to as *normalized waves*, and differ from our previous definition of wave variables uniquely for the scaling factor $\sqrt{\Gamma_i}$. It is straightforward to see that normalized pressure waves are scattered as $\mathbf{p}^- = \mathbf{A}\mathbf{p}^+$, where

$$\mathbf{A} = \begin{bmatrix} \frac{2\Gamma_1}{\Gamma_J} - 1, & \frac{2\sqrt{\Gamma_1\Gamma_2}}{\Gamma_J}, & \dots & \frac{2\sqrt{\Gamma_1\Gamma_N}}{\Gamma_J} \\ \frac{2\sqrt{\Gamma_2\Gamma_1}}{\Gamma_J}, & \frac{2\Gamma_2}{\Gamma_J} - 1, & \dots & \frac{2\sqrt{\Gamma_2\Gamma_N}}{\Gamma_J} \\ \vdots & & \ddots & \vdots \\ \frac{2\sqrt{\Gamma_N\Gamma_1}}{\Gamma_J}, & \frac{2\sqrt{\Gamma_N\Gamma_2}}{\Gamma_J}, & \dots & \frac{2\Gamma_N}{\Gamma_J} - 1 \end{bmatrix}, \quad \text{where } \Gamma_J = \sum_{l=1}^N \Gamma_l. \quad (4.51)$$

This normalized scattering matrix is immediately recognized as a Householder matrix:

$$\mathbf{A} = \frac{2}{\|\boldsymbol{\Gamma}\|} \boldsymbol{\Gamma}\boldsymbol{\Gamma}^T - \mathbf{I}, \quad \text{with } \boldsymbol{\Gamma} := [\sqrt{\Gamma_1}, \sqrt{\Gamma_2}, \dots, \sqrt{\Gamma_N}], \quad (4.52)$$

so we have this interesting geometrical interpretation: scattering of normalized pressure waves corresponds to a reflection around the vector $\boldsymbol{\Gamma}$.

If the acoustic bores are lossless and with ideal closed terminations, and if the length (in samples) of the i th bore is $m_i/2$, then the physical resonator of Fig. 4.15(a) can be modeled with the digital waveguide network given in Fig. 4.15(b). Now compare this scheme with the lossless FDN of Fig. 4.14: apart from the input signals $x_i[n]$, the two schemes implement the same computational structure. The incoming pressure waves $p_i^+[n]$ correspond to the output signals $y_i[n]$, and the outgoing pressure waves $p_i^-[n]$ correspond to the feedback signals generated by the feedback matrix.

4.4.2.2 General lossless scattering matrices

Showing the equivalence between DWNs and FDNs is more than a mere intellectual exercise: we can now design an entire new class of lossless FDN prototypes, in which the feedback matrix \mathbf{A} is given by Eq. (4.51) and have a straightforward physical interpretation.

Note that the matrix in Eq. (4.51) is still orthogonal (it is easy to verify that $\mathbf{A}\mathbf{A}^T = \mathbf{I}$). We can push the generalization further by generalizing our definition of losslessness, and consequently define new classes of lossless feedback matrices that are neither physical nor orthogonal. Consider a Hermitian, positive-definite $N \times N$ matrix $\boldsymbol{\Gamma}$ (we use this notation because we interpret $\boldsymbol{\Gamma}$ as a generalized junction admittance). This matrix induces a norm $\|\cdot\|_{\boldsymbol{\Gamma}}$, defined as follows: $\|\mathbf{x}\|_{\boldsymbol{\Gamma}} := \mathbf{x}^T \boldsymbol{\Gamma} \mathbf{x}$ for any real valued N -dimensional vector \mathbf{x} . We can then define a waveguide scattering matrix \mathbf{A} to be “lossless” if the scattering preserve the norm, i.e. the equality $\|\mathbf{p}^+\|_{\boldsymbol{\Gamma}} = \|\mathbf{p}^-\|_{\boldsymbol{\Gamma}}$ holds. This condition is clearly equivalent to the condition

$$\mathbf{A}^T \boldsymbol{\Gamma} \mathbf{A} = \boldsymbol{\Gamma} \quad (4.53)$$

for the scattering matrix \mathbf{A} . In the case $\boldsymbol{\Gamma} = \mathbf{I}$, the norm $\|\cdot\|_{\boldsymbol{\Gamma}}$ is the euclidean norm and the above equation reduces to the condition of \mathbf{A} being orthogonal. In the general case $\boldsymbol{\Gamma} \neq \mathbf{I}$ it can be shown that Eq. (4.53) holds if and only if \mathbf{A} has eigenvalues with modulus 1 and N linearly independent eigenvectors. We do not provide a proof of this characterization: intuitively it means that when such a feedback matrix is used in a lossless FDN prototype the system poles all have unit modulus and thus the system response consists of non-decaying eigenmodes.

Clearly orthogonal matrices are lossless in this sense, since they have unitary eigenvalues and pairwise orthogonal eigenvectors. Another class of matrices that satisfy this condition are triangular matrices: designing a triangular matrix with unitary eigenvalues is straightforward since we know from linear algebra that they lie on the diagonal. Additional care is required in order to ensure that the triangular matrix possesses N independent eigenvectors.

4.4.2.3 Waveguide meshes

So far we have seen DWNs in analogy with FDNs. In this section we discuss a new multidimensional waveguide structure, named *waveguide mesh*, that can be used to physically simulate resonating enclosures. What follows is only a quick and qualitative introduction to the subject, the interested reader can refer to the bibliography.

Consider again the N-D D'Alembert equation (4.1). Similarly to what we have done in the 1-D case (Chapter *Sound modeling: source based approaches*), we can simulate the traveling wave solution by using delay lines. In this case the delay lines are arranged in a mesh, that represents waves propagating in the x, y, z directions. At each node of the mesh continuity constraints must be satisfied, namely the pressure waves in each direction must provide the same pressure value.⁴ This means that at each node of the mesh the incoming pressure waves are scattered by a matrix identical to the matrix \mathbf{A} given in Chapter *Sound modeling: source based approaches*, in which all the incoming branches share the same impedance:

$$\mathbf{A} = \begin{bmatrix} \frac{2}{N} - 1 & \frac{2}{N} & \cdots & \frac{2}{N} \\ \frac{2}{N} & \frac{2}{N} - 1 & \cdots & \frac{2}{N} \\ \vdots & & \ddots & \vdots \\ \frac{2}{N} & \frac{2}{N} & \cdots & \frac{2}{N} - 1 \end{bmatrix}. \quad (4.54)$$

In order to clarify this idea, let us examine the 2-D case shown in Fig. 4.16. The outgoing pressure waves at each node are computed as $\mathbf{p}^- = \mathbf{A}\mathbf{p}^+$, i.e.

$$p_i^- [n] = p_J [n] - p_i^+ [n] \quad (i = 1 \dots 4) \quad \text{where} \quad p_J [n] = \frac{\sum_{i=1}^4 p_i^+ [n]}{2} \quad (4.55)$$

is the junction pressure. It can be shown that this rectangular waveguide mesh is equivalent to a finite-difference numerical solution of the the 2-D D'Alembert equation, in which the pressure at a certain node is expressed in terms of the pressures at its neighboring nodes one sample earlier, and itself two samples earlier.

The rectangular layout depicted in Fig. 4.16 is not the only possible one: other geometries may be used for assembling the mesh, like triangular, hexagonal, and so on. The choice of the geometry has a major influence on the *dispersion* error in the mesh, i.e the error in propagation speed as a function of frequency and direction along the mesh. It can be shown that the triangular waveguide mesh is the simplest 2-D mesh geometry with the least dispersion variation as a function of direction of propagation. In other words, the triangular mesh is closer to isotropic than all other known elementary geometries. Isotropy can be obtained also through interpolation, i.e. by using non integer propagation delays, but computational costs are higher. As far as frequency dispersion is concerned, frequency-warping methods can be used to minimize it in the mesh.

⁴In this section we are using waveguide meshes to simulate resonating enclosures and thus we work with pressure waves and consider parallel junctions. Waveguide meshes can also be used to simulate mechanical resonators, e.g membranes, and in that case it is natural to choose velocity waves and to consider series junctions at mesh nodes.

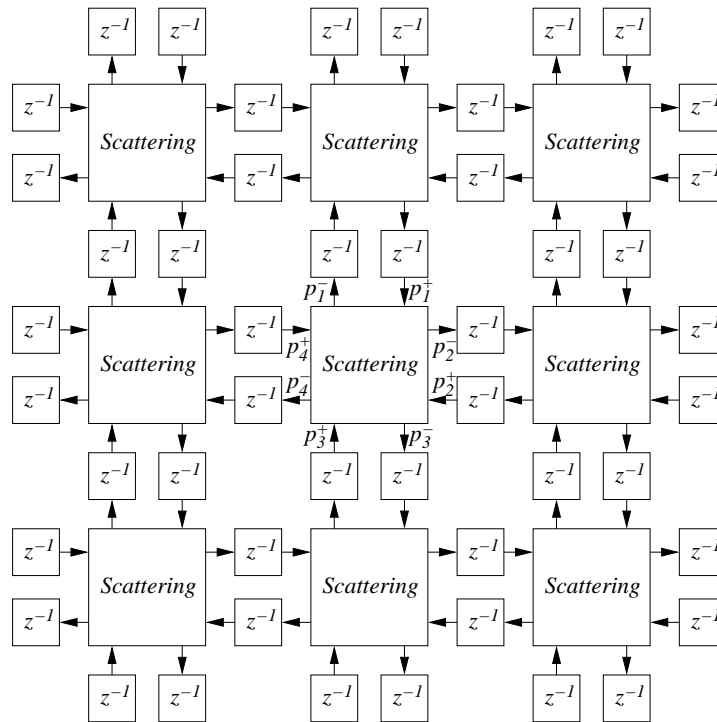


Figure 4.16: 2D rectilinear digital waveguide mesh.

The waveguide meshes analyzed so far simulate lossless propagation in an infinite medium. In order to model something similar to a real resonating enclosure we must add losses and boundary conditions into the structure. The techniques discussed in Chapter *Sound modeling: source based approaches* to simulate lossless in 1-D wave propagation can be extended to the waveguide mesh: the basic idea is once again that wave propagation during one sampling interval (in time) is associated with linear filtering by $G(z)$. The problem of modeling mesh boundaries is particularly important in the context of artificial reverberation: in order to obtain high temporal reflection densities, maximally *diffusing* boundaries have to be modeled.

As efficient solutions are found to deal with the above mentioned problems, 3-D waveguide meshes are being more and more used for the simulation of acoustic spaces.

4.5 Spatial hearing

Sound is transformed by the pinnae (the visible portion of the outer ear) and proximate parts of the body such as the shoulder and head. Following this are the effects of the meatus (or “ear canal”) that leads to the eardrum.

Our assumption is that the sound pressure at the two eardrums is a sufficient stimulus. Producing the same sound pressure will produce the same auditory perception. Caveats: Bone conduction, Adaptation, Conflicting visual cues, Conflicting expectations

Exact reproduction of the sound pressure is not necessary for producing the same auditory perception. The limitations of neural responses allow different (and simpler) stimuli to produce the same response. Examples: Bandwidth (20 Hz to 20 kHz) Amplitude (1-dB resolution) Monaural phase

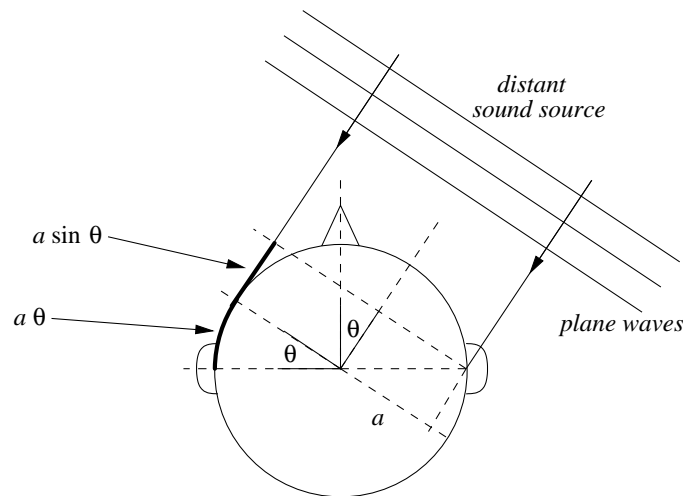


Figure 4.17: Estimate of ITD in the case of a distant sound source (plane waves) and spherical head..

(2-ms resolution) Latency (10-ms resolution) Spectral fine structure (critical bands, $Q = 8$)

4.5.1 The sound field at the eardrum

Spatial attributes of the sound field are coded into temporal and spectral attributes via this filtering effect of head, external ear, torso and shoulders.

4.5.1.1 Head

Our ears are not isolated objects in space. They are located, at the same height, on opposite sides of an acoustically rigid object: the head. This acts as an obstacle to the free propagation of sound and has two main effects: (1) it introduces an *interaural time difference (ITD)*, because a sound wave has to travel an extra distance in order to reach the farthest ear, and (2) it introduces an *interaural level difference (ILD)* because the farthest ear is acoustically “shadowed” by the presence of the head.

As one may expect, the ILD is highly frequency dependent: at low frequencies (i.e., for wavelengths that are long relative to the head diameter) there is hardly any difference in sound pressure at the two ears, while at high frequencies differences can be up to 20 dB or more. On the other hand an approximate yet quite accurate description of the ITD can be derived using a few simplifying assumptions, in particular by considering the case of “distant” sound sources and a spherical head: this situation is depicted in Fig. 4.17.

The first assumption implies that the sound waves that strike the head are plane waves. Then the extra-distance Δx needed for a sound ray to reach the farthest ear is estimated from elementary geometrical considerations, as shown in Fig. 4.17, and the ITD is simply $\Delta x/c$. Therefore

$$\text{ITD} \sim \frac{a}{c}(\theta + \sin \theta), \quad (4.56)$$

where a is the head radius and θ is the *azimuth* angle that defines the direction of the incoming sound on the horizontal plane. This formula shows that the ITD is zero when the source is directly ahead

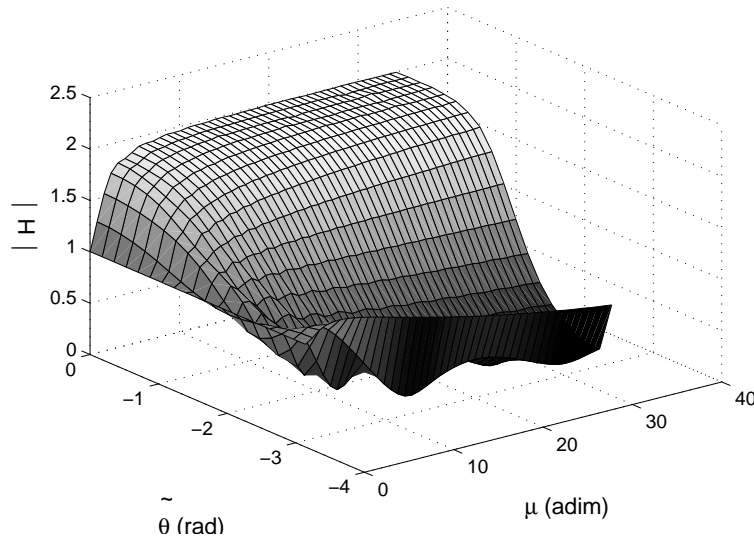


Figure 4.18: Magnitude response $|H(\infty, \mu, \tilde{\theta})|$ of a sphere for an infinitely distant source.

($\theta = 0$), and is a maximum of $a/c(\pi/2 + 1)$ when the source is off to one side ($\theta = \pi/2$). This represents an ITD of more than 0.6 ms for a head radius $a = 8.5$ cm, which is a realistic value.

Take a sphere of radius a , a point sound source at a distance $r > a$ from the center of the sphere, and a point on the sphere. It is customary to use the normalized variables $\mu = \omega a/c$ (normalized frequency) and $\rho = r/a$ (normalized distance). Then the diffraction of an acoustic wave by the sphere seen on the chosen point is expressed with the transfer function

$$H(\rho, \mu, \tilde{\theta}) = -\frac{\rho}{\mu} e^{-i\mu\rho} \sum_{m=0}^{+\infty} (2m+1) P_m(\cos \tilde{\theta}) \frac{h_m(\mu\rho)}{h'_m(\mu)}, \quad (4.57)$$

where P_m and h_m are the m th order Legendre polynomial and spherical Hankel function, respectively, and $\tilde{\theta}$ is the angle of incidence, i.e. the angle between the ray from the center of the sphere to the source and the ray to the measurement point on the surface of the sphere.⁵ Normal incidence corresponds to $\tilde{\theta} = 0$, while the sphere point opposite to the source is at $\tilde{\theta} = \pi$.

It is known that the Hankel function $h_m(x)$ admits an asymptotic approximation as the argument x goes to infinity. By exploiting this approximation one can study the behavior of the transfer function $H(\infty, \mu, \tilde{\theta})$ as the distance r between the source and the sphere becomes arbitrarily large. The approximate solution $|H(\infty, \mu, \tilde{\theta})|$ is plotted in Fig. 4.18.

At low frequencies the transfer function is not directionally dependent and the magnitude $|H|$ is essentially unity for any angle $\tilde{\theta}$. When μ exceeds 1 the dependence on $\tilde{\theta}$ becomes noticeable. The response increases around the front of the sphere, and in particular exhibits a 6 dB boost at high frequencies near the front of the sphere ($|H(\infty, \infty, 0)| = 2$), consistently with the requirement that in this limit the solution must reduce to that of a plane wave normally incident on a rigid plane surface. $|H|$ is approximately flat when $\tilde{\theta}$ is around 100 degrees, and progressively decreases around the back

⁵We are using a different notation with respect to the azimuth angle θ used previously, in order to avoid confusion. Given a 2-D reference system like that in Fig. 4.17, the transfer functions (4.57) at the right and left ear will use the angles $\tilde{\theta}^{(r)} = \theta - \pi/2$ and $\tilde{\theta}^{(l)} = \theta + \pi/2$, respectively.

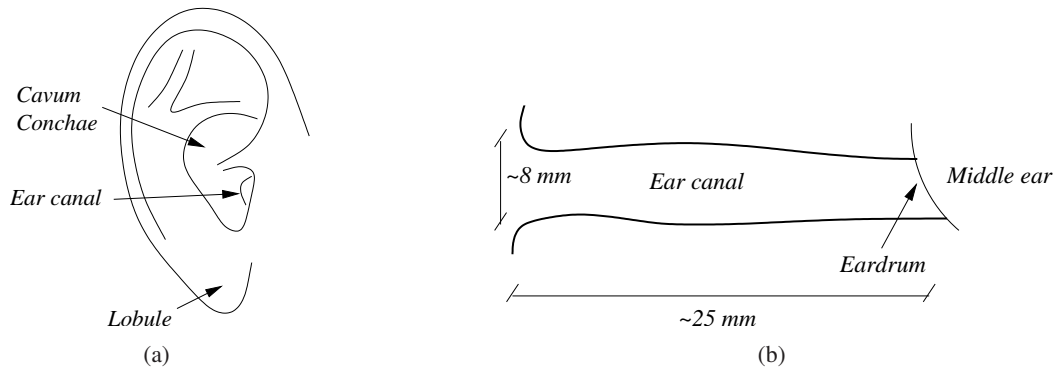


Figure 4.19: External ear: (a) pinna, and (b) ear canal.

of the sphere. Note however that the minimum response does not occur at the very back ($\tilde{\theta} = \pi$). Instead, this point exhibits a so-called “bright spot” effect, which is due to the fact that all the waves propagating around the sphere arrive at that point in phase. At very high frequencies the bright-spot lobe becomes extremely narrow, and the back of the sphere is effectively in a sound shadow. Finally, note that interference effects caused by waves propagating in various directions around the sphere introduce ripples in the response that are quite prominent on the shadowed side.

4.5.1.2 The external ear

External ear consists of the pinna and ear canal until the eardrum. Then middle ear and internal ear. Here we are interested in the external ear only, in Chapter *Auditory based processing* we will study the middle and internal ear.

Pinna: Fig. 4.19(a). It has a characteristic “bas-relief” form with features that differ greatly from one individual to another (just look at people’s ears). The pinna is connected to the *ear canal*: Fig. 4.19(b). It can be approximately described as a tube of constant width, with walls of high acoustic impedance. At the end opposite to the pinna, the ear canal is terminated by the eardrum diaphragm.

At a first approximation the acoustic behaviour of the ear canal is easily understood: it behaves like a one-dimensional resonator. On the other hand the pinna has much more complex effects, as it basically acts like an acoustic antenna. Its resonant cavities amplify some frequencies, and its geometry leads to interference effects that attenuate other frequencies. Moreover, its frequency response is directionally dependent. Acoustically it acts like a filter whose transfer function depends in general on the distance and direction of the sound source relative to the ear. Like for any other resonator, we can interpret these filtering effect either in the frequency domain or by looking at reflections of sound rays.

First approach: measurements of frequency responses using an imitation pinna and a ear canal with high impedance termination. The measurements give results like those depicted in Fig. 4.20(a). First resonance is that of a open-closed tube $\sim 33\%$ longer than the ear canal: the pinna acts as a prolongation of the ear canal with an aperture effect. Second resonance is a resonance of the *cavum concha* alone: the pressure distribution is similar to what would be obtained if the canal were plugged. The higher resonances instead are again associated to longitudinal standing waves: these are not very widely spaced and are quite dependent on the individual, therefore it can combine in a single broad peak of the magnitude response.

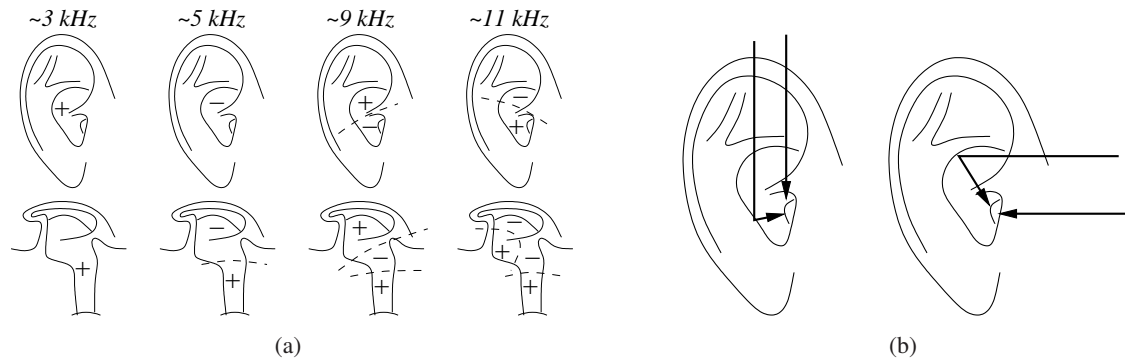


Figure 4.20: Effects of pinna: (a) resonances, and (b) direction-dependent reflections

Second approach: external ear as a sound reflector. Figure 4.20(b) shows two different directions of arrival. In each case there are two paths from the source to the ear canal—a direct path and a longer path following a reflection from the pinna. At moderately low frequencies, the pinna essentially collects additional sound energy, and the signals from the two paths arrive in phase. However, at high frequencies, the delayed signal is out of phase with the direct signal, and destructive interference occurs. The greatest interference occurs when the difference in path length is a half wavelength: this produces a “pinna notch”. Since the pinna is a more effective reflector for sounds coming from the front than for sounds from above, the resulting notch is much more pronounced for sources in front than for sources above. In addition, the path length difference changes with elevation.

The synthetic conclusion of this section is then that the pinna and the ear canal form a systems of acoustic resonators, whose resonances are excited to different extents depending on the direction and distance of the sound source.

4.5.1.3 Torso and shoulders

In the discussion up to now we have not considered a third element that, together with the head and the external ear, contributes to the shaping of the sound field at the eardrum: the torso. Torso and shoulders affect incident sound waves in two main respects. First, they provide additional reflections that sum up with the direct sound. Second, they have a shadowing effect for sound rays coming from below.

The geometry of the torso is quite complicated. However a simplified description can be derived by considering an ellipsoidal torso below a spherical head. These kind of approximate descriptions are sometimes called “snowman models”, for obvious reasons. Figure 4.21(a) depicts a snowman model and shows the main effects of the ellipsoidal torso on the sound field at the ear.

Reflections: Fig. 4.21(a). If we measured the impulse response at the right ear for the sound source locations depicted in Fig. 4.21(a) we would see that the initial pulse is followed by a series of subsequent pulses, whose delays increase and then decrease with elevation. These additional pulses are caused by reflections on the torso.

We could exploit the simplified geometry of the snowman model to compute analytically the delay of the reflected rays, given the model parameters and the sound source position. However some important remarks can already be made from a qualitative analysis. First, the delay between the direct sound and the reflected ray does not vary much if the sound source moves on a circumference in the horizontal plane (especially if its radius is large compared to the head radius). Second, the delay varies

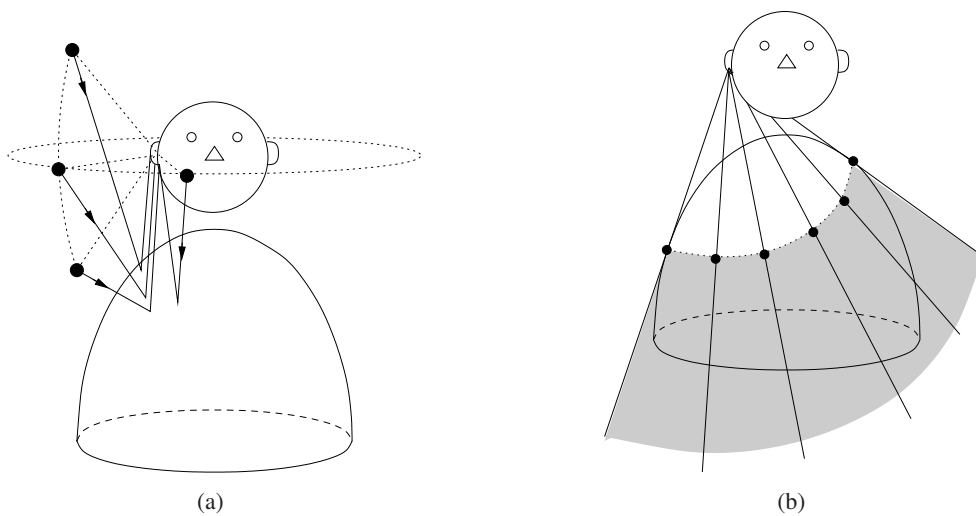


Figure 4.21: Effects of torso: (a) reflections, and (b) shadowing.

considerably if the sound source moves vertically, and in particular the reflected pulses are maximally delayed for sound source locations right above the listener. If we consider that the distance from the ear canal to the shoulder is roughly 16 cm, then a reflected ray from a source right above the subject has to travel an extra distance of approximately 32 cm, which corresponds to a delay of almost 1 ms.

In the frequency domain the torso reflections act as a comb filter, introducing periodic notches in the spectrum. The frequencies at which the notches occur are inversely related to the delays, and thus produce a pattern that varies with the elevation of the source. The lowest notch frequency corresponds to the longest delay. Delays longer than a sixth of a millisecond will produce one or more notches below 3 kHz, which is approximately the lowest frequency where pinna effects start to be noticeable.

Modeling the effects of the torso as specular reflections means accounting for only a part of the story. First, reflection is a high frequency concept. Second, and perhaps more important, as the source descends in elevation, a point of grazing incidence is reached, below which torso reflections disappear and *torso shadowing* emerges. As shown in Fig. 4.21(b), rays drawn from the ear to points of tangency around the upper torso define a torso-shadow cone. Clearly, the specular reflection model does not apply within the torso shadow cone. Instead, diffraction and scattering produce a qualitatively different behavior, characterized by a stronger attenuation for high frequencies (i.e. for wavelength comparable to or smaller than the size of the torso).

Although the acoustic effects of torso and shoulders are not as strong as those introduced by the pinna, they are important because they appear at lower frequencies, where typical sound signals have most of their energy and where the response of the pinna is essentially flat. In terms of frequency ranges the effects provided by the torso are therefore complementary to those provided by the pinna.

4.5.1.4 Head-related transfer functions

In the preceding sections we have investigated the influence of hear, torso and external ear on the sound field at the eardrum. All the effects that we have examined are linear, which means that (1) they can be described by means of transfer functions, and (2) they combine additively. Therefore the sound pressure produced by an arbitrary sound source at the eardrum is uniquely determined by the impulse

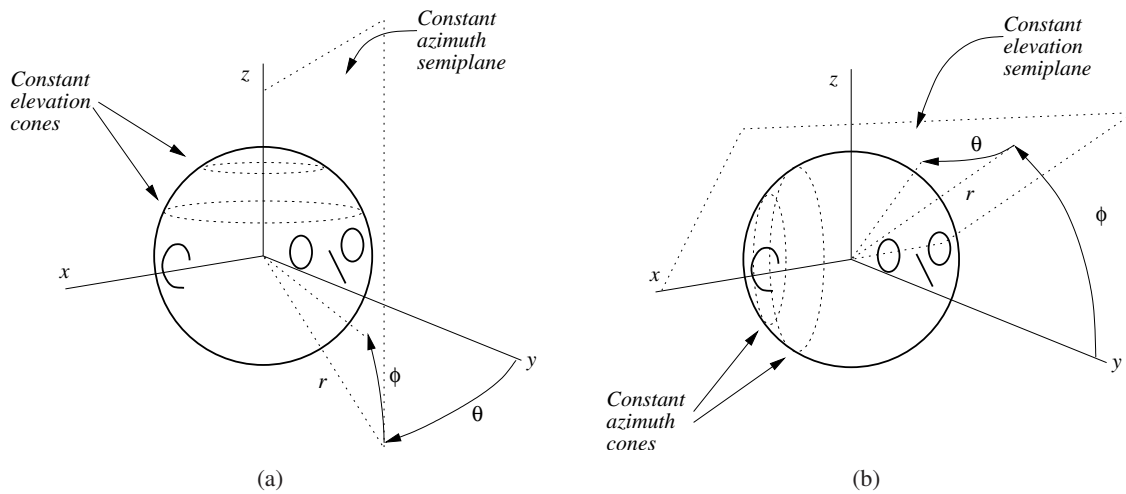


Figure 4.22: Spherical coordinate systems: (a) vertical-polar coordinate system, and (b) interaural-polar coordinate system.

response from the source to the eardrum. This is called *Head-Related Impulse Response (HRIR)*, and its Fourier transform is called *Head Related Transfer Function (HRTF)*. The HRTF captures all of the physical effects that we have examined separately in the previous sections.

The HRTF is a function of three spatial coordinates and frequency. Given the approximately spherical shape of the head, it is customary to use spherical coordinates. The angular coordinates are named *azimuth* and *elevation* and noted as θ and ϕ , respectively, while the radial coordinate is named *range* and noted as r . Note however that more than one choice of spherical coordinates is available. Figure 4.22(a) show the most popular one, sometimes called *vertical polar* coordinate system: in this system the azimuth is measured as the angle from the yz plane to a vertical plane containing the source and the z axis, and the elevation is measured as the angle up from the xy plane. With this choice, surfaces of constant azimuth are planes through the z axis, and surfaces of constant elevation are cones concentric about the z axis.⁶

In alternative the so-called *interaural-polar* coordinate system, shown in Fig. 4.22(b), is sometimes used. In this case the elevation is measured as the angle from the xy plane to a plane containing the source and the x axis, and the azimuth is then measured as the angle from the yz plane. With this choice, surfaces of constant elevation are planes through the x axis, and surfaces of constant azimuth are cones concentric with the x axis. One advantage of this system is that it makes it significantly simpler to express interaural differences at all elevations (in particular the constant-azimuth cones are the loci of points that share equals ILD and ITD values for a spherical head).

In the remainder of this chapter we will specify, when necessary, whether we are using the vertical-polar or the interaural-polar coordinate system. In any case we will indicate the HRTFs as $H_{l,r}(r, \theta, \phi, \omega)$ or, in the limiting case of a “distant” sound source (in most practical applications for $r > 1$ m), as $H_{l,r}(\theta, \phi, \omega)$. The subscripts l, r will indicate the HRTF at the left and right ear, respectively. In the hypothesis of a perfectly symmetrical geometry we will simply write $H(\theta, \phi, \omega)$, with $H_r(\theta, \phi, \omega) = H(\theta, \phi, \omega)$ and $H_l(\theta, \phi, \omega) = H(-\theta, \phi, \omega)$.

The HRTF is a surprisingly complicated function.

⁶This is the coordinate system that we have already introduced in Chapter *Sound modeling: source based approaches*.

(a) (b) (c) (d)

Figure 4.23: *HRIRs and HRTFs*

In spherical coordinates, for distances greater than about one meter, the source is said to be in the far field. Most HRTF measurements are made in the far field, which essentially reduces the HRTF to a function of azimuth, elevation and frequency.

We formally define the HRTF at one ear as the frequency-dependent ratio between the sound pressure level (SPL) $\Phi_{l,r}(\theta, \phi, \omega)$ at the corresponding eardrum and the free-field SPL at the center of the head $\Phi_f(\omega)$ as if the listener were absent:

$$H_l(\theta, \phi, \omega) = \frac{\Phi_l(\theta, \phi, \omega)}{\Phi_f(\omega)}, \quad H_r(\theta, \phi, \omega) = \frac{\Phi_r(\theta, \phi, \omega)}{\Phi_f(\omega)}. \quad (4.58)$$

4.5.2 Perception of sound source location

This is a complicated matter. Many competing and interfering effects can influence auditory perception of sound source location. In this section we provide a brief summary, but we warn the reader to be cautious when dealing with this matter and always to be aware of limitations and simplifying hypotheses.

4.5.2.1 Azimuth perception

The horizontal placement of the ears maximizes differences for sound events occurring around the listener, rather than from below or above, enabling audition of sound sources at the terrain level and outside the visual field of view. The ITD and the ILD are considered to be the key parameters for azimuth perception, in what is sometimes referred to as the *duplex theory* of localization.

For the sake of clarity, consider a sine wave reaching the left and right ear. At low frequencies the ITD shifts the waveform a fraction of a cycle, which is easily detected: see Fig 4.24(a). Qualitatively one can say that if the half wavelength is larger than the size of the head, then it is possible for the auditory system to detect the phase of these waveforms unambiguously, and the ITD cue can function. On the other hand, at high frequencies there is ambiguity in the ITD, since there can be several cycles of shift: see Fig 4.24(b). Qualitatively, we can consider the critical point to be the point where the half wavelength becomes shorter than the head size: for shorter wavelengths, the phase information in relation to relative time of arrival at the ears can no longer convey which is the leading wavefront. The critical point in frequency is usually assumed to be a value around 1.5 kHz.

If we now look at the ILD the situation is reversed. As we have seen in Sec. 4.5.1 (see in particular Fig. 4.18), at low frequencies the head transfer function is essentially flat and therefore there is little ILD information. On the other hand, at high frequencies the ILD is more marked and can become very large. For this reason the Duplex Theory asserts that the ILD and the ITD are complementary cues to azimuth perception, and that taken together they provide azimuth perception throughout the audible frequency range.

This is not completely true, though. In fact timing information can be exploited for azimuth perception also in the high frequency range because the timing differences in amplitude envelopes are detected. Again, for the sake of clarity consider a sine wave that is modulated in amplitude as in Fig. 4.24(c). Then an ITD envelope cue, sometimes referred to as *Interaural Envelope Difference*

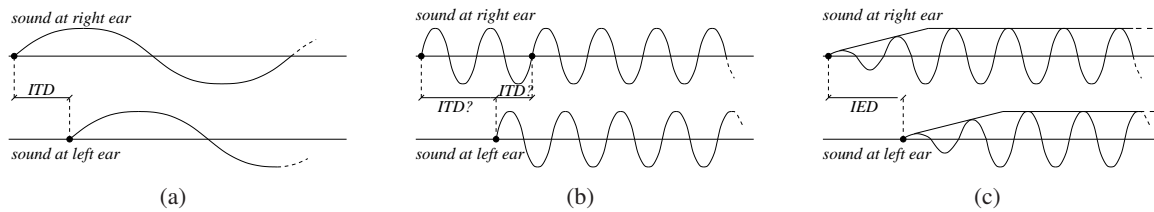


Figure 4.24: Time differences at the ears; (a) non ambiguous ITD, (b) ambiguous ITD, and (c) IED.

(IED) can be exploited, based on the hearing system’s extraction of the timing differences from the transients of amplitude envelopes, rather than from the timing of the waveform within the envelope. This is demonstrated by the so-called Franssen Effect. If a sine wave is suddenly turned on and a high-pass-filtered version is sent to Loudspeaker A while a low-pass filtered version is sent to Loudspeaker B, most listeners will localize the sound at Loudspeaker A. This is true even if the frequency of the sine wave is sufficiently low that in steady state most of the energy is coming from Loudspeaker B.

The information provided by ITD and ILD can be ambiguous. If we assume the spherical geometry of Fig. 4.17, a sound source located in front of the listener at a certain θ , and a second one located at the rear, at $\pi - \theta$, provide identical ITD and ILD values. In reality ITD and ILD will not be exactly identical at θ and $\pi - \theta$ because (1) human heads are not spherical, (2) there are asymmetries and other facial features, and (3) ears are not positioned as in Fig. 4.22 but lie below and behind the x axis. Nonetheless the values will be very similar, and *front-back confusion* is in fact often observed experimentally: listeners operate *reversals* in azimuth judgements, erroneously locating sources at the rear instead of at the front, or viceversa. The former reversal occurs more often than the latter. Some argue that this asymmetry may originate from a sort of ancestral “survival mechanism”, according to which if something (a predator?) can be heard but not seen then it must be at the rear (danger!).

The Duplex Theory essentially works in anechoic conditions. But in everyday conditions reverberation can severely degrade especially ITD information. As we know, in a typical room reflections begin to arrive a few milliseconds after the direct sound. Below a certain sound frequency, the first reflections reach the ear before one oscillation period is completed. Before the auditory system estimates the frequency of the incoming sound wave, and consequently infers the ITD, the number of reflections at the ear has increased exponentially and the auditory system is not able to estimate the ITD. Therefore sounds that possess energy in the low-frequency range only (indicatively below 250 Hz) are essentially impossible to localize in a reverberant environment.⁷ Instead the IED is used, because the starting transient provides unambiguous localization information, while the steady-state signal is very difficult to localize. In conclusion we can state –with some risk of oversimplification– that high-frequency energy only is important for localization in reverberant environments.

4.5.2.2 Lateralization and externalization

In Sec. 4.6 we will see that the simplest systems for spatial sound rendering are based on manipulation of the interaural cues examined above, and on headphone-based auditory display. These systems can be used in applications where only two-dimensional localization –in the horizontal plane– is required.

In this context, the term *lateralization* is typically used to indicate a special case of localization, where the spatial percept is heard inside the head, mostly along the interaural axis (the x of Fig. 4.22),

⁷This is why surround systems use many small loudspeakers for high frequencies and one subwoofer for low frequencies.

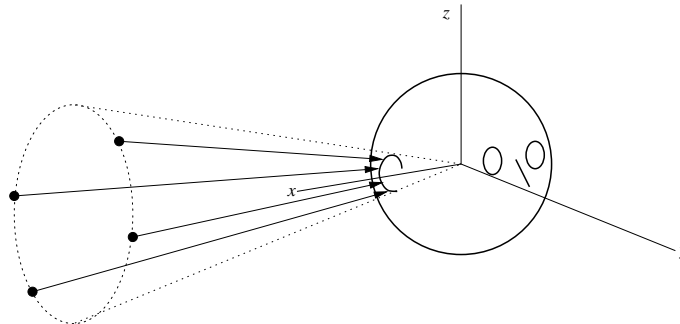


Figure 4.25: *Cone of confusion.*

and the means of producing the percept involves manipulation of ITD and/or ILD over headphones. Lateralization illustrates a fundamental example of virtual, as opposed to actual, sound source position. When identical monaural sounds are delivered from stereo headphones, the listener does not hear two distinct sounds coming from the transducers, and instead perceives a single virtual sound source which appears to be positioned at the center of the head. As ITD and ILD are increased, the perceived position of the virtual sound source will start to shift toward one of the ears, along an imaginary line. Once a critical value of the ITD or the ILD is reached, the perceived sound source will stop moving along the interaural axis and will be located at one of the ears. This effect is sometimes termed *inside-the-head localization* (IHL). Having knowledge of this effect is important since headphone playback is otherwise superior to loudspeakers for transmitting virtual acoustic imagery in three dimensions.

Achieving *externalization* of the sound (i.e. in removing the IHL effect) is in many respects the “sacred graal” of headphone-based spatial audio systems. It is not completely clear what additional cues are most effective in producing sound externalization. However it has been observed by many that externalization increases as the stimulation approximates more closely a stimulation that is natural and that especially reverberation, either natural or artificial, can enhance dramatically externalization. In general, IHL is not an inevitable consequence of headphone listening, simply because externalized sounds can be heard through headphones in many instances.

4.5.2.3 Elevation perception

While the relevant cues for the localization of a sound source in the horizontal plane are relatively well understood, things become more complicated when we consider non-null elevations

Figure 4.25 show that sound sources located anywhere on a conical surface extending out from the ear of a spherical head produce identical values of ITD and ILD. These surfaces are often referred to as *cones of confusion*, and extend the concept of front/back confusion that we have examined above. Of course this situation is only theoretical: in reality ITD and ILD will never be completely identical on the cone of Fig. 4.25, because of the facial features and asymmetries already mentioned above. Nonetheless, when ITD and ILD cues are maximally similar between two locations, a potential for confusion between the positions exists in the absence of other spatial cues.

The directional effects of the pinnae can disambiguate this confusion, and are considered to be particularly important for vertical localization. The role of the pinnae in improving vertical localization can be evaluated experimentally e.g. by comparing judgments made under normal conditions to a condition where the pinnae are bypassed or occluded. In fact vertical localization can be achieved even when one ear is completely occluded. This evidence support the idea that the spectral cues

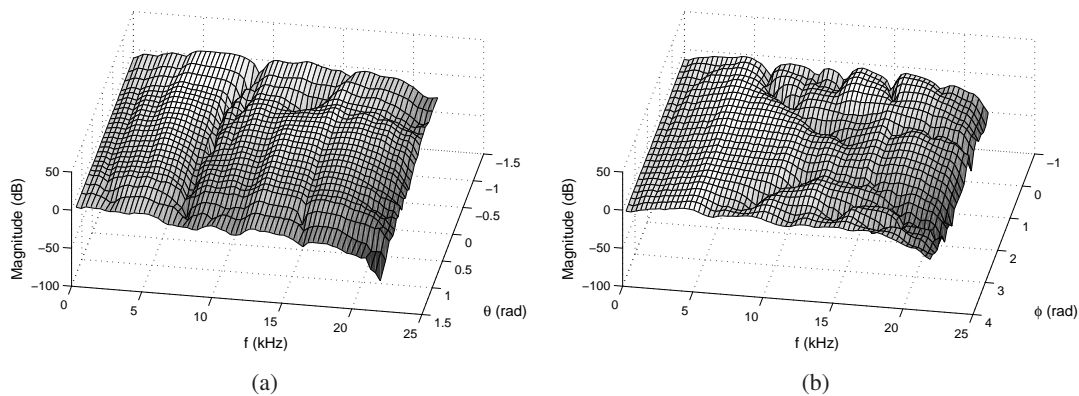


Figure 4.26: HRTFs with varying (a) azimuth and (b) elevation (interaural polar coordinates are used).

provided by the pinnae work mainly monaurally.

There are many theories about the role of pinnae spectral cues. Very roughly, all of them suggest that a major cue for elevation involves movement of spectral notches and/or peaks, that change as a function of source and listener orientation. Figure 4.26(a) provides a plot of measured HRTFs (magnitude response) for a sound source on the horizontal xy plane, as functions of θ . One can notice the movement in the center frequency of two spectral notches: this changes could contribute to the disambiguation of front-back source positions on a cone of confusion. Another way of appreciating the pinnae spectral cues is to examine the special case of sound sources along the yz plane of the listener: note that this is the locus of the points where not only IID and ITD are null, but also spectral differences between the left and right HRTFs are null as long as the left and right pinnae are identical. Figure 4.26(b) provides a plot of measured HRTFs (magnitude response) for a sound source on the yz plane, as functions of ϕ . Again, a moving spectral notch can be noticed, that is thought to be important for elevation perception.

In general it is difficult without extensive psychoacoustic evaluation to ascertain how importantly these changes function as spatial cues. In particular, it is unclear if localization cues are derived from a particular spectral feature such as a peak or a notch, or from the overall spectral shape. Also, it is generally considered that a sound source has to contain substantial energy in the high-frequency range for accurate judgment of elevation, because the pinna has limited dimensions in space and wavelengths longer than the size of the pinna are not affected (see also Fig. 4.20(a)). One could roughly state that the pinnae have a relatively little effect below 3 kHz.

While the role of the pinna in vertical localization has been extensively studied, the role of the torso is less well understood. We have seen in Sec. 4.5.1 that the torso disturbs incident sound waves at frequencies lower than those affected by the pinna. However, the effects of the torso are relatively weak, and experiments to establish the perceptual importance of low-frequency cues have produced mixed results.

4.5.2.4 Distance perception

It is an unanimous claim that auditory estimation of azimuth is more accurate than elevation estimation, and that distance estimation is the most difficult task. Similarly, the cues for azimuth are quite

well understood, those for elevation are less well understood, and those for distance are least well understood. Distance perception involve a process of integrating multiple cues, the most important being *loudness*, ratio of *direct to reverberant sound*, cognitive *familiarity*, and distance dependent *spectral effects*. Additional effects are produced in the so-called *near field* case, i.e. when the distance between the sound source and the listener is less than approximately 1 m. Any of these cues can be rendered ineffective by the summed result of other potential cues.

In the absence of other information, the *intensity* of a sound source is the primary distance cue used by listeners, who learn from experience to correlate the physical displacement of sound sources with corresponding increases or reductions in intensity. Under anechoic conditions, one can use the inverse square law to predict sound intensity reduction with increasing distance. Given a reference intensity and distance, an omnidirectional sound source's intensity will fall approximately 6 dB for each distance doubling (we have already remarked this point in Sec. 4.2.2 when discussing the clarity index parameter). However this law is not well motivated perceptually: intensity expresses the ratio of a sound source's intensity to a reference level, whereas the *perceived* magnitude of intensity is called *loudness*. Thus a mapping where the relative estimation of doubled distance follows "half-loudness" rather than "half-intensity" seems preferable: the two scales are different. Without entering into details, we can say that experimental results show that, for most sounds, an increase of 10 dB is roughly equivalent to a doubling of loudness.⁸

However loudness (or intensity) increments can only operate effectively as distance cues in the absence of other information, in particular reverberation. When reverberation is present the overall loudness at a listener's ear does not change much for very close and very distant sources: the distance-dependent scaling applies only to the direct sound whereas the *reflected energy* remains approximately constant. Reverberation is often not included in distance perception studies, thereby giving subjects an incomplete and non-realistic information with respect everyday listening situations. In particular, estimation of distance with anechoic stimuli is usually worse than in experiments with "optimal" reverberation conditions. As an example, many experimental results show an overall underestimation of the apparent distance of a sound source in an anechoic environment, which may be explained by the absence of reverberation in the stimulus. Many studies report that in a reverberant context the change in the proportion of reflected to direct energy, the so-called *R/D ratio*, functions as a stronger cue for distance than intensity scaling. In particular a sensation of changing distance can occur if the overall loudness remains constant but the the R/D ratio is altered. Note however that in some contexts the possible R/D ratio variation can be limited by the size of the particular environmental context, causing the cue to be less robust (e.g. in a small, acoustically treated room, the ratio would vary between smaller limits than in a large room like a gymnasium). It can be said that reverberation provides the "spatiality" that allows listeners to move from the domain of loudness inferences to the domain of distance inferences, i.e. from an analytic listening attitude to an *everyday listening* attitude.

Distance perception is also affected by expectation or *familiarity* with the sound source. If the sound source is completely synthetic (e.g., pulsed white noise), then a listener will typically concentrate on parametric changes in loudness and other cues that occur for different simulated distances (in this case loudness probably plays a more important role than reverberation effects). On the other hand, if the sound source is cognitively associated with a typical distance range, that range will be more easily perceived than unexpected or unfamiliar distances. This is especially true for speech: as an example, it is easier to simulate a whispering voice 20 cm away from your ear than it is to simulate the same whisper 5 m away.

Distance-dependent spectral effects can also affect distance perception, although to a lesser extent

⁸We will return on the concept of loudness in Chapter *Auditory based processing*.

than the cues discussed above. One effect is due to the influence of atmospheric conditions and air absorption: with increasing distance, higher frequencies of a complex sound are increasingly attenuated by air humidity and temperature. There is little experimental evidence this spectral effect is actually used by listeners in forming the distance of an auditory event, although some experimental results suggest that, in the absence of other cues, a low-frequency emphasis applied to a stimulus would be interpreted as “more distant” compared to an untreated stimulus. A second distance-dependent spectral effect is produced in the so-called *near field*, i.e. for distances less than approximately 1 m. Within this range it is not possible to assume the sound wavefronts to be planar, and the effect of their curvature must be taken into account. The relatively simple sphere transfer function given in Eq. (4.57) already shows what these effects are: by computing the response for the same angle of incidence $\hat{\theta}$ but with varying normalized distances ρ , one would note that as the source approach the sphere emphasis is added to lower frequencies. This phenomenon corresponds to the “darkening” of tone color that occurs as a sound source is moved very close to one’s ear.

Note that all the cues discussed above are essentially monaural cues. An open question is whether binaural listening improves the perception of distance. This could indeed be the case again in the near field limit. The spherical head model shows that in this limit both the ILD and the ITD at low frequencies are emphasized, especially for very lateralized sound sources ($\theta \sim \pm\pi$). This effect is sometimes termed *auditory parallax*, and has been interpreted by some to mean that the accuracy of estimation of a sound from the side should be improved when compared to distance perception on the median plane. There are numerous discrepancies in the literature, however, and the question of binaural cues to distance is therefore still unresolved.

4.5.2.5 Dynamic cues

So far in this section we have examined sound source perception in the implicit assumption of static conditions, i.e. with both listener and source not moving. However in everyday perception we use also *dynamic* cues in addition to static ones to reinforce localization. These arise from active, sometimes unconscious, motions of listeners, who change their position relative to the source. When we hear a sound that we want to localize, we move e.g. in order to minimize the interaural differences, using our head as a sort of “pointer”. Animals use movable pinnae for the same purpose (think of a cat).

Several studies have shown that allowing listeners to move their head can improve localization ability. Listeners apparently integrate some combination of the changes in ITD, ILD, and movement of spectral notches and peaks that occur with head movement over time, and subsequently use this information. Perhaps the most clear example in this respect is represented by front/back confusions: while these are common in static listening tests (see our discussion about cones of confusion), they can be disambiguated and disappear when listeners are allowed to turn their heads during the sound source localization task. A sound source located in the horizontal plane at $\theta = 30^\circ$ could potentially be confused with a source at $\theta = 150^\circ$. A listener who is trying to localize this source will probably attempt to center the auditory image by moving his head to the right, since ITD and ILD cues suggest that the source is somewhere to the right, in spite of front-back ambiguity. If the sound source becomes increasingly centered –i.e. interaural differences are minimized– consequently to rightward head motion, then it must be in the front. If instead it becomes increasingly lateralized –i.e., the sound becomes louder and arrives sooner at the right ear relative to the left– then it must be to the rear.

A second relevant example of the importance of dynamic cues is about externalization. We have seen previously that IHL can occur in lateralization with headphone reproduction. It has been observed that this effect is less likely to occur when head movement is allowed, probably for the same reason that front/back confusion is avoided: dynamic cues arising from head motion can be used to

disambiguate locations, while static conditions can potentially lead to judgments at a “default” position inside or at the edge of the head. An even worse situation is when the sound scene is presented through headphones without tracking of head/body motion, *and* the listener can move. In this case dynamic cues are absent and the scene rotates together with the user, creating discomfort and preventing externalization. The situation changes completely if visual cues are supplied, e.g. if one can move in a fully immersive virtual environment and can see the virtual sound source. In this case it is quite likely that the combination of vestibular and visual cues will enable externalization. In fact externalization can occur even when listening to a television with a single earpiece: this is because vision is more reliable than audition in spatial location, and therefore our brain “trusts” visual rather than auditory feedback (the general mechanism underlying this phenomenon is known as “visual capture”).

Finally, active motion of the listener can provide useful cues for distance perception. If a listener translates his or her head, the azimuth will undergo small or large variations depending on the sound source distance. For sources that are very close, a small shift causes a large change in azimuth, while for sources that are distant the azimuth change will be small. In the limit of infinite distance, there will be no change in azimuth irrespective of the amount of head shift. This dynamic cue is sometimes referred to as *motion parallax*, and is in many respects similar to its visual counterpart (a large, distant sphere and a small, near sphere look the same, but if we move the different changes in perspective reveal the different distances).

4.6 Algorithms for 3-D sound rendering

Before examining processing algorithms for 3-D sound rendering we have to understand that the techniques to be developed depend on the type of system that is going to be used: the type of the effectors (e.g. loudspeakers vs. headphones), as well as their number and geometric arrangement (e.g. stereo systems vs. 5 + 1 surround systems, etc.).

Stereo is the simplest system involving “spatial” sound. In order to place a sound on the left or to the right, its signal is sent to the corresponding loudspeaker. If the same signal is sent to both speakers, the speakers are wired “in phase”, and the listener is approximately equidistant from the speakers, then the listener will perceive a “phantom source” located midway between the two loudspeakers. By crossfading the signal from one speaker to the other, one can create the impression of the source moving continuously between the two loudspeaker positions. With this technique however the perceived source will never move outside the line segment between the two speakers.

Multichannel systems are the next step in complexity. The idea is to have a separate channel for every desired direction, possibly including above and below. Commercial home-theater systems are based on this idea. In typically reverberant environments, one can exploit the limitations of our perception (see in Sec. 4.5.2 our discussion about azimuth perception in reverberant environments) and use small loudspeakers everywhere, except for one large speaker (the “subwoofer”) that provides the nondirectional, low-frequency content.

Headphone-based systems have some disadvantages compared to loudspeakers: headphones are invasive and can be uncomfortable to wear for long periods of time; they have non-flat frequency responses that can severely compromise spatialization effects; they tend to provide the impression of too close sources, and do not compensate for listener motion unless a tracking system is used. On the other hand there are two main advantages in using headphones: first, they eliminate reverberation of the listening space; second, and probably more important, they allow to deliver distinct signals to each ear, which greatly simplifies the design of 3-D sound rendering techniques. On the contrary loudspeaker based systems suffer from “cross-talk”, i.e. the sound emitted by one loudspeaker will be

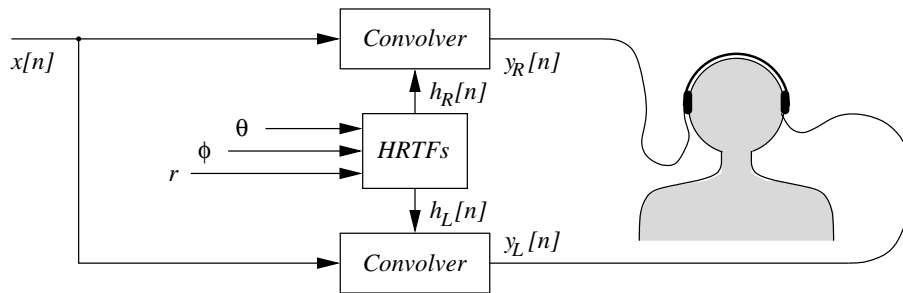


Figure 4.27: Block scheme of a headphone 3-D audio rendering system based on HRTFs.

always heard by both ears. If one ignores the effects of the listening environment, headphone listening conditions can be roughly approximated from stereo loudspeakers using *cross-talk cancellation* techniques, which try to pre-process the stereo signals in such a way that the sound emitted from one loudspeaker is cancelled at the opposite ear. Using these techniques the phantom source can be placed significantly outside of the line segment between the two loudspeakers and in particular elevation effects can be produced. The main problem is that the result will depend on where the listener is relative to the speakers: proper effects are obtained only near the so-called “sweet spot”, a specific listener location assumed by the system.

In this section we will focus on techniques for headphone-based systems.

4.6.1 HRTF-based rendering

The general idea in HRTF-based 3-D audio systems is to use measured HRIRs and HRTFs. Given an anechoic signal and a desired virtual sound source position (θ, ϕ) , a left and right signals are synthesized by convolving the anechoic signal with the corresponding left and right head-related impulse responses. A synthetic block scheme is given in Fig. 4.27. In the remainder of this section we summarize the main steps involved in the development of a HRTF-based 3-D audio system, including HRTF measurement and processing, approximation through synthetic HRTFs, and interpolation.

4.6.1.1 Measuring HRTFs

The typical setting for HRTF measurement is the following: an anechoic chamber, a set of speakers mounted on a geodesic sphere (with a radius of at least one meter in order to avoid near-field effects), at fixed intervals in azimuth and elevation. The listener is at the center of the sphere, with microphones placed in each ear. HRIRs are then measured by playing an analytic signal and recording the corresponding signals produced at the ears, for each desired virtual position.⁹ The listener and the speakers do not need to be moved, facilitating the collection of the measurements. Microphone placing is an issue: it can be placed at the entrance of a plugged ear canal, or near the eardrum to account for the response of the ear canal. Techniques and equipment are explored in order to minimize measurement variability, to improve the signal-to-noise ratio of the measurement hardware used and to determine the optimal placement of the measurement microphone.

In most 3-D sound applications one typically wants to use a single set of HRTFs for every user. One approach might be to use the features of a person who has “desirable” HRTFs, based on some

⁹There is a plethora of sophisticated techniques for Impulse Response estimation, which we do not discuss here

criteria. A set of HRTFs from a good localizer could be used if the criterion were localization performance. An alternative approach is to construct *generalized HRTFs*, that represent the common features of a number of individuals. Binaural impulse responses from many individuals can be “spectrally averaged” in the Fourier domain. However this can cause the resultant HRTF to have diminished spectral features relative to any particular individual’s spectral features. In the extreme case, one person has a 20 dB spectral notch at 8 kHz, and another has a 20 dB spectral peak – the average is no spectral feature at all.

Generalized HRTFs can also be obtained through the use of so-called “dummy heads”, which are mannequins constructed from averaged anthropometric measures and represent standardized heads with average pinnae and torso. The most widely used one is probably the *KEMAR* head (Knowles Electronics Manikin for Auditory Research), although many others are commercially available. Measurements with dummy heads are usually easier, since they are often part of integrated measurement and analysis systems. The low frequency response will be better than with probe mics, since the mic is built into the head; the results will be more replicable since the mic and head remain fixed in position. Moreover, 3-D sound systems based on dummy head HRTFs will be closely matched to recordings made by the same binaural head, allowing compatibility between the two different types of processing. One dummy head might sound more natural to a particular set of users than another, depending on the microphones, the technique used for simulating the ear canal, the head’s dimensions, and so on. The head size (and correspondingly, its diffraction effects and overall ITD) is a major component in the suitability of one dummy head versus another.

4.6.1.2 Post-processing of measured HRTFs

Measured HRTFs undergo a series of processing steps. First, the “blank” portion at the beginning of the impulse response, that results from the time needed for the sound to travel from the speaker to the microphone, is typically discarded. This can be applied to all the impulse responses by investigating the case involving the shortest path, i.e., the measurement position with the ear nearest the loudspeaker. One can even customize the delay inherent to each HRIR pair: by inserting or subtracting blank samples at the start of the impulse response corresponding to the ear furthest from the analytic signal, the overall ITD cue can be customized to a particular head size, or even exaggerated. A second typical procedure is post-equalization of HRTFs to eliminate potential spectral nonlinearities originated from the loudspeaker, the measuring microphone, and the headphones used for playback. As an example, probe microphone are usually small and are especially inefficient at low (< 400 Hz) frequencies, making high-pass filtering or “bass boosting” a fairly common HRTF postequalization procedure. A frequency curve approximating the ear canal resonance, usually derived from some standard equalization, can be applied if it was not part of the impulse response measurement. Since the ear canal resonance is almost independent on the angle of incidence, this needs to be done only once. For most applications, the listener’s own ear canal resonance will be present during headphone listening; this requires removal of the ear canal resonance that may have been present in the original measurement, to avoid a “double resonance”.

A final post-processing procedure is often applied to reduce redundancy in HRTF data. Spectral features that are common to raw HRTFs at all locations do not contain important directional cues, and do not need to be encoded in each single HRTF. Therefore a so-called *Common Transfer Function (CTF)* is often estimated, by computing the mean log-magnitude of the HRTFs measured at several spatial locations. The CTF will include the direction-independent spectral features shared by all HRTFs (e.g., the ear canal resonance). It will also include systematic measurement artifacts, if any. During postprocessing, the CTF can be removed from the raw HRTFs to yield the *Directional*

Transfer Function (DTF). The DTF is a function of θ and ϕ , and is the quantity that contains spectral cues responsible for spatial hearing. Let $C(\omega)$ be the known CTF and $D_{l,r}(\theta, \phi, \omega)$ be the unknown left and right ear DTFs respectively. Then $D_{l,r}$ are estimated from $H_{l,r}$ and C from the equality

$$H_{l,r}(\theta, \phi, \omega) = C(\omega)D_{l,r}(\theta, \phi, \omega). \quad (4.59)$$

The CTF captures the overall structure and dynamic range of the HRTFs, allowing each DTF to operate over a smaller dynamic range. This division allows us to vary a smaller parameter set (corresponding to only the DTF) to achieve space-varying HRTF approximations. Many of the algorithms described in the next sections can be applied either to the “raw” HRTFs or to the DTFs.

It is also known that it is not really necessary to preserve phase information in the interpolated HRTF, as humans are sensitive mostly to the magnitude spectrum for the localization purposes [40] and the measured phase is likely to be contaminated anyway due to difficulties of measuring it accurately because of sampling and other problems.

[40] A. Kulkarni, S. K. Isabelle, and H. S. Colburn (1999). Sensitivity of human subjects to head-related transfer-function phase spectra, *J. Acoust. Soc. Am.*, vol. 105, no. 5, pp. 2821-2840.

Having acquired HRTF magnitude responses, one can design low-order filters that approximate the original HRTFs, in a perceptually motivated way. The resulting filters are *synthetic HRTFs*, and should be perceptually undistinguishable from the measured ones while providing significant computational advantages. Convolution of the sound stimuli with a low-order filter requires little computational resources, while the direct use of measured HRTFs requires a convolution with long FIR filters. The reported duration of measured HRIRs varies across studies: assuming an average duration of ~ 10 ms, the corresponding FIR filter length is ~ 440 samples for a sampling rate of 44.1 kHz. Despite the ever increasing computational power at our disposal, such filter sizes can make it difficult to synthesize complex acoustic environments in real time, particularly when multiple sound sources and reverberant environments have to be rendered.

A perceptually appropriate low-order representation of the HRTFs may also provide insight into sound localization mechanisms. The usefulness of various cues embodied in the HRTF is incompletely understood, and identifying an appropriate simple representation can be used to study attributes that lead to directional perceptions. Moreover, an appropriate low-order model can be used to study the physical mechanisms that produce certain features in the HRTF. Gaining this insight could result in computational methods for generating HRTFs that would not rely upon making empirical measurements from individuals.

We can schematically divide the techniques for deriving synthetic HRTFs into two families. In *pole-zero models* the modeling problem is viewed as one of system identification, which has several classical solutions. One drawback is that the coefficients are usually complicated functions of azimuth and elevation, and have to be tabulated, which destroys the usefulness of the model. *Series expansions* let one represent the HRTF as a weighted sum of simpler basis functions. While this is useful for inspecting the data, the run-time complexity of such models severely limits their usefulness. Below we discuss both approaches.

4.6.1.3 Synthetic HRTFs: pole-zero models

We now briefly discuss the pole-zero modeling approach. For a given direction (θ, ϕ) , we want to approximate the corresponding HRTF, $H(z)$, with a rational transfer function $\hat{H}(z)$ defined as

$$\hat{H}(z) = \frac{b_0 + \sum_{k=1}^q b_k z^{-k}}{1 - \sum_{k=1}^p a_k z^{-k}} = \frac{B(z)}{A(z)}. \quad (4.60)$$

For simplicity, here and in the following we omit in the notation any dependence on (θ, ϕ) : in particular the coefficient vectors $\mathbf{b} = \{b_k\}$, $\mathbf{a} = \{a_k\}$ will depend on θ, ϕ . We will call this a *pole-zero* model (or an *ARMA* model)¹⁰ of the HRTFs. In the particular case $q = 0$, Eq. (4.60) is an *all-pole* model: we have already seen in Chapter *Sound modeling: signal based approaches* that linear prediction can be used in this case to estimate the coefficients $\{a_k\}$ that allow \hat{H} to best approximate H . In the general case $q \geq 1$, we can still re-state the problem as a problem of minimizing some error function. In the hypothesis that the head related impulse response $h[k]$ have length m , the most straightforward choice is to minimize the energy of the difference signal (the Least-Squares Error):

$$E\{h - \hat{h}\} = \sum_{k=0}^m \left(h[k] - \hat{h}[k] \right)^2 = \frac{1}{2\pi} \int_{-\pi}^{\pi} \left| H(e^{j\omega}) - \frac{B(e^{j\omega})}{A(e^{j\omega})} \right|^2 d\omega, \quad (4.61)$$

In practical applications the desired response $H(\omega_k)$ is specified on a set of L “design frequencies” $\omega_k = 2k\pi/LF_s$ and the error to be minimized will have e.g. the form

$$E\{h - \hat{h}\} \sim \frac{1}{L} \sum_{k=0}^{L-1} \left(H(\omega_k) - \hat{H}(\omega_k) \right)^2. \quad (4.62)$$

Minimizing the error $E\{h - \hat{h}\}$ means finding the coefficient vectors \mathbf{b} , \mathbf{a} for which the gradient of $E\{h - \hat{h}\}$ is null, that is solving the set of equations

$$\nabla_{\mathbf{a}} E\{h - \hat{h}\} = \nabla_{\mathbf{b}} E\{h - \hat{h}\} = \mathbf{0}, \quad (4.63)$$

where the notation $\nabla_{\mathbf{x}} E$ stands for the gradient of E with respect to the vector \mathbf{x} . We do not enter into the mathematics involved in writing and solving these equations and refer the reader to the literature on linear Least-Squares Error estimation.

Instead we note that, since our goal is to derive a fit to the perceptually salient features in the HRTF set using a minimal number of model coefficients, an error metric that utilizes absolute LS error on a linear scale is not the best choice, whereas an error criteria based on the ratio of the approximated to desired magnitude of the HRTF across frequency (e.g., the difference in log magnitude) might be perceptually more appropriate. Since both spectral peaks and spectral notches provide relevant information about the sound source location, minimizing the error on a log scale ensures that the solution is not biased toward peaks relative to notches. An example of such a perceptually motivated error criterion is

$$E_{\log}\{h - \hat{h}\} = \frac{1}{L} \sum_{k=0}^{L-1} \left(\ln |H(\omega_k)| - \ln |\hat{H}(\omega_k)| \right)^2, \quad (4.64)$$

where many other variants have been proposed in the literature. A drawback of this kind of log-magnitude response errors is that determining the pole-zero model parameters is a nonlinear problem. A possible approach to solve the resulting equations is using “quasi-Newton” gradient search algorithms. Note also that the error function (4.64) does not consider errors in the phase response. This is not a major issue since listener’s localization accuracy is not significantly degraded when appropriate interaural time delays are used and HRTFs are represented by their minimum phase responses (although it is also true that listeners can hear differences when they are asked to discriminate between signals passed through measured and approximated HRTFs that differ only in their phase spectrum).

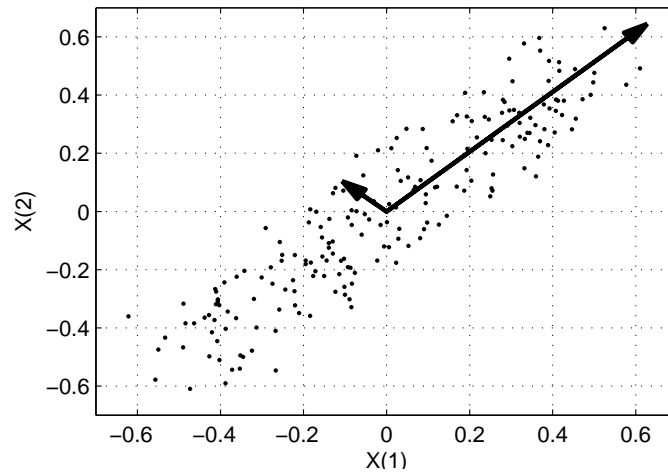


Figure 4.28:

4.6.1.4 Synthetic HRTFs: series expansions

Based on the notions given in Sec. 4.5.1, one can argue on a physical basis that HRTFs should be completely determined by a relatively small number of physical parameters: the average head radius, head eccentricity, maximum pinna diameter, etc. This suggests that the intrinsic dimensionality of the HRTFs might be small, and that their complexity primarily reflects the fact that we are not viewing them correctly.

Among the statistical procedures used to provide a “simpler” representation of a set of correlated measures, a powerful and very popular one is *Principal component analysis (PCA)*, also known as Karhunen-Loève transformation. The central idea of PCA is to reduce the dimensionality of a data set in which there are a large number of interrelated measures, while retaining as much as possible of the variation present in the data. A small set of *basis functions* is derived, and these are used to compute the *principal components*, i.e. the sets of weights that reflect the relative contributions of each basis function to the original data.

Assume we wish to represent M N -dimensional vectors $\mathbf{x}_1 \dots \mathbf{x}_M$ with a 1-dimensional projection (a line) through the sample mean. The line can be written as

$$\mathbf{x} = \mathbf{m} + a\mathbf{e}, \quad (4.65)$$

where \mathbf{e} is a unit vector in the direction of the line, a is a constant coefficient that indicates the distance of any \mathbf{x} from the sample mean $\mathbf{m} = 1/M \sum_{k=1}^M \mathbf{x}_k$. The k th vector \mathbf{x}_k is represented as $\mathbf{m} + a_k\mathbf{e}$, where the optimal coefficient a_k can be obtained by minimizing the “squared error criterion function”

$$E(a_1 \dots, a_k, \mathbf{e}) = \sum_{k=1}^M \|\mathbf{m} - a_k\mathbf{e} - \mathbf{x}_k\|^2. \quad (4.66)$$

For a given direction \mathbf{e} , the optimal coefficients are clearly $a_k = \mathbf{e}^T(\mathbf{x}_k - \mathbf{m})$, i.e. they are obtained by projecting the data vectors onto the line \mathbf{e} that passes through the sample mean. The question is

¹⁰See linear prediction in Chapter *Sound modeling: signal based approaches*.

now: what is the optimal direction e ? By exploiting the expression written above for the optimal a_k 's, the error E can be rewritten after some straightforward algebra as

$$E(a_1 \dots, a_k, e) = \sum_{k=1}^M \|(\mathbf{m} - a_k e) - \mathbf{x}_k\|^2 = \dots = -e^T \mathbf{S} e + \sum_{k=1}^M \|\mathbf{x}_k - \mathbf{m}\|^2, \quad (4.67)$$

where $\mathbf{S} = \sum_{k=1}^M (\mathbf{x}_k - \mathbf{m})(\mathbf{x}_k - \mathbf{m})^T$ is the $N \times N$ *scattering matrix* of the data (which coincides with the covariance matrix except for a multiplying factor $1/(N - 1)$). Therefore minimizing E means maximizing the function $f(e) = e^T \mathbf{S} e$, with the constraint $\|e\| = 1$. This can be done using Lagrange multipliers.¹¹ For our PCA problem we have $\mathcal{L}(e, \lambda) = e^T \mathbf{S} e - \lambda(1 - e^T e)$, and $\nabla_e \mathcal{L}(e, \lambda) = 2\mathbf{S} e - 2\lambda e$. In conclusion the points e that maximize $f(e)$ are those for which

$$\mathbf{S} e = \lambda e, \quad (4.68)$$

i.e. are the eigenvectors of \mathbf{S} corresponding to the eigenvalue λ . The single “best” line that represents the data is found by picking the eigenvector corresponding to the largest eigenvalue of \mathbf{S} so to ensure that $e^T \mathbf{S} e = \lambda$ is maximized.

This can be readily extended to larger dimensions. If we wish to represent the \mathbf{x}_k 's on a q -dimensional hyperplane through the sample mean, written as

$$\mathbf{x} = \mathbf{m} + \sum_{i=1}^q a_i \mathbf{e}_i, \quad (4.69)$$

then we project the data onto the q eigenvectors of \mathbf{S} corresponding to the q largest eigenvalues. If we choose to use all eigenvectors, that is project the data to all eigenvectors and then add them, we will get the original data back (with no dimensionality reduction). From a geometrical standpoint, the eigenvectors represent the principal axes along which the data (and hence the covariance matrix) show largest variance. The weight coefficients a_i are called the *principal components*. Note also that the basis functions are derived in such a way that the first function and its weights capture the majority of common variation present in the data and that the remaining functions and weights reflect decreasing common variation and increasing unique variation. The number q of basis functions required to provide an adequate representation of the data is largely a function of the amount of redundancy or correlation present in the data. The greater the redundancy, the smaller the number of basis functions needed.

Now suppose we have measured directional transfer functions D , on M directions (θ_k, ϕ_k) and on N frequency points: $D(\theta_k, \phi_k, \omega_j)$, $k = 1 \dots M$, $j = 1 \dots N$. We can apply PCA to the particular set of $M N$ -dimensional vectors \mathbf{x}_k constructed as $x_{k,j} = \log |D_k(\theta_k, \phi_k, \omega_j)|$, i.e. we work on the log magnitudes of the DTFs (as already remarked, approximation of log-magnitudes is perceptually more appropriate than approximation of linear magnitudes). The result is a set of q basis vectors \mathbf{e}_i (where $e_{i,j} = e_i(\omega_j)$), such that for the k th direction (θ_k, ϕ_k) the DTF can be approximated as

$$\log |D(\theta_k, \phi_k, \omega_j)| \sim \sum_{i=1}^q a_i(\theta_k, \phi_k) e_i(\omega_j). \quad (4.70)$$

Accurate evaluation of the procedure sketched above would show that the first five basis functions ($q = 5$) are sufficient to accurately represent the magnitudes of the DTF set, and listening tests would

¹¹Reminder: to find the extremum of a function $f(x)$ subject to a constraint $g(x) = 0$, one constructs the Lagrange function $\mathcal{L}(x, \lambda) = f(x) + \lambda g(x)$ and looks for a zero of the gradient $\nabla_x \mathcal{L}(x, \lambda)$.

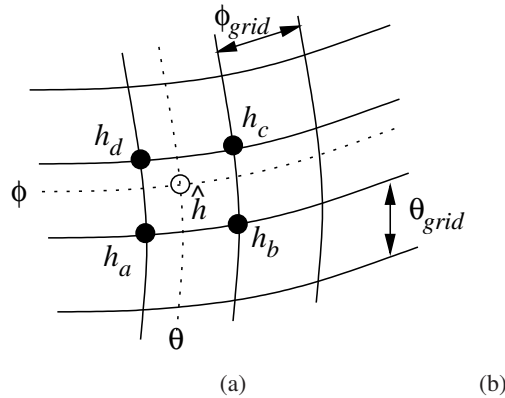


Figure 4.29: (a) Bilinear interpolation, and (b) Interpolation of zeros

show a high correlation between responses to the synthesized and measured conditions. Furthermore, it is possible to relate series expansions to anthropometric measurements and to scale them to account for individual differences .

A similar procedure can also be applied to the complex HRTF rather than to the log magnitude. Expanding the log-magnitude can be viewed as a cascade model, and is most natural for representing head diffraction, ear-canal resonance, and other operations that occur in sequence. Expanding the complex HRTF can be viewed as a parallel model, and is most natural for representing shoulder echoes, pinna echoes, and other multipath phenomena. One drawback of this kind of representations is that they require significant computation for real-time synthesis when head motion or source motion is involved, because the weights a_i are relatively complex functions of azimuth and elevation that must be tabulated. This remark leads us to the topic of HRTF interpolation.

4.6.1.5 Interpolation

HRTF measurements can only be made a finite set of locations, and when a sound source at an intermediate location must be rendered, the HRTF must be *interpolated*. If interpolation is not applied (e.g.. if a nearest neighbor approach is used) audible artifacts like clicks and noise are generated in the sound spectrum when the source position changes.

A straightforward way to perform interpolation directly on the HRIR samples is the bilinear method, which simply consists of computing the response at a given point (θ, ϕ) as a weighted mean of the measured responses associated with the four nearest points. More precisely, if the corresponding set of HRIRs has been measured over a spherical grid with steps θ_{grid} and ϕ_{grid} , the estimate of the HRIR at an arbitrary (θ, ϕ) can be obtained as (see Fig. 4.29(a))

$$\hat{h}[n] = (1 - c_\theta)(1 - c_\phi)h_a[n] + c_\theta(1 - c_\phi)h_b[n] + c_\theta c_\phi h_c[n] + (1 - c_\theta)c_\phi h_d[n], \quad (4.71)$$

where $h_\alpha[n]$ ($\alpha = a, b, c, d$) are the HRIRs associated with the four nearest points to the desired position. The parameters c_θ and c_ϕ are computed as

$$c_\theta = \frac{\theta \bmod \theta_{grid}}{\theta_{grid}}, \quad c_\phi = \frac{\phi \bmod \phi_{grid}}{\phi_{grid}}. \quad (4.72)$$

Several refinements can be applied to this simple technique, in order to improve efficiency. In particular, reduced-order HRIR such as those described earlier in this section can be used, and interpolation

can be performed using only three grid points (those which form a triangle around the desired position). Since some HRTF features arise due to coherent addition or cancelation of reflected and diffracted waves, interpolation may not preserve these features and produce perceptually poor results. Moreover, the interpolating filters are required to be minimum-phase: if this requirement is not satisfied, severe comb-filtering effects in the frequency domain can be produced when the phase delays of the interpolating filters vary considerably. Also, to capture fine details of the HRTF the sampling must be fine enough, i.e. satisfy a Nyquist criterion. Interpolation can be performed in the frequency domain as well (i.e. estimate the DFT of \hat{h} by interpolating the DFTs of the h_α 's). Besides linear approaches, geometric and spline interpolation can be used as well.

If synthetic HRTFs in the form of pole-zero filters are being used, interpolation can be performed on the poles and the zeros themselves. The case of an all-zero filter is relatively straightforward. Suppose that we want to interpolate between two transfer functions $H_\alpha(z)$ ($\alpha = a, b$) of the form

$$H_\alpha(z) = 1 + \sum_{k=1}^q b_{\alpha,k} z^{-k} = \prod_{k=0}^q (1 - c_{\alpha,k} z^{-1}), \quad \alpha = a, b, \quad (4.73)$$

where $b_{\alpha,0} = 1$ without loss of generality, and where we are assuming that the zeros of both filters are sorted according to their phases. Then an interpolated filter $\hat{H}(z) = \prod_{k=0}^q (1 - \hat{c}_k z^{-1})$ can be obtained by (1) pairing the zeros according to angular proximity, and (2) computing $\forall k$ the interpolated zero $\hat{c}_k = (1 - \rho)c_{a,k} + \rho c_{b,k}$. Note that if the H_α are minimum-phase the interpolated filter is also minimum-phase (see also Fig. 4.29(b))

If we use pole-zero synthetic HRTFs, i.e. of the form 4.60 with $p > 0$, then interpolation becomes more complicated. One can still use convex combinations of pole and zero values from neighbouring DTF approximations (note in particular that linear combination of stable poles is guaranteed to be stable). However a naive realization of this approach can result in erratic and occasionally large errors of the interpolated filters. In order to achieve regularity in the interpolation, more refined algorithms are needed that provide pairing and ordering on the entire HRTF database.

Reconstruction of the underlying continuous coefficient functions from the samples obtained is an inherently ill-posed problem because the samples do not uniquely define the functions in the absence of additional assumptions. Furthermore, the samples are usually corrupted by the presence of noise. Regularization theory [Tikhonov and Arsenin 1977] offers a general framework for transforming ill-posed problems to well posed problems through the use of smoothness constraints. Here smoothness constraints imply that a small change in θ, ϕ induces a small change in the coefficients: the HRTFs originate from a physical system of limited spatial extent. Various methods, including linear interpolation, could be used to reconstruct continuous coefficients. Alternatively, spline models [Gu (1989)] can be used.

4.6.2 Structural models

A synthetic block scheme is given in Fig. 4.30.

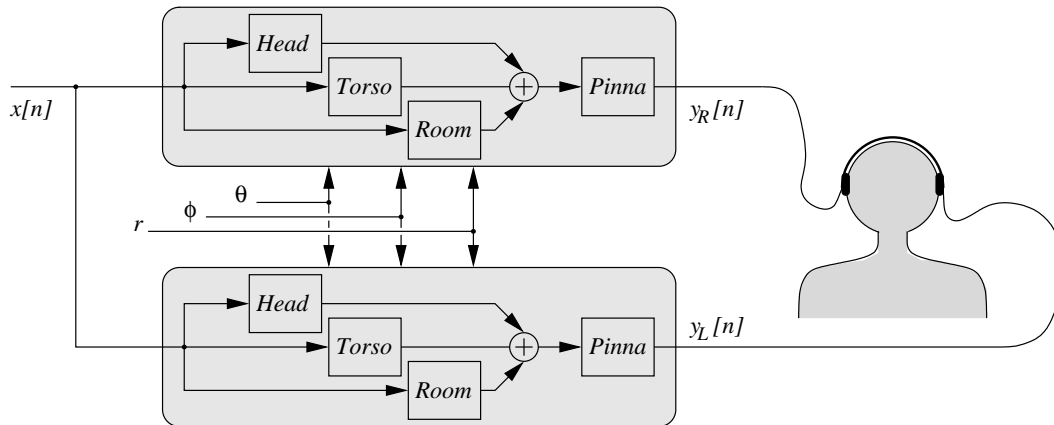


Figure 4.30: Block scheme of a headphone 3-D audio rendering system based on a structural model.

4.7 Commented bibliography

Room acoustics: Wallace C. Sabine has in a way invented the science of concert hall acoustics in the early '900s. For a review of his work and early literature on concert hall acoustics see [Sabine, 1939]. Note that the Paul E. Sabine author of this paper is the cousin of Wallace. A very complete discussion of physical aspects of room acoustics is provided by Kuttruff [1991]: Section 4.2.1 is almost entirely based on this book.

Concerning the research on perceptual attributes of reverberation, the tutorial paper by Beranek [1992] summarizes the main results obtained up to 1992. Research at IRCAM tried to provide a minimal set of independent parameters that give an exhaustive characterization of room acoustic quality [Jot, 1999]. These parameters are divided into three categories, that relate to room perception, source/room interaction, and source perception, respectively.

The first artificial reverberator was proposed by Manfred Schroeder in the early '60's. The reverberator realized in our example M-4.3 is in fact the Schroeder [1962]reverberator. Schroeder also provided a method for measuring the reverberation time [Schroeder, 1965], which can be used to realize the code in example M-4.1. Moreover, Schroeder [1970] proposed the combination of early reflections and late reverberation depicted in our Fig. 4.11(a).

An extensive experimentation on structures for artificial reverberation was conducted by Andy Moorer in the late seventies. He extended the work done by Schroeder in relating some basic computational structures (e.g., tapped delay lines, comb and allpass filters) with the physical behavior of actual rooms. The reverberator realized in our example M-4.4 is in fact the Moorer [1979] reverberator. He also proposed the combination of early reflections and late reverberation depicted in our Fig. 4.11(b).

Gardner [1998] has explored the use of structures based on all-pass and nested all-pass filters (see in particular Figs. 4.9 and 4.10). This reference, together with [Rocchesso, 2002], also provides an extensive overview of reverberation algorithms.

Feedback Delay Networks were first suggested for artificial reverberation by Gerzon [1971, 1972], who noted that several comb filters could “sound good” when cross-coupled. He proposed an orthogonal matrix feedback around a parallel bank of delay lines, as a means of maximizing cross-coupling. Some years later Stautner and Puckette [1982] independently suggested similar ideas and proposed a four-channel FDN reverberator based on the feedback matrix given in our Eq. (4.38). Jot [Jot and

Chaigne, 1991, Jot, 1991, 1997] developed a systematic FDN design methodology allowing largely independent setting of reverberation time in different frequency bands. Rocchesso and Smith [1997] have provided further insights about the structures of feedback matrices in FDNs, and discussed analogies between FDNs and DWNs. General discussions of the use of FDNs for artificial reverberation are provided by Gardner [1998], Rocchesso [2002], Smith [2006]

Waveguide meshes were first studied by Van Duyne and Smith [1993, 1995]. Since then many studies have focused on techniques for reducing dispersion errors. Savioja and Välimäki [2000, 2003] have proposed interpolation and frequency-warping techniques to reduce dispersion as function of both frequency and propagation direction. Fontana and Rocchesso [1998, 2001] have focused on 2-D meshes, and provided results both about applications to membrane modeling and about general numerical aspects: they compared square, triangular, and hexagonal meshes in terms of sampling efficiency and dispersion error. Bilbao [2004] has also investigated in details many numerical and computational properties of the waveguide mesh, in particular he analyzed dispersion properties of various mesh topologies using von Neumann analysis and he provided a unified view of the digital waveguide mesh and wave digital filters as particular classes of energy invariant finite difference schemes. Finally, another topic addressed in the literature is the design of mesh boundaries, with a special focus on modeling diffusion. This problem was addressed by Laird et al. [1999], and later by Lee and Smith [2004], who used quadratic residue sequences to design maximally diffusing boundaries.

References

- Leo L. Beranek. Concert hall acoustics. *J. Acoust. Soc. Am.*, 92(1):1–39, July 1992.
- Stefan Bilbao. *Wave and Scattering Methods for Numerical Simulation*. John Wiley and Sons, Inc., New York, 2004.
- Federico Fontana and Davide Rocchesso. Signal-Theoretic Characterization of Waveguide Mesh Geometries for Models of Two-Dimensional Wave Propagation in Elastic Media. *IEEE Trans. Speech Audio Process.*, 9(2):152–161, Feb. 2001.
- Federico Fontana and Davide Rocchesso. Physical modeling of membranes for percussion instruments. *Acta Acustica united with Acustica*, 84(3):529–542, May 1998.
- William G. Gardner. Reverberation algorithms. In Mark Kahrs and Karl-Heinz Brandenburg, editors, *Applications of Digital Signal Processing to Audio and Acoustics*, pages 85–131. Kluwer Academic Publishers, New York, Mar. 1998.
- Michael A. Gerzon. Synthetic stereo reverberation, Part I. *Studio Sound*, 13:632–635, Dec. 1971.
- Michael A. Gerzon. Synthetic stereo reverberation, Part II. *Studio Sound*, 14:24–28, Jan. 1972.
- Jean-Marc Jot. An analysis/synthesis approach to real-time artificial reverberation. In *Proc. IEEE Int. Conf. Acoust. Speech and Signal Process.*, volume 2, pages 221–224, S. Francisco, Feb. 1991.
- Jean-Marc Jot. Efficient models for reverberation and distance rendering in computer music and virtual audio reality. In *Proc. Int. Computer Music Conf.*, pages 236–243, Thessaloniki, 1997.
- Jean-Marc Jot. Real-time spatial processing of sounds for music, multimedia, and interactive human-computer interfaces. *Multimedia Systems*, 7(1):55–69, Jan. 1999.
- Jean-Marc Jot and Antoine Chaigne. Digital delay networks for designing artificial reverberators. In *Proc. Audio Engineering Society Convention*, Paris, Feb. 1991. Preprint 3030.
- Heinrich Kuttruff. *Room Acoustics*. Elsevier Applied Science, London and New York, 3rd edition, 1991.
- Joel Laird, Paul Masri, and Nishan Canagarajah. Modelling diffusion at the boundary of a digital waveguide mesh. In *Proc. Int. Computer Music Conf.*, pages 492–495, Beijing, Oct. 1999.



- Kyogu Lee and Julius O. Smith. Implementation of a highly diffusing 2-D digital waveguide mesh with a quadratic residue diffuser. In *Proc. Int. Computer Music Conf.*, Miami, Nov. 2004.
- Jame A. Moorer. About this reverberation business. *Computer Music J.*, 3(2):13–18, Summer 1979.
- Davide Rocchesso. Spatial effects. In Udo Zölzer, editor, *Digital Audio Effects*, pages 137–200. John Wiley & Sons, Chirchester Sussex, UK, 2002.
- Davide Rocchesso and Julius O. Smith. Circulant and elliptic feedback delay networks for artificial reverberation. *IEEE Trans. Speech Audio Process.*, 5(1):51–63, Jan. 1997.
- Paul E. Sabine. Architectural acoustics: Its past and its possibilities. *J. Acoust. Soc. Am.*, 11(1):21–28, July 1939.
- Lauri Savioja and Vesa Välimäki. Reducing the dispersion error in the digital waveguide mesh using interpolation and frequency-warping techniques. *IEEE Trans. Speech Audio Process.*, 8(2):184–194, Mar. 2000.
- Lauri Savioja and Vesa Välimäki. Interpolated rectangular 3-d digital waveguide mesh algorithms with frequency warping. *IEEE Trans. Speech Audio Process.*, 11(6):783–790, Nov. 2003.
- Manfred R. Schroeder. Natural-sounding artificial reverberation. *J. Audio Eng. Soc.*, 10(3):219–233, July 1962.
- Manfred R. Schroeder. New method of measuring reverberation time. *J. Acoust. Soc. Am.*, 37(6):1187–1188, June 1965.
- Manfred R. Schroeder. Digital simulation of sound transmission in reverberant spaces. *J. Acoust. Soc. Am.*, 47(2):424–431, Feb. 1970.
- Julius O. Smith. *Physical Audio Signal Processing: for Virtual Musical Instruments and Digital Audio Effects, August 2006 Edition*. <http://ccrma.stanford.edu/~jos/pasp/>, 2006.
- John Stautner and Miller Puckette. Designing multichannel reverberators. *Computer Music J.*, 3(2):52–65, 1982. Reprinted in *The Music Machine*, Curtis Roads (Ed.). Cambridge, The MIT Press, 1989. (pp. 569–582).
- Scott A. Van Duyne and Julius O. Smith. Physical modeling with the 2-d digital waveguide mesh. In *Proc. Int. Computer Music Conf.*, pages 40–47, Tokio, 1993.
- Scott A. Van Duyne and Julius O. Smith. The tetrahedral digital waveguide mesh. In *Workshop Appl. Signal Process. to Audio and Acoust.*, pages 234–237, Mohonk, Oct. 1995.



Contents

4	Sound in space	4.1
4.1	Introduction	4.1
4.2	Reverberation: physical and perceptual background	4.2
4.2.1	Basics of room acoustics	4.2
4.2.1.1	Sound waves in a closed space	4.2
4.2.1.2	Sound sources and room impulse responses	4.4
4.2.1.3	Reverberation time	4.5
4.2.1.4	Geometrical room acoustics	4.7
4.2.2	Perceptual reverberation parameters	4.9
4.2.2.1	Reverberance	4.10
4.2.2.2	Early reflections and spatial impression	4.11
4.2.2.3	Clarity	4.12
4.2.2.4	Other perceptually relevant parameters	4.13
4.3	Algorithms for synthetic reverberation: the perceptual approach	4.14
4.3.1	Approximating late reverberation	4.15
4.3.1.1	Recirculating delays	4.15
4.3.1.2	Tuning the parameters	4.15
4.3.2	Improved structures	4.17
4.3.2.1	Low-pass combs	4.17
4.3.2.2	Nested all-pass filters	4.17
4.3.2.3	Adding early reflections	4.19
4.4	Multidimensional reverberation structures	4.20
4.4.1	Feedback delay networks	4.20
4.4.1.1	A n-D generalization of the recursive comb filter	4.20
4.4.1.2	A general FDN reverberators	4.22
4.4.1.3	Designing the lossless prototype	4.23
4.4.1.4	Designing lossy components	4.24
4.4.2	Digital waveguide networks	4.26
4.4.2.1	The link between FDNs and DWNs	4.26
4.4.2.2	General lossless scattering matrices	4.27
4.4.2.3	Waveguide meshes	4.28
4.5	Spatial hearing	4.29
4.5.1	The sound field at the eardrum	4.30
4.5.1.1	Head	4.30
4.5.1.2	The external ear	4.32
4.5.1.3	Torso and shoulders	4.33

4.5.1.4	Head-related transfer functions	4.34
4.5.2	Perception of sound source location	4.36
4.5.2.1	Azimuth perception	4.36
4.5.2.2	Lateralization and externalization	4.37
4.5.2.3	Elevation perception	4.38
4.5.2.4	Distance perception	4.39
4.5.2.5	Dynamic cues	4.41
4.6	Algorithms for 3-D sound rendering	4.42
4.6.1	HRTF-based rendering	4.43
4.6.1.1	Measuring HRTFs	4.43
4.6.1.2	Post-processing of measured HRTFs	4.44
4.6.1.3	Synthetic HRTFs: pole-zero models	4.45
4.6.1.4	Synthetic HRTFs: series expansions	4.47
4.6.1.5	Interpolation	4.49
4.6.2	Structural models	4.50
4.7	Commented bibliography	4.51