

1 Introduction

... things inanimate have mov'd,
And, as with living Souls, have been inform'd,
By Magick Numbers and persuasive Sound.

—William Congreve (1697) *The Mourning Bride*

The ear is a most complex and beautiful organ. It is the most perfect *acoustic*, or hearing instrument, with which we are acquainted, and the ingenuity and skill of man would be in vain exercised to imitate it.

—John Frost (1838), *The Class Book of Nature: Comprising Lessons on the Universe, the Three Kingdoms of Nature, and the Form and Structure of the Human Body*

Would it truly be in vain to exercise our ingenuity to imitate the ear? It would have been, in the 1800s—but now we are beginning to do so, using the “magick” of numbers. Machines imitating the ear already perform useful services for us: answering our queries, telling us what music is playing, locating gunshots, and more. By imitating ears more faithfully, we will be able to make machines hear even better. The goal of this book is to teach readers how to do so.

Understanding how humans hear is the primary strategy in designing machines that hear. Like the study of vision, the study of human hearing is ancient, and has enjoyed impressive advances in the last few centuries. The idea of *machines* that can see and hear also dates back more than a century, though the computational power to build such machines has become available only in recent decades. It is now, as they say in the computer business, a simple matter of programming. Well, not quite—there is still work to be done to firm up our understanding of sound analysis in the ear, and yet more to be done to understand the enormous capabilities of the human brain, and to abstract these understandings to better support machine hearing. So let's get started.

Humans tend to take hearing for granted. We are so aware of what's going on around us, largely by extracting information from sound, yet so unable to describe or appreciate how we do it. Can we make machines do as well at interpreting their world, and ours, through sound? We can, if we leverage scientific knowledge of how humans process sound.

Being able to produce and analyze sound waves is a prerequisite to developing a better understanding of hearing. Early progress in the field was made with the help of analytical instruments such as Helmholtz's resonators and recording devices, like the waveform drawing device in Figure 1.1, and controlled sound production instruments such as Seebeck's siren, shown in Figure 1.2. Representing such waves as electrical

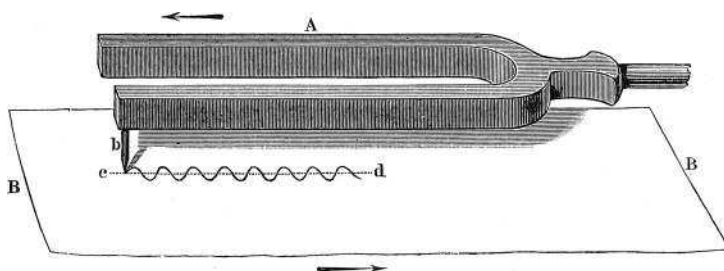


Figure 1.1 Helmholtz explained the idea of a sound's waveform via this diagram of a tuning fork with a stylus point attached, drawing its vibration on a moving piece of paper.

signals has been routine since the invention of the telephone. We now have a myriad of machines that help us generate, compress, communicate, store, reproduce, and modify sound signals, in ways tuned to how we hear. For most of these applications, though, the machines remain “deaf,” in that they get very little meaning out of the sounds they process.

What if you had a device at home, always listening to what's going on? Could it tell what interesting things it heard while you were out? Could it tell you the refrigerator sounds like it's wearing out? Would it understand if you asked it a question? Could it find you some music to listen to if you described your mood? Could it listen to you and determine your mood itself? Could it say where a mouse might be hiding because it heard it run there? Could it distinguish between normal household sounds and an anomaly in the dead of night? Could it also be your intelligent answering machine, and tell you who called, and why, based on hearing their voice? Of course it could.

Who might make such a machine? What crazy functionality might they give a machine that could hear and understand sounds? Have we chosen the best path through the complex web of theories about hearing? Can we do better on some tasks by modifying the approach? What advances in the study of human hearing might we discover while trying to put our theories to the test of real use? These are the kinds of ideas and questions about sound and hearing that have been going around in my head for decades—and that we are getting some answers on recently. I've worked on spatial effects in music and games, and on machines to synthesize and recognize speech and music, and on other fun things to do with sound. Where most others deal with sounds by various conventional or ad hoc methods, I keep coming back to how the ear would do it—and this approach has proved fruitful.

There is enough known about how the ear and hearing work that we have gotten serious about putting this knowledge to practical uses. Starting with the anatomy, we model the structure and function of the ear and the auditory nervous system; using physiological and psychophysical techniques, we figure out what the brain gets from the ear, and how it deals with the information to perform meaningful tasks. Then we program computing machines to do similarly, based on this knowledge. In essence, we mimic the biology.



Figure 1.2 A make-it-yourself acoustic siren, much like August Seebeck's, as shown by Alfred M. Mayer (1878). The spinning disk, driven from a crank via string and pulleys, interrupts a stream of air from the tube to make waves of sound pressure that we hear as a tone. Different tones can be made by moving the tube to a different row of holes, or by changing the disk to one with a different pattern of holes. August Seebeck and Hermann von Helmholtz were among the nineteenth-century scientists who used such devices in their research that contributed to connecting the physical and perceptual properties of musical tones to the mechanisms of human hearing—though their theories were somewhat in opposition to each other.

Today we have access to massive quantities of sound, to analyze, organize, index, and learn from. The soundtracks of YouTube videos alone have hundreds of millions of hours of sound, and so far our computers are rather ignorant of what those soundtracks are trying to communicate. Imagine what value there might be in having our machines just listen to them and understand. Speech, music, laughing babies, sounds of interesting events, activities, places, and personalities—it's all there to be discovered, categorized, indexed, summarized, remembered, and retrieved.

The full scope of machine hearing will reveal itself as people discover that it is relatively easy to have machines understand sounds of all sorts, and people find

imaginative uses for such machines. Elephant infrasound hearing and bat ultrasound hearing and echolocation suggest that the same basic strategies have been put to many purposes by other mammals. We might include other sonic applications—such as medical imaging—that use sound waves but don't rely on anything about sound perception. At Schlumberger Research in the 1980s, we experimented with hearing techniques applied to the analysis of underground sonic waves. Any far-out infrasound through ultrasound applications that can benefit from the use of techniques like those evolved by humans fall within the scope of what we're trying to teach via this book.

As we get more people engaged in machine hearing, there will be more good ideas and more things we can take on. The potential is enormous, and the scope broad.

1.1 On Vision and Hearing *à la* David Marr

The pioneering vision scientist David Marr was a big influence on my approach to modeling hearing. When I visited him at MIT in 1979 to show him what I was working on, he was very encouraging of the approach. Twisting his words, from vision to hearing, illustrates how his thinking influenced mine:

What does it mean to hear? The plain man's answer (and Aristotle's, too) would be, to know what is where by listening. In other words, hearing is the *process* of discovering from sounds what is present in the world and where it is.

Hearing is therefore, first and foremost, an information processing task, but we cannot think of it just as a process. For if we are capable of knowing what is where in the world, our brains must somehow be capable of representing this information—in all its profusion of color and form, beauty, motion, and detail.—modified from *Vision*, David Marr (1982)

I honor Marr's introduction to his ground-breaking book *Vision* in the quotation above, having changed *see* to *hear*, *looking* to *listening*, *vision* to *hearing*, and *images* to *sounds*. I've left the last phrase unchanged, as I believe that "*color and form, beauty, motion, and detail*" is a much more apt description of what our brains extract and represent about sound than the usual more pedestrian properties of *loudness*, *pitch*, and *timbre*.

Marr's computational and representational approach to vision helped to define the vibrant field of computer vision, or machine vision as it's also called, more than thirty years ago. My book is motivated by the feeling that something along these lines is still needed in the hearing field. It's a daunting challenge to try to live up to David Marr, even if I've had a few extra decades to prepare, but it's time to give it a shot.

Compared to other mammals, humans have put vision to some very special applications, like reading written language, and analogously have put hearing to use in spoken language and in music. These pinnacle applications should not exclusively drive the study of vision and hearing, however, and perhaps are best addressed only after low-level preliminaries are well understood, and more general applications are under control. Therefore, we focus on these more general and lower-level aspects, and on broader

applications of hearing, as Marr focused on the more general aspects of vision. At the end, we come back and touch on applications in speech and music.

David Mellinger (1991) should be credited with helping drive this approach via his dissertation, pointing out that “Advances in machine vision have long stemmed from a physiological approach where researchers have been heavily influenced by Marr’s computational theory. Perhaps the same transfer will begin to happen more in machine hearing.” But this transfer has been incomplete, so we need to drive it some more.

Martin Cooke (1993) has provided an excellent review of Marr’s approach to vision and its influence on work in speech and hearing. Marr’s identification of three levels at which the sensory system is to be understood—*function*, *process*, and *mechanism*, also described as *computation*, *algorithm*, and *implementation*—certainly does help us organize our study of hearing. In an interesting twist, Peter Dallos (1973) used a similar division of concerns into function, mode of operation, and anatomy to describe the auditory periphery, before Marr’s work. His scheme is still used this way and credited in current hearing books (Yost, 2007), as shown in Figure 1.3.

Cooke reviews several applications of Marr’s levels and principles to speech processing, but provides relatively little connection to hearing. The repurposing of Marr’s *primal sketch* concept into a *speech sketch*, by Green and Wood (1986), points up a disconnect: Marr didn’t go from primitive images directly to reading, and we shouldn’t go from primitive sound representations straight to speech; *primal* should imply a much lower level. A sketch is a “sparsified” version of an image, which may be used as part of a feature extraction strategy at the input to a learning system, as described in Section 25.7.

In vision, objects and images must be analyzed at many different scales. Referring to Marr, Andy Witkin (1983) said, “The problem of scale has emerged consistently as a fundamental source of difficulty, because the events we perceive and find meaningful vary enormously in size and extent. The problem is not so much to eliminate fine-scale noise, as to separate events at different scales arising from distinct physical processes.” In hearing, we have the same issue, especially in the temporal dimension, where sounds have periodicities and structure on all time scales.

The idea of an “auditory primal sketch” has been introduced by Neil Todd (1994) as a way to represent the rhythm and temporal structure of music and speech. I had published a related idea on multiscale temporal analysis, as part of a speech recognition approach (Lyon, 1987). Both of these are based on Witkin’s scale-space filtering, which was descended from Marr. Both fall far short of a comprehensive framework for machine hearing, but help to inspire some of the sorts of representations that we will be working with.

Albert Bregman (1990), in his book *Auditory Scene Analysis: The Perceptual Organization of Sound*, discusses how aspects of hearing are valued from an evolutionary perspective, yielding certain advantages of hearing over vision. The auditory system evolved in a context in which better understanding of meaning from an auditory scene—better answers to *what* and *where*—led to a better chance of survival. When I refer to *human hearing* in my title, I mean to include the cortical-level processing

Gross division	Outer ear	Middle ear	Inner ear	Central auditory nervous system
Anatomy				
Mode of operation	<i>Air vibration</i>	<i>Mechanical vibration</i>	<i>Mechanical, Hydrodynamic, Electrochemical</i>	<i>Electrochemical</i>
Function	<i>Protection, Amplification, Localization</i>	<i>Impedance matching, Selective oval window stimulation, Pressure equalization</i>	<i>Filtering distribution, Transduction</i>	<i>Information processing</i>

Figure 1.3 Ear diagram by Yost (2007). While the anatomy and modes of operation are important, we are most interested in emulating the *function*, described in the bottom row. The *information processing* in the central nervous system—the bit where meaning is extracted—is the part that remains most open to exploration and speculation. [Figure 6.1 (Yost, 2007) reproduced with permission of William A. Yost.]

systems that have evolved to handle speech, music, and other big-brain functions; but I do not mean to diminish the importance of the lower levels of auditory processing—in the ear, the brainstem, and the midbrain—that underlie the exquisite hearing capabilities of our pets (and pests), and that form the basis for robust representations of sound from which actionable information can be extracted. Even animals that don't normally use speech can learn to reliably recognize their own names, and discriminate

them against other speech sounds; for example, Shepherd (1911) taught four raccoons that their names were Jack, Jim, Tom, and Dolly.

We can question Marr's insistence that a symbolic representation or *description* be generated (Hacker, 1991). Some approaches to machine hearing systems successfully use representations that remain completely abstract and nameless until the final output—the information that the system is trained to extract—with intermediate steps being subsumed in the learning system. Other approaches will use explicit and named concepts, such as objects, events, musical instruments, notes, talkers, and so forth, that artificial intelligence systems can reason with. Different theories of mind, or different computational frameworks that we have available, will bias our machine hearing applications one way or the other. We are not yet in a position to say which way is likely to be more fruitful for any given area, and hope to encourage exploration in all such directions.

Comments on hearing's analogy with vision are not new. For example, in 1797, the effect of auditory masking on sensitivity was observed and compared to visual masking effects in "annotations" on Perrole's "Philosophical Memoir" on sound transmission (Perrole, 1797):

Sounds seem more intense, and are heard to a greater distance, by night than by day. . . . It is a practical question of some importance to ascertain whether this difference may arise from the different state of the air, the greater acuteness of the organ, or the absence of the ordinary noises produced in the day. By attentive listening to the vibrations of a clock in the night, and remarking the difference between the time when no other noise was heard, and when a coach passed along, it has appeared clear to me that this difference arises from the greater or less stillness only, and that no voluntary effort or attention can render the near sound much more audible, while another noise acts upon the organ. In this situation the ear is nearly in the state of the eye, which cannot perceive the stars in the day time, nor an object behind a candle.

In that memoir, Perrole also introduced the term *timbre* from the French to explain what he meant by *tone* in English: "The tone (*timbre*) was changed in the water in a striking manner." This "catch-all" term, as it has been called, captures everything about what a sound "sounds like," except for its pitch and loudness—sort of like *texture* in vision, which captures much of what shape, size, and brightness don't. It is the job of our machine hearing systems to map timbre (along with pitch and loudness and direction, and their evolution and rhythm over time) into useful information about what the sound represents, be it speech, music, environmental noises, or evidence of mundane or exceptional events.

1.2 Top-Down versus Bottom-Up Analysis

Top-down processing evaluates sensory evidence in support of hypothesized interpretations (meaning), while bottom-up processing converts sensory input to ever-higher-level representations that drive interpretation. Real systems are not necessarily at either extreme, but the distinction can be useful.

Marr says, with respect to general-to-specific (or coarse-to-fine) stereo matching approaches (Marr, 1982),

Nomenclature: What to Call This Endeavor

The terms *computer vision* and *machine vision* are in wide use, not quite interchangeably, the former having a more computer-science connotation, and the latter a more industrial or applications connotation. Terms like *computer hearing*, *computational hearing*, and *computer listening* seem awkward to me, especially since I spent a lot of years building analog electronic models of hearing, probably not qualifying as computers. And what about *listening* or *audition* as a better analogy to *vision*? Several of these terms have overloaded meanings: we can convene a hearing, or perform in an audition, or plant listening devices. The term *machine listening* is sometimes used, but mostly in connection with music listening and performance.

The term *machine hearing* has a strong history at Stanford's computer music lab, CCRMA. In their 1992 progress report, Bernard Mont-Reynaud (1992) wrote a section on machine hearing, which noted that "The purpose of this research is to design a model of Machine Hearing and implement it in a collection of computer programs that capture essential aspects of human hearing including source formation and selective attention to one source (the 'cocktail party problem') without tying the model closely to speech, music, or other domain of sound interpretation."

We hope that by calling the space of computer applications of sound analysis *machine hearing*, following Mont-Reynaud, we will leverage this good name and good direction, and help the field build around a good framework, as Marr did with what we refer to as *machine vision*.

This type of approach is typical of the so-called top-down school of thought, which was prevalent in machine vision in the 1960s and early 1970s, and our present approach was developed largely in reaction to it. Our general view is that although some top-down information is sometimes used and necessary, it is of only secondary importance in early visual processing.

Here we totally agree. Although I have nothing but respect for the strong case for the power of top-down information and expectations in human hearing (Slaney, 1998; Huron, 2006), and though there are prominent "descending" pathways at all levels of the auditory nervous system (Schofield, 2010), my understanding is that the more extensive and complex feedback is within the cortical levels of the central nervous system, and that early audition, like early vision, is best conceived as a modular set of mostly feed-forward bottom-up processing modules. There is feedback, to be sure, but its function can often be treated as secondary, as Marr says. At some levels, feedback may be about parameter learning and optimization; from cortex to thalamus, top-down projections may be about attention. These are important, but not where we start, especially in "early" layers as Marr says.

In the mammalian brain, these early hearing modules include the periphery (the ear) as well as auditory structures in the brainstem and midbrain, and maybe even some stages of cortical processing, such as primary auditory cortex. These levels were successful stable subsystems long before the evolution of the big neocortex that led to

speech and music. The “near decomposability” condition (Simon, 1981) is what allows complex systems to evolve. That’s why we rely so much on data from bottom-up experiments in animals to help us understand human hearing; we accept that the amazing abilities of humans evolved on top of these stable mammalian subsystems, which are themselves not so different from reptilian, bird, and even fish auditory systems.

Like Marr, we are partly reacting to an overreliance on top-down information in sound processing systems. For example, automatic speech recognition (ASR) systems have been gradually improved over the years by reliance on larger and more complex language models and by statistical models that can capture complex prior distributions, while their front-end processing remains relatively stagnant, stuck with spectro-temporal approaches that have no way to improve in terms of robustness to noise and interference, since they don’t represent the aspects of sound that help our auditory systems tease sound mixtures apart. Such problems demand that we understand hearing better, and build systems that can hear and understand multiple sounds at once; how else can we expect a speech recognizer to give us a transcript of a boisterous meeting? Of course, good prior distributions from top-down information will continue to play an important role, too.

Is the auditory system *complex*? Herb Simon (1981) characterizes a complex system this way:

In such systems, the whole is more than the sum of the parts, not in an ultimate, metaphysical sense, but in the important pragmatic sense that, given the properties of the parts and the laws of their interaction, it is not a trivial matter to infer the properties of the whole.

I think this applies to the auditory system as a whole, when the cortex is included, especially in a living organism in which the auditory system is interacting with visual, motor, and other systems, with strong top-down and feedback effects. But for the various bottom-up modules of lower-level auditory processing, perhaps the system is merely *complicated*, but not so complex that we can’t describe its function, and its process, in terms of its mechanisms. I think this is how Marr saw early vision, too. Otherwise, it would be hard to be optimistic about our ability to assemble machines to do similar jobs.

1.3 The Neuromimetic Approach

A strategic element of our machine hearing approach is to respect the representation of sounds on the auditory nerve, which involves both a *tonotopic* (arranged by frequency) organization and detailed temporal structure, as extracted by the rather nonlinear inner ear. At this level, the approach can be said to be *neuromimetic* (Jutten et al., 1988), or *neuromorphic* (Mead, 1990), in the sense that we may be building a copy of a complicated neural system, mimicking its function—or mimicking its structure when we can’t quite describe the function. In the neuromorphic case, copying the structure of the neural system, the expectation is that the structure will have an appropriate *emergent*

behavior and therefore a useful information-processing function. Here *emergent* means that the behavior is not explicitly designed in, but *emerges* from the simpler behaviors of the lower-level elements as a consequence of the structural pattern of interconnection of those elements (Bar-Yam, 1997).

This neuromimetic approach is somewhat distinct from the Marr approach, but sometimes a useful supplement. When a system built this way is found to have a useful function due to its emergent behavior, it can sometimes be further analyzed, and the important parts of its function abstracted, described, and reengineered more efficiently. I believe we are part of the way through this process with neuromimetic hearing front ends. At the level of the cochlea, for example, the function is largely understood, but the description is still as much structural as functional. We do not have the clean separation of function, process, and mechanism that Marr recommended, but we do have a structure for which we can understand the function.

Beyond the cochlea, we still have a mixed structural and functional view, though it is somewhat speculative, of what the function is—the little “information processing” box in the lower right corner of Bill Yost’s diagram, Figure 1.3, is where we ultimately extract meaning. We have pretty good ideas from physiological data about what kinds of auditory images are formed in the brainstem. The main thing we use that is neuromorphic is the very idea of an auditory image: a neural pathway with two spatial dimensions, like the optic nerve from the retina, projecting a time-varying pattern to a two-dimensional sheet of cortical tissue, the primary auditory cortex, for further processing.

An early proponent of a neuromimetic, or *bionic*, approach to machine hearing systems was John L. Stewart (1963), who published a number of reports, papers, patents, and a book on the topic in the 1960s and 1970s. He explains the reasoning behind this approach (Stewart, 1979):

The model becomes an intermediary—a surrogate reality. . . . It is my belief that effective explanations for the traits of living organisms demand the construction of models which behave as do their living counterparts. For, in no other way can the research be disciplined to produce an effective holistic theory!

Stewart (1979) anticipated much of our current approach, including a cochlear transmission-line analog with nonlinearities, a “neural-like analyzer” stage following the cochlea (Stewart, 1966), and the idea of efferent (feedback) adaptation to conditions, via coupled frequency-dependent gain control (Stewart, 1967).

1.4 Auditory Images

In our approach to hearing, we incorporate the notion of an *auditory image*: a presumed representation developed in the subcortical parts of the auditory nervous system (cochlea, brainstem, and midbrain), projecting to primary auditory cortex in the same way that the retinal image projects to primary visual cortex. This approach brings together the strategies of Marr with the two-dimensional neural circuits of the *place*