

man and *cop* differ in formality.) No one really knows why languages are so stingy with words and profligate with meanings, but children seem to expect it (or perhaps it is this expectation that causes it!), and that helps them further with the *gavagai* problem. If a child already knows a word for a kind of thing, then when another word is used for it, he or she does not take the easy but wrong way and treat it as a synonym. Instead, the child tries out some other possible concept. For example, Markman found that if you show a child a pair of pewter tongs and call it *biff*, the child interprets *biff* as meaning tongs in general, showing the usual bias for middle-level objects, so when asked for "more biffs," the child picks out a pair of plastic tongs. But if you show the child a pewter cup and call it *biff*, the child does not interpret *biff* as meaning "cup," because most children already know a word that means "cup," namely, *cup*. Loathing synonyms, the children guess that *biff* must mean something else, and the stuff the cup is made of is the next most readily available concept. When asked for more *biffs*, the child chooses a pewter spoon or pewter tongs.

Many other ingenious studies have shown how children home in on the correct meanings for different kinds of words. Once children know some syntax, they can use it to sort out different kinds of meaning. For example, the psychologist Roger Brown showed children a picture of hands kneading a mass of little squares in a bowl. If he asked them, "Can you see any sibbing?," the children pointed to the hands. If instead he asked them, "Can you see a sib?," they point to the bowl. And if he asked, "Can you see any sib?," they point to the stuff inside the bowl. Other experiments have uncovered great sophistication in children's understanding of how classes of words fit into sentence structures and how they relate to concepts and kinds.

So what's in a name? The answer, we have seen, is, a great deal. In the sense of a morphological product, a name is an intricate structure, elegantly assembled by layers of rules and lawful even at its quirk-iest. And in the sense of a listeme, a name is a pure symbol, part of a cast of thousands, rapidly acquired because of a harmony between the mind of the child, the mind of the adult, and the texture of reality.

6



The Sounds of Silence

When I was a student I worked in a laboratory at McGill University that studied auditory perception. Using a computer, I would synthesize trains of overlapping tones and determine whether they sounded like one rich sound or two pure ones. One Monday morning I had an odd experience: the tones suddenly turned into a chorus of screaming munchkins. Like this: (beep boop-boop) (beep boop-boop) (beep boop-boop) HUMPTY-DUMPTY-HUMPTY-DUMPTY-HUMPTY-DUMPTY (beep boop-boop) (beep boop-boop) HUMPTY-DUMPTY-HUMPTY-DUMPTY-HUMPTY-HUMPTY-DUMPTY-DUMPTY (beep boop-boop) (beep boop-boop) (beep boop-boop) HUMPTY-DUMPTY (beep boop-boop) HUMPTY-HUMPTY-HUMPTY-DUMPTY (beep boop-boop). I checked the oscilloscope: two streams of tones, as programmed. The effect had to be perceptual. With a bit of effort I could go back and forth, hearing the sound as either beeps or munchkins. When a fellow student entered, I recounted my discovery, mentioning that I couldn't wait to tell Professor Bregman, who directed the laboratory. She offered some advice: don't tell anyone, except perhaps Professor Poser (who directed the psychopathology program).

Years later I discovered what I had discovered. The psychologists

Robert Remez, David Pisoni, and their colleagues, braver men than I am, published an article in *Science* on “sine-wave speech.” They synthesized three simultaneous wavering tones. Physically, the sound was nothing at all like speech, but the tones followed the same contours as the bands of energy in the sentence. “Where were you a year ago?” Volunteers described what they heard as “science fiction sounds” or “computer bleeps.” A second group of volunteers was told that the sounds had been generated by a bad speech synthesizer. They were able to make out many of the words, and a quarter of them could write down the sentence perfectly. The brain can hear speech content in sounds that have only the remotest resemblance to speech. Indeed, sine-wave speech is how mynah birds fool us. They have a valve on each bronchial tube and can control them independently, producing two wavering tones which we hear as speech.

Our brains can flip between hearing something as a bleep and hearing it as a word because phonetic perception is like a sixth sense. When we listen to speech the actual sounds go in one ear and out the other; what we perceive is *language*. Our experience of words and syllables, of the “b”-ness of *b* and the “ee”-ness of *ee*, is as separable from our experience of pitch and loudness as lyrics are from a score. Sometimes, as in sine-wave speech, the senses of hearing and phonetics compete over which gets to interpret a sound, and our perception jumps back and forth. Sometimes the two senses simultaneously interpret a single sound. If one takes a tape recording of *da*, electronically removes the initial chirplike portion that distinguishes the *da* from *ga* and *ka*, and plays the chirp to one ear and the residue to the other, what people hear is a chirp in one ear and *da* in the other—a single clip of sound is perceived simultaneously as *d*-ness and a chirp. And sometimes phonetic perception can transcend the auditory channel. If you watch an English-subtitled movie in a language you know poorly, after a few minutes you may feel as if you are actually understanding the speech. In the laboratory, researchers can dub a speech sound like *ga* onto a close-up video of a mouth articulating *va*, *ba*, *tha*, or *da*. Viewers literally *hear* a consonant like the one they see the mouth

making—an astonishing illusion with the pleasing name “McGurk effect,” after one of its discoverers.

Actually, one does not need electronic wizardry to create a speech illusion. All speech is an illusion. We hear speech as a string of separate words, but unlike the tree falling in the forest with no one to hear it, a word boundary with no one to hear it has no sound. In the speech sound wave, one word runs into the next seamlessly; there are no little silences between spoken words the way there are white spaces between written words. We simply hallucinate word boundaries when we reach the edge of a stretch of sound that matches some entry in our mental dictionary. This becomes apparent when we listen to speech in a foreign language: it is impossible to tell where one word ends and the next begins. The seamlessness of speech is also apparent in “oronyms,” strings of sound that can be carved into words in two different ways:

The good can decay many ways.

The good candy came anyways.

The stuffy nose can lead to problems.

The stuff he knows can lead to problems.

Some others I've seen.

Some mothers I've seen.

Oronyms are often used in songs and nursery rhymes:

I scream,
You scream,
We all scream
For ice cream.

Mairzey doats and dozey doats
And little lamsey divey,
A kiddley-divey do,
Wouldn't you?

Fuzzy Wuzzy was a bear,
Fuzzy Wuzzy had no hair.

Fuzzy Wuzzy wasn't fuzzy,
Was he?

In fir tar is,
In oak none is.
In mud eel is,
In clay none is.
Goats eat ivy.
Mares eat oats.

And some are discovered inadvertently by teachers reading their students' term papers and homework assignments:

Jose can you see by the donzerly light? [Oh say can you see
by the dawn's early light?]
It's a doggy-dog world. [dog-eat-dog]
Eugene O'Neill won a Pullet Surprise. [Pulitzer Prize]
My mother comes from Pencil Vanea. [Pennsylvania]
He was a notor republic. [notary public]
They played the Bohemian Rap City. [Bohemian Rhapsody]

Even the sequence of sounds we think we hear within a word is an illusion. If you were to cut up a tape of someone's saying *cat*, you would not get pieces that sounded like *k*, *a*, and *t* (the units called "phonemes" that correspond roughly to the letters of the alphabet). And if you spliced the pieces together in the reverse order, they would be unintelligible, not *tack*. As we shall see, information about each component of a word is smeared over the entire word.

Speech perception is another one of the biological miracles making up the language instinct. There are obvious advantages to using the mouth and ear as a channel of communication, and we do not find any hearing community opting for sign language, though it is just as expressive. Speech does not require good lighting, face-to-face contact, or monopolizing the hands and eyes, and it can be shouted over long distances or whispered to conceal the message. But to take advantage of the medium of sound, speech has to overcome the problem that the ear is a narrow informational bottleneck. When engineers

first tried to develop reading machines for the blind in the 1940s, they devised a set of noises that corresponded to the letters of the alphabet. Even with heroic training, people could not recognize the sounds at a rate faster than good Morse code operators, about three units a second. Real speech, somehow, is perceived an order of magnitude faster: ten to fifteen phonemes per second for casual speech, twenty to thirty per second for the man in the late-night Veg-O-Matic ads, and as many as forty to fifty per second for artificially sped-up speech. Given how the human auditory system works, this is almost unbelievable. When a sound like a click is repeated at a rate of twenty times a second or faster, we no longer hear it as a sequence of separate sounds but as a low buzz. If we can hear forty-five phonemes per second, the phonemes cannot possibly be consecutive bits of sound; each moment of sound must have several phonemes packed into it that our brains somehow unpack. As a result, speech is by far the fastest way of getting information into the head through the ear.

No human-made system can match a human in decoding speech. It is not for lack of need or trying. A speech recognizer would be a boon to quadriplegics and other disabled people, to professionals who have to get information into a computer while their eyes or hands are busy, to people who never learned to type, to users of telephone services, and to the growing number of typists who are victims of repetitive-motion syndromes. So it is not surprising that engineers have been working for more than forty years to get computers to recognize the spoken word. The engineers have been frustrated by a tradeoff. If a system has to be able to listen to many different people, it can recognize only a tiny number of words. For example, telephone companies are beginning to install directory assistance systems that can recognize anyone saying the word *yes*, or, in the more advanced systems, the ten English digits (which, fortunately for the engineers, have very different sounds). But if a system has to recognize a large number of words, it has to be trained to the voice of a single speaker. No system today can duplicate a person's ability to recognize both many words and many speakers. Perhaps the state of the art is a system called Dragon-Dictate, which runs on a personal computer and can recognize 30,000

words. But it has severe limitations. It has to be trained extensively on the voice of the user. You . . . have . . . to . . . talk . . . to . . . it . . . like . . . this, with quarter-second pauses between the words (so it operates at about one-fifth the rate of ordinary speech). If you have to use a word that is not in its dictionary, like a name, you have to spell it out using the “Alpha, Bravo, Charlie” alphabet. And the program still garbles words about fifteen percent of the time, more than once per sentence. It is an impressive product but no match for even a mediocre stenographer.

The physical and neural machinery of speech is a solution to two problems in the design of the human communication system. A person might know 60,000 words, but a person’s mouth cannot make 60,000 different noises (at least, not ones that the ear can easily discriminate). So language has exploited the principle of the discrete combinatorial system again. Sentences and phrases are built out of words, words are built out of morphemes, and morphemes, in turn, are built out of phonemes. Unlike words and morphemes, though, phonemes do not contribute bits of meaning to the whole. The meaning of *dog* is not predictable from the meaning of *d*, the meaning of *o*, the meaning of *g*, and their order. Phonemes are a different kind of linguistic object. They connect outward to speech, not inward to mentalese: a phoneme corresponds to an act of making a sound. A division into independent discrete combinatorial systems, one combining meaningless sounds into meaningful morphemes, the others combining meaningful morphemes into meaningful words, phrases, and sentences, is a fundamental design feature of human language, which the linguist Charles Hockett has called “duality of patterning.”

But the phonological module of the language instinct has to do more than spell out the morphemes. The rules of language are discrete combinatorial systems: phonemes snap cleanly into morphemes, morphemes into words, words into phrases. They do not blend or melt or coalesce: *Dog bites man* differs from *Man bites dog*, and believing in God is different from believing in Dog. But to get these structures out of one head and into another, they must be converted to audible signals. The audible signals people can produce are not a series of crisp

beeps like on a touch-tone phone. Speech is a river of breath, bent into hisses and hums by the soft flesh of the mouth and throat. The problems Mother Nature faced are digital-to-analog conversion when the talker encodes strings of discrete symbols into a continuous stream of sound, and analog-to-digital conversion when the listener decodes continuous speech back into discrete symbols.

The sounds of language, then, are put together in several steps. A finite inventory of phonemes is sampled and permuted to define words, and the resulting strings of phonemes are then massaged to make them easier to pronounce and understand before they are actually articulated. I will trace out these steps for you and show you how they shape some of our everyday encounters with speech: poetry and song, slips of the ear, accents, speech recognition machines, and crazy English spelling.

One easy way to understand speech sounds is to track a glob of air through the vocal tract into the world, starting in the lungs.

When we talk, we depart from our usual rhythmic breathing and take in quick breaths of air, then release them steadily, using the muscles of the ribs to counteract the elastic recoil force of the lungs. (If we did not, our speech would sound like the pathetic whine of a released balloon.) Syntax overrides carbon dioxide: we suppress the delicately tuned feedback loop that controls our breathing rate to regulate oxygen intake, and instead we time our exhalations to the length of the phrase or sentence we intend to utter. This can lead to mild hyperventilation or hypoxia, which is why public speaking is so exhausting and why it is difficult to carry on a conversation with a jogging partner.

The air leaves the lungs through the trachea (windpipe), which opens into the larynx (the voice-box, visible on the outside as the Adam’s apple). The larynx is a valve consisting of an opening (the glottis) covered by two flaps of retractable muscular tissue called the vocal folds (they are also called “vocal cords” because of an early anatomist’s error; they are not cords at all). The vocal folds can close off the glottis tightly, sealing the lungs. This is useful when we want to stiffen our upper body, which is a floppy bag of air. Get up from your

chair without using your arms; you will feel your larynx tighten. The larynx is also closed off in physiological functions like coughing and defecation. The grunt of the weightlifter or tennis player is a reminder that we use the same organ to seal the lungs and to produce sound.

The vocal folds can also be partly stretched over the glottis to produce a buzz as the air rushes past. This happens because the high-pressure air pushes the vocal folds open, at which point they spring back and get sucked together, closing the glottis until air pressure builds up and pushes them open again, starting a new cycle. Breath is thus broken into a series of puffs of air, which we perceive as a buzz, called "voicing." You can hear and feel the buzz by making the sounds *ssssssss*, which lacks voicing, and *zzzzzzzz*, which has it.

The frequency of the vocal folds' opening and closing determines the pitch of the voice. By changing the tension and position of the vocal folds, we can control the frequency and hence the pitch. This is most obvious in humming or singing, but we also change pitch continuously over the course of a sentence, a process called intonation. Normal intonation is what makes natural speech sound different from the speech of robots in old science fiction movies and of the Coneheads on *Saturday Night Live*. Intonation is also controlled in sarcasm, emphasis, and an emotional tone of voice such as anger or cheeriness. In "tone languages" like Chinese, rising or falling tones distinguish certain vowels from others.

Though voicing creates a sound wave with a dominant frequency of vibration, it is not like a tuning fork or a test of the Emergency Broadcasting System, a pure tone with that frequency alone. Voicing is a rich, buzzy sound with many "harmonics." A male voice is a wave with vibrations not only at 100 cycles per second but also at 200 cps, 300 cps, 400 cps, 500 cps, 600 cps, 700 cps, and so on, all the way up to 4000 cps and beyond. A female voice has vibrations at 200 cps, 400 cps, 600 cps, and so on. The richness of the sound source is crucial—it is the raw material that the rest of the vocal tract sculpts into vowels and consonants.

If for some reason we cannot produce a hum from the larynx, any rich source of sound will do. When we whisper, we spread the

vocal folds, causing the air stream to break apart chaotically at the edges of the folds and creating a turbulence or noise that sounds like hissing or radio static. A hissing noise is not a neatly repeating wave consisting of a sequence of harmonics, as we find in the periodic sound of a speaking voice, but a jagged, spiky wave consisting of a hodgepodge of constantly changing frequencies. This mixture, though, is all that the rest of the vocal tract needs for intelligible whispering. Some laryngectomy patients are taught "esophageal speech," or controlled burping, which provides the necessary noise. Others place a vibrator against their necks. In the 1970s the guitarist Peter Frampton funneled the amplified sound of his electric guitar through a tube into his mouth, allowing him to articulate his twangings. The effect was good for a couple of hit records before he sank into rock-and-roll oblivion.

The richly vibrating air then runs through a gantlet of chambers before leaving the head: the throat or "pharynx" behind the tongue, the mouth region between the tongue and palate, the opening between the lips, and an alternative route to the external world through the nose. Each chamber has a particular length and shape, which affects the sound passing through by the phenomenon called "resonance." Sounds of different frequencies have different wavelengths (the distance between the crests of the sound wave); higher pitches have shorter wavelengths. A sound wave moving down the length of a tube bounces back when it reaches the opening at the other end. If the length of the tube is a certain fraction of the wavelength of the sound, each reflected wave will reinforce the next incoming one; if it is of a different length, they will interfere with one another. (This is similar to how you get the best effect pushing a child on a swing if you synchronize each push with the top of the arc.) Thus a tube of a particular length amplifies some sound frequencies and filters out others. You can hear the effect when you fill a bottle. The noise of the sloshing water gets filtered by the chamber of air between the surface and the opening: the more water, the smaller the chamber, the higher the resonant frequency of the chamber, and the tinnier the gurgle.

What we hear as different vowels are the different combinations of amplifications and filtering of the sound coming up from the larynx. These combinations are produced by moving five speech organs around in the mouth to change the shapes and lengths of the resonant cavities that the sound passes through. For example, *ee* is defined by two resonances, one from 200 to 350 cps produced mainly by the throat cavity, and the other from 2100 to 3000 cps produced mainly by the mouth cavity. The range of frequencies that a chamber filters is independent of the particular mixture of frequencies that enters it, so we can hear an *ee* as an *ee* whether it is spoken, whispered, sung high, sung low, burped, or twanged.

The tongue is the most important of the speech organs, making language truly the "gift of tongues." Actually, the tongue is three organs in one: the hump or body, the tip, and the root (the muscles that anchor it to the jaw). Pronounce the vowels in *bet* and *butt* repeatedly, *e-uh*, *e-uh*, *e-uh*. You should feel the body of your tongue moving forwards and backwards (if you put a finger between your teeth, you can feel it with the finger). When your tongue is in the front of your mouth, it lengthens the air chamber behind it in your throat and shortens the one in front of it in your mouth, altering one of the resonances: for the *bet* vowel, the mouth amplifies sounds near 600 and 1800 cps; for the *butt* vowel, it amplifies sounds near 600 and 1200. Now pronounce the vowels in *beet* and *bat* alternately. The body of your tongue will jump up and down, at right angles to the *bet-butt* motion; you can even feel your jaw move to help it. This, too, alters the shapes of the throat and mouth chambers, and hence their resonances. The brain interprets the different patterns of amplification and filtering as different vowels.

The link between the postures of the tongue and the vowels it sculpts gives rise to a quaint curiosity of English and many other languages called phonetic symbolism. When the tongue is high and at the front of the mouth, it makes a small resonant cavity there that amplifies some higher frequencies, and the resulting vowels like *ee* and *i* (as in *bit*) remind people of little things. When the tongue is low and to the back, it makes a large resonant cavity that amplifies some lower

frequencies, and the resulting vowels like *a* in *father* and *o* in *core* and in *cot* remind people of large things. Thus mice are *teeny* and *squeak*, but elephants are *humongous* and *roar*. Audio speakers have small *tweeters* for the high sounds and large *woofers* for the low ones. English speakers correctly guess that in Chinese *ch'ing* means light and *ch'ung* means heavy. (In controlled studies with large numbers of foreign words, the hit rate is statistically above chance, though just barely.) When I questioned our local computer wizard about what she meant when she said she was going to *frob* my workstation, she gave me this tutorial on hackerese. When you get a brand-new graphic equalizer for your stereo and aimlessly slide the knobs up and down to hear the effects, that is *frobbing*. When you move the knobs by medium-sized amounts to get the sound to your general liking, that is *twiddling*. When you make the final small adjustments to get it perfect, that is *tweaking*. The *ob*, *id*, and *eak* sounds perfectly follow the large-to-small continuum of phonetic symbolism.

And at the risk of sounding like Andy Rooney on *Sixty Minutes*, have you ever wondered why we say *fiddle-faddle* and not *faddle-fiddle*? Why is it *ping-pong* and *pitter-patter* rather than *pong-ping* and *patter-pitter*? Why *dribs and drabs*, rather than vice versa? Why can't a kitchen be *span and spic*? Whence *riff-raff*, *mish-mash*, *flim-flam*, *chit-chat*, *tit for tat*, *knick-knack*, *zig-zag*, *sing-song*, *ding-dong*, *King Kong*, *criss-cross*, *shilly-shally*, *see-saw*, *hee-haw*, *flip-flop*, *hippity-hop*, *tick-tock*, *tic-tac-toe*, *eeny-meeny-miney-moe*, *bric-a-brac*, *clickety-clack*, *hickory-dickory-dock*, *kit and kaboodle*, and *bibbity-bobbity-boo*? The answer is that the vowels for which the tongue is high and in the front always come before the vowels for which the tongue is low and in the back. No one knows why they are aligned in this order, but it seems to be a kind of syllogism from two other oddities. The first is that words that connote me-here-now tend to have higher and fronter vowels than verbs that connote distance from "me": *me* versus *you*, *here* versus *there*, *this* versus *that*. The second is that words that connote me-here-now tend to come before words that connote literal or metaphorical distance from "me" (or a prototypical generic speaker): *here and there* (not *there and here*), *this and that*, *now and then*, *father*

and son, man and machine, friend or foe, the Harvard-Yale game (among Harvard students), the Yale-Harvard game (among Yalies), Serbo-Croatian (among Serbs), Croat-Serbian (among Croats). The syllogism seems to be: “me” = high front vowel; me first; therefore, high front vowel first. It is as if the mind just cannot bring itself to flip a coin in ordering words; if meaning does not determine the order, sound is brought to bear, and the rationale is based on how the tongue produces the vowels.

Let’s look at the other speech organs. Pay attention to your lips when you alternate between the vowels in *boot* and *book*. For *boot*, you round the lips and protrude them. This adds an air chamber, with its own resonances, to the front of the vocal tract, amplifying and filtering other sets of frequencies and thus defining other vowel contrasts. Because of the acoustic effects of the lips, when we talk to a happy person over the phone, we can literally hear the smile.

Remember your grade-school teacher telling you that the vowel sounds in *bat*, *bet*, *bit*, *bottle*, and *butt* were “short,” and the vowel sounds in *bait*, *beet*, *bite*, *boat*, and *boot* were “long”? And you didn’t know what she was talking about? Well, forget it; her information is five hundred years out of date. Older stages of English differentiated words by whether their vowels were pronounced quickly or were drawn out, a bit like the modern distinction between *bad* meaning “bad” and *baaaad* meaning “good.” But in the fifteenth century English pronunciation underwent a convulsion called the Great Vowel Shift. The vowels that had simply been pronounced longer now became “tense”: by advancing the tongue root (the muscles attaching the tongue to the jaw), the tongue becomes tense and humped rather than lax and flat, and the hump narrows the air chamber in the mouth above it, changing the resonances. Also, some tense vowels in modern English, like in *bite* and *brow*, are “diphthongs,” two vowels pronounced in quick succession as if they were one: ba-eet, bra-oh.

You can hear the effects of the fifth speech organ by drawing out the vowel in *Sam* and *sat*, postponing the final consonant indefinitely. In most dialects of English, the vowels will be different: the vowel in *Sam* will have a twangy, nasal sound. That is because the soft palate

or velum (the fleshy flap at the back of the hard palate) is opened, allowing air to flow out through the nose as well as through the mouth. The nose is another resonant chamber, and when vibrating air flows through it, yet another set of frequencies gets amplified and filtered. English does not differentiate words by whether their vowels are nasal or not, but many languages, like French, Polish, and Portuguese, do. English speakers who open their soft palate even when pronouncing *sat* are said to have a “nasal” voice. When you have a cold and your nose is blocked, opening the soft palate makes no difference, and your voice is the opposite of nasal.

So far we have just discussed the vowels—sounds where the air has clear passage from the larynx to the world. When some barrier is put in the way, one gets a consonant. Pronounce *sssss*. The tip of your tongue—the sixth speech organ—is brought up almost against the gum ridge, leaving a small opening. When you force a stream of air through the opening, the air breaks apart turbulently, creating noise. Depending on the size of the opening and the length of the resonant cavities in front of it, the noise will have some of its frequencies louder than others, and the peak and range of frequencies define the sound we hear as *s*. This noise-making comes from the friction of moving air, so this kind of sound is called a fricative. When rushing air is squeezed between the tongue and palate, we get *sh*; between the tongue and teeth, *th*; and between the lower lip and teeth, *f*. The body of the tongue, or the vocal folds of the larynx, can also be positioned to create turbulence, defining the various “ch” sounds in languages like German, Hebrew, and Arabic (*Bach*, *Chanukah*, and so on).

Now pronounce a *t*. The tip of the tongue gets in the way of the airstream, but this time it does not merely impede the flow; it stops it entirely. When the pressure builds up, you release the tip of the tongue, allowing the air to pop out (flutists use this motion to demarcate musical notes). Other “stop” consonants can be formed by the lips (*p*), by the body of the tongue pressed against the palate (*k*), and by the larynx (in the “glottal” consonants in *uh-oh*). What a listener hears when you produce a stop consonant is the following. First,

nothing, as the air is dammed up behind the stoppage: stop consonants are the sounds of silence. Then, a brief burst of noise as the air is released; its frequency depends on the size of the opening and the resonant cavities in front of it. Finally, a smoothly changing resonance, as voicing fades in while the tongue is gliding into the position of whatever vowel comes next. As we shall see, this hop-skip-and-jump makes life miserable for speech engineers.

Finally, pronounce *m*. Your lips are sealed, just like for *p*. But this time the air does not back up silently; you can say *mmmmm* until you are out of breath. That is because you have also opened your soft palate, allowing all of the air to escape through your nose. The voicing sound is now amplified at the resonant frequencies of the nose and of the part of the mouth behind the blockage. Releasing the lips causes a sliding resonance similar in shape to what we heard for the release in *p*, except without the silence, noise burst, and fade-in. The sound *n* works similarly to *m*, except that the blockage is created by the tip of the tongue, the same organ used for *d* and *s*. So does the *ng* in *sing*, except that the body of the tongue does the job.

Why do we say *razzle-dazzle* instead of *dazzle-razzle*? Why *super-duper*, *helter-skelter*, *harum-scarum*, *hocus-pocus*, *willy-nilly*, *hully-gully*, *roly-poly*, *holy moly*, *herky-jerky*, *walkie-talkie*, *namby-pamby*, *mumbo-jumbo*, *loosey-goosey*, *wing-ding*, *wham-bam*, *hobnob*, *razza-matazz*, and *rub-a-dub-dub*? I thought you'd never ask. Consonants differ in "obstruency"—the degree to which they impede the flow of air, ranging from merely making it resonate, to forcing it noisily past an obstruction, to stopping it up altogether. The word beginning with the less obstruent consonant always comes before the word beginning with the more obstruent consonant. Why ask why?

Now that you have completed a guided tour up the vocal tract, you can understand how the vast majority of sounds in the world's languages are created and heard. The trick is that a speech sound is not a single gesture by a single organ. Every speech sound is a *combination* of gestures, each exerting its own pattern of sculpting of the sound wave, all executed more or less simultaneously—that is one of the

reasons speech can be so rapid. As you may have noticed, a sound can be nasal or not, and produced by the tongue body, the tongue tip, or the lips, in all six possible combinations:

	Nasal (Soft Palate Open)	Not Nasal (Soft Palate Closed)
Lips	<i>m</i>	<i>p</i>
Tongue tip	<i>n</i>	<i>t</i>
Tongue body	<i>ng</i>	<i>k</i>

Similarly, voicing combines in all possible ways with the choice of speech organ:

	Voicing (Larynx Hums)	No Voicing (Larynx Doesn't Hum)
Lips	<i>b</i>	<i>p</i>
Tongue tip	<i>d</i>	<i>t</i>
Tongue body	<i>g</i>	<i>k</i>

Speech sounds thus nicely fill the rows and columns and layers of a multidimensional matrix. First, one of the six speech organs is chosen as the major articulator: the larynx, soft palate, tongue body, tongue tip, tongue root, or lips. Second, a manner of moving that articulator is selected: fricative, stop, or vowel. Third, configurations of the other speech organs can be specified: for the soft palate, nasal or not; for the larynx, voiced or not; for the tongue root, tense or lax; for the lips, rounded or unrounded. Each manner or configuration is a symbol for a set of commands to the speech muscles, and such symbols are called features. To articulate a phoneme, the commands must be executed with precise timing, the most complicated gymnastics we are called upon to perform.

English multiplies out enough of these combinations to define 40 phonemes, a bit above the average for the world's languages. Other languages range from 11 (Polynesian) to 141 (Khoisan or "Bushman"). The total inventory of phonemes across the world numbers in the thousands, but they are all defined as combinations of the

six speech organs and their shapes and motions. Other mouth sounds are not used in any language: scraping teeth, clucking the tongue against the floor of the mouth, making raspberries, and squawking like Donald Duck, for instance. Even the unusual Khoisan and Bantu clicks (similar to the sound of *tsk-tsk* and made famous by the Xhosa pop singer Miriam Makeba) are not miscellaneous phonemes added to those languages. Clicking is a manner-of-articulation feature, like stop or fricative, and it combines with all the other features to define a new layer of rows and columns in the language's table of phonemes. There are clicks produced by the lips, tongue tip, and tongue body, any of which can be nasalized or not, voiced or not, and so on, as many as 48 click sounds in all!

An inventory of phonemes is one of the things that gives a language its characteristic sound pattern. For example, Japanese is famous for not distinguishing *r* from *l*. When I arrived in Japan on November 4, 1992, the linguist Masaaki Yamanashi greeted me with a twinkle and said, "In Japan, we have been very interested in Clinton's erection."

We can often recognize a language's sound pattern even in a speech stream that contains no real words, as with the Swedish chef on *The Muppets* or John Belushi's samurai dry cleaner. The linguist Sarah G. Thomason has found that people who claim to be channeling back to past lives or speaking in tongues are really producing gibberish that conforms to a sound pattern vaguely reminiscent of the claimed language. For example, one hypnotized channeler, who claimed to be a nineteenth-century Bulgarian talking to her mother about soldiers laying waste to the countryside, produced generic pseudo-Slavic gobbledygook like this:

Ovishta reshta rovishta. Vishna beretishti? Ushna barishta
dashto. Na darishnoshto. Korapshnoshashit darishtoy.
Aobashni bedetpa.

And of course, when the words in one language are pronounced with the sound pattern of another, we call it a foreign accent, as in the following excerpt from a fractured fairy tale by Bob Belviso:

GIACCHE ENNE BINNESTAUCCHE

Uans appona taim uase disse boi. Neimmese Giacche.
Naise boi. Live uite ise mamma. Mainde da cao.

Uane dei, di spaghetti ise olle ronne aute. Dei goine feinte
fromme no fudde. Mamma soi orais, "Oreie Giacche, teicche
da cao enne traide erra forre bocchese spaghetti enne somme
uaine."

Bai enne bai commese omme Giacche. I garra no fudde, i
garra no uaine. Meichese misteicche, enne traidesse da cao forre
bonce binnese.

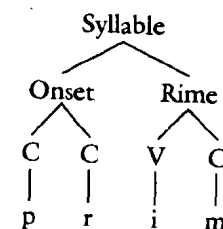
Giacchasse!

What defines the sound pattern of a language? It must be more than just an inventory of phonemes. Consider the following words:

ptak	thale	hlad
plaft	sram	mgla
vlas	flutch	dnom
rtut	roasp	nyip

All of the phonemes are found in English, but any native speaker recognizes that *thale*, *plaft*, and *flutch* are not English words but could be, whereas the remaining ones are not English words and could not be. Speakers must have tacit knowledge about how phonemes are strung together in their language.

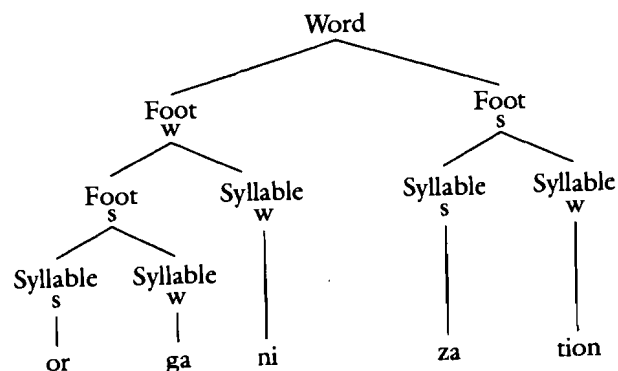
Phonemes are not assembled into words as one-dimensional left-to-right strings. Like words and phrases, they are grouped into units, which are then grouped into bigger units, and so on, defining a tree. The group of consonants (C) at the beginning of a syllable is called an onset; the vowel (V) and any consonants coming after it are called the rime:



The rules generating syllables define legal and illegal kinds of words in a language. In English an onset can consist of a cluster of consonants, like *flit*, *thrive*, and *spring*, as long as they follow certain restrictions. (For example, *vlit* and *sring* are impossible.) A rime can consist of a vowel followed by a consonant or certain clusters of consonants, as in *toast*, *lift*, and *sixths*. In Japanese, in contrast, an onset can have only a single consonant and a rime must be a bare vowel; hence *strawberry ice cream* is translated as *sutoroberi aisukurimo*, *girl-friend* as *garufurendo*. Italian allows some clusters of consonants in an onset but no consonants at the end of a rime. Belviso used this constraint to simulate the sound pattern of Italian in the Giacche story; *and* becomes *enne*, *from* becomes *fromme*, *beans* becomes *binnese*.

Onsets and rimes not only define the possible sounds of a language; they are the pieces of word-sound that are most salient to people, and thus are the units that get manipulated in poetry and word games. Words that rhyme share a rime; words that alliterate share an onset (or just an initial consonant). Pig Latin, eggy-peggy, aygo-paygo, and other secret languages of children tend to splice words at onset-rime boundaries, as does the Yinglish construction in *fancy-shmancy* and *Oedipus-Shmoedipus*. In the 1964 hit song "The Name Game" ("Noam Noam Bo-Boam, Bonana Fana Fo-Foam, Fee Fi Mo Moam, Noam"), Shirley Ellis could have saved several lines in the stanza explaining the rules if she had simply referred to onsets and rimes.

Syllables, in turn, are collected into rhythmic groups called feet:



Syllables and feet are classified as strong (s) and weak (w) by other rules, and the pattern of weak and strong branches determines how much stress each syllable will be given when it is pronounced. Feet, like onsets and rhymes, are salient chunks of word that we tend to manipulate in poetry and wordplay. Meter is defined by the kind of feet that go into a line. A succession of feet with a strong-weak pattern is a trochaic meter, as in *Mary had a little lamb*; a succession with a weak-strong pattern is iambic, as in *The rain in Spain falls mainly in the plain*. An argot popular among young ruffians contains forms like *fan-fuckin-tastic*, *abso-bloody-lutely*, *Phila-fuckin-delphia*, and *Kalama-fuckin-zoo*. Ordinarily, expletives appear in front of an emphatically stressed word; Dorothy Parker once replied to a question about why she had not been at the symphony lately by saying "I've been too fucking busy and vice versa." But in this lingo they are placed inside a single word, always in front of a stressed foot. The rule is followed religiously: *Philadel-fuckin-phia* would get you launched out of the pool hall.

The assemblies of phonemes in the morphemes and words stored in memory undergo a series of adjustments before they are actually articulated as sounds, and these adjustments give further definition to the sound pattern of a language. Say the words *pat* and *pad*. Now add the inflection *-ing* and pronounce them again: *patting*, *padding*. In many dialects of English they are now pronounced identically; the original difference between the *t* and the *d* has been obliterated. What obliterated them is a phonological rule called flapping: if a stop consonant produced with the tip of the tongue appears between two vowels, the consonant is pronounced by flicking the tongue against the gum ridge, rather than keeping it there long enough for air pressure to build up. Rules like flapping apply not only when two morphemes are joined, like *pat* and *-ing*; they also apply to one-piece words. For many English speakers *ladder* and *latter*, though they "feel" like they are made out of different sounds and indeed are represented differently in the mental dictionary, are pronounced the same (except in artificially

exaggerated speech). Thus when cows come up in conversation, often some wag will speak of an udder mystery, an udder success, and so on.

Interestingly, phonological rules apply in an ordered sequence, as if words were manufactured on an assembly line. Pronounce *write* and *ride*. In most dialects of English, the vowels differ in some way. At the very least, the *i* in *ride* is longer than the *i* in *write*. In some dialects, like the Canadian English of newscaster Peter Jennings, hockey star Wayne Gretzky, and yours truly (an accent satirized a few years back, eh, in the television characters Bob and Doug McKenzie), the vowels are completely different: *ride* contains a diphthong gliding from the vowel in *hot* to the vowel *ee*; *write* contains a diphthong gliding from the higher vowel in *hut* to *ee*. But regardless of exactly how the vowel is altered, it is altered in a consistent pattern: there are no words with long/low *i* followed by *t*, nor with short/high *i* followed by *d*. Using the same logic that allowed Lois Lane in her rare lucid moments to deduce that Clark Kent and Superman were the same, namely that they are never in the same place at the same time, we can infer that there is a single *i* in the mental dictionary, which is altered by a rule before being pronounced, depending on whether it appears in the company of *t* or *d*. We can even guess that the initial form stored in memory is like the one in *ride*, and that *write* is the product of the rule, rather than vice versa. The evidence is that when there is no *t* or *d* after the *i*, as in *rye*, and thus no rule disguising the underlying form, it is the vowel in *ride* that we hear.

Now pronounce *writing* and *riding*. The *t* and *d* have been made identical by the flapping rule. But the two *i*'s are still different. How can that be? It is only the difference between *t* and *d* that causes a difference between the two *i*'s, and that difference has been erased by the flapping rule. This shows that the rule that alters *i* must have applied *before* the flapping rule, while *t* and *d* were still distinct. In other words, the two rules apply in a fixed order, vowel-change before flapping. Presumably the ordering comes about because the flapping rule is in some sense there to make articulation easier and thus is farther downstream in the chain of processing from brain to tongue.

Notice another important feature of the vowel-altering rule. The

vowel *i* is altered in front of many different consonants, not just *t*. Compare:

prize	price
five	fife
jibe	hype
geiger	biker

Does this mean there are five different rules that alter *i*—one for *z* versus *s*, one for *v* versus *f*, and so on? Surely not. The change-triggering consonants *t*, *s*, *f*, *p*, and *k* all differ in the same way from their counterparts *d*, *z*, *v*, *b*, and *g*: they are unvoiced, whereas the counterparts are voiced. We need only one rule, then: change *i* whenever it appears before an *unvoiced* consonant. The proof that this is the real rule in people's heads (and not just a way to save ink by replacing five rules with one) is that if an English speaker succeeds in pronouncing the German *ch* in the *Third Reich*, that speaker will pronounce the *ei* as in *write*, not as in *ride*. The consonant *ch* is not in the English inventory, so English speakers could not have learned any rule specifically applying to it. But it is an unvoiced consonant, and if the rule applies to any unvoiced consonant, an English speaker knows exactly what to do.

This selectivity works not only in English but in all languages. Phonological rules are rarely triggered by a single phoneme; they are triggered by an entire class of phonemes that share one or more features (like voicing, stop versus fricative manner, or which organ is doing the articulating). This suggests that rules do not "see" the phonemes in a string but instead look right through them to the features they are made from.

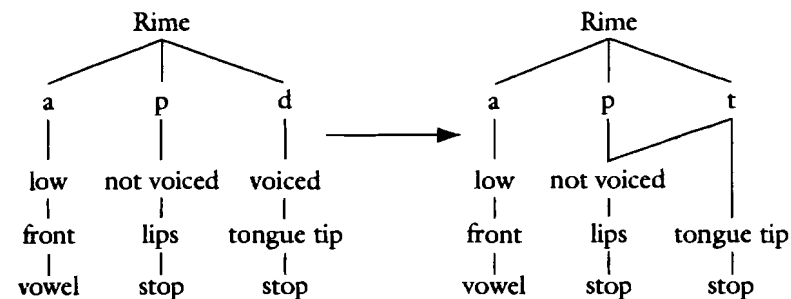
And it is features, not phonemes, that are manipulated by the rules. Pronounce the following past-tense forms:

walked	jogged
slapped	sobbed
passed	fizzed

In *walked*, *slapped*, and *passed*, the *-ed* is pronounced as a *t*; in *jogged*, *sobbed*, and *fizzed*, it is pronounced as a *d*. By now you can probably

figure out what is behind the difference: the *t* pronunciation comes after voiceless consonants like *k*, *p*, and *s*; the *d* comes after voiced consonants like *g*, *b*, and *z*. There must be a rule that adjusts the pronunciation of the suffix *-ed* by peering back into the final phoneme of the stem and checking to see if it has the voicing feature. We can confirm the hunch by asking people to pronounce *Mozart out-Bached Bach*. The *ver* to *out-Bach* contains the sound *ch*, which does not exist in English. Nonetheless everyone pronounces the *-ed* as a *t*, because the *ch* is unvoiced, and the rule puts a *t* next to any unvoiced consonant. We can even determine whether people store the *-ed* suffix as a *t* in memory and use the rule to convert it to a *d* for some words, or the other way around. Words like *play* and *row* have no consonant at the end, and everyone pronounces their past tenses like *plade* and *rode*, not *plate* and *rote*. With no stem consonant triggering a rule, we must be hearing the suffix in its pure, unaltered form in the mental dictionary, that is, *d*. It is a nice demonstration of one of the main discoveries of modern linguistics: a morpheme may be stored in the mental dictionary in a different form from the one that is ultimately pronounced.

Readers with a taste for theoretical elegance may want to bear with me for one more paragraph. Note that there is an uncanny pattern in what the *d*-to-*t* rule is doing. First, *d* itself is voiced, and it ends up next to voiced consonants, whereas *t* is unvoiced, and it ends up next to unvoiced consonants. Second, except for voicing, *t* and *d* are the same; they use the same speech organ, the tongue tip, and that organ moves in the same way, namely sealing up the mouth at the gum ridge and then releasing. So the rule is not just tossing phonemes around arbitrarily, like changing a *p* to an *l* following a high vowel or any other substitution one might pick at random. It is doing delicate surgery on the *-ed* suffix, adjusting it to be the same in voicing as its neighbor, but leaving the rest of its features alone. That is, in converting *slap* + *ed* to *slapt*, the rule is “spreading” the voicing instruction, packaged with the *p* at the end of *slap*, onto the *-ed* suffix like this:



The voicelessness of the *t* in *slapped* matches the voicelessness of the *p* in *slapped* because they are the same voicelessness; they are mentally represented as a single feature linked to two segments. This happens very often in the world's languages. Features like voicing, vowel quality, and tones can spread sideways or sprout connections to several phonemes in a word, as if each feature lived on its own horizontal "tier," rather than being tethered to one and only one phoneme.

So phonological rules "see" features, not phonemes, and they adjust features, not phonemes. Recall, too, that languages tend to arrive at an inventory of phonemes by multiplying out the various combinations of some set of features. These facts show that features, not phonemes, are the atoms of linguistic sound stored and manipulated in the brain. A phoneme is merely a bundle of features. Thus even in dealing with its smallest units, the features, language works by using a combinatorial system.

Every language has phonological rules, but what are they for? You may have noticed that they often make articulation easier. Flapping a *t* or a *d* between two vowels is faster than keeping the tongue in place long enough for air pressure to build up. Spreading voicelessness from the end of a word to its suffix spares the talker from having to turn the larynx off while pronouncing the end of the stem and then turn it back on again for the suffix. At first glance, phonological rules seem to be a mere summary of articulatory laziness. And from here it is a small step to notice phonological adjustments in some dialect other than one's own and conclude that they typify the slovenliness of the

speakers. Neither side of the Atlantic is safe. George Bernard Shaw wrote:

The English have no respect for their language and will not teach their children to speak it. They cannot spell it because they have nothing to spell it with but an old foreign alphabet of which only the consonants—and not all of them—have any agreed speech value. Consequently it is impossible for an Englishman to open his mouth without making some other Englishman despise him.

In his article “Howta Reckanize American Slurvian,” Richard Lederer writes:

Language lovers have long bewailed the sad state of pronunciation and articulation in the United States. Both in sorrow and in anger, speakers afflicted with sensitive ears wince at such mumblings as *guvmint* for *government* and *assessories* for *accessories*. Indeed, everywhere we turn we are assaulted by a slew of slurrings.

But if their ears were even more sensitive, these sorrowful speakers might notice that in fact there is no dialect in which sloppiness prevails. Phonological rules give with one hand and take away with the other. The same bumpkins who are derided for dropping *g*'s in *Nothin' doin'* are likely to enunciate the vowels in *pó-lice* and *accidens* that pointy-headed intellectuals reduce to a neutral “uh” sound. When the Brooklyn Dodgers pitcher Waite Hoyt was hit by a ball, a fan in the bleachers shouted, “Hurt's hoit!” Bostonians who *pahk* their cah in Hahvahd Yahd name their daughters Sheiler and Linder. In 1992 an ordinance was proposed that would have banned the hiring of any immigrant teacher who “speaks with an accent” in—I am not making this up—Westfield, Massachusetts. An incredulous woman wrote to the *Boston Globe* recalling how her native New England teacher defined “homonym” using the example *orphan* and *often*. Another amused reader remembered incurring the teacher's

wrath when he spelled “cuh-rée-uh” *k-o-r-e-a* and “cuh-rée-ur” *c-a-r-e-r*, rather than vice versa. The proposal was quickly withdrawn.

There is a good reason why so-called laziness in pronunciation is in fact tightly regulated by phonological rules, and why, as a consequence, no dialect allows its speakers to cut corners at will. Every act of sloppiness on the part of a speaker demands a compensating measure of mental effort on the part of the conversational partner. A society of lazy talkers would be a society of hard-working listeners. If speakers were to have their way, all rules of phonology would spread and reduce and delete. But if listeners were to have their way, phonology would do the opposite: it would enhance the acoustic differences between confusable phonemes by forcing speakers to exaggerate or embroider them. And indeed, many rules of phonology do that. (For example, there is a rule that forces English speakers to round their lips while saying *sh* but not while saying *s*. The benefit of forcing everyone to make this extra gesture is that the long resonant chamber formed by the pursed lips enhances the lower-frequency noise that distinguishes *sh* from *s*, allowing for easier identification of the *sh* by the listener.) Although every speaker soon becomes a listener, human hypocrisy would make it unwise to depend on the speaker's foresight and consideration. Instead, a single, partly arbitrary set of phonological rules, some reducing, some enhancing, is adopted by every member of a linguistic community when he or she acquires the local dialect as a child.

Phonological rules help listeners even when they do not exaggerate some acoustic difference. By making speech patterns predictable, they add redundancy to a language; English text has been estimated as being between two and four times as long as it has to be for its information content. For example, this book takes up about 900,000 characters on my computer disk, but my file compression program can exploit the redundancy in the letter sequences and squeeze it into about 400,000 characters; computer files that do not contain English text cannot be squished nearly that much. The logician Quine explains why many systems have redundancy built in:

It is the judicious excess over minimum requisite support. It is why a good bridge does not crumble when subjected to stress beyond what reasonably could have been foreseen. It is fallback and failsafe. It is why we address our mail to city and state in so many words, despite the zip code. One indistinct digit in the zip code would spoil everything. . . . A kingdom, legend tells us, was lost for want of a horseshoe nail. Redundancy is our safeguard against such instability.

Thanks to the redundancy of language, yxx cxn xndxrstxnd whxt x xm wrxtxng xvsn xf x rxplxcx xll thx vxwxls wxth xn "x" (t gts ltl hrdr f y dn't vn kn whr th vwls r). In the comprehension of speech, the redundancy conferred by phonological rules can compensate for some of the ambiguity in the sound wave. For example, a listener can know that "thisrip" must be *this rip* and not *the srip* because the English consonant cluster *sr* is illegal.

So why is it that a nation that can put a man on the moon cannot build a computer that can take dictation? According to what I have explained so far, each phoneme should have a telltale acoustic signature: a set of resonances for vowels, a noise band for fricatives, a silence-burst-transition sequence for stops. The sequences of phonemes are massaged in predictable ways by ordered phonological rules, whose effects could presumably be undone by applying them in reverse.

The reason that speech recognition is so hard is that there's many a slip 'twixt brain and lip. No two people's voices are alike, either in the shape of the vocal tract that sculpts the sounds, or in the person's precise habits of articulation. Phonemes also sound very different depending on how much they are stressed and how quickly they are spoken; in rapid speech, many are swallowed outright.

But the main reason an electric stenographer is not just around the corner has to do with a general phenomenon in muscle control called coarticulation. Put a saucer in front of you and a coffee cup a foot or so away from it on one side. Now quickly touch the saucer

and pick up the cup. You probably touched the saucer at the edge nearest the cup, not dead center. Your fingers probably assumed the handle-grasping posture while your hand was making its way to the cup, well before it arrived. This graceful smoothing and overlapping of gestures is ubiquitous in motor control. It reduces the forces necessary to move body parts around and lessens the wear and tear on the joints. The tongue and throat are no different. When we want to articulate a phoneme, our tongue cannot assume the target posture instantaneously; it is a heavy slab of meat that takes time to heft into place. So while we are moving it, our brains are anticipating the next posture in planning the trajectory, just like the cup-and-saucer maneuver. Among the range of positions in the mouth that can define a phoneme, we place the tongue in the one that offers the shortest path to the target for the next phoneme. If the current phoneme does not specify where a speech organ should be, we anticipate where the next phoneme wants it to be and put it there in advance. Most of us are completely unaware of these adjustments until they are called to our attention. Say *Cape Cod*. Until now you probably never noticed that your tongue body is in different positions for the two *k* sounds. In *horseshoe*, the first *s* becomes a *sh*; in *NPR*, the *n* becomes an *m*; in *month* and *width*, the *n* and *d* are articulated at the teeth, not the usual gum ridge.

Because sound waves are minutely sensitive to the shapes of the cavities they pass through, this coarticulation wreaks havoc with the speech sound. Each phoneme's sound signature is colored by the phonemes that come before and after, sometimes to the point of having nothing in common with its sound signature in the company of a different set of phonemes. That is why you cannot cut up a tape of the sound *cat* and hope to find a beginning piece that contains the *k* alone. As you make earlier and earlier cuts, the piece may go from sounding like *ka* to sounding like a chirp or whistle. This shingling of phonemes in the speech stream could, in principle, be a boon to an optimally designed speech recognizer. Consonant and vowels are being signaled simultaneously, greatly increasing the rate of phonemes per second, as I noted at the beginning of this chapter, and there are

many redundant sound cues to a given phoneme. But this advantage can be enjoyed only by a high-tech speech recognizer, one that has some kind of knowledge of how vocal tracts blend sounds.

The human brain, of course, is a high-tech speech recognizer, but no one knows how it succeeds. For this reason psychologists who study speech perception and engineers who build speech recognition machines keep a close eye on each other's work. Speech recognition may be so hard that there are only a few ways it could be solved in principle. If so, the way the brain does it may offer hints as to the best way to build a machine to do it, and how a successful machine does it may suggest hypotheses about how the brain does it.

Early in the history of speech research, it became clear that human listeners might somehow take advantage of their expectations of the kinds of things a speaker is likely to say. This could narrow down the alternatives left open by the acoustic analysis of the speech signal. We have already noted that the rules of phonology provide one sort of redundancy that can be exploited, but people might go even farther. The psychologist George Miller played tapes of sentences in background noise and asked people to repeat back exactly what they heard. Some of the sentences followed the rules of English syntax and made sense.

Furry wildcats fight furious battles.
Respectable jewelers give accurate appraisals.
Lighted cigarettes create smoky fumes.
Gallant gentlemen save distressed damsels.
Soapy detergents dissolve greasy stains.

Others were created by scrambling the words within phrases to create colorless-green-ideas sentences, grammatical but nonsensical:

Furry jewelers create distressed stains.
Respectable cigarettes save greasy battles.
Lighted gentlemen dissolve furious appraisals.
Gallant detergents fight accurate fumes.
Soapy wildcats give smoky damsels.

A third kind was created by scrambling the phrase structure but keeping related words together, as in

Furry fight furious wildcat battles.
Jewelers respectable appraisals accurate give.

Finally, some sentences were utter word salad, like

Furry create distressed jewelers stains.
Cigarettes respectable battles greasy save.

People did best with the grammatical sensible sentences, worse with the grammatical nonsense and the ungrammatical sense, and worst of all with the ungrammatical nonsense. A few years later the psychologist Richard Warren taped sentences like *The state governors met with their respective legislatures convening in the capital city*, excised the first *s* from *legislatures*, and spliced in a cough. Listeners could not tell that any sound was missing.

If one thinks of the sound wave as sitting at the bottom of a hierarchy from sounds to phonemes to words to phrases to the meanings of sentences to general knowledge, these demonstrations seem to imply that human speech perception works from the top down rather than just from the bottom up. Maybe we are constantly guessing what a speaker will say next, using every scrap of conscious and unconscious knowledge at our disposal, from how coarticulation distorts sounds, to the rules of English phonology, to the rules of English syntax, to stereotypes about who tends to do what to whom in the world, to hunches about what our conversational partner has in mind at that very moment. If the expectations are accurate enough, the acoustic analysis can be fairly crude; what the sound wave lacks, the context can fill in. For example, if you are listening to a discussion about the destruction of ecological habitats, you might be on the lookout for words pertaining to threatened animals and plants, and then when you hear speech sounds whose phonemes you cannot pick out like "eesees," you would perceive it correctly as *species*—unless you are Emily Litella, the hearing-impaired editorialist on *Saturday Night Live* who argued passionately against the campaign to protect endan-

gered feces. (Indeed, the humor in the Gilda Radner character, who also fulminated against saving Soviet jewelry, stopping violins in the streets, and preserving natural racehorses, comes not from her impairment at the bottom of the speech-processing system but from her ditziness at the top, the level that should have prevented her from arriving at her interpretations.)

The top-down theory of speech perception exerts a powerful emotional tug on some people. It confirms the relativist philosophy that we hear what we expect to hear, that our knowledge determines our perception, and ultimately that we are not in direct contact with any objective reality. In a sense, perception that is strongly driven from the top down would be a barely controlled hallucination, and that is the problem. A perceiver forced to rely on its expectations is at a severe disadvantage in a world that is unpredictable even under the best of circumstances. There is a reason to believe that human speech perception is, in fact, driven quite strongly by acoustics. If you have an indulgent friend, you can try the following experiment. Pick ten words at random out of a dictionary, phone up the friend, and say the words clearly. Chances are the friend will reproduce them perfectly, relying only on the information in the sound wave and knowledge of English vocabulary and phonology. The friend could not have been using any higher-level expectations about phrase structure, context, or story line because a list of words blurted out of the blue has none. Though we may call upon high-level conceptual knowledge in noisy or degraded circumstances (and even here it is not clear whether the knowledge alters perception or just allows us to guess intelligently after the fact), our brains seem designed to squeeze every last drop of phonetic information out of the sound wave itself. Our sixth sense may perceive speech as language, not as sound, but it *is* a sense, something that connects us to the world, and not just a form of suggestibility.

Another demonstration that speech perception is not the same thing as fleshing out expectations comes from an illusion that the columnist Jon Carroll has called the mondegreen, after his mis-hearing of the folk ballad "The Bonnie Earl O'Moray":

Oh, ye hielands and ye lowlands,
Oh, where hae ye been?
They have slain the Earl of Moray,
And laid him on the green.

He had always thought that the lines were "They have slain the Earl of Moray, And Lady Mondegreen." Mondegreens are fairly common (they are an extreme version of the Pullet Surprises and Pencil Vaneas mentioned earlier); here are some examples:

A girl with colitis goes by. [A girl with kaleidoscope eyes. From the Beatles song "Lucy in the Sky with Diamonds."]
Our father wishart in heaven; Harold be thy name . . . Lead us not into Penn Station.
Our father which art in Heaven; hallowed by thy name . . . Lead us not into temptation. From the Lord's Prayer.]
He is trampling out the vintage where the grapes are wrapped and stored. [. . . grapes of wrath are stored. From "The Battle Hymn of the Republic."]
Gladly the cross-eyed bear. [Gladly the cross I'd bear.]
I'll never be your pizza burnin'. [. . . your beast of burden. From the Rolling Stones song.]
It's a happy enchilada, and you think you're gonna drown. [It's a half an inch of water . . . From the John Prine song "That's the Way the World Goes 'Round."]

The interesting thing about mondegreens is that the mis-hearings are generally *less* plausible than the intended lyrics. In no way do they bear out any sane listener's general expectations of what a speaker is likely to say or mean. (In one case a student stubbornly misheard the Shocking Blue hit song "I'm Your Venus" as "I'm Your Penis" and wondered how it was allowed on the radio.) The mondegreens do conform to English phonology, English syntax (sometimes), and English vocabulary (though not always, as in the word *mondegreen* itself). Apparently, listeners lock in to some set of words that fit the sound and that hang together more or less as English

words and phrases, but plausibility and general expectations are ~~not~~ running the show.

The history of artificial speech recognizers offers a similar moral. In the 1970s a team of artificial intelligence researchers at Carnegie-Mellon University headed by Raj Reddy designed a computer program called HEARSAY that interpreted spoken commands to move chess pieces. Influenced by the top-down theory of speech perception, they designed the program as a “community” of “expert” subprograms cooperating to give the most likely interpretation of the signal. There were subprograms that specialized in acoustic analysis, in phonology, in the dictionary, in syntax, in rules for the legal moves of chess, even in chess strategy as applied to the game in progress. According to one story, a general from the defense agency that was funding the research came up for a demonstration. As the scientists sweated he was seated in front of a chessboard and a microphone hooked up to the computer. The general cleared his throat. The program printed “Pawn to King 4.”

The recent program DragonDictate, mentioned earlier in the chapter, places the burden more on good acoustic, phonological, and lexical analyses, and that seems to be responsible for its greater success. The program has a dictionary of words and their sequences of phonemes. To help anticipate the effects of phonological rules and coarticulation, the program is told what every English phoneme sounds like in the context of every possible preceding phoneme and every possible following phoneme. For each word, these phonemes-in-context are arranged into a little chain, with a probability attached to each transition from one sound unit to the next. This chain serves as a crude model of the speaker, and when a real speaker uses the system, the probabilities in the chain are adjusted to capture that person’s manner of speaking. The entire word, too, has a probability attached to it, which depends on its frequency in the language and on the speaker’s habits. In some versions of the program, the probability value for a word is adjusted depending on which word precedes it; this is the only top-down information that the program uses. All this knowledge allows the program to calculate which word is most likely

~~to~~ have come out of the mouth of the speaker given the input sound. Even then, DragonDictate relies more on expectancies than an able-bodied human does. In the demonstration I saw, the program had to be coaxed into recognizing *word* and *worm*, even when they were pronounced as clear as a bell, because it kept playing the odds and guessing higher-frequency *were* instead.

Now that you know how individual speech units are produced, how they are represented in the mental dictionary, and how they are rearranged and smeared before they emerge from the mouth, you have reached the prize at the bottom of this chapter: why English spelling is not as deranged as it first appears.

The complaint about English spelling, of course, is that it pretends to capture the sounds of words but does not. There is a long tradition of doggerel making this point, of which this stanza is a typical example:

Beware of heard, a dreadful word
That looks like beard and sounds like bird,
And dead: it’s said like bed, not bead—
For goodness’ sake don’t call it “deed”!
Watch out for meat and great and threat
(They rhyme with suite and straight and debt).

George Bernard Shaw led a vigorous campaign to reform the English alphabet, a system so illogical, he said, that it could spell *fish* as “ghoti”—*gh* as in *tough*, *o* as in *women*, *ti* as in *nation*. (“Mnompoute” for *minute* and “mnopspteiche” for *mistake* are other examples.) In his will Shaw bequeathed a cash prize to be awarded to the designer of a replacement alphabet for English, in which each sound in the spoken language would be recognizable by a single symbol: He wrote:

To realize the annual difference in favour of a forty-two letter phonetic alphabet . . . you must multiply the number of minutes in the year, the number of people in the world who are

continuously writing English words, casting types, manufacturing printing and writing machines, by which time the total figure will have become so astronomical that you will realize that the cost of spelling even one sound with two letters has cost us centuries of unnecessary labour. A new British 42 letter alphabet would pay for itself a million times over not only in hours but in moments. When this is grasped, all the useless twaddle about enough and cough and laugh and simplified spelling will be dropped, and the economists and statisticians will be set to work to gather in the orthographic Golconda.

My defense of English spelling will be halfhearted. For although language is an instinct, written language is not. Writing was invented a small number of times in history, and alphabetic writing, where one character corresponds to one sound, seems to have been invented only once. Most societies have lacked written language, and those that have it inherited it or borrowed it from one of the inventors. Children must be taught to read and write in laborious lessons, and knowledge of spelling involves no daring leaps from the training examples like the leaps we saw in Simon, Mayela, and the Jabba and *mice-eater* experiments in Chapters 3 and 5. And people do not uniformly succeed. Illiteracy, the result of insufficient teaching, is the rule in much of the world, and dyslexia, a presumed congenital difficulty in learning to read even with sufficient teaching, is a severe problem even in industrial societies, found in five to ten percent of the population.

But though writing is an artificial contraption connecting vision and language, it must tap into the language system at well-demarcated points, and that gives it a modicum of logic. In all known writing systems, the symbols designate only three kinds of linguistic structure: the morpheme, the syllable, and the phoneme. Mesopotamian cuneiform, Egyptian hieroglyphs, Chinese logograms, and Japanese kanji encode morphemes. Cherokee, Ancient Cypriot, and Japanese kana are syllable-based. All modern phonemic alphabets appear to be descended from a system invented by the Canaanites around 1700 B.C. No writing system has symbols for actual sound units that can be

identified on an oscilloscope or spectrogram, such as a phoneme as it is pronounced in a particular context or a syllable chopped in half.

Why has no writing system ever met Shaw's ideal of one symbol per sound? As Shaw himself said elsewhere, "There are two tragedies in life. One is not to get your heart's desire. The other is to get it." Just think back to the workings of phonology and coarticulation. A true Shavian alphabet would mandate different vowels in *write* and *ride*, different consonants in *write* and *writing*, and different spellings for the past-tense suffix in *slapped*, *sobbed*, and *sorted*. *Cape Cod* would lose its visual alliteration. A *horse* would be spelled differently from its *horseshoe*, and National Public Radio would have the enigmatic abbreviation *MPR*. We would need brand-new letters for the *n* in *month* and the *d* in *width*. I would spell *often* differently from *orphan*, but my neighbors here in the Hub would not, and their spelling of *career* would be my spelling of *Korea* and vice versa.

Obviously, alphabets do not and should not correspond to sounds; at best they correspond to the phonemes specified in the mental dictionary. The actual sounds are different in different contexts, so true phonetic spelling would only obscure their underlying identity. The surface sounds are predictable by phonological rules, though, so there is no need to clutter up the page with symbols for the actual sounds; the reader needs only the abstract blueprint for a word and can flesh out the sound if needed. Indeed, for about eighty-four percent of English words, spelling is completely predictable from regular rules. Moreover, since dialects separated by time and space often differ most in the phonological rules that convert mental dictionary entries into pronunciations, a spelling corresponding to the underlying entries, not the sounds, can be widely shared. The words with truly weird spellings (like *of*, *people*, *women*, *have*, *said*, *do*, *done*, and *give*) generally are the commonest ones in the language, so there is ample opportunity for everyone to memorize them.

Even the less predictable aspects of spelling bespeak hidden linguistic regularities. Consider the following pairs of words where the same letters get different pronunciations:

electric–electricity	declare–declaration
photograph–photography	muscle–muscular
grade–gradual	condemn–condemnation
history–historical	courage–courageous
revise–revision	romantic–romanticize
adore–adoration	industry–industrial
bomb–bombard	fact–factual
nation–national	inspire–inspiration
critical–criticize	sign–signature
mode–modular	malign–malignant
resident–residential	

Once again the similar spellings, despite differences in pronunciation, are there for a reason: they are identifying two words as being based on the same root morpheme. This shows that English spelling is not completely phonemic; sometimes letters encode phonemes, but sometimes a sequence of letters is specific to a morpheme. And a morphemic writing system is more useful than you might think. The goal of reading, after all, is to understand the text, not to pronounce it. A morphemic spelling can help a reader distinguishing homophones, like *meet* and *mete*. It can also tip off a reader that one word contains another (and not just a phonologically identical impostor). For example, spelling tells us that *overcome* contains *come*, so we know that its past tense must be *overcame*, whereas *succumb* just contains the sound “kum,” not the morpheme *come*, so its past tense is not *succame* but *succumbed*. Similarly, when something *recedes*, one has a *recession*, but when someone *re-seeds* a lawn, we have a *re-seeding*.

In some ways, a morphemic writing system has served the Chinese well, despite the inherent disadvantage that readers are at a loss when they face a new or rare word. Mutually unintelligible dialects can share texts (even if their speakers pronounce the words very differently), and many documents that are thousands of years old are readable by modern speakers. Mark Twain alluded to such inertia in our own Roman writing system when he wrote, “They spell it Vinci and pronounce it Vinchy; foreigners always spell better than they pronounce.”

Of course English spelling could be better than it is. But it is already much better than people think it is. That is because writing systems do not aim to represent the actual sounds of talking, which we do not hear, but the abstract units of language underlying them, which we do hear.