

## INTRODUCTION

*(Not) Like a Rock*

Here's how January 21, 2000 panned out for three different elements of the natural order.

### *Element 1: A Rock*

Here is a day in the life of a small, gray-white rock nestling amidst the ivy in my St. Louis backyard. It stayed put. Some things happened to it: there was rain, and it became wet and shiny; there was wind, and it was subtly eroded; my cat chased a squirrel nearby, and this made the rock sway. That's about it, really. There is no reason to believe the rock had any thoughts, or that any of this felt like anything to the rock. Stuff happened, but that was all.

### *Element 2: A Cat*

Lolo, my cat, had a rather different kind of day. About 80% of it was spent, as usual, asleep. But there were forays into the waking, wider world. Around 7 A.M. some inner stirring led Lolo to exit the house, making straight for the catflap from the warm perch of the living room sofa. Outside, bodily functions doubtless dominated, at least at first. Later, following a brief trip back inside (unerringly routed via the catflap and the food tray), squirrels were chased and dangers avoided. Other cats were dealt with in ways appropriate to their rank, station, girth, and meanness. There was a great deal of further sleeping.

### *Element 3: Myself*

My day was (I think) rather more like Lolo's than like the rock's. We both (Lolo and I) pursued food and warmth. But my day included, I suspect, rather more outright

*contemplation*. The kind of spiraling meta-contemplation, in fact, that has sometimes gotten philosophy a bad name. Martin Amis captured the spirit well:

I experienced thrilling self-pity. "What will that mind of your get up to next?" I said, recognizing the self-congratulation behind this thought and the self-congratulation behind that recognition, and the self-congratulation behind recognizing that recognition.

Steady on. (Martin Amis, *The Rachel Papers*, p. 96)

I certainly did some of that. I had thoughts, even "trains of thought" (reasonable sequences of thinkings such as "It's 1 P.M. Time to eat. What's in the fridge?" and so on). But there were also thoughts about thoughts, as I sat back and observed my own trains of thought, alert for colorful examples to import into this text.

What, then, distinguishes cat from rock, and (perhaps) person from cat? What are the mechanisms that make thought and feeling possible? And what further tricks or artifices give my own kind of mindfulness its peculiar self-aware tinge? Such questions seem to focus attention on three different types of phenomena:

1. The feelings that characterize daily experience (hunger, sadness, desire, and so on)
2. The flow of thoughts and reasons
3. The meta-flow of thoughts about thoughts (and thoughts about feelings), of reflection on reasons, and so on.

Most of the research programs covered in this text have concentrated on the middle option. They have tried to explain how my thought that it is 1 P.M. could lead to my thought about lunch, and how it could cause my subsequent lunch-seeking actions. All three types of phenomena are, however, the subject of what philosophers call "mentalistic discourse." A typical example of mentalistic discourse is the appeal to beliefs (and desires) to explain actions. The more technical phrase "propositional attitude psychology" highlights the standard shape of such explanations: such explanations pair mental attitudes (believing, hoping, fearing, etc.) with specific propositions ("that it is raining," "that the coffee is in the kitchen," "that the squirrel is up the tree," etc.) so as to explain intelligent action. Thus in a sentence such as "Pepa hopes that the wine is chilled," the that-construction introduces a proposition ("the wine is chilled") toward which the agent is supposed to exhibit some attitude (in this case, hoping). Other attitudes (such as believing, desiring, fearing, and so on) may, of course, be taken to the same proposition. Our everyday understandings of each other's behavior involve hefty doses of propositional attitude ascription: for example, I may explain Pepa's reluctance to open the wine by saying "Pepa believes that the wine is not yet chilled and desires that it remain in the fridge for a few more minutes."

Such ways of speaking (and thinking) pay huge dividends. They support a surprising degree of predictive success, and are the common currency of many of our social and practical projects. In this vein, the philosopher Jerry Fodor suggests that commonsense psychology is *ubiquitous*, almost *invisible* (because it works so well), and practically *indispensable*. For example, it enables us to make precise plans on the basis of someone's 2-month-old statement that they will arrive on flight 594 on Friday, November 20, 1999. Such plans often work out—a truly amazing fact given the number of physical variables involved. They work out (when they do) because the statement reflects an intention (to arrive that day, on that flight) that is somehow an active shaper of my behavior. I desire that I should arrive on time. You know that I so desire. And on that basis, with a little cooperation from the world at large, miracles of coordination can occur. Or as Fodor more colorfully puts it:

If you want to know where my physical body will be next Thursday, mechanics—our best science of middle-sized objects after all, and reputed to be pretty good in its field—is *no use to you at all*. Far the best way to find out (usually in practice, the only way to find out) is: *ask me!* (Fodor, 1987, p. 6, original emphasis)

Commonsense psychology thus works, and with a vengeance. But why? Why is it that treating each other as having beliefs, hopes, intentions, and the like allows us successfully to explain, predict, and understand so much daily behavior? Beliefs, desires, and so on are, after all, invisible. We see (what we take to be) their effects. But no one has ever actually seen a belief. Such things are (currently? permanently?) unobservable. Commonsense psychology posits these unobservables, and looks to be committed to a body of law-like relations involving them. For example, we explain Fred's jumping up and down by saying that he is happy because his sister just won the Nobel Prize. Behind this explanation lurks an implicit belief in a law-like regularity, viz. "if someone desires *x*, and *x* occurs, then (all other things being equal) they feel happy." All this makes commonsense psychology look like a theory about the invisible, but *causally potent*, roots of intelligent behavior. What, then, can be making the theory true (assuming that it is)? What is a belief (or a hope, or a fear) such that it can cause a human being (or perhaps a cat, dog, etc.) to act in an appropriate way?

Once upon a time, perhaps, it would have been reasonable to respond to the challenge by citing a special kind of spirit-substance: the immaterial but *causally* empowered seat of the mental [for some critical discussion, see Churchland (1984), pp. 7–22, and Appendix I of the present text]. Our concerns, however, lie squarely with attempts that posit nothing extra—nothing beyond the properties and organization of the material brain, body, and world. The goal is a fully materialistic story in which mindware emerges as *nothing but* the playing out of ordinary physical states and processes in the familiar physical world. Insofar as the mental is in any way *special*, according to these views, it is special because it depends on some

particular and unusual ways in which ordinary physical stuff can be built, arranged, and organized.

Views of this latter kind are broadly speaking *monistic*: that is to say, they posit only one basic *kind* of stuff (the material stuff) and attempt to explain the distinctive properties of mental phenomena in terms that are continuous with, or at least appropriately grounded in, our best understanding of the workings of the nonmental universe. A common, but still informative, comparison is with the once-lively (sic) debate between vitalists and nonvitalists. The vitalist held that living things were quite fundamentally different from the rest of inanimate nature, courtesy of a special extra force or ingredient (the “vital spark”), that was missing elsewhere. This is itself a kind of dualism. The demonstration of the fundamental unity of organic and inorganic chemistry (and the absence, in that fundament, of anything resembling a vital spark) was thus a victory—as far as we can tell—for a kind of monism. The animate world, it seems, is the result of *nothing but* the fancy combination of the same kinds of ingredients and forces responsible for inanimate nature. As it was with the animate, so materialists (which is to say, nearly all those working in contemporary cognitive science, the present author included) believe it must be with the mental. The mental world, it is anticipated, must prove to depend on nothing but the fancy combination and organization of ordinary physical states and processes.

Notice, then, the problem. The mental certainly *seems* special, unusual, and different. Indeed, as we saw, it *is* special, unusual, and different: thoughts give way to other thoughts and actions in a way that *respects reasons*: the thought that the forecast was sun (to adapt the famous but less upbeat example) causes me to apply sunscreen, to don a Panama hat, and to think “just another day in paradise.” And there is a qualitative feel, a “something it is like” to have a certain kind of mental life: I *experience* the stabbings of pain, the stirrings of desire, the variety of tastes, colors, and sounds. It is the burden of materialism to somehow get to grips with these various special features in a way that is continuous with, or appropriately grounded in, the way we get to grips with the rest of the physical world—by some understanding of material structure, organization, and causal flow. This is a tall order, indeed. But, as Jerry Fodor is especially fond of pointing out, there is at least one good idea floating around—albeit one that targets just one of the two special properties just mentioned: reason-respecting flow.

The idea, in a supercompressed nutshell, is that the power of a thought (e.g., that the forecast is sun) to cause further thoughts and actions (to apply sunscreen, to think “another day in paradise”) is fully explained by what are broadly speaking *structural* properties of the system in which the thought occurs. By a structural property I here mean simply a physical or organizational property: something whose nature is explicable *without* invoking the specific thought-content involved. An example will help. Consider the way a pocket calculator outputs the sum of two numbers given a sequence of button pushings that we interpret as inputting “2”

“+” “2.” The calculator need not (and does not) understand anything about numbers for this trick to work. It is simply structured so that those button pushings will typically lead to the output “4” as surely as a river will typically find the path of least resistance down a mountain. It is just that in the former case, but not the latter, there has been a process of design such that the physical stuff became organized *so as* its physical unfoldings would reflect the arithmetical constraints governing sensible (arithmetic-respecting) transitions in number space. Natural selection and lifetime learning, to complete the (supercompressed) picture, are then imagined to have sculpted our *brains* so that certain structure-based physical unfoldings respect the constraints on sensible sequences of thoughts and sensible thought-action transitions. Recognition of the predator thus causes running, hiding, and thoughts of escape, whereas recognition of the food causes eating, vigilance, and thoughts of where to find more. Our whole reason-respecting mental life, so the story goes, is just the unfolding of what is, at bottom, a physical and structural story. Mindfulness is just matter, nicely orchestrated.

(As to that *other* distinctive property, “qualitative feel,” let’s just say—and see Appendix II—that it’s a problem. Maybe that too is just a property of matter, nicely orchestrated. But how the orchestration *yields* the property is in this case much less clear, even in outline. So we’ll be looking where the light is.)

In the next eight chapters, I shall expand and pursue that simple idea of mindware (selected aspects!) as matter, nicely orchestrated. The chase begins with a notion of mind as a kind of souped-up pocket calculator (mind as a familiar kind of computer, but built out of meat rather than silicon). It proceeds to the vision of mind as dependent on the operation of a radically different *kind* of computational device (the kind known as artificial neural networks). And it culminates in the contemporary (and contentious) research programs that highlight the complex interactions among brains, bodies, and environmental surroundings (work on robotics, artificial life, dynamics, and situated cognition).

The narrative is, let it be said, biased. It reflects my own view of what we have learned in the past 30 or 40 years of cognitive scientific research. What we have learned, I suggest, is that there are many deeply different ways to put flesh onto that broad, materialistic framework, and that some once-promising incarnations face deep and unexpected difficulties. In particular, the simple notion of the brain as a kind of symbol-crunching computer is probably too simple, and too far removed from the neural and ecological realities of complex, time-critical interaction that sculpted animal minds. The story I tell is thus a story of (a kind of) *inner symbol flight*. But it is a story of progress, refinement, and renewal, not one of abandonment and decay. The sciences of the mind are, in fact, in a state of rude health, of exuberant flux. Time, then, to start the story, to seek the origins of mind in the whirr and buzz of well-orchestrated matter.

# MEAT MACHINES

## *Mindware as Software*



### 1.1 Sketches

### 1.2 Discussion

- A. Why Treat Thought as Computation?
- B. Is Software an Autonomous Level in Nature?
- C. *Mimicking, Modeling, and Behavior*
- D. Consciousness, Information, and Pizza

### 1.3 A Diversion

### 1.4 Suggested Readings

### 1.1 Sketches

The computer scientist Marvin Minsky once described the human brain as a meat machine—no more no less. It is, to be sure, an ugly phrase. But it is also a striking image, a compact expression of both the genuine scientific excitement and the rather gung-ho materialism that tended to characterize the early years of cognitive scientific research. Mindware—our thoughts, feelings, hopes, fears, beliefs, and intellect—is cast as nothing but the operation of the biological brain, the meat machine in our head. This notion of the brain as a meat *machine* is interesting, for it immediately in-

vites us to focus not so much on the material (the meat) as on the machine: the way the material is organized and the kinds of operation it supports. The same machine (see Box 1.1) can, after all, often be made of iron, or steel, or tungsten, or whatever. What we confront is thus both a rejection of the idea of mind as immaterial spirit-stuff and an affirmation that mind is best studied from a kind of *engineering perspective that reveals the nature of the machine that all that wet, white, gray, and sticky stuff happens to build.*

What exactly is meant by casting the brain as a machine, albeit one made out of meat? There exists a historical trend, to be sure, of trying to understand the workings of the brain by analogy with various currently fashionable technologies: the telegraph, the steam engine, and the telephone switchboard are all said to have had their day in the sun. But the "meat machine" phrase is intended, it should now be clear, to do more than hint at some rough analogy. For with regard to the very special class of machines known as computers, the claim is that the brain (and, by

## Box 1.1

## THE "SAME MACHINE"

In what sense can "the same machine" be made out of iron, or steel, or whatever? Not, obviously, in the strict sense of numerical identity. A set of steel darts and a set of tungsten ones cannot be the *very same* (numerically identical) set of darts. The relevant sense of sameness is, rather, some sense of *functional* sameness. You can make a perfectly good set of darts out of either material (though not, I suppose, out of jello), just as you can make a perfectly good corkscrew using a myriad (in this latter case quite radically) different designs and materials. In fact, what makes something a corkscrew is simply that it is designed as, and is capable of acting as, a cork-removing device. The notion of a brain as a meat machine is meant to embody a similar idea: that what matters about the brain is not the stuff it is made of but the way that stuff is organized so as to support thoughts and actions. The idea is that this capability depends on quite abstract properties of the physical device that could very well be duplicated in a device made, say, out of wires and silicon. Sensible versions of this idea need not claim then that *any* material will do: perhaps, for example, a certain stability over time (a tendency not to rapidly disorganize) is needed. The point is just that given that certain preconditions are met the same functionality can be pressed from multiple different materials and designs. For some famous opposition to this view, see Searle (1980, 1992).

not unproblematic extension, the mind) actually *is* some such device. It is not that the brain is somehow *like* a computer: everything is like everything else in some respect or other. It is that neural tissues, synapses, cell assemblies, and all the rest are just nature's rather wet and sticky way of building a hunk of honest-to-God computing machinery. Mindware, it is then claimed, is found "in" the brain in just the way that software is found "in" the computing system that is running it.

The attractions of such a view can hardly be overstated. It makes the mental special without making it ghostly. It makes the mental depend on the physical, but in a rather complex and (as we shall see) liberating way. And it provides a ready-made answer to a profound puzzle: how to get sensible, reason-respecting behavior out of a hunk of physical matter. To flesh out this idea of nonmysterious reason-respecting behavior, we next review some crucial developments<sup>1</sup> in the history (and prehistory) of artificial intelligence.

<sup>1</sup>The next few paragraphs draw on Newell and Simon's (1976) discussion of the development of the Physical Symbol Hypothesis (see Chapter 2 following), on John Haugeland's (1981a), and on Glymour, Ford, and Hayes' (1995).

One key development was the appreciation of the power and scope of formal logics. A decent historical account of this development would take us too far afield, touching perhaps on the pioneering efforts in the seventeenth century by Pascal and Leibniz, as well as on the twentieth-century contributions of Boole, Frege, Russell, Whitehead, and others. A useful historical account can be found in Glymour, Ford, and Hayes (1995). The idea that shines through the history, however, is the idea of finding and describing "laws of reason"—an idea whose clearest expression emerged first in the arena of formal logics. Formal logics are systems comprising sets of symbols, ways of joining the symbols so as to express complex propositions, and rules specifying how to legally derive new symbol complexes from old ones. The beauty of formal logics is that the steadfast application of the rules guarantees that you will never legally infer a false conclusion from true premises, even if you have no idea what, if anything, the strings of symbols actually mean. Just follow the rules and truth will be preserved. The situation is thus a little (just a little) like a person, incompetent in practical matters, who is nonetheless able to successfully build a cabinet or bookshelf by following written instructions for the manipulation of a set of preprovided pieces. Such building behavior can look as if it is rooted in a deep appreciation of the principles and laws of woodworking; but in fact, the person is just blindly making the moves allowed or dictated by the instruction set.

Formal logics show us how to preserve at least one kind of semantic (meaning-involving; see Box 1.2) property without relying on anyone's actually appreciating the meanings (if any) of the symbol strings involved. The seemingly ghostly and ephemeral world of meanings and logical implications is respected, and in a certain sense recreated, in a realm whose operating procedures do not rely on meanings at all! It is recreated as a realm of marks or "tokens," recognized by their physical ("syntactic") characteristics alone and manipulated according to rules that refer only to those physical characteristics (characteristics such as the shape of the symbol—see Box 1.2). As Newell and Simon comment:

Logic . . . was a game played with meaningless tokens according to certain purely syntactic rules. Thus progress was first made by walking away from all that seemed relevant to meaning and human symbols. (Newell and Simon, 1976, p. 43)

Or, to put it in the more famous words of the philosopher John Haugeland:

If you take care of the syntax, *the semantics will take care of itself.* (Haugeland, 1981a, p. 23, original emphasis)

This shift from meaning to form (from semantics to syntax if you will) also begins to suggest an attractive liberalism concerning actual physical structure. For what matters, as far as the identity of these formal systems is concerned, is not, e.g., the precise shape of the symbol for "and." The shape could be "AND" or "and" or "&" or "A" or whatever. All that matters is that the shape is used consistently and that the rules are set up so as to specify how to treat strings of symbols joined by that shape: to allow, for example, the derivation of "A" from the string "A and

## Box 1.2

## SYNTAX AND SEMANTICS

Semantic properties are the "meaning-involving" properties of words, sentences, and internal representations. Syntactic properties, at least as philosophers tend to use the term, are nonsemantic properties of, e.g., written or spoken words, or of any kinds of inscriptions of meaningful items (e.g., the physical states that the pocket calculator uses to store a number in memory). Two synonymous written words ("dog" and "chien") are thus semantically identical but syntactically distinct, whereas ambiguous words ("bank" as in river or "bank" as in high street) are syntactically identical but semantically distinct. The idea of a *token* is the idea of a *specific syntactic item* (e.g., this occurrence of the word "dog"). A pocket calculator manipulates physical tokens (inner syntactic states) to which the operation of the device is sensitive. It is by being sensitive to the *distinct syntactic features* of the inner tokens that the calculator manages to behave in an arithmetic-respecting fashion: it is set up *precisely* so that syntax-driven operations on inner tokens standing for numbers respect meaningful arithmetical relations between the numbers. Taking care of the syntax, in Haugeland's famous phrase, thus allows the semantics to take care of itself.

B." Logics are thus first-rate examples of *formal systems* in the sense of Haugeland (1981a, 1997). They are systems whose essence lies not in the *precise physical details* but in the web of legal moves and transitions.

Most games, Haugeland notes, are formal systems in exactly this sense. You can play chess on a board of wood or marble, using pieces shaped like animals, movie stars, or the crew of the star ship Enterprise. You could even, Haugeland suggests, play chess using helicopters as pieces and a grid of helipads on top of tall buildings as the board. All that matters is *again the web of legal moves* and the physical distinguishability of the tokens.

Thinking about formal systems thus liberates us in two very powerful ways at a single stroke. Semantic relations (such as *truth preservation*: if "A and B" is true, "A" is true) are seen to be respected in virtue of procedures that make no intrinsic reference to meanings. And the specific physical details of any such system are seen to be unimportant, *since what matters is the golden web of moves and transitions*. Semantics is thus made unmysterious without making it brute physical. Who says you can't have your cake and eat it?

The next big development was the formalization (Turing, 1936) of the notion of computation itself. Turing's work, which predates the development of the dig-

ital computer, introduced the foundational notion of (what has since come to be known as) the Turing machine. This is an imaginary device consisting of an infinite tape, a simple processor (a "finite state machine"), and a read/write head. The tape acts as data store, using some fixed set of symbols. The read/write head can read a symbol off the tape, move itself one square backward or forward on the tape, and write onto the tape. The finite state machine (a kind of central processor) has enough memory to recall what symbol was just read and what state it (the finite state machine) was in. These two facts together determine the next action, which is carried out by the read/write head, and determine also the next state of the finite state machine. What Turing showed was that some such device, performing a sequence of simple computations governed by the symbols on the tape, could compute the answer to any sufficiently well-specified problem (see Box 1.3).

We thus confront a quite marvelous confluence of ideas. Turing's work clearly suggested the notion of a physical machine whose *syntax-following properties* would enable it to solve any well-specified problem. Set alongside the earlier work on logics and formal systems, this amounted to nothing less than

... the emergence of a new level of analysis, independent of physics yet mechanistic in spirit . . . a science of structure and function divorced from material substance. (Pylyshyn, 1986, p. 68)

This was classical cognitive science conceived. The vision finally became flesh, however, only because of a third (and final) innovation: the actual construction of general purpose electronic computing machinery and the development of flexible, high-level programming techniques. The bedrock machinery (the digital computer) was designed by John von Neumann in the 1940s and with its advent all the pieces seemed to fall finally into place. For it was now clear that once realized in the physical medium of an electronic computer, a formal system could run *on its own*, without a human being sitting there deciding how and when to apply the rules to initiate the legal transformations. The well-programmed electronic computer, as John Haugeland nicely points out, is really just an automatic ("self-moving") formal system:

It is like a chess set that sits there and plays chess by itself, without any intervention from the players, or an automatic formal system that writes out its own proofs and theorems without any help from the mathematician. (Haugeland, 1981a, p. 10; also Haugeland, 1997, pp. 11–12)

Of course, the machine needs a program. And programs were, in those days (but see Chapter 4), written by good old-fashioned human beings. But once the program was in place, and the power on, the machine took care of the rest. The transitions between legal syntactic states (states that also, under interpretation, meant something) no longer required a human operator. The physical world suddenly included clear, nonevolved, nonorganic examples of what Daniel Dennett would later dub "syntactic engines"—quasiautonomous systems whose sheer physical make-

## Box 1.3

## A TURING MACHINE

To make the idea of Turing machine computation concrete, let us borrow an example from Kim (1996, pp. 80–85). Suppose the goal is to get a Turing machine to add positive numbers. Express the numbers to be added as a sequence of the symbols “#” (marking the beginning and end of numbers) “1” and “+.” So the sum  $3 + 2$  is encoded on the tape as shown in Figure 1.1. A neat program for adding the numbers (where “A” indicates the initial location and initial state of the read/write head) is as follows:

Instruction 1: If read-write head is in machine state A and encounters a “1,” it moves one square to the right, and the head stays in state A.

Instruction 2: If the head is in state A and encounters a “+,” it replaces it with a “1,” stays in state A, and moves one square to the right.

Instruction 3: If the head is in state A and it encounters a “#,” move one square left and go into machine state B.

Instruction 4: If the head is in machine state B and encounters a “1,” delete it, replace with a “#,” and halt.

You should be able to see how this works. Basically, the machine starts “pointed” at the leftmost “1.” It scans right seeking a “+,” which it replaces with a “1.” It continues scanning right until the “#” indicates the end of the sum, at which point it moves one square left, deletes a single “1,” and replaces it with a “#.” The tape now displays the answer to the addition problem in the same notation used to encode the question, as shown in Figure 1.2.

Similar set-ups (try to imagine how they work) can do subtraction, multiplication, and more (see Kim, 1996, pp. 83–85). But Turing’s most strik-

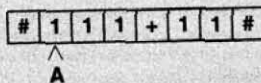


Figure 1.1 (After Kim, 1996, p. 81.)

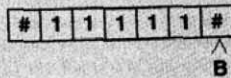


Figure 1.2 (After Kim, 1996, p. 81.)

ing achievement in this area was to show that you could then define a special kind of Turing machine (the aptly-named universal Turing machine) able to imitate any other Turing machine. The symbols on the tape, in this universal case, encode a description of the behavior of the other machine. The universal Turing machine uses this description to mimic the input-output function of any other such device and hence is itself capable of carrying out any sufficiently well-specified computation. (For detailed accounts see Franklin, 1995; Haugeland, 1985; Turing, 1936, 1950.)

The Turing machine affords a fine example of a simple case in which syntax-driven operations support a semantics-respecting (meaning-respecting) process. Notice also that you could build a simple Turing machine out of many different materials. It is the formal (syntactic) organization that matters for its semantic success.

up ensured (under interpretation) some kind of ongoing reason-respecting behavior. No wonder the early researchers were jubilant! Newell and Simon nicely capture the mood:

It is not my aim to surprise or shock you. . . . But the simplest way I can summarize is to say that there are now in the world machines that think, that learn and that create. Moreover, their ability to do these things is going to increase rapidly until—in a visible future—the range of problems they can handle will be co-extensive with the range to which the human mind has been applied. (Newell and Simon, 1958, p. 6, quoted in Dreyfus and Dreyfus, 1990, p. 312)

This jubilant mood deepened as advanced programming techniques<sup>2</sup> brought forth impressive problem-solving displays, while the broader theoretical and philosophical implications (see Box 1.4) of these early successes could hardly have been more striking. The once-mysterious realm of mindware (represented, admittedly, by just two of its many denizens: truth preservation and abstract problem solving) looked ripe for conquest and understanding. Mind was not ghostly stuff, but the operation of a formal, computational system implemented in the meatware of the brain.

Such is the heart of the matter. Mindware, it was claimed, is to the neural meat machine as software is to the computer. The brain may be the standard (local, earthly, biological) implementation—but cognition is a program-level thing. Mind

<sup>2</sup>For example, list-processing languages, as pioneered in Newell and Simon’s Logic Theorist program in 1956 and perfected in McCarthy’s LISP around 1960, encouraged the use of more complex “recursive programming” strategies in which symbols point to data structures that contain symbols pointing to further data structures and so on. They also made full use of the fact that the same electronic memory could store both program and data, a feature that allowed programs to be modified and operated on in the same ways as data. LISP even boasted a universal function, EVAL, that made it as powerful, modulo finite memory limitations, as a Universal Turing Machine.

## Box 1.4

## MACHINE FUNCTIONALISM

The leading philosophical offspring of the developments in artificial intelligence went by the name of machine functionalism, and it was offered as an answer to one of the deepest questions ever asked by humankind, viz. what is the essence (the deep nature) of the mental? What fundamental facts make it the case that some parts of the physical world have mental lives (thoughts, beliefs, feelings, and all the rest) and others do not? Substance dualists, recall, thought that the answer lay in the presence or absence of a special kind of mental *stuff*. Reacting against this idea (and against so-called philosophical behaviorism—see Appendix I). Mind-brain identity theorists, such as Smart (1959) (and again, see Appendix I), claimed that mental states *just are* processes going on in the brain. This bald identity claim, however, threatened to make the link between mental states and specific, material brain states a little too intimate. A key worry (e.g., Putnam, 1960, 1967) was that if it was really essential to being in a certain mental state that one be in a specific brain state, it would seem to follow that creatures lacking brains built just like ours (say, Martians or silicon-based robots) could not be in those very same mental states. But surely, the intuition went, creatures with very different brains from ours could, at least in principle, share, e.g., the belief that it is raining. Where, then, should we look for the commonality that could unite the robot, the Martian, and the Bostonian? The work in logic and formal systems, Turing machines, and electronic computation now suggested an answer: look not to the specific physical story (of neurons and wetware), nor to the surface behavior, but to the inner organization, that is to say, to the golden web: to the abstract, formal organization of the system. It is this organization—depicted by the machine functionalists as a web of links between possible inputs, inner computational states, and outputs (actions, speech)—that fixes the shape and contents of a mental life. The building materials do not matter: the web of transitions could be realized in flesh, silicon, or cream cheese (Putnam, 1975, p. 291). To be in such and such a mental state is simply to be a physical device, of whatever composition, that satisfies a specific formal description. Mindware, in humans, happens to run on a meat machine. But the very same mindware (as picked out by the web of legal state transitions) might run in some silicon device, or in the alien organic matter of a Martian.

is thus ghostly enough to float fairly free of the gory neuroscientific details. But it is not so ghostly as to escape the nets of more abstract (formal, computational) scientific investigation. This is an appealing story. But is it correct? Let's worry.

## 1.2 Discussion

(A brief note of reassurance: many of the topics treated below recur again and again in subsequent chapters. At this point, we lack much of the detailed background needed to really do them justice. But it is time to test the waters.)

## A. WHY TREAT THOUGHT AS COMPUTATION?

Why treat thought as computation? The principal reason (apart from the fact that it seems to work!) is that thinkers are physical devices whose behavior patterns are reason respecting. Thinkers act in ways that are usefully understood as sensitively guided by reasons, ideas, and beliefs. Electronic computing devices show us one way in which this strange "dual profile" (of physical substance and reason-respecting behavior) can actually come about.

The notion of reason-respecting behavior, however, bears immediate amplification. A nice example of this kind of behavior is given by Zenon Pylyshyn. Pylyshyn (1986) describes the case of the pedestrian who witnesses a car crash, runs to a telephone, and punches out 911. We could, as Pylyshyn notes, try to explain this behavior by telling a purely physical story (maybe involving specific neurons, or even quantum events, whatever). But such a story, Pylyshyn argues, will not help us understand the behavior in its *reason-guided* aspects. For example, suppose we ask: what would happen if the phone was dead, or if it was a dial phone instead of a touch-tone phone, or if the accident occurred in England instead of the United States? The neural story underlying the behavioral response will differ widely if the agent dials 999 (the emergency code in England) and not 911, or must run to find a working phone. Yet common sense psychological talk makes sense of all these options at a stroke by depicting the agent as seeing a crash and *wanting to get help*. What we need, Pylyshyn powerfully suggests, is a scientific story that remains in touch with this more abstract and reason-involving characterization. And the simplest way to provide one is to imagine that the agent's brain contains states ("symbols") that represent the event *as* a car crash and that the computational state-transitions occurring inside the system (realized as physical events in the brain) then lead to new sets of states (more symbols) whose proper interpretation is, e.g., "seek help," "find a telephone," and so on. The interpretations thus glue inner states to sensible real-world behaviors. Cognizers, it is claimed, "instantiate . . . representation physically as cognitive codes and . . . their behavior is a causal consequence of operations carried out on those codes" (Pylyshyn, 1986, p. xiii).

The same argument can be found in, e.g., Fodor (1987), couched as a point about content-determined transitions in trains of thought, as when the thought "it

is raining" leads to the thought "let's go indoors." This, for Fodor (but see Chapters 4 onward), is the essence of human rationality. How is such rationality mechanically possible? A good empirical hypothesis, Fodor suggests, is that there are neural symbols (inner states apt for interpretation) that mean, e.g., "it is raining" and whose physical properties lead in context to the generation of other symbols that mean "let's go indoors." If that is how the brain works then the brain is indeed a computer in exactly the sense displayed earlier. And if such were the case, then the mystery concerning reason-guided (content-determined) transitions in thought is resolved:

If the mind is a sort of computer, we begin to see how . . . there could be non-arbitrary content-relations among causally related thoughts. (Fodor, 1987, p. 19)

Such arguments aim to show that the mind *must* be understood as a kind of computer implemented in the wetware of the brain, on pain of failing empirically to account for rational transitions among thoughts. Reason-guided action, it seems, makes good scientific sense if we imagine a neural economy organized as a syntax-driven engine that tracks the shape of semantic space (see, e.g., Fodor, 1987, pp. 19–20).

#### B. IS SOFTWARE AN AUTONOMOUS LEVEL IN NATURE?

The mindware/software equation is as beguiling as it is, at times, distortive. One immediate concern is that all this emphasis on algorithms, symbols, and programs tends to promote a somewhat misleading vision of *crisp level distinctions in nature*. The impact of the theoretical independence of algorithms from hardware is an artifact of the long-term neglect of issues concerning *real-world action taking* and the time course of computations. For an algorithm or program as such is just a sequence of steps with no inbuilt relation to real-world timing. Such timing depends crucially on the particular way in which the algorithm is implemented on a real device. Given this basic fact, the theoretical independence of algorithm from hardware is unlikely to have made much of an impact on Nature. We must expect to find biological computational strategies closely tailored to getting useful real-time results from available, slow, wetware components. In practice, it is thus unlikely that we will be able to fully appreciate the formal organization of natural systems without some quite detailed reference to the nature of the neural hardware that provides the supporting implementation. In general, attention to the nature of real biological hardware looks likely to provide both important clues about and constraints on the kinds of computational strategy used by real brains. This topic is explored in more depth in Chapters 4 through 6.

Furthermore, the claim that mindware is software is—to say the least—merely schematic. For the space of possible types of explanatory story, all broadly computational (but see Box 1.5), is very large indeed. The comments by Fodor and by

#### Box 1.5

### WHAT IS COMPUTATION?

It is perhaps worth mentioning that the foundational notion of computation is itself still surprisingly ill understood. What do we really mean by calling some phenomenon "computational" in the first place? There is no current consensus at least (in the cognitive scientific community) concerning the answer to this question. It is mostly a case of "we know one when we see one." Nonetheless, there is a reasonable consensus concerning what I'll dub the "basic profile," which is well expressed by the following statement:

we count something as a computer because, and only when, its inputs and outputs can be usefully and systematically interpreted as representing the ordered pairs of some function that interests us. (Churchland and Sejnowski, 1992, p. 65)

Thus consider a pocket calculator. This physical device computes, on this account, because first, there is a reliable and systematic way of interpreting various states of the device (the marks and numerals on the screen and keyboard) as representing other things (numbers). And second, because the device is set up so that under that interpretation, its physical state changes mirror semantic (meaningful) transitions in the arithmetical domain. Its physical structure thus forces it to respect mathematical constraints so that inputs such as "4 × 3" lead to outputs such as "12" and so on.

A truly robust notion of the conditions under which some actual phenomenon counts as computational would require, however, some rather more *objective* criterion for determining when an encountered (nondesignated) physical process is actually implementing a computation—some criterion that does not place our interpretive activities and interests so firmly at center stage.

The best such account I know of is due to Dave Chalmers (1996, Chapter 9). Chalmers' goal is to give an "objective criterion for implementing a computation" (p. 319). Intuitively, a physical device 'implements' an abstract, formal computational specification just in case the physical device is set up to undergo state changes that march in step with those detailed in the specification. In this sense a specific word-processing program might, for example, constitute a formal specification that can (appropriately configured) be made to run on various kinds of physical device (MACS, PCs, etc.).

Chalmers' proposal, in essence, is that a physical device implements an abstract formal description (a specification of states and state-transition relations) just in case "the causal structure of the system mirrors the formal